

# Executive Summary of Artificial Intelligence and Emerging Technologies

Prepared by: Dr. Mejdal Alqahtani

Linked in.





## Table of contents

- 1. Overview on Al
- 2. Basic math and probabilities
- 3. Machine learning
  - 3.1. Supervised learning
  - 3.2. Unsupervised learning
- 4. Deep learning
- 5. Reinforcement learning
- 6. Programing languages
  - 6.1. Python language
  - 6.2. R language
  - 6.3. SQL language
- 7. Business Intelligence
  - 7.1. Tableau
  - 7.2. Power BI
  - 7.3. Data Visualization



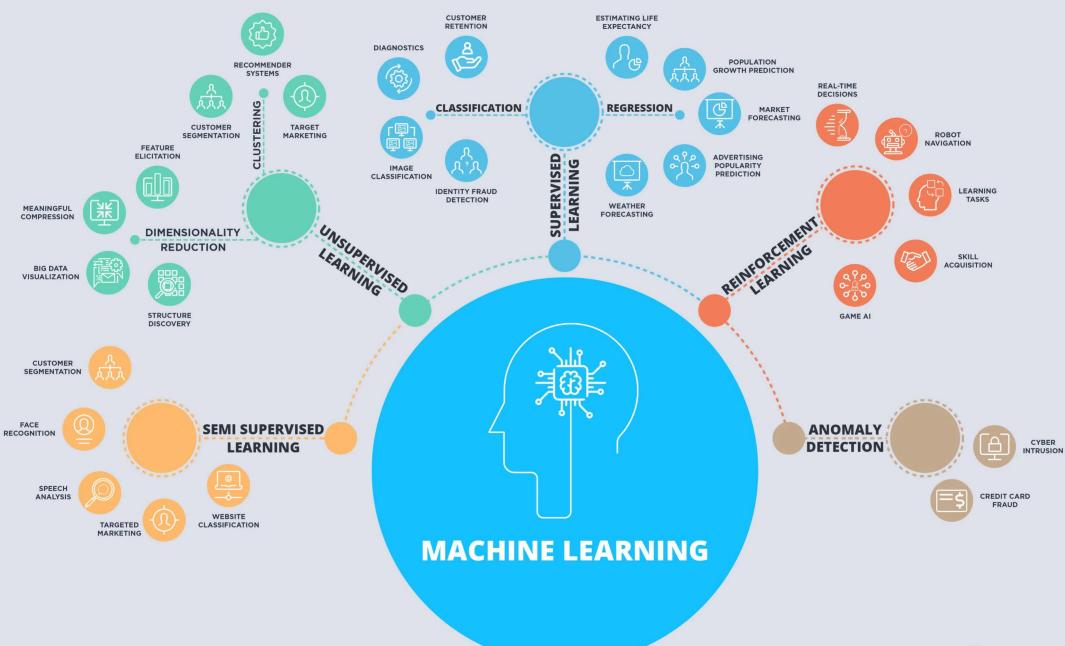
- 8.1. Artificial Intelligence (AI)
- 8.2. Internet of Things (IoT)
- 8.3. Robotic Process Automation (RPA)
- 8.4. Cloud Computing
- 8.5. Quantum Computing
- 8.6. 3D Printing
- 8.7. 5th Generation Network (5G)
- 8.8. Extended Reality (XR)
- 8.9. Blockchain
- 8.10. Cyber Security



1. Overview on Al



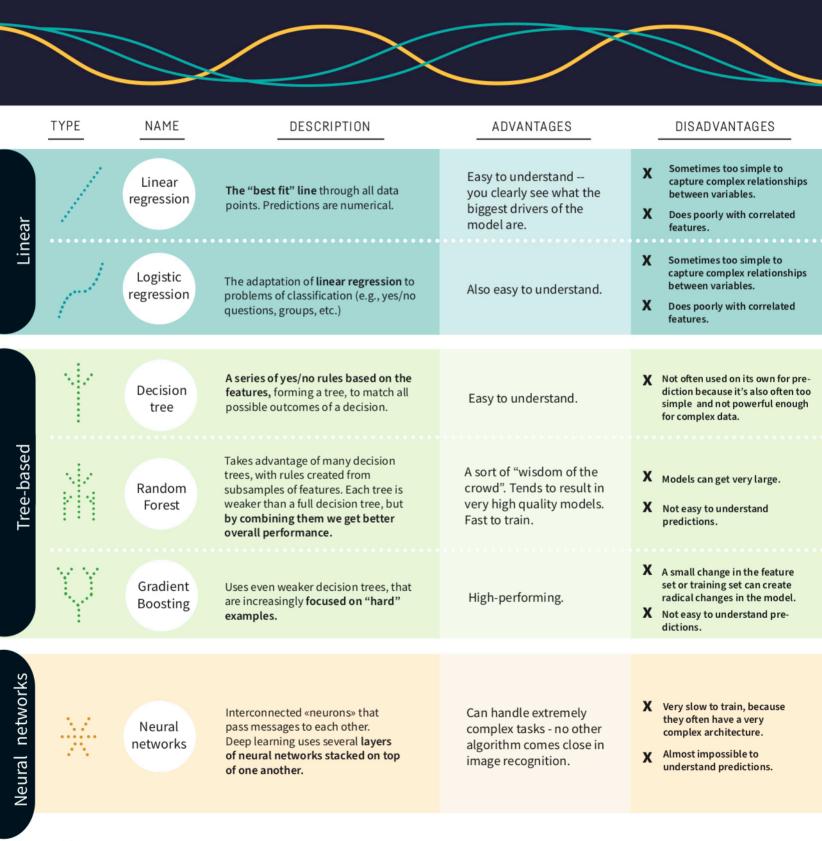




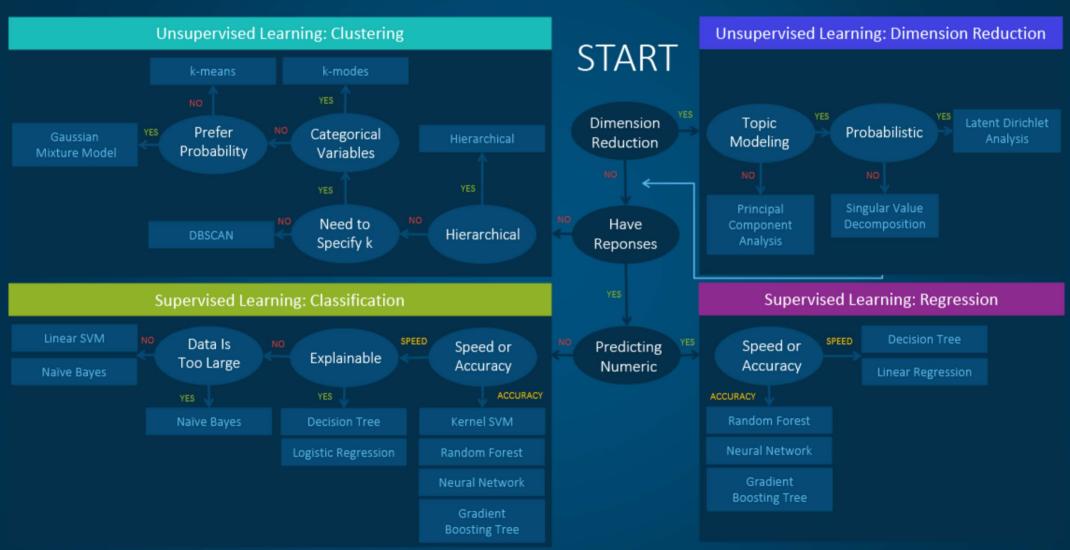
dataschencedojo



## TOP PREDICTION ALGORITHMS

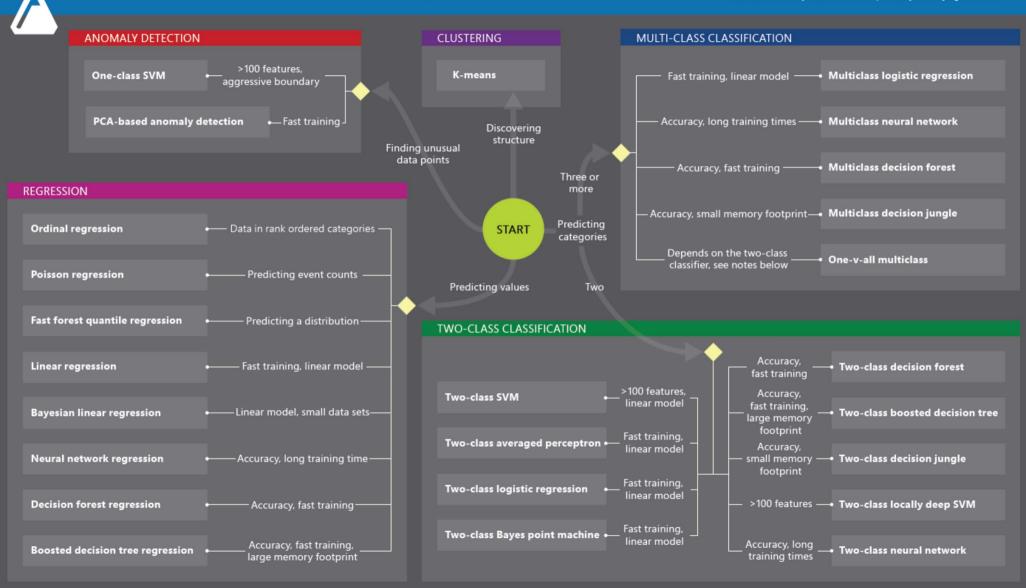


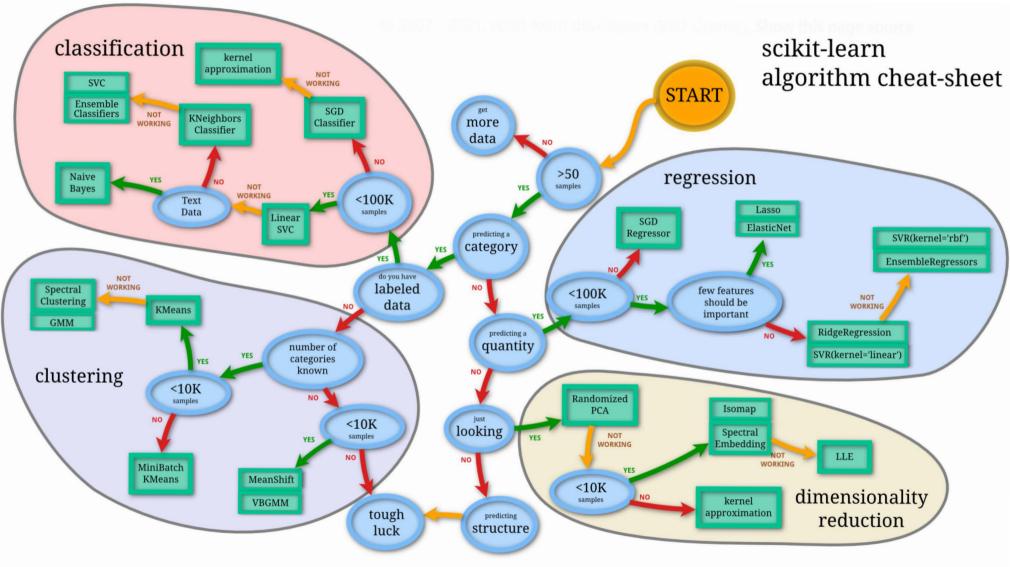
# Machine Learning Algorithms Cheat Sheet



## Microsoft Azure Machine Learning: Algorithm Cheat Sheet

This cheat sheet helps you choose the best Azure Machine Learning Studio algorithm for your predictive analytics solution. Your decision is driven by both the nature of your data and the question you're trying to answer.







# 2. Basic math and probabilities

## VIP Refresher: Linear Algebra and Calculus

Afshine AMIDI and Shervine AMIDI

General notations

 $\Box$  Vector – We note  $x \in \mathbb{R}^n$  a vector with n entries, where  $x_i \in \mathbb{R}$  is the  $i^{th}$  entry:

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^n$$

**D** Matrix – We note  $A \in \mathbb{R}^{m \times n}$  a matrix with *m* rows and *n* columns, where  $A_{i,j} \in \mathbb{R}$  is the entry located in the  $i^{th}$  row and  $j^{th}$  column:

$$A = \begin{pmatrix} A_{1,1} & \cdots & A_{1,n} \\ \vdots & & \vdots \\ A_{m,1} & \cdots & A_{m,n} \end{pmatrix} \in \mathbb{R}^{m \times n}$$

Remark: the vector x defined above can be viewed as a  $n \times 1$  matrix and is more particularly called a column-vector.

**D** Identity matrix – The identity matrix  $I \in \mathbb{R}^{n \times n}$  is a square matrix with ones in its diagonal and zero everywhere else:

$$I = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 1 \end{pmatrix}$$

Remark: for all matrices  $A \in \mathbb{R}^{n \times n}$ , we have  $A \times I = I \times A = A$ .

**Diagonal matrix** – A diagonal matrix  $D \in \mathbb{R}^{n \times n}$  is a square matrix with nonzero values in its diagonal and zero everywhere else:

$$D = \begin{pmatrix} d_1 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & d_n \end{pmatrix}$$

Remark: we also note D as  $diag(d_1,...,d_n)$ .

#### Matrix operations

□ Vector-vector multiplication – There are two types of vector-vector products:

• inner product: for  $x, y \in \mathbb{R}^n$ , we have:

$$x^T y = \sum_{i=1}^n x_i y_i \in \mathbb{R}$$

• outer product: for  $x \in \mathbb{R}^m, y \in \mathbb{R}^n$ , we have:

$$xy^{T} = \begin{pmatrix} x_{1}y_{1} & \cdots & x_{1}y_{n} \\ \vdots & & \vdots \\ x_{m}y_{1} & \cdots & x_{m}y_{n} \end{pmatrix} \in \mathbb{R}^{m \times n}$$

**D** Matrix-vector multiplication – The product of matrix  $A \in \mathbb{R}^{m \times n}$  and vector  $x \in \mathbb{R}^n$  is a vector of size  $\mathbb{R}^m$ , such that:

$$Ax = \begin{pmatrix} a_{r,1}^T x \\ \vdots \\ a_{r,m}^T x \end{pmatrix} = \sum_{i=1}^n a_{c,i} x_i \in \mathbb{R}^m$$

where  $a_{r,i}^T$  are the vector rows and  $a_{c,j}$  are the vector columns of A, and  $x_i$  are the entries of x.

**D** Matrix-matrix multiplication – The product of matrices  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{n \times p}$  is a matrix of size  $\mathbb{R}^{n \times p}$ , such that:

$$AB = \begin{pmatrix} a_{r,1}^T b_{c,1} & \cdots & a_{r,1}^T b_{c,p} \\ \vdots & \vdots \\ a_{r,m}^T b_{c,1} & \cdots & a_{r,m}^T b_{c,p} \end{pmatrix} = \sum_{i=1}^n a_{c,i} b_{r,i}^T \in \mathbb{R}^{n \times p}$$

where  $a_{r,i}^T, b_{r,i}^T$  are the vector rows and  $a_{c,j}, b_{c,j}$  are the vector columns of A and B respectively.

 $\Box$  Transpose – The transpose of a matrix  $A \in \mathbb{R}^{m \times n}$ , noted  $A^T$ , is such that its entries are flipped:

$$\forall i,j, \qquad A_{i,j}^T = A_{j,i}$$

Remark: for matrices A, B, we have  $(AB)^T = B^T A^T$ .

 $\square$  Inverse – The inverse of an invertible square matrix A is noted  $A^{-1}$  and is the only matrix such that:

$$AA^{-1} = A^{-1}A = I$$

Remark: not all square matrices are invertible. Also, for matrices A,B, we have  $(AB)^{-1} = B^{-1}A^{-1}$ 

 $\Box$  Trace – The trace of a square matrix A, noted tr(A), is the sum of its diagonal entries:

$$\operatorname{tr}(A) = \sum_{i=1}^{n} A_{i,i}$$

Remark: for matrices A,B, we have  $tr(A^T) = tr(A)$  and tr(AB) = tr(BA)

□ Determinant – The determinant of a square matrix  $A \in \mathbb{R}^{n \times n}$ , noted |A| or det(A) is expressed recursively in terms of  $A_{\backslash i, \backslash j}$ , which is the matrix A without its  $i^{th}$  row and  $j^{th}$  column, as follows:

$$\det(A) = |A| = \sum_{j=1}^{n} (-1)^{i+j} A_{i,j} |A_{\backslash i, \backslash j}|$$

Remark: A is invertible if and only if  $|A| \neq 0$ . Also, |AB| = |A||B| and  $|A^T| = |A|$ .

#### Matrix properties

**Symmetric decomposition** – A given matrix A can be expressed in terms of its symmetric and antisymmetric parts as follows:

$$A = \underbrace{\frac{A + A^T}{2}}_{\text{Symmetric}} + \underbrace{\frac{A - A^T}{2}}_{\text{Antisymmetric}}$$

**D** Norm – A norm is a function  $N: V \longrightarrow [0, +\infty[$  where V is a vector space, and such that for all  $x, y \in V$ , we have:

- $N(x+y) \leq N(x) + N(y)$
- N(ax) = |a|N(x) for a scalar
- if N(x) = 0, then x = 0

For  $x \in V$ , the most commonly used norms are summed up in the table below:

Norm	Notation	Definition	Use case
Manhattan, $L^1$	$  x  _1$	$\sum_{i=1}^{n}  x_i $	LASSO regularization
Euclidean, $L^2$	$  x  _{2}$	$\sqrt{\sum_{i=1}^n x_i^2}$	Ridge regularization
$p$ -norm, $L^p$	$  x  _p$	$\sum_{i=1}^n x_i^p \right)^{\frac{1}{p}}$	Hölder inequality
Infinity, $L^{\infty}$	$  x  _{\infty}$	$\max_i  x_i $	Uniform convergence

**Linearly dependence** – A set of vectors is said to be linearly dependent if one of the vectors in the set can be defined as a linear combination of the others.

Remark: if no vector can be written this way, then the vectors are said to be linearly independent.

**D** Matrix rank – The rank of a given matrix A is noted rank(A) and is the dimension of the vector space generated by its columns. This is equivalent to the maximum number of linearly independent columns of A.

**Dositive semi-definite matrix** – A matrix  $A \in \mathbb{R}^{n \times n}$  is positive semi-definite (PSD) and is noted  $A \succeq 0$  if we have:

$$\boxed{A = A^T} \quad \text{and} \quad \forall x \in \mathbb{R}^n, \quad x^T A x \ge 0$$

Remark: similarly, a matrix A is said to be positive definite, and is noted  $A \succ 0$ , if it is a PSD matrix which satisfies for all non-zero vector x,  $x^T A x > 0$ .

**D** Eigenvalue, eigenvector – Given a matrix  $A \in \mathbb{R}^{n \times n}$ ,  $\lambda$  is said to be an eigenvalue of A if there exists a vector  $z \in \mathbb{R}^n \setminus \{0\}$ , called eigenvector, such that we have:

 $Az = \lambda z$ 

**Spectral theorem** – Let  $A \in \mathbb{R}^{n \times n}$ . If A is symmetric, then A is diagonalizable by a real orthogonal matrix  $U \in \mathbb{R}^{n \times n}$ . By noting  $\Lambda = \text{diag}(\lambda_1, ..., \lambda_n)$ , we have:

 $\exists \Lambda \text{ diagonal}, \quad A = U \Lambda U^T$ 

□ Singular-value decomposition – For a given matrix A of dimensions  $m \times n$ , the singularvalue decomposition (SVD) is a factorization technique that guarantees the existence of  $U \ m \times m$ unitary,  $\Sigma \ m \times n$  diagonal and  $V \ n \times n$  unitary matrices, such that:

 $A = U \Sigma V^T$ 

#### Matrix calculus

**Gradient** – Let  $f : \mathbb{R}^{m \times n} \to \mathbb{R}$  be a function and  $A \in \mathbb{R}^{m \times n}$  be a matrix. The gradient of f with respect to A is a  $m \times n$  matrix, noted  $\nabla_A f(A)$ , such that:

$$\left(\nabla_A f(A)\right)_{i,j} = \frac{\partial f(A)}{\partial A_{i,j}}$$

Remark: the gradient of f is only defined when f is a function that returns a scalar.

**D** Hessian – Let  $f : \mathbb{R}^n \to \mathbb{R}$  be a function and  $x \in \mathbb{R}^n$  be a vector. The hessian of f with respect to x is a  $n \times n$  symmetric matrix, noted  $\nabla_x^2 f(x)$ , such that:

$$\left(\nabla_x^2 f(x)\right)_{i,j} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j}$$

Remark: the hessian of f is only defined when f is a function that returns a scalar.

**Gradient operations** – For matrices A, B, C, the following gradient properties are worth having in mind:

$$\nabla_A \operatorname{tr}(AB) = B^T \qquad \nabla_{A^T} f(A) = (\nabla_A f(A))^T$$
$$\nabla_A \operatorname{tr}(ABA^T C) = CAB + C^T A B^T \qquad \nabla_A |A| = |A| (A^{-1})^T$$

## VIP Refresher: Probabilities and Statistics

Afshine AMIDI and Shervine AMIDI

August 6, 2018

#### Introduction to Probability and Combinatorics

 $\Box$  Sample space – The set of all possible outcomes of an experiment is known as the sample space of the experiment and is denoted by S.

 $\Box$  Event – Any subset E of the sample space is known as an event. That is, an event is a set consisting of possible outcomes of the experiment. If the outcome of the experiment is contained in E, then we say that E has occurred.

 $\Box$  Axioms of probability – For each event E, we denote P(E) as the probability of event E occuring. By noting  $E_1, \ldots, E_n$  mutually exclusive events, we have the 3 following axioms:

(1) 
$$0 \leq P(E) \leq 1$$
 (2)  $P(S) = 1$  (3)  $P = \bigcup_{i=1}^{n} E_i = \sum_{i=1}^{n} P(E_i)$ 

**Dermutation** – A permutation is an arrangement of r objects from a pool of n objects, in a given order. The number of such arrangements is given by P(n, r), defined as:

$$P(n,r) = \frac{n!}{(n-r)!}$$

**Combination** – A combination is an arrangement of r objects from a pool of n objects, where the order does not matter. The number of such arrangements is given by C(n, r), defined as:

$$C(n,r) = \frac{P(n,r)}{r!} = \frac{n!}{r!(n-r)!}$$

Remark: we note that for  $0 \leq r \leq n$ , we have  $P(n,r) \geq C(n,r)$ .

#### **Conditional Probability**

**D** Bayes' rule – For events A and B such that P(B) > 0, we have:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Remark: we have  $P(A \cap B) = P(A)P(B|A) = P(A|B)P(B)$ .

□ **Partition** – Let  $\{A_i, i \in [\![1,n]\!]\}$  be such that for all  $i, A_i \neq \emptyset$ . We say that  $\{A_i\}$  is a partition if we have:

$$\forall i \neq j, A_i \cap A_j = \emptyset \quad \text{and} \quad \bigcup_{i=1}^n A_i = S$$

Remark: for any event B in the sample space, we have 
$$P(B) = \sum_{i=1}^{n} P(B|A_i)P(A_i)$$
.

**D** Extended form of Bayes' rule – Let  $\{A_i, i \in [\![1,n]\!]\}$  be a partition of the sample space. We have:

$$P(A_k|B) = \frac{P(B|A_k)P(A_k)}{\sum_{i=1}^{n} P(B|A_i)P(A_i)}$$

 $\Box$  Independence – Two events A and B are independent if and only if we have:

$$P(A \cap B) = P(A)P(B)$$

#### Random Variables

 $\Box$  Random variable – A random variable, often noted X, is a function that maps every element in a sample space to a real line.

□ Cumulative distribution function (CDF) – The cumulative distribution function F, which is monotonically non-decreasing and is such that  $\lim_{x\to-\infty} F(x) = 0$  and  $\lim_{x\to+\infty} F(x) = 1$ , is defined as:

$$F(x) = P(X \leqslant x)$$

Remark: we have  $P(a < X \leq B) = F(b) - F(a)$ .

 $\Box$  **Probability density function (PDF)** – The probability density function f is the probability that X takes on values between two adjacent realizations of the random variable.

 $\Box$  **Relationships involving the PDF and CDF** – Here are the important properties to know in the discrete (D) and the continuous (C) cases.

Case	CDF F	<b>PDF</b> $f$	Properties of PDF
(D)	$F(x) = \sum_{x_i \leqslant x} P(X = x_i)$	$f(x_j) = P(X = x_j)$	$0 \leqslant f(x_j) \leqslant 1$ and $\sum_j f(x_j) = 1$
(C)	$F(x) = \int_{-\infty}^{x} f(y) dy$	$f(x) = \frac{dF}{dx}$	$f(x) \ge 0$ and $\int_{-\infty}^{+\infty} f(x)dx = 1$

 $\Box$  Variance – The variance of a random variable, often noted Var(X) or  $\sigma^2$ , is a measure of the spread of its distribution function. It is determined as follows:

$$Var(X) = E[(X - E[X])^2] = E[X^2] - E[X]^2$$

 $\Box$  Standard deviation – The standard deviation of a random variable, often noted  $\sigma$ , is a measure of the spread of its distribution function which is compatible with the units of the actual random variable. It is determined as follows:

$$\sigma = \sqrt{\operatorname{Var}(X)}$$

value E[X], generalized expected value E[q(X)],  $k^{th}$  moment  $E[X^k]$  and characteristic function function  $f_{XY}$ , we have:  $(\omega)$  for the discrete and continuous cases:

C	ase	E[X]	E[g(X)]	$E[X^k]$	(ω)
(]	D)	$\sum_{i=1}^{n} x_i f(x_i)$	$\sum_{i=1}^{n} g(x_i) f(x_i)$	$\sum_{i=1}^{n} x_i^k f(x_i)$	$\sum_{i=1}^{n} f(x_i) e^{i\omega x_i}$
(	C)	$\int_{-\infty}^{+\infty} x f(x) dx$	$\int_{-\infty}^{+\infty} g(x)f(x)dx$	$\int_{-\infty}^{+\infty} x^k f(x) dx$	$\int_{-\infty}^{+\infty} f(x)e^{i\omega x}dx$

Remark: we have  $e^{i\omega x} = \cos(\omega x) + i\sin(\omega x)$ .

 $\Box$  Revisiting the  $k^{th}$  moment – The  $k^{th}$  moment can also be computed with the characteristic function as follows:

$$E[X^k] = \frac{1}{i^k} \left[ \frac{\partial^k \psi}{\partial \omega^k} \right]_{\omega = 0}$$

 $\Box$  Transformation of random variables – Let the variables X and Y be linked by some function. By noting  $f_X$  and  $f_Y$  the distribution function of X and Y respectively, we have:

$$f_Y(y) = f_X(x) \left| \frac{dx}{dy} \right|$$

 $\Box$  Leibniz integral rule – Let q be a function of x and potentially c, and a, b boundaries that may depend on c. We have:

$$\left| \frac{\partial}{\partial c} \left( \int_{a}^{b} g(x) dx \right) = \frac{\partial b}{\partial c} \cdot g(b) - \frac{\partial a}{\partial c} \cdot g(a) + \int_{a}^{b} \frac{\partial g}{\partial c}(x) dx \right|$$

 $\Box$  Chebyshev's inequality – Let X be a random variable with expected value  $\mu$  and standard deviation  $\sigma$ . For  $k, \sigma > 0$ , we have the following inequality:

$$P(|X - \mu| \ge k\sigma) \le \frac{1}{k^2}$$

Jointly Distributed Random Variables

 $\Box$  Conditional density – The conditional density of X with respect to Y, often noted  $f_{X|Y}$ , is defined as follows:

$$f_{X|Y}(x) = \frac{f_{XY}(x,y)}{f_Y(y)}$$

 $\Box$  Independence – Two random variables X and Y are said to be independent if we have:

$$f_{XY}(x,y) = f_X(x)f_Y(y)$$

🗆 Expectation and Moments of the Distribution – Here are the expressions of the expected 🛛 Marginal density and cumulative distribution – From the joint density probability

CaseMarginal densityCumulative function(D)
$$f_X(x_i) = \sum_j f_{XY}(x_i, y_j)$$
 $F_{XY}(x, y) = \sum_{x_i \leqslant x} \sum_{y_j \leqslant y} f_{XY}(x_i, y_j)$ (C) $f_X(x) = \int_{-\infty}^{+\infty} f_{XY}(x, y) dy$  $F_{XY}(x, y) = \int_{-\infty}^{x} \int_{-\infty}^{y} f_{XY}(x', y') dx' dy'$ 

 $\Box$  Distribution of a sum of independent random variables – Let  $Y = X_1 + ... + X_n$  with  $X_1, \ldots, X_n$  independent. We have:

$$_{Y}(\omega) = \prod_{k=1}^{n} \psi_{X_{k}}(\omega)$$

 $\Box$  Covariance – We define the covariance of two random variables X and Y, that we note  $\sigma_{XY}^2$ or more commonly  $\operatorname{Cov}(X,Y)$ , as follows:

$$\operatorname{Cov}(X,Y) \triangleq \sigma_{XY}^2 = E[(X - \mu_X)(Y - \mu_Y)] = E[XY] - \mu_X \mu_Y$$

 $\Box$  Correlation – By noting  $\sigma_X, \sigma_Y$  the standard deviations of X and Y, we define the correlation between the random variables X and Y, noted  $\rho_{XY}$ , as follows:

$$\rho_{XY} = \frac{\sigma_{XY}^2}{\sigma_X \sigma_Y}$$

Remarks: For any X, Y, we have  $\rho_{XY} \in [-1,1]$ . If X and Y are independent, then  $\rho_{XY} = 0$ .

**Main distributions** – Here are the main distributions to have in mind:

Type	Distribution	PDF	$(\omega)$	E[X]	$\operatorname{Var}(X)$
(D)	$\begin{aligned} X \sim \mathcal{B}(n,p) \\ \text{Binomial} \end{aligned}$	$\begin{split} P(X = x) &= \binom{n}{x} p^x q^{n-x} \\ x &\in \llbracket 0, n \rrbracket \end{split}$	$(pe^{i\omega}+q)^n$	np	npq
	$\begin{array}{l} X \sim \mathrm{Po}(\mu) \\ \mathrm{Poisson} \end{array}$	$P(X = x) = \frac{\mu^x}{x!}e^{-\mu}$ $x \in \mathbb{N}$	$e^{\mu(e^{i\omega}-1)}$	μ	$\mu$
	$X \sim \mathcal{U}(a, b)$ Uniform	$f(x) = \frac{1}{b-a}$ $x \in [a,b]$	$\frac{e^{i\omega b} - e^{i\omega a}}{(b-a)i\omega}$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$
(C)	$\begin{aligned} X \sim \mathcal{N}(\mu, \sigma) \\ \text{Gaussian} \end{aligned}$	$x \in [a,b]$ $f(x) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$ $x \in \mathbb{R}$	$e^{i\omega\mu - \frac{1}{2}\omega^2\sigma^2}$	μ	$\sigma^2$
	$X \sim \operatorname{Exp}(\lambda)$ Exponential	$f(x) = \lambda e^{-\lambda x}$ $x \in \mathbb{R}_+$	$\frac{1}{1 - \frac{i\omega}{\lambda}}$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$

### Parameter estimation

 $\Box$  Random sample – A random sample is a collection of *n* random variables  $X_1, ..., X_n$  that are independent and identically distributed with *X*.

 $\Box$  Estimator – An estimator  $\hat{\theta}$  is a function of the data that is used to infer the value of an unknown parameter  $\theta$  in a statistical model.

**D** Bias – The bias of an estimator  $\hat{\theta}$  is defined as being the difference between the expected value of the distribution of  $\hat{\theta}$  and the true value, i.e.:

$$\operatorname{Bias}(\hat{\theta}) = E[\hat{\theta}] - \theta$$

Remark: an estimator is said to be unbiased when we have  $E[\hat{\theta}] = \theta$ .

**D** Sample mean and variance – The sample mean and the sample variance of a random sample are used to estimate the true mean  $\mu$  and the true variance  $\sigma^2$  of a distribution, are noted  $\overline{X}$  and  $s^2$  respectively, and are such that:

$$\overline{X} = \frac{1}{n} \sum_{i=1}^{n} X_i \quad \text{and} \quad s^2 = \hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^{n} (X_i - \overline{X})^2$$

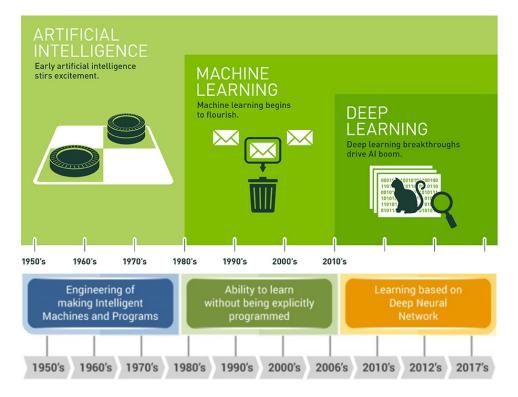
**Central Limit Theorem** – Let us have a random sample  $X_1, ..., X_n$  following a given distribution with mean  $\mu$  and variance  $\sigma^2$ , then we have:

$$\overline{X} \underset{n \to +\infty}{\sim} \mathcal{N}\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

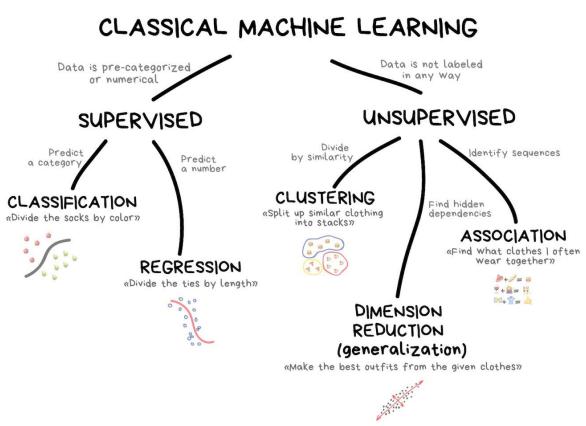


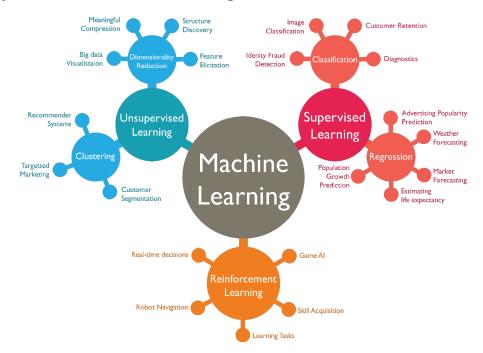
# 3. Machine learning

## • History of Al



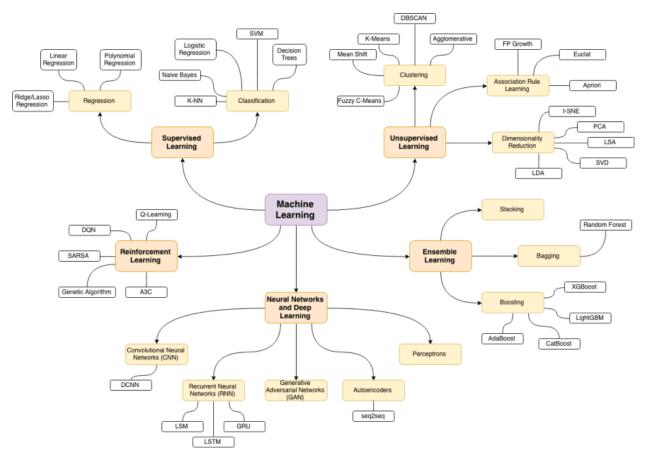
• Structure of Machine learning





• Applications of Machine learning

• Top Algorithms of machine and deep learning





# 3.1. Supervised learning

## VIP Cheatsheet: Supervised Learning

Afshine AMIDI and Shervine AMIDI

September 9, 2018

### Introduction to Supervised Learning

Given a set of data points  $\{x^{(1)}, ..., x^{(m)}\}$  associated to a set of outcomes  $\{y^{(1)}, ..., y^{(m)}\}$ , we want to build a classifier that learns how to predict y from x.

 $\square$  Type of prediction – The different types of predictive models are summed up in the table below:

	Regression	Classifier	
Outcome	Continuous	Class	
Examples	Linear regression	Logistic regression, SVM, Naive Bayes	

 $\square$  Type of model – The different models are summed up in the table below:

	Discriminative model	Generative model
Goal	Directly estimate $P(y x)$	Estimate $P(x y)$ to deduce $P(y x)$
What's learned	Decision boundary	Probability distributions of the data
Illustration		
Examples	Regressions, SVMs	GDA, Naive Bayes

### Notations and general concepts

**D** Hypothesis – The hypothesis is noted  $h_{\theta}$  and is the model that we choose. For a given input data  $x^{(i)}$ , the model prediction output is  $h_{\theta}(x^{(i)})$ .

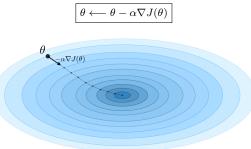
**D** Loss function – A loss function is a function  $L : (z,y) \in \mathbb{R} \times Y \mapsto L(z,y) \in \mathbb{R}$  that takes as inputs the predicted value z corresponding to the real data value y and outputs how different they are. The common loss functions are summed up in the table below:

Least squared	Logistic	Hinge	Cross-entropy
$\frac{1}{2}(y-z)^2$	$\log(1 + \exp(-yz))$	$\max(0,1-yz)$	$-\left[y\log(z) + (1-y)\log(1-z)\right]$
	y = -1	y = -1	y = 0 1 0 $y = 1$
Linear regression	Logistic regression	SVM	Neural Network

**Cost function** – The cost function J is commonly used to assess the performance of a model, and is defined with the loss function L as follows:

$$J(\theta) = \sum_{i=1}^{m} L(h_{\theta}(x^{(i)}), y^{(i)})$$

**Gradient descent** – By noting  $\alpha \in \mathbb{R}$  the learning rate, the update rule for gradient descent is expressed with the learning rate and the cost function J as follows:



Remark: Stochastic gradient descent (SGD) is updating the parameter based on each training example, and batch gradient descent is on a batch of training examples.

**D** Likelihood – The likelihood of a model  $L(\theta)$  given parameters  $\theta$  is used to find the optimal parameters  $\theta$  through maximizing the likelihood. In practice, we use the log-likelihood  $\ell(\theta) = \log(L(\theta))$  which is easier to optimize. We have:

$$\theta^{\mathrm{opt}} = \arg \max_{\theta} L(\theta)$$

**Dewton's algorithm** – The Newton's algorithm is a numerical method that finds  $\theta$  such that  $\ell'(\theta) = 0$ . Its update rule is as follows:

$$heta = heta - rac{\ell'( heta)}{\ell''( heta)}$$

 ${\it Remark: the multidimensional generalization, also known as the Newton-Raphson method, has the following update rule:}$ 

$$\theta = \theta - \left(\nabla_{\theta}^2 \ell(\theta)\right)^{-1} \nabla_{\theta} \ell(\theta)$$

Linear regression

We assume here that  $y|x; \theta \sim \mathcal{N}(\mu, \sigma^2)$ 

 $\Box$  Normal equations – By noting X the matrix design, the value of  $\theta$  that minimizes the cost function is a closed-form solution such that:

 $\theta = (X^T X)^{-1} X^T y$ 

 $\Box$  LMS algorithm – By noting  $\alpha$  the learning rate, the update rule of the Least Mean Squares (LMS) algorithm for a training set of m data points, which is also known as the Widrow-Hoff learning rule, is as follows:

$$\forall j, \quad \theta_j \leftarrow \theta_j + \alpha \sum_{i=1}^m \left[ y^{(i)} - h_\theta(x^{(i)}) \right] x_j^{(i)}$$

Remark: the update rule is a particular case of the gradient ascent.

**D** LWR – Locally Weighted Regression, also known as LWR, is a variant of linear regression that weights each training example in its cost function by  $w^{(i)}(x)$ , which is defined with parameter  $\tau \in \mathbb{R}$  as:

$$w^{(i)}(x) = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right)$$

#### Classification and logistic regression

 $\square$  Sigmoid function – The sigmoid function g, also known as the logistic function, is defined as follows:

$$\forall z \in \mathbb{R}, \quad g(z) = \frac{1}{1 + e^{-z}} \in ]0,1[$$

 $\square$  Logistic regression – We assume here that  $y|x;\theta\sim \text{Bernoulli}(\phi).$  We have the following form:

$$\phi = p(y = 1 | x; \theta) = \frac{1}{1 + \exp(-\theta^T x)} = g(\theta^T x)$$

Remark: there is no closed form solution for the case of logistic regressions.

□ Softmax regression – A softmax regression, also called a multiclass logistic regression, is used to generalize logistic regression when there are more than 2 outcome classes. By convention, we set  $\theta_K = 0$ , which makes the Bernoulli parameter  $\phi_i$  of each class *i* equal to:

$$\phi_i = \frac{\exp(\theta_i^T x)}{\displaystyle\sum_{j=1}^{K} \exp(\theta_j^T x)}$$

**Generalized Linear Models** 

**D** Exponential family – A class of distributions is said to be in the exponential family if it can be written in terms of a natural parameter, also called the canonical parameter or link function,  $\eta$ , a sufficient statistic T(y) and a log-partition function  $a(\eta)$  as follows:

$$p(y;\eta) = b(y) \exp(\eta T(y) - a(\eta))$$

Remark: we will often have T(y) = y. Also,  $\exp(-a(\eta))$  can be seen as a normalization parameter that will make sure that the probabilities sum to one.

Here are the most common exponential distributions summed up in the following table:

Distribution	$\eta$	T(y)	$a(\eta)$	b(y)
Bernoulli	$\log\left(\frac{\phi}{1-\phi}\right)$	y	$\log(1 + \exp(\eta))$	1
Gaussian	$\mu$	y	$\frac{\eta^2}{2}$	$\frac{1}{\sqrt{2\pi}}\exp\left(-\frac{y^2}{2}\right)$
Poisson	$\log(\lambda)$	y	$e^{\eta}$	$\frac{1}{y!}$
Geometric	$\log(1-\phi)$	y	$\log\left(\frac{e^{\eta}}{1-e^{\eta}}\right)$	1

**C** Assumptions of GLMs – Generalized Linear Models (GLM) aim at predicting a random variable y as a function fo  $x \in \mathbb{R}^{n+1}$  and rely on the following 3 assumptions:

(1) 
$$y|x; \theta \sim \text{ExpFamily}(\eta)$$
 (2)  $h_{\theta}(x) = E[y|x; \theta]$  (3)  $\eta = \theta^T x$ 

 ${\it Remark: \ ordinary \ least \ squares \ and \ logistic \ regression \ are \ special \ cases \ of \ generalized \ linear models.}$ 

### Support Vector Machines

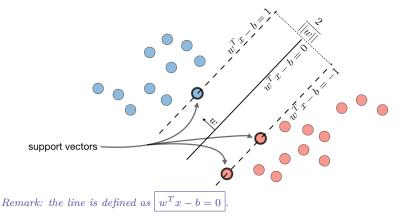
The goal of support vector machines is to find the line that maximizes the minimum distance to the line.

 $\Box$  Optimal margin classifier – The optimal margin classifier h is such that:

$$h(x) = \operatorname{sign}(w^T x - b)$$

where  $(w, b) \in \mathbb{R}^n \times \mathbb{R}$  is the solution of the following optimization problem:

$$\boxed{\min \frac{1}{2} ||w||^2} \quad \text{such that} \quad \boxed{y^{(i)}(w^T x^{(i)} - b) \ge 1}$$

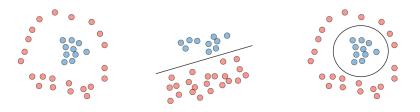


 $\square$  Hinge loss – The hinge loss is used in the setting of SVMs and is defined as follows:  $\boxed{L(z,y)=[1-yz]_+=\max(0,1-yz)}$ 

**G** Kernel – Given a feature mapping  $\phi$ , we define the kernel K to be defined as:

$$K(x,z) = \phi(x)^T \phi(z)$$

In practice, the kernel K defined by  $K(x,z) = \exp\left(-\frac{||x-z||^2}{2\sigma^2}\right)$  is called the Gaussian kernel and is commonly used.



Non-linear separability  $\implies$  Use of a kernel mapping  $\phi \implies$  Decision boundary in the original space

Remark: we say that we use the "kernel trick" to compute the cost function using the kernel because we actually don't need to know the explicit mapping  $\phi$ , which is often very complicated. Instead, only the values K(x,z) are needed.

**Lagrangian** – We define the Lagrangian  $\mathcal{L}(w,b)$  as follows:

$$\mathcal{L}(w,b) = f(w) + \sum_{i=1}^{l} \beta_i h_i(w)$$

Remark: the coefficients  $\beta_i$  are called the Lagrange multipliers.

#### Generative Learning

A generative model first tries to learn how the data is generated by estimating P(x|y), which we can then use to estimate P(y|x) by using Bayes' rule.

#### Gaussian Discriminant Analysis

 $\square$  Setting – The Gaussian Discriminant Analysis assumes that y and x|y=0 and x|y=1 are such that:

$$\boxed{\begin{array}{c} y \sim \text{Bernoulli}(\phi) \end{array}}$$
  
$$x|y = 0 \sim \mathcal{N}(\mu_0, \Sigma) \quad \text{and} \quad \boxed{x|y = 1 \sim \mathcal{N}(\mu_1, \Sigma)}$$

 $\square$  Estimation – The following table sums up the estimates that we find when maximizing the likelihood:

$$\label{eq:phi} \begin{array}{|c|c|c|c|c|} \hline \widehat{\phi} & \widehat{\mu_{j}} & (j=0,1) & \widehat{\Sigma} \\ \hline \\ \hline \frac{1}{m} \sum_{i=1}^{m} \mathbf{1}_{\{y^{(i)}=1\}} & \frac{\sum_{i=1}^{m} \mathbf{1}_{\{y^{(i)}=j\}} x^{(i)}}{\sum_{i=1}^{m} \mathbf{1}_{\{y^{(i)}=j\}}} & \frac{1}{m} \sum_{i=1}^{m} (x^{(i)} - \mu_{y^{(i)}}) (x^{(i)} - \mu_{y^{(i)}})^T \end{array}$$

#### Naive Bayes

 $\square$  Assumption – The Naive Bayes model supposes that the features of each data point are all independent:

$$P(x|y) = P(x_1, x_2, \dots | y) = P(x_1|y)P(x_2|y)\dots = \prod_{i=1}^n P(x_i|y)$$

 $\square$  Solutions – Maximizing the log-likelihood gives the following solutions, with  $k \in \{0,1\}, l \in [\![1,L]\!]$ 

$$P(y=k) = \frac{1}{m} \times \#\{j|y^{(j)} = k\} \quad \text{and} \quad P(x_i = l|y = k) = \frac{\#\{j|y^{(j)} = k \text{ and } x_i^{(j)} = l\}}{\#\{j|y^{(j)} = k\}}$$

Remark: Naive Bayes is widely used for text classification and spam detection.

#### Tree-based and ensemble methods

These methods can be used for both regression and classification problems.

 $\Box$  CART – Classification and Regression Trees (CART), commonly known as decision trees, can be represented as binary trees. They have the advantage to be very interpretable.

 $\square$  Random forest – It is a tree-based technique that uses a high number of decision trees built out of randomly selected sets of features. Contrary to the simple decision tree, it is highly uninterpretable but its generally good performance makes it a popular algorithm.

Remark: random forests are a type of ensemble methods.

 $\square$  **Boosting** – The idea of boosting methods is to combine several weak learners to form a stronger one. The main ones are summed up in the table below:

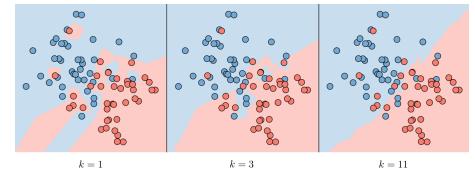
https://	/stanford.edu/	~shervine
----------	----------------	-----------

Adaptive boosting	Gradient boosting
<ul><li>High weights are put on errors to</li></ul>	- Weak learners trained
improve at the next boosting step <li>Known as Adaboost</li>	on remaining errors

#### Other non-parametric approaches

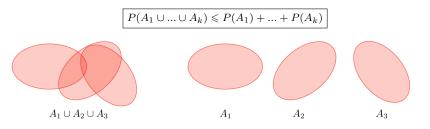
 $\Box$  k-nearest neighbors – The k-nearest neighbors algorithm, commonly known as k-NN, is a non-parametric approach where the response of a data point is determined by the nature of its k neighbors from the training set. It can be used in both classification and regression settings.

Remark: The higher the parameter k, the higher the bias, and the lower the parameter k, the higher the variance.



### Learning Theory

**D** Union bound – Let  $A_1, ..., A_k$  be k events. We have:



**D** Hoeffding inequality – Let  $Z_1, ..., Z_m$  be *m* iid variables drawn from a Bernoulli distribution of parameter  $\phi$ . Let  $\hat{\phi}$  be their sample mean and  $\gamma > 0$  fixed. We have:

$P( \phi - \widehat{\phi}  > \gamma) \leqslant 2 \exp(-2\gamma^2 m)$
--

Remark: this inequality is also known as the Chernoff bound.

 $\Box$  Training error – For a given classifier h, we define the training error  $\hat{\epsilon}(h)$ , also known as the empirical risk or empirical error, to be as follows:

$$\widehat{\epsilon}(h) = \frac{1}{m} \sum_{i=1}^{m} \mathbf{1}_{\{h(x^{(i)}) \neq y^{(i)}\}}$$

 $\Box$  **Probably Approximately Correct (PAC)** – PAC is a framework under which numerous results on learning theory were proved, and has the following set of assumptions:

- the training and testing sets follow the same distribution
- the training examples are drawn independently

□ Shattering – Given a set  $S = \{x^{(1)}, ..., x^{(d)}\}$ , and a set of classifiers  $\mathcal{H}$ , we say that  $\mathcal{H}$  shatters S if for any set of labels  $\{y^{(1)}, ..., y^{(d)}\}$ , we have:

$$\exists h \in \mathcal{H}, \quad \forall i \in \llbracket 1, d \rrbracket, \quad h(x^{(i)}) = y^{(i)}$$

**Dpper bound theorem** – Let  $\mathcal{H}$  be a finite hypothesis class such that  $|\mathcal{H}| = k$  and let  $\delta$  and the sample size *m* be fixed. Then, with probability of at least  $1 - \delta$ , we have:

$\epsilon(\widehat{h}) \leqslant \left(\min_{h \in \mathcal{H}} \epsilon(h)\right) + 2\sqrt{\frac{1}{2m} \log\left(\frac{2k}{\delta}\right)}$
---

 $\Box$  VC dimension – The Vapnik-Chervonenkis (VC) dimension of a given infinite hypothesis class  $\mathcal{H}$ , noted VC( $\mathcal{H}$ ) is the size of the largest set that is shattered by  $\mathcal{H}$ .

Remark: the VC dimension of  $\mathcal{H} = \{set \ of \ linear \ classifiers \ in \ 2 \ dimensions\}$  is 3.



**D** Theorem (Vapnik) – Let  $\mathcal{H}$  be given, with  $VC(\mathcal{H}) = d$  and m the number of training examples. With probability at least  $1 - \delta$ , we have:

$$\epsilon(\widehat{h}) \leqslant \left(\min_{h \in \mathcal{H}} \epsilon(h)\right) + O\left(\sqrt{\frac{d}{m} \log\left(\frac{m}{d}\right) + \frac{1}{m} \log\left(\frac{1}{\delta}\right)}\right)$$

## VIP Cheatsheet: Machine Learning Tips

Afshine AMIDI and Shervine AMIDI

September 9, 2018

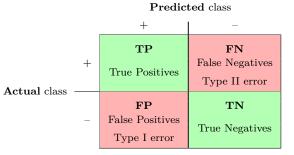
### Metrics

Given a set of data points  $\{x^{(1)}, ..., x^{(m)}\}$ , where each  $x^{(i)}$  has n features, associated to a set of outcomes  $\{y^{(1)}, ..., y^{(m)}\}$ , we want to assess a given classifier that learns how to predict y from x.

### Classification

In a context of a binary classification, here are the main metrics that are important to track to assess the performance of the model.

 $\Box$  Confusion matrix – The confusion matrix is used to have a more complete picture when assessing the performance of a model. It is defined as follows:

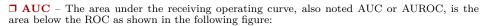


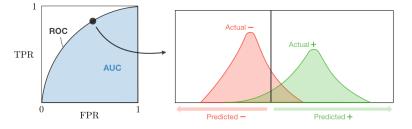
**D** Main metrics – The following metrics are commonly used to assess the performance of classification models:

Metric	Formula	Interpretation
Accuracy	$\frac{\mathrm{TP} + \mathrm{TN}}{\mathrm{TP} + \mathrm{TN} + \mathrm{FP} + \mathrm{FN}}$	Overall performance of model
Precision	$\frac{\mathrm{TP}}{\mathrm{TP}+\mathrm{FP}}$	How accurate the positive predictions are
Recall Sensitivity	$\frac{\mathrm{TP}}{\mathrm{TP} + \mathrm{FN}}$	Coverage of actual positive sample
Specificity	$\frac{\mathrm{TN}}{\mathrm{TN}+\mathrm{FP}}$	Coverage of actual negative sample
F1 score	$\frac{2\mathrm{TP}}{2\mathrm{TP}+\mathrm{FP}+\mathrm{FN}}$	Hybrid metric useful for unbalanced classes

 $\Box$  **ROC** – The receiver operating curve, also noted ROC, is the plot of TPR versus FPR by varying the threshold. These metrics are are summed up in the table below:

Metric	Formula	Equivalent
True Positive Rate TPR	$\frac{\mathrm{TP}}{\mathrm{TP} + \mathrm{FN}}$	Recall, sensitivity
False Positive Rate FPR	$\frac{\rm FP}{\rm TN+FP}$	1-specificity





### Regression

**D** Basic metrics – Given a regression model f, the following metrics are commonly used to assess the performance of the model:

Total su	m of squares	Explained sum of squares	Residual sum of squares
$SS_{tot} =$	$\sum_{i=1}^m (y_i - \overline{y})^2$	$SS_{reg} = \sum_{i=1}^{m} (f(x_i) - \overline{y})^2$	$SS_{res} = \sum_{i=1}^{m} (y_i - f(x_i))^2$

**Coefficient of determination** – The coefficient of determination, often noted  $R^2$  or  $r^2$ , provides a measure of how well the observed outcomes are replicated by the model and is defined as follows:

$$R^2 = 1 - \frac{\mathrm{SS}_{\mathrm{res}}}{\mathrm{SS}_{\mathrm{tot}}}$$

**D** Main metrics – The following metrics are commonly used to assess the performance of regression models, by taking into account the number of variables n that they take into consideration:

Mallow's Cp	AIC	BIC	Adjusted $R^2$
$\frac{\mathrm{SS}_{\mathrm{res}} + 2(n+1)\widehat{\sigma}^2}{m}$	$2\Big[(n+2) - \log(L)\Big]$	$\log(m)(n+2) - 2\log(L)$	$1 - \frac{(1 - R^2)(m - 1)}{m - n - 1}$

where L is the likelihood and  $\widehat{\sigma}^2$  is an estimate of the variance associated with each response.

### Model selection

 $\square$  Vocabulary – When selecting a model, we distinguish 3 different parts of the data that we have as follows:

Training set	Validation set	Testing set
<ul><li>Model is trained</li><li>Usually 80% of the dataset</li></ul>	<ul><li>Model is assessed</li><li>Usually 20% of the dataset</li><li>Also called hold-out</li><li>or development set</li></ul>	<ul><li>Model gives predictions</li><li>Unseen data</li></ul>

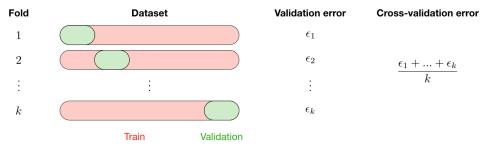
Once the model has been chosen, it is trained on the entire dataset and tested on the unseen test set. These are represented in the figure below:



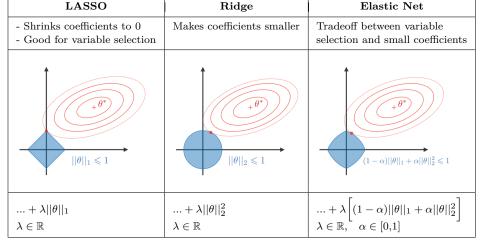
 $\Box$  Cross-validation – Cross-validation, also noted CV, is a method that is used to select a model that does not rely too much on the initial training set. The different types are summed up in the table below:

k-fold	Leave- <i>p</i> -out
- Training on $k-1$ folds and	- Training on $n - p$ observations and
assessment on the remaining one	assessment on the $p$ remaining ones
- Generally $k = 5$ or $10$	- Case $p = 1$ is called leave-one-out

The most commonly used method is called k-fold cross-validation and splits the training data into k folds to validate the model on one fold while training the model on the k-1 other folds, all of this k times. The error is then averaged over the k folds and is named cross-validation error.



**D** Regularization – The regularization procedure aims at avoiding the model to overfit the data and thus deals with high variance issues. The following table sums up the different types of commonly used regularization techniques:



□ Model selection – Train model on training set, then evaluate on the development set, then pick best performance model on the development set, and retrain all of that model on the whole training set.

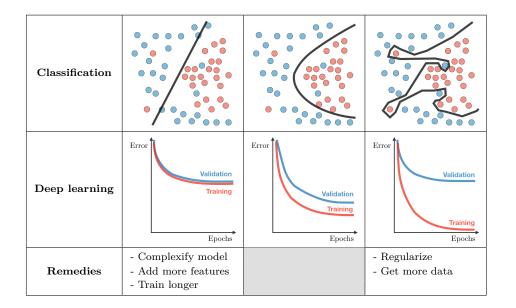
#### Diagnostics

 $\square$  Bias – The bias of a model is the difference between the expected prediction and the correct model that we try to predict for given data points.

 $\square$  **Variance** – The variance of a model is the variability of the model prediction for given data points.

 $\square$  Bias/variance tradeoff – The simpler the model, the higher the bias, and the more complex the model, the higher the variance.

	Underfitting	Just right	Overfitting
Symptoms	<ul> <li>High training error</li> <li>Training error close to test error</li> <li>High bias</li> </ul>	- Training error slightly lower than test error	<ul> <li>Low training error</li> <li>Training error much</li> <li>lower than test error</li> <li>High variance</li> </ul>
Regression			Myst



 $\square$  Error analysis – Error analysis is analyzing the root cause of the difference in performance between the current and the perfect models.

 $\square$  Ablative analysis – Ablative analysis is analyzing the root cause of the difference in performance between the current and the baseline models.

## Cheat Sheet – Regression Analysis

### What is Regression Analysis?

Fitting a function f(.) to datapoints  $y_i = f(x_i)$  under some error function. Based on the estimated function and error, we have the following types of regression

 $min_{\beta} \sum_{i} \|y_i - f_{\beta}^{linear}(x_i)\|^2$ 

 $f_{\beta}^{inear}(x_i) = \beta_0 + \beta_1 x_i$  $min_{\beta} \sum_{i=0}^{m} \|y_i - f_{\beta}^{poly}(x_i)\|^2$ 

 $f_{\beta}^{poly}(x_i) = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \ldots + \beta_k x_i^k$ 

 $min_{\beta} \sum \|y_i - \mathcal{N}\left(f_{\beta}(x_i), \sigma^2\right)\|^2$ 

 $f_{\beta}(x_i) \stackrel{i}{=} f_{\beta}^{poly}(x_i) \text{ or } f_{\beta}^{linear}(x_i)$ 

 $min_{\beta} \sum_{i=0}^{m} \|y_i - f_{\beta}(x_i)\|^2 + \sum_{i=0}^{k} \beta_j^2$ 

 $f_{\beta}(x_i) = f_{\beta}^{poly}(x_i) \text{ or } f_{\beta}^{linear}(x_i)$ 

 $min_{\beta} \sum_{i=0}^{m} \|y_i - f_{\beta}(x_i)\|^2 + \sum_{i=0}^{k} |\beta_j|$ 

 $f_{\beta}(x_i) = f_{\beta}^{poly}(x_i) \text{ or } f_{\beta}^{linear}(x_i)$ 

 $f_{\beta}(x_i) = f_{\beta}^{poly}(x_i) \text{ or } f_{\beta}^{linear}(x_i)$ 

 $\sigma(t) = \frac{1}{1 + e^{-t}}$ 

### 1. Linear Regression:

Fits a line minimizing the sum of mean-squared error for each datapoint.

2. Polynomial Regression:

Fits a polynomial of order k (k+1 unknowns) minimizing the sum of mean-squared error for each datapoint.

## 3. Bayesian Regression:

For each datapoint, fits a gaussian distribution by minimizing the mean-squared error. As the number of data points x<sub>i</sub> increases, it converges to point  $\mathcal{N}(\mu, \sigma^2) \to \text{Gaussian}$  with mean  $\mu$  and variance  $\sigma^2$ estimates i.e.  $n \to \infty, \sigma^2 \to 0$ 

### 4. Ridge Regression:

Can fit either a line, or polynomial minimizing the sum of mean-squared error for each datapoint and the weighted L2 norm of the function parameters beta.

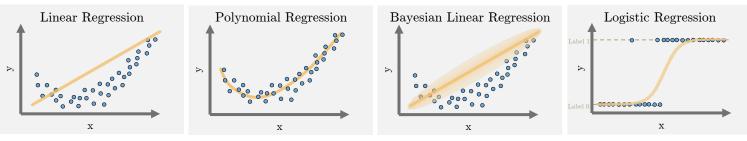
## 5. LASSO Regression:

Can fit either a line, or polynomial minimizing the the sum of mean-squared error for each datapoint and the weighted L1 norm of the function parameters beta.

6. Logistic Regression:

 $\min_{\beta} \sum_{i} -y_{i} \log\left(\sigma\left(f_{\beta}(x_{i})\right)\right) - (1-y_{i}) \log\left(1-\sigma\left(f_{\beta}(x_{i})\right)\right)$ Can fit either a line, or polynomial with sigmoid activation minimizing the binary cross-entropy loss for each datapoint. The labels y are binary class labels.

## Visual Representation:



### Summary:

	What does it fit?	Estimated function	Error Function
Linear	A line in n dimensions	$f_{\beta}^{linear}(x_i) = \beta_0 + \beta_1 x_i$	$\sum_{i=0}^m \ y_i - f_\beta(x_i)\ ^2 \cdot$
Polynomial	A polynomial of order <b>k</b>	$f_{\beta}^{poly}(x_i) = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots$	$\sum_{i=0}^m \ y_i - f_eta(x_i)\ ^2 \cdot$
Bayesian Linear	Gaussian distribution for each point	$\mathcal{N}\left(f_{\beta}(x_i),\sigma^2\right)$	$\sum_{i} \ y_i - \mathcal{N}\left(f_eta(x_i), \sigma^2 ight)\ ^2$
Ridge	Linear/polynomial	$f_{\beta}^{poly}(x_i) \ or \ f_{\beta}^{linear}(x_i)$	$\sum_{\substack{i=0\\m}}^{m} \ y_i - f_\beta(x_i)\ ^2 + \sum_{\substack{j=0\\n}}^{n} \beta_j^2$
LASSO	Linear/polynomial	$f_{\beta}^{poly}(x_i) \text{ or } f_{\beta}^{linear}(x_i)$	$\sum_{i=0}^{m} \ y_i - f_{\beta}(x_i)\ ^2 + \sum_{i=0}^{n}  \beta_j $
Logistic	Linear/polynomial with sigmoid	$\sigma(f_eta(x_i)) \qquad min_eta\sum\limits_i -y_i log$	$q\left(\sigma\left(f_{\beta}(x_{i})\right)\right) - (1 - y_{i})log\left(1 - \sigma\left(f_{\beta}(x_{i})\right)\right)$

Source: https://www.cheatsheets.aqeel-anwar.com Tutorial: Click here

## Cheat Sheet – Regularization in ML

## What is Regularization in ML?

- Regularization is an approach to address over-fitting in ML.
- Overfitted model fails to generalize estimations on test data
- When the underlying model to be learned is low bias/high variance, or when we have small amount of data, the estimated model is prone to over-fitting.
- Regularization reduces the variance of the model

## Types of Regularization:

## 1. Modify the loss function:

• L2 Regularization: Prevents the weights from getting too large (defined by L2 norm). Larger the weights, more complex the model is, more chances of overfitting.

$$loss = error(y, \hat{y}) + \lambda \sum_{j} \beta_{j}^{2} \quad \lambda \ge 0, \ \lambda \propto model \ bias, \ \lambda \propto \frac{1}{model \ variance}$$

• L1 Regularization: Prevents the weights from getting too large (defined by L1 norm). Larger the weights, more complex the model is, more chances of overfitting. L1 regularization introduces sparsity in the weights. It forces more weights to be zero, than reducing the the average magnitude of all weights

$$loss = error(y, \hat{y}) + \lambda \sum_{j} |\beta_{j}| \quad \lambda \ge 0, \ \lambda \propto model \ bias, \ \lambda \propto \frac{1}{model \ variance}$$

• Entropy: Used for the models that output probability. Forces the probability distribution towards uniform distribution.

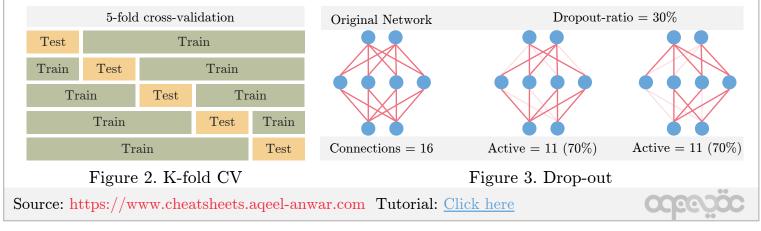
$$loss = error(p, \hat{p}) - \lambda \sum_{i} \hat{p}_{i} log(\hat{p}_{i}) \qquad \lambda \geq 0, \ \lambda \propto model \ bias, \ \lambda \propto \frac{1}{model \ variance}$$

## 2. Modify data sampling:

- Data augmentation: Create more data from available data by randomly cropping, dilating, rotating, adding small amount of noise etc.
- K-fold Cross-validation: Divide the data into k groups. Train on (k-1) groups and test on 1 group. Try all k possible combinations.

## 3. Change training approach:

- Injecting noise: Add random noise to the weights when they are being learned. It pushes the model to be relatively insensitive to small variations in the weights, hence regularization
- **Dropout:** Generally used for neural networks. Connections between consecutive layers are randomly dropped based on a dropout-ratio and the remaining network is trained in the current iteration. In the next iteration, another set of random connections are dropped.



Under-fittingJust RightOver-fittingPreferred if size<br/>of dataset is smallPreferred if size<br/>of dataset is large

Figure 1. Overfitting



## Cheat Sheet – Bias-Variance Tradeoff

## What is Bias?

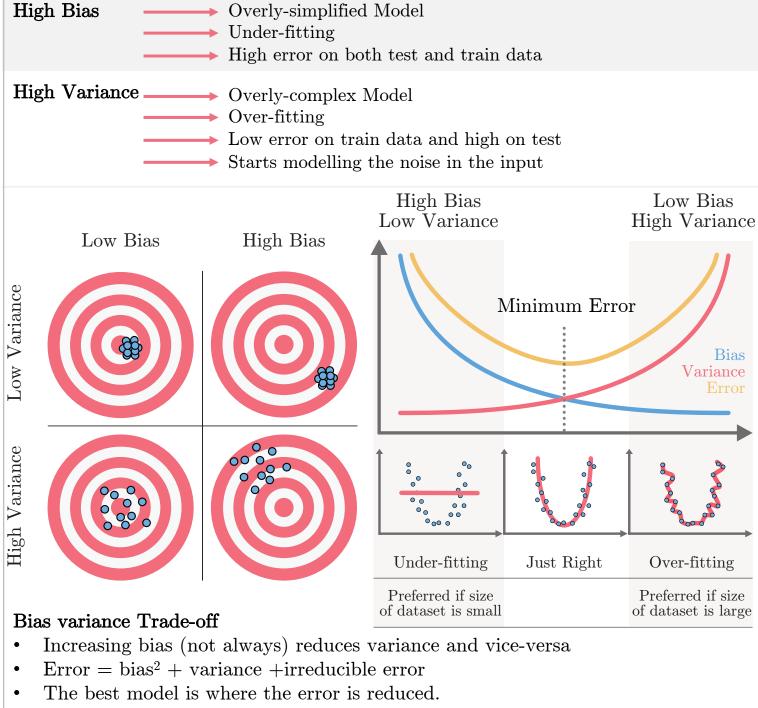
- Error between average model prediction and ground truth
- The bias of the estimated function tells us the capacity of the underlying model to predict the values

## What is Variance?

$$eriance = \mathbb{E}\left[\left(f'(x) - \mathbb{E}[f'(x)]\right)^2\right]$$

 $bias = \mathbb{E}[f'(x)] - f(x)$ 

- Average variability in the model prediction for the given dataset  $variance = \mathbb{E}$
- The variance of the estimated function tells you how much the function can adjust to the change in the dataset



• Compromise between bias and variance

Source: https://www.cheatsheets.aqeel-anwar.com Tutorial: Click here



## Cheat Sheet – Bayes Theorem and Classifier

## What is Bayes' Theorem?

• Describes the probability of an event, based on prior knowledge of conditions that might be related to the event.

$$P(A|B) = \frac{P(B|A)(likelihood) \times P(A)(prior)}{P(B)(evidence)}$$

• How the probability of an event changes when we have knowledge of another event

$$(A) \longrightarrow P(A|B)$$

$$(A) \longrightarrow Usually, a better estimate than P(A)$$

## Example

• Probability of fire P(F) = 1%

Ρ

- Probability of smoke P(S) = 10%
- Prob of smoke given there is a fire P(S|F) = 90%
- What is the probability that there is a fire given we see a smoke P(F|S)?

$$P(F|S) = \frac{P(S|F) \times P(F)}{P(S)} = \frac{0.9 \times 0.01}{0.1} = 9\%$$

## Maximum Aposteriori Probability (MAP) Estimation

The MAP estimate of the random variable y, given that we have observed iid  $(x_1, x_2, x_3, ...)$ , is given by. We try to accommodate our prior knowledge when estimating.

$$\hat{y}_{MAP} = argmax_y \quad P(y) \prod_i P(x_i|y)$$

y that maximizes the product of prior and likelihood

## Maximum Likelihood Estimation (MLE)

The MAP estimate of the random variable y, given that we have observed iid  $(x_1, x_2, x_3, ...)$ , is given by. We assume we don't have any prior knowledge of the quantity being estimated.

$$\hat{y}_{MLE} = argmax_y \prod_i P(x_i|y)$$
 y that maximizes only the likelihood

MLE is a special case of MAP where our prior is uniform (all values are equally likely)

## Naïve Bayes' Classifier (Instantiation of MAP as classifier)

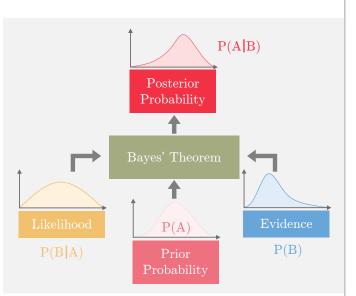
Suppose we have two classes,  $y=y_1$  and  $y=y_2$ . Say we have more than one evidence/features ( $x_1$ ,  $x_2$ ,  $x_3$ , ... ), using Bayes' theorem

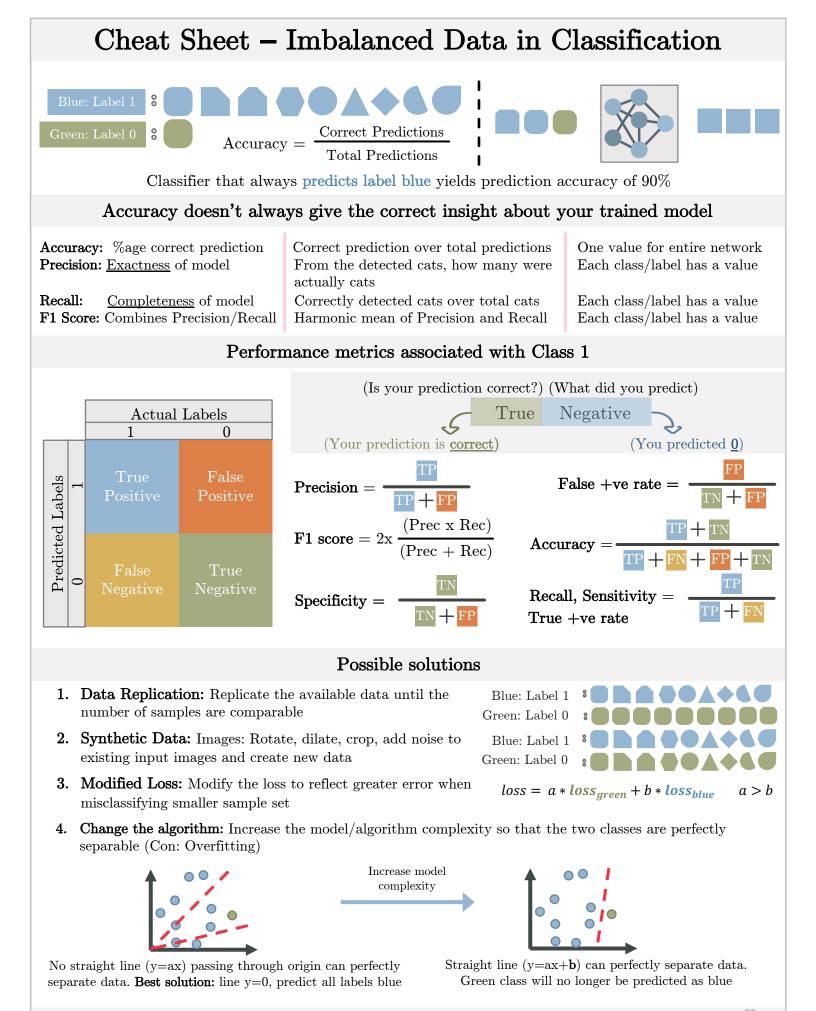
$$P(y|x_1, x_2, x_3, \ldots) = \frac{P(x_1, x_2, x_3, \ldots | y) \times P(y)}{P(x_1, x_2, x_3, \ldots)}$$

Naïve Bayes' theorem assumes the features  $(x_1, x_2, ...)$  are i.i.d. i.e  $P(x_1, x_2, x_3, ... | y) = \prod_{i=1}^{n} P(x_i | y)$ 

$$P(y|x_1, x_2, x_3, \ldots) = \prod_i P(x_i|y) \frac{P(y)}{P(x_1, x_2, x_3, \ldots)}$$
$$\hat{y} = y_1 \ if \ \frac{P(y_1|x_1, x_2, x_3, \ldots)}{P(y_2|x_1, x_2, x_3, \ldots)} > 1 \ else \ \hat{y} = y_2$$

Source: https://www.cheatsheets.aqeel-anwar.com





Source: https://www.cheatsheets.aqeel-anwar.com Tutorial: Click here

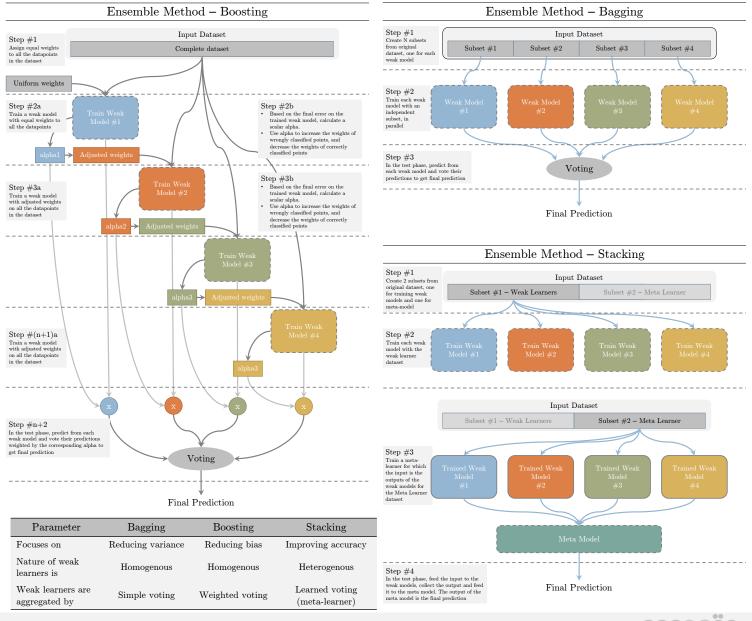
## Cheat Sheet – Ensemble Learning in ML

### What is Ensemble Learning? Wisdom of the crowd

Combine multiple weak models/learners into one predictive model to reduce bias, variance and/or improve accuracy.

## Types of Ensemble Learning: N number of weak learners

- **1.Bagging:** Trains N different weak models (usually of same types homogenous) with N non-overlapping subset of the input dataset in parallel. In the test phase, each model is evaluated. The label with the greatest number of predictions is selected as the prediction. Bagging methods reduces variance of the prediction
- **2.Boosting:** Trains N different weak models (usually of same types homogenous) with the complete dataset in a sequential order. The datapoints wrongly classified with previous weak model is provided more weights to that they can be classified by the next weak leaner properly. In the test phase, each model is evaluated and based on the test error of each weak model, the prediction is weighted for voting. Boosting methods decreases the bias of the prediction.
- **3.Stacking:** Trains N different weak models (usually of different types heterogenous) with one of the two subsets of the dataset in parallel. Once the weak learners are trained, they are used to trained a meta learner to combine their predictions and carry out final prediction using the other subset. In test phase, each model predicts its label, these set of labels are fed to the meta learner which generates the final prediction.



### The block diagrams, and comparison table for each of these three methods can be seen below.

Source: https://www.cheatsheets.aqeel-anwar.com Tutorial: <u>Click here</u>



# 3.2. Unsupervised learning

## VIP Cheatsheet: Unsupervised Learning

Afshine AMIDI and Shervine AMIDI

September 9, 2018

### Introduction to Unsupervised Learning

**D** Motivation – The goal of unsupervised learning is to find hidden patterns in unlabeled data  $\{x^{(1)}, \dots, x^{(m)}\}.$ 

**]** Jensen's inequality – Let f be a convex function and X a random variable. We have the following inequality:

E[f(X)]	$\geq f(E[X])$
---------	----------------

#### Expectation-Maximization

**D** Latent variables – Latent variables are hidden/unobserved variables that make estimation problems difficult, and are often denoted z. Here are the most common settings where there are latent variables:

Setting	Latent variable $z$	x z	Comments
Mixture of $k$ Gaussians	$\mathrm{Multinomial}(\phi)$	$\mathcal{N}(\mu_j, \Sigma_j)$	$\mu_j \in \mathbb{R}^n, \phi \in \mathbb{R}^k$
Factor analysis	$\mathcal{N}(0,I)$	$\mathcal{N}(\mu + \Lambda z, \psi)$	$\mu_j \in \mathbb{R}^n$

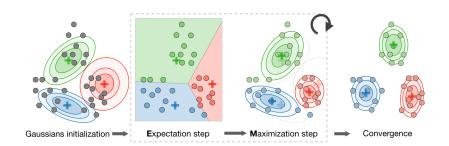
 $\Box$  Algorithm – The Expectation-Maximization (EM) algorithm gives an efficient method at estimating the parameter  $\theta$  through maximum likelihood estimation by repeatedly constructing a lower-bound on the likelihood (E-step) and optimizing that lower bound (M-step) as follows:

• <u>E-step</u>: Evaluate the posterior probability  $Q_i(z^{(i)})$  that each data point  $x^{(i)}$  came from a particular cluster  $z^{(i)}$  as follows:

$$Q_i(z^{(i)}) = P(z^{(i)}|x^{(i)};\theta)$$

• M-step: Use the posterior probabilities  $Q_i(z^{(i)})$  as cluster specific weights on data points  $\overline{x^{(i)}}$  to separately re-estimate each cluster model as follows:

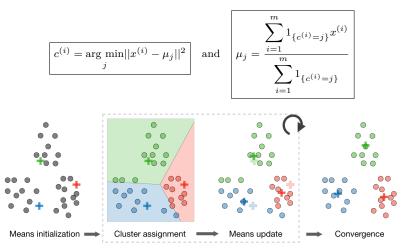
$$\theta_i = \operatorname*{argmax}_{\theta} \sum_{i} \int_{z^{(i)}} Q_i(z^{(i)}) \log\left(\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}\right) dz^{(i)}$$



#### k-means clustering

We note  $c^{(i)}$  the cluster of data point *i* and  $\mu_j$  the center of cluster *j*.

 $\Box$  Algorithm – After randomly initializing the cluster centroids  $\mu_1, \mu_2, ..., \mu_k \in \mathbb{R}^n$ , the k-means algorithm repeats the following step until convergence:



 $\square$  **Distortion function** – In order to see if the algorithm converges, we look at the distortion function defined as follows:

$$J(c,\mu) = \sum_{i=1}^{m} ||x^{(i)} - \mu_{c^{(i)}}||^2$$

#### Hierarchical clustering

 $\square$  Algorithm – It is a clustering algorithm with an agglomerative hierarchical approach that build nested clusters in a successive manner.

 $\Box$  **Types** – There are different sorts of hierarchical clustering algorithms that aims at optimizing different objective functions, which is summed up in the table below:

Ward linkage	Average linkage	Complete linkage	
Minimize within cluster	Minimize average distance	Minimize maximum distance	
distance	between cluster pairs	of between cluster pairs	

#### **Clustering assessment metrics**

In an unsupervised learning setting, it is often hard to assess the performance of a model since we don't have the ground truth labels as was the case in the supervised learning setting.

**Solution** Silhouette coefficient – By noting a and b the mean distance between a sample and all other points in the same class, and between a sample and all other points in the next nearest cluster, the silhouette coefficient s for a single sample is defined as follows:

$$s = \frac{b-a}{\max(a,b)}$$

 $\Box$  Calinski-Harabaz index – By noting k the number of clusters,  $B_k$  and  $W_k$  the between and within-clustering dispersion matrices respectively defined as

$$B_k = \sum_{j=1}^k n_{c(i)} (\mu_{c(i)} - \mu) (\mu_{c(i)} - \mu)^T, \qquad W_k = \sum_{i=1}^m (x^{(i)} - \mu_{c(i)}) (x^{(i)} - \mu_{c(i)})^T$$

the Calinski-Harabaz index s(k) indicates how well a clustering model defines its clusters, such that the higher the score, the more dense and well separated the clusters are. It is defined as follows:

$$s(k) = \frac{\operatorname{Tr}(B_k)}{\operatorname{Tr}(W_k)} \times \frac{N-k}{k-1}$$

#### Principal component analysis

It is a dimension reduction technique that finds the variance maximizing directions onto which to project the data.

**D** Eigenvalue, eigenvector – Given a matrix  $A \in \mathbb{R}^{n \times n}$ ,  $\lambda$  is said to be an eigenvalue of A if there exists a vector  $z \in \mathbb{R}^n \setminus \{0\}$ , called eigenvector, such that we have:

 $Az = \lambda z$ 

**Spectral theorem** – Let  $A \in \mathbb{R}^{n \times n}$ . If A is symmetric, then A is diagonalizable by a real orthogonal matrix  $U \in \mathbb{R}^{n \times n}$ . By noting  $\Lambda = \text{diag}(\lambda_1, ..., \lambda_n)$ , we have:

$$\exists \Lambda \text{ diagonal}, \quad A = U \Lambda U^T$$

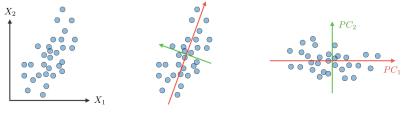
Remark: the eigenvector associated with the largest eigenvalue is called principal eigenvector of matrix A.

 $\Box$  Algorithm – The Principal Component Analysis (PCA) procedure is a dimension reduction technique that projects the data on k dimensions by maximizing the variance of the data as follows:

• Step 1: Normalize the data to have a mean of 0 and standard deviation of 1.

$$x_{j}^{(i)} \leftarrow \frac{x_{j}^{(i)} - \mu_{j}}{\sigma_{j}} \quad \text{where} \quad \boxed{\mu_{j} = \frac{1}{m} \sum_{i=1}^{m} x_{j}^{(i)}} \quad \text{and} \quad \boxed{\sigma_{j}^{2} = \frac{1}{m} \sum_{i=1}^{m} (x_{j}^{(i)} - \mu_{j})^{2}}$$

- <u>Step 2</u>: Compute  $\Sigma = \frac{1}{m} \sum_{i=1}^{m} x^{(i)} x^{(i)^T} \in \mathbb{R}^{n \times n}$ , which is symmetric with real eigenvalues.
- Step 3: Compute  $u_1, ..., u_k \in \mathbb{R}^n$  the k orthogonal principal eigenvectors of  $\Sigma$ , i.e. the orthogonal eigenvectors of the k largest eigenvalues.
- Step 4: Project the data on  $\operatorname{span}_{\mathbb{R}}(u_1,...,u_k)$ . This procedure maximizes the variance among all k-dimensional spaces.



Data in feature space 🛛 🛶 Find principal components 🛶 Data in principal components space

#### Independent component analysis

It is a technique meant to find the underlying generating sources.

**D** Assumptions – We assume that our data x has been generated by the *n*-dimensional source vector  $s = (s_1, ..., s_n)$ , where  $s_i$  are independent random variables, via a mixing and non-singular matrix A as follows:

$$x = As$$

The goal is to find the unmixing matrix  $W = A^{-1}$  by an update rule.

 $\square$  Bell and Sejnowski ICA algorithm – This algorithm finds the unmixing matrix W by following the steps below:

• Write the probability of  $x = As = W^{-1}s$  as:

$$p(x) = \prod_{i=1}^{n} p_s(w_i^T x) \cdot |W|$$

- Write the log likelihood given our training data  $\{x^{(i)}, i \in [\![1,m]\!]\}$  and by noting g the sigmoid function as:

$$l(W) = \sum_{i=1}^m \left( \sum_{j=1}^n \log \left( g'(w_j^T x^{(i)}) \right) + \log |W| \right)$$

Therefore, the stochastic gradient ascent learning rule is such that for each training example  $x^{(i)}$ , we update W as follows:

$$W \longleftarrow W + \alpha \left( \begin{pmatrix} 1 - 2g(w_1^T x^{(i)}) \\ 1 - 2g(w_2^T x^{(i)}) \\ \vdots \\ 1 - 2g(w_n^T x^{(i)}) \end{pmatrix} x^{(i)^T} + (W^T)^{-1} \right)$$

# Cheat Sheet – PCA Dimensionality Reduction

## What is PCA?

- Based on the dataset find a new set of orthogonal feature vectors in such a way that the data spread is maximum in the direction of the feature vector (or dimension)
- Rates the feature vector in the decreasing order of data spread (or variance)
- The datapoints have maximum variance in the first feature vector, and minimum variance in the last feature vector
- The variance of the datapoints in the direction of feature vector can be termed as a measure of information in that direction.

## Steps

- 1. Standardize the datapoints
- 2. Find the covariance matrix from the given datapoints
- 3. Carry out eigen-value decomposition of the covariance matrix
- 4. Sort the eigenvalues and eigenvectors

## Dimensionality Reduction with PCA

• Keep the first m out of n feature vectors rated by PCA. These m vectors will be the best m vectors preserving the maximum information that could have been preserved with m vectors on the given dataset

## Steps:

- 1. Carry out steps 1-4 from above
- 2. Keep first m feature vectors from the sorted eigenvector matrix
- $V_{reduced} = V[:, 0:m]$

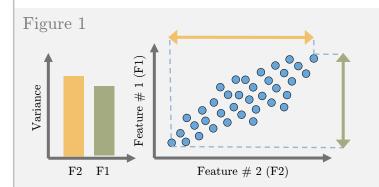
 $X_{new} = \frac{X - mean(X)}{std(X)}$ 

 $\Sigma_{sort} = sort(\Sigma) \ V_{sort} = sort(V, \Sigma_{sort})$ 

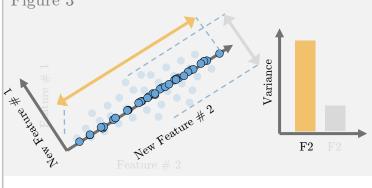
 $C[i,j] = cov(x_i, x_j)$ 

 $C = V \Sigma V^{-1}$ 

- 3. Transform the data for the new basis (feature vectors)  $X_{reduced} = X_{new} \times V_{reduced}$
- 4. The importance of the feature vector is proportional to the magnitude of the eigen value







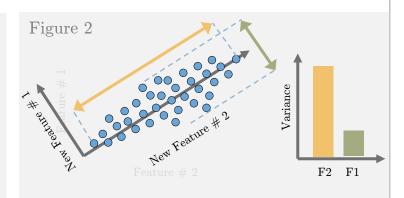


Figure 1: Datapoints with feature vectors as x and y-axis

Figure 2: The cartesian coordinate system is rotated to maximize the standard deviation along any one axis (new feature # 2) Figure 3: Remove the feature vector with minimum standard deviation of datapoints (new feature # 1) and project the data on new feature # 2

Source: https://www.cheatsheets.aqeel-anwar.com





4. Deep learning

# VIP Cheatsheet: Deep Learning

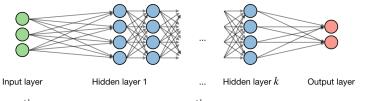
Afshine AMIDI and Shervine AMIDI

## September 15, 2018

#### **Neural Networks**

Neural networks are a class of models that are built with layers. Commonly used types of neural networks include convolutional and recurrent neural networks.

 $\square$  Architecture – The vocabulary around neural networks architectures is described in the figure below:



By noting i the  $i^{th}$  layer of the network and j the  $j^{th}$  hidden unit of the layer, we have:

$$z_{j}^{[i]} = w_{j}^{[i]T}x + b_{j}^{[i]}$$

where we note w, b, z the weight, bias and output respectively.

<b>Activation function</b> – Activation functions are used at the end of a hidden unit to introduce	
non-linear complexities to the model. Here are the most common ones:	

Sigmoid	Tanh	ReLU	Leaky ReLU
$g(z) = \frac{1}{1 + e^{-z}}$	$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$	$g(z) = \max(0, z)$	$g(z) = \max(\epsilon z, z)$ with $\epsilon \ll 1$

**Cross-entropy loss** – In the context of neural networks, the cross-entropy loss L(z,y) is commonly used and is defined as follows:

$$L(z,y) = -\left[y \log(z) + (1-y) \log(1-z)\right]$$

**D** Learning rate – The learning rate, often noted  $\eta$ , indicates at which pace the weights get updated. This can be fixed or adaptively changed. The current most popular method is called Adam, which is a method that adapts the learning rate.

**Backpropagation** – Backpropagation is a method to update the weights in the neural network by taking into account the actual output and the desired output. The derivative with respect to weight w is computed using chain rule and is of the following form:

$\partial L(z,y)$	$\partial L(z,y)$	) da	$\partial z$
$-\frac{\partial w}{\partial w} =$	$\partial a$	$\times \overline{\partial z} \times$	$\overline{\partial w}$

As a result, the weight is updated as follows:

$w \leftarrow$	$w - \eta \frac{\partial L(z,y)}{\partial w}$
	Ow

□ Updating weights – In a neural network, weights are updated as follows:

- Step 1: Take a batch of training data.
- Step 2: Perform forward propagation to obtain the corresponding loss.
- Step 3: Backpropagate the loss to get the gradients.
- Step 4: Use the gradients to update the weights of the network.

 $\Box$  Dropout – Dropout is a technique meant at preventing overfitting the training data by dropping out units in a neural network. In practice, neurons are either dropped with probability p or kept with probability 1 - p.

#### **Convolutional Neural Networks**

**Convolutional layer requirement** – By noting W the input volume size, F the size of the convolutional layer neurons, P the amount of zero padding, then the number of neurons N that fit in a given volume is such that:

$$N = \frac{W - F + 2P}{S} + 1$$

 $\Box$  Batch normalization – It is a step of hyperparameter  $\gamma, \beta$  that normalizes the batch  $\{x_i\}$ . By noting  $\mu_B, \sigma_B^2$  the mean and variance of that we want to correct to the batch, it is done as follows:

$$x_i \longleftarrow \gamma \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

It is usually done after a fully connected/convolutional layer and before a non-linearity layer and aims at allowing higher learning rates and reducing the strong dependence on initialization.

#### **Recurrent Neural Networks**

 $\square$  Types of gates – Here are the different types of gates that we encounter in a typical recurrent neural network:

Input gate	Forget gate	Output gate	Gate
Write to cell or not?	Erase a cell or not?	Reveal a cell or not?	How much writing?

 $\square$  LSTM – A long short-term memory (LSTM) network is a type of RNN model that avoids the vanishing gradient problem by adding 'forget' gates.

#### **Reinforcement Learning and Control**

The goal of reinforcement learning is for an agent to learn how to evolve in an environment.

**D** Markov decision processes – A Markov decision process (MDP) is a 5-tuple  $(S, A, \{P_{sa}\}, \gamma, R)$  where:

- ${\mathcal S}$  is the set of states
- ${\mathcal A}$  is the set of actions
- $\{P_{sa}\}$  are the state transition probabilities for  $s \in S$  and  $a \in A$
- $\gamma \in [0,1[$  is the discount factor
- $R: S \times A \longrightarrow \mathbb{R}$  or  $R: S \longrightarrow \mathbb{R}$  is the reward function that the algorithm wants to maximize

**D** Policy – A policy  $\pi$  is a function  $\pi : S \longrightarrow A$  that maps states to actions.

Remark: we say that we execute a given policy  $\pi$  if given a state s we take the action  $a = \pi(s)$ .

 $\Box$  Value function – For a given policy  $\pi$  and a given state s, we define the value function  $V^{\pi}$  as follows:

$$V^{\pi}(s) = E \bigg[ R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots | s_0 = s, \pi \bigg]$$

**Bellman equation** – The optimal Bellman equations characterizes the value function  $V^{\pi^*}$  of the optimal policy  $\pi^*$ :

$$V^{\pi^{*}}(s) = R(s) + \max_{a \in \mathcal{A}} \gamma \sum_{s' \in S} P_{sa}(s') V^{\pi^{*}}(s')$$

Remark: we note that the optimal policy  $\pi^*$  for a given state s is such that:

$$\pi^*(s) = \operatorname*{argmax}_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} P_{sa}(s') V^*(s')$$

 $\square$  Value iteration algorithm – The value iteration algorithm is in two steps:

• We initialize the value:

 $V_0(s)=0$ 

• We iterate the value based on the values before:

$$V_{i+1}(s) = R(s) + \max_{a \in \mathcal{A}} \left[ \sum_{s' \in \mathcal{S}} \gamma P_{sa}(s') V_i(s') \right]$$

**Maximum likelihood estimate** – The maximum likelihood estimates for the state transition probabilities are as follows:

$$P_{sa}(s') = \frac{\# \text{times took action } a \text{ in state } s \text{ and got to } s'}{\# \text{times took action } a \text{ in state } s}$$

 $\Box$  Q-learning – Q-learning is a model-free estimation of Q, which is done as follows:

 $Q(s,a) \leftarrow Q(s,a) + \alpha \left[ R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a) \right]$ 

# Super VIP Cheatsheet: Deep Learning

Afshine AMIDI and Shervine AMIDI

November 25, 2018

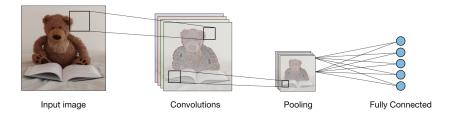
## Contents

1	Con	volutional Neural Networks	<b>2</b>
	1.1	Overview	2
	1.2	Types of layer	2
	1.3	Filter hyperparameters	2
	1.4	Tuning hyperparameters	3
	1.5	Commonly used activation functions	3
	1.6	Object detection	4
		1.6.1 Face verification and recognition	5
		1.6.2 Neural style transfer	5
		1.6.3 Architectures using computational tricks	6
<b>2</b>	Rec	urrent Neural Networks	7
	2.1	Overview	7
	2.2	Handling long term dependencies	8
	2.3	Learning word representation	9
		2.3.1 Motivation and notations	9
		2.3.2 Word embeddings	9
	2.4	Comparing words	9
	2.5		10
	2.6	0 0	10
	2.7		10
3	Dee	p Learning Tips and Tricks	11
	3.1		11
	3.2		12
	-	3	12
		3.2.2 Finding optimal weights	12
	3.3		12
			12
			$12^{$
	3.4		13
	3.5		13

## **1** Convolutional Neural Networks

1.1 Overview

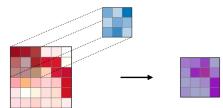
**Architecture of a traditional CNN** – Convolutional neural networks, also known as CNNs, are a specific type of neural networks that are generally composed of the following layers:



The convolution layer and the pooling layer can be fine-tuned with respect to hyperparameters that are described in the next sections.

#### 1.2Types of layer

**Convolutional layer (CONV)** – The convolution layer (CONV) uses filters that perform convolution operations as it is scanning the input I with respect to its dimensions. Its hyperparameters include the filter size F and stride S. The resulting output O is called *feature map* or activation map.

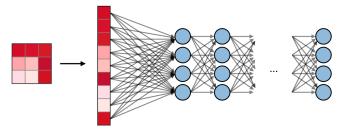


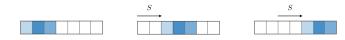
Remark: the convolution step can be generalized to the 1D and 3D cases as well.

**Pooling (POOL)** – The pooling layer (POOL) is a downsampling operation, typically applied after a convolution layer, which does some spatial invariance. In particular, max and average pooling are special kinds of pooling where the maximum and average value is taken, respectively.

	Max pooling	Average pooling
Purpose	Each pooling operation selects the maximum value of the current view	Each pooling operation averages the values of the current view
Illustration		avg
Comments	<ul><li> Preserves detected features</li><li> Most commonly used</li></ul>	- Downsamples feature map - Used in LeNet

**\Box** Fully Connected (FC) – The fully connected layer (FC) operates on a flattened input where each input is connected to all neurons. If present, FC layers are usually found towards the end of CNN architectures and can be used to optimize objectives such as class scores.



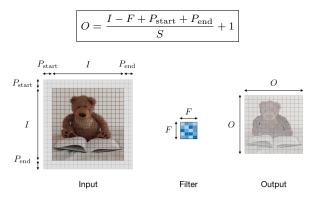


**D** Zero-padding – Zero-padding denotes the process of adding P zeroes to each side of the boundaries of the input. This value can either be manually specified or automatically set through one of the three modes detailed below:

	Valid	Same	Full
Value	P = 0	$P_{\text{start}} = \left\lfloor \frac{S \lceil \frac{I}{S} \rceil - I + F - S}{2} \right\rfloor$ $P_{\text{end}} = \left\lceil \frac{S \lceil \frac{I}{S} \rceil - I + F - S}{2} \right\rceil$	$P_{\text{start}} \in \llbracket 0, F - 1  rbracket$ $P_{\text{end}} = F - 1$
Illustration			
Purpose	- No padding - Drops last convolution if dimensions do not match	- Padding such that feature map size has size $\left\lceil \frac{I}{S} \right\rceil$ - Output size is mathematically convenient - Also called 'half' padding	<ul> <li>Maximum padding such that end convolutions are applied on the limits of the input</li> <li>Filter 'sees' the input end-to-end</li> </ul>

#### 1.4 Tuning hyperparameters

**D** Parameter compatibility in convolution layer – By noting I the length of the input volume size, F the length of the filter, P the amount of zero padding, S the stride, then the output size O of the feature map along that dimension is given by:

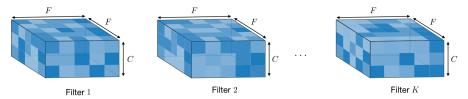


Remark: often times,  $P_{start} = P_{end} \triangleq P$ , in which case we can replace  $P_{start} + P_{end}$  by 2P in the formula above.

## 1.3 Filter hyperparameters

The convolution layer contains filters for which it is important to know the meaning behind its hyperparameters.

**Dimensions of a filter** – A filter of size  $F \times F$  applied to an input containing C channels is a  $F \times F \times C$  volume that performs convolutions on an input of size  $I \times I \times C$  and produces an output feature map (also called activation map) of size  $O \times O \times 1$ .



Remark: the application of K filters of size  $F \times F$  results in an output feature map of size  $O \times O \times K$ .

 $\Box$  Stride – For a convolutional or a pooling operation, the stride S denotes the number of pixels by which the window moves after each operation.

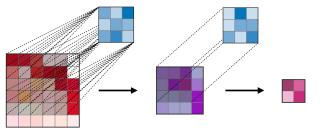
**Understanding the complexity of the model** – In order to assess the complexity of a model, it is often useful to determine the number of parameters that its architecture will have. In a given layer of a convolutional neural network, it is done as follows:

	CONV	POOL	FC
Illustration	$F \xrightarrow{F} \times K$	F max	N <sub>in</sub> N <sub>out</sub>
Input size	$I \times I \times C$	$I \times I \times C$	$N_{ m in}$
Output size	$O \times O \times K$	$O \times O \times C$	$N_{ m out}$
Number of parameters	$(F \times F \times C + 1) \cdot K$	0	$(N_{\rm in}+1) \times N_{\rm out}$
Remarks	<ul> <li>One bias parameter per filter</li> <li>In most cases, S &lt; F</li> <li>A common choice for K is 2C</li> </ul>	- Pooling operation done channel-wise - In most cases, $S = F$	<ul> <li>Input is flattened</li> <li>One bias parameter per neuron</li> <li>The number of FC neurons is free of structural constraints</li> </ul>

 $\Box$  Receptive field – The receptive field at layer k is the area denoted  $R_k \times R_k$  of the input that each pixel of the k-th activation map can 'see'. By calling  $F_j$  the filter size of layer j and  $S_i$  the stride value of layer i and with the convention  $S_0 = 1$ , the receptive field at layer k can be computed with the formula:

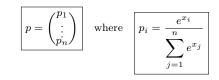
$$R_k = 1 + \sum_{j=1}^k (F_j - 1) \prod_{i=0}^{j-1} S_i$$

In the example below, we have  $F_1 = F_2 = 3$  and  $S_1 = S_2 = 1$ , which gives  $R_2 = 1+2 \cdot 1+2 \cdot 1 = 5$ .



# ReLULeaky ReLUELU $g(z) = \max(0,z)$ $g(z) = \max(\epsilon z, z)$ <br/>with $\epsilon \ll 1$ $g(z) = \max(\alpha(e^z - 1), z)$ <br/>with $\alpha \ll 1$ 1111011101110111011101110111011101110111011101110111011101110111011101110111

 $\Box$  Softmax – The softmax step can be seen as a generalized logistic function that takes as input a vector of scores  $x \in \mathbb{R}^n$  and outputs a vector of output probability  $p \in \mathbb{R}^n$  through a softmax function at the end of the architecture. It is defined as follows:



#### 1.6 Object detection

**Types of models** – There are 3 main types of object recognition algorithms, for which the nature of what is predicted is different. They are described in the table below:

Image classification	Classification w. localization	Detection
Teddy bear	Teddy bear	Teddy bear
- Classifies a picture	- Detects object in a picture - Predicts probability of	- Detects up to several objects in a picture
- Predicts probability of object	object and where it is located	- Predicts probabilities of objects and where they are located
Traditional CNN	Simplified YOLO, R-CNN	YOLO, R-CNN

#### 1.5 Commonly used activation functions

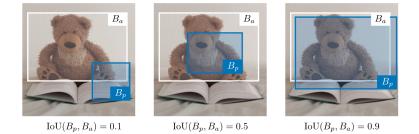
 $\Box$  Rectified Linear Unit – The rectified linear unit layer (ReLU) is an activation function g that is used on all elements of the volume. It aims at introducing non-linearities to the network. Its variants are summarized in the table below:

 $\Box$  Detection – In the context of object detection, different methods are used depending on whether we just want to locate the object or detect a more complex shape in the image. The two main ones are summed up in the table below:

Bounding box detection	Landmark detection
Detects the part of the image where the object is located	<ul><li>Detects a shape or characteristics of an object (e.g. eyes)</li><li>More granular</li></ul>
$b_h$ $(b_x, b_y)$ $b_w$	$\begin{array}{c}(l_{1x}, l_{1y}) & (l_{2x}, l_{2y}) \\(l_{4x}, l_{4y}) & (l_{7x}, l_{7y}) \\(l_{5x}, l_{5y}) & (l_{3x}, l_{3y}) & (l_{8x}, l_{8y}) \\(l_{6x}, l_{6y}) & (l_{9x}, l_{9y})\end{array}$
Box of center $(b_x, b_y)$ , height $b_h$ and width $b_w$	Reference points $(l_{1x}, l_{1y}),, (l_{nx}, l_{ny})$

 $\Box$  Intersection over Union – Intersection over Union, also known as IoU, is a function that quantifies how correctly positioned a predicted bounding box  $B_p$  is over the actual bounding box  $B_a$ . It is defined as:

$$IoU(B_p, B_a) = \frac{B_p \cap B_a}{B_p \cup B_a}$$

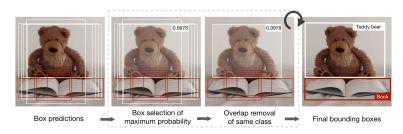


Remark: we always have  $IoU \in [0,1]$ . By convention, a predicted bounding box  $B_p$  is considered as being reasonably good if  $IoU(B_p, B_a) \ge 0.5$ .

 $\Box$  Anchor boxes – Anchor boxing is a technique used to predict overlapping bounding boxes. In practice, the network is allowed to predict more than one box simultaneously, where each box prediction is constrained to have a given set of geometrical properties. For instance, the first prediction can potentially be a rectangular box of a given form, while the second will be another rectangular box of a different geometrical form.

□ Non-max suppression – The non-max suppression technique aims at removing duplicate overlapping bounding boxes of a same object by selecting the most representative ones. After having removed all boxes having a probability prediction lower than 0.6, the following steps are repeated while there are boxes remaining:

- Step 1: Pick the box with the largest prediction probability.
- Step 2: Discard any box having an IoU  $\geqslant 0.5$  with the previous box.



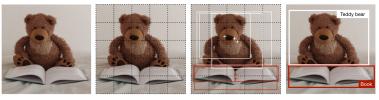
 $\square$  YOLO – You Only Look Once (YOLO) is an object detection algorithm that performs the following steps:

- Step 1: Divide the input image into a  $G \times G$  grid.
- Step 2: For each grid cell, run a CNN that predicts y of the following form:

$$y = \left[\underbrace{p_c, b_x, b_y, b_h, b_w, c_1, c_2, \dots, c_p}_{\text{repeated } k \text{ times}}, \dots\right]^T \in \mathbb{R}^{G \times G \times k \times (5+p)}$$

where  $p_c$  is the probability of detecting an object,  $b_x, b_y, b_h, b_w$  are the properties of the detected bouding box,  $c_1, \ldots, c_p$  is a one-hot representation of which of the p classes were detected, and k is the number of anchor boxes.

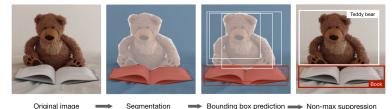
• Step 3: Run the non-max suppression algorithm to remove any potential duplicate overlapping bounding boxes.



Original image  $\implies$  Division in  $G \times G$  grid  $\implies$  Bounding box prediction  $\implies$  Non-max suppression

Remark: when  $p_c = 0$ , then the network does not detect any object. In that case, the corresponding predictions  $b_x, ..., c_p$  have to be ignored.

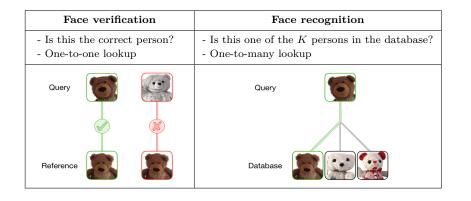
 $\Box$  **R-CNN** – Region with Convolutional Neural Networks (R-CNN) is an object detection algorithm that first segments the image to find potential relevant bounding boxes and then run the detection algorithm to find most probable objects in those bounding boxes.



Remark: although the original algorithm is computationally expensive and slow, newer architectures enabled the algorithm to run faster, such as Fast R-CNN and Faster R-CNN.

#### 1.6.1 Face verification and recognition

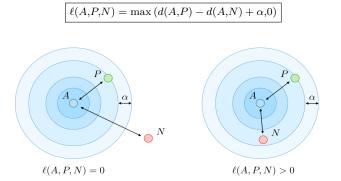
□ **Types of models** – Two main types of model are summed up in table below:



**Done Shot Learning** – One Shot Learning is a face verification algorithm that uses a limited training set to learn a similarity function that quantifies how different two given images are. The similarity function applied to two images is often noted d(image 1, image 2).

**Siamese Network** – Siamese Networks aim at learning how to encode images to then quantify how different two images are. For a given input image  $x^{(i)}$ , the encoded output is often noted as  $f(x^{(i)})$ .

 $\Box$  Triplet loss – The triplet loss  $\ell$  is a loss function computed on the embedding representation of a triplet of images A (anchor), P (positive) and N (negative). The anchor and the positive example belong to a same class, while the negative example to another one. By calling  $\alpha \in \mathbb{R}^+$  the margin parameter, this loss is defined as follows:



#### 1.6.2 Neural style transfer

 $\Box$  Motivation – The goal of neural style transfer is to generate an image G based on a given content C and a given style S.



 $\Box$  Activation – In a given layer l, the activation is noted  $a^{[l]}$  and is of dimensions  $n_H \times n_w \times n_c$ 

**Content cost function** – The content cost function  $J_{\text{content}}(C,G)$  is used to determine how the generated image G differs from the original content image C. It is defined as follows:

$$J_{\text{content}}(C,G) = \frac{1}{2} ||a^{[l](C)} - a^{[l](G)}||^2$$

 $\Box$  Style matrix – The style matrix  $G^{[l]}$  of a given layer l is a Gram matrix where each of its elements  $G^{[l]}_{kk'}$  quantifies how correlated the channels k and k' are. It is defined with respect to activations  $a^{[l]}$  as follows:

$$G_{kk'}^{[l]} = \sum_{i=1}^{n_H^{[l]}} \sum_{j=1}^{n_w^{[l]}} a_{ijk}^{[l]} a_{ijk'}^{[l]}$$

Remark: the style matrix for the style image and the generated image are noted  $G^{[l](S)}$  and  $G^{[l](G)}$  respectively.

**Style cost function** – The style cost function  $J_{\text{style}}(S,G)$  is used to determine how the generated image G differs from the style S. It is defined as follows:

$$J_{\text{style}}^{[l]}(S,G) = \frac{1}{(2n_H n_w n_c)^2} ||G^{[l](S)} - G^{[l](G)}||_F^2 = \frac{1}{(2n_H n_w n_c)^2} \sum_{k,k'=1}^{n_c} \left(G_{kk'}^{[l](S)} - G_{kk'}^{[l](G)}\right)^2$$

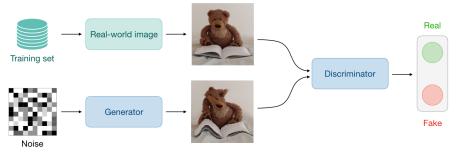
**Overall cost function** – The overall cost function is defined as being a combination of the content and style cost functions, weighted by parameters  $\alpha,\beta$ , as follows:

 $J(G) = \alpha J_{\text{content}}(C,G) + \beta J_{\text{style}}(S,G)$ 

Remark: a higher value of  $\alpha$  will make the model care more about the content while a higher value of  $\beta$  will make it care more about the style.

#### 1.6.3 Architectures using computational tricks

**Generative Adversarial Network** – Generative adversarial networks, also known as GANs, are composed of a generative and a discriminative model, where the generative model aims at generating the most truthful output that will be fed into the discriminative which aims at differentiating the generated and true image.



Remark: use cases using variants of GANs include text to image, music generation and synthesis.

**ResNet** – The Residual Network architecture (also called ResNet) uses residual blocks with a high number of layers meant to decrease the training error. The residual block has the following characterizing equation:

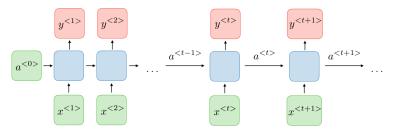
$$a^{[l+2]} = g(a^{[l]} + z^{[l+2]})$$

**D** Inception Network – This architecture uses inception modules and aims at giving a try at different convolutions in order to increase its performance. In particular, it uses the  $1 \times 1$  convolution trick to lower the burden of computation.

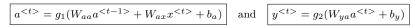
#### 2 Recurrent Neural Networks

#### 2.1 Overview

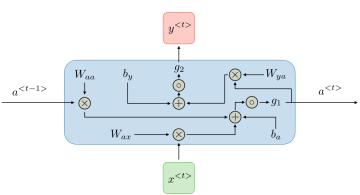
□ Architecture of a traditional RNN – Recurrent neural networks, also known as RNNs, are a class of neural networks that allow previous outputs to be used as inputs while having hidden states. They are typically as follows:



For each timestep t, the activation  $a^{\langle t \rangle}$  and the output  $y^{\langle t \rangle}$  are expressed as follows:



where  $W_{ax}, W_{aa}, W_{ya}, b_a, b_y$  are coefficients that are shared temporally and  $g_1, g_2$  activation functions



The pros and cons of a typical RNN architecture are summed up in the table below:

Advantages	Drawbacks
- Possibility of processing input of any length	- Computation being slow
- Model size not increasing with size of input	- Difficulty of accessing information
- Computation takes into account	from a long time ago
historical information	- Cannot consider any future input
- Weights are shared across time	for the current state

 $\Box$  Applications of RNNs – RNN models are mostly used in the fields of natural language processing and speech recognition. The different applications are summed up in the table below:

\* \* \*

Type of RNN	Illustration	Example
One-to-one $T_x = T_y = 1$	$\hat{y}$ $a^{<0>} \rightarrow \uparrow$ 1 x	Traditional neural network
One-to-many $T_x = 1, T_y > 1$	$a^{<0>} \rightarrow \begin{array}{c} y^{<1>} \\ 1 \\ x \end{array} \begin{array}{c} y^{<2>} \\ 1 \\ y \end{array} \begin{array}{c} y^{<2>} \\ 1 \\ y \end{array} \begin{array}{c} y^{<2>} \\ y^{<2>} \\ y \end{array} \begin{array}{c} y^{<2>} \\ y^{<2>} \\ y \end{array} \begin{array}{c} y^{<2>} \\ y^{<2>} \\ y^{<2>} \end{array} \begin{array}{c} y^{<2>} \\ y^{<2>} \\ y^{<2>} \\ y^{<2>} \end{array} \begin{array}{c} y^{<2>} \\ y^{<2>} \\ y^{<2>} \\ y^{<2>} \end{array} \begin{array}{c} y^{<2>} \\ y^{<2>} \\ y^{<2>} \\ y^{<2>} \end{array} \begin{array}{c} y^{<2>} \\ y^{<2>} \\ y^{<2>} \\ y^{<2>} \\ y^{<2>} \end{array} \begin{array}{c} y^{<2>} \\ $	Music generation
Many-to-one $T_x > 1, T_y = 1$	$ \begin{array}{c}                                     $	Sentiment classification
Many-to-many $T_x = T_y$	$a^{<0>} \rightarrow \begin{array}{c} \hat{y}^{<1>} & \hat{y}^{<2>} & \hat{y}^{} \\ \uparrow & \uparrow & \uparrow \\ x^{<1>} & x^{<2>} & x^{} \end{array}$	Name entity recognition
Many-to-many $T_x \neq T_y$	$ \begin{array}{c}  y^{<1>} \\  y^{<1>} \\  x^{<1>} \end{array} $	Machine translation

 $\Box$  Loss function – In the case of a recurrent neural network, the loss function  $\mathcal{L}$  of all time steps is defined based on the loss at every time step as follows:

$$\mathcal{L}(\widehat{y}, y) = \sum_{t=1}^{T_y} \mathcal{L}(\widehat{y}^{< t>}, y^{< t>})$$

$$\frac{\partial \mathcal{L}^{(T)}}{\partial W} = \sum_{t=1}^{T} \left. \frac{\partial \mathcal{L}^{(T)}}{\partial W} \right|_{(t)}$$

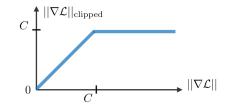
#### 2.2 Handling long term dependencies

**Commonly used activation functions** – The most common activation functions used in RNN modules are described below:

Sigmoid	Tanh	RELU
$g(z) = \frac{1}{1 + e^{-z}}$	$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$	$g(z) = \max(0, z)$

□ Vanishing/exploding gradient – The vanishing and exploding gradient phenomena are often encountered in the context of RNNs. The reason why they happen is that it is difficult to capture long term dependencies because of multiplicative gradient that can be exponentially decreasing/increasing with respect to the number of layers.

 $\square$  **Gradient clipping** – It is a technique used to cope with the exploding gradient problem sometimes encountered when performing backpropagation. By capping the maximum value for the gradient, this phenomenon is controlled in practice.



 $\Box$  Types of gates – In order to remedy the vanishing gradient problem, specific gates are used in some types of RNNs and usually have a well-defined purpose. They are usually noted  $\Gamma$  and are equal to:

$$\Gamma = \sigma(Wx^{} + Ua^{} + b)$$

timestep T, the derivative of the loss  $\mathcal{L}$  with respect to weight matrix W is expressed as follows: are summed up in the table below:

 $\Box$  Backpropagation through time – Backpropagation is done at each point in time. At where W, U, b are coefficients specific to the gate and  $\sigma$  is the sigmoid function. The main ones

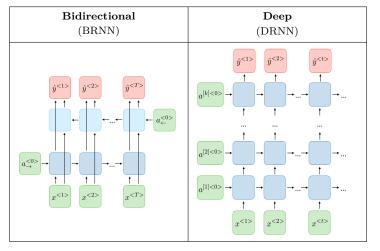
Type of gate	Role	Used in
Update gate $\Gamma_u$	How much past should matter now?	GRU, LSTM
Relevance gate $\Gamma_r$	Drop previous information?	GRU, LSTM
Forget gate $\Gamma_f$	Erase a cell or not?	LSTM
Output gate $\Gamma_o$	How much to reveal of a cell?	LSTM

 $\Box$  **GRU/LSTM** – Gated Recurrent Unit (GRU) and Long Short-Term Memory units (LSTM) deal with the vanishing gradient problem encountered by traditional RNNs, with LSTM being a generalization of GRU. Below is a table summing up the characterizing equations of each architecture:

	Gated Recurrent Unit (GRU)	Long Short-Term Memory (LSTM)
$\tilde{c}^{}$	$\tanh(W_c[\Gamma_r \star a^{}, x^{}] + b_c)$	$\tanh(W_c[\Gamma_r \star a^{}, x^{}] + b_c)$
$c^{\langle t \rangle}$	$\Gamma_u \star \tilde{c}^{\langle t \rangle} + (1 - \Gamma_u) \star c^{\langle t-1 \rangle}$	$\Gamma_u \star \tilde{c}^{} + \Gamma_f \star c^{}$
$a^{}$	$c^{\langle t \rangle}$	$\Gamma_o \star c^{}$
Dependencies	$c^{} \xrightarrow{\tilde{c}^{}} c^{}$ $a^{} \xrightarrow{\tilde{c}^{}} a^{}$	$c^{} \xrightarrow{\overline{C} < t>} c^{}$ $a^{} \xrightarrow{\Gamma_{f}} (\Gamma_{u}) (\Gamma_{r}) \xrightarrow{\Gamma_{o}} a^{}$ $a^{} \xrightarrow{T_{f}} x^{}$

Remark: the sign  $\star$  denotes the element-wise multiplication between two vectors.

 $\square$  Variants of RNNs – The table below sums up the other commonly used RNN architectures:

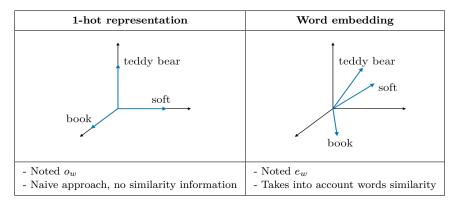


#### 2.3 Learning word representation

In this section, we note V the vocabulary and |V| its size.

#### 2.3.1 Motivation and notations

 $\square$  Representation techniques – The two main ways of representing words are summed up in the table below:



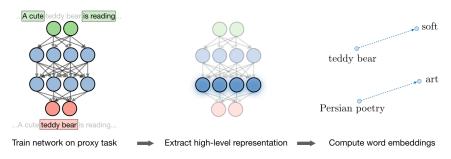
**D** Embedding matrix – For a given word w, the embedding matrix E is a matrix that maps its 1-hot representation  $o_w$  to its embedding  $e_w$  as follows:

 $e_w = Eo_w$ 

Remark: learning the embedding matrix can be done using target/context likelihood models.

## 2.3.2 Word embeddings

**Word2vec** – Word2vec is a framework aimed at learning word embeddings by estimating the likelihood that a given word is surrounded by other words. Popular models include skip-gram, negative sampling and CBOW.



**Skip-gram** – The skip-gram word2vec model is a supervised learning task that learns word embeddings by assessing the likelihood of any given target word t happening with a context word c. By noting  $\theta_t$  a parameter associated with t, the probability P(t|c) is given by:

$$P(t|c) = \frac{\exp(\theta_t^T e_c)}{\sum_{j=1}^{|V|} \exp(\theta_j^T e_c)}$$

Remark: summing over the whole vocabulary in the denominator of the softmax part makes this model computationally expensive. CBOW is another word2vec model using the surrounding words to predict a given word.

**D** Negative sampling – It is a set of binary classifiers using logistic regressions that aim at assessing how a given context and a given target words are likely to appear simultaneously, with the models being trained on sets of k negative examples and 1 positive example. Given a context word c and a target word t, the prediction is expressed by:

$$P(y=1|c,t) = \sigma(\theta_t^T e_c)$$

Remark: this method is less computationally expensive than the skip-gram model.

**GloVe** – The GloVe model, short for global vectors for word representation, is a word embedding technique that uses a co-occurence matrix X where each  $X_{i,j}$  denotes the number of times that a target *i* occurred with a context *j*. Its cost function *J* is as follows:

$$J(\theta) = \frac{1}{2} \sum_{i,j=1}^{|V|} f(X_{ij})(\theta_i^T e_j + b_i + b'_j - \log(X_{ij}))^2$$

here f is a weighting function such that  $X_{i,j} = 0 \Longrightarrow f(X_{i,j}) = 0$ .

Given the symmetry that e and  $\theta$  play in this model, the final word embedding  $e_w^{(\text{final})}$  is given by:

$$e_w^{\text{(final)}} = \frac{e_w + \theta_w}{2}$$

 ${\it Remark: the individual \ components \ of \ the \ learned \ word \ embeddings \ are \ not \ necessarily \ interpretable.}$ 

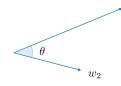
#### 2.4 Comparing words

 $\Box$  Cosine similarity – The cosine similarity between words  $w_1$  and  $w_2$  is expressed as follows:

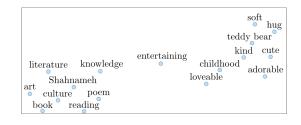
similarity = 
$$\frac{w_1 \cdot w_2}{||w_1|| \cdot ||w_2||} = \cos(\theta)$$

 $w_1$ 

Remark:  $\theta$  is the angle between words  $w_1$  and  $w_2$ .



 $\Box$  *t*-SNE – *t*-SNE (*t*-distributed Stochastic Neighbor Embedding) is a technique aimed at reducing high-dimensional embeddings into a lower dimensional space. In practice, it is commonly used to visualize word vectors in the 2D space.

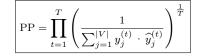


#### 2.5 Language model

 $\Box$  Overview – A language model aims at estimating the probability of a sentence P(y).

 $\Box$  *n*-gram model – This model is a naive approach aiming at quantifying the probability that an expression appears in a corpus by counting its number of appearance in the training data.

**Derplexity** – Language models are commonly assessed using the perplexity metric, also known as PP, which can be interpreted as the inverse probability of the dataset normalized by the number of words T. The perplexity is such that the lower, the better and is defined as follows:



Remark: PP is commonly used in t-SNE.

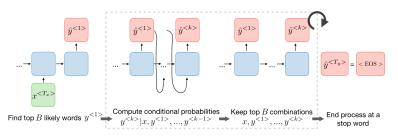
#### 2.6 Machine translation

**Overview** – A machine translation model is similar to a language model except it has an encoder network placed before. For this reason, it is sometimes referred as a conditional language model. The goal is to find a sentence y such that:

$$y = \arg\max_{y^{<1>},...,y^{}} P(y^{<1>},...,y^{}|x)$$

 $\Box$  Beam search – It is a heuristic search algorithm used in machine translation and speech recognition to find the likeliest sentence y given an input x.

- Step 1: Find top B likely words  $y^{<1>}$
- Step 2: Compute conditional probabilities  $y^{\leq k \geq} |x, y^{\leq 1 \geq}, \dots, y^{\leq k-1 \geq}$
- Step 3: Keep top B combinations  $x, y^{<1>}, \dots, y^{<k>}$



#### Remark: if the beam width is set to 1, then this is equivalent to a naive greedy search.

**Deam width** – The beam width B is a parameter for beam search. Large values of B yield to better result but with slower performance and increased memory. Small values of B lead to worse results but is less computationally intensive. A standard value for B is around 10.

**Length normalization** – In order to improve numerical stability, beam search is usually applied on the following normalized objective, often called the normalized log-likelihood objective, defined as:

Objective = 
$$\frac{1}{T_y^{\alpha}} \sum_{t=1}^{T_y} \log \left[ p(y^{}|x,y^{<1>},...,y^{}) \right]$$

Remark: the parameter  $\alpha$  can be seen as a softener, and its value is usually between 0.5 and 1.

**D** Error analysis – When obtaining a predicted translation  $\hat{y}$  that is bad, one can wonder why we did not get a good translation  $y^*$  by performing the following error analysis:

Case	$P(y^* x) > P(\widehat{y} x)$	$P(y^* x) \leqslant P(\widehat{y} x)$
Root cause	Beam search faulty	RNN faulty
Remedies	Increase beam width	<ul><li>Try different architecture</li><li>Regularize</li><li>Get more data</li></ul>

Remark: the attention scores are commonly used in image captioning and machine translation.





A cute teddy bear is reading Persian literature

A cute teddy bear is reading Persian literature

□ Attention weight – The amount of attention that the output  $y^{<t>}$  should pay to the activation  $a^{<t'>}$  is given by  $\alpha^{<t,t'>}$  computed as follows:

$$\boxed{\alpha^{< t, t'>} = \frac{\exp(e^{< t, t'>})}{\sum_{t''=1}^{T_x} \exp(e^{< t, t''>})}}$$

Remark: computation complexity is quadratic with respect to  $T_x$ .

\* \* \*

**D** Bleu score – The bilingual evaluation understudy (bleu) score quantifies how good a machine translation is by computing a similarity score based on n-gram precision. It is defined as follows:

bleu score = 
$$\exp\left(\frac{1}{n}\sum_{k=1}^{n}p_{k}\right)$$

where  $p_n$  is the bleu score on *n*-gram only defined as follows:

$$p_n = \frac{\sum_{\substack{\text{n-gram} \in \widehat{y} \\ \text{n-gram} \in \widehat{y}}} \text{count}_{\text{clip}}(\text{n-gram})}{\sum_{\substack{\text{n-gram} \in \widehat{y} \\ \text{out}}} \text{count}(\text{n-gram})}$$

Remark: a brevity penalty may be applied to short predicted translations to prevent an artificially inflated bleu score.

#### 2.7 Attention

□ Attention model – This model allows an RNN to pay attention to specific parts of the input that is considered as being important, which improves the performance of the resulting model in practice. By noting  $\alpha^{<t,t'>}$  the amount of attention that the output  $y^{<t>}$  should pay to the activation  $a^{<t'>}$  and  $c^{<t>}$  the context at time t, we have:

$$\boxed{c^{} = \sum_{t'} \alpha^{} a^{} = 1$$

#### 3 Deep Learning Tips and Tricks

#### 3.1 Data processing

**Data augmentation** – Deep learning models usually need a lot of data to be properly trained. It is often useful to get more data from the existing ones using data augmentation techniques. The main ones are summed up in the table below. More precisely, given the following input image, here are the techniques that we can apply:

Original	Flip	Rotation	Random crop
- Image without any modification	- Flipped with respect to an axis for which the meaning of the image is preserved	<ul> <li>Rotation with</li> <li>a slight angle</li> <li>Simulates incorrect</li> <li>horizon calibration</li> </ul>	<ul> <li>Random focus on one part of the image</li> <li>Several random crops can be done in a row</li> </ul>

Color shift	Noise addition	Information loss	Contrast change
<ul> <li>Nuances of RGB</li> <li>is slightly changed</li> <li>Captures noise</li> <li>that can occur</li> <li>with light exposure</li> </ul>	- Addition of noise - More tolerance to quality variation of inputs	<ul> <li>Parts of image</li> <li>ignored</li> <li>Mimics potential</li> <li>loss of parts of image</li> </ul>	<ul> <li>Luminosity changes</li> <li>Controls difference</li> <li>in exposition due</li> <li>to time of day</li> </ul>

**D** Batch normalization – It is a step of hyperparameter  $\gamma, \beta$  that normalizes the batch  $\{x_i\}$ . By noting  $\mu_B, \sigma_B^2$  the mean and variance of that we want to correct to the batch, it is done as follows:

$$x_i \longleftarrow \gamma \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

It is usually done after a fully connected/convolutional layer and before a non-linearity layer and aims at allowing higher learning rates and reducing the strong dependence on initialization.

#### 3.2 Training a neural network

#### 3.2.1 Definitions

 $\square$  **Epoch** – In the context of training a model, epoch is a term used to refer to one iteration where the model sees the whole training set to update its weights.

**Mini-batch gradient descent** – During the training phase, updating weights is usually not based on the whole training set at once due to computation complexities or one data point due to noise issues. Instead, the update step is done on mini-batches, where the number of data points in a batch is a hyperparameter that we can tune.

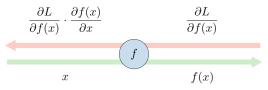
**D** Loss function – In order to quantify how a given model performs, the loss function L is usually used to evaluate to what extent the actual outputs y are correctly predicted by the model outputs z.

 $\Box$  Cross-entropy loss – In the context of binary classification in neural networks, the crossentropy loss L(z,y) is commonly used and is defined as follows:

$$L(z,y) = -\left[y \log(z) + (1-y) \log(1-z)\right]$$

#### 3.2.2 Finding optimal weights

**Backpropagation** – Backpropagation is a method to update the weights in the neural network by taking into account the actual output and the desired output. The derivative with respect to each weight w is computed using the chain rule.

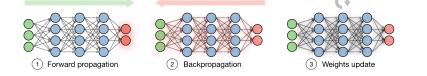


Using this method, each weight is updated with the rule:

$w \longleftarrow w - \alpha \frac{\partial L(z,y)}{\partial w}$	
--	--

□ Updating weights – In a neural network, weights are updated as follows:

- Step 1: Take a batch of training data and perform forward propagation to compute the loss.
- Step 2: Backpropagate the loss to get the gradient of the loss with respect to each weight.
- Step 3: Use the gradients to update the weights of the network.



#### 3.3 Parameter tuning

#### 3.3.1 Weights initialization

**Xavier initialization** – Instead of initializing the weights in a purely random manner, Xavier initialization enables to have initial weights that take into account characteristics that are unique to the architecture.

**Transfer learning** – Training a deep learning model requires a lot of data and more importantly a lot of time. It is often useful to take advantage of pre-trained weights on huge datasets that took days/weeks to train, and leverage it towards our use case. Depending on how much data we have at hand, here are the different ways to leverage this:

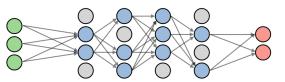
Training size	Illustration	Explanation
Small		Freezes all layers, trains weights on softmax
Medium		Freezes most layers, trains weights on last layers and softmax
Large		Trains weights on layers and softmax by initializing weights on pre-trained ones

Method	Explanation	$\mathbf{Update of}\ w$	$\mathbf{Update} \ \mathbf{of} \ b$
Momentum	<ul><li>Dampens oscillations</li><li>Improvement to SGD</li><li>2 parameters to tune</li></ul>	$w - \alpha v_{dw}$	$b - lpha v_{db}$
RMSprop	<ul> <li>Root Mean Square propagation</li> <li>Speeds up learning algorithm</li> <li>by controlling oscillations</li> </ul>	$w - \alpha \frac{dw}{\sqrt{s_{dw}}}$	$b \longleftarrow b - \alpha \frac{db}{\sqrt{s_{db}}}$
Adam	<ul> <li>Adaptive Moment estimation</li> <li>Most popular method</li> <li>4 parameters to tune</li> </ul>	$w - \alpha \frac{v_{dw}}{\sqrt{s_{dw}} + \epsilon}$	$b \longleftarrow b - \alpha \frac{v_{db}}{\sqrt{s_{db}} + \epsilon}$

Remark: other methods include Adadelta, Adagrad and SGD.

#### 3.4 Regularization

**Dropout** – Dropout is a technique used in neural networks to prevent overfitting the training data by dropping out neurons with probability p > 0. It forces the model to avoid relying too much on particular sets of features.



*Remark:* most deep learning frameworks parametrize dropout through the 'keep' parameter 1-p.

**Weight regularization** – In order to make sure that the weights are not too large and that the model is not overfitting the training set, regularization techniques are usually performed on the model weights. The main ones are summed up in the table below:

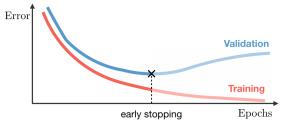
LASSO	Ridge	Elastic Net
<ul><li>Shrinks coefficients to 0</li><li>Good for variable selection</li></ul>	Makes coefficients smaller	Tradeoff between variable selection and small coefficients
$  \theta  _1 \leq 1$	$  \theta  _2 \leqslant 1$	$(1-\alpha)  \theta  _1 + \alpha   \theta  _2^2 \leqslant 1$
$egin{array}{lll} \ldots + \lambda    heta  _1 \ \lambda \in \mathbb{R} \end{array}$	$egin{aligned} & \ldots + \lambda     heta   _2^2 \ & \lambda \in \mathbb{R} \end{aligned}$	$\dots + \lambda \left[ (1 - \alpha)   \theta  _1 + \alpha   \theta  _2^2 \right]$ $\lambda \in \mathbb{R}, \alpha \in [0, 1]$

#### 3.3.2 Optimizing convergence

**D** Learning rate – The learning rate, often noted  $\alpha$  or sometimes  $\eta$ , indicates at which pace the weights get updated. It can be fixed or adaptively changed. The current most popular method is called Adam, which is a method that adapts the learning rate.

 $\Box$  Adaptive learning rates – Letting the learning rate vary when training a model can reduce the training time and improve the numerical optimal solution. While Adam optimizer is the most commonly used technique, others can also be useful. They are summed up in the table below:

 $\Box$  Early stopping – This regularization technique stops the training process as soon as the validation loss reaches a plateau or starts to increase.



#### **3.5** Good practices

 $\Box$  Overfitting small batch – When debugging a model, it is often useful to make quick tests to see if there is any major issue with the architecture of the model itself. In particular, in order to make sure that the model can be properly trained, a mini-batch is passed inside the network to see if it can overfit on it. If it cannot, it means that the model is either too complex or not complex enough to even overfit on a small batch, let alone a normal-sized training set.

**Gradient checking** – Gradient checking is a method used during the implementation of the backward pass of a neural network. It compares the value of the analytical gradient to the numerical gradient at given points and plays the role of a sanity-check for correctness.

	Numerical gradient	Analytical gradient
Formula	$\frac{df}{dx}(x) \approx \frac{f(x+h) - f(x-h)}{2h}$	$\frac{df}{dx}(x) = f'(x)$
Comments	<ul> <li>Expensive; loss has to be computed two times per dimension</li> <li>Used to verify correctness of analytical implementation</li> <li>Trade-off in choosing h not too small (numerical instability) nor too large (poor gradient approx.)</li> </ul>	<ul><li>'Exact' result</li><li>Direct computation</li><li>Used in the final implementation</li></ul>

\* \* \*

# Cheat Sheet – Convolutional Neural Network

## **Convolutional Neural Network:**

The data gets into the CNN through the input layer and passes through various hidden layers before getting to the output layer. The output of the network is compared to the actual labels in terms of loss or error. The partial derivatives of this loss w.r.t the trainable weights are calculated, and the weights are updated through one of the various methods using backpropagation.

## **CNN** Template:

3.

a.

b.

c.

0 O

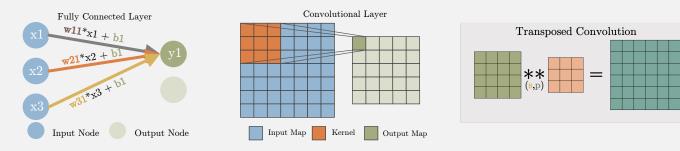
Sigmoid

0

 $\alpha = 10.0$ 

Most of the commonly used hidden layers (not all) follow a pattern

- 1. Layer function: Basic transforming function such as convolutional or fully connected layer.
- a. Fully Connected: Linear functions between the input and the
- a. Octividuational Layers: These layers are applied to 2D (3D) input feature maps. The trainable weights are a 2D (3D) kernel/filter that moves across the input feature map, generating dot products with the overlapping region of the input feature map.
- b.Transposed Convolutional (DeConvolutional) Layer: Usually used to increase the size of the output feature map (Upsampling) The idea behind the transposed convolutional layer is to undo (not exactly) the convolutional layer



**Pooling:** Non-trainable layer to change the size of the feature map 2.

Activation: Introduce non-linearity so

efficiently map non-linear complex mapping.

- Max/Average Pooling: Decrease the spatial size of the input layer based on a. selecting the maximum/average value in receptive field defined by the kernel
- **b**. **UnPooling:** A non-trainable layer used to increase the spatial size of the input layer based on placing the input pixel at a certain index in the receptive field of the output defined by the kernel.
- **3.** Normalization: Usually used just before the activation functions to limit the unbounded activation from increasing the output layer values too high
- Output Local Response Normalization LRN: A non-trainable layer that square-normalizes the pixel values in a feature map ล. within a local neighborhood.
- **Batch Normalization:** A trainable approach to normalizing the data by learning scale and shift variable during training. b.

can

CNN

5. Loss function: Quantifies how far off the CNN prediction is from the actual labels.

Type: max'pool - Stride: 1 Padding: 1

0

0

0

0 0

Regression Loss Functions: MAE, MSE, Huber loss a.

4.3 5 12

12 12 6

8.5 8.4

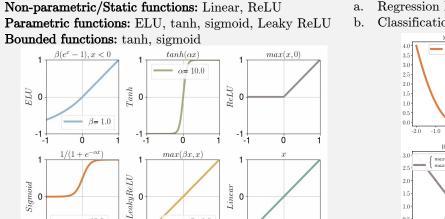
8.3 5.8 9.7

7.6 6 10

3.9 11 5.7 3.6 11

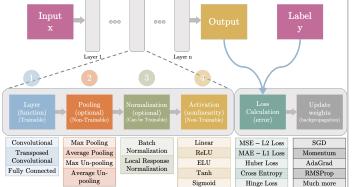
13 7.1

- Classification Loss Functions: Cross entropy, Hinge loss
  - MAE Loss Huber L MSE Los  $mse = (x - \hat{x})^2$  $\begin{cases} \frac{1}{2}(x - \hat{x})^2 \\ \gamma |x - \hat{x}| - \frac{1}{2}\gamma^2 \end{cases}$ 1.75 3.0 1.3 1.25 1.25 2.0 1.0 1.0 0.75 0.7 1.0 0.5 0.5 0.5 0.25 0.25 0.0 0.0 0.0 0.0 1.0 2.0 1.0 0.0 1.0 Hinge Loss Cross Entropy Los 3.0 -ylog(p) - (1 - y)log(1 - p) $ax(0, 1 - \hat{x})$  : x = 1 $ax(0, 1 + \hat{x})$  : x = -18.0 2.3 2.0 6.0 1.3 4.0 1.0 2.0 0.5



Source: https://www.cheatsheets.ageel-anwar.com Tutorial: Click here

 $\beta = 1.0$ 



Dataset: (x, y)

# Cheat Sheet – Famous CNNs

## AlexNet - 2012

Why: AlexNet was born out of the need to improve the results of the ImageNet challenge.

What: The network consists of 5 Convolutional (CONV) layers and 3 Fully Connected (FC) layers. The activation used is the Rectified Linear Unit (ReLU).

**How:** Data augmentation is carried out to reduce over-fitting, Uses Local response localization.

## VGGNet - 2014

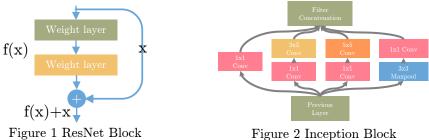
**Why:** VGGNet was born out of the need to reduce the # of parameters in the CONV layers and improve on training time

What: There are multiple variants of VGGNet (VGG16, VGG19, etc.) How: The important point to note here is that all the conv kernels are of size 3x3 and maxpool kernels are of size 2x2 with a stride of two.

## ResNet - 2015

Why: Neural Networks are notorious for not being able to find a simpler mapping when it exists. ResNet solves that.

What: There are multiple versions of ResNetXX architectures where 'XX' denotes the number of layers. The most used ones are ResNet50 and ResNet101. Since the vanishing gradient problem was taken care of (more about it in the How part), CNN started to get deeper and deeper How: ResNet architecture makes use of shortcut connections do solve the vanishing gradient problem. The basic building block of ResNet is a Residual block that is repeated throughout the network.



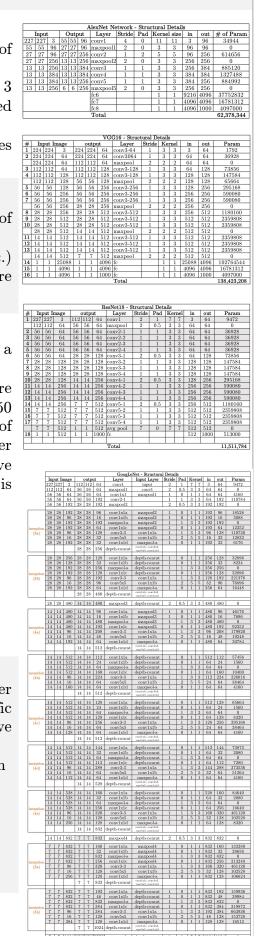
## Inception -2014

Why: Lager kernels are preferred for more global features, on the other hand, smaller kernels provide good results in detecting area-specific features. For effective recognition of such a variable-sized feature, we need kernels of different sizes. That is what Inception does.

What: The Inception network architecture consists of several inception modules of the following structure. Each inception module consists of four operations in parallel, 1x1 conv layer, 3x3 conv layer, 5x5 conv layer, max pooling

**How:** Inception increases the network space from which the best network is to be chosen via training. Each inception module can capture salient features at different levels.

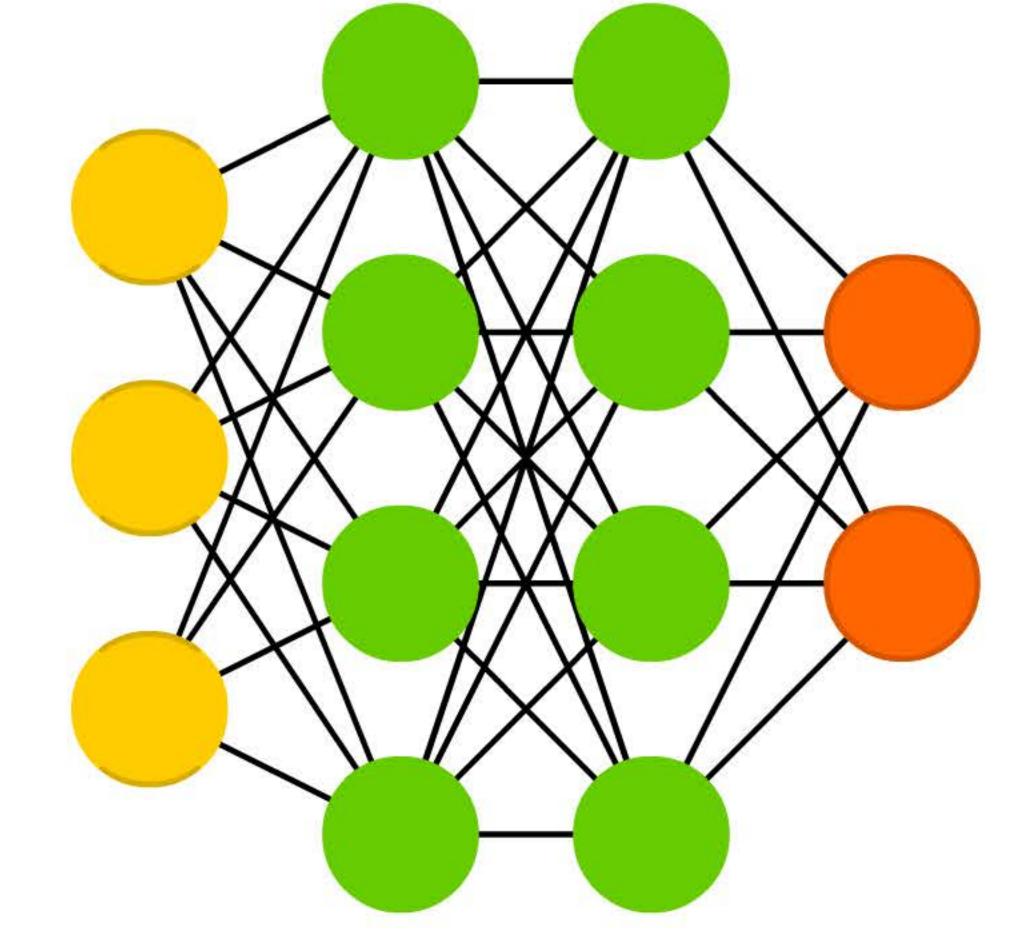
Comparison					
Network	Year	Salient Feature	top5 accuracy	Parameters	FLOP
AlexNet	2012	Deeper	84.70%	62M	1.5B
VGGNet	2014	Fixed-size kernels	92.30%	138M	19.6B
Inception	2014	Wider - Parallel kernels	93.30%	6.4M	2B
ResNet-152	2015	Shortcut connections	95.51%	60.3M	11B



Source: https://www.cheatsheets.aqeel-anwar.com Tutorial: Click here

# A mostly complete chart of Neural Networks ©2019 Fjodor van Veen & Stefan Leijnen asimovinstitute.org Radial Basis Network (RBF) Feed Forward (FF) Perceptron (P) Recurrent Neural Network (RNN) Long / Short Term Memory (LSTM)

Deep Feed Forward (DFF)



Spiking Hidden Cell

Capsule Cell

Input Cell

Backfed Input Cell

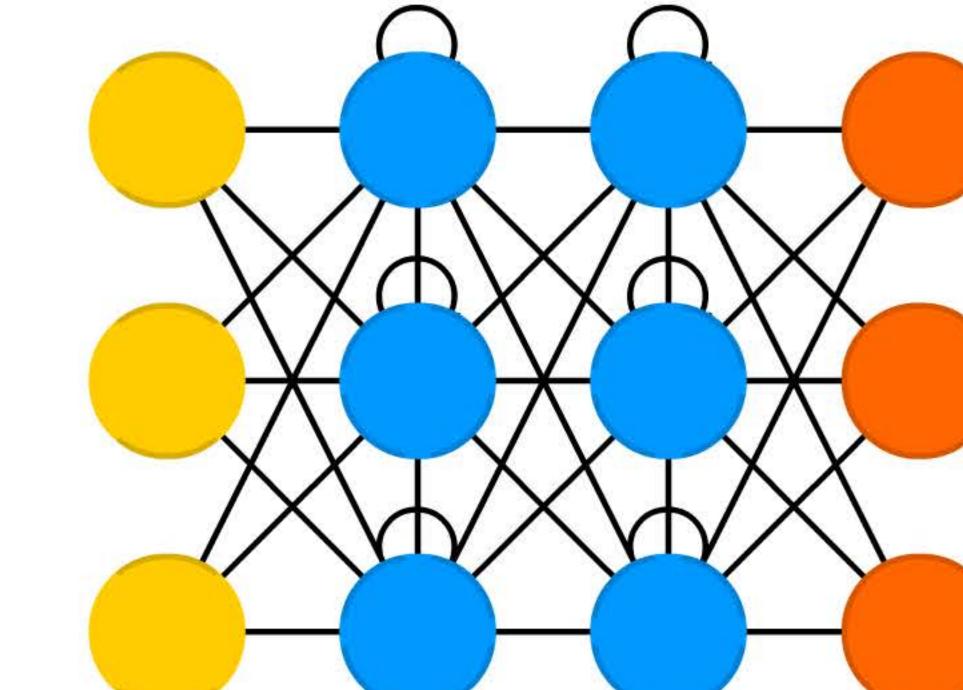
Noisy Input Cell

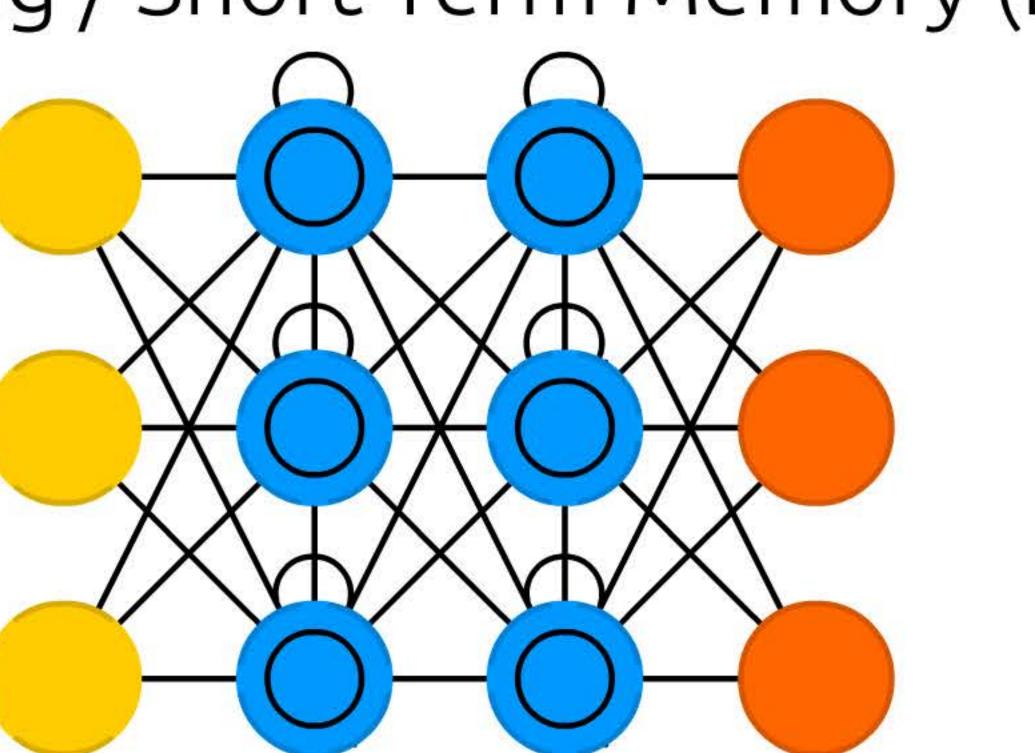
Hidden Cell

Output Cell

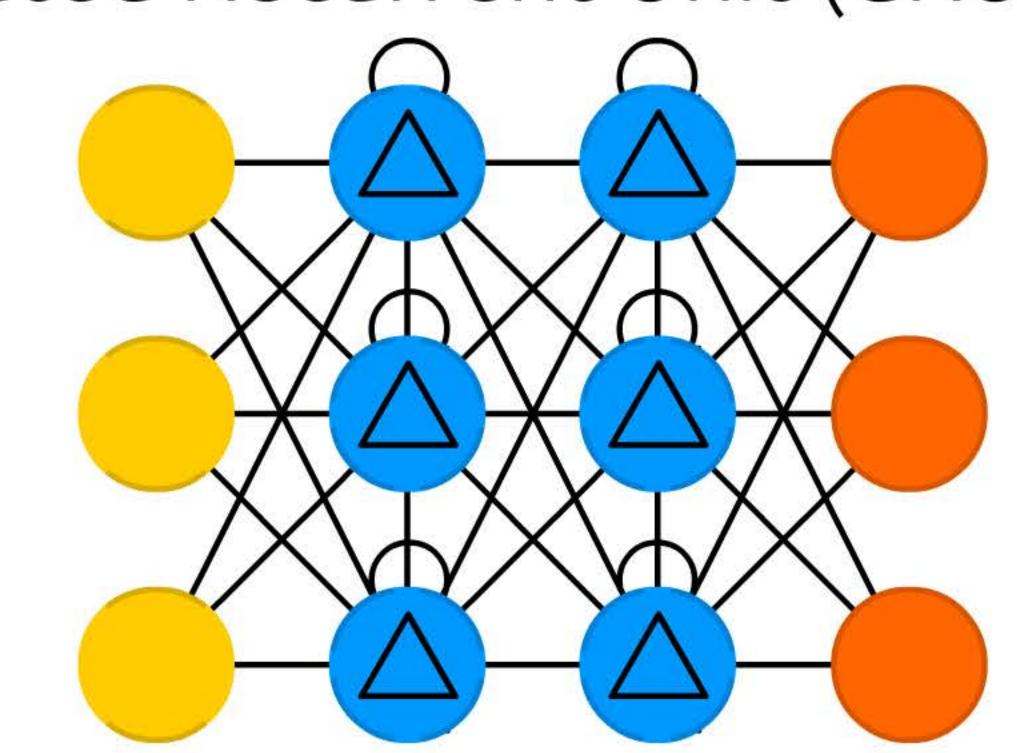
Match Input Output Cell

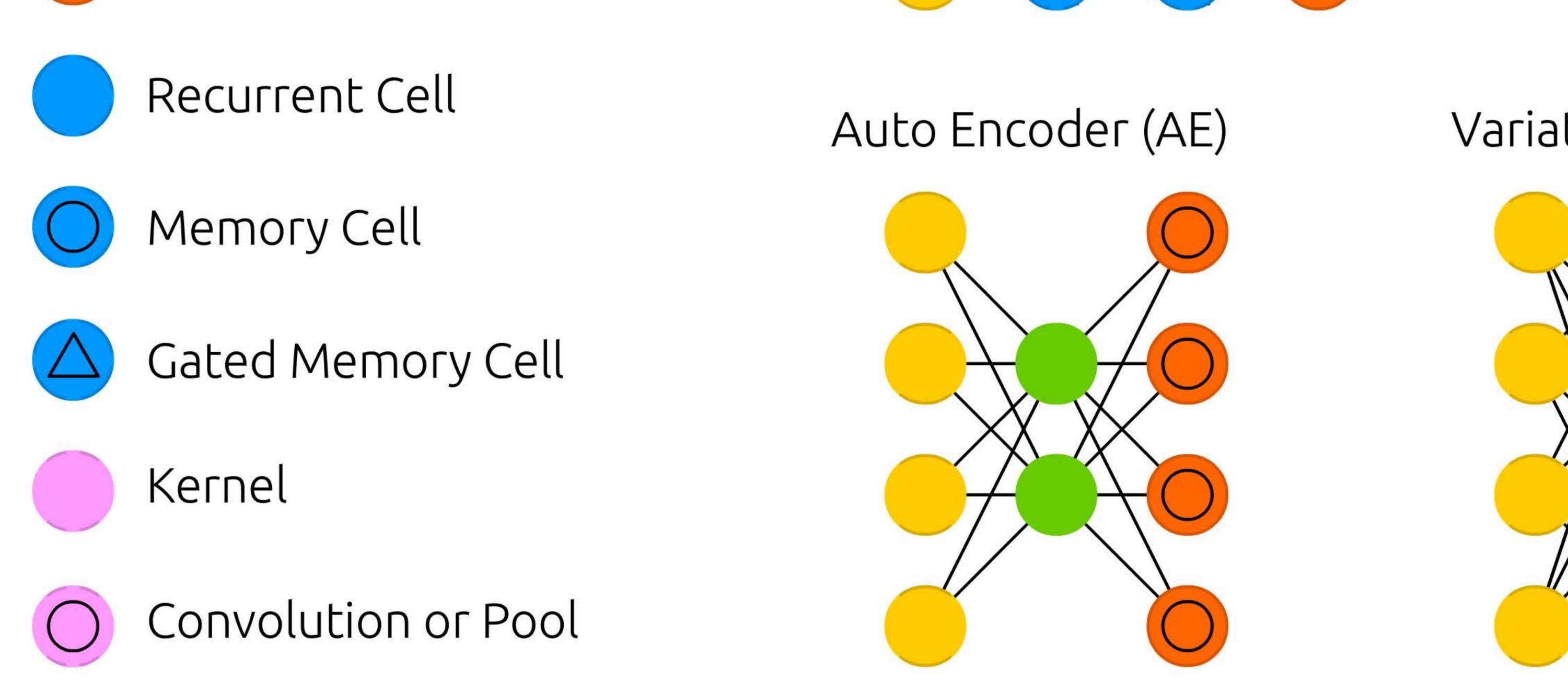
Probablistic Hidden Cell

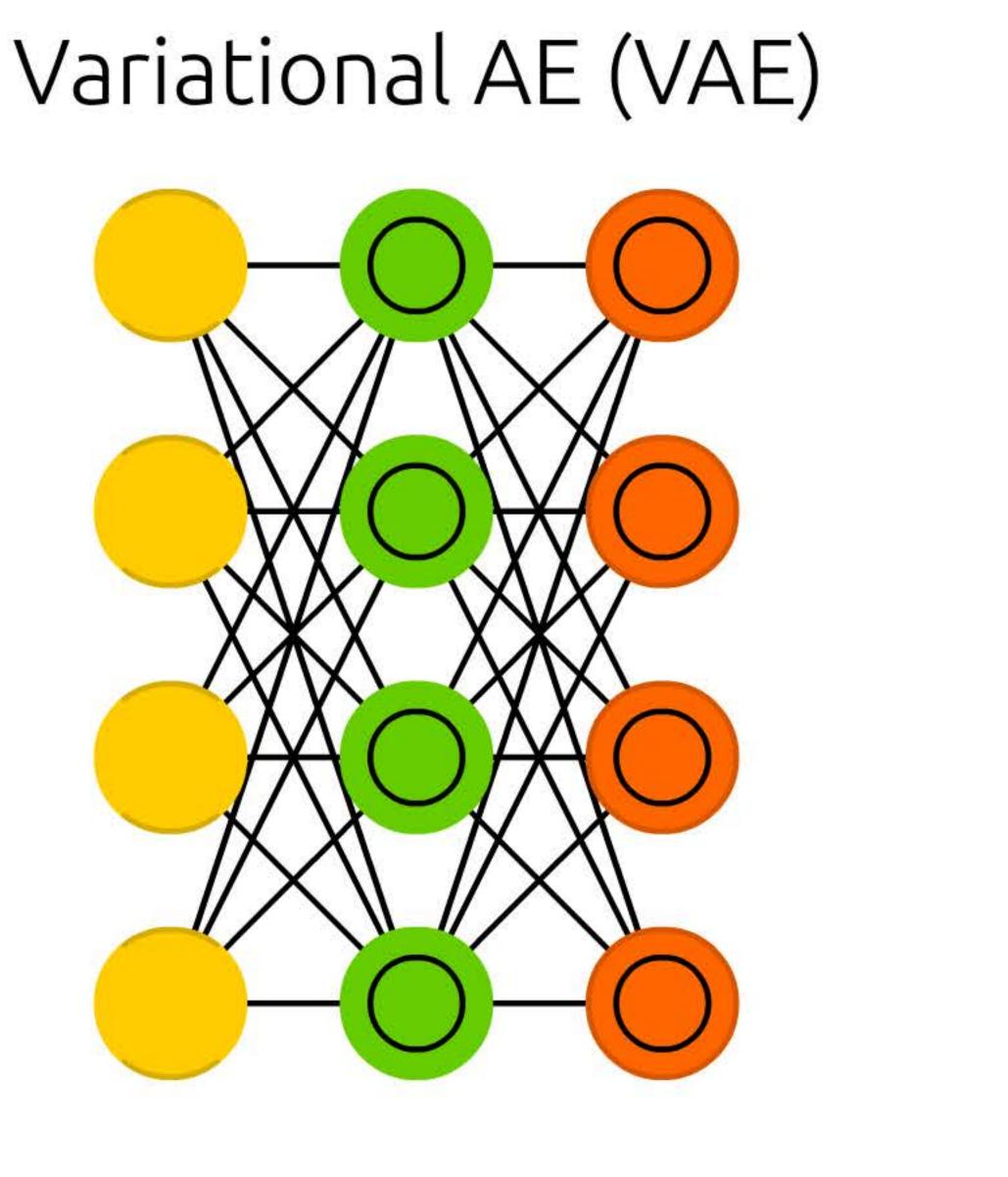


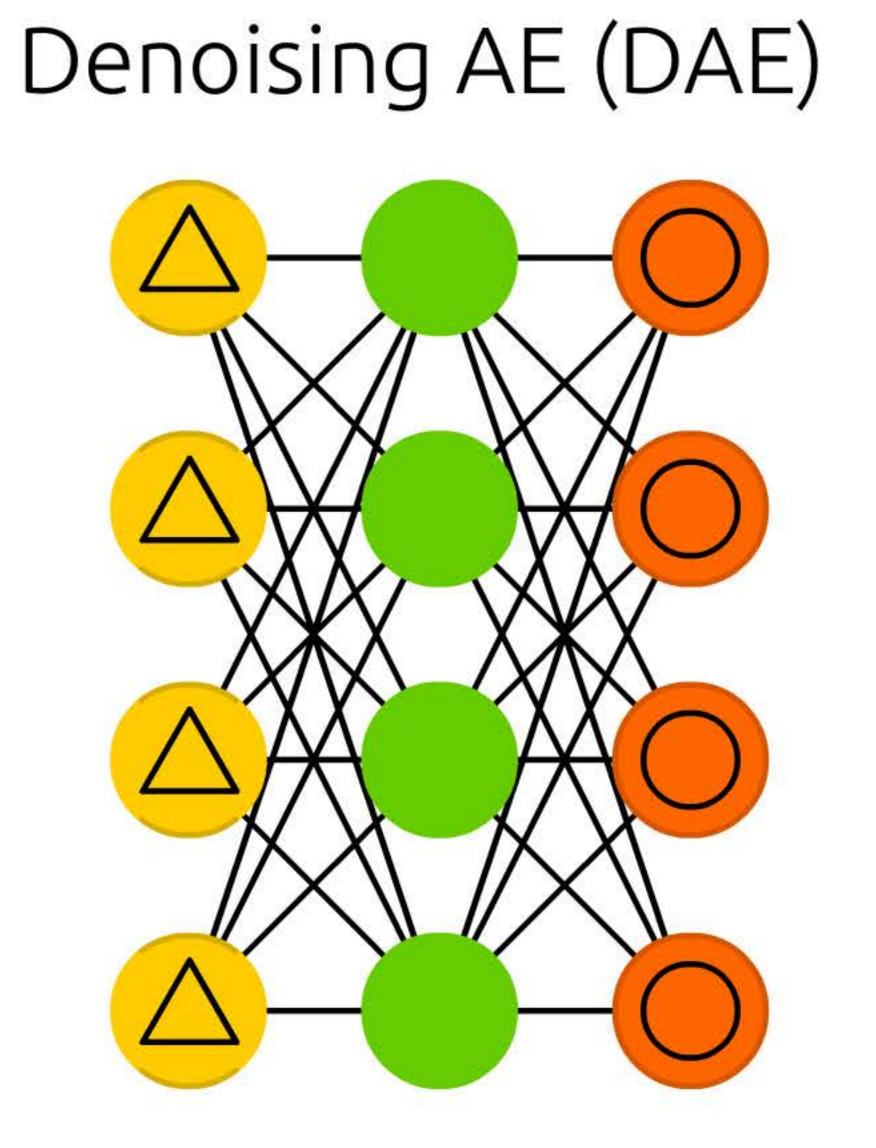


Gated Recurrent Unit (GRU)

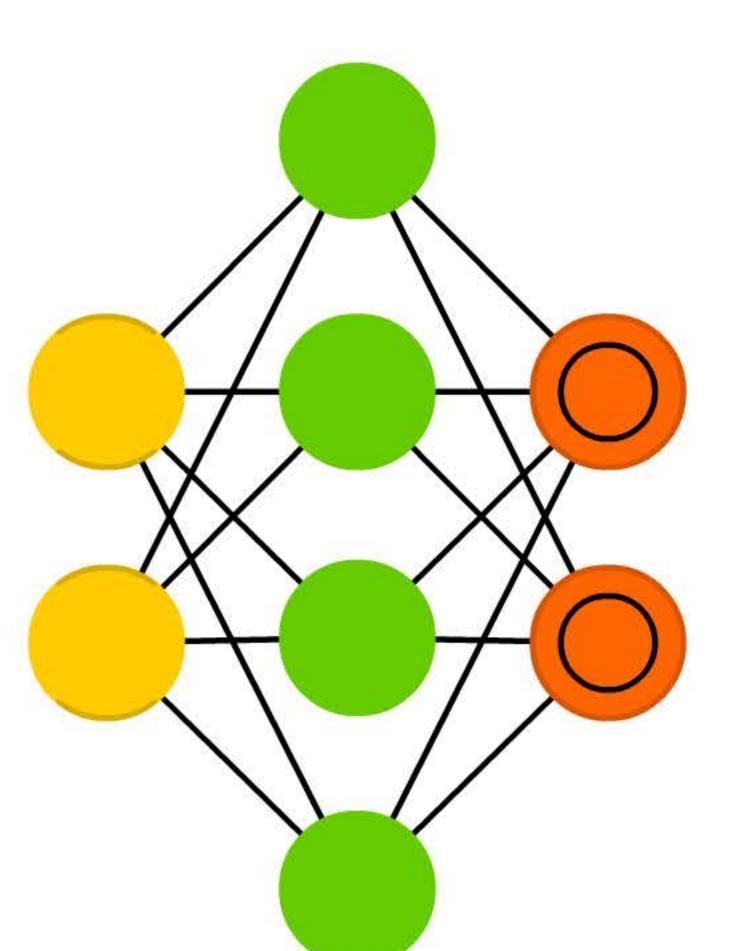








Sparse AE (SAE)

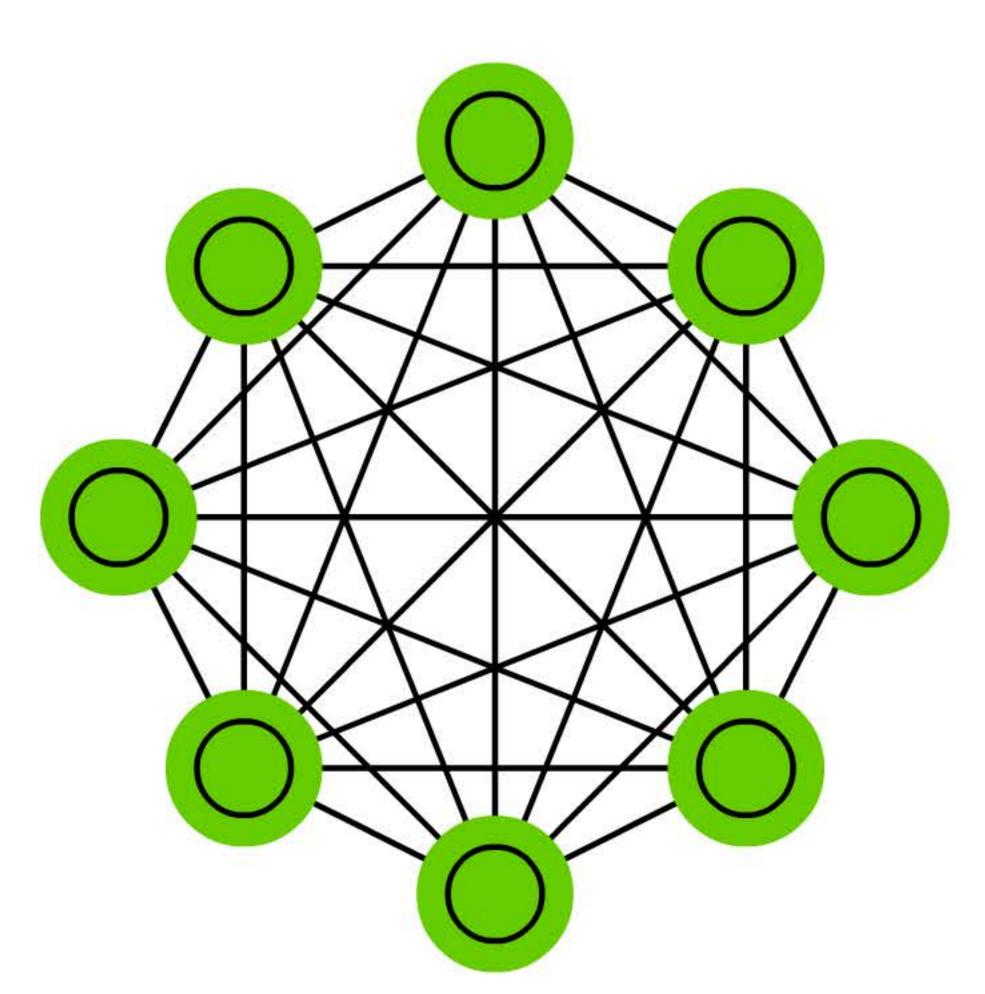


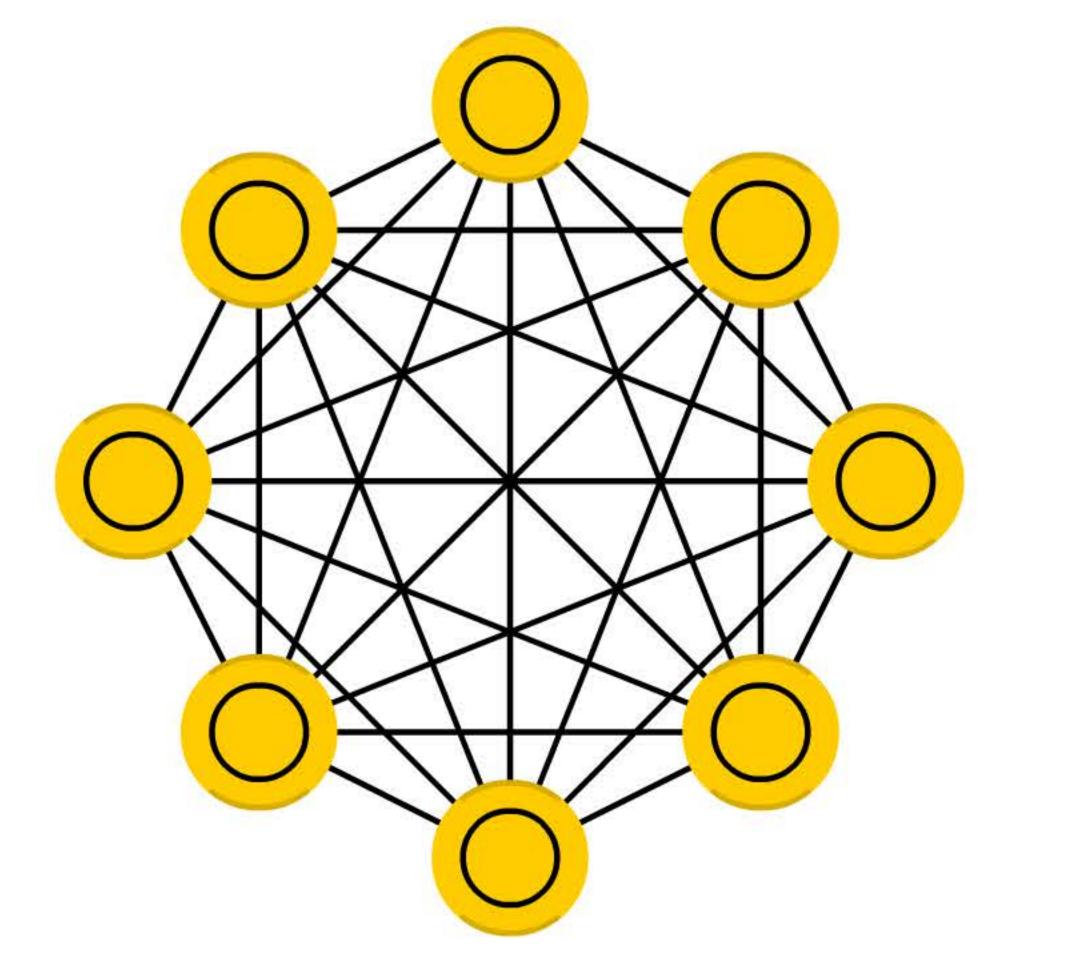
Markov Chain (MC)

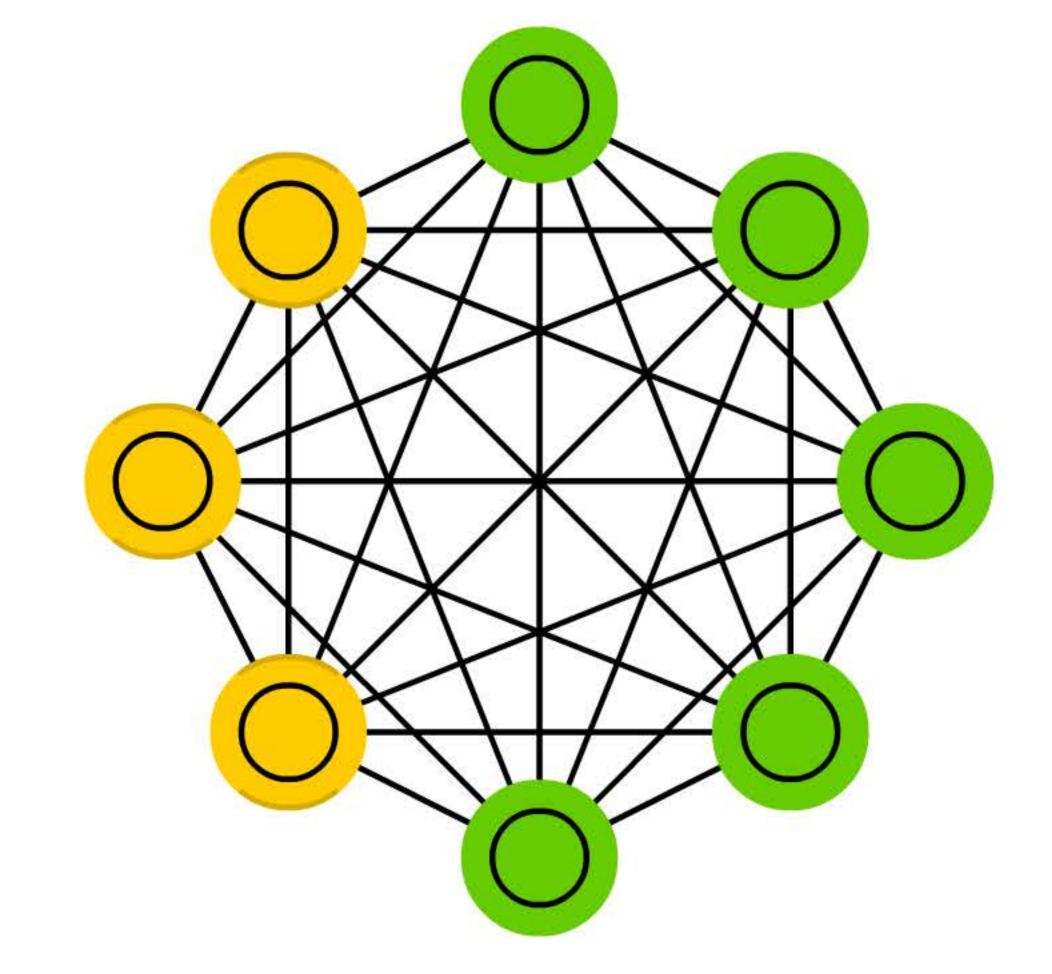
Hopfield Network (HN) Boltzmann Machine (BM)

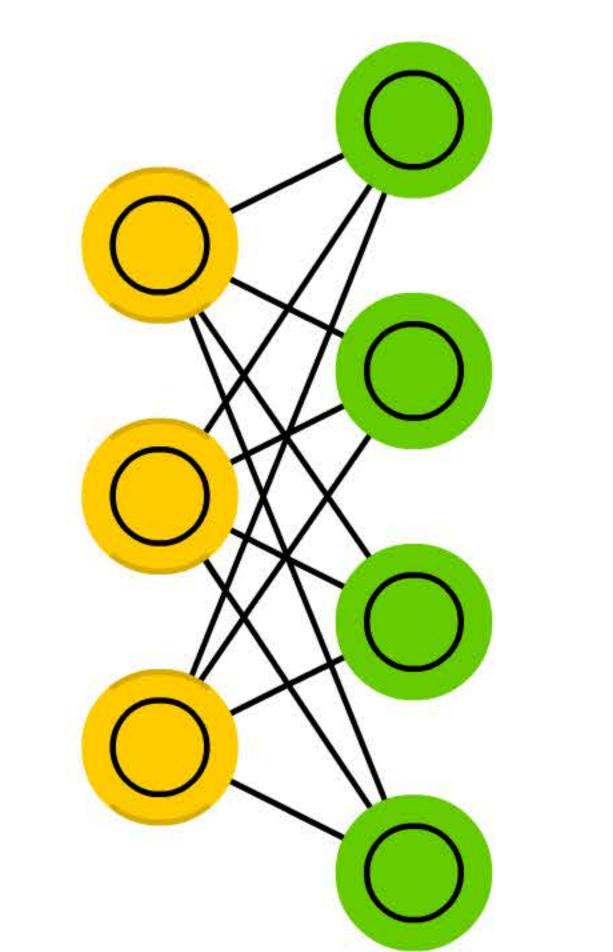
Restricted BM (RBM)

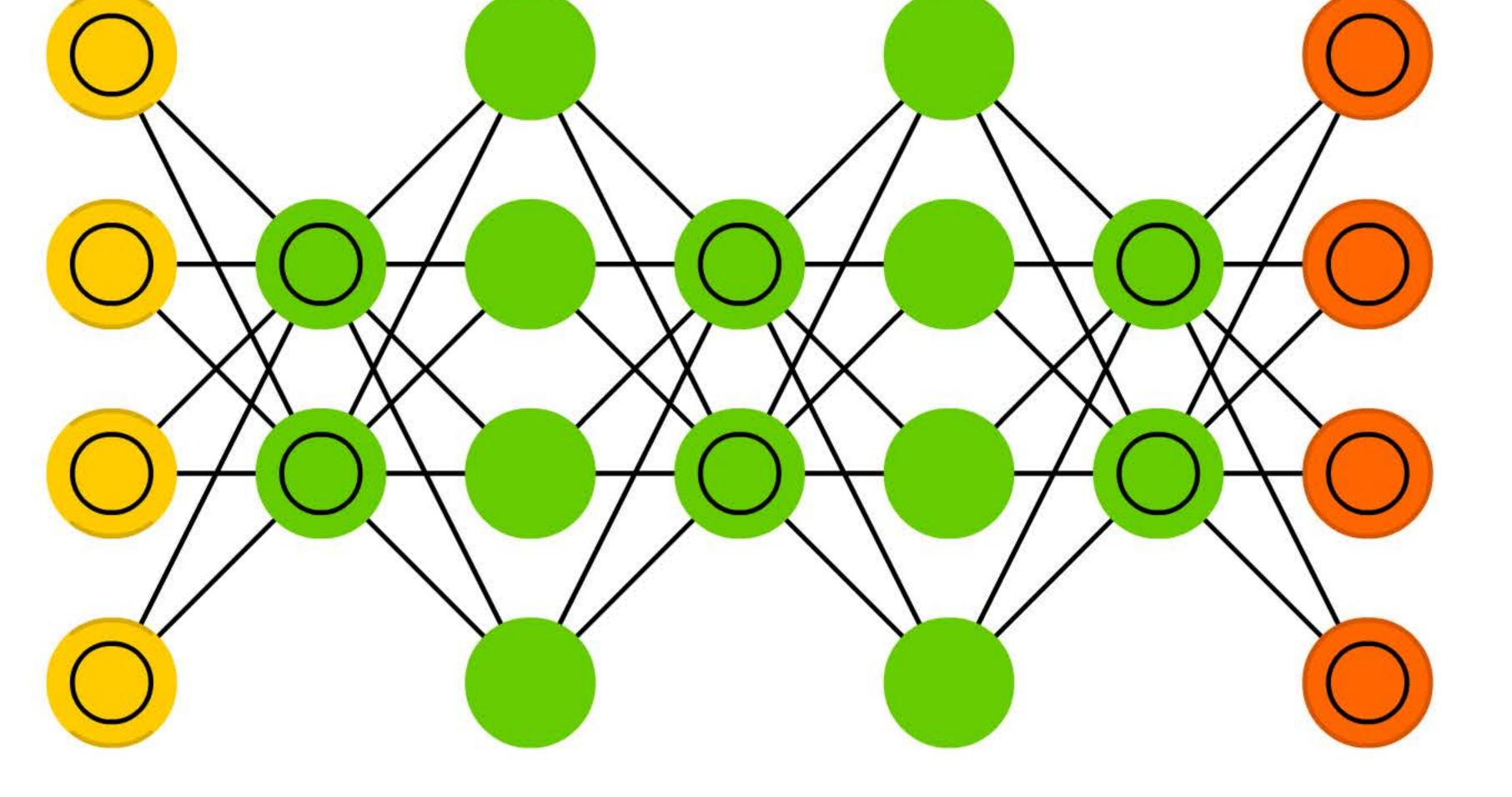
Deep Belief Network (DBN)

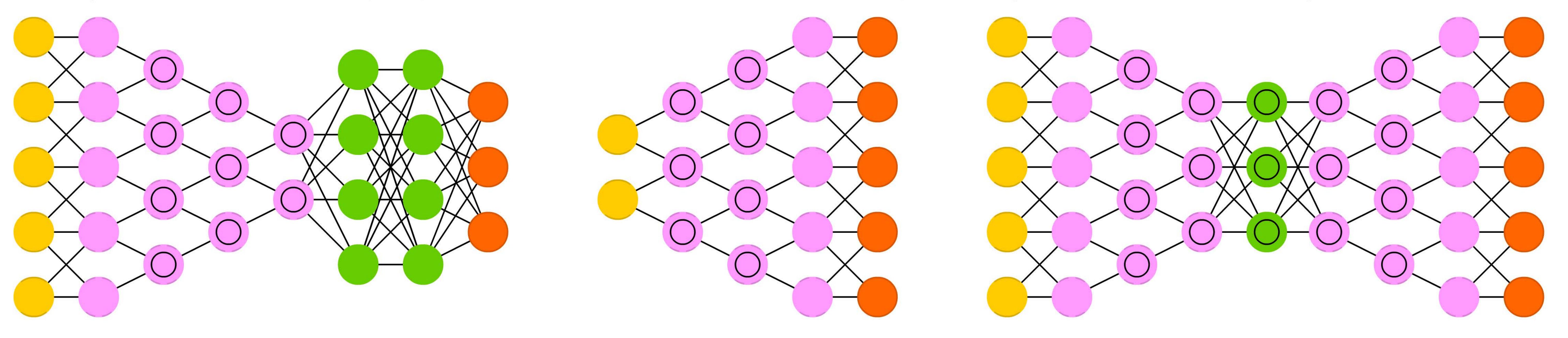




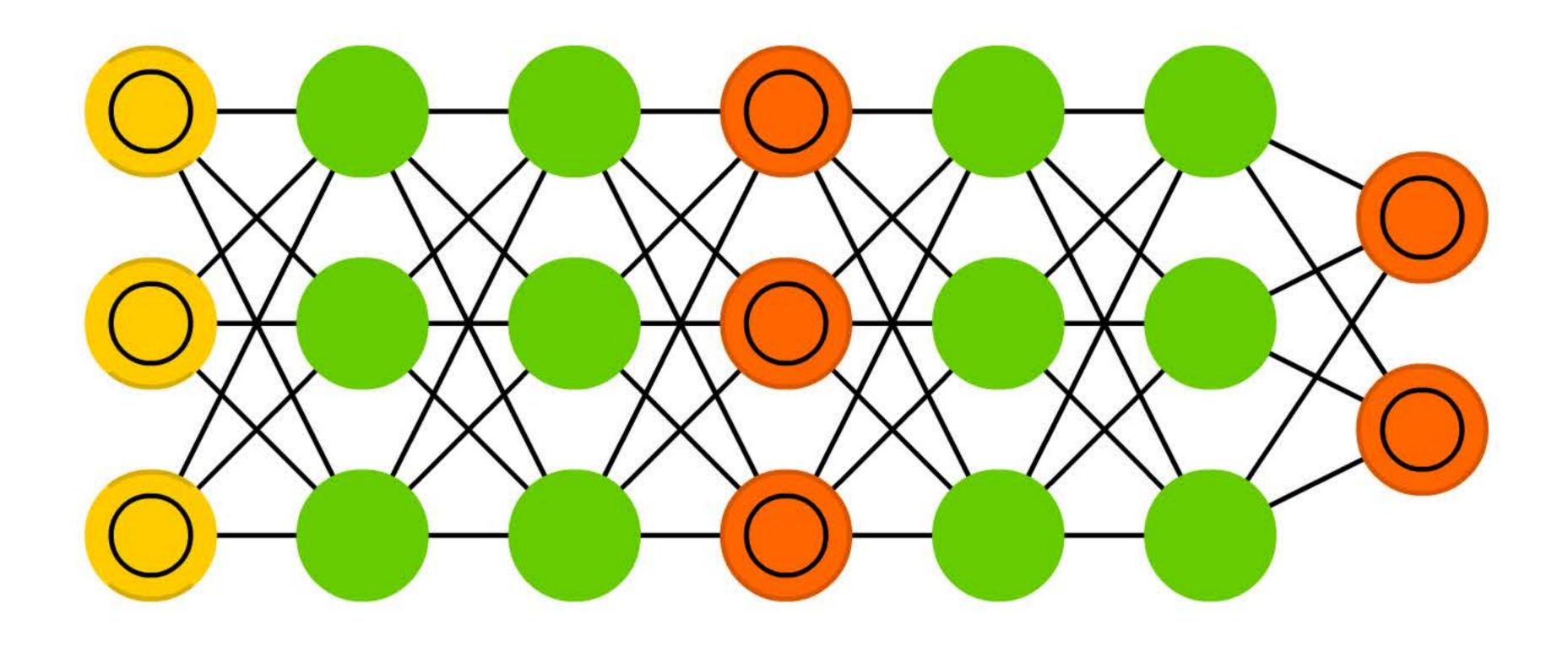




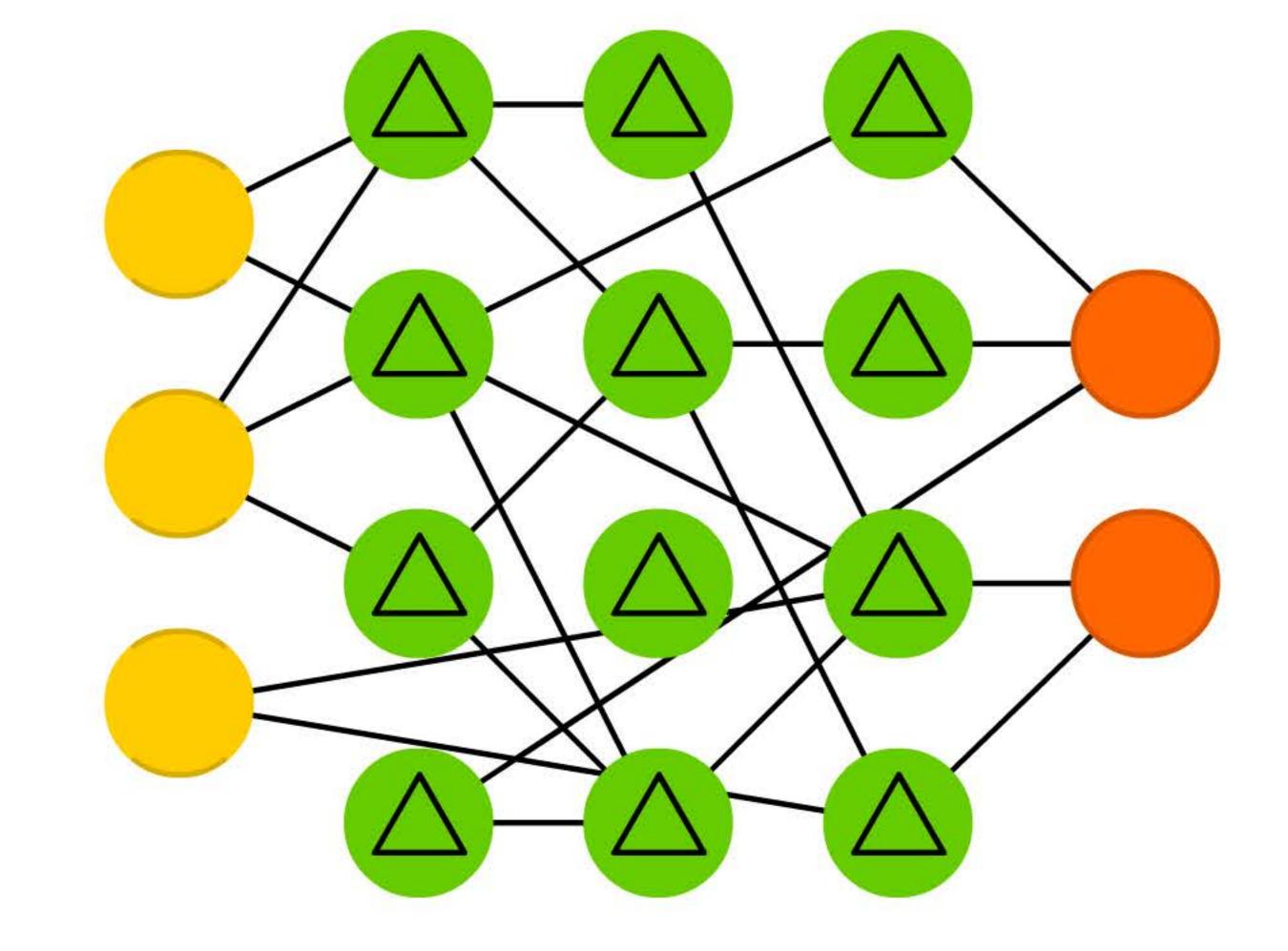


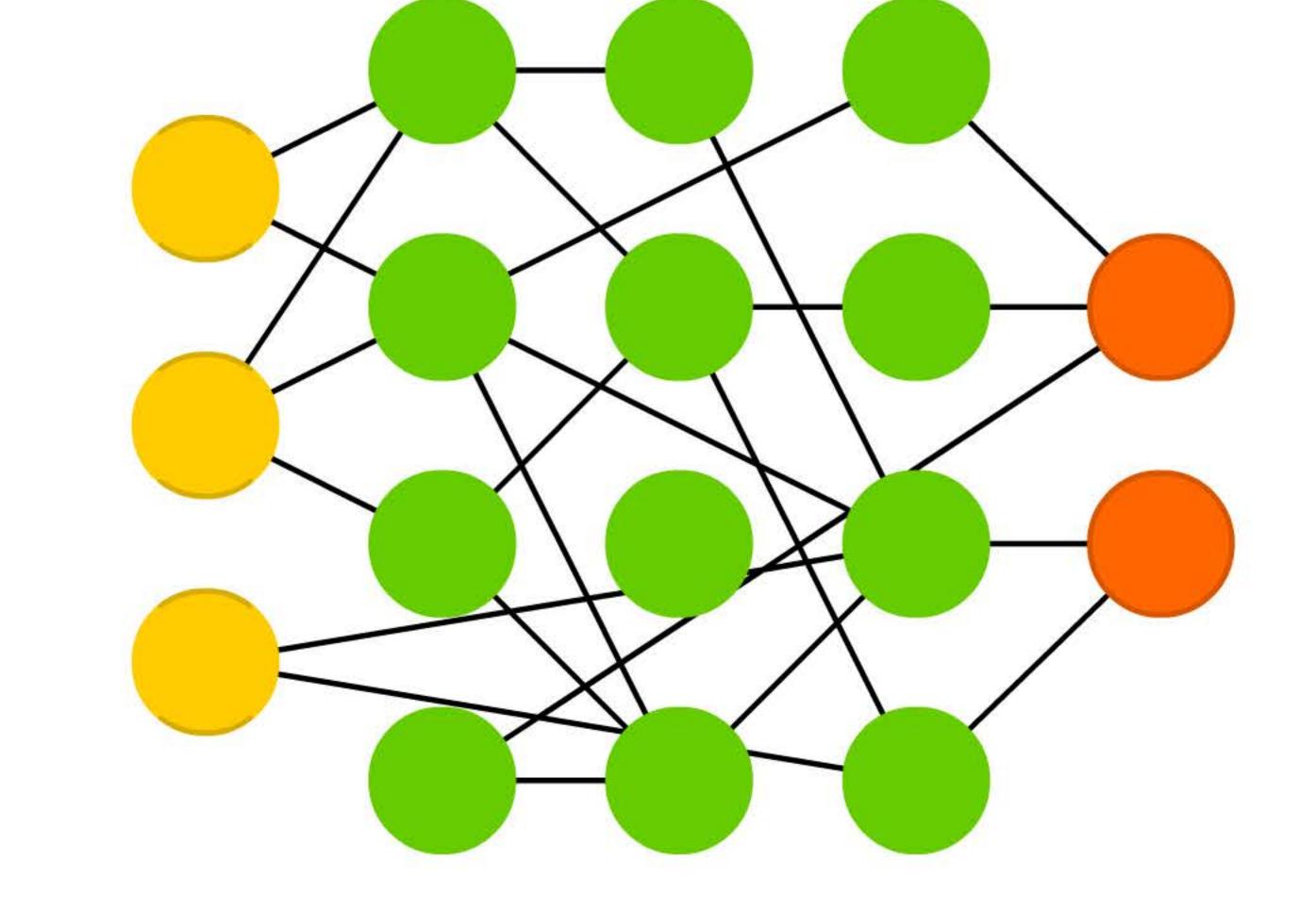


Generative Adversarial Network (GAN)

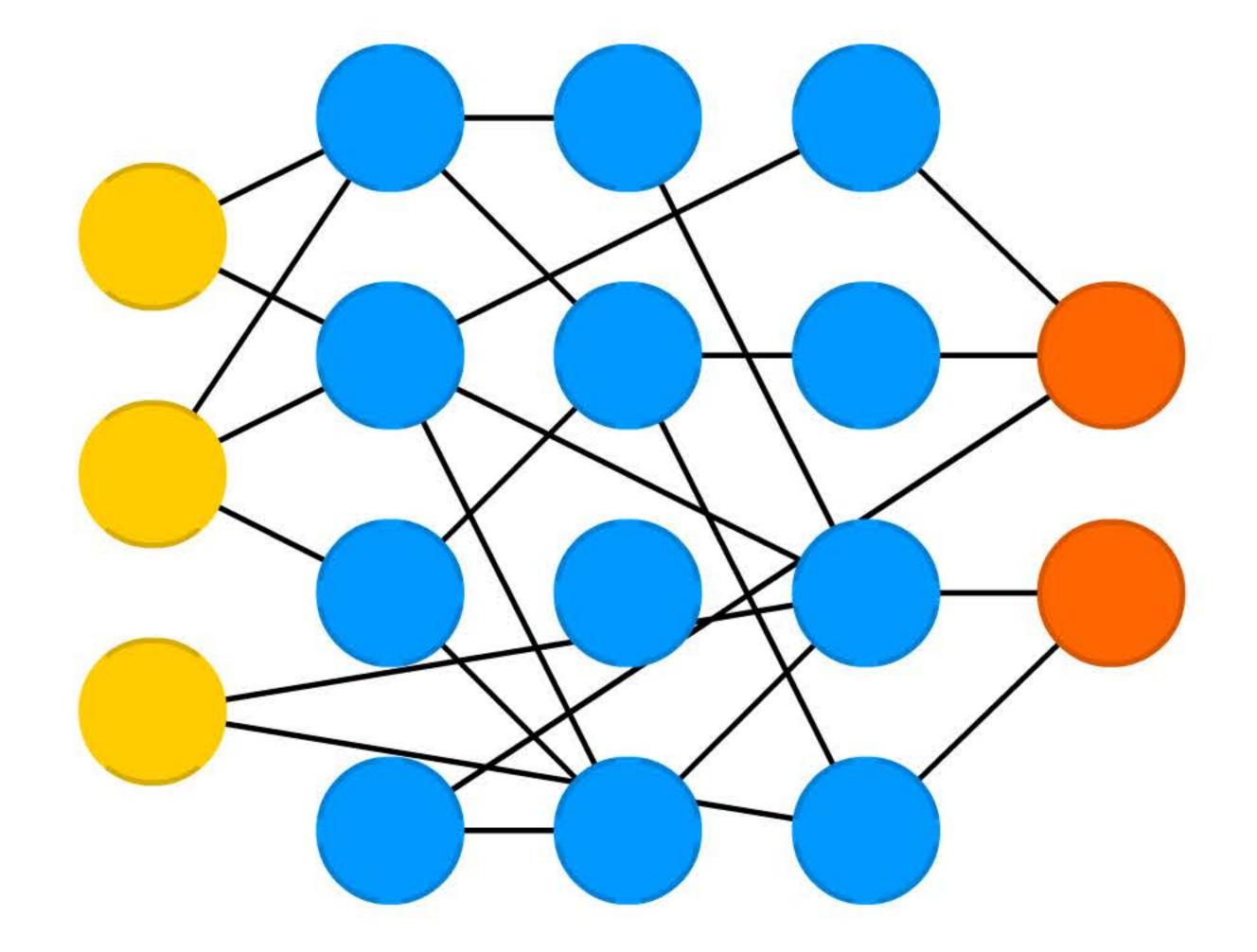


Liquid State Machine (LSM) Extreme Learning Machine (ELM)

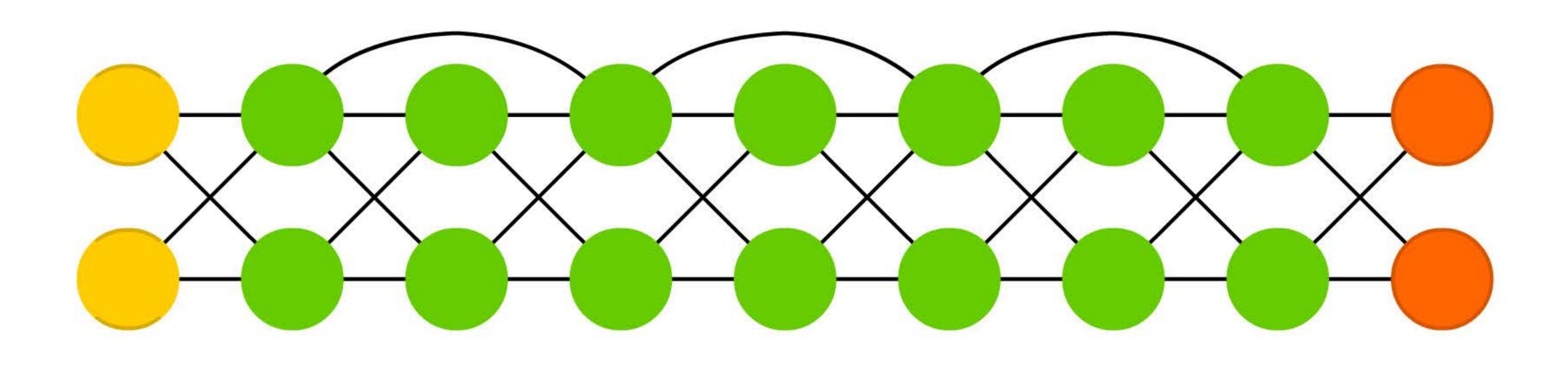


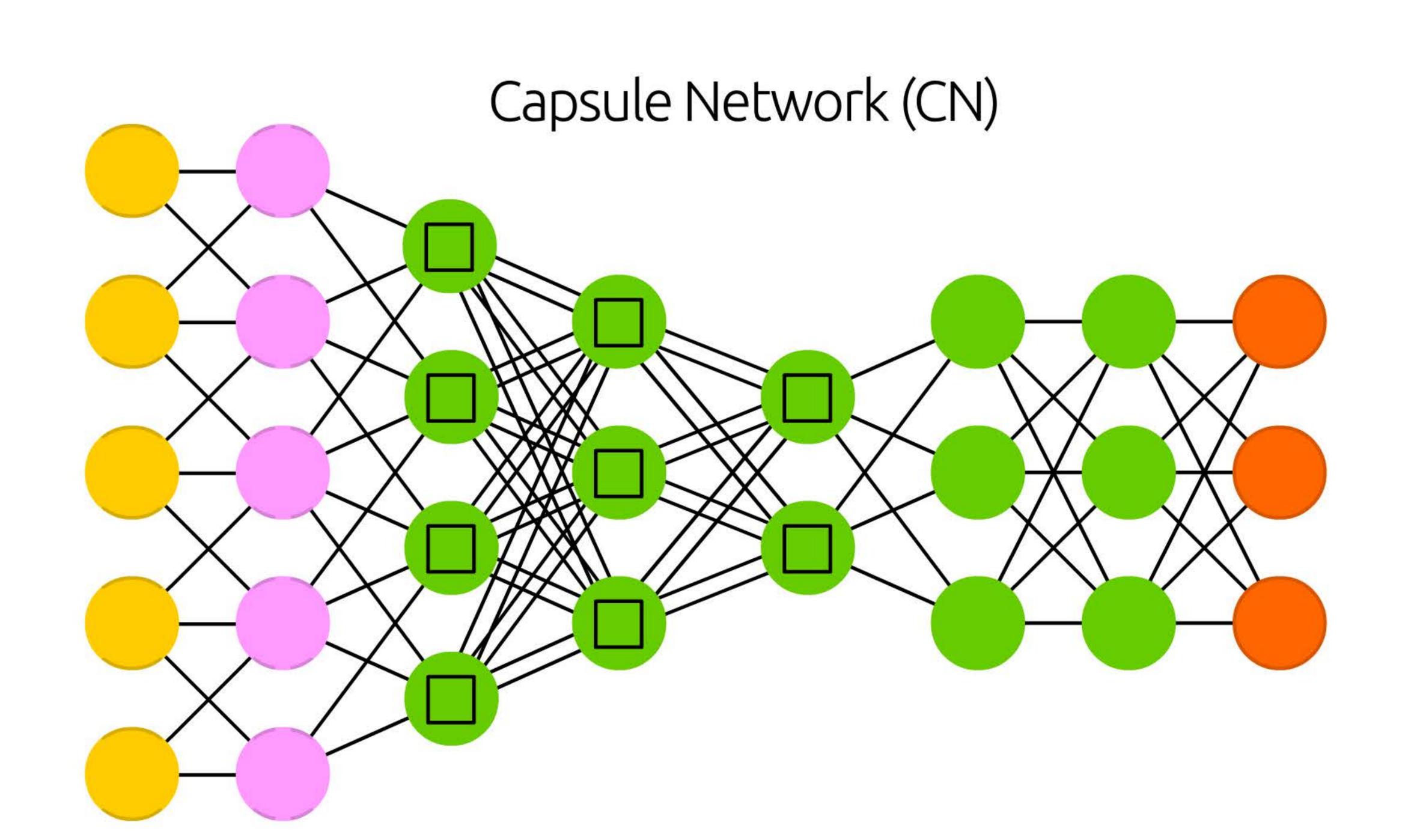


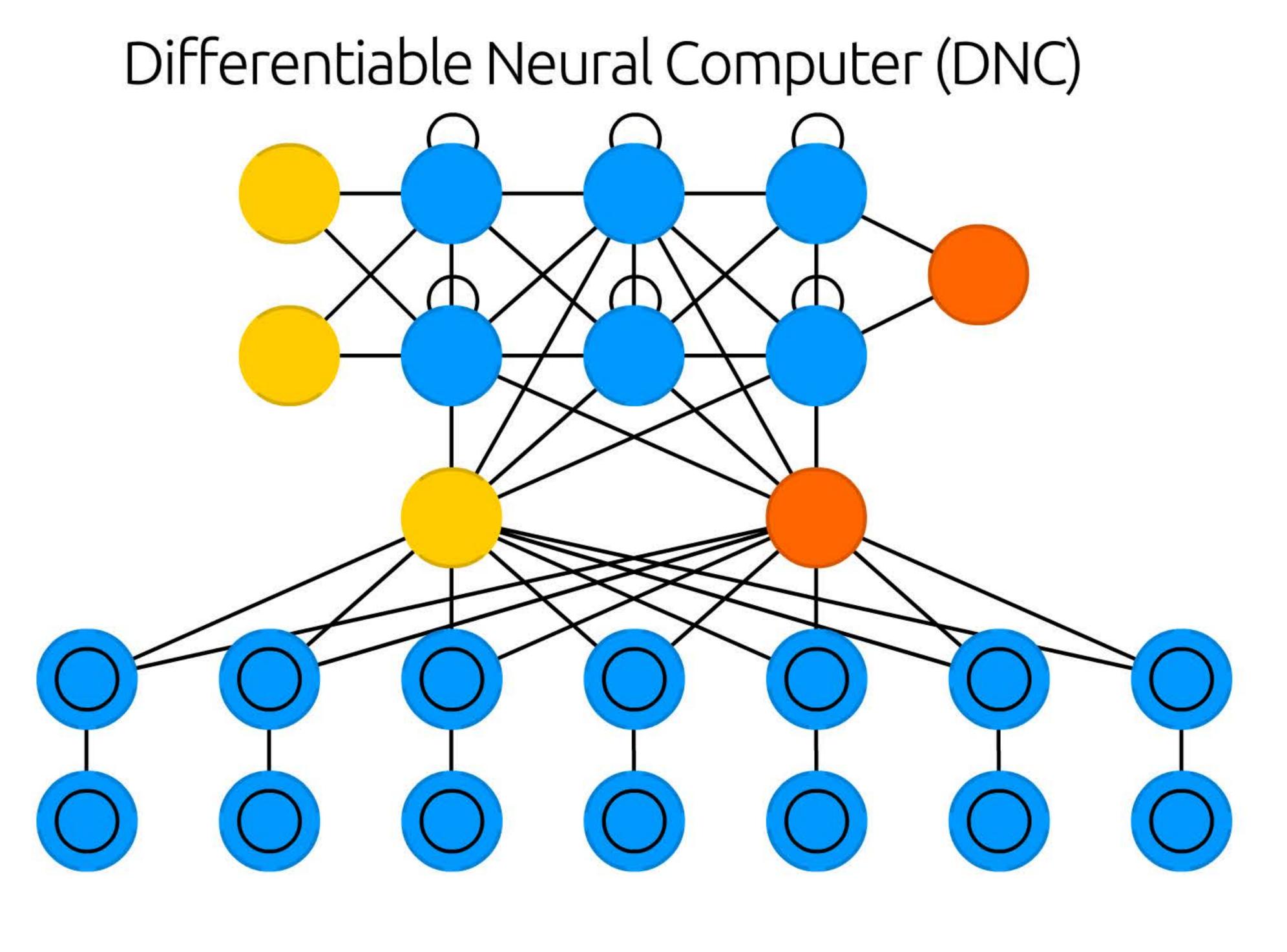
Echo State Network (ESN)



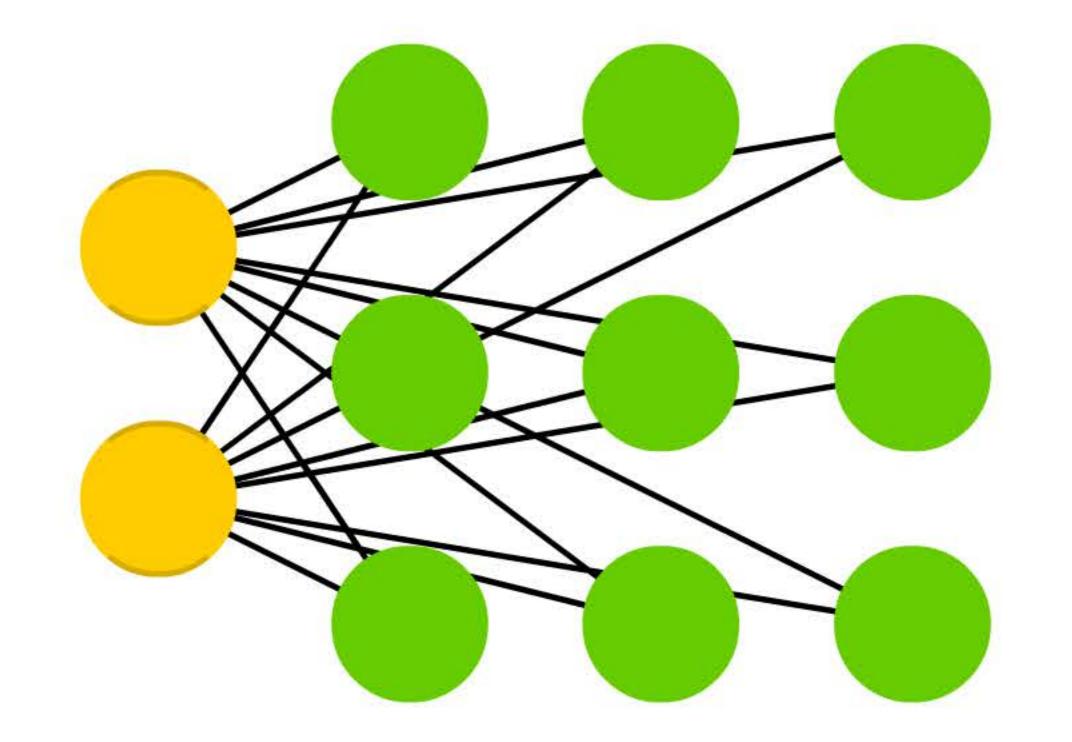




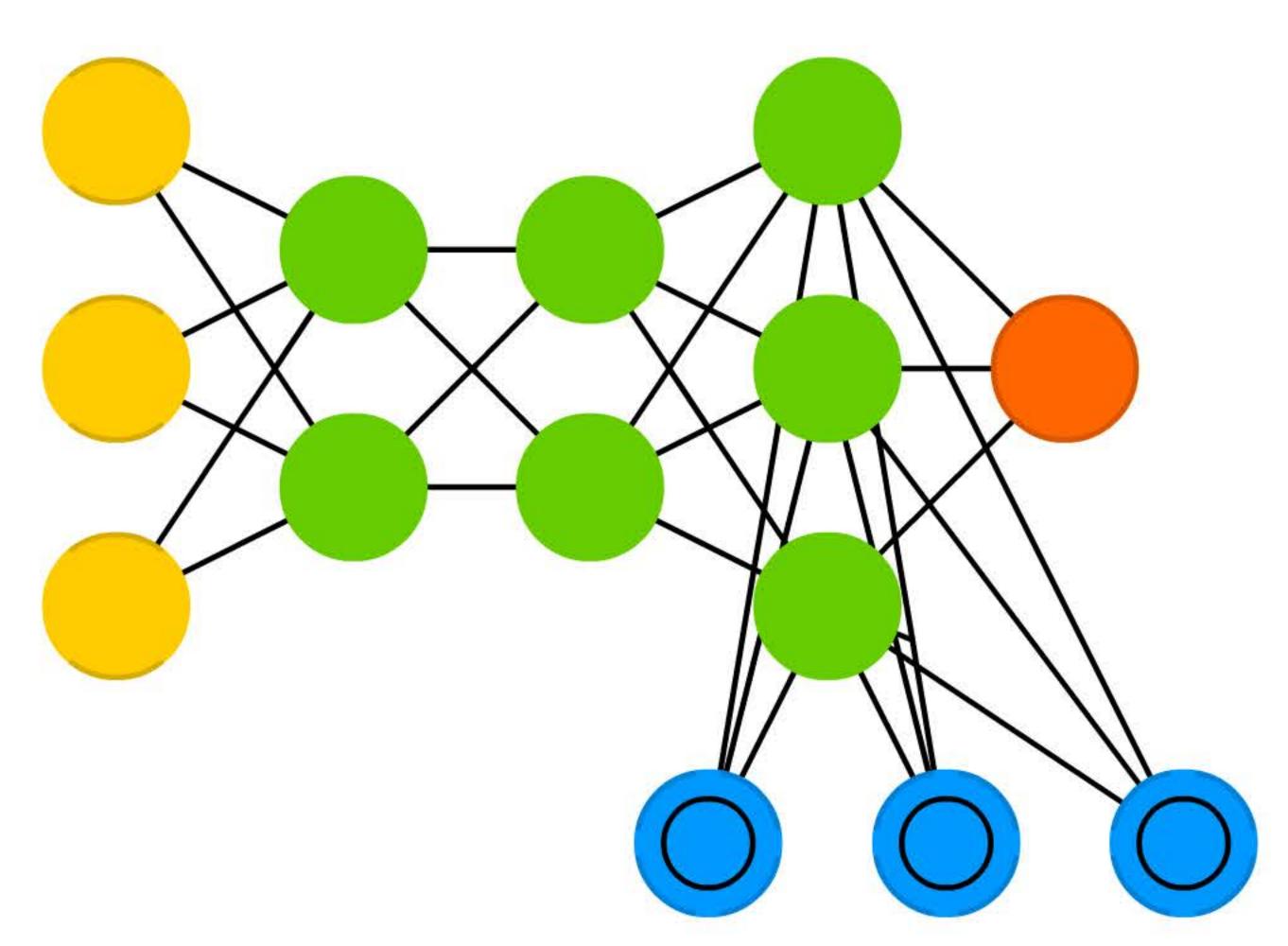




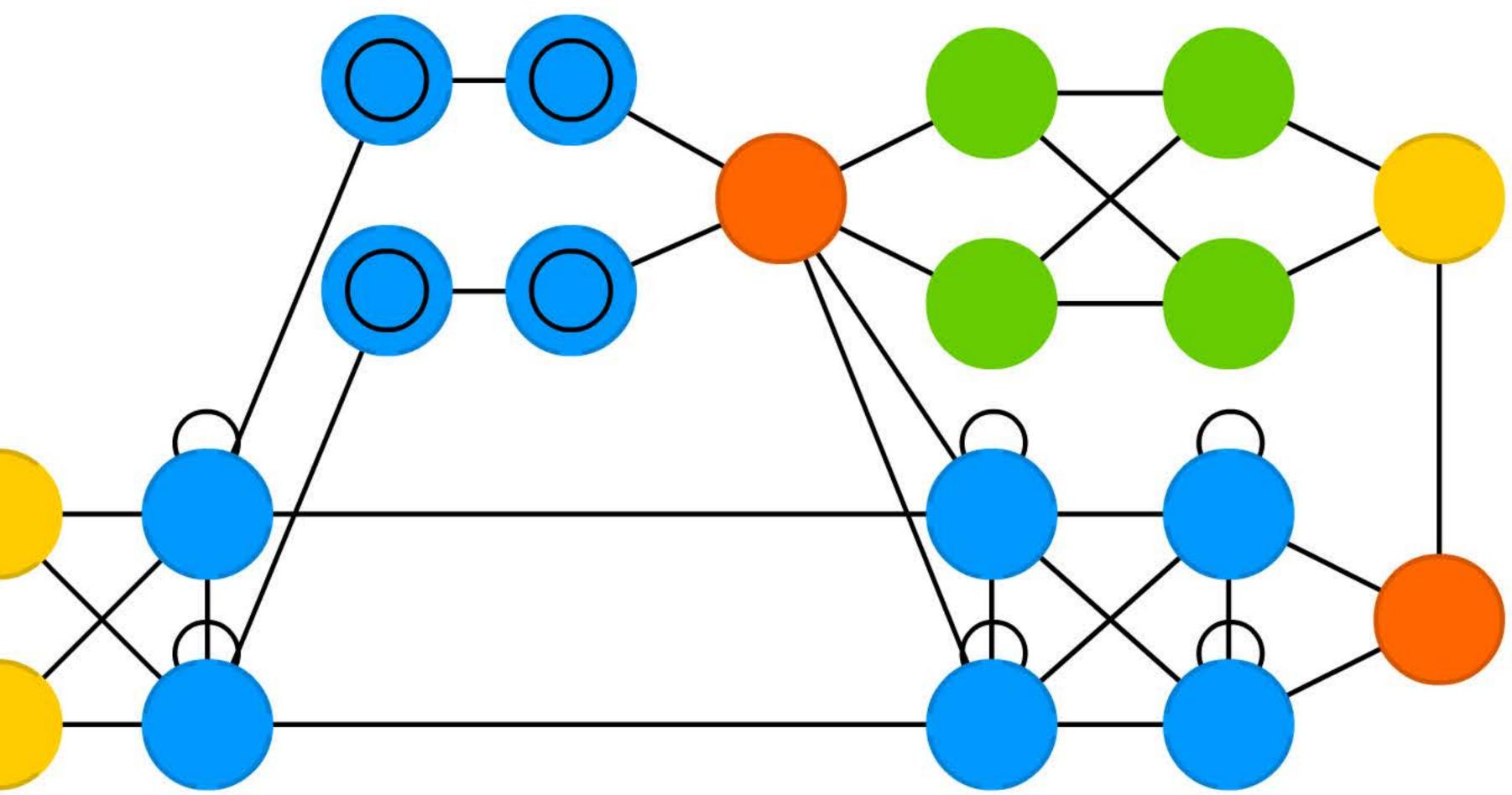
Kohonen Network (KN)



Neural Turing Machine (NTM)



Attention Network (AN)

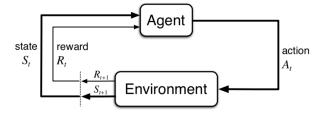




# 5. Reinforcement learning

## Reinforcement Learning Cheat Sheet Optimal

## Agent-Environment Interface



The Agent at each step t receives a representation of the environment's state,  $S_t \in S$  and it selects an action  $A_t \in A(s)$ . Then, as a consequence of its action the agent receives a reward,  $R_{t+1} \in R \in \mathbb{R}$ .

## Policy

A policy is a mapping from a state to an action

(1)

(2)

That is the probability of select an action  $A_t = a$  if  $S_t = s$ .

 $\pi_t(s|a)$ 

## Reward

The total *reward* is expressed as:

$$G_t = \sum_{k=0}^{H} \gamma^k r_{t+k+1}$$

Where  $\gamma$  is the *discount factor* and *H* is the *horizon*, that can be infinite.

#### Markov Decision Process

A **Markov Decision Process**, MPD, is a 5-tuple  $(S, A, P, R, \gamma)$  where:

finite set of states:  

$$s \in S$$
  
finite set of actions:  
 $a \in A$   
state transition probabilities:  
 $p(s'|s,a) = Pr\{S_{t+1} = s'|S_t = s, A_t = a\}$   
expected reward for state-action-nexstate:  
 $r(s',s,a) = \mathbb{E}[R_{t+1}|S_{t+1} = s', S_t = s, A_t = a]$ 
(3)

## Value Function

Value function describes how good is to be in a specific state s under a certain policy  $\pi$ . For MDP:

$$V_{\pi}(s) = \mathbb{E}[G_t | S_t = s] \tag{4}$$

Informally, is the expected return (expected cumulative discounted reward) when starting from s and following  $\pi$ 

$$V_*(s) = \max V_\pi(s)$$

## Action-Value (Q) Function

We can also denoted the expected reward for state, action pairs.

$$q_{\pi}(s,a) = \mathbb{E}_{\pi} \left[ G_t | S_t = s, A_t = a \right]$$
(6)

#### Optimal

The optimal value-action function:

$$q_*(s,a) = \max_{\pi} q^{\pi}(s,a) \tag{7}$$

Clearly, using this new notation we can redefine  $V^*$ , equation 5, using  $q^*(s, a)$ , equation 7:

$$V_*(s) = \max_{a \in A(s)} q_{\pi*}(s, a)$$

Intuitively, the above equation express the fact that the value of a state under the optimal policy **must be equal** to the expected return from the best action from that state.

## **Bellman Equation**

An important recursive property emerges for both Value (4) and Q (6) functions if we expand them.

## Value Function

$$V_{\pi}(s) = \mathbb{E}_{\pi} \left[ G_t | S_t = s \right]$$
  
=  $\mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right]$   
=  $\mathbb{E}_{\pi} \left[ R_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} | S_t = s \right]$   
=  $\sum_{a} \pi(a|s) \sum_{s'} \sum_{r} p(s', r|s, a)$   
Sum of all probabilities  $\forall$  possible  $r$   
 $\left[ r + \gamma \underbrace{\mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} | S_{t+1} = s' \right]}_{\text{Expected reward from } s_{t+1}} \right]$   
=  $\sum_{a} \pi(a|s) \sum_{s'} \sum_{r} p(s', r|s, a) \left[ r + \gamma V_{\pi}(s') \right]$ 

Similarly, we can do the same for the Q function:

$$(s,a) = \mathbb{E}_{\pi} \left[ G_t | S_t = s, A_t = a \right]$$

$$= \mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right]$$

$$= \mathbb{E}_{\pi} \left[ R_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} | S_t = s, A_t = a \right]$$

$$= \sum_{s',r} p(s',r|s,a) \left[ r + \gamma \mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} | S_{t+1} = s' \right] \right]$$

$$= \sum_{s',r} p(s',r|s,a) \left[ r + \gamma V_{\pi}(s') \right]$$
(10)

## **Dynamic Programming**

Taking advantages of the subproblem structure of the V and Q function we can find the optimal policy by just *planning* 

#### **Policy Iteration**

(5)

(8)

(9)

 $q_{\pi}$ 

We can now find the optimal policy 1. Initialisation  $V(s) \in \mathbb{R}$ , (e.g V(s) = 0) and  $\pi(s) \in A$  for all  $s \in S$ ,  $\Delta \leftarrow 0$ 2. Policy Evaluation while  $\Delta > \theta$  (a small positive number) do foreach  $s \in S$  do  $v \leftarrow V(s)$  $V(s) \leftarrow \sum_{a}^{\cdot} \pi(a|s) \sum_{s',r} p(s',r|s,a) \left[r + \gamma V(s')\right]$  $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ end end 3. Policy Improvement policy-stable  $\leftarrow true$ foreach  $s \in S$  do old-action  $\leftarrow \pi(s)$  $\pi(s) \leftarrow \operatorname*{argmax}_{a} \sum_{s',r} p(s',r|s,a) \left[r + \gamma V(s')\right]$ policy-stable  $\leftarrow$  old-action =  $\pi(s)$ end if *policy-stable* return  $V \approx V_*$  and  $\pi \approx \pi_*$ , else go to 2 Algorithm 1: Policy Iteration

#### Value Iteration

We can avoid to wait until V(s) has converged and instead do policy improvement and truncated policy evaluation step in one operation

Initialise  $V(s) \in \mathbb{R}, e.gV(s) = 0$   $\Delta \leftarrow 0$ while  $\Delta \ge \theta$  (a small positive number) do foreach  $s \in S$  do  $v \leftarrow V(s)$   $V(s) \leftarrow \max_{a} \sum_{s',r} p(s', r|s, a) [r + \gamma V(s')]$   $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ end end ouput: Deterministic policy  $\pi \approx \pi_*$  such that  $\pi(s) = \underset{a}{\operatorname{argmax}} \sum_{s',r} p(s', r|s, a) [r + \gamma V(s')]$ 

## Algorithm 2: Value Iteration

## Monte Carlo Methods

Monte Carlo (MC) is a *Model Free* method, It does not require complete knowledge of the environment. It is based on **averaging sample returns** for each state-action pair. The following algorithm gives the basic implementation

```
Initialise for all s \in S, a \in A(s):
  Q(s, a) \leftarrow \text{arbitrary}
  \pi(s) \leftarrow \text{arbitrary}
  Returns(s, a) \leftarrow empty list
while forever do
     Choose S_0 \in S and A_0 \in A(S_0), all pairs have
      probability > 0
     Generate an episode starting at S_0, A_0 following \pi
       for each pair s, a appearing in the episode do
         G \leftarrow return following the first occurrence of s, a
         Append G to Returns(s, a))
         Q(s, a) \leftarrow average(Returns(s, a))
     end
     foreach s in the episode do
         \pi(s) \leftarrow \operatorname{argmax} Q(s, a)
     end
end
```

#### Algorithm 3: Monte Carlo first-visit

For non-stationary problems, the Monte Carlo estimate for, e.g, V is:

$$V(S_t) \leftarrow V(S_t) + \alpha \left[ G_t - V(S_t) \right]$$
(11)

Where  $\alpha$  is the learning rate, how much we want to forget about past experiences.

## Sarsa

6

Sarsa (State-action-reward-state-action) is a on-policy TD control. The update rule:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t) \right]$$

#### *n*-step Sarsa

Define the n-step Q-Return

$$q^{(n)} = R_{t+1} + \gamma Rt + 2 + \ldots + \gamma^{n-1} R_{t+n} + \gamma^n Q(S_{t+n})$$

n-stepSarsa updateQ(S,a) towards the n-step Q-return

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ q_t^{(n)} - Q(s_t, a_t) \right]$$

Forward View  $Sarsa(\lambda)$ 

$$q_t^{\lambda} = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} q_t^{(n)}$$

Forward-view  $Sarsa(\lambda)$ :

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ q_t^{\lambda} - Q(s_t, a_t) \right]$$

 $\begin{array}{c|c} \mbox{Initialise } Q(s,a) \mbox{ arbitrarily and } \\ Q(terminal - state,) = 0 \\ \mbox{foreach } episode \in episodes \mbox{ do} \\ \hline \mbox{Choose } a \mbox{ from } s \mbox{ using policy derived from } Q \mbox{ (e.g., } \\ \epsilon \mbox{-greedy}) \\ \mbox{ while } s \mbox{ is not terminal } \mbox{ do} \\ \hline \mbox{ Take action } a, \mbox{ observer } r, s' \\ \mbox{ Choose } a' \mbox{ from } s' \mbox{ using policy derived from } Q \\ \mbox{ (e.g., } \epsilon \mbox{-greedy}) \\ \mbox{ Q(s, a) } \leftarrow Q(s, a) + \alpha \left[ r + \gamma Q(s', a') - Q(s, a) \right] \\ \mbox{ s } \leftarrow a' \\ \mbox{ end } \\ \hline \mbox{ end } \end{array}$ 

Algorithm 4:  $Sarsa(\lambda)$ 

## Temporal Difference - Q Learning

Temporal Difference (TD) methods learn directly from raw experience without a model of the environment's dynamics. TD substitutes the expected discounted reward  $G_t$  from the episode with an estimation:

$$V(S_t) \leftarrow V(S_t) + \alpha \left[ R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \right]$$
(12)

The following algorithm gives a generic implementation. Initialise Q(s, a) arbitrarily and Q(terminal - state, ) = 0foreach episode  $\in$  episodes do while s is not terminal do Choose a from s using policy derived from Q(e.g.,  $\epsilon$ -greedy) Take action a, observer r, s'  $Q(s, a) \leftarrow$   $Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a)\right]$   $s \leftarrow s'$ end end

......

## Algorithm 5: Q Learning

## Deep Q Learning

Created by DeepMind, Deep Q Learning, DQL, substitutes the Q function with a deep neural network called Q-network. It also keep track of some observation in a memory in order to use them to train the network.

$$L_{i}(\theta_{i}) = \mathbb{E}_{(s,a,r,s') \sim U(D)} \begin{bmatrix} (r + \gamma \max_{a} Q(s',a';\theta_{i-1}) - Q(s,a;\theta_{i}))^{2} \\ \underbrace{(r + \gamma \max_{a} Q(s',a';\theta_{i-1}) - Q(s,a;\theta_{i}))^{2}}_{\text{target}} \end{bmatrix}$$
(13)

Where  $\theta$  are the weights of the network and U(D) is the experience replay history.

```
Initialise replay memory D with capacity N
Initialise Q(s, a) arbitrarily
foreach episode \in episodes do
     while s is not terminal do
           With probability \epsilon select a random action
             a \in A(s)
           otherwise select a = \max_a Q(s, a; \theta)
           Take action a, observer r, s'
           Store transition (s, a, r, s') in D
           Sample random minibatch of transitions
            (s_i, a_i, r_i, s'_i) from D
           Set y_i \leftarrow
                                               for terminal s'_i
           \begin{cases} \bar{r_j} + \gamma \max_a Q(s', a'; \theta) & \text{for non-terminal } s'_j \\ \text{Perform gradient descent step on} \end{cases} 
           (y_j - Q(s_j, a_j; \Theta))^2
s \leftarrow s'
     end
end
```

#### Algorithm 6: Deep Q Learning

Copyright © 2018 Francesco Saverio Zuppichini https://github.com/FrancescoSaverioZuppichini/Reinforcement-Learning-Cheat-Sheet



# 6. Programing languages



6.1 Python language



# **Python 3 Beginner's Reference Cheat Sheet**

## Main data types

<b>boolean</b> = True / False
integer = 10
float = 10.01
string = "123abc"
list = [ value1, value2, ]
<pre>dictionary = { key1:value1, key2:value2,}</pre>

Numeric operators			Comparison operators		
+ - * / ** % //	<ul> <li>subtraction</li> <li>subtraction</li> <li>multiplication</li> <li>division</li> <li>exponent</li> <li>modulus</li> </ul>		== != > < >= <=	equal different higher lower higher or equal lower or equal	
Boolean operators					ecial acters
and or not	logical AND logical OR logical NOT		# \n \ <ch< th=""><th>iar&gt;</th><th>coment new line scape char</th></ch<>	iar>	coment new line scape char

## String operations

string[i]	retrieves character at position i
string[-1]	retrieves last character
string[i:j]	retrieves characters in range i to j

## List operations

list = []	defines an empty list		
list[i] = x	stores x with index i		
list[i]	retrieves the item with index I		
list[-1]	retrieves last item		
list[i:j]	retrieves items in the range i to j		
del list[i]	removes the item with index i		

## Dictionary operations

dict = {}	defines an empty dictionary
dict[k] = x	stores x associated to key k
dict[k]	retrieves the item with key k
del dict[k]	removes the item with key k

## String methods

string.upper()	cor
string.lower()	cor
string.count(x)	cou
	tim
string.find(x)	pos
	000
string.replace(x,y)	rep
string.strip(x)	ret
	del
string.join(L)	ret
	val
string.format(x)	ret
	inc

converts to uppercase
converts to lowercase
counts how many
times x appears
position of the x first
occurrence
replaces x for y
returns a list of values
delimited by x
returns a string with L
values joined by string
returns a string that
includes formatted x

## List methods

list.append(x)	adds x to the end of the list
list.extend(L)	appends L to the end of the list
list.insert(i,x)	inserts x at i position
list.remove(x)	removes the first list item whose value is x
list.pop(i)	removes the item at position i and returns its value
list.clear()	removes all items from the list
list.index(x)	returns a list of values delimited by x
list.count(x)	returns a string with list values joined by S
list.sort()	sorts list items
list.reverse()	reverses list elements
list.copy()	returns a copy of the list

## **Dictionary methods**

dict.keys()	returns a list of keys
dict.values()	returns a list of values
dict.items()	returns a list of pairs (key,value)
dict.get(k)	returns the value associtated to
	the key k
dict.pop()	removes the item associated to
	the key and returns its value
dict.update(D)	adds keys-values (D) to dictionary
dict.clear()	removes all keys-values from the
	dictionary
dict.copy()	returns a copy of the dictionary

Legend: x,y stand for any kind of data values, s for a string, n for a number, L for a list where i,j are list indexes, D stands for a dictionary and k is a dictionary key.

# **Python 3 Beginner's Reference Cheat Sheet**

E	Built-in functions	Conditional statements	Loops	Functions
print(x, sep='y') input(s)	prints x objects separated by y prints s and waits for an input	if <condition> : <code></code></condition>	<b>while</b> <condition>: <code></code></condition>	<b>def</b> function( <params>): <code></code></params>
	that will be returned	else if <condition> : <code></code></condition>	for <variable> in <list>:</list></variable>	<b>return</b> <data></data>
len(x)	returns the length of x (s, L or D)	 else:	<code></code>	Modules
min(L)	returns the minimum value in L	<code></code>	for <variable> in</variable>	
max(L) sum(L)	returns the maximum value in L returns the sum of the values in L	if <value> in <list>:</list></value>	<b>range</b> (start,stop,step): <code></code>	<pre>import module module.function()</pre>
range(n1,n2,n)	returns a sequence of numbers from n1 to n2 in steps of n	Data validation	for key, value in dict.items(): <code></code>	<b>from</b> module <b>import</b> * function()
abs(n)	returns the absolute value of n	try:		Reading and
round(n1,n)	returns the n1 number rounded to n digits	<code> <b>except</b> <error>: <code></code></error></code>	Loop control statements	writing files
type(x)	returns the type of x (string, float, list, dict)	else: <code></code>	break finishes loop execution	f = open( <path>,'r') f.read(<size>) f.readline(<size>)</size></size></path>
str(x)	converts x to string		<b>continue</b> jumps to next iteration	f.close()
list(x)	converts x to a list	Working with files and folders	pass does nothing	f = open( <path>,'r')</path>
int(x)	converts x to a integer number	import os	Running external	for line in f: <code></code>
float(x)	converts x to a float number	os.getcwd()	programs	f.close()
help(s)	prints help about x	os.makedirs( <path>) os.chdir(<path>)</path></path>	import os	f = open( <path>,'w') f.write(<str>)</str></path>
map(function, L)	Applies function to values in L	os.listdir( <path>)</path>	os.system( <command/> )	f.close()

Legend: x,y stand for any kind of data values, s for a string, n for a number, L for a list where i,j are list indexes, D stands for a dictionary and k is a dictionary key.

**Python Basics** 

Learn More Python for Data Science Interactively at www.datacamp.com



## Variables and Data Types

## Variable Assignment

>>>	x=5

>>> x

5

## **Calculations With Variables**

>>>	x+2	Sum of two variables
7		
>>>	x-2	Subtraction of two variables
3		
>>>	x*2	Multiplication of two variables
10		
>>>	x**2	Exponentiation of a variable
25		
>>>	x%2	Remainder of a variable
1		
	x/float(2)	Division of a variable
///	x/110al(2)	DIVISION OF a VANADIE
2.5		

## Types and Type Conversion

str()	'5', '3.45', 'True'	Variables to strings
int()	5, 3, 1	Variables to integers
float()	5.0, 1.0	Variables to floats
bool()	True, True, True	Variables to booleans

## **Asking For Help**

>>> help(str)

## Strings

>>> my string = 'thisStringIsAwesome' >>> my string 'thisStringIsAwesome'

## **String Operations**

>>> my string \* 2 'thisStringIsAwesomethisStringIsAwesome' >>> my string + 'Innit' 'thisStringIsAwesomeInnit' >>> 'm' in my string

#### Lists Also see NumPy A >>> a = 'is' >>> b = 'nice'

>>> my list = ['my', 'list', a, b] >>> my list2 = [[4,5,6,7], [3,4,5,6]]

Selecting List Elements	Index starts at o
<pre>Subset &gt;&gt;&gt; my_list[1]</pre>	Select item at index 1
>>> my_list[-3] Slice	Select 3rd last item
>>> my_list[1:3] >>> my list[1:]	Select items at index 1 and 2 Select items after index 0
>>> my_list[:3] >>> my_list[:]	Select items before index 3 Copy my list
Subset Lists of Lists	
>>> my_list2[1][0] >>> my_list2[1][:2]	my_list[list][itemOfList]
List Operations	

#### List Operations

>>> my list + my list ['my', 'list', 'is', 'nice', 'my', 'list', 'is', 'nice'] >>> my list \* 2 ['my', 'list', 'is', 'nice', 'my', 'list', 'is', 'nice'] >>> my list2 > 4

## List Methods

>>> my list.index(a) >>> my list.count(a) >>> my list.append('!') >>> my list.remove('!') >>> del(my\_list[0:1]) >>> my list.reverse() >>> my list.extend('!') >>> my list.pop(-1) >>> my list.insert(0,'!') >>> my list.sort()

String Operations	Index starts at
>>> my_string[3] >>> my_string[4:9]	
String Methods	
>>> my_string.upper()	String to uppercase
—	String to uppercase String to lowercase
>>> my_string.lower()	5
<pre>&gt;&gt;&gt; my_string.upper() &gt;&gt;&gt; my_string.lower() &gt;&gt;&gt; my_string.count('w') &gt;&gt;&gt; my_string.replace('e', 'i')</pre>	String to lowercase

Also see NumPy Arrays	Libraries	
', a, b] [3,4,5,6]]	<pre>Import libraries &gt;&gt;&gt; import numpy &gt;&gt;&gt; import numpy as np Selective import &gt;&gt;&gt; from math import pi</pre>	pandas       Image: Second Secon
Index starts at O	Install Dython	
ct item at index 1 ct 3rd last item ct items at index 1 and 2 ct items after index 0 ct items before index 3 y my_list	5 1	DE that is included with Anaconda
list[list][itemOfList]		Also see Lists
	<pre>Numpy Arrays &gt;&gt;&gt; my_list = [1, 2, 3, 4] &gt;&gt;&gt; my_array = np.array(my_ &gt;&gt;&gt; my 2darray = np.array([</pre>	_list)
:', 'is', 'nice']	Selecting Numpy Array Elem	ents Index starts at o
', 'is', 'nice']	Subset >>> my_array[1]	Select item at index 1
	<pre>Slice &gt;&gt;&gt; my_array[0:2]</pre>	Select items at index 0 and 1
Get the index of an item Count an item Append an item at a time Remove an item	<pre>array([1, 2]) Subset 2D Numpy arrays &gt;&gt;&gt; my_2darray[:,0] array([1, 4])</pre>	my_2darray[rows, columns]
Remove an item	Numpy Array Operations	
Reverse the list Append an item Remove an item Insert an item Sort the list	<pre>&gt;&gt;&gt; my_array &gt; 3 array([False, False, False, &gt;&gt;&gt; my_array * 2 array([2, 4, 6, 8]) &gt;&gt;&gt; my_array + np.array( array([6, 8, 10, 12])</pre>	
	Numpy Array Functions	
Index starts at o	<pre>&gt;&gt;&gt; my_array.shape &gt;&gt;&gt; np.append(other_arra &gt;&gt;&gt; np.insert(my_array, &gt;&gt;&gt; np.delete(my_array,[ &gt;&gt;&gt; np.mean(my_array)</pre>	1, 5) Insert items in an array

>>> np.median(my array)

>>> my array.corrcoef()

>>> np.std(my array)

DataCamp Learn Python for Data Science Interactively

Median of the array

Standard deviation

Correlation coefficient

0

**Importing Data** 

Learn Python for data science Interactively at www.DataCamp.com



## **Importing Data in Python**

Most of the time, you'll use either **NumPy** or **pandas** to import your data:

>>> import numpy as np >>> import pandas as pd

## Help

>>> np.info(np.ndarray.dtype)
>>> help(pd.read\_csv)

## Text Files

## **Plain Text Files**

>>>	filename = 'huck finn.txt'	
>>>	<pre>file = open(filename, mode='r')</pre>	Open the file for
>>>	<pre>text = file.read()</pre>	Read a file's cont
>>>	print(file.closed)	Check whether fi
>>>	file.close()	Close file
>>>	print(text)	

```
pen the file for reading
ead a file's contents
heck whether file is closed
lose file
```

#### Using the context manager with

ı open('huck	_finn.txt',	'r')	as	file:	
print(file.r	eadline())				Read a single line
print(file.r	eadline())				
print(file.r	eadline())				
	print(file.r print(file.r	<pre>h open('huck_finn.txt', print(file.readline()) print(file.readline()) print(file.readline())</pre>	<pre>print(file.readline()) print(file.readline())</pre>	<pre>print(file.readline()) print(file.readline())</pre>	print(file.readline())

## Table Data: Flat Files

#### (Importing Flat Files with numpy)

#### Files with one data type

skipro	ter=',',String used to separate valuesbws=2,Skip the first 2 lines.s=[0,2],Read the 1st and 3rd column
Files with mixed data types	

>>> filename = 'titanic.csv'		
>>> data = np.genfromtxt(filename,		
delimiter=',',		
names=True,	Look for column header	
dtype=None)		
	J	

>>> data\_array = np.recfromcsv(filename)

The default dtype of the <code>np.recfromcsv()</code> function is <code>None</code>.

#### Importing Flat Files with pandas

## Excel Spreadsheets

parse\_cols=[0],
 skiprows=[0],
 names=['Country'])

#### To access the sheet names, use the sheet names attribute:

>>> data.sheet\_names

## SAS Files

## **Stata Files**

>>> data = pd.read\_stata('urbanpop.dta')

## **Relational Databases**

>>> from sqlalchemy import create\_engine
>>> engine = create\_engine('sqlite://Northwind.sqlite')

#### Use the table names () method to fetch a list of table names:

>>> table\_names = engine.table\_names()

#### Querying Relational Databases

>>> con = engine.connect()
>>> rs = con.execute("SELECT \* FROM Orders")
>>> df = pd.DataFrame(rs.fetchall())
>>> df.columns = rs.keys()
>>> con.close()

#### Using the context manager with

>>> with engine.connect() as con: rs = con.execute("SELECT OrderID FROM Orders") df = pd.DataFrame(rs.fetchmany(size=5)) df.columns = rs.keys()

#### Querying relational databases with pandas

>>> df = pd.read\_sql\_query("SELECT \* FROM Orders", engine)

Exploring Your Data	
NumPy Arrays	
>>> data_array.dtype >>> data_array.shape >>> len(data_array)	Data type of array elements Array dimensions Length of array
pandas DataFrames	
<pre>&gt;&gt;&gt; df.head() &gt;&gt;&gt; df.tail() &gt;&gt;&gt; df.index &gt;&gt;&gt; df.columns &gt;&gt;&gt; df.info() &gt;&gt;&gt; dat_array = data.values</pre>	Return first DataFrame rows Return last DataFrame rows Describe index Describe DataFrame columns Info on DataFrame Convert a DataFrame to an a NumPy array

## **Pickled Files**

## **HDF5** Files

>>> import h5py
>>> filename = 'H-H1\_LOSC\_4\_v1-815411200-4096.hdf5'
>>> data = h5py.File(filename, 'r')

## **Matlab Files**

>>> import scipy.io
>>> filename = 'workspace.mat'

>>> mat = scipy.io.loadmat(filename)

## **Exploring Dictionaries**

## Accessing Elements with Functions

<pre>&gt;&gt;&gt; print(mat.keys()) &gt;&gt;&gt; for key in data.keys():</pre>	
print(key)	
meta	
quality	
strain	
>>> pickled_data.values()	
>>> print(mat.items())	

Return dictionary values Returns items in list format of (key, value)

0

Print dictionary keys Print dictionary keys

tuple pairs

## Accessing Data Items with Keys

>>> for key in data ['meta'].keys() print(key)	Explore the HDF5 structure
Description	
DescriptionURL	
Detector	
Duration	
GPSstart	
Observatory	
Туре	
UTCstart	
>>> print(data['meta']['Description'].value)	Retrieve the value for a key

## Navigating Your FileSystem

Magic Commands	
!ls %cd %pwd	List directory contents of files and directories Change current working directory Return the current working directory path
os Library	

	import os	
>>>	path = "/usr/tmp"	
>>>	wd = os.getcwd()	Store the name of current directory in a strin
>>>	os.listdir(wd)	Output contents of the directory in a list
>>>	os.chdir(path)	Change current working directory
>>>	os.rename("test1.txt",	Rename a file
	"test2.txt")	
>>>	<pre>os.remove("test1.txt")</pre>	Delete an existing file
>>>	os.mkdir("newdir")	Create a new directory

DataCamp

Learn R for Data Science Interactively

Pandas Basics

Learn Python for Data Science Interactively at www.DataCamp.com

# $\bigcirc$

## Pandas

The **Pandas** library is built on NumPy and provides easy-to-use data structures and data analysis tools for the Python programming language. pandas

Use the following import convention:

>>> import pandas as pd

>>> xlsx = pd.ExcelFile('file.xls') >>> df = pd.read excel(xlsx, 'Sheet1')

Pandas Data Structures	'Belgium'	column
	>>> df.iat([0],[0]) 'Belgium'	
Series A one-dimensional labeled array a 3	By Label	
	<pre>&gt;&gt;&gt; df.loc[[0], ['Country']]     'Belgium'</pre>	Select single value by row & column labels
Index c 7 d 4	<pre>&gt;&gt;&gt; df.at([0], ['Country'])     'Belgium'</pre>	
>>> s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])	By Label/Position	Select single row of
DataFrame	Country Brazil Capital Brasilia Population 207847528	subset of rows
	<pre>&gt;&gt;&gt; df.ix[:,'Capital'] 0 Brussels</pre>	Select a single column of subset of columns
0     Belgium     Brussels     11190846     data structure with columns of potentially different types	1 New Delhi 2 Brasília	
Index 1 India New Deini 1303171035 2 Brazil Brasília 207847528	<pre>&gt;&gt;&gt; df.ix[1,'Capital']     'New Delhi'</pre>	Select rows and columns
<pre>&gt;&gt;&gt; data = {'Country': ['Belgium', 'India', 'Brazil'],</pre>	Boolean Indexing	Series s where value is not >1
	>>> $s[(s < -1)   (s > 2)]$	$_{\rm S}$ where value is <-1 or >2
<pre>'Population': [11190846, 1303171035, 207847528]} &gt;&gt;&gt; df = pd.DataFrame(data,</pre>	<pre>&gt;&gt;&gt; df[df['Population']&gt;120000000] Setting</pre>	Use filter to adjust DataFrame
<pre>columns=['Country', 'Capital', 'Population'])</pre>	>>> s['a'] = 6	Set index a of Series s to 6

**Asking For Help** 

Selection

Getting

-5

>>> s['b']

>>> df[1:]

Country

**By Position** 

>>> df.iloc[[0],[0]]

>>> help(pd.Series.loc)

Capital Population

1 India New Delhi 1303171035 2 Brazil Brasília 207847528

Selecting, Boolean Indexing & Setting

Read and Write to CSV	Read and Write to SQL Query or Database Table	
<pre>&gt;&gt;&gt; pd.read_csv('file.csv', header=None, nrows=5) &gt;&gt;&gt; df.to_csv('myDataFrame.csv')</pre>	<pre>&gt;&gt;&gt; from sqlalchemy import create_engine &gt;&gt;&gt; engine = create_engine('sqlite:///:memory:')</pre>	
Read and Write to Excel	<pre>&gt;&gt;&gt; pd.read_sql("SELECT * FROM my_table;", engine) &gt;&gt;&gt; pd.read_sql_table('my_table', engine)</pre>	
<pre>&gt;&gt;&gt; pd.read_excel('file.xlsx')</pre>	<pre>&gt;&gt;&gt; pd.read_sql_query("SELECT * FROM my_table;", engine)</pre>	
<pre>&gt;&gt;&gt; pd.to_excel('dir/myDataFrame.xlsx', sheet_name='Sheet1') Read multiple sheets from the same file</pre>	<pre>read_sql() is a convenience wrapper around read_sql_table() and</pre>	

L table() and read sql query()

>>> pd.to sql('myDf', engine)

		Dropping		
	Also see NumPy Arrays	<pre>&gt;&gt;&gt; s.drop(['a', 'c']) &gt;&gt;&gt; df.drop('Country', axis=1) Drop values from columns(axis=1)</pre>		
Get o	ne element	Sort & Rank		
	ubset of a DataFrame	<pre>&gt;&gt;&gt; df.sort_index() &gt;&gt;&gt; df.sort_values(by='Country') &gt;&gt;&gt; df.rank()</pre> Sort by labels along an axis Assign ranks to entries		
		Retrieving Series/DataFrame Information		
ting		Basic Information		
	elect single value by row & olumn	>>> df.shape(rows,columns)>>> df.indexDescribe index>>> df.columnsDescribe DataFrame columns>>> df.info()Info on DataFrame>>> df.count()Number of non-NA values		
		Summary		
	elect single value by row & olumn labels	<pre>&gt;&gt;&gt; df.sum() &gt;&gt;&gt; df.cumsum() &gt;&gt;&gt; df.min()/df.max() &gt;&gt;&gt; df.idxmin()/df.idxmax() &gt;&gt;&gt; df.idxmin()/df.idxmax() &gt;&gt;&gt; df.describe() &gt;&gt;&gt; df.mean() </pre> Sum of values Summary statistics Adf.median() Median of values		
	elect single row of ubset of rows	Applying Functions		
	elect a single column of ubset of columns	>>> f = lambda x: x*2         >>> df.apply(f)         >>> df.applymap(f)    Apply function Apply function element-wise		
		Data Alignment		
Se	elect rows and columns	Internal Data Alignment		
00] U	eries s where value is not >1 where value is <-1 or >2 Ise filter to adjust DataFrame et index a of Series s to 6	<pre>NA values are introduced in the indices that don't overlap: &gt;&gt;&gt; s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd']) &gt;&gt;&gt; s + s3     a         10.0     b     NaN     c         5.0     d         7.0</pre>		
		Arithmetic Operations with Fill Methods		
Datab	oase Table	You can also do the internal data alignment yourself with the help of the fill methods:		
	ngine ///:memory:') able;", engine)	<pre>&gt;&gt;&gt; s.add(s3, fill_value=0) a 10.0 b -5.0 </pre>		

7.0 >>> s.sub(s3, fill value=2) >>> s.div(s3, fill value=4)

> DataCamp Learn Python for Data Science Interactively

>>> s.mul(s3, fill\_value=3)

5.0

С d

9

## **Pandas**

Learn Python for Data Science Interactively at <a href="http://www.DataCamp.com">www.DataCamp.com</a>



#### **Reshaping Data** Pivot >>> df3= df2.pivot(index='Date', Spread rows into columns columns='Type', values='Value') Date Type Value 11.432 2016-03-01 Type b с 13.031 2016-03-02 h Date 2016-03-01 20.784 с 2016-03-01 11.432 NaN 20.784 2016-03-03 99,906 а 1.303 13.031 NaN 2016-03-02 2016-03-02 1.303 а 2016-03-03 99,906 NaN 20.784 2016-03-03 20.784 **Pivot Table** >>> df4 = pd.pivot table(df2, Spread rows into columns values='Value'. index='Date', columns='Type'] Stack / Unstack >>> stacked = df5.stack() Pivot a level of column labels >>> stacked.unstack() Pivot a level of index labels 1 5 0 0.233482 0 1 5 0.233482 1 0.390959 0.390959 2 4 0.184713 2 4 0 0.237102 0.184713 3 3 0.433522 0.429401 0.237102 Unstacked 3 3 0 0.433522 0.429401 1 Stacked Melt Gather columns into rows >>> pd.melt(df2, id vars=["Date"], value\_vars=["Type", "Value"], value name="Observations") Date Variable Observations Date Type Value 0 2016-03-01 Туре а 2016-03-01 11.432 а 1 2016-03-02 Туре h 2016-03-02 13.031 2 2016-03-01 Туре с 3 2016-03-03 Туре 2016-03-01 20.784 а 4 2016-03-02 Туре а 2016-03-03 99,906 2016-03-03 5 Туре С 2016-03-02 1.303 6 2016-03-01 Value 11.432 20.784 2016-03-03 7 2016-03-02 Value 13.031 8 Value 20,784 2016-03-01 Value 9 2016-03-03 99.906 10 2016-03-02 Value 1.303 11 2016-03-03 Value 20.784 Iteration (Column-index, Series) pairs >>> df.iteritems()

(Row-index, Series) pairs

>>> df.iterrows()

Advanced Indexing	Also see NumPy Arrays
<pre>Selecting &gt;&gt;&gt; df3.loc[:,(df3&gt;1).any()] &gt;&gt;&gt; df3.loc[:,(df3&gt;1).all()] &gt;&gt;&gt; df3.loc[:,df3.isnull().any()] &gt;&gt;&gt; df3.loc[:,df3.notnull().all() Indexing With isin &gt;&gt;&gt; df[(df.Country.isin(df2.Type) &gt;&gt;&gt; df3.filter(items="a","b"]) &gt;&gt;&gt; df.select(lambda x: not x%5) Where &gt;&gt;&gt; s.where(s &gt; 0) Query &gt;&gt;&gt; df6.query('second &gt; first')</pre>	Select cols with any vals >1 Select cols with vals >1 Select cols with NaN Select cols without NaN
Setting/Resetting Index	
<pre>&gt;&gt;&gt; df.set_index('Country') &gt;&gt;&gt; df4 = df.reset_index() &gt;&gt;&gt; df = df.rename(index=str,</pre>	Set the index Reset the index Rename DataFrame ":"cptl", tion":"ppltn"})
Reindexing	
>>> s2 = s.reindex(['a','c','d','e	e','b'])
Forward Filling	Backward Filling
<pre>&gt;&gt;&gt; df.reindex(range(4),</pre>	<pre>&gt;&gt;&gt; s3 = s.reindex(range(5),</pre>
MultiIndexing	
<pre>&gt;&gt;&gt; arrays = [np.array([1,2,3]),</pre>	<pre>uples(tuples, names=['first', 'second']) .rand(3, 2), index=index)</pre>
Duplicate Data	
<pre>&gt;&gt;&gt; s3.unique() &gt;&gt;&gt; df2.duplicated('Type') &gt;&gt;&gt; df2.drop_duplicates('Type', kg &gt;&gt;&gt; df.index.duplicated()</pre>	eep='last') Return unique values Check duplicates Drop duplicates Check index duplicates
•	

#### Aggregation >>> df2.groupby(by=['Date','Type']).mean() >>> df4.groupby(level=0).sum() >>> df4.groupby(level=0).agg({'a':lambda x:sum(x)/len(x), 'b': np.sum}) Transformation >>> customSum = lambda x: (x+x%2) >>> df4.groupby(level=0).transform(customSum)

## **Missing Data**

- >>> df.dropna() >>> df3.fillna(df3.mean()) >>> df2.replace("a", "f")
- Drop NaN values Fill NaN values with a predetermined value Replace values with others

5	Combining Data
s >1 N	data1     data2       X1     X2     X1     X3       a     11.432     a     20.784       b     1.303     b     NaN       c     99.906     d     20.784
s	Merge
	<pre>&gt;&gt;&gt; pd.merge(data1,</pre>
	<pre>&gt;&gt;&gt; pd.merge(data1,</pre>
	>>> pd.merge(data1, data2, how='inner', on='X1') X1 X2 X3 a 11.432 20.784 b 1.303 NaN
1')	>>> pd.merge(data1, data2, how='outer', on='X1')
	Join >>> data1.join(data2, how='right')
1)	<pre>Concatenate Vertical &gt;&gt;&gt; s.append(s2) Horizontal/Vertical &gt;&gt;&gt; pd.concat([s,s2],axis=1, keys=['One','Two']) &gt;&gt;&gt; pd.concat([data1, data2], axis=1, join='inner')</pre>
	Dates
25	<pre>&gt;&gt;&gt; df2['Date']= pd.to_datetime(df2['Date']) &gt;&gt;&gt; df2['Date']= pd.date_range('2000-1-1',</pre>
	Visualization Also see Matplotlib
	>>> import matplotlib.pyplot as plt
	>>> s.plot() >>> plt.show() >>> plt.show()

DataCamp

Learn Python for Data Science Interactively

9

**NumPy Basics** 

Learn Python for Data Science Interactively at www.DataCamp.com

# Q

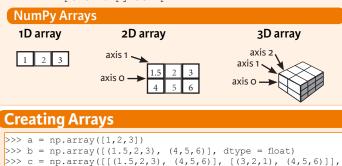
## NumPy

The **NumPy** library is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays.

NumPy

Use the following import convention:

```
>>> import numpy as np
```



#### **Initial Placeholders**

>>>	np.zeros((3,4))	Create an array of zeros
>>>	np.ones((2,3,4),dtype=np.int16)	Create an array of ones
>>>	<pre>d = np.arange(10,25,5)</pre>	Create an array of evenly
		spaced values (step value)
>>>	np.linspace(0,2,9)	Create an array of evenly
		spaced values (number of samples)
>>>	e = np.full((2,2),7)	Create a constant array
>>>	f = np.eye(2)	Create a 2X2 identity matrix
>>>	np.random.random((2,2))	Create an array with random values
>>>	np.empty((3,2))	Create an empty array

dtype = float)

## I/O

## Saving & Loading On Disk

>>> np.save('my\_array', a)
>>> np.savez('array.npz', a, b)
>>> np.load('my\_array.npy')

## Saving & Loading Text Files

>>> np.loadtxt("myfile.txt")
>>> np.genfromtxt("my\_file.csv", delimiter=',')
>>> np.savetxt("myarray.txt", a, delimiter=" ")

## Data Types

- >>> np.int64
  >>> np.float32
  >>> np.complex
  >>> np.bool
  >>> np.object
  >>> np.string
  >>> np.unicode
- Signed 64-bit integer types Standard double-precision floating point Complex numbers represented by 128 floats Boolean type storing TRUE and FALSE values Python object type Fixed-length string type Fixed-length unicode type

Inspecting Your Array				
>>> a.shape       Array dimension         >>> len(a)       Length of array         >>> b.ndim       Number of array         >>> e.size       Number of array         >>> b.dtype       Data type of arr         >>> b.dtype.name       Name of data type		dimensions elements ay elements		
Asking For Help	e)			
Array Mathematics				
Arithmetic Operations				
>>> g = a - b array([[-0.5, 0., 0.],		Subtraction		
<pre>[-3., -3., -3.]]) &gt;&gt;&gt; np.subtract(a,b) &gt;&gt;&gt; b + a array([[ 2.5, 4., 6.],</pre>		Subtraction Addition		
<pre>[ 5. , 7. , 9. ]]) &gt;&gt;&gt; np.add(b,a) &gt;&gt;&gt; a / b array([[ 0.666666667, 1. [ 0.25] , 0.4</pre>	, 1. ], , 0.5 ]])	Addition Division		
<pre>&gt;&gt;&gt; np.divide(a,b) &gt;&gt;&gt; a * b array([[ 1.5, 4., 9.</pre>	Division Multiplication			
<pre>&gt;&gt;&gt; np.multiply(a,b) &gt;&gt;&gt; np.exp(b) &gt;&gt;&gt; np.sqrt(b) &gt;&gt;&gt; np.sin(a) &gt;&gt;&gt; np.log(a) &gt;&gt;&gt; e.dot(f) array([[7., 7.],         [7., 7.]])</pre>	,	Multiplication Exponentiation Square root Print sines of an array Element-wise cosine Element-wise natural logarithm Dot product		
Comparison				
<pre>&gt;&gt;&gt; a == b array([[False, True, True]</pre>		Element-wise comparison		
<pre>[False, False, False]], dtype=bool) &gt;&gt;&gt; a &lt; 2 array([True, False, False], dtype=bool)</pre>		Element-wise comparison		
>>> np.array_equal(a, b)		Array-wise comparison		
Aggregate Functions				
<pre>&gt;&gt;&gt; a.sum() &gt;&gt;&gt; a.min() &gt;&gt;&gt; b.max(axis=0) &gt;&gt;&gt; b.cumsum(axis=1) &gt;&gt;&gt; a.mean() &gt;&gt;&gt; b.median() &gt;&gt;&gt; a.corrcoef() &gt;&gt;&gt; np.std(b)</pre>	Maximu Cumulat Mean Median Correlat	ise sum ise minimum value m value of an array row ive sum of the elements ion coefficient d deviation		

Create a view of the array with the same data

Sort the elements of an array's axis

Create a copy of the array

Sort an array

Create a deep copy of the array

## Copying Arrays

>>> h = a.view() >>> np.copy(a) >>> h = a.copy()

## **Sorting Arrays**

>>> a.sort() >>> c.sort(axis=0)

Subsetting, Slicing, Indexing Also see Lists Subsetting 1 2 3 Select the element at the 2nd index >>> a[2] 3 1.5 2 3 >>> b[1,2] Select the element at row 1 column 2 4 5 6 (equivalent to b[1] [2]) 6.0 Slicina >>> a[0:2] Select items at index 0 and 1 1 2 3 array([1, 2]) 3 >>> b[0:2,1] Select items at rows 0 and 1 in column 1 5 6 array([ 2., 5.]) 4 3 Select all items at row o >>> b[:1] 4 5 6 (equivalent to b[0:1, :]) array([[1.5, 2., 3.]]) >>> c[1,...] Same as [1, :, :] array([[[ 3., 2., 1.], [ 4., 5., 6.]]]) >>> a[ : :-1] Reversed array a array([3, 2, 1]) **Boolean Indexing** >>> a[a<2] Select elements from a less than 2 1 2 3 array([1]) Fancy Indexing >>> b[[1, 0, 1, 0], [0, 1, 2, 0]] Select elements (1,0), (0,1), (1,2) and (0,0) array([ 4. , 2. , 6. , 1.5]) Select a subset of the matrix's rows >>> b[[1, 0, 1, 0]][:,[0,1,2,0]] array([[4, ,5, ,6, ,4,], [1.5, 2. , 3, , 1.5], [4. , 5. , 6. , 4.], [1.5, 2. , 3. , 1.5]]) and columns Array Manipulation Transposing Array >>> i = np.transpose(b) Permute array dimensions Permute array dimensions >>> i.T

Changing Array Shape
>>> b.ravel()
>>> g.reshape(3,-2)

Adding/Removing Elements

>>> h.resize((2,6))
>>> np.append(h,g)
>>> np.insert(a, 1, 5)
>>> np.delete(a,[1])

#### Combining Arrays

>>> np.concatenate((a,d),axis=0)
array([1, 2, 3, 10, 15, 20])
>>> np.vstack((a,b))
array([[1, 2, 3, 3, ],
 [1.5, 2, 3, 3, ],
 [4, 5, 2, 6, ]])
>>> np.r\_[e,f]
>>> np.n\_r[e,f]
>>> np.hstack((e,f))
array([[7, 7, 1, 1, 0,],
 [7, 7, 0, 1, 1]])
>>> np.column\_stack((a,d))
array([[1, 10],
 [2, 15],
 [3, 20]])
>>> np.c\_[a,d]

## **Splitting Arrays**

>>> np.hsplit(a,3)
[array([1]),array([2]),array([3])]
>>> np.vsplit(c,2)
[array([[ 1.5, 2., 1.],
 [ 4., 5., 6. ]]]),
 array([[ 3., 2., 3.],
 [ 4., 5., 6.]])]

Stack arrays vertically (row-wise) Stack arrays vertically (row-wise) Stack arrays horizontally (column-wise)

Flatten the array

Reshape, but don't change data

Append items to an array

Delete items from an array

Insert items in an array

Concatenate arrays

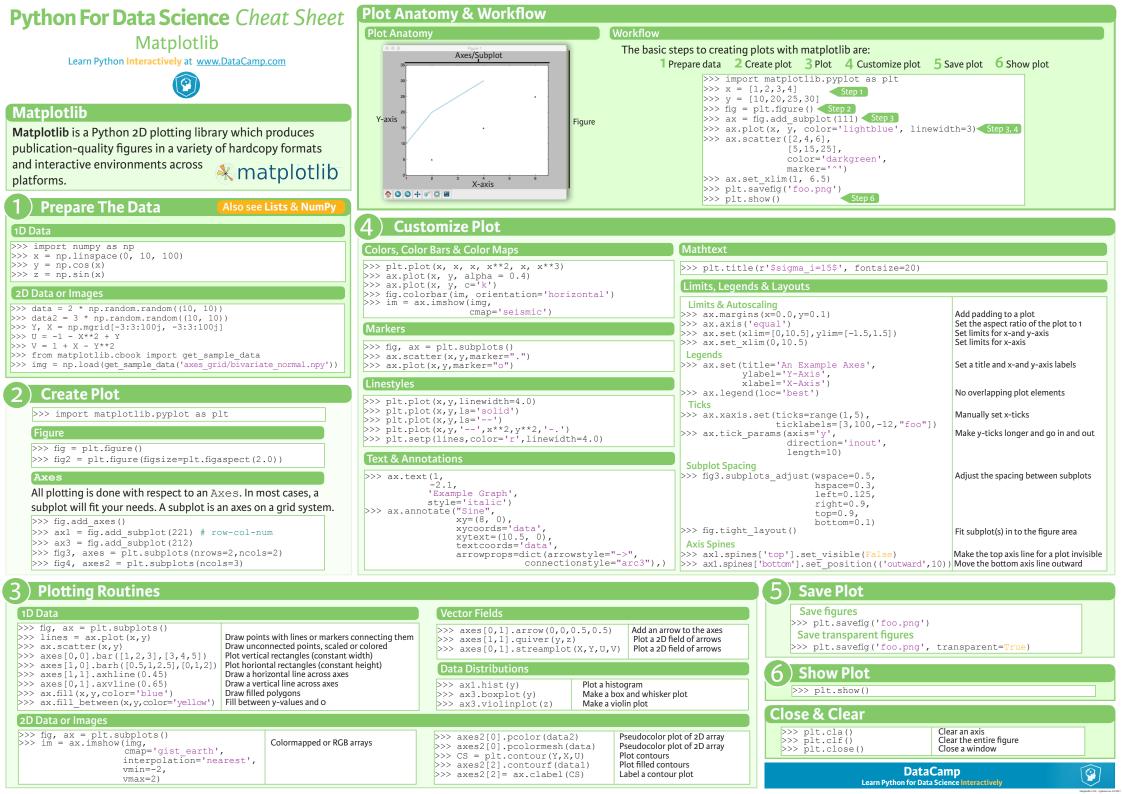
Return a new array with shape (2,6)

Create stacked column-wise arrays

Create stacked column-wise arrays

Split the array horizontally at the 3rd index Split the array vertically at the 2nd index

DataCamp Learn Python for Data Science Interactively  $\bigcirc$ 



# Python For Data Science Cheat Sheet 3 Plotting With Seaborn

Seaborn

Learn Data Science Interactively at www.DataCamp.com



## Statistical Data Visualization With Seaborn

The Python visualization library **Seaborn** is based on matplotlib and provides a high-level interface for drawing attractive statistical graphics.

## Make use of the following aliases to import the libraries:

ne basic steps to creating   1. Prepare some data	plots with Seaborn are:	
2. Control figure aesthetic	rs in the second s	
3. Plot with Seaborn		
4. Further customize your	rplot	
> import matplotlib.pypl	ot as plt	
> import seaborn as sns		
> tips = sns.load_datase > sns.set style("whitegr		
> g = sns.lmplot(x="tip"		
	l bill",	
data=ti		
aspect=	-	
<pre>&gt; g = (g.set_axis_tables t(xlim=(0,10),ylim=(0,10)</pre>	("Tip", "Total bill(USD)").	
<pre>&gt; plt.title("title")</pre>	Step 4	
· · · · ·	tep 5	
Data	Also see Lists, NumPy & Pan	das

Seaborn also offers built-in data sets:			
	'y':np.random.normal(0,4,100)})		
>>>	<pre>data = pd.DataFrame({'x':np.arange(1,101),</pre>		
	uniform data = np.random.rand(10, 12)		
	import numpy as np		
	import pandas as pd		

>>> titanic = sns.load\_dataset("titanic")
>>> iris = sns.load\_dataset("iris")

2) Figure Aesthetics			Also see Matplotlib
		Context Functions	
>>> f, ax = plt.subplots(figsize=(5,6))	Create a figure and one subplot	>>> sns.set context("talk")	Set context to "talk"
Seaborn styles		<pre>&gt;&gt;&gt; sns.set_context("notebook",</pre>	Set context to "notebook", scale font elements and override param mapping
<pre>&gt;&gt;&gt; sns.set() (Re)set the seaborn default &gt;&gt;&gt; sns.set_style("whitegrid") Set the matplotlib parameters</pre>		Color Palette	
{"xtick.major.size":8, "ytick.major.size":8}) >>> sns.axes style("whitegrid")	Set the matplotlib parameters Return a dict of params or use with Atch to temporarily set the style	<pre>&gt;&gt;&gt; sns.color_palette("husl") Use with wi &gt;&gt;&gt; flatui = ["#9b59b6","#3498db","#95a5a6","#e74c</pre>	color palette th to temporarily set palette 3c", "#34495e", "#2ecc71"] n color palette

Axis Grids			
<pre>&gt;&gt;&gt; g = sns.FacetGrid(titanic,</pre>	Subplot grid for plotting conditional relationships Draw a categorical plot onto a Facetgrid Plot data and regression model fits across a FacetGrid	<pre>&gt;&gt;&gt; h = h.map(plt.scatter) &gt;&gt;&gt; sns.pairplot(iris) &gt;&gt;&gt; i = sns.JointGrid(x="x",</pre>	Subplot grid for plotting pairwise relationships Plot pairwise bivariate distributions Grid for bivariate plot with marginal univariate plots Plot bivariate distribution
Categorical Plots		Regression Plots	
<pre>Scatterplot &gt;&gt;&gt; sns.stripplot(x="species",     y="petal_length",     data=iris)</pre>	Scatterplot with one categorical variable	<pre>&gt;&gt;&gt; sns.regplot(x="sepal_width",</pre>	Plot data and a linear regression model fit
>>> sns.swarmplot(x="species",	Categorical scatterplot with non-overlapping points	Distribution Plots	
y="petal_length", data=iris) Bar Chart	non-ovenapping points	>>> plot = sns.distplot(data.y, kde=False, color="b")	
<pre>&gt;&gt;&gt; sns.barplot(x="sex",</pre>	Show point estimates and confidence intervals with	Matrix Plots	
hue="class", data=titanic)	scatterplot glyphs	>>> sns.heatmap(uniform_data,vmin=	=0,vmax=1) Heatmap
<pre>Count Plot &gt;&gt;&gt; sns.countplot(x="deck",</pre>	Show count of observations	4) Further Customizatio	Ons Also see Matplotlik
<pre>palette="Greens_d") Point Plot</pre>		Axisgrid Objects	
<pre>&gt;&gt;&gt; sns.pointplot(x="class", y="survived", hue="sex", data=titanic, palette={"male":"g",</pre>		<pre>&gt;&gt;&gt; g.despine(left=True) &gt;&gt;&gt; g.set_ylabels("Survived") &gt;&gt;&gt; g.set_axticklabels(rotation &gt;&gt;&gt; g.set_axis_labels("Survive</pre>	
<pre>Boxplot &gt;&gt;&gt; sns.boxplot(x="alive",</pre>	Boxplot	yticks=[0,2.5,5])	
<pre>y="age", hue="adult_male", data=titanic) &gt;&gt;&gt; sns.boxplot(data=iris,orient="h") Violinplot &gt;&gt;&gt; sns.violinplot(x="age", y="sex", hue="survived", data=titanic)</pre>	Boxplot with wide-form data Violin plot	<pre>&gt;&gt;&gt; plt.xlabel("Sex") &gt;&gt;&gt; plt.ylim(0,100) &gt;&gt;&gt; plt.xlim(0,10) &gt;&gt;&gt; plt.setp(ax,yticks=[0,5]) &gt;&gt;&gt; plt.tight_layout()</pre>	Add plot title Adjust the label of the y-axis Adjust the label of the x-axis Adjust the limits of the y-axis Adjust the limits of the x-axis Adjust a plot property Adjust subplot params
	Also see Matplotlib	5) Show or Save Plot	Also see Matplotlik
Context Functions	Also see Matpiotilb	<pre>&gt;&gt;&gt; plt.show() &gt;&gt;&gt; plt.savefig("foo.png")</pre>	Show the plot Save the plot as a figure
<pre>&gt;&gt;&gt; sns.set_context("talk")</pre>	Set context to "talk"	>>> plt.saveng("foo.png") >>> plt.savefig("foo.png", transparent=Th	Save transparent figure

**Close & Clear** 

>>> plt.cla()
>>> plt.clf()
>>> plt.close()

Clear an axis Clear an entire figure Close a window

9

DataCamp

Learn Python for Data Science Interactively

## Bokeh

Learn Bokeh Interactively at www.DataCamp.com, taught by Bryan Van de Ven, core contributor

## **Plotting With Bokeh**

The Python interactive visualization library **Bokeh** enables high-performance visual presentation of large datasets in modern web browsers.



Bokeh's mid-level general purpose bokeh.plotting interface is centered around two main components: data and glyphs.



The basic steps to creating plots with the **bokeh.plotting** interface are:

1. Prepare some data:

Python lists, NumPy arrays, Pandas DataFrames and other sequences of values

- 2. Create a new plot
- 3. Add renderers for your data, with visual customizations
- 4. Specify where to generate the output

#### 5. Show or save the results

## ) Data

## Also see Lists, NumPy & Pandas

Under the hood, your data is converted to Column Data

#### Sources. You can also do this manually:

>>> from bokeh.models import ColumnDataSource
>>> cds df = ColumnDataSource(df)

## 2 Plotting

## **Renderers & Visual Customizations**

#### Glyphs

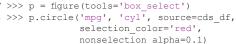
#### Scatter Markers

color='blue', size=1)

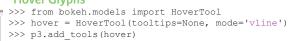
#### Line Glyphs

#### **Customized Glyphs**

#### Selection and Non-Selection Glyphs



## Hover Glyphs



#### Colormapping

legend='Origin')

#### Legend Location

#### Inside Plot Area

>>> p.legend.location = 'bottom\_left'

#### **Outside Plot Area**

## Legend Orientation

>>> p.legend.orientation = "horizontal"
>>> p.legend.orientation = "vertical"

## Legend Background & Border

>>> p.legend.border\_line\_color = "navy"
>>> p.legend.background\_fill\_color = "white"

## Rows & Columns Layout

Rows
>>> from bokeh.layouts import row
>>> layout = row(p1,p2,p3)

Columns
>>> from bokeh.layouts import columns
>>> layout = column(p1,p2,p3)

Nesting Rows & Columns
>>>layout = row(column(p1,p2), p3)

## **Grid Layout**

>>> from bokeh.layouts import gridplot
>>> row1 = [p1,p2]
>>> row2 = [p3]
>>> layout = gridplot([[p1,p2],[p3]])

## **Tabbed Layout**

>>> from bokeh.models.widgets import Panel, Tabs
>>> tab1 = Panel(child=p1, title="tab1")
>>> tab2 = Panel(child=p2, title="tab2")
>>> layout = Tabs(tabs=[tab1, tab2])

## Linked Plots

#### Linked Axes

>>> p2.x\_range = p1.x\_range
>>> p2.y\_range = p1.y\_range

#### Linked Brushing

## 1 ) Output & Export

## Notebook

>>> from bokeh.io import output\_notebook, show
>>> output\_notebook()

## HTML

#### Standalone HTML

>>> from bokeh.embed import file\_html
>>> from bokeh.resources import CDN
>>> html = file html(p, CDN, "my plot")

>>> from bokeh.io import output\_file, show
>>> output\_file('my\_bar\_chart.html', mode='cdn')

#### Components

>>> from bokeh.embed import components >>> script, div = components(p)

## PNG

>>> from bokeh.io import export\_png
>>> export png(p, filename="plot.png")

## SVG

- >>> from bokeh.io import export\_svgs
- >>> p.output\_backend = "svg"
- >>> export\_svgs(p, filename="plot.svg")

## Show or Save Your Plots

>>> show(p1) >>> save(p1)

>>> show(layout) >>> save(layout)

9

### SciPy - Linear Algebra

Learn More Python for Data Science Interactively at www.datacamp.com

### 0

### SciPy

The SciPy library is one of the core packages for scientific computing that provides mathematical S SciPv algorithms and convenience functions built on the NumPy extension of Python.

### Interacting With NumPy

>>> import numpy as np >>> a = np.array([1, 2, 3]) >>> b = np.array([(1+5j,2j,3j), (4j,5j,6j)]) >>> c = np.array([[(1.5,2,3), (4,5,6)], [(3,2,1), (4,5,6)]])

#### Index Tricks

>>> np.mgrid[0:5,0:5] >>> np.ogrid[0:2,0:2] >>> np.r [[3,[0]\*5,-1:1:10j] >>> np.c [b,c]

Create a dense meshgrid Create an open meshgrid Stack arrays vertically (row-wise) Create stacked column-wise arrays

Create a polynomial object

Return the real part of the array elements

Cast object to a data type

Return the imaginary part of the array elements

Return a real array if complex parts close to o

Also see NumPv

### Shape Manipulation

>>> np.transpose(b) Permute array dimensions >>> b.flatten() Flatten the array >>> np.hstack((b,c)) Stack arrays horizontally (column-wise) Stack arrays vertically (row-wise) >>> np.vstack((a,b)) Split the array horizontally at the 2nd index >>> np.hsplit(c,2) Split the array vertically at the 2nd index >>> np.vpslit(d,2)

### Polvnomials

>>> from numpy import poly1d >>> p = poly1d([3, 4, 5])

### **Vectorizing Functions**

>>> def myfunc(a): if a < 0:return a\*2 else. return a/2 >>> np.vectorize(myfunc) Vectorize functions

### Type Handling

>>> np.real(c) >>> np.imag(c) >>> np.real if close(c,tol=1000) >>> np.cast['f'](np.pi)

### Other Useful Functions

	<pre>np.angle(b,deg=True)</pre>	Return the angle of the complex argument
>>>	<pre>g = np.linspace(0,np.pi,num=5)</pre>	Create an array of evenly spaced values
>>>	g [3:] += np.pi	(number of samples)
>>>	np.unwrap(g)	Unwrap
>>>	np.logspace(0,10,3)	Create an array of evenly spaced values (log scale)
>>>	np.select([ $c<4$ ],[ $c*2$ ])	Return values from a list of arrays depending on
		conditions
>>>	misc.factorial(a)	Factorial
>>>	<pre>misc.comb(10,3,exact=True)</pre>	Combine N things taken at k time
>>>	<pre>misc.central_diff_weights(3)</pre>	Weights for Np-point central derivative
>>>	misc.derivative(myfunc,1.0)	Find the n-th derivative of a function at a point

### Linear Algebra

You'll use the linalg and sparse modules. Note that scipy.linalg contains and expands on numpy.linalg.

>>> from scipy import linalg, sparse

### Creating Matrices

>>> A = np.matrix(np.random.random((2,2))) >>> B = np.asmatrix(b) >>> C = np.mat(np.random.random((10,5)))

Inverse

Inverse

Trace

Tranpose matrix

Frobenius norm

Matrix rank

Determinant

equation

(SVD)

Conjugate transposition

L1 norm (max column sum)

L inf norm (max row sum)

Solver for dense matrices

Solver for dense matrices

(least-squares solver)

Sparse matrix exponential

Least-squares solution to linear matrix

Compute the pseudo-inverse of a matrix

Compute the pseudo-inverse of a matrix

>>> D = np.mat([[3,4], [5,6]])

### **Basic Matrix Routines**

#### Inverse

>>> A.I >>> linalg.inv(A) >>> A.T >>> A.H >>> np.trace(A)

#### Norm

>>> linalg.norm(A) >>> linalg.norm(A,1) >>> linalg.norm(A,np.inf)

Rank >>> np.linalg.matrix rank(C)

Determinant >>> linalg.det(A) Solving linear problems

>>> linalg.solve(A,b) >>> E = np.mat(a).T >>> linalg.lstsq(D,E)

### Generalized inverse

>>> linalg.pinv(C)

>>> linalq.pinv2(C)

### **Creating Sparse Matrices**

- >>> F = np.eye(3, k=1) Create a 2X2 identity matrix >>> G = np.mat(np.identity(2)) Create a 2x2 identity matrix >>> C[C > 0.5] = 0>>> H = sparse.csr matrix(C) Compressed Sparse Row matrix >>> I = sparse.csc matrix(D) Compressed Sparse Column matrix Dictionary Of Keys matrix >>> J = sparse.dok matrix(A) >>> E.todense() Sparse matrix to full matrix Identify sparse matrix >>> sparse.isspmatrix csc(A) Sparse Matrix Routines
- Inverse >>> sparse.linalg.inv(I) Inverse Norm Norm >>> sparse.linalg.norm(I) Solving linear problems
- >>> sparse.linalg.spsolve(H,I) Solver for sparse matrices

### Sparse Matrix Functions

>>> sparse.linalq.expm(I)

Asking For Help >>> help(scipy.linalg.diagsvd) >>> np.info(np.matrix)

### Matrix Functions

#### Addition >>> np.add(A,D)

Subtraction >>> np.subtract(A,D) Division

#### >>> np.divide(A,D) Multiplication

>>> np.multiply(D,A) >>> np.dot(A,D) >>> np.vdot(A,D) >>> np.inner(A,D) >>> np.outer(A,D) >>> np.tensordot(A,D) >>> np.kron(A,D)

#### **Exponential Functions**

>>> linalg.expm(A) >>> linalg.expm2(A) >>> linalq.expm3(D)

#### Logarithm Function >>> linalg.logm(A)

**Trigonometric Tunctions** >>> linalg.sinm(D) >>> linalg.cosm(D) >>> linalg.tanm(A)

#### Hyperbolic Trigonometric Functions >>> linalq.sinhm(D) >>> linalq.coshm(D) >>> linalg.tanhm(A)

**Matrix Sign Function** >>> np.sigm(A)

#### Matrix Square Root >>> linalg.sqrtm(A)

**Arbitrary Functions** >>> linalg.funm(A, lambda x: x\*x)

### Matrix sign function Matrix square root

Hypberbolic matrix sine

Hyperbolic matrix cosine

Hyperbolic matrix tangent

Evaluate matrix function

### Decompositions

**Eigenvalues and Eigenvectors** Solve ordinary or generalized >>> la, v = linalg.eig(A) eigenvalue problem for square matrix >>> 11, 12 = la Unpack eigenvalues First eigenvector >>> v[:,0] >>> v[:,1] Second eigenvector Unpack eigenvalues >>> linalg.eigvals(A) **Singular Value Decomposition** Singular Value Decomposition (SVD) >>> U,s,Vh = linalq.svd(B) >>> M.N = B.shape Construct sigma matrix in SVD >>> Sig = linalg.diagsvd(s,M,N) LU Decomposition >>> P,L,U = linalg.lu(C) LU Decomposition

### Sparse Matrix Decompositions

>>> la, v = sparse.linalg.eigs(F,1) **Eigenvalues and eigenvectors** SVD >>> sparse.linalg.svds(H, 2)

Also see NumPy

Addition

Division

Subtraction

Multiplication

Inner product

Outer product

decomposition)

Matrix sine

Matrix cosine Matrix tangent

Vector dot product

Tensor dot product

Kronecker product

Matrix exponential

Matrix logarithm

Matrix exponential (Taylor Series)

Matrix exponential (eigenvalue

Dot product

Scikit-Learn

Learn Python for data science Interactively at www.DataCamp.com



### Scikit-learn

Scikit-learn is an open source Python library that implements a range of machine learning, preprocessing, cross-validation and visualization algorithms using a unified interface.

### A Basic Example

- >>> from sklearn import neighbors, datasets, preprocessing
  >>> from sklearn.model\_selection import train\_test\_split
  >>> from sklearn.metrics import accuracy\_score
  >>> iris = datasets.load\_iris()
  >>> X, y = iris.data[:, :2], iris.target
  >>> X\_train,X\_test, y\_train, y\_test= train\_test\_split(X, y, random\_state=33)
  >>> scaler = preprocessing.StandardScaler().fit(X\_train)
  >>> X\_train = scaler.transform(X\_train)
  >>> k\_n = neighbors.KNeighborsClassifier(n\_neighbors=5)
- >>> knn.fit(X\_train, y\_train)
  >>> y pred = knn.predict(X test)
- >>> y\_pred = knn.predict(X\_test)
  >>> accuracy score(y test, y pred)

### **Loading The Data**

#### Also see NumPy & Pandas

Your data needs to be numeric and stored as NumPy arrays or SciPy sparse matrices. Other types that are convertible to numeric arrays, such as Pandas DataFrame, are also acceptable.

>>> import numpy as np

### **Training And Test Data**

### **Preprocessing The Data**

### Standardization

### >>> from sklearn.preprocessing import StandardScaler

- >>> scaler = StandardScaler().fit(X\_train)
- >>> standardized X = scaler.transform(X train) >>> standardized X test = scaler.transform(X test)
- >>> stanuarurzeu\_A\_test = scater.transform(X\_tes

### Normalization

- >>> from sklearn.preprocessing import Normalizer
- >>> scaler = Normalizer().fit(X\_train)
- >>> normalized\_X = scaler.transform(X\_train)
  >>> normalized\_X test = scaler.transform(X test)
- /// normallzed\_x\_test = scaler.transform(X\_tes

### Binarization

>>> from sklearn.preprocessing import Binarizer
>>> binarizer = Binarizer(threshold=0.0).fit(X)
>>> binary\_X = binarizer.transform(X)

### **Create Your Model**

#### Supervised Learning Estimators

#### Linear Regression

>>> from sklearn.linear\_model import LinearRegression
>>> lr = LinearRegression(normalize=True)

#### Support Vector Machines (SVM)

>>> from sklearn.svm import SVC
>>> svc = SVC(kernel='linear')

#### **Naive Bayes**

>>> from sklearn.naive\_bayes import GaussianNB >>> gnb = GaussianNB()

KNN

>>> from sklearn import neighbors
>>> knn = neighbors.KNeighborsClassifier(n\_neighbors=5)

### Unsupervised Learning Estimators

#### Principal Component Analysis (PCA)

>>> from sklearn.decomposition import PCA
>>> pca = PCA(n\_components=0.95)

K Means

>>> from sklearn.cluster import KMeans
>>> k\_means = KMeans(n\_clusters=3, random\_state=0)

### **Model Fitting**

### Supervised learning

>>> lr.fit(X, y)
>>> knn.fit(X\_train, y\_train)
>>> svc.fit(X\_train, y\_train)
Unsupervised Learning

#### >>> k\_means.fit(X\_train)

>>> pca\_model = pca.fit\_transform(X\_train) Fit to data, then transform it

Fit the model to the data

Fit the model to the data

Predict labels

Predict labels

Estimate probability of a label

Predict labels in clustering algos

### Prediction

# Supervised Estimators >>> y\_pred = svc.predict(np.random.random((2,5))) >>> y\_pred = lr.predict(X\_test) >>> y pred = knn.predict proba(X test)

Unsupervised Estimators

>>> y\_pred = k\_means.predict(X\_test)

### **Encoding Categorical Features**

>>> from sklearn.preprocessing import LabelEncoder
>>> enc = LabelEncoder()
>>> y = enc.fit\_transform(y)

### Imputing Missing Values

>>> from sklearn.preprocessing import Imputer
>>> imp = Imputer(missing\_values=0, strategy='mean', axis=0)
>>> imp.fit\_transform(X\_train)

### **Generating Polynomial Features**

>>> from sklearn.preprocessing import PolynomialFeatures
>>> poly = PolynomialFeatures(5)
>>> poly.fit\_transform(X)

### **Evaluate Your Model's Performance**

### **Classification Metrics**

### Accuracy Score

>>> knn.score(X\_test, y\_test) Estimator score method >>> from sklearn.metrics import accuracy\_score Metric scoring functions

>>> accuracy\_score(y\_test, y\_pred)
Classification Report
>>> from sklearn.metrics import classification report Precision, recall, fi-score

>>> print(classification\_report(y\_test, y\_pred))
and support
Confusion Matrix

>>> from sklearn.metrics import confusion\_matrix
>>> print(confusion\_matrix(y\_test, y\_pred))

### **Regression Metrics**

### Mean Absolute Error

>>> from sklearn.metrics import mean\_absolute\_error
>>> y\_true = [3, -0.5, 2]
>>> mean\_absolute\_error(y\_true, y\_pred)

#### Mean Squared Error >>> from sklearn.metrics import mean squared\_error >>> mean squared error(y test, y pred)

R<sup>2</sup> Score
>>> from sklearn.metrics import r2\_score
>>> r2 score(y true, y pred)

### **Clustering Metrics**

### **Adjusted Rand Index**

>>> from sklearn.metrics import adjusted\_rand\_score
>>> adjusted\_rand\_score(y\_true, y\_pred)

#### Homogeneity

>>> from sklearn.metrics import homogeneity\_score
>>> homogeneity\_score(y\_true, y\_pred)

V-measure

>>> from sklearn.metrics import v\_measure\_score
>>> metrics.v\_measure\_score(y\_true, y\_pred)

### **Cross-Validation**

>>> from sklearn.cross\_validation import cross\_val\_score
>>> print(cross\_val\_score(knn, X\_train, y\_train, cv=4))
>>> print(cross\_val\_score(lr, X, y, cv=2))

### Tune Your Model

### **Grid Search**

- >>> grid.nt(x\_train, y\_trai >>> print(grid.best score )
- >>> print(grid.best\_stimator .n neighbors)

### Randomized Parameter Optimization

DataCamp

Learn Python for Data Science Interactively

9

>>> rsearch.fit(X\_train, y\_train)
>>> print(rsearch.best\_score\_)

Keras

Learn Python for data science Interactively at www.DataCamp.com

### 0

### Keras

Keras is a powerful and easy-to-use deep learning library for Theano and TensorFlow that provides a high-level neural networks API to develop and evaluate deep learning models.

#### A Basic Example

>>> import numpy as np >>> from keras.models import Sequential >>> from keras.layers import Dense >>> data = np.random.random((1000,100)) >>> labels = np.random.randint(2,size=(1000,1)) >>> model = Sequential() >>> model.add(Dense(32, activation='relu', input dim=100)) >>> model.add(Dense(1, activation='sigmoid')) >>> model.compile(optimizer='rmsprop', loss='binary crossentropy', metrics=['accuracy']) >>> model.fit(data,labels,epochs=10,batch size=32) >>> predictions = model.predict(data)

#### Data

Also see NumPy, Pandas & Scikit-Learn

Your data needs to be stored as NumPy arrays or as a list of NumPy arrays. Ideally, you split the data in training and test sets, for which you can also resort

to the train test split module of sklearn.cross validation.

### Keras Data Sets

>>> from keras.datasets import boston\_housing, mnist. cifar10, imdb

- >>> (x\_train,y\_train),(x\_test,y\_test) = mnist.load data() >>> (x train2, y train2), (x test2, y test2) = boston housing.load data() >>> (x\_train3,y\_train3),(x\_test3,y\_test3) = cifar10.load\_data() >>> (x train4,y train4),(x test4,y test4) = imdb.load data(num words=20000)
- >>> num classes = 10

#### Other

>>> from urllib.request import urlopen >>> data = np.loadtxt(urlopen("http://archive.ics.uci.edu/ ml/machine-learning-databases/pima-indians-diabetes/ pima-indians-diabetes.data"),delimiter=",") >>> X = data[:,0:8] >>> y = data [:,8]

### Preprocessing

### Sequence Padding

>>> from keras.preprocessing import sequence >>> x train4 = sequence.pad sequences(x train4,maxlen=80) >>> x test4 = sequence.pad sequences(x test4,maxlen=80)

### One-Hot Encoding

- >>> from keras.utils import to categorical
- >>> Y train = to categorical(y train, num classes) >>> Y test = to categorical(y test, num classes)
- >>> Y\_train3 = to\_categorical(y\_train3, num\_classes) >>> Y\_test3 = to\_categorical(y\_test3, num\_classes)

Model Architecture

### Sequential Model

>>> from keras.models import Sequential >>> model = Sequential() >>> model2 = Seguential() >>> model3 = Sequential()

### ( Multilayer Perceptron (MLP)

#### **Binary Classification**

>>> from keras.layers import Dense >>> model.add(Dense(12, input dim=8, kernel initializer='uniform', activation='relu')) >>> model.add(Dense(8,kernel initializer='uniform',activation='relu')) >>> model.add(Dense(1,kernel initializer='uniform',activation='sigmoid'))

#### Multi-Class Classification

>>> from keras.layers import Dropout >>> model.add(Dense(512,activation='relu',input shape=(784,)))

- >>> model.add(Dropout(0.2))
- >>> model.add(Dense(512,activation='relu'))
- >>> model.add(Dropout(0.2))
- >>> model.add(Dense(10,activation='softmax'))

#### Regression

>>> model.add(Dense(64,activation='relu',input dim=train data.shape[1])) >>> model.add(Dense(1))

### Convolutional Neural Network (CNN)

>>> from keras.layers import Activation,Conv2D,MaxPooling2D,Flatten >>> model2.add(Conv2D(32,(3,3),padding='same',input shape=x train.shape[1:])) >>> model2.add(Activation('relu')) >>> model2.add(Conv2D(32,(3,3))) >>> model2.add(Activation('relu'))

- >>> model2.add(MaxPooling2D(pool size=(2,2)))
- >>> model2.add(Dropout(0.25))
- >>> model2.add(Conv2D(64,(3,3), padding='same'))
- >>> model2.add(Activation('relu'))
- >>> model2.add(Conv2D(64,(3, 3)))
- >>> model2.add(Activation('relu'))
- >>> model2.add(MaxPooling2D(pool size=(2,2))) >>> model2.add(Dropout(0.25))
- >>> model2.add(Flatten()) >>> model2.add(Dense(512))
- >>> model2.add(Activation('relu'))
- >>> model2.add(Dropout(0.5))
- >>> model2.add(Dense(num classes))
- >>> model2.add(Activation('softmax'))

#### Recurrent Neural Network (RNN)

>>> from keras.klayers import Embedding,LSTM

- >>> model3.add(Embedding(20000,128))
- >>> model3.add(LSTM(128,dropout=0.2,recurrent\_dropout=0.2))
- >>> model3.add(Dense(1,activation='sigmoid'))

### Also see NumPy & Scikit-Learn

### **Train and Test Sets**

>>> from sklearn.model selection import train test split >>> X train5,X test5,y train5,y test5 = train test split(X,

> test size=0 33. random state=42)

### Standardization/Normalization

- >>> from sklearn.preprocessing import StandardScaler
- >>> scaler = StandardScaler().fit(x train2)
- >>> standardized X = scaler.transform(x train2)
- >>> standardized X test = scaler.transform(x test2)

### Inspect Model

>>> model.output shape >>> model.summary() >>> model.get config() >>> model.get weights()

Model output shape Model summary representation Model configuration List all weight tensors in the model

### **Compile Model**

#### **MLP: Binary Classification**

>>> model.compile(optimizer='adam', loss='binary crossentropy', metrics=['accuracy'])

#### MLP: Multi-Class Classification

>>> model.compile(optimizer='rmsprop', loss='categorical crossentropy', metrics=['accuracy'])

**MLP: Regression** 

>>> model.compile(optimizer='rmsprop', loss='mse', metrics=['mae'])

#### Recurrent Neural Network

>>> model3.compile(loss='binary crossentropy', optimizer='adam', metrics=['accuracy'])

### **Model Training**

>>> model3.fit(x train4, y Train4, batch size=32, epochs=15, verbose=1, validation data=(x test4, y test4))

### **Evaluate Your Model's Performance**

>>> score = model3.evaluate(x test, y\_test, batch size=32)

### Prediction

>>> model3.predict(x test4, batch size=32) >>> model3.predict classes(x test4,batch size=32)

### Save/ Reload Models

>>> from keras.models import load model >>> model3.save('model file.h5') >>> my model = load model('my model.h5')

### **Model Fine-tuning**

### Optimization Parameters

>>> from keras.optimizers import RMSprop >>> opt = RMSprop(1r=0.0001, decay=1e-6) >>> model2.compile(loss='categorical crossentropy', optimizer=opt, metrics=['accuracy'])

### Early Stopping

>>> from keras.callbacks import EarlyStopping >>> early stopping monitor = EarlyStopping(patience=2) >>> model3.fit(x train4, y train4, batch size=32, epochs=15, validation data=(x test4, y test4),

callbacks=[early\_stopping\_monitor])

9

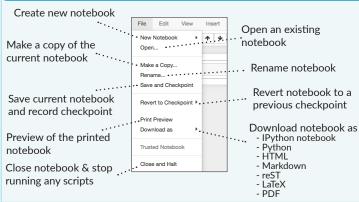
DataCamp Learn Python for Data Science Interactively

Jupyter Notebook

Learn More Python for Data Science Interactively at www.DataCamp.com

### 0

#### Saving/Loading Notebooks



#### Writing Code And Text

Code and text are encapsulated by 3 basic cell types: markdown cells, code cells, and raw NBConvert cells.

#### Edit Cells Cut currently selected cells Copy cells from to clipboard clipboard to current cursor position Paste cells from clipboard above Paste cells from View current cell clipboard below Cut Cells current cell Paste cells from Conv Cells clipboard on top . Delete current cells of current cel Paste Cells & Re Split up a cell from Delete Cells Revert "Delete Cells" Indo Delete Cell current cursor invocation Split Cell position Merge Cell Ab Merge current cell .... Merge current cell Merge Cell Be with the one above with the one below Move Cell Up Move Cell Down Move current cell Move current cell up down Adjust metadata underlying the Find and replace Find and Replace . current notebook in selected cells Cut Cell Attachments Copy Cell Attachments . Remove cell Copy attachments of aste Cell Attachment attachments current cell insert Image . . . . . Paste attachments of Insert image in current cell selected cells Insert Cells Cell Insert Kern Add new cell below the Add new cell above the Insert Cell Above current one current one Insert Cell Below

#### Working with Different Programming Languages

Kernels provide computation and communication with front-end interfaces like the notebooks. There are three main kernels:



Installing Jupyter Notebook will automatically install the IPython kernel.



#### **Command Mode:**

IP[y]:

IPython



IJ[••]

Uulia

Widgets

JavaScript.

Download serialized

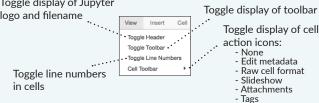
state of all widget

models in use ...

### In [ ]:

di	+	ΝЛ	~	А	~	
	ι	ואו	υ	u	e	

Edit Mode:	1. Save and checkpoint       9. Interrupt kernel         2. Insert cell below       10. Restart kernel         3. Cut cell       11. Display characteristics         4. Copy cell(s)       12. Open command palette         5. Paste cell(s) below       13. Current kernel
Executing Cells         Run selected cell(s)         Run current cells down and create a new one	6. Move cell up14. Kernel status7. Move cell down15. Log out from notebook server9. Dup current cell15. Log out from notebook server
Cell     Kernel     Widgets     below       Run current cells down and create a new one above     • Run Cells Run Cells and Select Below • Run Cells and Insert Below Run All     • Run all cells	Asking For Help Walk through a UI tour
current cell       Run All Below       the current cell         Change the cell type of current cell	Edit the built-in keyboard shortcuts       User Interface Tour Keyboard Shortcuts       shortcuts         Description of markdown available       Notebook Help       Information on
scrolling and clear all output View Cells	in notebook Python help topics NumPy help topics IPython IPyth
Toggle display of Jupyter logo and filename	NumPy help topics       IPython       IPython       SciPy         NumPy       SciPy       SciPy help topics         Matplotlib help topics       Matplotlib       SymPy help topics



Save notebook

with interactive

widgets

widgets

Embed current

Notebook widgets provide the ability to visualize and control changes

You can use them to build interactive GUIs for your notebooks or to

synchronize stateful and stateless information between Python and

Help

Save Notebook with Widgets

Download Widget State

Embed Widgets ....

in your data, often as a control like a slider, textbox, etc.

Widgets

DataCamp arn Python for Data Science Interactively	

2

ß

About Jupyter Notebook

SymPy

pandas

About

Pandas help topics ....

### Jupyter Notebook Markdown Cheatsheet

From <u>SqlBak.com</u> with
-----------------------------

#⊔Header 1 Header 1	Header 1
=======	
##⊔Header 2	Header 2
	neauer 2
Header 2	
###⊔Header 3	Header 3
####⊔Header 4	Header 4
#####uHeader 5	Header 5
*italics*	italics
_italics_	nanoo
<pre>\*literal asterisks\*</pre>	*literal asterisks*
**bold** bold	bold
~~strikethrough~~	strikethrough
1.⊔First item 2.⊔Second item ⊔1.⊔Subitem	<ol> <li>First item</li> <li>Second item         <ul> <li>A. Subitem</li> </ul> </li> </ol>
*uItem 1 uIndent	Item 1     Indent
-uItem 2	Item 2
u+uItem 3	■ Item 3
- [x] Done - [ ] To do	<ul> <li>☑ Done</li> <li>□To do</li> </ul>
A Lineuu	A Line
Break	Break
* * *	

<a id="anchor"></a> [Go to anchor](#anchor) #⊔Top Header [Go to header](#Top-Header)	<u>Go to anchor</u>
<pre>https://sqlbak.com [Link](https://sqlbak.com "optional title") Click [here][id] [id]:https://sqlbak.com</pre>	<u>Link</u>
> blockquote text	blockquote text
<pre>```python print('hello'); ```` `inline_code();`</pre>	<pre>print('hello');</pre>
Left  Center Right   : ::   1  A  C    2  B  D	LeftCenterRight1AC2BD
<pre>![alt text](logo.png "Title") ![][id] [id]:logo.png "Title"</pre>	Ø
<pre>\$\$\sqrt{k}\$\$ Inline: \$\sqrt{k}\$</pre>	$\sqrt{k}$
<pre>[![Img Alt Text](http://img.youtube.com/vi/ aZCXOw707nc/0.jpg)](https://yout u.be/aZCXOw707nc "Video Title")</pre>	YouTube

### Cheatography

### Natural Language Processing with Python & nltk Cheat Sheet by RJ Murray (murenei) via cheatography.com/58736/cs/15485/

Hendline Text			Contours Dension	
Handling Text			Sentence Parsing	
text='Some words'	assign string		g=nltk.data.load('grammar.c	fg') Load a grammar from a file
list(text)	Split text into	character tokens	g=nltk.CFG.fromstring("""	. " " Manually define grammar
set(text)	Unique token	S	")	
len(text)	Number of ch	naracters	<pre>parser=nltk.ChartParser(g)</pre>	Create a parser out of the grammar
Accessing corpora and lexic	al resources		trees=parser.parse_all(text	.)
from nltk.corpus import	import Corp	usReader object	for tree in trees: prin	t tree
brown		·	from nltk.corpus import treebank	
brown.words(text_id)	Returns pre of words	tokenised document as list	<pre>treebank.parsed_sents('wsj_ 1.mrg')</pre>	000 Treebank parsed sentences
brown.fileids()		n Brown corpus	i.mig /	
brown.categories() Lists categories in Brown corpus		Text Classification		
			from sklearn.feature_extrac	
Tokenization			CountVectorizer, TfidfVecto	rizer
text.split(" ")	Split b	y space	<pre>vect=CountVectorizer().fit(</pre>	_ 6
<pre>nltk.word_tokenizer(tex</pre>	t) nltk in-	built word tokenizer	ain)	data
<pre>nltk.sent_tokenize(doc)</pre>	nltk in-	built sentence tokenizer	<pre>vect.get_feature_names()</pre>	Get features
			<pre>vect.transform(X_train)</pre>	Convert to doc-term matrix
Lemmatization & Stemming			Entity Recognition (Chunking/Chi	nkina)
input="List listed list	s listing	Different suffixes	<pre>g="NP: {<dt>?<jj>*<nn>}"</nn></jj></dt></pre>	Regex chunk grammar
listings"				
words=input.lower().spl	lit(' ')	Normalize (lowercase) words	cp=nltk.RegexpParser(g)	Parse grammar
porter=nltk.PorterStemm	ner	Initialise Stemmer	ch=cp.parse(pos_sent)	Parse tagged sent. using grammar
[porter.stem(t) for t i	n wordsl	Create list of stems	print(ch)	Show chunks
WNL=nltk.WordNetLemmatizer()		Initialise WordNet	ch.draw()	Show chunks in IOB tree
		lemmatizer	cp.evaluate(test_sents)	Evaluate against test doc
[WNL.lemmatize(t) for t	in words]	Use the lemmatizer	<pre>sents=nltk.corpus.treebank.</pre>	tagged_sents()
			<pre>print(nltk.ne_chunk(sent))</pre>	Print chunk tree
Part of Speech (POS) Taggin	g			
<pre>nltk.help.upenn_tagset</pre>	Lookup defini	tion for a POS tag		
('MD')				

nltk.pos\_tag(words)

nltk in-built POS tagger

<use an alternative tagger to illustrate ambiguity>

By **RJ Murray** (murenei) cheatography.com/murenei/ tutify.com.au Published 28th May, 2018. Last updated 29th May, 2018. Page 1 of 2. Sponsored by CrosswordCheats.com Learn to solve cryptic crosswords! http://crosswordcheats.com

### Cheatography

RegEx with Pandas & Named Groups
<pre>df=pd.DataFrame(time_sents, columns=['text'])</pre>
<pre>df['text'].str.split().str.len()</pre>
<pre>df['text'].str.contains('word')</pre>
df['text'].str.count(r'\d')
<pre>df['text'].str.findall(r'\d')</pre>
<pre>df['text'].str.replace(r'\w+day\b', '???')</pre>
<pre>df['text'].str.replace(r'(\w)', lambda x: x.groups() [0][:3])</pre>
df['text'].str.extract(r'(\d?\d):(\d\d)')
<pre>df['text'].str.extractall(r'((\d?\d):(\d\d) ? ([ap]m))')</pre>
<pre>df['text'].str.extractall(r'(?P<digits>\d)')</digits></pre>



By **RJ Murray** (murenei) cheatography.com/murenei/ tutify.com.au

Published 28th May, 2018. Last updated 29th May, 2018. Page 2 of 2. Sponsored by CrosswordCheats.com Learn to solve cryptic crosswords! http://crosswordcheats.com

### Cheatography

### spaCy Cheat Sheet by Nuozhi via cheatography.com/122797/cs/22963/

#### Init

from spacy.lang.en import
English
nlp = English()

#### Basic

doc = nlp("SOME TEXTS")
span = doc[i:j]
token = doc[i]

### **Pre-trained Model**

nlp =
spacy.load('en\_core\_web\_sm')
doc = nlp(MY\_TEXT)

#### Name entity

doc.ents

.text

.label

spacy.	tokens
Doc	Doc(nlp.vocab, words=-
	words, spaces = spaces)
Span	Span(doc, i, j, label="-
	PERSON")

index: i, j

*words*: a collection of words *spaces*: a collecture of booleans

#### Matcher

matcher =

```
spacy.matcher.Matcher(nlp.vocab)
matches = matcher(doc)
```

[(id, start, end)]

#### Add pattern to matcher

```
pattern = [ { key: value } ]
matcher.add("PATTERN_NAME",
None, pattern)
```

Two types of key:

- 1. regex pattern
- 2. label (i.e. POS, entity)



### By Nuozhi cheatography.com/nuozhi/

#### **Phrase matching**

#### matcher =

```
spacy.matcher.PhraseMatcher(nlp.vocab)
pattern = nlp("Golden Retriever")
matcher.add("DOG", None, pattern)
for match_id, start, end in matche-
r(doc):
```

span = doc[start:end]

# Similarity token.vector source of the sector of the sector

	Pipeline
Text	Text - tokenizer tageer ner - Doc

nlp.pipe\_names
nlp.pipeline

### Add pipeline component

```
def fn(doc):
    # function body
    return doc
nlp.add_pipe(fn, last, first, before,
after)
```

### Set custom attributes add doc.\_.ATTR = "ATTRIBUTE

metadata NAME"
register Doc.set\_extension("ATTR",
globally default=None)

set to doc, tokens, spans access property via .\_\_

Not published yet. Last updated 24th May, 2020. Page 1 of 1.

#### **Extension attribute types**

attribute	Token.set_extensi-
	on("ATTR", defaut-
	=Bool)
property	Span.set_extensio-
	n("PROP", getter=fn)
method	Doc.set_extension-
	("METHOD", method-
	=fn)

#### Boost up

nlp.pipe(DATA)

#### **Passing in context**

data = [ ("SOME TEXTS",
{"KEY": "VAL"}), (...), ]
# Method 1
for doc, ctx in nlp.pipe(data, as\_tuple=True):
 print( doc.ATTR,
 ctx[KEY] )
# Method 2
Doc.set\_extension("KEY",
default=None)
for doc, ctx in nlp.pipe(data, as\_tuples=True):
 doc.\_.KEY = ctx["KEY"]

### Using tokenizer only

```
# Method 1
doc = nlp.make_doc("SOME
TEXTS")
# Method 2
with nlp.disable_pipes("tag-
ger", "parser"):
    doc = nlp(text)
```

Sponsored by **ApolloPad.com** Everyone has a novel in them. Finish Yours! https://apollopad.com

### matple

#### Ouick start

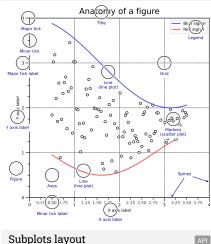
import numpy as np import matplotlib as mpl import matplotlib.pyplot a

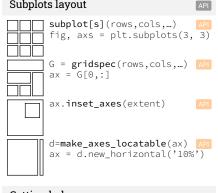
X = np.linspace(0, 2\*np.p Y = np.cos(X)

fig, ax = plt.subplots() ax.plot(X, Y, color='green

fig.savefig("figure.pdf") fig.show()

#### Anatomy of a figure





- Getting help
- matplotlib.org
- **Q** github.com/matplotlib/matplotlib/issues
- **D** discourse.matplotlib.org
- ▲ stackoverflow.com/questions/tagged/matplotlib ₩ gitter.im/matplotlib
- ♥ twitter.com/matplotlib
- Matplotlib users mailing list

+lih	Basic p	lots
Version 3.5.0		<pre>plot([X],Y,[fmt],) X,Y,fmt, color, marker, linestyle</pre>
API		scatter(X,Y,) X,Y, [s]izes, [c]olors, marker, cmap
as plt		bar[h] (x,height,)x, height, width, bottom, align, color
oi, 100)		imshow(Z,) Z, cmap, interpolation, extent, origin
en')		<b>contour[f]</b> ([X],[Y],Z,) AF X, Y, <b>Z</b> , levels, colors, extent, origin
		<b>pcolormesh</b> ([X],[Y],Z,) X, Y, <b>Z</b> , vmin, vmax, cmap
Bli (signal Re(signal Legend		<b>quiver</b> ([X],[Y],U,V,) AF X, Y, <b>U</b> , <b>V</b> , C, units, angles
Grid		pie(X,) Z, explode, labels, colors, radius
\$ \$ \$ \$ \$ \$	TEXT	text(x,y,text,) AF x, y, text, va, ha, size, weight, transform
o o o (scatter plot)		fill[_between][x]()AFX, Y1, Y2, color, where
Spines		

### Advanced plots

step(X,Y,[fmt],...) X, Y, fmt, color, marker, where boxplot(X,...) X, notch, sym, bootstrap, widths errorbar(X,Y,xerr,yerr,...) API X, Y, xerr, yerr, fmt hist(X, bins, ...)

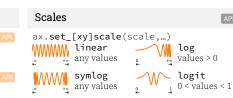
> **X**, bins, range, density, weights violinplot(D,...)

D, positions, widths, vert barbs([X],[Y], U, V, ...)

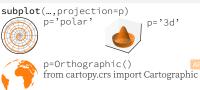
X, Y, U, V, C, length, pivot, sizes

eventplot(positions,...) positions, orientation, lineoffsets

hexbin(X,Y,C,...) X, Y, C, gridsize, bins



#### Projections



### Lines

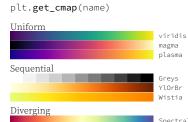
API

linestyle or ls .... n\_ ... (0, (0, 01, 2))capstyle or dash\_capstyle "butt' "round "projecting"

Ma	irke	rs									API
	0		÷	83	☆	$\bigcirc$	$\diamond$	$\triangleleft$	$\triangleright$	$\triangle$	$\nabla$
1.1	'0'	's'	'P'	'X'	'*'	'p'	'D'	'<'	'>'	1.4.1	'v'
$\mathbf{Y}$	+	$\prec$	$\succ$	+	$\times$	1	—	<	>	$\wedge$	$\sim$
'1'	'2'	'3'	'4'	'+'	'x'	2P	1.1	4	5	6	7
	٠	•	٠	$\rightarrow$	←	1	4	${}^{\bullet}$	0	$\Theta$	$\bigcirc$
'\$ <b></b> \$\$'	'\$ <b>*</b> \$'	'\$♥\$'	\$\$\$'	'\$→\$'	'\$←\$'	'\$1\$'	'\$↓\$'	\$0\$	\$0\$	'\$ <del>Q</del> \$'	'\$ <del>⊜</del> \$ '
mar	kev	erv									
~		÷.	~			_		-00	<u> </u>		
	10			[0, -1]			(25, 5)		[0	, 25, -	1]

Co	lor	s								
C⊗	C1	C2	C3	C4	C5	C6	C7	C8	C9	'Cn'
ь	g					у		k	w	,×,
Dark	Red			k Crii	nson		nRed	Sal		'name'
	0,0)						0			(R,G,B[,A])
#FF	6986									'#RRGGBB[AA
6.0	3.1	0.2	6.3	6.4 6	.5 0.				1.0	'x.v'

#### Colormaps



#### Spectral coolwarm RdGv Oualitative tab10

tah20 Cyclic twilight

#### **Tick locators**

from matplotlib import ticker ax.[xy]axis.set [minor|major] locator(locator)

#### ticker.NullLocator() ticker.MultipleLocator(0.5)

0.0 0.5 1.0 1.5 2.0 2.5 3.0 3.5 4.0 4.5 ticker.FixedLocator([0, 1, 5]) ticker.LinearLocator(numticks=3) 2.5 ticker.IndexLocator(base=0.5, offset=0.25) 4.25 ticker.AutoLocator() ticker.MaxNLocator(n=4) 1.5 ticker.LogLocator(base=10, numticks=15)

#### Tick formatters

API

API

API

from matplotlib import ticker ax.[xy]axis.set\_[minor|major]\_formatter(formatter)

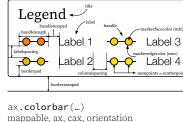
ticker.NullFormatter()

tick	er.EixedFor	matter(['	', '0', '1'	1)	
6 0	25 0.50 1 0.75	0.25 2 0.50 0.	75 3 0.25 0.50	0.75 4 0.	25 0.50 5
tick	er.FuncForm	natter(lam	bda x, pos:	"[%.2f]	<u>" % x</u> )
.00]	[1.00]	[2.00]	[3.00]	[4.00]	[5.00
tick	er.FormatSt	rFormatte	r('>%d<')		
-0<	>1<	>2<	>3<	>4<	>5+
tick	er.ScalarFo	rmatter()			
ò	i	ż	3	4	5
tick	er.StrMetho	dFormatte	r('{x}')		
0.0	1.0	2.0	3.0	4.0	5.0

ticker.PercentFormatter(xmax=5)

#### Ornaments

ax.legend(...) handles, labels, loc, title, frameon



0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9

ax.annotate(...) text, xy, xytext, xycoords, textcoords, arrowprops

text xytext textcoords

#### Event handling

fig, ax = plt.subplots() def on\_click(event): print(event) fig.canvas.mpl\_connect( 'button\_press\_event', on\_click)

#### Animation

API

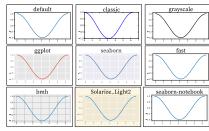
API

import matplotlib.animation as mpla

T = np.linspace(0, 2\*np.pi, 100) S = np.sin(T)line, = plt.plot(T, S) def animate(i): line.set\_ydata(np.sin(T+i/50)) anim = mpla.FuncAnimation( plt.gcf(), animate, interval=5) plt.show()

#### Styles

plt.style.use(style)



ax.patch.set\_alpha(0) ax.set\_[xy]lim(vmin, vmax) ax.set\_[xy]label(label) ax.set\_[xy]ticks(list) ax.set\_[xy]ticklabels(list) ax.set\_[sup]title(title) ax.tick\_params(width=10, ...) ax.set\_axis\_[on|off]()

plt.gcf(), plt.gca()
mpl.rc('axes', linewidth=1, ...) fig.patch.set alpha(0) text=r'\$\frac{-e^{i\pi}}{2^n}\$'

#### Keyboard shortcuts

ctrl + s Save	ctrl + w Close plot
r Reset view	f Fullscreen 0/1
f View forward	b View back
P Pan view	<ul> <li>Zoom to rect</li> </ul>
🗙 X pan/zoom	y Y pan/zoom
g Minor grid 0/1	G Major grid 0/1
I X axis log/linear	L Y axis log/linear

#### Ten simple rules

.

xy xycoords

API

- 1. Know Your Audience
- 2. Identify Your Message 3. Adapt the Figure
- 4. Captions Are Not Optional
- 5. Do Not Trust the Defaults
- 6. Use Color Effectively
- 7. Do Not Mislead the Reader
- 8. Avoid "Chartiunk"
- 9. Message Trumps Beauty
- 10. Get the Right Tool

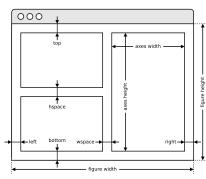
## Quick reminder

ax.grid()

fig.tight\_layout()

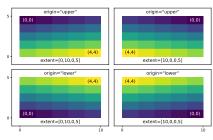
#### Axes adjustments

#### plt.subplots\_adjust( ... )

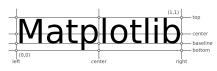


#### Extent & origin

ax.imshow( extent=..., origin=... )







#### Text parameters

ax.text(..., family=..., size=..., weight=...) ax.text(..., fontproperties=...)

The quick brown fox	xx-large	(1.73)
The quick brown fox	x-large	(1.44)
The quick brown fox	large	(1.20)
The quick brown fox	medium	(1.00)
The quick brown fox	small	(0.83)
The quick brown fox	x-small	(0.69)
The quick brown fox	xx-small	(0.58)

#### The quick brown fox jumps over the lazy dog The quick brown fox jumps over the lazy dog The quick brown fox jumps over the lazy dog The quick brown fox jumps over the lazy dog The quick brown fox jumps over the lazy dog

The sudah beaus for dimonstructure the last day	
The quick brown fox jumps over the lazy dog	monospace
The quick brown fox jumps over the lazy dog	serif
The quick brown fox jumps over the lazy dog	sans
The quick brown fox jumps over the lazy dog	cursive
The quick brown fox jumps over the lazy dog	italic

The quick brown fox jumps over the lazy dog The quick brown fox jumps over the lazy dog

THE QUICK BROWN FOX JUMPS OVER THE LAZY DOG The quick brown fox jumps over the lazy dog

# Uniform colormaps

API

API

API

API

black (900)

bold (700)

normal

normal

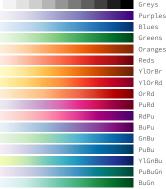
small-caps

semibold (600)

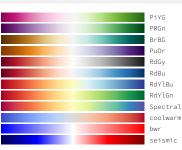
ultralight (100)

normal (400)

### Sequential colormaps



#### **Diverging colormaps**



#### Qualitative colormaps

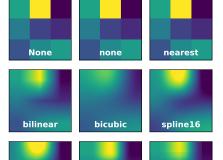
Pastel1
Pastel2
Paired
Accent
Dark2
Set1
Set2
Set3
tab10
tab20
tab20b
tab20c

#### Miscellaneous colormaps



darksalmon     darksgreen     mediumpurple       coral     g     mediumpurple       orangered     green     blueviolet       lightsalmon     imme     midgo       ispitsalmon     imme     midgo       saashell     mediumseagreen     darkviolet       chocolate     springgreen     mediumscholet       sandybrown     mediumscholet     purple       peachpuff     mediumaquamarine     violet       bique     lightsagreen     mediumscholet       bique     gaturquoise     fichtsia       bique     gaturquoise     mediumscholet       bightsagreen     mediumscholet     mediumscholet       bightsagreen     mediumscholet     mediumscholet       bightsagreen     mediumscholet     mediumscholet       bightsagreen     mediumscholet <th>Color names</th> <th></th> <th>API</th>	Color names		API
papayawhip teal lavenderblush moccasin darkcyan palevioletred	<ul> <li>black</li> <li>dimgray</li> <li>dimgray</li> <li>dimgrey</li> <li>gray</li> <li>grey</li> <li>grey<td>darkgöldernod goldernod cornsilk politik politik darkkhaki darkkhaki beige lightyellow lightyellow lightyellow lightyellow darkölwegreen darkölwegreen darkölwegreen darkölwegreen honeydew darkölwegreen honeydew darksegreen lightgreen lightgreen darksegreen mediumsegreen darksegreen lime segreen darksegreen mediumsegreen darksegreen da</td><td><ul> <li>cadetbiue</li> <li>powderbiue</li> <li>ightbiue</li> <li>ightbiue</li> <li>ightbiue</li> <li>ightbiue</li> <li>ightbiue</li> <li>ightbiue</li> <li>ightbiategray</li> <li>islategray</li> <li>islategray</li> <li>islategray</li> <li>islategray</li> <li>islategray</li> <li>islatebiue</li> <li>darkbiue</li> <li>darkbiue</li></ul></td></li></ul>	darkgöldernod goldernod cornsilk politik politik darkkhaki darkkhaki beige lightyellow lightyellow lightyellow lightyellow darkölwegreen darkölwegreen darkölwegreen darkölwegreen honeydew darkölwegreen honeydew darksegreen lightgreen lightgreen darksegreen mediumsegreen darksegreen lime segreen darksegreen mediumsegreen darksegreen da	<ul> <li>cadetbiue</li> <li>powderbiue</li> <li>ightbiue</li> <li>ightbiue</li> <li>ightbiue</li> <li>ightbiue</li> <li>ightbiue</li> <li>ightbiue</li> <li>ightbiategray</li> <li>islategray</li> <li>islategray</li> <li>islategray</li> <li>islategray</li> <li>islategray</li> <li>islatebiue</li> <li>darkbiue</li> <li>darkbiue</li></ul>
vheat aqua pink oldiace cyan lightpink	papayawhip moccasin orange wheat	teal darkcyan c aqua	lavenderblush palevioletred crimson pink

Image int	terpola	tion	











#### Legend placement K J 9 1 Α 6 10 7 οH В 8 G C Ē D F ax.legend(loc="string", bbox\_to\_anchor=(x,y))

10: center

Annotation connection styles

9: upper center 1: upper right

8: lower center 4: lower right

B: center right / (-0.1,0.5)

D: upper left / (0.1,-0.1)

F: upper right / (0.9,-0.1)

H: center left / (1.1,0.5)

J: lower right / (0.9,1.1)

L: lower left / (0.1,1.1)

7: center right

API

2: upper left

6: center left

3: lower left

API

A: upper right / (-0.1,0.9)

C: lower right / (-0.1,0.1)

G: lower left / (1.1,0.1)

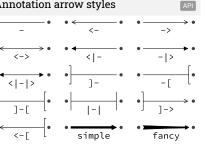
I: upper left / (1.1,0.9)

K: lower center / (0.5,1.1)

E: upper center / (0.5,-0.1)

arc3, rad=0	arc3, rad=0.3	angle3, angle4, angle8=90
anglo, anglo, anglo≣188, rad=0	angles-90, angles-90, angle5-180, rad-25	arc[eAz=96, ang[eAz=96, argAz=0, argAz=0, rad=0, rad=0,
bar, fraction=0.3	bar, fraction=-0.3	bar, angle=189, fraction=-0.2

#### Annotation arrow styles



#### How do I ..

- ... resize a figure?  $\rightarrow$  fig.set\_size\_inches(w, h)
- ... save a figure?  $\rightarrow$  fig.savefig("figure.pdf")
- ... save a transparent figure?
- → fig.savefig("figure.pdf", transparent=True) ... clear a figure/an axes?
- $\rightarrow$  fig.clear()  $\rightarrow$  ax.clear()
- ... close all figures?  $\rightarrow$  plt.close("all")
- ... remove ticks?
- $\rightarrow$  ax.set\_[xy]ticks([]) ... remove tick labels ?
- $\rightarrow$  ax.set\_[xy]ticklabels([]) ... rotate tick labels ?
- $\rightarrow$  ax.set\_[xy]ticks(rotation=90) ... hide top spine?
- $\rightarrow$  ax.spines['top'].set\_visible(False) ... hide legend border?
- $\rightarrow$  ax.legend(frameon=False)
- ... show error as shaded region?  $\rightarrow$  ax.fill\_between(X, Y+error, Y-error)
- ... draw a rectangle?
- $\rightarrow$  ax.add\_patch(plt.Rectangle((0, 0), 1, 1)) ... draw a vertical line?
- $\rightarrow$  ax.axvline(x=0.5) ... draw outside frame?
- $\rightarrow$  ax.plot(..., clip\_on=False) ... use transparency?
- $\rightarrow$  ax.plot(..., alpha=0.25) ... convert an RGB image into a gray image?
- $\rightarrow$  grav = 0.2989\*R + 0.5870\*G + 0.1140\*B
- ... set figure background color? → fig.patch.set\_facecolor("grey") ... get a reversed colormap?
- $\rightarrow$  plt.get\_cmap("viridis\_r")
- ... get a discrete colormap?
- $\rightarrow$  plt.get\_cmap("viridis", 10) ... show a figure for one second?
- $\rightarrow$  fig.show(block=False), time.sleep(1)

#### Performance tips

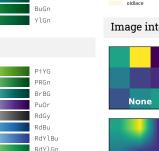
<mark>catter(X, Y)</mark>	<mark>slow</mark>
lot(X, Y, marker="o", ls="")	fast
<pre>or i in range(n): plot(X[i]) lot(sum([x+[None] for x in X],[]))</pre>	<mark>slow</mark> fast
la(), imshow(…), canvas.draw()	<mark>slow</mark>
m.set_data(…), canvas.draw()	fast

### **Beyond Matplotlib**

Seaborn: Statistical Data Visualization **Cartopy**: Geospatial Data Processing yt: Volumetric data Visualization mpld3: Bringing Matplotlib to the browser Datashader: Large data processing pipeline plotnine: A Grammar of Graphics for Python

Matplotlib Cheatsheets Copyright (c) 2021 Matplotlib Development Team Released under a CC-BY 4.0 International License





viridis

plasma

inferno

cividis

magma

### Matplotlib for beginners

Matplotlib is a library for making 2D plots in Python. It is designed with the philosophy that you should be able to create simple plots with just a few commands:

### 1 Initialize

import numpy as np import matplotlib.pyplot as plt

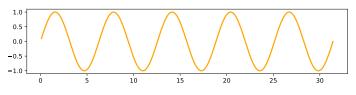
### 2 Prepare

X = np.linspace(0, 4\*np.pi, 1000) Y = np.sin(X)

### 3 Render

```
fig, ax = plt.subplots()
ax.plot(X, Y)
fig.show()
```

### 4 Observe



### Choose

Matplotlib offers several kind of plots (see Gallery):

```
X = np.random.uniform(0, 1, 100)
Y = np.random.uniform(0, 1, 100)
ax.scatter(X, Y)
```

X = np.arange(10)Y = np.random.uniform(1, 10, 10)ax.bar(X, Y)

Z = np.random.uniform(0, 1, (8,8))

ax.**imshow**(Z)

```
Z = np.random.uniform(0, 1, (8,8))
```

```
ax.contourf(Z)
```

Z = np.random.uniform(0, 1, 4)

ax.**pie**(Z)

Z = np.random.normal(0, 1, 100)

### ax.**hist**(Z)

X = np.arange(5)Y = np.random.uniform(0, 1, 5)ax.errorbar(X, Y, Y/4)

Z = np.random.normal(0, 1, (100,3))

ax.boxplot(Z)

### Tweak

You can modify pretty much anything in a plot, including limits, colors, markers, line width and styles, ticks and ticks labels, titles, etc.

```
X = np.linspace(0, 10, 100)
Y = np.sin(X)
ax.plot(X, Y, color="black")
```

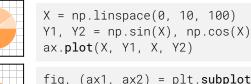
```
X = np.linspace(0, 10, 100)
Y = np.sin(X)
ax.plot(X, Y, linestyle="--")
```

X = np.linspace(0, 10, 100)Y = np.sin(X)ax.plot(X, Y, linewidth=5)

X = np.linspace(0, 10, 100)Y = np.sin(X)ax.plot(X, Y, marker="o")

### Organize

You can plot several data on the the same figure, but you can also split a figure in several subplots (named Axes):





fig, (ax1, ax2) = plt.subplots((2,1))ax1.plot(X, Y1, color="C1") ax2.plot(X, Y2, color="C0")



### Label (everything)

ax.plot(X, Y) fig.suptitle(None) ax.set\_title("A Sine wave")



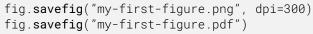
ax.plot(X, Y) ax.set vlabel(None) ax.set\_xlabel("Time")



### Explore

Figures are shown with a graphical user interface that allows to zoom and pan the figure, to navigate between the different views and to show the value under the mouse

**Save** (bitmap or vector format)





Matplotlib 3.5.0 handout for beginners. Copyright (c) 2021 Matplotlib Development Team. Released under a CC-BY 4.0 International License. Supported by NumFOCUS.





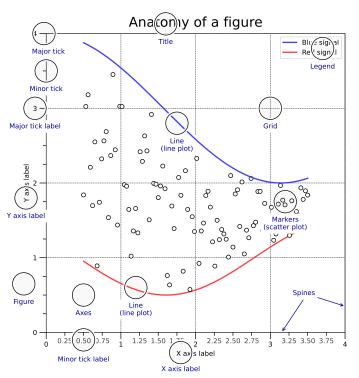






### Matplotlib for intermediate users

A matplotlib figure is composed of a hierarchy of elements that forms the actual figure. Each element can be modified.



### Figure, axes & spines



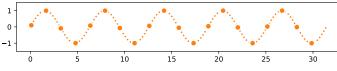
### Ticks & labels

from mpl.ticker import MultipleLocator as ML
from mpl.ticker import ScalarFormatter as SF
ax.xaxis.set\_minor\_locator(ML(0.2))
ax.xaxis.set\_minor\_formatter(SF())
ax.tick\_params(axis='x',which='minor',rotation=90)

### Lines & markers

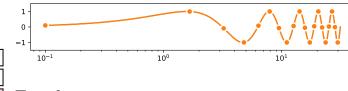
X = np.linspace(0.1, 10\*np.pi, 1000) Y = np.sin(X)

ax.plot(X, Y, "C1o:", markevery=25, mec="1.0")

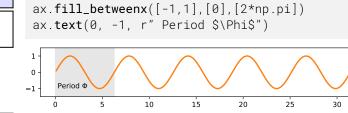


### **Scales & projections**

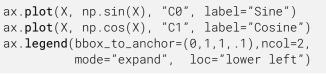
```
fig, ax = plt.subplots()
ax.set_xscale("log")
ax.plot(X, Y, "Clo-", markevery=25, mec="1.0")
```

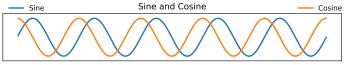


### Text & ornaments



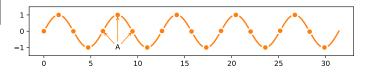
### Legend





### Annotation

ax.annotate("A", (X[250],Y[250]),(X[250],-1), ha="center", va="center",arrowprops = {"arrowstyle" : "->", "color": "C1"})



### Colors

 Any color can be used, but Matplotlib offers sets of colors:

 c0
 C1
 C2
 C3
 C4
 C5
 C6
 C7
 C8
 C9

 0.0
 0.1
 0.2
 0.3
 0.4
 0.5
 0.6
 0.7
 0.8
 0.9
 1.0

### Size & DPI

Consider a square figure to be included in a two-columns A4 paper with 2cm margins on each side and a column separation of 1cm. The width of a figure is (21 - 2\*2 - 1)/2 = 8cm. One inch being 2.54cm, figure size should be  $3.15 \times 3.15$  in.

fig = plt.figure(figsize=(3.15,3.15), dpi=50)
plt.savefig("figure.pdf", dpi=600)

Matplotlib 3.5.0 handout for intermediate users. Copyright (c) 2021 Matplotlib Development Team. Released under a CC-BY 4.0 International License. Supported by NumFOCUS.

### Matplotlib tips & tricks

### Transparency

Scatter plots can be enhanced by using transparency (alpha) in order to show area with higher density. Multiple scatter plots can be used to delineate a frontier.

X = np.random.normal(-1, 1, 500) Y = np.random.normal(-1, 1, 500) ax.scatter(X, Y, 50, "0.0", lw=2) # optional ax.scatter(X, Y, 50, "1.0", lw=0) # optional ax.scatter(X, Y, 40, "C1", lw=0, alpha=0.1)



### Rasterization

If your figure has many graphical elements, such as a huge scatter, you can rasterize them to save memory and keep other elements in vector format.

```
X = np.random.normal(-1, 1, 10_000)
Y = np.random.normal(-1, 1, 10_000)
ax.scatter(X, Y, rasterized=True)
fig.savefig("rasterized-figure.pdf", dpi=600)
```

### Offline rendering

Use the Agg backend to render a figure directly in an array.

```
from matplotlib.backends.backend_agg import FigureCanvas
canvas = FigureCanvas(Figure()))
... # draw som stuff
canvas.draw()
Z = np.array(canvas.renderer.buffer_rgba())
```

### **Range of continuous colors**

You can use colormap to pick from a range of continuous colors.

X = np.random.randn(1000, 4)
cmap = plt.get\_cmap("Oranges")
colors = cmap([0.2, 0.4, 0.6, 0.8])

ax.hist(X, 2, histtype='bar', color=colors)

### Text outline

Use text outline to make text more visible.

```
import matplotlib.patheffects as fx
text = ax.text(0.5, 0.1, "Label")
text.set_path_effects([
    fx.Stroke(linewidth=3, foreground='1.0'),
    fx.Normal()])
```

### Multiline plot

You can plot several lines at once using None as separator.

```
X,Y = [], []
for x in np.linspace(0, 10*np.pi, 100):
```

```
X.extend([x, x, None]), Y.extend([0, sin(x), None])
ax.plot(X, Y, "black")
```



### **Dotted lines**

To have rounded dotted lines, use a custom linestyle and modify dash\_capstyle.

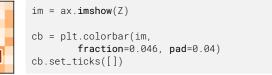
### 

### **Combining axes**

You can use overlaid axes with different projections.

### Colorbar adjustment

You can adjust a colorbar's size when adding it.





### Taking advantage of typography

You can use a condensed font such as Roboto Condensed to save space on tick labels.

```
for tick in ax.get_xticklabels(which='both'):
    tick.set_fontname("Roboto Condensed")
```

### 0 0.2 0.4 0.6 0.8 1 1.2 1.4 1.6 1.8 2 2.2 2.4 2.6 2.8 3 3.2 3.4 3.6 3.8 4 4.2 4.4 4.6 4.8 5

### Getting rid of margins

Once your figure is finished, you can call tight\_layout() to remove white margins. If there are remaining margins, you can use the pdfcrop utility (comes with TeX live).

### Hatching

You can achieve a nice visual effect with thick hatch patterns.

```
cmap = plt.get_cmap("Oranges")
plt.rcParams['hatch.color'] = cmap(0.2)
plt.rcParams['hatch.linewidth'] = 8
ax.bar(X, Y, color=cmap(0.6), hatch="/")
```



### **Read the documentation**

Matplotlib comes with an extensive documentation explaining the details of each command and is generally accompanied by examples. Together with the huge online gallery, this documentation is a gold-mine.

Matplotlib 3.5.0 handout for tips & tricks. Copyright (c) 2021 Matplotlib Development Team. Released under a CC-BY 4.0 International License. Supported by NumFOCUS.



Labe



### 6.2. R language

### Base R Cheat Sheet

	C(2, 4, 0)
Getting Help	2:6
Accessing the help files	seq(2, 3, by
<b>?mean</b> Get help of a particular function.	rep(1:2, tim
help.search('weighted mean') Search the help files for a word or phrase.	rep(1:2, eac
<pre>help(package = 'dplyr') Find help for a package.</pre>	
More about an object	<pre>sort(x)</pre>
<pre>str(iris) Get a summary of an object's structure. class(iris)</pre>	Return x sorte <b>table(x)</b> See counts o
Find the class an object belongs to.	Sele
Using Packages	
<b>install.packages('dplyr')</b> Download and install a package from CRAN.	x [ 4
<b>library(dplyr)</b> Load the package into the session, making all its functions available to use.	x[-4 x[2:
<b>dplyr::select</b> Use a particular function from a package.	x[-( <mark>2</mark> :
<b>data(iris)</b> Load a built-in dataset into the environment.	x[c(1,
Working Directory	x[x ==
getwd() Find the current working directory (where	x[x <
Find the current working directory (where inputs are found and outputs are sent).	x[x %
<pre>setwd('C://file/path') Change the current working directory.</pre>	c(1, 2,

Use projects in RStudio to set the working directory to the folder you are working in.

Vactors					
	Vectors Creating Vectors				
c(2, 4, 6)	246	Join elements into a vector	for (va		
2:6	23456	An integer sequence	Do so		
seq(2, 3, by=0.5)	2.0 2.5 3.0	A complex sequence			
rep(1:2, times=3)	121212	Repeat a vector	for (i j <-		
rep(1:2, each=3)	1 1 1 2 2 2	Repeat elements of a vector	prin <sup>.</sup> }		
Vecto	or Function	s			
<pre>sort(x) Return x sorted. table(x) See counts of values</pre>	<b>rev(x)</b> Return x re <b>unique(</b> . See uniqu	(x)	if (con Do so } else Do so		
Selecting	Vector Eler	nents	}		
By Positionx[4]The fourth element.x[-4]All but the fourth.			if (i > print } else - print		
x[-4] x[2:4]	}				
x[-(2:4)]	All element two to		Input		
x[c(1, 5)]	Elements of five	df <- r			
В	y Value				
x[x == 10]		ts which Ial to 10.	df <-		
x[x < 0]		ents less zero.			
x[x %in% c(1, 2, 5)]	Elements 1, 2	lo			
Nan	ned Vectors				
x['apple']	Elemer name 'a		Condition		

Programming							
For Loop	Tiog		5	Whil		20	
-						ob	
<pre>for (variable in sequence){</pre>			e ( <mark>condit</mark>				
Do something			o somethi	ng			
}		}					
Example					ample	9	
for (i in 1:4){			e(i < 5)	1			
j <- i + 10			rint(i)				
<pre>print(j)</pre>			<- i + 1				
}		}					
If Statements				Fun	ctior	IS	
if (condition){		func	tion_name	e <- f	unct	ion(var){	
<pre>Do something } else {</pre>		D	o somethi	ing			
Do something different }		r }	eturn(new	v_vari	able)	)	
Example				Exa	ample	9	
if (i > 3){		squa	re <- fun	oction	(x){		
<pre>print('Yes') } else {</pre>		S	quared <-	- x*x			
<pre>print('No')</pre>		r	eturn(squ	uared)			
}		}					
Reading and Writing Data Also see the readr package.							
							aonagor
Input	Ouput			De	escript	ion	
<pre>df &lt;- read.table('file.txt')</pre>	write.t	able(df,	file.txt	<mark>')</mark> R	Read an	d write a delir file.	nited text
					Pood	and write a co	mma
df <- read.csv('file.csv')	write.	csv(df, '	file.csv'	)	separa	ted value file.	This is a
					specia	al case of read write.table.	.table/
<pre>load('file.RData')</pre>	<pre>save(df,</pre>	file = 't	file.Rdata	a') <sup>F</sup>		nd write an R d type special f	
a == b Are equal	a > b	Greater than	a >= b	Greate or equ		is.na(a)	Is missing
Conditions a != b Not equal	a < b	Less than	a <= b	Less the	han or	is.null(a)	Is null
				Uque			

### Types

Converting between common data types in R. Can always go from a higher value in the table to a lower value.

as.logical	TRUE, FALSE, TRUE	Boolean values (TRUE or FALSE).
as.numeric	1, 0, 1	Integers or floating point numbers.
as.character	'1', '0', '1'	Character strings. Generally preferred to factors.
as.factor	'1', '0', '1', levels: '1', '0'	Character strings with preset levels. Needed for some statistical models.

### **Maths Functions**

log(x)	Natural log.	sum(x)	Sum.
exp(x)	Exponential.	mean(x)	Mean.
max(x)	Largest element.	<pre>median(x)</pre>	Median.
min(×)	Smallest element.	<pre>quantile(x)</pre>	Percentage quantiles.
round(x, n)	Round to n decimal places.	rank(x)	Rank of elements.
<pre>signif(x, n)</pre>	Round to n significant figures.	var(x)	The variance.
cor(x, y)	Correlation.	sd(x)	The standard deviation.

### Variable Assignment

<- 'apple' > a > a [1] 'apple'

	vironment
ls()	List all variables in the environment.
rm(x)	Remove x from the environment.
rm(list = ls())	Remove all variables from the environment.

Matrices	
<pre>m &lt;- matrix(x, nrow = 3, Create a matrix from x</pre>	
<b>m[2,</b> ] - Select a row	t(m) Transpose
m[, 1] - Select a column	<b>m %*% n</b> Matrix Multiplication
m[2, 3] - Select an element	<b>solve(m, n)</b> Findxin:m*x=n

### Lists

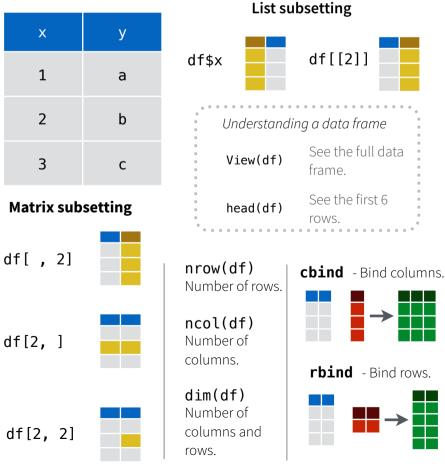
l <- list(x = 1:5, y = c('a', 'b'))</pre> A list is a collection of elements which can be of different types.

l[[2]]	l[1]	l\$x	l['y']
Second element of l.	New list with only the first element.	Element named x.	New list with only element named y.

### Also see the **dplyr** package.

df <- data.frame(x = 1:3, y = c('a', 'b', 'c'))</pre> A special case of a list where all elements are the same length.

**Data Frames** 



Strings	Also see the <b>stringr</b> package.
<pre>paste(x, y, sep = ' ')</pre>	Join multiple vectors together.
<pre>paste(x, collapse = ' ')</pre>	Join elements of a vector together.
<pre>grep(pattern, x)</pre>	Find regular expression matches in x.
<pre>sub(pattern, replace, x)</pre>	Replace matches in x with a string.
toupper(x)	Convert to uppercase.
tolower(x)	Convert to lowercase.
nchar(x)	Number of characters in a string.
Fact	ors
factor(x)	cut(x, breaks = 4)

Turn a vector into a factor. Can set the levels of the factor and the order.

q

cut(x, breaks Turn a numeric vector into a factor by 'cutting' into sections.

### **Statistics**

 $lm(y \sim x, data=df)$ prop.test t.test(x, y) Linear model. Perform a t-test for difference between  $glm(y \sim x, data=df)$ means. Generalised linear model. pairwise.t.test summary Perform a t-test for

Get more detailed information out a model.

proportions. paired data.

aov Analysis of variance.

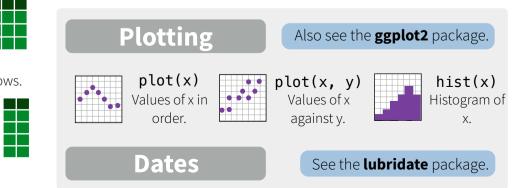
Test for a

difference

between

### **Distributions**

	Random Variates	Density Function	Cumulative Distribution	Quantile
Normal	rnorm	dnorm	pnorm	qnorm
Poisson	rpois	dpois	ppois	qpois
Binomial	rbinom	dbinom	pbinom	qbinom
Uniform	runif	dunif	punif	qunif



RStudio® is a trademark of RStudio, Inc. • CC BY Mhairi McNeill • mhairihmcneill@gmail.com • 844-448-1212 • rstudio.com

Learn more at web page or vignette • package version • Updated: 3/15

# **Data Import** with readr, tibble, and tidyr

Cheat Sheet



R's tidyverse is built around tidy data stored in tibbles. an enhanced version of a data frame.



The reverse side shows how to create tibbles with **tibble** and to layout tidy data with **tidyr**.

to read text files into R with **readr**.

The front side of this sheet shows how

### **Other types of data**

Try one of the following packages to import other types of files

- haven SPSS, Stata, and SAS files •
- readxl excel files (.xls and .xlsx) •
- **DBI** databases •
- isonlite ison •
- xml2 XMI •
- httr - Web APIs
- rvest HTML (Web Scraping) •

### Write functions

Save **x**, an R object, to **path**, a file path, with:

write\_csv(x, path, na = "NA", append = FALSE, col names = !append)

Tibble/df to comma delimited file.

write\_delim(x, path, delim = " ", na = "NA", append = FALSE, col\_names = !append) Tibble/df to file with any delimiter.

write\_excel\_csv(x, path, na = "NA", append = FALSE, col names = !append)

Tibble/df to a CSV for excel

write\_file(x, path, append = FALSE) String to file.

write lines(x, path, na = "NA", append = FALSE)

String vector to file, one element per line.

- write\_rds(x, path, compress = c("none", "gz", "bz2", "xz"), ...**)** Object to RDS file.
- write tsv(x, path, na = "NA", append = FALSE, col\_names = !append) Tibble/df to tab delimited files.

### **Read functions**

### **Read tabular data to tibbles**

Reads Semi-colon delimited files.

read delim("file.txt", delim = "|")

read\_csv2("file2.csv")

### These functions share the common arguments:

read\_\*(file, col\_names = TRUE, col\_types = NULL, locale = default\_locale(), na = c("", "NA"), quoted na = TRUE, comment = "", trim ws = TRUE, skip = 0, n max = Inf, guess max = min(1000, n max), progress = interactive())

read csv2()







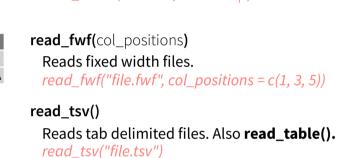


a,b,c

1,2,3

4,5,NA

A B C



read\_delim(delim, quote = "\"", escape\_backslash = FALSE,

escape double = TRUE) Reads files with any delimiter.

### **Useful arguments**



write\_csv (path = "file.csv", *x* = *read\_csv("a,b,c\n1,2,3\n4,5,NA")*)

#### No header 1 2 3 read\_csv("file.csv", 4 5 NA

col names = FALSE)

#### x y z **Provide header**

- A B C read csv("file.csv", 1 2 3
- col\_names = c("x", "y", "z")) 4 5 NA

### **Read non-tabular data**

read file(file, locale = default locale()) Read a file into a single string.

read\_file\_raw(file)

- Read a file into a raw vector.
- read\_lines(file, skip = 0, n\_max = -1L, locale = default\_locale(), na = character(), progress = interactive()) Read each line into its own string.

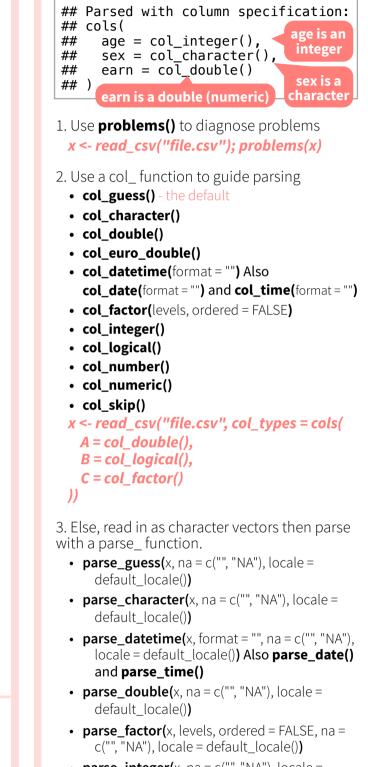
read\_lines\_raw(file, skip = 0, n\_max = -1L, progress = interactive()) Read each line into a raw vector.

**read log(**file, col names = FALSE, col types = NULL, skip = 0, n\_max = -1, progress = interactive()) Apache style log files.

### **Parsing data types**

readr functions guess the types of each column and convert types when appropriate (but will NOT convert strings to factors automatically).

A message shows the type of each column in the result.



- parse\_integer(x, na = c("", "NA"), locale = default locale())
- **parse\_logical(**x, na = c("", "NA"), locale = default\_locale())
- **parse\_number(**x, na = c("", "NA"), locale = default locale())
- x\$A <- parse\_number(x\$A)

Learn more at browseVignettes(package = c("readr", "tibble", "tidyr")) • readr 1.1.0 • tibble 1.2.12 • tidyr 0.6.0 • Updated: 2017-01

4 5 NA ABC Read in a subset 1 2 3 read\_csv("file.csv",

2 3

**Skip lines** 

**skip = 1**)

*n max* = 1)

**Missing Values** 

read csv("file.csv",

na = c("4", "5", "."))

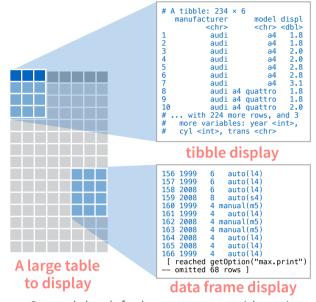
read\_csv("file.csv",



### Tibbles - an enhanced data frame

The **tibble** package provides a new S3 class for storing tabular data, the tibble. Tibbles inherit the data frame class, but improve two behaviors:

- Display When you print a tibble, R provides a concise view of the data that fits on one screen.
- **Subsetting** [ always returns a new tibble, [[ and \$ always return a vector.
- No partial matching You must use full column names when subsetting



- Control the default appearance with options:
  - options(tibble.print\_max = n, tibble.print\_min = m, tibble.width = Inf)
- View entire data set with **View**(x, title) or **glimpse(**x, width = NULL, ...)
- Revert to data frame with **as.data.frame()** (required for some older packages)

#### **Construct a tibble in two ways** tibble(...) Both make Construct by columns. *tibble*(*x* = 1:3. this tibble *y* = *c*("*a*", "*b*", "*c*")) tribble(...) A tibble: $3 \times 2$ Construct by rows. <int> <dbl> tribble( а ~y, "a" 2 b 1, 3 3 С 2, "b", 3. "c") as tibble(x,...) Convert data frame to tibble. enframe(x, name = "name", value = "value")

Converts named vector to a tibble with a names column and a values column. is tibble(x) Test whether x is a tibble.

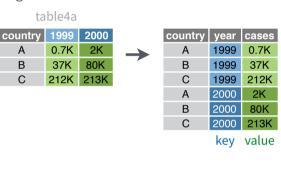
Tidy data is a way to organize tabular data. It provides a consistent data structure across packages. A table is tidy if: Tidy data: R B С C Each variable is in Each **observation**. or Makes variables easy case, is in its own row its own **column** to access as vectors vectorized operations

### Reshape Data - change the layout of values in a table

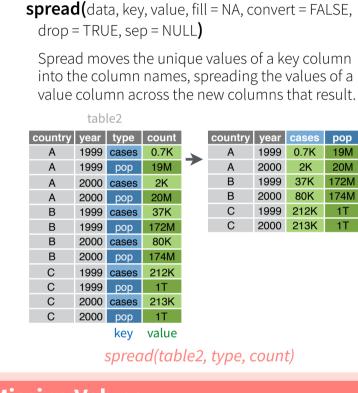
Use gather() and spread() to reorganize the values of a table into a new layout. Each uses the idea of a key column: value column pair.

**gather(**data, key, value, ..., na.rm = FALSE, convert = FALSE, factor key = FALSE

Gather moves column names into a key column, gathering the column values into a single value column.



qather(table4a, `1999`, `2000`, kev = "year", value = "cases")



### **Handle Missing Values**

drop_na(data,)	<pre>fill(data,, .direction = c("down", "up"))</pre>	replace_na(data,
Drop rows containing NA's in columns.	Fill in NA's in columns with most recent non-NA values.	replace = list(), <b>)</b> Replace NA's by column. ×
$\begin{array}{c} x1 & x2 \\ A & 1 \\ B & NA \\ C & NA \\ D & 3 \\ E & NA \end{array} \xrightarrow{\begin{array}{c} x1 & x2 \\ A & 1 \\ D & 3 \\ \end{array}} x1 \\ x2 \\ A \\ D \\ x1 \\ x2 \\ A \\ x1 \\ x2 \\ x1 \\ x2 \\ x2 \\ x2 \\ x1 \\ x2 \\ x2$	$\begin{array}{c} x1 & x2 \\ A & 1 \\ B & NA \\ C & NA \\ D & 3 \\ E & NA \end{array} \xrightarrow{\begin{array}{c} x1 & x2 \\ A & 1 \\ B & 1 \\ C & 1 \\ D & 3 \\ E & 3 \end{array}}$	$\begin{array}{c} x1 & x2 \\ A & 1 \\ B & NA \\ C & NA \\ D & 3 \\ E & NA \end{array} \xrightarrow{\begin{array}{c} x1 & x2 \\ A & 1 \\ B & 2 \\ C & 2 \\ D & 3 \\ E & 2 \end{array}}$
drop_na(x, x2)	fill(x, x2)	replace_na(x,list(x2 = 2), x2)
Expand Tabl	<b>es</b> - quickly create tables with com	binations of values

### complete(data, ..., fill = list())

d

Adds to the data missing combinations of the values of the variables listed in ... complete(mtcars, cyl, gear, carb)

### expand(data, ...)

Create new tibble with all possible combinations of the values of the variables listed in ... expand(mtcars, cyl, gear, carb)

### **Split and Combine Cells**

Use these functions to split or combine cells into individual, isolated values.

separate(data, col, into, sep = "[^[:alnum:]]+", remove = TRUE, convert = FALSE. extra = "warn", fill = "warn", ...)

Separate each cell in a column to make several columns.

	table	5					
country	year	rate		country	year	cases	рор
Α	1999	0.7K <b>/</b> 19M		Α	1999	0.7K	19M
Α	2000	2K/20M	$\rightarrow$	Α	2000	2K	20M
В	1999	37K <b>/172M</b>		В	1999	37K	172
В	2000	80K/174M		В	2000	80K	174
С	1999	212K/1T		С	1999	212K	1T
С	2000	213K/1T		С	2000	213K	1T

+ 0 0 0 1 0 1

### separate\_rows(table3, rate, *into* = *c*("*cases*", "*pop*"))

### separate\_rows(data, ..., sep = "[^[:alnum:].]+", convert = FALSE)

Separate each cell in a column to make several rows. Also separate rows ()

1	table3					
country	year	rate		country	year	rate
А	1999	0.7K <b>/</b> 19M		А	1999	0.7K
А	2000	2K <b>/20M</b>	$\rightarrow$	А	1999	19M
В	1999	37K/172M		А	2000	2K
В	2000	80K/174M		А	2000	20M
С	1999	212K <b>/1T</b>		В	1999	37K
С	2000	213K <b>/1T</b>		В	1999	172M
				В	2000	80K
				В	2000	174M
				С	1999	212K
				С	1999	1T
				С	2000	213K
				С	2000	1T

separate\_rows(table3, rate)

### **unite**(data, col, ..., sep = "\_", remove = TRUE)

Collapse cells across several columns to make a single column.



unite(table5, century, year, *col* = "*year*", *sep* = "")

### **Tidy Data with tidyr**

A

A \* B -> C

Preserves cases during

С

20M

172M

174M

1T

\* B

### **Data Wrangling** with dplyr and tidyr

Cheat Sheet

Studio

### **Syntax** - Helpful conventions for wrangling

### dplyr::tbl df(iris)

Converts data to tbl class. tbl's are easier to examine than data frames. R displays only the data that fits onscreen:

Source: lo	ocal data f	rame [150 x 5	5]
		l.Width Peta	-
1	5.1	3.5	1.4
2	4.9	3.0	1.4
3	4.7	3.2	1.3
4	4.6	3.1	1.5
5	5.0	3.6	1.4
Variables Species		Petal.Width	(dbl),

### dplyr::glimpse(iris)

Information dense summary of tbl data.

### utils::View(iris)

View data set in spreadsheet-like display (note capital V).

	iris ×				
$\langle \diamondsuit \rangle$	🖒 🙇 🛛 🖓 Fil	ter		Q,	
	Sepal.Length $\stackrel{\diamond}{=}$	Sepal.Width $^{\diamond}$	Petal.Length ≑	Petal.Width 🗘	Species $^{\diamond}$
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
4	4.6	3.1	1.5	0.2	setosa
5	5.0	3.6	1.4	0.2	setosa
6	5.4	3.9	1.7	0.4	setosa
7	4.6	3.4	1.4	0.3	setosa
8	5.0	3.4	1.5	0.2	setosa

### dplyr::%>%

Passes object on left hand side as first argument (or . argument) of function on righthand side.

> $x \gg f(y)$  is the same as f(x, y)y %>% f(x, ., z) is the same as f(x, y, z)

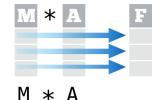
"Piping" with %>% makes code more readable, e.g.

iris %>% group\_by(Species) %>% summarise(avg = mean(Sepal.Width)) %>% arrange(avg)

### Tidy Data - A foundation for wrangling in R



Each **observation** is saved in its own **row**  Tidy data complements R's vectorized **operations**. R will automatically preserve observations as you manipulate variables. No other format works as intuitively with R.



### Reshaping Data - Change the layout of a data set



tidyr::gather(cases, "year", "n", 2:4) Gather columns into rows.

Fach variable is saved

in its own **column** 

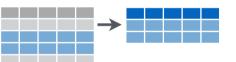
In a tidv

data set:



tidyr::separate(storms, date, c("y", "m", "d")) Separate one column into several.

### Subset Observations (Rows)



dplyr::filter(iris, Sepal.Length > 7)

Extract rows that meet logical criteria.

### dplyr::distinct(iris)

Remove duplicate rows.

dplyr::sample\_frac(iris, 0.5, replace = TRUE)

Randomly select fraction of rows.

dplyr::sample\_n(iris, 10, replace = TRUE)

Randomly select n rows.

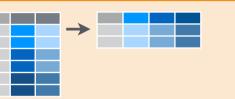
dplyr::slice(iris, 10:15)

Select rows by position.

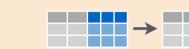
### dplyr::top\_n(storms, 2, date)

Select and order top n entries (by group if grouped data).

	Logic in R - ?	Comparison, ?base	::Logic
<	Less than	!=	Not equal to
>	Greater than	%in%	Group membership
==	Equal to	is.na	Is NA
<=	Less than or equal to	!is.na	Is not NA
>=	Greater than or equal to	&, ,!,xor,any,all	Boolean operators



tidyr::spread(pollution, size, amount) Spread rows into columns.



tidyr::unite(data, col, ..., sep) Unite several columns into one.

### dplyr::data frame(a = 1:3, b = 4:6)Combine vectors into data frame (optimized).

dplyr::arrange(mtcars, mpg)

Order rows by values of a column (low to high).

dplyr::arrange(mtcars, desc(mpg))

Order rows by values of a column (high to low).

dplyr::rename(tb, y = year) Rename the columns of a data frame

### Subset Variables (Columns)

dplyr::select(iris, Sepal.Width, Petal.Length, Species)

Select columns by name or helper function.

### Helper functions for select - ?select select(iris, contains(".")) Select columns whose name contains a character string.

select(iris, ends\_with("Length"))

Select columns whose name ends with a character string.

select(iris, everything()) Select every column.

### select(iris, matches(".t."))

Select columns whose name matches a regular expression.

### select(iris, num\_range("x", 1:5))

Select columns named x1, x2, x3, x4, x5. select(iris, one\_of(c("Species", "Genus")))

Select columns whose names are in a group of names.

### select(iris, starts\_with("Sepal"))

Select columns whose name starts with a character string.

select(iris, Sepal.Length:Petal.Width)

Select all columns between Sepal.Length and Petal.Width (inclusive). select(iris, -Species) Select all columns except Species.

RStudio® is a trademark of RStudio, Inc. • CC BY RStudio • info@rstudio.com • 844-448-1212 • rstudio.com

devtools::install\_github("rstudio/EDAWR") for data sets

Learn more with browseVignettes(package = c("dplyr", "tidyr")) • dplyr 0.4.0• tidyr 0.2.0 • Updated: 1/15

### **Summarise Data**

 $\rightarrow$ 

dplyr::summarise(iris, avg = mean(Sepal.Length)) Summarise data into single row of values.

dplyr::summarise each(iris, funs(mean))

Apply summary function to each column.

### dplyr::count(iris, Species, wt = Sepal.Length)

Count number of rows with each unique value of variable (with or without weights).

> summary function

Summarise uses **summary functions**, functions that take a vector of values and return a single value, such as:

dplyr::first	min
First value of a vector.	Minimum value in a vector.
dplyr::last	max
Last value of a vector.	Maximum value in a vector.
dplyr::nth	mean
Nth value of a vector.	Mean value of a vector.
dplyr:: <b>n</b>	median
# of values in a vector.	Median value of a vector.
dplyr::n_distinct	var
# of distinct values in	Variance of a vector.
a vector.	sd
IQR	Standard deviation of a
IQR of a vector.	vector.

### **Group Data**

### dplyr::group\_by(iris, Species)

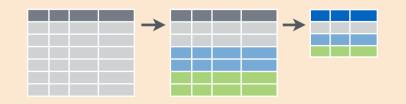
Group data into rows with the same value of Species.

dplyr::ungroup(iris)

Remove grouping information from data frame.

### iris %>% group\_by(Species) %>% summarise(...)

Compute separate summary row for each group.



### **Make New Variables**



dplyr::mutate(iris, sepal = Sepal.Length + Sepal. Width)

Compute and append one or more new columns.

dplyr::mutate each(iris, funs(min rank))

Apply window function to each column.

dplyr::lead

dplyr::lag

dplyr::ntile

dplvr::between

dplyr::transmute(iris, sepal = Sepal.Length + Sepal. Width)

Compute one or more new columns. Drop original columns.

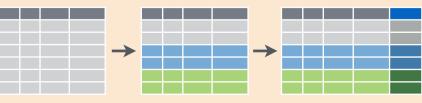


Mutate uses window functions, functions that take a vector of values and return another vector of values, such as:

### dplyr::cumall Copy with values shifted by 1. Cumulative **all** dplyr::cumany Copy with values lagged by 1. Cumulative **any** dplyr::dense\_rank dplyr::cummean Ranks with no gaps. Cumulative **mean** dplyr::min\_rank cumsum Ranks. Ties get min rank. Cumulative **sum** dplyr::percent\_rank cummax Ranks rescaled to [0, 1]. Cumulative **max** dplyr::row\_number cummin Cumulative **min** Ranks. Ties got to first value. cumprod Cumulative **prod** Bin vector into n buckets. pmax Are values between a and b? dplyr::cume dist pmin Cumulative distribution. Element-wise **min**

iris %>% group\_by(Species) %>% mutate(...)

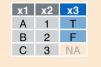
Compute new variables by group.



### **Combine Data Sets**

#### b а x1 x2 x1 x3 Α 1 Т В 2 C 3

### Mutating Joins



Join matching rows from b to a.

dplyr::left join(a, b, by = "x1")

XI	XJ	XZ
Α	Т	1
В	F	2
D	Т	NA

dplyr::right\_join(a, b, by = "x1") Join matching rows from a to b.

x1	x2	х3	d
А	1	Т	
В	2	F	

### plyr::inner\_join(a, b, by = "x1") Join data. Retain only rows in both sets.

x1	x2	x3
А	1	Т
В	2	F
С	3	NA
D	NA	Т

### dplyr::full\_join(a, b, by = "x1")

Join data. Retain all values, all rows.

### Filtering Joins

x1	x2	<pre>dplyr::semi_join(a, b, by = "x1")</pre>
Α	1	• • • • • • •
В	2	All rows in a that have a match ir
x1	x2	<pre>dplyr::anti_join(a, b, by = "x1")</pre>
С		

rows in a that have a match in b. /r::anti\_join(a, b, by = "x1") All rows in a that do not have a match in b.

	у			Z	
x1	x2		<b>x1</b>	x2	
А	1	1.1	В	2	_
В	2		С	3	
С	3	_	D	4	

### Set Operations

B 2	x1	x2
0 0	В	2
C 3	С	3

x1 x2

A 1

в

2

Rows that appear in both y and z.

### dplyr::union(y, z)

dplyr::intersect(y, z)

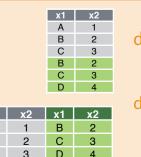
Rows that appear in either or both y and z.

### dplyr::setdiff(y, z)

Rows that appear in y but not z.

### Binding

С



dplyr::bind\_rows(y, z) Append z to y as new rows.

### dplyr::bind\_cols(y, z)

Append z to y as new columns. Caution: matches rows by position.

RStudio® is a trademark of RStudio, Inc. • CC BY RStudio • info@rstudio.com • 844-448-1212 • rstudio.com

devtools::install\_github("rstudio/EDAWR") for data sets

Learn more with **browseVignettes(package = c("dplyr", "tidyr"))** • dplyr 0.4.0• tidyr 0.2.0 • Updated: 1/15

C 3 D 4 Element-wise **max** x1 x2 A 1

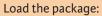
### **R For Data Science** *Cheat Sheet* data.table

Learn R for data science Interactively at www.DataCamp.com



### data.table

data.table is an R package that provides a high-performance version of base R's data.frame with syntax and feature enhancements for ease of use, convenience and programming speed. 



> library(data.table)

### Creating A data.table

set.seed(45L) Create a data.table and call it DT > DT <- data.table(V1=c(1L,2L), V2=LETTERS[1:3], V3=round(rnorm(4),4), V4=1:12)

### Subsetting Rows Using i

> DT[3:5,]	Select 3rd to 5th row	
> DT[3:5]	Select 3rd to 5th row	
> DT[V2=="A"]	Select all rows that have value A in column v2	
> DT[V2 %in% c("A","C")]	Select all rows that have value A or C in column v2	

### Manipulating on Columns in j

Г		ā
>	> DT[,V2]	Return v2 as a vector
	[1] "A" "B" "C" "A" "B" "C"	
>	DT[,.(V2,V3)]	Return v2 and v3 as a data.table
>	<pre>DT[,sum(V1)]</pre>	Return the sum of all elements of v1 in a
	[1] 18	vector
>	<pre>DT[,.(sum(V1),sd(V3))]</pre>	Return the sum of all elements of v1 and the
	V1 V2	std.dev.of V3 in a data.table
	1: 18 0.4546055	
>	<pre>DT[,.(Aggregate=sum(V1),</pre>	The same as the above, with new names
	Sd.V3=sd(V3))]	
	Aggregate Sd.V3	
	1: 18 0.4546055	
>	<pre>DT[,.(V1,Sd.V3=sd(V3))]</pre>	Select column v2 and compute std. dev. of v3,
		which returns a single value and gets recycled
>	<pre>DT[,.(print(V2),</pre>	Print column v2 and plot v3
	plot(V3),	
	NULL)]	
L		

### Doing j by Group

> DT[,.(V4.Sum=sum(V4)),by=V1]	Calculate sum of v4 for every group in v1	> DT[c("A","C"),
V1 V4.Sum		 sum(V4),
1: 1 36		 by=.EACHI]
2: 2 42		V2 V1
> DT[,.(V4.Sum=sum(V4)),	Calculate sum of v4 for every group in v1	1: A 22
by=.(V1,V2)]	and v2	 2: C 30
> DT[,.(V4.Sum=sum(V4)),	Calculate sum of v4 for every group in	 > setkey(DT,V1,V2)
by=sign(V1-1)]	sign(V1-1)	> DT[.(2,"C")]
	Sign(VI-I)	 V1 V2 V3 V4
sign V4.Sum 1: 0 36		1: 2 C 0.3262 6
2: 1 42		2: 2 C -1.6148 12
> DT[,.(V4.Sum=sum(V4)),	The same as the above, with new name	> DT[.(2,c("A","C"))]
	for the variable you're grouping by	V1 V2 V3 V4
> DT[1:5,.(V4.Sum=sum(V4)),	Calculate sum of $v_4$ for every group in $v_1$	1: 2 A -1.6148 4
by=V1]	after subsetting on the first 5 rows	2: 2 A 0.3262 10
> DT[,.N,by=V1]	Count number of rows for every group in	3: 2 C 0.3262 6
[,.Μ,Υ]	v1	4: 2 C -1.6148 12

General form: DT[i, j,	by] —, (			Advanced Data Table O	perations
"Take DT, subset rows us		culate j grouped by by"		> DT[.N-1] > DT[,.N]	Return the penultimate row of the DT Return the number of rows
Adding/Updating Colu	ımns By Ref	erence in 🕆 Using :=		<pre>&gt; DT[,.(V2,V3)] &gt; DT[,list(V2,V3)]</pre>	Return v2 and v3 as a data.table Return v2 and v3 as a data.table
<pre>&gt; DT[,V1:=round(exp(V1),2)]</pre>	initis by Ref	v1 is updated by what is after :=		<pre>&gt; DT[,mean(V3),by=.(V1,V2)]</pre>	Return the result of $j$ , grouped by all possible combinations of groups specified in by
> DT V1 V2 V3 V4		Return the result by calling $\mathbb{D}\mathbb{T}$		1: 1 A 0.4053 2: 1 B 0.4053 3: 1 C 0.4053	
1: 2.72 A -0.1107 1 2: 7.39 B -0.1427 2				4: 2 A -0.6443 5: 2 B -0.6443	
3: 2.72 C -1.8893 3 4: 7.39 A -0.3571 4				6: 2 C -0.6443	
 > DT[,c("V1","V2"):=list(rou	nd(exp(V1),2),	Columns v1 and v2 are updated by		.SD & .SDcols	
	TERS[4:6])]	what is after := Alternative to the above one. With [],		<pre>&gt; DT[,print(.SD),by=V2] &gt; DT[,.SD[c(1,.N)],by=V2]</pre>	Look at what $.$ SD contains Select the first and last row grouped by $v_2$
V2=LETTERS[4:6])] V1 V2 V3 V4		you print the result to the screen		> DT[,lapply(.SD,sum),by=V2]	Calculate sum of columns in . SD grouped by V2
1: 15.18 D -0.1107 1 2: 1619.71 E -0.1427 2				<pre>&gt; DT[,lapply(.SD,sum),by=V2, .SDcols=c("V3","V4")]</pre>	Calculate sum of $v_3$ and $v_4$ in $.s_D$ grouped by $v_2$
3: 15.18 F -1.8893 3 4: 1619.71 D -0.3571 4				V2 V3 V4 1: A -0.478 22 2: B -0.478 26	
<pre>&gt; DT[,V1:=NULL] &gt; DT[,c("V1","V2"):=NULL]</pre>		Remove v1 Remove columns v1 and v2		<pre>3: c -0.478 30 &gt; DT[,lapply(.SD,sum),by=V2,</pre>	Calculate sum of $v_3$ and $v_4$ in $. s_D$ grouped by
<pre>&gt; Cols.chosen=c("A","B") &gt; DT[,Cols.Chosen:=NULL]</pre>		Delete the column with column name		.SDcols=paste0("V",3:4)	V2
> DT[,(Cols.Chosen):=NULL]		Cols.chosen Delete the columns specified in the		Chaining	
		variable Cols.chosen		> DT <- DT[,.(V4.Sum=sum(V4)),	Calculate sum of $v_4$ , grouped by $v_1$
Indexing And Keys				by=V1] V1 V4.Sum	
> setkey(DT,V2)		utput is returned invisibly		1: 1 36 2: 2 42	
> DT["A"] V1 V2 V3 V4	the value A	re the key column (set to v2) has	>	<pre>&gt; DT[V4.Sum&gt;40] &gt; DT[,.(V4.Sum=sum(V4)),</pre>	Select that group of which the sum is >40 Select that group of which the sum is >40
1: 1 A -0.2392 1 2: 2 A -1.6148 4				by=V1][V4.Sum>40] V1 V4.Sum	(chaining)
3: 1 A 1.0498 7 4: 2 A 0.3262 10 > DT[c("A","C")]	Poturn all rows who	re the key column (v2) has value A or C	>	1: 2 42 > DT[,.(V4.Sum=sum(V4)),	Calculate sum of $v_4$ , grouped by $v_1$ ,
> DT["A", mult="first"]		ll rows that match value A in key		by=V1][order(-V1)] V1 V4.Sum 1: 2 42	ordered on V1
<pre>&gt; DT["A",mult="last"]</pre>		l rows that match value A in key		1: 2 42 2: 1 36	
> DT[c("A", "D")] v1 v2 v3 v4		re key column V2 has value A or D		cot() Family	
1: 1 A -0.2392 1 2: 2 A -1.6148 4				set()-Family	
3: 1 A 1.0498 7 4: 2 A 0.3262 10				set()	
5: NA D NA NA > DT[c("A", "D"), nomatch=0]	Return all rows whe	re key column v2 has value A or D	5	> rows <- list(3:4,5:6)	t(DT, row, column, new value)
V1 V2 V3 V4 1: 1 A -0.2392 1				> cols <- 1:2 > for(i in seq along(rows))	Sequence along the values of rows, and
2: 2 A -1.6148 4 3: 1 A 1.0498 7 4: 2 A 0.3262 10				<pre>{set(DT,     i=rows[[i]],</pre>	for the values of cols, set the values of those elements equal to NA (invisible)
<pre>&gt; DT[c("A", "C"), sum(V4)]</pre>	Return total sum of have values A or c	v4, for rows of key column v2 that		j=cols[i], value=NA)}	
<pre>&gt; DT[c("A", "C"),</pre>	Return sum of colun	nn v4 for rows of v2 that have value A, r rows of v2 that have value C		setnames()	
by=.EACHI]			Syntax: setnames	(DT, "old", "new") []	
1: A 22 2: C 30			>	<pre>&gt; setnames(DT,"V2","Rating") &gt; setnames(DT,</pre>	Set name of V2 to Rating (invisible) Change 2 column names (invisible)
<pre>&gt; setkey(DT,V1,V2) &gt; DT[.(2,"C")]</pre>		by v2 within each group of v1 (invisible) ve value 2 for the first key (v1) and the		c("V2","V3"), c("V2.rating","V3.I	) (C"))
V1 V2 V3 V4 1: 2 C 0.3262 6	value c for the seco			<pre>setnames()</pre>	
2: 2 C -1.6148 12 > DT[.(2,c("A","C"))]	Select rows that have	e value 2 for the first key (v1) and within		Syntax: setcolor	der(DT, "neworder")
V1 V2 V3 V4 1: 2 A -1.6148 4		e A or C for the second key (V2)	>	<pre>&gt; setcolorder(DT,</pre>	Change column ordering to contents , "V3")) of the specified vector (invisible)
2: 2 A 0.3262 10 3: 2 C 0.3262 6					- ^
4: 2 C -1.6148 12				Data Learn Python for Data	

# Data Transformation with data.table :: **CHEAT SHEET**

### **Basics**

data.table is an extremely fast and memory efficient package for transforming data in R. It works by converting R's native data frame objects into data.tables with new and enhanced functionality. The basics of working with data.tables are:

### dt[i, j, by]

Take data.table **dt**. subset rows using i and manipulate columns with i. grouped according to **by**.

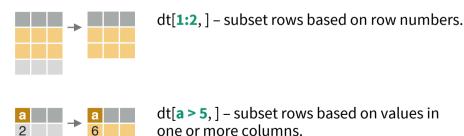
data.tables are also data frames - functions that work with data frames therefore also work with data.tables.

### Create a data.table

**data.table**(a = c(1, 2), b = c("a", "b")) – create a data.table from scratch. Analogous to data.frame().

setDT(df)\* or as.data.table(df) - convert a data frame or a list to a data.table.

### Subset rows using i



5

dt[a > 5, ] – subset rows based on values in one or more columns.

### LOGICAL OPERATORS TO USE IN i

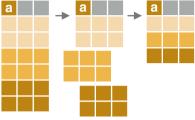
<	<=	is.na()	%in%	
>	>=	!is.na()	!	&

### %like% %between%



dt[, **b** := as.integer(b)] – convert the type of a column using as.integer(), as.numeric(), as.character(), as.Date(), etc..

### Group according to by



dt[, j, by = .(a)] – group rows by values in specified columns.

dt[, j, keyby = .(a)] – group and simultaneously sort rows by values in specified columns.

### **COMMON GROUPED OPERATIONS**

dt[, .(c = sum(b)), by = a] - summarize rows within groups.

dt[, c := sum(b), by = a] – create a new column and compute rows within groups.

dt[, .SD[1], by = a] - extract first row of groups.

dt[, .SD[.N], by = a] – extract last row of groups.

### Chaining

dt[...][...] – perform a sequence of data.table operations by *chaining* multiple "[]".

### Functions for data.tables

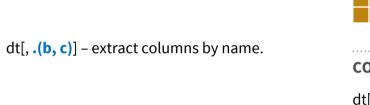
### REORDER

а	b		а	b	
1	2		1	2	
2	2		1	1	
1	1		2	2	

**setorder**(dt, a, -b) – reorder a data.table according to specified columns. Prefix column names with "-" for descending order.

### \* SET FUNCTIONS AND :=

data.table's functions prefixed with "set" and the operator ":=" work without "<-" to alter data without making copies in memory. E.g., the more efficient "setDT(df)" is analogous to "df <- as.data.table(df)".



dt[, c(2)] – extract columns by number. Prefix

column numbers with "-" to drop.

### **SUMMARIZE**

**EXTRACT** 



dt[, .(x = sum(a))] – create a data.table with new columns based on the summarized values of rows.

Manipulate columns with j

Summary functions like mean(), median(), min(), max(), etc. can be used to summarize rows.

### **COMPUTE COLUMNS\***



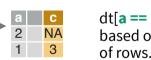
C d

**CONVERT COLUMN TYPE** 

1 2

1 2

dt[, c := 1 + 2] – compute a column based on an expression.

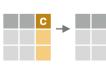


dt[a == 1, c := 1 + 2] - compute a columnbased on an expression but only for a subset of rows.

dt[, := (c = 1, d = 2)] - compute multiplecolumns based on separate expressions.

dt[, **c** := NULL] – delete a column.

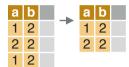
### **DELETE COLUMN**



2



### **UNIQUE ROWS**



**unique(**dt, by = c("a", "b")**)** – extract unique rows based on columns specified in "by". Leave out "by" to use all columns.

**uniqueN(**dt, by = c("a", "b")) – count the number of unique rows based on columns specified in "by".

### **RENAME COLUMNS**



**setnames(**dt, c("a", "b"), c("x", "y")**)** – rename columns.

### SET KEYS

**setkey(**dt, a, b) – set keys to enable fast repeated lookup in specified columns using "dt[.(value), ]" or for merging without specifying merging columns using "dt\_a[dt\_b]".

### Combine data.tables

а	b		X	у		а	b	X
1	С		3	b	_	3	b	3
2	а	+	2	С	—	1	С	2
3	b		1	а		2	а	1

dt\_a[dt\_b, **on = .(b = y)**] – join data.tables on rows with equal values.

dt\_a[dt\_b, on = .(b = y, c > z)] -

ioin data.tables on rows with

equal and unequal values.

а	b	С		X	У	z		а	b	С	X
1	С	7		3	b	4	_	3	b	4	3
2	а	5	+	2	С	5	_	1	С	5	2
3	b	6		1	а	8		NA	а	8	1

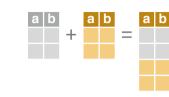
### **ROLLING JOIN**

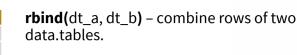
a <mark>id</mark>	date		b	id	date		a	id	d	ate	b
1 A	01-01-2010	+	1	А	01-01-2013	3 =	2	А	01-0	1-2013	1
2 A	01-01-2012		1	В	01-01-2013	3	2	В	01-0	1-2013	1
3 A	01-01-2014										
1 B	01-01-2010										
2 B	01-01-2012										

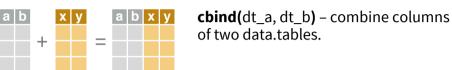
dt\_a[dt\_b, **on = .(id = id, date = date), roll = TRUE**] – join

data.tables on matching rows in id columns but only keep the most recent preceding match with the left data.table according to date columns. "roll = -Inf" reverses direction.

### BIND

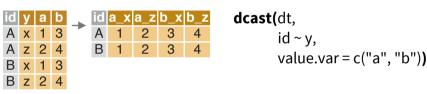






### Reshape a data.table

### **RESHAPE TO WIDE FORMAT**



### Reshape a data.table from long to wide format.

dt A data.table.

id ~ y Formula with a LHS: ID columns containing IDs for multiple entries. And a RHS: columns with values to spread in column headers.

value.var Columns containing values to fill into cells.

### **RESHAPE TO LONG FORMAT**

<b>melt(</b> dt,					b_z	b_x	a_z	a_x	id
id.v	3	1	1	А	4	3	2	1	А
	3	1	1	В	4	3	2	1	В
mea	4	2	2	А					
vari	4	2	2	В					
vali				_					

### id.vars = c("id"), measure.vars = patterns("^a", "^b"),

variable.name = "y",

value.name = c("a", "b")**)** 

### Reshape a data.table from wide to long format.

dt	A data.table.
id.vars	ID columns with IDs for multiple entries.
measure.vars	Columns containing values to fill into cells (often in
	pattern form).
variable.name,	Names of new columns for variables and values
value.name	derived from old headers.

### Apply function to cols.



а	b		а	b	
1	4		2	5	
2	5				
3	6				

<mark>a a\_m</mark>

2

2

2

1

2

3

dt[, **lapply(.SD, mean), .SDcols = c("a", "b")**] – apply a function – e.g. mean(), as.character(), which.max() – to columns specified in .SDcols with lapply() and the .SD symbol. Also works with groups.

### cols <- c("a")

dt[, paste0(cols, "\_m") := lapply(.SD, mean),

**.SDcols = cols**] – apply a function to specified columns and assign the result with suffixed variable names to the original data.

### Sequential rows

### **ROW IDS**

1

2

3

a	b	а	b	С	
1	а	1	а	1	
2	а	2	а	2	
3	b	3	b	1	

dt[, **c** := 1:.N, by = b] – within groups, compute a column with sequential row IDs.

### LAG & LEAD

a	b	a	b	С
1	а	1	а	NA
2	а	2	а	1
3	b	3	b	NA
4	b	4	b	3
5	b	5	b	4

dt[, **c** := **shift(a**, **1)**, **by** = **b**] – within groups, duplicate a column with rows *lagged* by specified amount.

dt[, **c** := **shift(a**, **1**, **type** = "**lead**"), **by** = **b**] – within groups, duplicate a column with rows *leading* by specified amount.

### read & write files

### IMPORT

fread("file.csv") - read data from a flat file such as .csv or .tsv into R.

**fread(**"file.csv", select = c("a", "b")) – read specified columns from a flat file into R.

### EXPORT

fwrite(dt, "file.csv") - write data to a flat file from R.

### Machine Learning Modelling in R : : CHEAT SHEET

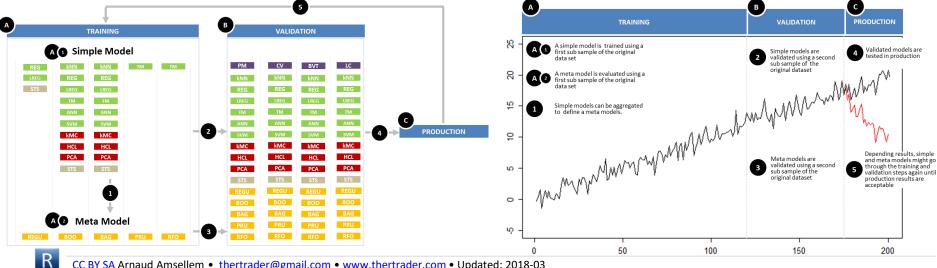
ALGORITHM	DESCRIPTION	R PACKAGE::FUNCTION	SAMPLE CODE
Naïve Bayes classifier	A classification technique based on Bayes' Theorem with an assumption of independence among predictors. In simple terms, a Navie Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature	e1071::naiveBayes	naiveBayes(class ~ ., data = x)
k-Nearest Neighbours	A non-parametric method used for classification and regression. In both cases, the input consists of the k closest training examples in the feature space. The output depends on whether k-NN is used for classification or regression	class::knn	knn(train, test, cl, k = 1, l = 0, prob = FALSE, use.all = TRUE)
RES Linear Regression	Model the linear relationship between a scalar dependant variable Y and one or more explanatory variables (or independent variables) denoted X	stats::lm	Im(dist ~ speed, data=cars)
Hest Cogistic Regression	Used to predict a binary outcome (1 / 0, Yes / No, True / False) given a set of independent variables.	stats::glm	$\operatorname{glm}(Y^{\sim}$ ., family = binomial (link = 'logit'), data = X)
Tree-Based Models	The idea is to consecutively divide (branch) the training dataset based on the input features until an assignment criterion with respect to the target variable into a "data bucket" (leaf) is reached	rpart::rpart	rpart(Kyphosis~ Age + Number + Start, data = kyphosis)
Artificial Neural Network	Neural networks are built from units called perceptrons. Perceptrons have one or more inputs, an activation function and an output. An ANM model is built up by combining perceptrons in structured layers.	neuralnet::neuralnet	neuralnet(f,data=train_,hidden=c(5,3),linear.output=T)
Support Vector Machine	A data classification method that separates data using hyperplanes	e1071::svm	svm(formula, data = NULL,, subset, na.action = na.omit, scale = TRUE)
Principal Component Analysis	A procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components.	stats::pricomp stats::princomp FactoMineR::PCA ade4:;dudi.pca amap::acp	<pre>stats : prcomp(formula, data = NULL, subset, na.action,) stats : princomp(formula, data = NULL, subset, na.action,) FactoMineR : PCA(decathlon, quanti.sup = 11:12, quali.sup=13) ade4 : dudi.pca(deugStab, center = deugScent, scale = FALSE, scan = FALSE amap : acg(bubsch)</pre>
k-Mean Clustering	Aims at partitioning $\boldsymbol{n}$ observations into $k$ clusters in which each observation belongs to the cluster with the nearest mean	stats::kmeans	kmeans(x, centers, iter.max = 10, nstart = 1, algorithm = c("Hartigan-Wong", "Lloyd", "Forgy", "MacQueen"), trace=FALSE)
HCL Hierarchical Clustering	An approach which builds a hierarchy from the bottom-up, and doesn't require the number of clusters to be specified beforehand.	stats::hclust	hclust(d, method = "complete", members = NULL)

### **Meta-Algorithm, Time Series & Model Validation**

	ALGORITHM	DESCRIPTION	R PACKAGE::FUNCTION	SAMPLE CODE
	L1 (Lasso) L2 (Ridge)	Regularization adds a penalty on the different parameters of a model to reduce the freedom of the model. Hence, the model will be less likely to fit the noise of the training data and will improve the generalization abilities of the model	glmnet::glmnet	$\label{eq:L1:glmnet(myMatrixA, myMatrixB, family = "gaussian", alpha = 1) \\ L2:glmnet(myMatrixA, myMatrixB, family = "gaussian", alpha = 0) \\ \end{cases}$
	B00 Boosting	A process of iteratively refining, e.g. by reweighting, of estimated regression and classification functions (though it has primarily been applied to the latter), in order to improve predictive ability.	Paramatric model - mboost::glmboost	glmboost(Yen ~ ., data = curr1[trnidxs,])
= TRUE)	Bagging	Bagging is a way to increase the power of a predictive statistical model by taking multiple random samples (with replacement) of the training data set, and using each of them to construct a separate model and separate predictions for the original test set	All models: foreach Tree models: ipred::bagging	<pre>foreach: d &lt;&lt; data.frame(x=1:10, y=rnorm(10)) s &lt; foreach(d=iter(d, by='row'), .combine=rbind) %dopar%d identical(s, d) ipred i bagging(formula, data, subset, na.action=na.rpart, \dots)</pre>
	Pruning	Pruning is a technique that reduces the size of decision tree by removing sections of the tree that provide little power to classify instances. Pruning reduces the complexity of the final classifier and hence improves predictive accuracy by reducing overfitting	rpart::prune	prune(x, cp = 0.1)
iosis)	Random Forrest	An ensemble learning method for classification, regression and other tasks, that operate by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression)	randomForest::randomForest	randomForest(X ~ . , data = Y , subset = mySub)
ut=T) a.omit,	Lead-lag analysis, Auto-correlation, Spectral analysis, Time series clustering, Seasonality, Trend	Random sampling of observations for training and testing a model can be an issue when faced with a times dimension. Random sampling may either destroy serial correlation properties in the data which we would like to exploit	stats tts forecast spectral TTR 	Auto-correlation: acf(x, lag.max = NULL, type = c("correlation", "covariance", "partial")) Spectral Analysis: spec.pgram((my (s, spans = NULL) Seasonal Decomposition of Time Series st((s, s.window = 7, t.window = 50, t.jump = 1) 
tion,) action,	Phil Performance metrics	Depends on the problem: • Regression: squared errors, outliers, error rate • Classification: Accuracy, precision, recall, F-score	Regression-stats::outlierTest, stats:: qqPlot Classification-ROCR:: Tree: caret:: confusionMatrix	Regresion: fit <- Im(Y~X,data=myData) outlierTest(fit) qqPlot(fit, main="QQ Plot")
le =	Biais-Variance Tradeoff	<ul> <li>Simple models with few parameters are easier to compute but may lead to poorer fits (high bias).</li> <li>Complex models may provide more accurate fits but may over-fit the data (high variance)</li> </ul>	Tailored to the analysis	Tailored to the analysis
:hm = '),	Cross validation	Cross validation compares the test performances of different model realisations with different sets or values of parameters	caret::createDataPartition caret::createFolds	createDataPartition(classes, p = 0.8, list = FALSE)
	Learning Curves	Learning curves plot a model's training and test errors, or the chosen performance metric, depending on the training set size	caret::learing_curve_dat	learing_curve_dat(dat, outcome = NULL, proportion = (1:10)/10, test_prop = 0, verbose = TRUE,)

**Time Series View** 

### **Standard Modelling Workflow**



# Machine Learning with R

### Introduction

mlr offers a unified interface for the basic building blocks of machine learning: tasks, learners, hyperparameters, etc.

Tasks contain a description of a task (classification, regression, clustering, etc.) and a data set.

Learners specify a machine learning algorithm (GLM, SVM, xgboost, etc.) and its parameters.

Hyperparameters are learner settings that can be specified directly or tuned. A parameter set lists the possible hyperparameters for a given learner.

Wrapped Models are learners that have been trained on a task and can be used to make predictions.

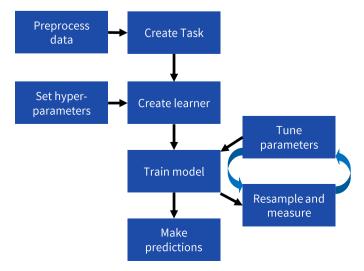
**Predictions** are the results of applying a model to either new data or the original training data.

Measures control how learner performance is evaluated, e.g. RMSE, LogLoss, AUC, etc.

**Resampling** estimates generalization performance by separating training data from test data. Common strategies include holdout and cross-validation.

### Links: Tutorial | CRAN | Github

### mlr workflow



### Setup

### **Preprocessing data**

createDummyFeatures(obj=,target=,method=,cols=) Creates (0,1) flags for each non-numeric variable excluding target. Can be applied to entire dataset or only specific cols

normalizeFeatures(obj=,target=,method=,cols=, range=.on.constant=)

Normalizes numerical features according to specified method: • "center" (subtract mean)

- "scale" (divide by std. deviation)
- "standardize" (center and scale)
- "range" (linear scale to given range, default range=c(0,1))

mergeSmallFactorLevels(task=,cols=,min.perc=) Combine infrequent factor levels into a single merged level

summarizeColumns(obi=) where obi is a data.frame or task. Provides type, NA, and distributional data about each column

See also capLargeValues dropFeatures removeConstantFeatures summarizeLevels



BC

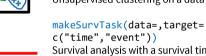
makeClassifTask(data=,target=) Classification of a target variable, with optional positive class positive



Regression on a target variable makeMultilabelTask(data=,target=)

Classification where the target can belong to more than one class per observation





c("time","event")) Survival analysis with a survival time column and an event column

makeCostSensTask(data=,costs=) Cost-sensitive classification where each observation-cost pair has a specified cost

Other arguments that can be passed to a task:

В

• weights = Weighting vector to apply to observations • blocking= Factor vector where each level indicates a block of observations that will not be split up in resampling

### Making a learner

makeLearner(cl=,predict.type=,...,par.vals=) Choose an algorithm class to perform the task and determine what that algorithm will predict

- cl=name of algorithm, e.g. "classif.xgboost" "regr.randomForest" "cluster.kmeans"
- predict.type="response" returns a prediction type that matches the source data; "prob" returns a predicted probability for classification problems only; "se" returns the a standard error of the prediction for regression problems only. Only certain learners can return "prob" and "se"
- par.vals = takes a list of hyperparameters and passes them to the learner; parameters can also be passed directly (...) You can make multiple learners at once with makeLearners()
- mlr has integrated over 170 different learning algorithms
- Full list: View(listLearners()) shows all learners Available learners for a task: View(listLearners(task))
- Filtered list: View(listLearners("classif", properties=c("prob", "factors"))) shows all classification learners "classif" which can predict probabilities "prob" and handle factor inputs "factors"
- See also getLearnerProperties()

mlr cheatsheet <u>CC BY</u> Aaron Cooley • aaronsama@gmail.com • Learn more about mlr at <u>https://mlr-org.github.io/mlr-tutorial/devel/html/</u>

### **Training & Testing**

#### **Setting hyperparameters** setHyperPars(learner=,...) Set the hyperparameters (settings) for each mtry / learner, if you don't want to use the defaults. You can also specify hyperparameters in the makeLearner() call



### Train a model and predict

getParamSet(learner=)



train(learner=,task=) Train a model (WrappedModel) by applying a learner to a task. By default. the model will train on all observations. The underlying model can be extracted with getLearnerModel()

Show the possible universe of parameters for

your learner; can take a learner directly, or a

text string such as "classif.qda"

predict(object=,task=,newdata=) Use a trained model to make predictions on a task or dataset. The resulting pred object can be viewed with View(pred) or accessed by as.data.frame(pred)

### **Measuring performance**

performance(pred=, measures=)

Calculate performance of predictions according to one or more of several measures (use listMeasures() for full list):

- classif acc auc bac ber brier[.scaled] f1 fdr fn fnr fp fpr gmean multiclass[.au1u .aunp .aunu .brier] npv ppv qsr ssr tn tnr tp tpr wkappa
- regrarsq expvar kendalltau mae mape medae medse mse msle rae rmse rmsle rrse rsg sae spearmanrho sse
- cluster db dunn G1 G2 silhouette
- multilabelmultilabel[.fl .subset01 .tpr .ppv .acc .hamloss]
- costsens mcp meancosts
- surv cindex • other featperc timeboth timepredict timetrain

#### For detailed performance data on classification tasks, use: calculateConfusionMatrix(pred=)

calculateROCMeasures(pred=)

### **Resampling a learner**

makeResampleDesc(method=,...,stratify=)

- method must be one of the following:
- "CV" (cross-validation, for number of folds use iters=)
- "LOO" (leave-one-out cross-validation, for folds use iters=) • "RepCV" (repeated cross-validation, for number of repetitions
- use reps=, for folds use folds=)
- "Subsample" (aka Monte-Carlo cross-validation, for iterations use iters=, for train % use split=)
- "Bootstrap" (out-of-bag bootstrap, uses iters=)
- "Holdout" (for train % use split=) stratify keeps target proportions consistent across samples.

makeResampleInstance(desc=,task=) can reduce noise by ensuring the resampling is done identically every time.

resample(learner=,task=,resampling=,measures=) Train and test model according to specified resampling strategy.

mlrincludes several pre-specified resample descriptions: cv2 (2fold cross-validation), cv3, cv5, cv10, hout (holdout with split 2/3 for training, 1/3 for testing).

Convenience functions also exist to resample() with a specific strategy:crossval(),repcv(),holdout(),subsample(), bootstrapOOB(),bootstrapB632(),bootstrapB632plus()

### **Refining** Performance

### **Tuning hyperparameters**

- Set search space using makeParamSet(make<type>Param())
- makeNumericParam(id=,lower=,upper=,trafo=)
- makeIntegerParam(id=,lower=,upper=,trafo=) makeIntegerVectorParam(id=,len=,lower=,upper=,
- trafo=)
  - makeDiscreteParam(id=,values=c(...))(can also be used to test discrete values of numeric or integer parameters)

trafo transforms the parameter output using a specified function, e.g. lower=-2, upper=2, trafo=function(x) 10^x would test values between 0.01 and 100, scaled exponentially Other acceptable parameter types include Logical LogicalVector CharacterVector DiscreteVector

- Set a search algorithm with makeTuneControl<type>()
- Grid(resolution=10L) Grid of all possible points
- Random(maxit=100) Randomly sample search space .
- MBO (budget=) Use Bayesian model-based optimization Irace(n.instances=) Iterated racing process
- Other types: CMAES, Design, GenSA

Tune using tuneParams(learner=,task=,resampling=, measures=,par.set=,control=)

### Quickstart

Prepare data for training and testing library(mlbench)

data(Soybean)

soy = createDummyFeatures(Soybean,target="Class") tsk = makeClassifTask(data=soy,target="Class") ho = makeResampleInstance("Holdout",tsk) tsk.train = subsetTask(tsk,ho\$train.inds[[1]]) tsk.test = subsetTask(tsk,ho\$test.inds[[1]]) Convert the factor inputs in the Soybean dataset into (0,1) dummy features which can be used by the XGboost algorithm. Create a task to precict the "Class" column. Create a train set with 2/3 of data and a test set with the remaining 1/3 (default).

Create learner and evaluate performance lrn = makeLearner("classif.xgboost",nrounds=10)

### cv = makeResampleDesc("CV",iters=5)

res = resample(lrn,tsk.train,cv,acc) Create an XGboost learner which will build 10 trees. Then test performance using 5-fold cross-validation. Accuracy should be between 0.90-0.92.

Tune hyperparameters and retrain model ps = makeParamSet(makeNumericParam("eta",0,1), makeNumericParam("lambda",0,200),

makeIntegerParam("max\_depth",1,20)) tc = makeTuneControlMBO(budget=100) tr = tuneParams(lrn,tsk.train,cv5,acc,ps,tc) lrn = setHyperPars(lrn,par.vals=tr\$x) Tune hyperparameters eta, lambda, and max depth by defining a search space and using Model Based Optimization (MBO) to control the search. Then perform 100 rounds of 5-fold crossvalidation, improving accuracy to ~0.93. Update the XGboost

Train the model on the train set and make predictions on the test

model on the full set to use on new data. You are now ready to go

set. Show performance as a confusion matrix. Finally, re-train

out into the real world and make 93% accurate predictions!

Legend for functions (not all parameters shown): Function(required\_parameters=,optional\_parameters=)

#### mdl = train(lrn,tsk.train) prd = predict(mdl,tsk.test) calculateConfusionMatrix(prd) mdl = train(lrn,tsk)

learner with the tuned hyperparameters.

### Configuration

mlr's default settings can be changed using configureMlr():

- show.infoWhether to show verbose output by default when training, tuning, resampling, etc. (TRUE)
- on.learner.error How to handle a learner error. "stop" halts execution, "warn" returns NAs and displays a warning, "quiet" returns NAs with no warning ("stop")
- on.learner.warning How to handle a learner warning. "warn" displays a warning, "quiet" supresses it ("warn") • on.par.without.desc How to handle a parameter with no
- description. "stop", "warn", "quiet" ("stop") on.par.out.of.bounds How to handle a parameter with an out-of-bounds value. "stop", "warn", "quiet" ("stop")
- on.measure.not.applicable How to handle a measure not applicable to a learner. "stop", "warn", "quiet" ("stop")
- show.learner.output Whether to show learner output to the console during training (TRUE)
- on.error.dump Whether to create an error dump for crashed learners if on.learner.error is not set to "stop" (TRUE)

Use getMlrOptions() to see current settings

### Parallelization

mlr works with the parallelMap package to take advantage of multicore and cluster computing for faster operations. mlr automatically detects which operations are able to run in parallel.

#### To begin parallel operation use:

parallelStart(mode=,cpus=,level=) mode determines how the parallelization is performed:

- "local" no parallelization applied, simply uses mapply • "multicore" multicore execution on a single machine, uses
- parallel::mclapply.Not available in Windows. • "socket" multicore execution in socket mode
- "mpi" Snow MPI cluster on one or multiple machines using parallel::makeCluster and parallel::clusterMap "BatchJobs" Batch queuing HPC clusters using
- BatchJobs::batchMap • cpus determines how many logical cores will be used
- level controls parallelization: "mlr.benchmark"
- "mlr.resample" "mlr.selectFeatures" "mlr.tuneParams" "mlr.ensemble"

To end parallelization, use parallelStop()

### Imputation

impute(obj=,target=,cols=,dummy.cols=,dummy.type=) Applies specified logic to data frame or task containing NAs and returns an imputation description which can be used on new data

- obj=data frame or task on which to perform imputation
- target=specify target variable which will not be imputed
- cols=column names and logic for imputation\*
- dummy.cols=column names to create a NA (T/F) column\* • dummy.type=set to "numeric" to use (0,1) instead of (T/F)
- \*Can also use classes and dummy.classes in place of cols

Imputation logic is passed to cols or classes via a list, e.g.: cols=list(V1=imputeMean()) where V1 is the column to which to apply the imputation, and imputeMean() is the imputation method. Available imputation methods include: imputeConst(const=) imputeMedian() imputeMode() imputeMin(multiplier=) imputeMax(multiplier=) imputeNormal(mean=,sd=) imputeHist(breaks=,use.mids=) imputeLearner(learner=,features=) impute returns a list containing the imputed dataset or task as well as an imputation description that can be used to reapply the same imputation to new data using reimpute

reimpute(obj=, desc=) Imputes missing values on a task or dataset (obj) using a description (desc) created by impute

mlr cheatsheet <u>CC BY</u> Aaron Cooley • aaronsama@gmail.com • Learn more about mlr at <u>https://mlr-org.github.io/mlr-tutorial/devel/html</u>/

### **Feature Extraction**

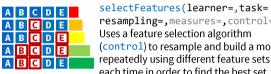
### Feature filtering

filterFeatures(task=,method=, perc=,abs=,threshold=) Uses a learner-agnostic feature evaluation method to rank feature importance, then includes only features in the top n percent ACDBFE (perc=), top n (abs=), or which meet a set performance threshold (threshold=).

Outputs a task with features that failed the test omitted. method defaults to "randomForestSRC.rfsrc", but can be set to: "anova.test" "carscore" "cforest.importance" "chi.squared" "gain.ratio" "information.gain" "kruskal.test" "linear.correlation" "mrmr" "oneR" "permutation.importance" "randomForest.importance" "randomForestSRC.rfsrc" "randomForestSRC.var.select" "rank.correlation" "relief"

"symmetrical.uncertainty" "univariate.model.score" "variance"

### Feature selection



resampling=,measures=,control=) Uses a feature selection algorithm (control) to resample and build a model repeatedly using different feature sets each time in order to find the best set.

Available controls include:

- makeFeatSelControlExhaustive(max.features=)Try every combination of features up to optional max.features makeFeatSelControlRandom(maxit=,prob=,
- max.features=) Randomly sample features with probability prob (default 0.5) until maxit (default 100) iterations; return the best one found
- makeFeatSelControlSequential(method=, maxit=, max.features=,alpha=,beta=) Perform an iterative search using a method from the following: "sfs" forward search, "sbs" backward search, "sffs" floating forward search, "sfbs" floating backward search. alpha indicates minimum improvement required to add a feature; beta indicates minimum required to remove a feature
- makeFeatSelControlGA(maxit=,max.features=,mu=, lambda=, crossover.rate=, mutation.rate=) Genetic algorithm trains on random feature vectors, then uses crossover on the best performers to produce 'offspring', repeated over generations. mu is size of parent population, lambda is size of children population. crossover.rate is probability of choosing a bit from first parent, mutation.rate is probability of flipping a bit (on or off)

selectFeatures returns a FeatSelResult object which contains optimal features and an optimization path. To apply feature selection result (fsr) to your task (tsk), use: tsk = subsetTask(tsk,features=fsr\$x)

### **Benchmarking**

benchmark(learners=,tasks=,resamplings=,measures=) Allows easy comparison of multiple learners on a single task, a single learner on multiple tasks, or multiple learners on multiple tasks. Returns a benchmark result object.

Benchmark results can be accessed with a variety of functions beginning with getBMR<object>: AggrPerformance FeatSelResults FilteredFeatures LearnerIds LeanerShortNames Learners MeasureIds Measures Models Performances Predictions TaskDescs TaskIds TuneResults

mlr contains several toy tasks which are useful for benchmarking: agri.task bc.task bh.task costiris.task iris.task lung.task mtcars.task pid.task sonar.task wpbc.task yeast.task

### Visualization

### Performance

generateThreshVsPerfData(obj=,measures=) Measure performance at different probability cutoffs to determine optimal decision threshold for binary classification problems plotThreshVsPerf(obj) Plot visual

representation of threshold curve(s) from ThreshVsPerfData plotROCCurves(obj) Plot receiver

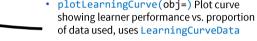
operating characteristic (ROC) curve from ThreshVsPerfData. Must set measures=list(fpr,tpr)

Residuals

 plotResiduals(obj=) Plots residuals for Prediction or BenchmarkResult

### Learning curve

generateLearningCurveData(learners=,task=, resampling=, percs=, measures=) Measure performance of learner(s) trained on different percentages of task data plotLearningCurve(obj=) Plot curve



### **Feature importance**

generateFilterValuesData(task=,method=) Get feature importance rankings using specified filter method

• plotFilterValues(obj=) Plot bar chart of feature importance based on filter method using FilterValuesData

### Hyperparameter tuning

generateHyperParsEffectData(tune.result=) Get the impact of different hyperparameter settings on model performance



t.data=, x=, y=, z=) Create a plot showing hyperparameter impact on performance using HyperParsEffectData

plotHyperParsEffect(hyperpars.effec

### See also:

- plotOptPath(op=) Display details of optimization process. Takes <obj>\$opt.path, where <obj> is an object of class tuneResult or featSelResult
- plotTuneMultiCritResult(res=) Show pareto front for results of tuning to multiple performance measures

### **Partial dependence**

generatePartialDependenceData(obj=,input=)Get partial dependence of model (obj) prediction over each feature of data (input)

plotPartialDependence(obj=) Plots partial dependence of model using PartialDependenceData

### Benchmarking

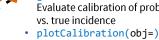
- plotBMRBoxplots (bmr=) Distribution of performances
- plotBMRSummary(bmr=) Scatterplot of avg. performances
- plotBMRRanksAsBarChart(bmr=) Rank learners in bar plot

### Other

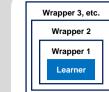
 generateCritDifferencesData(bmr=, measure=,p.value=,test=)Perform critical-differences test using either the Bonferroni-Dunn ("bd") or "Nemenyi" test з plotCritDifferences(obj=)



generateCalibrationData(obj=) Evaluate calibration of probability predictions



### Wrappers



Wrappers fuse a learner with additional functionality. mlr treats a learner with wrappers as a single learner, and hyperparameters of wrappers can be tuned jointly with underlying model parameters. Models trained with wrappers will apply them to new data.

### Preprocessing and imputation

makeDummvFeaturesWrapper(learner=) makeImputeWrapper(learner=, classes=, cols=) makePreprocWrapper(learner=,train=,predict=) makePreprocWrapperCaret(learner=,...) makeRemoveConstantFeaturesWrapper(learner=)

### **Class imbalance**

makeOverBaggingWrapper(learner=) makeSMOTEWrapper(learner=) makeUndersampleWrapper(learner=) makeWeightedClassesWrapper(learner=)

### **Cost-sensitive learning**

makeCostSensClassifWrapper(learner=) makeCostSensRegrWrapper(learner=) makeCostSensWeightedPairsWrapper(learner=)

### Multilabel classification

makeMultilabelBinaryRelevanceWrapper(learner=) makeMultilabelClassifierChainsWrapper(learner=) makeMultilabelDBRWrapper(learner=) makeMultilabelNestedStackingWrapper(learner=) makeMultilabelStackingWrapper(learner=)

### Other

makeBaggingWrapper(learner=) makeConstantClassWrapper(learner=) makeDownsampleWrapper(learner=,dw.perc=) makeFeatSelWrapper(learner=, resampling=, control=) makeFilterWrapper(learner=,fw.perc=,fw.abs=, fw.threshold=makeMultiClassWrapper(learner=)

makeTuneWrapper(learner=, resampling=, par.set=, control=)

### **Nested Resampling**

mlr supports nested resampling for complex operations such as tuning and feature selection through wrappers. In order to get a good estimate of generalization performance and avoid data leakage, both an outer (for tuning/feature selection) and an inner (for the base model) resampling process are advised.

• Outer resampling can be specified in resample or benchmark • Inner resampling can be specified in makeTuneWrapper, makeFeatSelWrapper,etc.

### **Ensembles**

"stack.nocv", "stack.cv" train super learner on results

• "compress" with a neural network for faster performance

makeStackedLearner(base.learners=, super.learner=, method=) Combines multiple learners to create an ensemble • **base.learners**=learners to use for initial predictions

• super.learner=learner to use for final prediction

method=how to combine base learner predictions:

• "average" simple average of all base learners

of base learners, with or without cross-validation

"hill.climb" search for optimal weighted average

### **R For Data Science** *Cheat Sheet*

xts

Learn R for data science Interactively at www.DataCamp.com



### xts

eXtensible Time Series (xts) is a powerful package that provides an extensible time series class, enabling uniform handling of many R time series classes by extending zoo.

Load the package as follows:

> library(xts)

#### xts Objects

xts objects have three main components:

- coredata: always a matrix for xts objects, while it could also be a vector for zoo objects
- index: vector of any Date, POSIXct, chron, yearmon, yeargtr, or DateTime classes
- xtsAttributes: arbitrary attributes

### Creating xts Objects

xts1 <- xts(x=1:10, order.by=Sys.Date()-1:10)</pre> data <- rnorm(5) dates <- seg(as.Date("2017-05-01"),length=5,by="days")</pre> xts2 <- xts(x=data, order.by=dates)</pre> > xts3 <- xts(x=rnorm(10),</pre> order.by=as.POSIXct(Sys.Date()+1:10), born=as.POSIXct("1899-05-08")) > xts4 <- xts(x=1:10, order.by=Sys.Date()+1:10)</pre> Convert To And From xts

> data (AirPassengers) > xts5 <- as.xts(AirPassengers)</p>

#### Import From Files

> dat <- read.csv(tmp file)</p> > xts(dat, order.by=as.Date(rownames(dat),"%m/%d/%Y")) dat zoo <- read.zoo(tmp file, index.column=0, sep=",",
format="%m/%d/%Y") > dat zoo <- read.zoo(tmp,sep=",",FUN=as.yearmon)</pre>

> dat xts <- as.xts(dat zoo)

### Inspect Your Data

<pre>&gt; core_data &lt;- coredata(xts2) &gt; index(xts1)</pre>	Extract core data of objects Extract index of objects	Periods, Periodicity & Timest	amps	Other Useful Functions	
Class Attributes <pre>&gt; indexClass(xts2) &gt; indexClass(convertIndex(xts,'POSIXct')) &gt; indexTZ(xts5) &gt; indexFormat(xts5) &lt;- "%Y-%m-%d" Time Zones</pre>	Get index class Replacing index class Get index class Change format of time display	<pre>&gt; periodicity(xts5) &gt; to.yearly(xts5) &gt; to.monthly(xts3) &gt; to.quarterly(xts5) &gt; to.period(xts5,period="quarters") &gt; to.period(xts5,period="years") &gt; nmonths(xts5) &gt; nquarters(xts5) &gt; nyears(xts5) &gt; make.index.unique(xts3,eps=1e-4)</pre>	Convert to yearly OHLC Count the months in xts5 Count the quarters in xts5 Count the years in xts5	<pre>&gt; .index(xts4) &gt; .indexwday(xts3) &gt; .indexhour(xts3) &gt; start(xts3) &gt; end(xts4) &gt; str(xts3) &gt; time(xts1) &gt; head(xts2) &gt; tail(xts2)</pre>	Extract raw numeric index of xts1 Value of week(day), starting on Sunday, in index of xts3 Value of hour in index of xts3 Extract first observation of xts3 Extract last observation of xts4 Display structure of xts3 Extract raw numeric index of xts1 First part of xts2 Last part of xts2
<pre>&gt; tzone(xts1) &lt;- "Asia/Hong_Kong" &gt; tzone(xts1)</pre>	Change the time zone Extract the current time zone			Data	

### Export xts Objects

data xts <- as.xts(matrix)</pre> tmp <- tempfile()</pre> write.zoo(data xts, sep=", ", file=tmp)

### **Replace & Update**

> xts2[dates] <- 0 > xts5["1961"] <- NA > xts2["2016-05-02"] <- NA

### Applying Functions

> ep1 <- endpoints(xts4,on="weeks",k=2) Take index values by time 0 5 10 > ep2 <- endpoints(xts5,on="years") [1] 0 12 24 36 48 60 72 84 96 108 120 132 144 Calculate the yearly mean > period.apply(xts5, INDEX=ep2, FUN=mean) > xts5 yearly <- split(xts5,f="years")</pre> Split xts5 by year Create a list of yearly means > lapply(xts5 yearly,FUN=mean) Find the last observation in > do.call(rbind, lapply(split(xts5, "years"), each year in xts5 function(w) last(w, n="1 month"))) Calculate cumulative annual do.call(rbind, lapply(split(xts5, "years"), passengers cumsum)) > rollapply(xts5, 3, sd) Apply sd to rolling margins of xts5

Replace values in xts2 on dates with 0

Replace the value at 1 specific index with NA

Index of weekend days

Extract weekend days of xts1

Replace dates from 1961 with NA

### Selecting, Subsetting & Indexing

> index <- which(.indexwday(xts1)==0|.indexwday(xts1)==6)</pre>

xts1[index]

#### Select > mar55 <- xts5["1955-03"] Get value for March 1955 Subset > xts5 1954 <- xts5["1954"] Get all data from 1954 xts5 janmarch <- xts5["1954/1954-03"] Extract data from Jan to March '54 xts5 janmarch <- xts5["/1954-03"] Get all data until March '54 xts4[ep1] Subset xts4 using ep2 first() and last() > first(xts4, '1 week') Extract first 1 week > first(last(xts4, '1 week'), '3 days') Get first 3 days of the last week of data Indexing > xts2[index(xts3)] Extract rows with the index of xts3 days <- c("2017-05-03", "2017-05-23") Extract rows using the vector days xts3[days] > xts2[as.POSIXct(days,tz="UTC")] Extract rows using days as POSIXct

### **Missing Values** na.omit(xts5) > xts last <- na.locf(xts2)</pre> xts last <- na.locf(xts2,</pre> fromLast=TRUE) na.approx(xts2)

#### Omit NA values in xts5 Fill missing values in xts2 using last observation Fill missing values in xts2 using next observation Interpolate NAs using linear approximation

### **Arithmetic Operations**

#### coredata()Of as.numeric()

xts3 + as.numeric(xts2) Addition > xts3 \* as.numeric(xts4) Multiplication > coredata(xts4) - xts3 Subtraction > coredata(xts4) / xts3 Division

#### Shifting Index Values

xts5 - lag(xts5) Period-over-period differences diff(xts5, lag=12, differences=1) Lagged differences

### Reindexing

el 2017-05-04 5.231538 2017-05-05 5.829257	
201/-03-03 3.02923/	
2017-05-06 4.000000	
2017-05-07 3.000000 2017-05-08 2.000000	
2017-05-09 1.000000	
> xts1 - merge(xts2,index(xts1),fill=na.locf) Subtraction	1
e1 2017-05-04 5.231538	
2017-05-04 5.231538	
2017-05-06 4.829257	
2017-05-07 3.829257	
2017-05-08 2.829257 2017-05-09 1.829257	

### Meraina

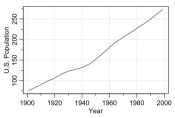
<b></b>	
<pre>&gt; merge(xts2,xts1,join='inner')</pre>	Inner join of xts2 and xts1
2017-05-05 -0.8382068 10	
<pre>&gt; merge(xts2,xts1,join='left',fill=0)</pre>	Left join of xts2 and xts1, fill empty spots with 0
xts2 xts1 2017-05-01 1.7482704 0 2017-05-02 -0.2314678 0 2017-05-03 0.1685517 0 2017-05-04 1.1685649 0 2017-05-05 -0.8382068 10	
> rbind(xts1, xts4)	Combine xts1 and xts4 by rows

Learn R for Data Science Interactively

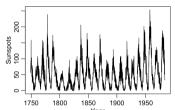
# **Time Series Cheat Sheet**

### **Plot Time Series**

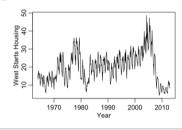
1. tsplot(x=time, y=data)



2. plot(ts(data, start=start\_time, frequency=gap))



3. ts.plot(ts(data, start=start\_time, frequency=gap))



### Simulation

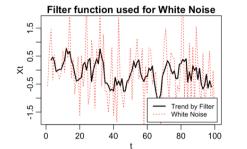
 $\begin{aligned} &\textbf{Autoregression of Order p} \\ &X_t = \phi_1 X_{t-1} + \phi_2 X_{t-2} + \ldots + \phi_p X_{t-p} + W_t \\ &\textbf{Moving Average of Order q} \\ &X_t = Z_t + \theta_1 Z_{t-1} + \theta_2 Z_{t-2} + \ldots + \theta_q Z_{t-p} \\ &\textbf{ARMA (p, q)} \\ &X_t = \phi_1 X_{t-1} + \phi_2 X_{t-2} + \ldots + \phi_p X_{t-p} + \\ &Z_t + \theta_1 Z_{t-1} + \theta_2 Z_{t-2} + \ldots + \theta_q Z_{t-p} \end{aligned}$ 

### Simulation of ARMA (p, q)

arima.sim(model=list(ar=c( $\phi_1, \dots, \phi_p$ ), ma=c( $\theta_1, \dots, \theta_q$ )), n=n)

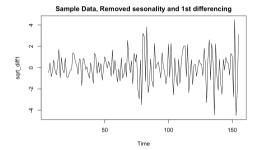
### Filters

- Linear Filter: filter()
- filter(data, filter=filter\_coefficients, sides=2, method="convolution", circular=F)



Differencing Filter: diff()

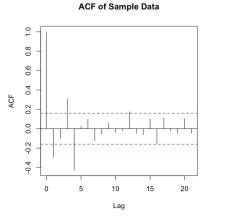
diff(data, lag=4, differences=1)



### Auto-correlation

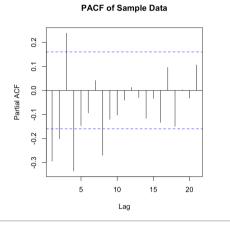
Use ACF and PACF to detect model

(Complete) Auto-correlation function: acf() acf(data, type='correlation', na.action=na.pass)



Partial Auto-correlation function: pacf() pacf(data, na.action=na.pass)

**OR:** acf(data, type='partial', na.action=na.pass)



### **Parameter Estimation**

it an ARMA time series model to the data

**ar():** To estimate parameters of an AR model ar(x=data, aic=T, order.max = NULL,

c("yule-walker", "burg", "ols", "mle", "yw"))

Call: ar(x = sqrt1, aic = TRUE, order.max = NULL, method = c("yule-walker", "burg", "ols", "mle", "y Coefficients: 1 2 3 4 5 6 7 8 9 -0.3066 -0.1903 0.0793 -0.5065 -0.1873 -0.1149 -0.0580 -0.3031 -0.1207 Order selected 9 sigma^2 estimated as 1.52

# **arima():** To estimate parameters of an AM or ARMA model, and build model

arima(data, order=c(p, o, q),method=c('ML'))

```
Call:
arima(x = sqrt1, order = c(2, 0, 6), method = c("ML"))
```

Coefficients: ar1 ar2 ma1 ma2 ma3 ma4 ma5 ma6 intercept -0.1658 -0.7951 -0.1536 0.7524 -0.2101 -0.7680 0.0605 -0.6812 -0.0048 s.e. 0.0918 0.0754 0.0916 0.1047 0.0794 0.0743 0.0812 0.0895 0.0062 sigma^2 estimated as 1.193: log likelihood = -230.78, aic = 481.55

```
AICc(): Compare models using AICC
```

AICc(fittedModel)

# Forecasting

Forecasting future observations given a fitted ARMA model

**predict():** Predict future observations given a fitted ARMA model

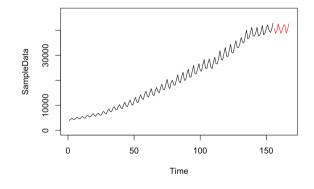
predict(arima\_model, number\_to\_predict)

### Plot Predicted values and Confidence Interval:

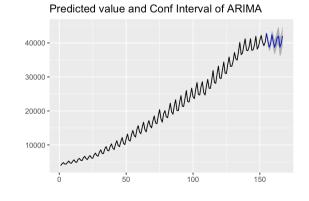
fit<-predict(arima\_model, number\_to\_predict)
ts.plot(data,</pre>

xlim=c(1, length(data)+number\_to\_predict), ylim=c(0, max(fit\$pred+1.96\*fit\$se))) lines(length(data)+1:length(data)+

number\_to\_predict, fit\$pred)



**OR:** autoplot(forecast(arima\_model, level=c(95), h=number\_to\_predict))



# Deep Learning with Keras :: **CHEAT SHEET**



### Intro

**DEFINE A MODEL** 

a linear stack of layers

**COMPILE A MODEL** 

**FIT A MODEL** 

(iterations)

samples

generator

by-batch by a generator

**EVALUATE A MODEL** 

NULL) Evaluate a Keras model

Studio

GPUs

keras\_model() Keras Model

Keras is a high-level neural networks API developed with a focus on enabling fast experimentation. It supports multiple backends, including TensorFlow, CNTK and Theano.

TensorFlow is a lower level mathematical library for building deep neural network architectures. The keras R package makes it easy to use Keras and TensorFlow in R.

Working with keras models

keras\_model\_sequential() Keras Model composed of

multi\_gpu\_model() Replicates a model on different

compile(object, optimizer, loss, metrics = NULL)

**fit**(object, x = NULL, y = NULL, batch size = NULL,

Train a Keras model for a fixed number of epochs

fit\_generator() Fits the model on data yielded batch-

train on batch() test on batch() Single gradient

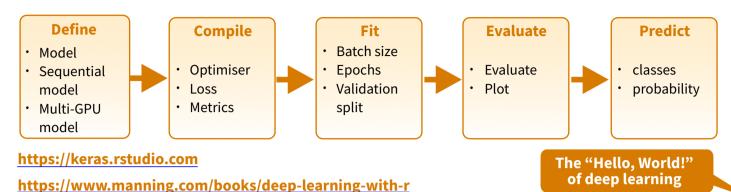
evaluate(object, x = NULL, y = NULL, batch\_size =

evaluate\_generator() Evaluates the model on a data

update or model evaluation over one batch of

epochs = 10, verbose = 1, callbacks = NULL, ...)

Configure a Keras model for training



### INSTALLATION

The keras R package uses the Python keras library. You can install all the prerequisites directly from R. https://keras.rstudio.com/reference/install keras.html

library(keras) See ?install\_keras install\_keras() for GPU instructions

This installs the required libraries in an Anaconda environment or virtual environment 'r-tensorflow'.

### TRAINING AN IMAGE RECOGNIZER ON MNIST DATA

#### *#* input layer: use MNIST images Ч 0 mnist <- dataset mnist()</pre>

x\_train <- mnist\$train\$x; y\_train <- mnist\$train\$y</pre>

x\_test <- mnist\$test\$x; y\_test <- mnist\$test\$y</pre>

### # reshape and rescale

 $x_train <- array_reshape(x_train, c(nrow(x_train), 784))$ x test <- array reshape(x test, c(nrow(x test), 784))x train <- x train / 255; x test <- x test / 255

y\_train <- to\_categorical(y\_train, 10)</pre> y test <- to categorical(y test, 10)

model <- keras model sequential() model %>% layer\_dense(units = 256, activation = 'relu', input shape = c(784)) %>% layer\_dropout(rate = 0.4) %>% layer dense(units = 128, activation = 'relu') %>% layer dense(units = 10, activation = 'softmax')

### # compile (define loss and optimizer)

model %>% compile( loss = 'categorical\_crossentropy', optimizer = optimizer\_rmsprop(), metrics = c('accuracy')

### # train (fit)

model %>% fit( x train, y train, epochs = 30, batch size = 128, validation\_split = 0.2

model %>% evaluate(x\_test, y\_test) model %>% predict\_classes(x\_test)

### PREDICT

predict() Generate predictions from a Keras model

predict proba() and predict classes() Generates probability or class probability predictions

input samples from a data generator

### **OTHER MODEL OPERATIONS**

**pop\_layer()** Remove the last layer in a model

save\_model\_hdf5(); load\_model\_hdf5() Save/ Load models using HDF5 files

serialize\_model(); unserialize\_model() Serialize a model to an R object

clone\_model() Clone a model instance

freeze\_weights(); unfreeze\_weights() Freeze and unfreeze weights

for the input samples

**predict on batch()** Returns predictions for a single batch of samples

predict\_generator() Generates predictions for the

summary() Print a summary of a Keras model

export\_savedmodel() Export a saved model

get\_layer() Retrieves a layer based on either its name (unique) or index



**CORE LAYERS** 





















activity





f(x)



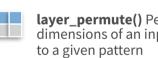




activation function to an output

layer\_dropout() Applies Dropout

layer permute() Permute the dimensions of an input according to a given pattern



to the input

layer repeat vector() Repeats the input n times

layer\_input() Input layer

**layer\_lambda**(object, f) Wraps arbitrary expression as a layer

layer\_activity\_regularization() Layer that applies an update to the cost function based input

layer\_masking() Masks a sequence by using a mask value to skip timesteps

layer\_flatten() Flattens an input

RStudio<sup>®</sup> is a trademark of RStudio, Inc. • CC BY SA RStudio • info@rstudio.com • 844-448-1212 • rstudio.com • Learn more at keras.rstudio.com • keras 2.1.2 • Updated: 2017-12

### # defining the model and layers

### More layers

### **CONVOLUTIONAL LAYERS**



**layer\_conv\_1d()** 1D, e.g. temporal convolution

**layer\_conv\_2d\_transpose()** Transposed 2D (deconvolution)

**layer\_conv\_2d()** 2D, e.g. spatial convolution over images



layer\_conv\_3d\_transpose()
Transposed 3D (deconvolution)
layer\_conv\_3d() 3D, e.g. spatial
convolution over volumes

layer\_conv\_lstm\_2d() Convolutional LSTM

**layer\_separable\_conv\_2d()** Depthwise separable 2D



layer\_upsampling\_1d()
layer\_upsampling\_2d()
layer\_upsampling\_3d()
Upsampling layer



layer\_zero\_padding\_1d() layer\_zero\_padding\_2d() layer\_zero\_padding\_3d() Zero-padding layer

layer\_cropping\_1d()

layer\_cropping\_2d()

### layer\_cropping\_3d() Cropping layer

### **POOLING LAYERS**

layer\_max\_pooling\_1d() layer\_max\_pooling\_2d() layer\_max\_pooling\_3d() Maximum pooling for 1D to 3D





Studio

layer\_global\_max\_pooling\_1d() layer\_global\_max\_pooling\_2d() layer\_global\_max\_pooling\_3d() Global maximum pooling

layer\_global\_average\_pooling\_1d() layer\_global\_average\_pooling\_2d() layer\_global\_average\_pooling\_3d() Global average pooling

### **ACTIVATION LAYERS**

layer\_activation(object, activation) Apply an activation function to an output

- layer\_activation\_leaky\_relu() Leaky version of a rectified linear unit
- layer\_activation\_parametric\_relu()
  Parametric rectified linear unit

layer\_activation\_thresholded\_relu() Thresholded rectified linear unit

layer\_activation\_elu() Exponential linear unit

### **DROPOUT LAYERS**

layer\_dropout() Applies dropout to the input

layer\_spatial\_dropout\_1d()
layer\_spatial\_dropout\_2d()
layer\_spatial\_dropout\_3d()
Spatial 1D to 3D version of dropout

### **RECURRENT LAYERS**

layer\_simple\_rnn() Fully-connected RNN where the output is to be fed back to input

**layer\_gru()** Gated recurrent unit - Cho et al

**layer\_cudnn\_gru()** Fast GRU implementation backed by CuDNN

**layer\_lstm()** Long-Short Term Memory unit -Hochreiter 1997

**layer\_cudnn\_lstm()** Fast LSTM implementation backed by CuDNN

### LOCALLY CONNECTED LAYERS

**layer\_locally\_connected\_1d() layer\_locally\_connected\_2d()** Similar to convolution, but weights are not shared, i.e. different filters for each patch

### Preprocessing

### SEQUENCE PREPROCESSING

**pad\_sequences()** Pads each sequence to the same length (length of the longest sequence)

**skipgrams()** Generates skipgram word pairs

**make\_sampling\_table()** Generates word rank-based probabilistic sampling table

### **TEXT PREPROCESSING**

text\_tokenizer() Text tokenization utility

fit\_text\_tokenizer() Update tokenizer internal
vocabulary

save\_text\_tokenizer(); load\_text\_tokenizer()
Save a text tokenizer to an external file

texts\_to\_sequences();
texts\_to\_sequences\_generator()
Transforms each text in texts to sequence of integers

texts\_to\_matrix(); sequences\_to\_matrix()
Convert a list of sequences into a matrix

text\_one\_hot() One-hot encode text to word indices

**text\_hashing\_trick()** Converts a text to a sequence of indexes in a fixedsize hashing space

text\_to\_word\_sequence()
Convert text to a sequence of words (or tokens)

### **IMAGE PREPROCESSING**

image\_load() Loads an image into PIL format.

flow\_images\_from\_data()
flow\_images\_from\_directory()
Generates batches of augmented/normalized data
from images and labels, or a directory

**image\_data\_generator()** Generate minibatches of image data with real-time data augmentation.

fit\_image\_data\_generator() Fit image data generator internal statistics to some sample data

generator\_next() Retrieve the next item

image\_to\_array(); image\_array\_resize()
image\_array\_save() 3D array representation



### **Pre-trained models**

Keras applications are deep learning models that are made available alongside pre-trained weights. These models can be used for prediction, feature extraction, and fine-tuning.

application\_xception() xception\_preprocess\_input() Xception v1 model

**application\_inception\_v3() inception\_v3\_preprocess\_input()** Inception v3 model, with weights pre-trained on ImageNet

application\_inception\_resnet\_v2() inception\_resnet\_v2\_preprocess\_input() Inception-ResNet v2 model, with weights trained on ImageNet

application\_vgg16(); application\_vgg19()
VGG16 and VGG19 models

application\_resnet50() ResNet50 model

application\_mobilenet()
mobilenet\_preprocess\_input()
mobilenet\_decode\_predictions()
mobilenet\_load\_model\_hdf5()
MobileNet model architecture

IMAGENET ImageNet is a large database of images with labels, extensively used for deep learning

imagenet\_preprocess\_input()
imagenet\_decode\_predictions()
Preprocesses a tensor encoding a batch of
images for ImageNet, and decodes predictions

### Callbacks

visualizations

A callback is a set of functions to be applied at given stages of the training procedure. You can use callbacks to get a view on internal states and statistics of the model during training.

callback\_early\_stopping() Stop training when
a monitored quantity has stopped improving
callback\_learning\_rate\_scheduler() Learning
rate scheduler
callback\_tensorboard() TensorBoard basic

### Need to Know

Pattern arguments in stringr are interpreted as regular expressions after any special characters have been parsed.

In R, you write regular expressions as strings, sequences of characters surrounded by quotes ("") or single quotes(").

Some characters cannot be represented directly in an R string. These must be represented as **special characters**, sequences of characters that have a specific meaning., e.g.

Special Character	Represents			
//	\			
\"	н			
∖n	new line			
Run ?""" to see a complete list				

Because of this, whenever a \ appears in a regular expression, you must write it as \\ in the string that represents the regular expression.

Use writeLines() to see how R views your string after all special characters have been parsed.

writeLines("\\.") #\.

writeLines("\\ is a backslash") #\is a backslash

### **INTERPRETATION**

Patterns in stringr are interpreted as regexs. To change this default, wrap the pattern in one of:

regex(pattern, ignore\_case = FALSE, multiline = FALSE, comments = FALSE, dotall = FALSE, ...) Modifies a regex to ignore cases, match end of lines as well of end of strings, allow R comments within regex's, and/or to have . match everything including \n. str\_detect("I", regex("i", TRUE))

fixed() Matches raw bytes but will miss some characters that can be represented in multiple ways (fast). str\_detect("\u0130", fixed("i"))

**coll()** Matches raw bytes and will use locale specific collation rules to recognize characters that can be represented in multiple ways (slow). str\_detect("\u0130", coll("i", TRUE, locale = "tr"))

**boundary()** Matches boundaries between characters, line\_breaks, sentences, or words. str\_split(sentences, boundary("word"))



<b>Regular Expressions -</b>	Regular expressions, or <i>regexps</i> , are a concise language for describing patterns in strings.
------------------------------	---

MATCH C	HARACTERS	see <- function(rx)	str_view_all("abc AB	3C 123\t.!?\\(){}\n", rx)
string	regexp	matches	example	
(type this)	(to mean this)	(which matches this)		
	a (etc.)	a (etc.)	see("a")	abc ABC 123 .!?\(){}
//.	\.		see("\\.")	abc ABC 123 <mark>.</mark> !?\(){}
\\!	\!	!	see("\\!")	abc ABC 123 <mark>.!</mark> ?\(){}
\\?	\?	?	see("\\?")	abc ABC 123 .! <mark>?</mark> \(){}
1111	W		see("\\\\")	abc ABC 123 .!?\(){}
\\(	\ <b>(</b>	(	see("\\(")	abc ABC 123 .!?\ <mark>(</mark> ){}
\\ <b>)</b>	V	)	see("\\)")	abc ABC 123 .!?\( <mark>)</mark> {}
\\ <b>{</b>	\{ {	{	see("\\{")	abc ABC 123 .!?\() <mark>{</mark> }
\\}	\}	}	see( "\\}")	abc ABC 123 .!?\(){ <mark>}</mark>
\\ <b>n</b>	\n	new line (return)	see("\\n")	abc ABC 123 .!?\(){}
\\ <b>t</b>	\t	tab	see("\\t")	abc ABC 123 .!?\(){}
\ <b>\s</b>	\s	any whitespace (\ <b>S</b> for non-whitespaces)	see("\\s")	abc ABC 123 .!?\(){}
//d	\d	any digit (\ <b>D</b> for non-digits)	see("\\d")	abc ABC <mark>123</mark> .!?\(){}
\\w	\w	any word character (\W for non-word chars)	see("\\w")	abc ABC 123 .!?\(){}
\\b	\ <b>b</b>	word boundaries	see("\\b")	abc ABC 123 .!?\(){}
	[:digit:] <sup>1</sup>	digits	see("[:digit:]")	abc ABC <mark>123</mark> .!?\(){}
	[:alpha:] <sup>1</sup>	letters	see("[:alpha:]")	abc ABC 123 .!?\(){}
	[:lower:] <sup>1</sup>	lowercase letters	see("[:lower:]")	abc ABC 123 .!?\(){}
	[:upper:] <sup>1</sup>	uppercase letters	see("[:upper:]")	abc <mark>ABC</mark> 123 .!?\(){}
	[:alnum:] <sup>1</sup>	letters and numbers	see("[:alnum:]")	abc ABC 123 .!?\(){}
	[:punct:] <sup>1</sup>	punctuation	see("[:punct:]")	abc ABC 123 <mark>.!?\(){}</mark>
	[:graph:]	letters, numbers, and punctuation	see("[:graph:]")	abc <mark>ABC</mark> 123 .!?\(){}
	[:space:]	space characters (i.e. \s)	see("[:space:]")	abc <mark>ABC</mark> 123 .!?\(){}
	[:blank:] <sup>1</sup>	space and tab (but not new line)	see("[:blank:]")	abc ABC 123 .!?\(){}
	•	every character except a new line	see(".")	abc ABC 123 .!?\(){}
		<sup>1</sup> Many base R functions require classes to be v	vrapped in a second se	et of [ ], e.g. [[ <b>:digit:</b> ]]

uase κ runctions require classes to be wrapped in a second set of [], e.g. [[:digit:]]

ALTERNATES		alt <- function	n(rx) str_view_all("a	abcde", rx)
	regexp	matches	example	
	abd	or	alt("ab d")	abcde
	abe	one of	alt("[abe]")	abcde
	[^abe]	anything but	alt("[^abe]")	abcde
	a-c	range	alt("[a-c]")	abcde
ANCHORS		anchor <- funct	ion(rx) str_view_all	("aaa", rx)
	regexp	matches	example	
	<b>∧</b> a	start of string	anchor("^a")	aaa
	a\$	end of string	anchor("a\$")	aaa j

QUANTIFIERS	quant <- function(rx) str_view_all(".a.aa.aaa", rx)				
	regexp	matches	example		
	a <b>?</b>	zero or one	quant("a?")	.a.aa.aaa	
	a*	zero or more	quant("a*")	.a.aa.aaa	
	a <b>+</b>	one or more	quant("a+")	. <mark>a.aa.aaa</mark>	
1-2n-	a{n}	exactly <b>n</b>	quant("a{2}")	.a. <mark>aa</mark> .aaa	
1 2 n -	a{n, }	<b>n</b> or more	quant("a{2,}")	.a. <mark>aa</mark> .aaa	
- nm-	a{n, m}	between <b>n</b> and <b>m</b>	quant("a{2,4}")	.a. <mark>aa</mark> .aaa	

GROUPS	5	ref <- f	unction(rx) str_view_al	l("abbaab", rx)
Use parer	ntheses to set p	precedent (order of ev	aluation) and create gro	oups
	regexp (ab d)e	matches sets precedence	example alt("(ab d)e")	abc <mark>de</mark>
		to refer to and duplic r to each group by its	ate parentheses groups order of appearance	s that occur
string (type this)	regexp (to mean this)	matches (which matches this)	example (the result is the same as	ref("abba"))
\\1	<b>\1</b> (etc.)	first () group, etc.	ref("(a)(b)\\2\\1")	<mark>abba</mark> ab

LOOK AROUNDS		look <- function(rx) str_view_all("bacad", rx)					
	regexp	matches	example		Use an es		
	a(?=c)	followed by	look("a(?=c)")	b <mark>a</mark> cad	earlier in		
	a(?!c)	not followed by	look("a(?!c)")	bac <mark>a</mark> d	string		
	(?<=b)a	preceded by	look("(?<=b)a")	b <mark>a</mark> cad	(type this)		
	(? b)a</th <th>not preceded by</th> <th>look("(?<!--b)a")</th--><th>bac<mark>a</mark>d</th><th>\\1</th></th>	not preceded by	look("(? b)a")</th <th>bac<mark>a</mark>d</th> <th>\\1</th>	bac <mark>a</mark> d	\\1		

[:space:]			23C					
[:blank:] space tab								
[:graph:]								
[:punct:] [:symbol:]								
	/ * @ # { } ( )	` = ~ < >	• \$					
[:	alnum:]							
	[:digit:] 0 1 2 3 4 5 6 7 8 9							
I	:alpha:]							
[:lower:]	[:	upper:]						
a b c d e	f A B	C D E	F					
ghijk		IJK						
mnopq		ΟΡQ						
stuvw		UVW	Х					
y z	ΥZ							

# Data visualization with ggplot2 :: CHEAT SHEET

Each function returns a layer.

### Basics

ggplot2 is based on the grammar of graphics, the idea that you can build every graph from the same components: a data set, a coordinate system, and geoms-visual marks that represent data points.



To display values, map variables in the data to visual properties of the geom (aesthetics) like size, color, and x and **v** locations.



Complete the template below to build a graph.

ggplot (data = <data>) +</data>	required
<pre><geom_function>(mapping = aes(<mappin)< pre=""></mappin)<></geom_function></pre>	GS>),
stat = <b><stat></stat></b> , position = <b><position></position></b> ) +	Not
<coordinate_function> +</coordinate_function>	required, sensible
<facet_function> +</facet_function>	defaults
<scale_function> +</scale_function>	supplied
<theme_function></theme_function>	

ggplot(data = mpg, aes(x = cty, y = hwy)) Begins a plot that you finish by adding layers to. Add one geom function per laver.

last\_plot() Returns the last plot.

ggsave("plot.png", width = 5, height = 5) Saves last plot as 5' x 5' file named "plot.png" in working directory. Matches file type to file extension.

### Aes Common aesthetic values.

**color** and **fill** - string ("red", "#RRGGBB")

**linetype** - integer or string (0 = "blank", 1 = "solid", 2 = "dashed", 3 = "dotted", 4 = "dotdash", 5 = "longdash", 6 = "twodash")

**lineend** - string ("round", "butt", or "square")

linejoin - string ("round", "mitre", or "bevel")

0 1 2 3 4 5 6 7 8 9 10 11 12 **size** - integer (line width in mm) shape - integer/shape name or 13 14 15 16 17 18 19 20 21 22 23 24 25 

# Studio

**GRAPHICAL PRIMITIVES** a <- ggplot(economics, aes(date, unemploy))  $b \le ggplot(seals, aes(x = long, y = lat))$ 

a + geom blank() and a + expand limits() Ensure limits include values across all plots.

**b**+geom\_curve(aes(yend = lat + 1, xend = long + 1), curvature = 1) - x, xend, y, yend, alpha, angle, color, curvature, linetype, size

a + geom path(lineend = "butt", linejoin = "round", linemitre = 1) x, y, alpha, color, group, linetype, size

a + geom\_polygon(aes(alpha = 50)) - x, y, alpha, color, fill, group, subgroup, linetype, size

**b** + geom\_rect(aes(xmin = long, ymin = lat, xmax = long + 1, ymax = lat + 1)) - xmax, xmin. ymax, ymin, alpha, color, fill, linetype, size

**a + geom ribbon(**aes(ymin = unemploy - 900, ymax = unemploy + 900)) - x, ymax, ymin, alpha, color, fill, group, linetype, size

### LINE SEGMENTS

Geoms

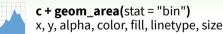
common aesthetics: x, y, alpha, color, linetype, size

**b + geom\_abline**(aes(intercept = 0, slope = 1)) **b** + geom\_hline(aes(yintercept = lat)) **b** + geom\_vline(aes(xintercept = long))

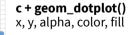
**b + geom\_segment(**aes(yend = lat + 1, xend = long + 1)) **b + geom\_spoke(**aes(angle = 1:1155, radius = 1))

### **ONE VARIABLE** continuous

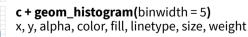
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)

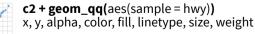


c + geom\_density(kernel = "gaussian") x, y, alpha, color, fill, group, linetype, size, weight



c + geom\_freqpoly() x, y, alpha, color, group, linetype, size





### discrete

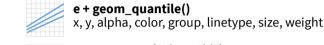
d <- ggplot(mpg, aes(fl)) d + geom bar() x, alpha, color, fill, linetype, size, weight

#### **TWO VARIABLES** both continuous e <- ggplot(mpg, aes(cty, hwy))

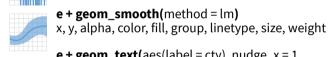
Use a geom function to represent data points, use the geom's aesthetic properties to represent variables.

**e + geom\_label(**aes(label = cty), nudge\_x = 1,  $nudge_y = 1$ ) - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

e + geom\_point() x, y, alpha, color, fill, shape, size, stroke

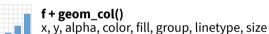


e + geom\_rug(sides = "bl") x, y, alpha, color, linetype, size



**e + geom text(**aes(label = cty), nudge x = 1, C  $nudge_y = 1$ ) - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust A<sub>B</sub>

### one discrete. one continuous f <- ggplot(mpg, aes(class, hwy))



### f + geom\_boxplot()

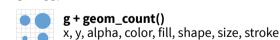
x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight

f + geom dotplot(binaxis = "y", stackdir = "center") x, y, alpha, color, fill, group



f + geom\_violin(scale = "area") x, y, alpha, color, fill, group, linetype, size, weight

#### both discrete g <- ggplot(diamonds, aes(cut, color))

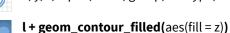


e + geom\_jitter(height = 2, width = 2) x, y, alpha, color, fill, shape, size er 197

### **THREE VARIABLES**

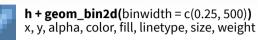
seals\$z <- with(seals, sqrt(delta\_long^2 + delta\_lat^2)); l <- ggplot(seals, aes(long, lat))</pre>

l + geom\_contour(aes(z = z)) x, y, z, alpha, color, group, linetype, size, weight



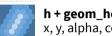
x, y, alpha, color, fill, group, linetype, size, subgroup

### continuous bivariate distribution h <- ggplot(diamonds, aes(carat, price))



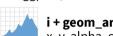
h + geom density 2d()x, y, alpha, color, group, linetype, size

ggplot



h + geom hex() x, y, alpha, color, fill, size

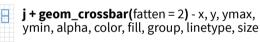
### continuous function i <- ggplot(economics, aes(date, unemploy))



- i + geom area() x, y, alpha, color, fill, linetype, size
- i + geom\_line() x, y, alpha, color, group, linetype, size
- **i + geom** step(direction = "hv") x, y, alpha, color, group, linetype, size

### visualizing error

df < - data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))



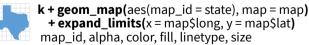
j + geom\_errorbar() - x, ymax, ymin, alpha, color, group, linetype, size, width Also geom\_errorbarh().



**j + geom\_pointrange() -** x, y, ymin, ymax, alpha, color, fill, group, linetype, shape, size

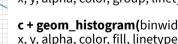
### maps

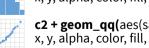
data <- data.frame(murder = USArrests\$Murder, state = tolower(rownames(USArrests))) map <- map\_data("state")</pre> k < - ggplot(data, aes(fill = murder))



l + geom\_raster(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE) x, y, alpha, fill

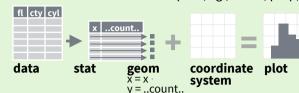
l + geom\_tile(aes(fill = z)) x, y, alpha, color, fill, linetype, size, width





### Stats An alternative way to build a layer.

A stat builds new variables to plot (e.g., count, prop).



Visualize a stat by changing the default stat of a geom function, geom\_bar(stat="count") or by using a stat function, stat\_count(geom="bar"), which calls a default geom to make a layer (equivalent to a geom function). Use ...name.. syntax to map stat variables to aesthetics.



**c + stat bin(**binwidth = 1, boundary = 10) **x**, **y** | ...count.., ..ncount.., ..density.., ..ndensity.. c + stat\_count(width = 1) x, y | ...count..., ...prop...

c + stat density(adjust = 1, kernel = "gaussian") **x, y** | ...count.., ..density.., ...scaled..

e + stat\_bin\_2d(bins = 30, drop = T) **x, y, fill** ...count.., ..density..

e + stat\_bin\_hex(bins = 30) x, y, fill | ...count.., ..density..

e + stat\_density\_2d(contour = TRUE, n = 100) x, y, color, size ... level..

e + stat\_ellipse(level = 0.95, segments = 51, type = "t")

**l** + stat\_contour(aes(z = z)) x, y, z, order | ...level...

l + stat\_summary\_hex(aes(z = z), bins = 30, fun = max) x, y, z, fill | ..value..

l + stat\_summary\_2d(aes(z = z), bins = 30, fun = mean) x, y, z, fill | ..value..

f + stat\_boxplot(coef = 1.5)

**x**, **y** | ..lower.., ..middle.., ..upper.., ..width.. , ..ymin.., ..ymax.. f + stat\_ydensity(kernel = "gaussian", scale = "area") x, y

...density.., ..scaled.., ..count.., ..n.., ..violinwidth.., ..width..

**e + stat\_ecdf(**n = 40) **x, y** | ...x., ...y..

e + stat\_quantile(quantiles = c(0.1, 0.9), formula =  $y \sim log(x)$ , method = "rq") **x**, **y** | ...quantile...

e + stat\_smooth(method = "lm", formula = y ~ x, se = T, level = 0.95) **x**, **y** | ...se.., ...x.., ...y.., ...ymin..., ...ymax...

ggplot() + xlim(-5, 5) + stat\_function(fun = dnorm, n = 20, geom = "point") x | ...x.., ...y..

ggplot() + stat\_qq(aes(sample = 1:100)) x, y, sample | ...sample.., ..theoretical..

e + stat\_sum() x, y, size | ..n., ..prop..

e + stat summary(fun.data = "mean cl boot")

**h** + stat summary bin(fun = "mean", geom = "bar")

e + stat\_identity()

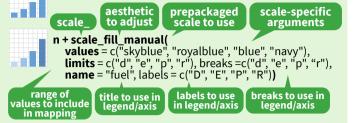
e + stat\_unique()





**Scales** map data values to the visual values of an aesthetic. To change a mapping, add a new scale.

 $n \le d + geom bar(aes(fill = fl))$ 



### **GENERAL PURPOSE SCALES**

### Use with most aesthetics

scale\_\*\_continuous() - Map cont' values to visual ones. scale \* discrete() - Map discrete values to visual ones. **scale** \* **binned()** - Map continuous values to discrete bins. scale\_\*\_identity() - Use data values as visual ones. scale\_\*\_manual(values = c()) - Map discrete values to manually chosen visual ones. scale\_\*\_date(date\_labels = "%m/%d"), date\_breaks = "2 weeks") - Treat data values as dates. scale\_\*\_datetime() - Treat data values as date times. Same as scale\_\*\_date(). See ?strptime for label formats.

### **X & Y LOCATION SCALES**

Use with x or y aesthetics (x shown here) scale\_x\_log10() - Plot x on log10 scale. scale\_x\_reverse() - Reverse the direction of the x axis. scale\_x\_sqrt() - Plot x on square root scale.

### **COLOR AND FILL SCALES (DISCRETE)**

n + scale\_fill\_brewer(palette = "Blues") For palette choices:

RColorBrewer::display.brewer.all() **n + scale\_fill\_grey(**start = 0.2,

end = 0.8, na.value = "red") 

### **COLOR AND FILL SCALES (CONTINUOUS)**

o <- c + geom\_dotplot(aes(fill = ..x..))</pre>

..... o + scale\_fill\_distiller(palette = "Blues")

o + scale fill gradient(low="red", high="vellow")

o + scale\_fill\_gradient2(low = "red", high = "blue", mid = "white", midpoint = 25)

o + scale\_fill\_gradientn(colors = topo.colors(6)) Also: rainbow(), heat.colors(), terrain.colors(), cm.colors(), RColorBrewer::brewer.pal()

### SHAPE AND SIZE SCALES

•

p <- e + geom\_point(aes(shape = fl, size = cyl))</pre>

p + scale\_shape() + scale\_size()  $\langle \rangle$ p + scale\_shape\_manual(values = c(3:7))

 $+_{X}$ 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25

p + scale\_radius(range = c(1,6)) 

p + scale\_size\_area(max\_size = 6)

### **Coordinate Systems**

### $r < -d + geom_bar()$



- **r + coord fixed(**ratio = 1/2) ratio, xlim, ylim - Cartesian coordinates with fixed aspect ratio between x and y units.
  - ggplot(mpg, aes(y = fl)) + geom\_bar() Flip cartesian coordinates by switching x and y aesthetic mappings.
  - r + coord\_polar(theta = "x", direction=1) theta, start, direction - Polar coordinates.
- **r + coord\_trans(**y = "sqrt") x, y, xlim, ylim Transformed cartesian coordinates. Set xtrans and ytrans to the name of a window function.

### $\pi$ + coord\_quickmap()

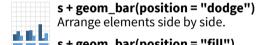
 $\pi$  + coord\_map(projection = "ortho", orientation  $\frac{1}{9} = c(41, -74, 0)$ ) - projection, xlim, ylim

Map projections from the mapproj package (mercator (default), azegualarea, lagrange, etc.).

### **Position Adjustments**

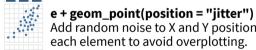
Position adjustments determine how to arrange geoms that would otherwise occupy the same space.

s <- ggplot(mpg, aes(fl, fill = drv))</pre>



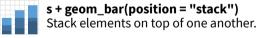
•A

s + geom\_bar(position = "fill") Stack elements on top of one another, normalize height.



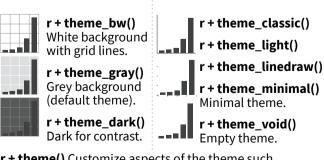
Add random noise to X and Y position of each element to avoid overplotting. e + geom\_label(position = "nudge")

●B Nudge labels away from points.



Each position adjustment can be recast as a function with manual width and height arguments: s + geom\_bar(position = position\_dodge(width = 1))

### Themes



r + theme() Customize aspects of the theme such as axis, legend, panel, and facet properties. r + ggtitle("Title") + theme(plot.title.postion = "plot" r + theme(panel.background = element\_rect(fill = "blue"))

RStudio® is a trademark of RStudio, PBC • CC BY SA RStudio • info@rstudio.com • 844-448-1212 • rstudio.com • Learn more at ggplot2.tidyverse.org • ggplot2 3.3.5 • Updated: 2021-08

### Faceting

Facets divide a plot into subplots based on the values of one or more discrete variables.

t <- ggplot(mpg, aes(cty, hwy)) + geom\_point()

ggplot

fl: r

<b>t + facet_grid(cols = vars(fl))</b> Facet into columns based on fl.
<b>t + facet_grid(rows = vars(year))</b> Facet into rows based on year.
<pre>t + facet_grid(rows = vars(year), cols = vars(fl)) Facet into both rows and columns.</pre>
 <b>t + facet_wrap(vars(fl))</b> Wrap facets into a rectangular layout.

Set scales to let axis limits vary across facets.

- t + facet\_grid(rows = vars(drv), cols = vars(fl), scales = "free") x and y axis limits adjust to individual facets: "free x" - x axis limits adjust
  - "free\_y" y axis limits adjust

### Set labeller to adjust facet label:

t + facet\_grid(cols = vars(fl), labeller = label\_both)

fl: d fl: e fl: p fl: c

t + facet\_grid(rows = vars(fl), labeller = label\_bquote(alpha ^ .(fl)))

 $\alpha^c$   $\alpha^d$   $\alpha^e$   $\alpha^p$   $\alpha^r$ 

### Labels and Legends

Use **labs()** to label the elements of your plot. **t** + **labs**(**x** = "New x axis label", **y** = "New y axis label", **title** ="Add a title above the plot", subtitle = "Add a subtitle below title", **caption** = "Add a caption below plot",

**alt** = "Add alt text to the plot",

<AES> = "New <AES> legend title")

**t + annotate(**geom = "text", x = 8, y = 9, label = "A") Places a geom with manually selected aesthetics.

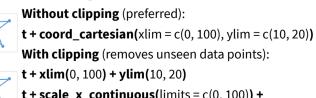
**p** + guides(x = guide\_axis(n.dodge = 2)) Avoid crowded or overlapping labels with guide\_axis(n.dodge or angle).

**n + guides(**fill = "none") Set legend type for each aesthetic: colorbar, legend, or none (no legend).

**n + theme(**legend.position = "bottom") Place legend at "bottom", "top", "left", or "right".

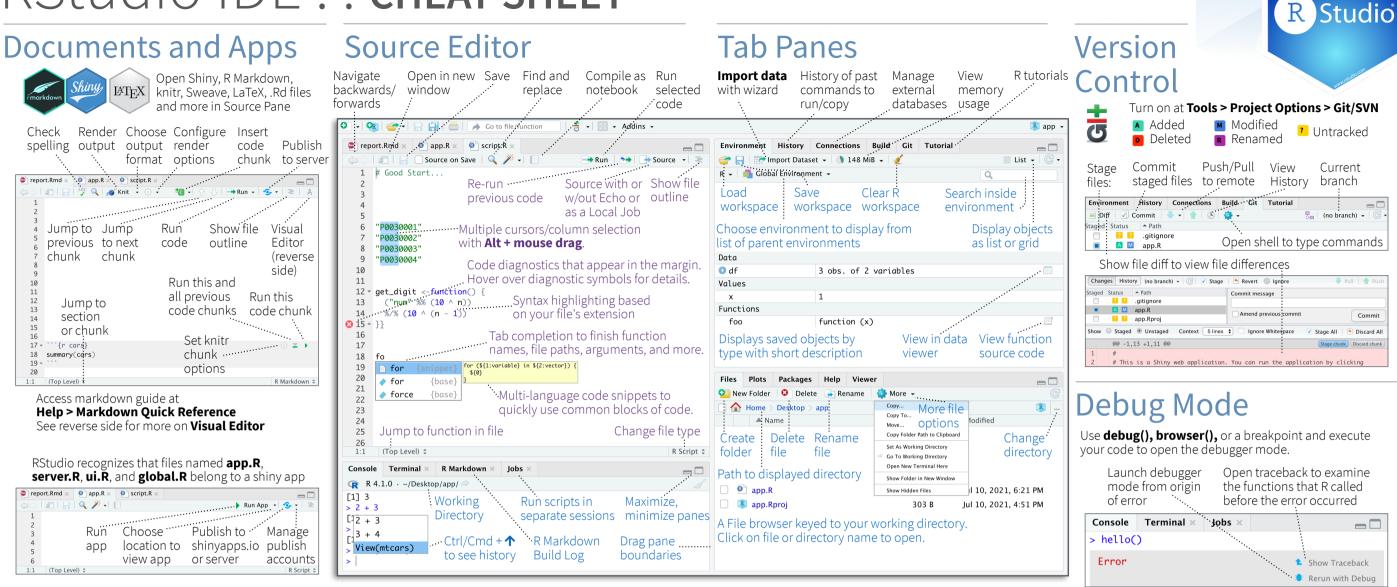
**n + scale\_fill\_discrete(**name = "Title", labels = c("A", "B", "C", "D", "E")**)** Set legend title and labels with a scale function.

### Zooming



t + scale\_x\_continuous(limits = c(0, 100)) + scale\_y\_continuous(limits = c(0, 100))

# RStudio IDE : : CHEAT SHEET



### **Package Development**

Create a new package with File > New Project > New Directory > R Package Enable roxygen documentation with Tools > Project Options > Build Tools

### Roxygen guide at **Help > Roxygen Quick Reference**

See package information in the **Build Tab** 

Install package and restart R Environment History Connecti	Run devtools:: and reload cha	
🔊 Install and Restart 🛛 🔽 Check	🍄 More 👻	đ
Run R CMD check	Clean All ^☆L Clean and Rebuild ······· Test Package ☆ #T.	Clear output and rebuild
Customize package build	Check Package 企業E Build Source Package Build Binary Package	<sup></sup> Run package
options	•• 💸 Configure Build Tools	tests

### RStudio opens plots in a dedicated **Plots** pane

Files	Plots	Packages	Help	Viewer			
Image:	🔎 Zo	oom 🛛 - 🚬 I	Export -	Θ 🥑	·····	4	🍹 Publish 🕞 🤆
: Nav rece	rigate ent pl	ots w	)pen i /indo\	n Ex N <b>P</b>	kport lot	Delete plot	<sup>…</sup> Delete all plots

### GUI Package manager lists every installed package

	•	0	2		0
Files Plots	Packages	Help Viewer			
可 Install 🔰 🤇	Update		Q,		
Name	De	scription		Version	
Install Package	Upda s Pack			owse ckage :	site
tibble	Sir	nple Data Frames		3.1.2	
🖌 tidyr	Tie	dy Messy Data		1.1.3	• 8
Click to l <b>library(</b> package	). Unclicl	kage with < to detach t <b>ach()</b> .	Package version installec	. fr	elete om brary

### RStudio opens documentation in a dedicated **Help** pane

	0 0 0 0 1 1 0 1 1 0 1 1 0 0 1			ea <b>e</b> p pane
Files Plots Packages	Help Viewer			
🦛 🔿 🏠 🔊		Q	•	🕝 Refresh Help Topic
R: Render R Markdown 🗸 🛛 Fi	nd in Topic			
: Home page of helpful links	·Search help file			arch for p file

### **Viewer** pane displays HTML content, such as Shiny apps, RMarkdown reports, and interactive visualizations

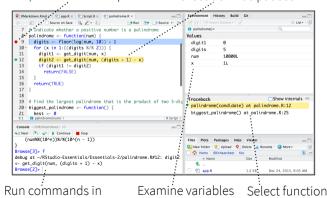
	-					
Files	Plots	Packages	Help	Viewer		
0	1 🔊					💁 Publish 👻 🥝
: Sto app	p Shi	ny			n to shinyapps.io, RSConnect,	: Refresh

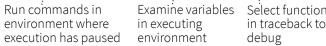
### View(<data>) opens spreadsheet like view of data set

i 🖒   🔊   🖓 Fil	ter						Q	:	
^	mpg 🍦	cyl 🗘	disp 🍦	hp 🗘	drat 🕴 🗘	wt 🌣	qsec 🌼	vs 🌣	am 🍦
Mazda RX4	21.0	6	160.0	110	3.90	2.620	16.46	0	1
Mazda RX4 Wag	21.0	6	160.0	110	3.90	2.875	17.02	0	1
i Filter or val			alue		Sor valı	t by Jes		Beard or va	

Click next to line number to add/remove a breakpoint.

### Highlighted line shows where execution has paused







### Koyboard Shortcute

		ICUTS		
RUN CODE Search comma Interrupt curre Clear console	-	Windows/Linux Ctrl+↑ Esc Ctrl+L	Mac Cmd+ <b>↑</b> Esc Ctrl+L	<b>DOCUMENTS AND A</b> Knit Document (knit Insert chunk (Sweave Run from start to cur
NAVIGATE COD Go to File/Func		Ctrl+.	Ctrl+.	MORE KEYBOARD SI Keyboard Shortcuts
Insert %>% (pi	gnment operator) pe operator)	Ctrl+Shift+M	Tab or Ctrl+Space Option+- Cmd+Shift+M	Show Command Pale View the Keyboard Shortd Reference with <b>Tools &gt; Ke</b> <b>Shortcuts</b> or <b>Alt/Option</b>
(Un)Comment : MAKE PACKAG Load All (devto Test Package (E Document Pack	<b>ES</b> ols) Desktop)	Ctrl+Shift+C Windows/Linux Ctrl+Shift+L Ctrl+Shift+T Ctrl+Shift+D	Cmd+Shift+C Mac Cmd+Shift+L Cmd+Shift+T Cmd+Shift+D	Tabs     Source       Tabs     #P1       Trabs     #F10       Trabs     #F10       Trabs     #E10       Trabs     #E10       Trabs     #E10       Trabs     #E10       Trabs     #E10       Trabs     #E10       Tabs     #E10       Tabs     #E10       Tabs     #E10       Tabs     #E10       Tabs     #E10       Tabs     #E100       Tabs
ſ	Cł sp Preport.Rmd ×	neck Render out belling output for	pose Choose put output mat location	Insert Jump to Ju code previous to chunk chunk ch
Block format	Cł sp Preport.Rmd ×	neck Render out belling output for ABC Q	put output mat location	code previous to chunk chunk ch to to to to to to to to to to to to to

### APPS

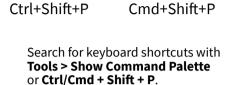
tr) ve & Knitr) irrent line

### **SHORTCUTS**

Help lette

tcut Quick (eyboard + Shift + K

Keyboard Shortcut Quick <b>Tabs</b>	
"^⊕⊊um Switch to Tab	#F9 Back
· ^ uth Next Tab	<b>%F10</b> Forward
^⊕- Previous Tab	\™U Find Usages
^oF11 First Tab	Selection for Find
^#F12 Last Tab	#F Find
	^G Find Next
Panes	r≋G Find Previous
	Replace and Find
^1 Move Focus to Source	^. Go To File/Function



Ctrl+Shift+K

Ctrl+Alt+I

Ctrl+Alt+B

Alt+Shift+K

Cmd+Shift+K

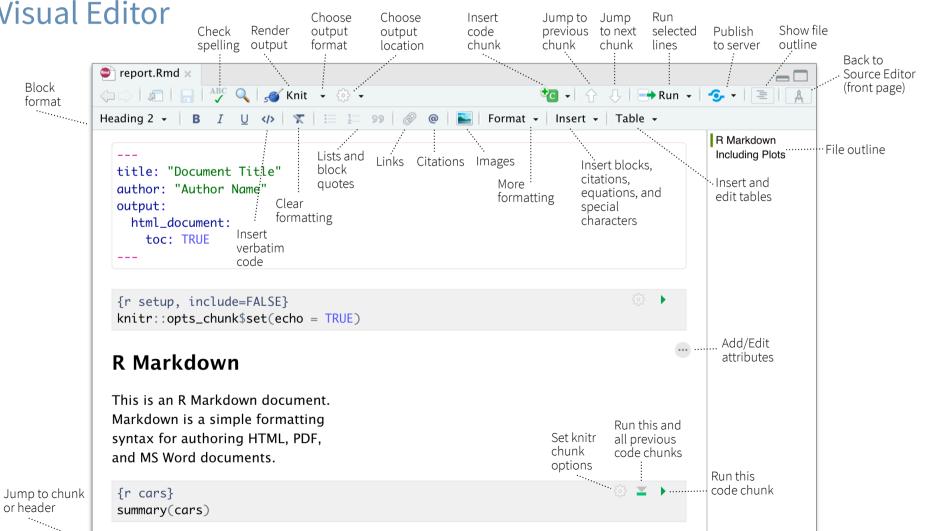
Cmd+Option+I

Cmd+Option+B

Option+Shift+K

History Send Command to Console Create a New R Script Ctrl Shift N

Create a new R Markdown document



### **RStudio** Workbench

### WHY RSTUDIO WORKBENCH?

Extend the open source server with a commercial license, support, and more:

- open and run multiple R sessions at once
- tune your resources to improve performance
- administrative tools for managing user sessions
- collaborate real-time with others in shared projects
- switch easily from one version of R to a different version
- integrate with your authentication, authorization, and audit practices

R Studio

• work in the RStudio IDE, JupyterLab, Jupyter Notebooks, or VS Code

Download a free 45 day evaluation at

### www.rstudio.com/products/workbench/evaluation/

### **Share Projects**

### File > New Project

RStudio saves the call history, workspace, and working Close R Session Start new R Session directory associated with a in project in current project project. It reloads each when you re-open a project. T H J garrett 🕞 Sessions - 🔅 ۲ 🔳 IDEcheatsheet 🕶 🛛 R 3.2.2 🕶 🤏 New Project... ✓ R version 3.2.2 Active shared R version 3.1.3 🕣 Open Project.. collaborators Close Proiect R version 3.0.3 Name of R version 2.15.3 🚨 Share Project.. current IDEcheatsheet RStudio-Essential project Select Essentials Ð **R** Version A Share Project shiny-examples with Collaborators Clear Proiect List Project Options. Ω₩.

### **Run Remote Jobs**

Run R on remote (Kubernetes/Sluri Job Launcher	010.000.0	Run Source Source with Echo	Source ・ 企業S 企業・
Monitor launcher job	Launch a job	A Source as Launcher Source as Local Job	
Console Terminal	Jobs × Launcher ×		
↓ fast.R	Running	Local 0:09	<b>()</b>
sleepy.R	Succeeded 11:22 AM	Local 0:41	۵
sleepy.R	Idle	KubernetesX Waiting	<b>•• &gt;</b>
	Run launcher jobs remotely		

udio

# R Markdown 🗘

R Markdown 🗘



### 6.3 SQL language

# Study Guide: Data Retrieval with SQL

Afshine AMIDI and Shervine AMIDI

#### August 21, 2020

Category	Operator	Command
	Equality / non-equality	= / !=, <>
	Inequalities	>=, >, <, <=
General	Belonging	IN (val_1,, val_n)
General	And / or	AND / OR
	Check for missing value	IS NULL
	Between bounds	BETWEEN val_1 AND val_2
Strings	Pattern matching	LIKE '%val%'

#### General concepts

 $\square$  Structured Query Language – Structured Query Language, abbreviated as SQL, is a language that is largely used in the industry to query data from databases.

#### **Query structure** – Queries are usually structured as follows:

COL	
SQL	
Select fields SELECT	mandatory
col_1,	
col_2, ,	
col_n	
Source of data	mandatory
FROM table t	
	optional
JOIN other_table ot ON (t.key = ot.key)	
Conditions	optional
WHERE some_condition(s)	optional
Aggregating	optional
GROUP BY column_group_list	
Sorting values	optional
ORDER BY column_order_list	
Restricting aggregated values	optional
HAVING some_condition(s)	
Limiting number of rows	optional
LIMIT some_value	

Remark: the SELECT DISTINCT command can be used to ensure not having duplicate rows.

 $\square$  Condition – A condition is of the following format:

 $\mathbf{SQL}$ 

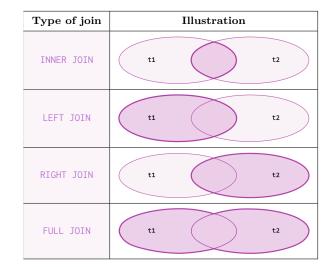
some\_col some\_operator some\_col\_or\_value

where <code>some\_operator</code> can be among the following common operations:

 $\square$  Joins – Two tables table\_1 and table\_2 can be joined in the following way:



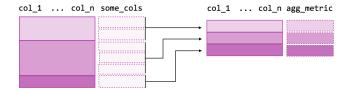
where the different type\_of\_join commands are summarized in the table below:



Remark: joining every row of table 1 with every row of table 2 can be done with the CROSS JOIN command, and is commonly known as the cartesian product.

#### Aggregations

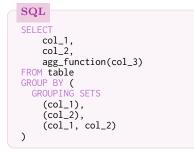
**Grouping data** – Aggregate metrics are computed on grouped data in the following way:



#### The SQL command is as follows:



**Grouping sets** – The GROUPING SETS command is useful when there is a need to compute aggregations across different dimensions at a time. Below is an example of how all aggregations across two dimensions are computed:



□ Aggregation functions – Th	he table below s	summarizes the	main	aggregate i	functions th	hat
can be used in an aggregation que	ery:					

Category	Operation	Command
	Mean	AVG(col)
	Percentile	<pre>PERCENTILE_APPROX(col, p)</pre>
Values	Sum / $\#$ of instances	<pre>SUM(col) / COUNT(col)</pre>
	Max / min	MAX(col) / MIN(col)
	Variance / standard deviation	VAR(col) / STDEV(col)
Arrays	Concatenate into array	<pre>collect_list(col)</pre>

Remark: the median can be computed using the PERCENTILE\_APPROX function with p equal to 0.5.

 $\square$  Filtering – The table below highlights the differences between the <code>WHERE</code> and <code>HAVING</code> commands:

WHERE	HAVING	
<ul><li>Filter condition applies to individual rows</li><li>Statement placed right after FROM</li></ul>	<ul><li>Filter condition applies to aggregates</li><li>Statement placed right after GROUP BY</li></ul>	

Remark: if WHERE and HAVING are both in the same query, WHERE will be executed first.

#### Window functions

 $\square$  **Definition** – A window function computes a metric over groups and has the following structure:

col_1 col_n	some_cols	col_1 col_n	win_metric

The SQL command is as follows:

$\mathbf{SQL}$				

some\_window\_function() OVER(PARTITION BY some\_col ORDER BY another\_col)

Remark: window functions are only allowed in the SELECT clause.

**Row numbering** – The table below summarizes the main commands that rank each row across specified groups, ordered by a specific column:

Command	Description	Example
ROW_NUMBER()	Ties are given different ranks	1, 2, 3, 4
RANK ()	Ties are given same rank and skip numbers	1, 2, 2, 4
DENSE_RANK ()	Ties are given same rank and don't skip numbers	1, 2, 2, 3

 $\square$  Values – The following window functions allow to keep track of specific types of values with respect to the partition:

Command	Description
FIRST_VALUE(col)	Takes the first value of the column
LAST_VALUE(col)	Takes the last value of the column
LAG(col, n)	Takes the $n^{\text{th}}$ previous value of the column
LEAD(col, n)	Takes the $n^{\text{th}}$ following value of the column
NTH_VALUE(col, n)	Takes the $n^{\text{th}}$ value of the column

#### Advanced functions

 $\square$  SQL tips – In order to keep the query in a clear and concise format, the following tricks are often done:

Operation	Command	Description
Renaming columns	SELECT operation_on_column AS col_name	New column names shown in query results
Abbreviating tables	FROM table_1 t1	Abbreviation used within query for simplicity in notations
Simplifying group by	GROUP BY col_number_list	Specify column position in SELECT clause instead of whole column names
Limiting results	LIMIT n	Display only n rows

 $\square$  Sorting values – The query results can be sorted along a given set of columns using the following command:

#### SQL ... [query] ... ORDER BY col\_list

Remark: by default, the command sorts in ascending order. If we want to sort it in descending order, the DESC command needs to be used after the column.

 $\Box$  Column types – In order to ensure that a column or value is of one specific data type, the following command is used:

#### $\mathbf{SQL}$

CAST(some\_col\_or\_value AS data\_type)

where data\_type is one of the following:

Data type	Description	Example	
INT	Integer	2	
DOUBLE	Numerical value	2.0	
STRING	String	'teddy bear'	
VARCHAR			
DATE	Date	2020-01-01'	
TIMESTAMP	Timestamp	'2020-01-01 00:00:00.000'	

Remark: if the column contains data of different types, the TRY\_CAST() command will convert unknown types to NULL instead of throwing an error.

 $\square$  Column manipulation – The main functions used to manipulate columns are described in the table below:

Category	Operation	Command
	Take first non-NULL value	COALESCE(col_1, col_2,, col_n)
General	Create a new column combining existing ones	CONCAT(col_1,, col_n)
Value	Round value to n decimals	ROUND(col, n)
	Converts string column to lower / upper case	LOWER(col) / UPPER(col)
	Replace occurrences of old in col to new	REPLACE(col, old, new)
String	Take the substring of col, with a given start and length	SUBSTR(col, start, length)
	Remove spaces from the left / right / both sides	LTRIM(col) / RTRIM(col) / TRIM(col)
	Length of the string	LENGTH(col)
Date	Truncate at a given granularity (year, month, week)	DATE_TRUNC(time_dimension, col_date)
	Transform date	DATE_ADD(col_date, number_of_days)

 $\Box$  Conditional column – A column can take different values with respect to a particular set of conditions with the CASE when command as follows:

#### SQL CASE WHEN some\_condition THEN some\_value

- WHEN some\_other\_condition THEN some\_other\_value
- ELSE some\_other\_value\_n END

 $\square$  Combining results – The table below summarizes the main ways to combine results in queries:

Category	Command	Remarks
	UNION	Guarantees distinct rows
Union	UNION ALL	Potential newly-formed duplicates are kept
Intersection	INTERSECT	Keeps observations that are in all selected queries

 $\Box$  Common table expression – A common way of handling complex queries is to have temporary result sets coming from intermediary queries, which are called common table expressions (abbreviated CTE), that increase the readability of the overall query. It is done thanks to the WITH ... AS ... command as follows:

 $\mathbf{SQL}$ 

WITH cte\_1 AS ( SELECT ... ),

cte_n AS ( SELECT )
SELECT FROM

#### Table manipulation

**Table creation** – The creation of a table is done as follows:

$\mathbf{SQL}$	
	E [table_type] TABLE [creation_type] table_name( _1 data_type_1,
col. ) [optic	_n data_type_n
Lober	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,

where [table\_type], [creation\_type] and [options] are one of the following:

Category	Command	Description
	Blank	Default table
Table type	EXTERNAL TABLE	External table
Creation type	Blank	Creates table and overwrites current one if it exists
	IF NOT EXISTS	Only creates table if it does not exist
Options	<pre>location 'path_to_hdfs_folder'</pre>	Populate table with data from hdfs folder
	stored as data_format	Stores the table in a specific data format, e.g. parquet, orc or avro

 $\square$  Data insertion – New data can either append or overwrite already existing data in a given table as follows:

SQL	
WITH	optional
<pre>INSERT [insert_type] table_name</pre>	mandatory
SELECT;	mandatory

where  $\verb[insert_type]$  is among the following:

Command	Description
OVERWRITE	Overwrites existing data
INTO	Appends to existing data

**Dropping table** – Tables are dropped in the following way:

SQI	
-----	--

DROP TABLE table\_name;

 $\square$  View – Instead of using a complicated query, the latter can be saved as a view which can then be used to get the data. A view is created with the following command:

#### $\mathbf{SQL}$

CREATE VIEW view\_name AS complicated\_query;

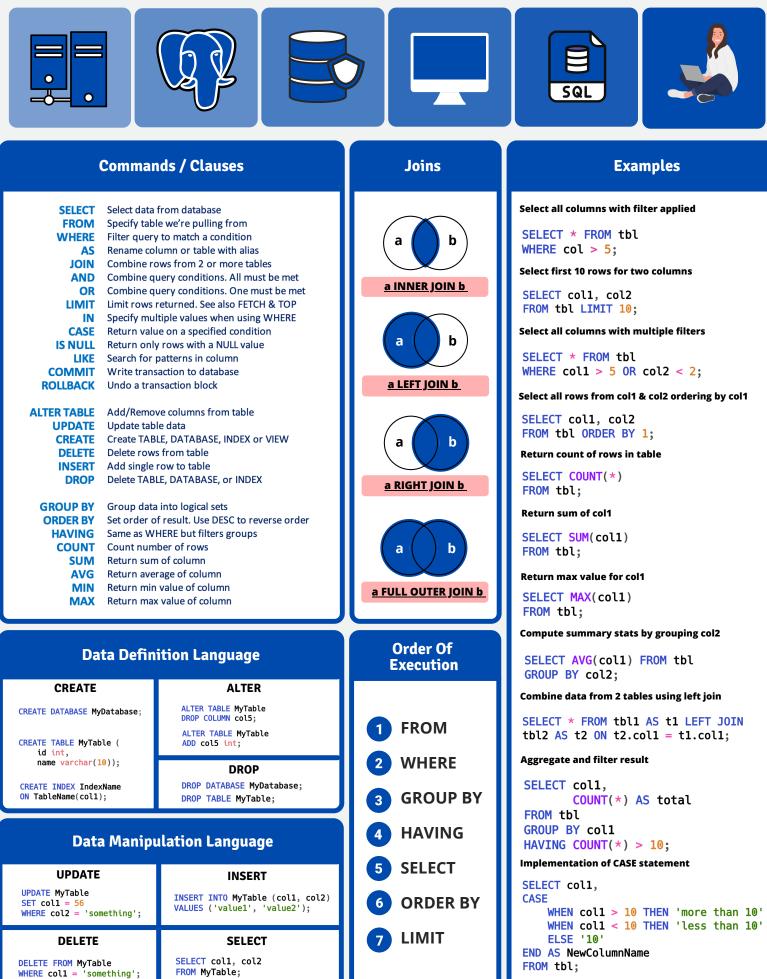
Remark: a view does not create any physical table and is instead seen as a shortcut.

#### Sources

- W3Schools.com
- DataQuest.io







# SQL cheat sheet

For more awesome cheat sheets **REBELLA** Visit rebellabs.org!

## **Basic Queries**

-- filter your columns SELECT col1, col2, col3, ... FROM table1

-- filter the rows WHERE col4 = 1 AND col5 = 2

-- aggregate the data GROUP by ...

-- limit aggregated data
 HAVING count(\*) > 1
 -- order of the results

ORDER BY col2

#### Useful keywords for SELECTS:

**DISTINCT** - return unique results **BETWEEN** a **AND** b - limit the range, the values can be numbers, text, or dates **LIKE** - pattern search within the column text **IN** (a, b, c) - check if the value is contained among given.

## **Data Modification**

-- update specific data with the **WHERE** clause **UPDATE** table1 **SET** col1 = 1 **WHERE** col2 = 2

-- insert values manually
 INSERT INTO table1 (ID, FIRST\_NAME, LAST\_NAME)
 VALUES (1, 'Rebel', 'Labs');

-- or by using the results of a query

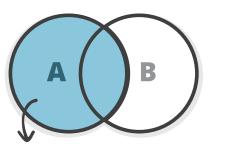
INSERT INTO table1 (ID, FIRST\_NAME, LAST\_NAME) SELECT id, last\_name, first\_name FROM table2

## Views

A **VIEW** is a virtual table, which is a result of a query. They can be used to create virtual tables of complex queries.

CREATE VIEW view1 AS SELECT col1, col2 FROM table1 WHERE ...

## The Joy of JOINs



**LEFT OUTER JOIN -** all rows from table A, even if they do not exist in table B

## **Updates on JOINed Queries**

You can use **JOIN**s in your **UPDATE**s **UPDATE** t1 **SET** a = 1 **FROM** table1 t1 **JOIN** table2 t2 **ON** t1.id = t2.t1\_id **WHERE** t1.col1 = 0 **AND** t2.col2 **IS NULL**;

NB! Use database specific syntax, it might be faster!

# Semi JOINs

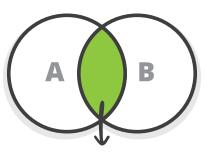
You can use subqueries instead of JOINs:

SELECT col1, col2 FROM table1 WHERE id IN (SELECT t1\_id FROM table2 WHERE date > CURRENT\_TIMESTAMP)

## Indexes

If you query by a column, index it! **CREATE INDEX** index1 **ON** table1 (col1)

<u>Don't forget:</u> Avoid overlapping indexes Avoid indexing on too many columns Indexes can speed up **DELETE** and **UPDATE** operations



*INNER JOIN -* fetch the results that exist in both tables

AB

**RIGHT OUTER JOIN -** all rows from table B, even if they do not exist in table A

## **Useful Utility Functions**

-- convert strings to dates:

- TO\_DATE (Oracle, PostgreSQL), STR\_TO\_DATE (MySQL)
- -- return the first non-NULL argument: **COALESCE** (col1, col2, "default value") -- return current time:
- CURRENT\_TIMESTAMP
- -- compute set operations on two result sets **SELECT** col1, col2 **FROM** table1 **UNION / EXCEPT / INTERSECT SELECT** col3, col4 FROM table2;

 Union - returns data from both queries
 Except - rows from the first query that are not present in the second query
 Intersect - rows that are returned from both queries

# Reporting

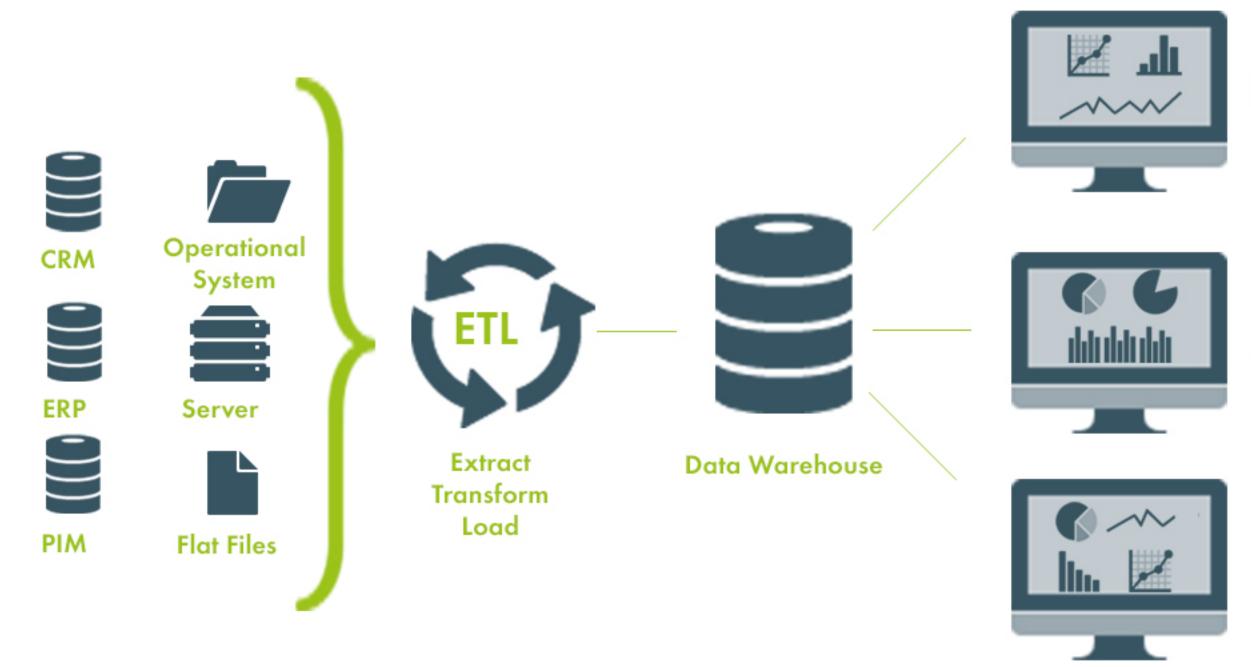
Use aggregation functions

COUNT - return the number of rows SUM - cumulate the values AVG - return the average for the group MIN / MAX - smallest / largest value





# 7. Business Intelligence



# **Business Inteligence**

# **Data Mining**

**Business Analytics** 

# **Big Data**

**Predictive Analysis** 

Reporting

# Planning

**Data Visualisation** 



# 7.1 Tableau

# **Tableau for Business Intelligence Tableau Basics Cheat Sheet**

Learn Tableau online at www.DataCamp.com

# What is Tableau?

Tableau is a business intelligence tool that allows you to effectively report insights through easy-to-use customizable visualizations and dashboards

			A 🛪 🛛 🍕
Home > Salesforce > Open_Pipeline > O	pen Pipeline 🔆 🖯 1		^
- Undo → Redo  ← Revert 🔒	Refresh C	🔄 View: Original 🛕 Alert 🖾 Subsoribe 🎤 Edit 👳	Share 🖓 Download 🖓 Comments 💢 Full Scr
Open Pipeline	Ellie Rogers'	steam	
Total Expected Amoun	nt Avg. Deal Size	Number of Open Opportunities	Avg. Age of Opportunity
\$3,093,00	0 \$19,453	159	75 days
Use Options Below to Filter	Overall Opportunities Closing   by Day	Number of Open Opportunities and Average Days in Stage	* Click a Stage to Filter
licse Date 10/1/2017 10/31/2017 	10 5 0 0115 0125 0125	Outlification Prospecting Value Proposition Needs Analysis On Tal (16) Exclorem Makers	
	Breakdown by Stage   by Day * Click a Day to Filter Qualification	Perception Analysis Proposal/Price Quote	
(All) *	Prespecting	Negotiation/Review 0 10 20 Opportunities F	30 0 20 40 60 Days in Stage
(All) •	Value Proposition	Top 10 Opportunities   by Expected Amount	
Iling Country All)	Needs Analysis	Id Sapien Cras Corporation 5 Sogittis Felis Donec PC 8	\$75,000 \$72,000
Illing State/Province (All)	Id. Decision Makers	Mi Ltd 13 Inceptos Hymenaeos Mauris Associates 4	\$28,000 \$15,000
dustry	Perception Analysis Proposal/Price	Nec Tempus Corp. 13 Ligula Aenean PC 2	\$10,000 \$4,000
(AII) *	Quote	Et Netus Corporation 10	\$3,500

# Why use Tableau?

**<>** Easy to use—no coding involved

>

Integrates seamlessly with any data source

4 Fast and can handle large datasets

# **Tableau Versions** >

There are two main versions of Tableau

# TABLEAU PUBLIC

A free version of Tableau that lets you connect to limited data sources, create visualizations and dashboards, and publish dashboards online

# TABLEAU DESKTOP

A paid version of tableau which lets you connect to all types of data sources, allows you to save work locally, and unlimited data sizes

# Getting started with Tableau >

When working with Tableau, you will work with Workbooks. Workbooks contain sheets, dashboards, and stories. Similar to Microsoft Excel, a Workbook can contain multiple sheets. A sheet can be any of the following and can be accessed on the bottom left of a workbook



# WORKSHEET

A worksheet is a single view in a workbook. You can add shelves, cards, legends, visualizations, and more in a worksheet



DASHBOARD A collection of multiple worksheets used to display multiple views simultaneously



# STORY

A story is a collection of multiple dashboards and/ or sheets that describe a data story

# >

The Anatomy of a Worksheet

When opening a worksheet, you will work with a variety of tools and interfaces

# The Sidebar

- In the sidebar, you'll find useful panes for working with data
- 1. Data: The data pane on the left-hand side contains all of the fields in the currently selected data source
- 2. Analytics: The analytics pane on the left-hand side lets you add useful insights like trend lines, error bars, and other useful summaries to visualizations

# **Tableau Data Definitions**

When working with data in Tableau, there are multiple definitions to be mindful of

- 1. Fields: Fields are all of the different columns or values in a data source or that are calculated in the workbook. They show up in the data pane and can either be dimension or measure fields
- 2. **Dimensions:** A dimension is a type of field that contains qualitative values (e.g. locations, names, and departments). Dimensions dictate the amount of granularity in visualizations and help reveal nuanced details in the data

3. Measures: A measure is a type of field that contains quantitative values (e.g. revenue, costs, and market sizes). When dragged into a view, this data is aggregated, which is determined by the dimensions in the view

4. Data types: Every field has a data type which is determined by the type of information it contains. The available data types in Tableau include text, date values, date & time values, numerical values, boolean values, geographical values, and cluster groups

# The Canvas

The canvas is where you'll create data visualizations

1. Tableau Canvas: The canvas takes up most of the screen on Tableau and is where you can add visualizations 2. Rows and columns: Rows and columns dictate how the data is displayed in the canvas. When dimensions are placed, they create headers for the rows or columns while measures add quantitative values 3. Marks card: The marks card allows users to add visual details such as color, size, labels, etc. to rows and columns. This is done by dragging fields from the data pane into the marks card

# Visualizing Your First Dataset

Upload a dataset to Tableau

1. Launch Tableau

>

- 2. In the Connect section, under To a File, press on the file format of your choice
- 3. For selecting an Excel file, select .xlsx or .xlsx

## Creating your first visualization

- 1. Once your file is uploaded, open a Worksheet and click on the Data pane on the left-hand side
- 2. Drag and drop at least one field into the *Columns* section, and one field into the *Rows* section at the top of the canvas
- 3. To add more detail, drag and drop a dimension into the *Marks* card (e.g. drag a dimension over the color square in the marks card to color visualization components by that dimension)
- 4. To a summary insight like a trendline, click on the *Analytics* pane and drag the trend line into your visualization
- 5. You can change the type of visualization for your data by clicking on the Show Me button on the top right

# Data Visualizations in Tableau

Tableau provides a wide range of data visualizations to use. Here is a list of the most useful visualizations you have in Tableau

- **Bar Charts:** Horizontal bars used for comparing specific values across categories (*e.g. sales by region*)
- **Stacked Bar Chart:** Used to show categorical data within a bar chart (*e.g., sales by region and department*)
- Side-by-Side Bar Chart: Used to compare values across categories in a bar chart format (e.g., sales by region comparing product types)

Eine Charts: Used for looking at a numeric value over time (e.g., revenue over time)

- $o_{o+}^{+o}$  Scatter Plot: Used to identify patterns between two continuous variables (e.g., profit vs. sales volume)
- **Histogram:** Used to show a distribution of data (e.g., Distribution of monthly revenue)
- Box-and-Whisker Plot: Used to compare distributions between categorical variables (e.g., distribution of revenue by region)
- Heat Map: Used to visualize data in rows and columns as colors (e.g., revenue by marketing channel)
- Highlight Table: Used to show data values with conditional color formatting (e.g., site-traffic by marketing) channel and year)
- 🛚 🟋 Symbol Map: Used to show geographical data (e.g., Market size opportunity by state)
- 🛚 💱 🗺 Map: Used to show geographical data with color formatting (e.g., Covid cases by state)
- **Treemap:** Used to show hierarchical data (e.g., Show how much revenue subdivisions generate relative to the whole department within an organization)
- **Dual Combination:** Used to show two visualizations within the same visualization (e.g., profit for a store each month as a bar chart with inventory over time as a line chart)

# **Customizing Visualizations with Tableau**

Tableau provides a deep ability to filter, format, aggregate, customize, and highlight specific parts of your data visualizations

## Filtering data with highlights

1. Once you've created a visual, click and drag your mouse over the specific portion you want to highlight

2. Once you let go, you will have the option to  $\checkmark$  Keep Only or  $\times$  Exclude the data

3. Open the *Data* pane on the side bar. Then, you can drag-and-drop a field into the fitlers card just to the

## Filtering data with filters

left of the pane.

- 1. Open the *Data* pane on the left-hand-side
- 2. Drag-and-drop a field you want to filter on and add it to the Filters card
- 3. Fill out in the modal how you would like your visuals to be filtered on the data

# Aggregating data

When data is dragged into the Rows and Columns on a sheet, it is aggregated based on the dimensions in the sheet. This is typically a summed value. The default aggregation can be changed using the steps below:

# Changing colors

# Changing fonts use the following steps



>

- 1. Launch Tableau
- 3. Select your file 4. Click the **M** New Sheet at the bottom to create a new sheet
- 5. Create a visualization in the sheet by following the steps in the previous sections of this cheat sheet
- 6. Repeat steps 4 and 5 untill you have created all the visualizations you want to include in your dashboard 7. Click the **H** New Dashboard at the bottom of the screen

+;+ • •	bleau			
Home > On	acle Eloqua ,	Account_Er	ngagement >	Account E
← Undo ·	→ Redo	Revert	🕃 Refresh	Gi Par
Enga	igeme	nt A	nalysi	s fo
Day of Activit November 15				
	Control	Activity Ty	pe	
PageVie		ebVisit	FormSub	-t
			Ac	counts
Con	npany		Total Ad	tivity
		Sing	er Lumber	
		Si	mply Save	
	Environ	Architectu	ral Design	
			Food Barn	
	Mu	lti Tech De	velopment	
			0	50
				Timelir
20		1		



- 1. Click the **U** *New Story* at the bottom of the screen
- 2. Change the size of the story to the desired size in the bottom left-hand corner of the screen under Size 3. Edit the title of the story by renaming the story. To do this, right-click on the story sheet at the bottom and press Rename

- 4. A story is made of story points, which lets you cycle through different visualizations and dashboards 5. To begin adding to the story, add a story point from the left-hand side. You can add a blank story point 6. To add a summary text to the story, click *Add a caption* and summarize the story point
- 7. Add as many story points as you would like to finalize your data story

# A Case Study on Tables



Most Profitable Tables





1. Right-click on a measure field in the *Data* pane

2. Go down to Default properties, Aggregation, and select the aggregation you would like to use

Color is a critical component of visualizations. It draws attention to details. Attention is the most important component of strong storytelling. Colors in a graph can be set using the marks card.

1. Create a visualization by dragging fields into the *Rows* and *Columns* section at the top of the screen 2. Drag dimensions into the *Marks* field, specifically into the *Color* square 3. To change from the default colors, go to the upper-right corner of the color legend and select *Edit Colors*. This

will bring up a dialog that allows you to select a different palette

Fonts can help with the aesthetic of the visualization or help with consistent branding. To change the workbook's font,

1. In the Format menu on the top ribbon, press on Select Workbook. This will replace the Data pane and allow you to make formatting decisions for the Workbook

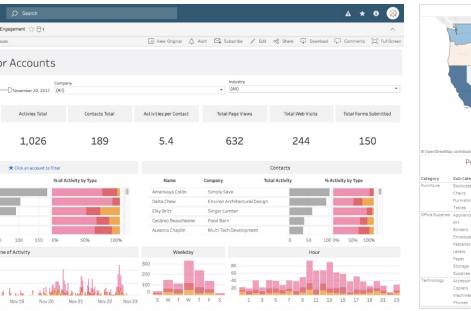
2. From here, select the font, font size, and color

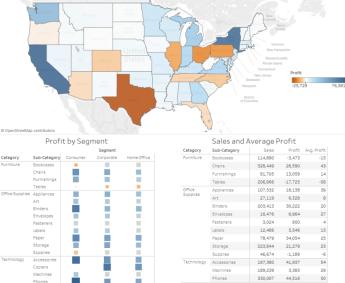
# Creating dashboards with Tableau

Dashboards are an excellent way to consolidate visualizations and present data to a variety of stakeholders. Here is a step by step process you can follow to create a dashboard

2. In the *Connect* section under *To A File*, press on your desired file type

- 8. On the left-hand side, you will see all your created sheets. Drag sheets into the dashboard
- 9. Adjust the layout of your sheets by dragging and dropping your visualizations

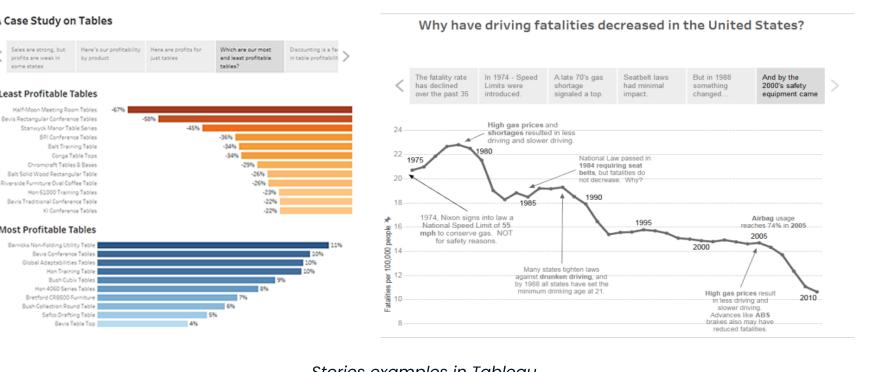




Dashboard examples in Tableau

# **Creating stories with Tableau**

A story is a collection of multiple dashboards and/or sheets that describe a data story



Stories examples in Tableau

# Learn Data Skills Online at <a href="https://www.DataCamp.com">www.DataCamp.com</a>

# Tableau - Desktop

# CHEAT SHEET

#### **Data Sources**

File Systems	CSV, Excel, etc.
Relational Systems	Oracle, Sql Server, DB2, etc.
Cloud Systems	Windows Azure, Google BigQuery, etc.
Other Sources	ODBC

**Data Blending** 

• Blending the Data

Preparing Data for Blending

Adding Secondary Data Source

#### **Data Extract**

- Extraction of data is done by following ٠ Menu  $\rightarrow$  Data  $\rightarrow$  Extract Data.
- Applying Extract Filters to create subset of data .
- To add more data for an already created extract  $Data \rightarrow Extract \rightarrow Append Data from File$
- Extract History ٠ Menu - Data → Extract History

	Data	Joining
--	------	---------

- Creating a Join
- Editing a Join Type
- Editing Join Fields

#### **Operators**

- General Operators Relational Operators •
- Arithmetic Operators Logical Operators •

#### **LOD Expressions**

• Fixed LOD , Include LOD and Exclude LOD

#### Sorting

- Computed Sorting: Directly applied on an axis using the sort dialog button.
- Manual Sorting: Rearrange the order of dimension fields by dragging them next to each other.

## **Data Sources**

Types	Work
Filter Dimensions	Applied on the dimension fields.
Filter Measures	Applied on the measure fields.
Filter Dates	Applied on the date fields.
Single Value (List)	Select one value at a time in a list.
Single Value (Dropdown)	Select a single value in a drop-down list.

#### **Tableau Charts**

Туре	Description				
Text Table (Crosstab)	To see your data in rows and columns.				
Heat Map	Just like Crosstab, but it uses size and color as visual cues to describe the data.				
Highlight Table	Just like Excel table, but the cells here are colored.				
Symbol Map	Visualize and highlight geographical data.				
Filled Map	Color filled geographical data visualization.				
Pie Chart	Represents data as slices of a circle with different sizes and colors.				
Horizontal Bar Chart	Represents data in horizontal bars, visually digestible.				
Stacked Bar Chart	Visualize data of a category having sub-categories.				
Side-by-Side Bar Chart	Side by side comparison of data, vertical representation.				
Treemap	Similar to a heat map, but the boxes are grouped by items that are close in hierarchy.				
Circle View	Shows the different values that are within the categories.				
Side-by-Side Circle View fields)	Combination of Circle view and Side-by-Side Bar Chart				
Line Chart (Continuous)	Several number of lines in the view to show continuous flow of data, must have a date.				
Line Chart (Discrete)	This allows slicing and dicing of the graph, graph not continuous.				
Dual Line Chart	Comparing two measures over a period.				
Scatter Plot	Scatter plot shows many points scattered in the Cartesian plane				
Histogram	A histogram represents the frequencies of values of a variable bucketed into ranges				
Gantt Chart	It illustrates a project schedule.				
Bullet Graph	Two bars drawn upon one another to indicate their individual values at the same position in the graph				
Waterfall Chart	It shows where a value starts, ends and how it gets there incrementally				

Work Types Multiple Values (List) Select one or more values in a list. Multiple Values (Dropdown) Select one or more values in a drop-down list. Multiple Values (Custom List) Search and select one or more values. Single Value (Slider) Drag a horizontal slider to select a single value. Wildcard Match Select values containing the specified characters.



🕂 + a b | e a u

**FURTHERMORE:** Tableau Training and Certification - Tableau 10 Desktop Course

# Tableau-Desktop Shortcuts & Terminologies CHEAT SHEET

#### What is Tableau?

A powerful Data visualization and Business intelligence tool with a strong and intuitive interface. No coding knowledge or experience needed to work with Tableau.

#### **Data Short Cuts:**

Data	Operation			
ALT+D+A	Automatic updates			
ALT+D+D	Connect to data			
CTRL+D	Connect to data source			
ALT+D+C+D	Duplicate data connection			
ALT+D+A	Extract			
ALT+D+C+P	Properties of data connection			
ALT+D+R	Refresh data			
F5	Refreshes the data source			
F9	Run query			
ALT+D+U	Toggle use extract			
F10	Toggles automatic updates on and off			
ALT+D+P+D	Update dashboard			
ALT+D+P+Q	Update quick filters			
ALT+D+P+W	Update worksheet			

#### **File Short Cuts:**

File	Operation
ALT+F4	Closes the current workbook
ALT+F+E+I	Export to image
ALT+F+E+P	Export to packaged workbook
CTRL+N	New workbook
CTRL+O	Open file
CTRL+P	Print
CTRL+S	Save file
	-

#### File Short Cuts:

Analysis	Operation
ALT+A+A	Aggregate measures
ALT+A+C	Create calculated field
ALT+A+B	Describe trend model
ALT+A+U	Edit calculated field
ALT+A+L	Edit trend lines
ALT+A+F	Filter
CTRL+1	Show me!
ALT+A+S	Sort
ALT+A+M+F	Stack marks off
ALT+A+M+O	Stack marks on

#### Get started with Tableau

Tableau is BI software

 Allows users to connect to data, visualize data and create interactive and sharable dashboards.

#### **Design Flow**

#### **Tableau Terminologies**

Alias: Refers to a field or to a dimension member

Bin: User-defined grouping of measures in the data source

Bookmark: .tbm file in the Bookmarks folder contains a single worksheet of the Tableau repository

Calculated Field: New field created by using a formula to modify the existing fields in data source

Crosstab: Text table view to display the numbers associated with dimension members

Dashboard: Use dashboards to compare and monitor a variety of data simultaneously

Data Pane: Displays the fields (divided into dimensions and measures) of the data sources to which Tableau is connected

Data Source Page: A page to set up your data source consists of – left pane, join area, preview area, and metadata area

Dimension: Categorical data field holds discrete data such as hierarchies and members that cannot be aggregated Extract: A saved subset of a data source that can be used to improve performance and analyze offline.

Filters Shelf: Used to exclude data from a view by filtering it using measures and dimensions

Format Pane: Contains formatting settings that control the entire worksheet, as well as individual fields in the view Level of Detail (LOD) Expression: A syntax that supports aggregation at dimensionalities other than the view level. Marks: A part of the view that visually represents one or more rows in a data source. A mark can be, for example, a bar, line, or square. You can control the type, color, and size of marks

Marks Card: A card to the left of the view, where you can drag fields to control mark properties such as type, color, size, shape, label, tooltip, and detail

**ntelliPaat** 

Pages Shelf: Used to split a view into a sequence of pages based on the members and values in a discrete or continuous field

Rows Shelf: Used to create the rows of a data table, also accepts any number of dimensions and measures

Worksheet: A sheet to build views of the data

Workbook: Contains one or more worksheets

Connect to Data Source → Build Data Views → Enhance Data Views → Worksheets → Create and Organize Dashboards → Story Telling

FURTHERMORE: Tableau Training and Certification - Tableau 10 Desktop Course



# 7.2 Power Bl

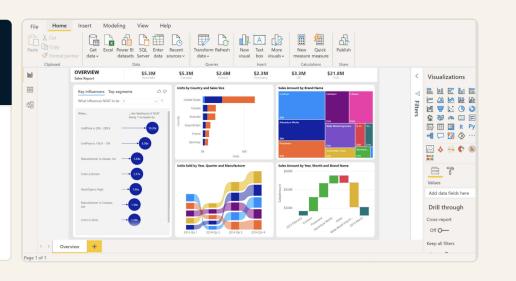
# datacamp Microsoft **Power BI for Business Intelligence**

# Power BI Cheat Sheet

Learn Power BI online at www.DataCamp.com

# What is Power BI?

Power BI is a business intelligence tool that allows you to effectively report insights through easy-to-use customizable visualizations and dashboards.



# Why use Power BI?

**<>** Easy to use—no coding involved

>

- Integrates seamlessly with any data source
- 4 Fast and can handle large datasets

# Power BI Components >

There are three components to Power BI-each of them serving different purposes

# **POWER BI DESKTOP**

Free desktop application that provides data analysis and creation tools.

# **POWER BI SERVICE**

Cloud-based version of Power Bl with report editing and publishing features.

# **POWER BI MOBILE**

A mobile app of Power Bl, which allows you to author, view, and share reports on the go.

# Getting started with Power Bl >

There are three main views in Power Bl

# REPORT VIEW

>

This view is the default view, where you can visualize data and create reports



# DATA VIEW

This view lets you examine datasets associated with your reports



# MODEL VIEW

This view helps you establish different relationships between datasets

# Visualizing your first dataset

# Upload datasets into Power BI

- Underneath the Home tab, click on Get Data
- Choose any of your datasets and double click
- Click on *Load* if not prior data needs processing
- If you need to transform the data, click *Transform* which will launch Power Query. Keep reading this cheat sheet for how to apply transformations in Power Query.
- Inspect your data by clicking on the 🎹 Data View

# Create relationships in Power BI

Sales Performance	
SalesPersonID	
Employee Database	
EmployeeID	

- If you have different datasets you want to connect. First, upload them into Power Bl
- Click on the 唱 Model View from the left-hand pane
- · Connect key columns from different datasets by dragging one to another (e.g., EmployeeID to e.g., SalespersonID)

# Create your first visualization

- Click on the 📶 Report View and go to the Visualizations pane on the right-hand side
- Select the type of visualization you would like to plot your data on. Keep reading this cheat to learn different visualizations available in Power Bl.
- Under the *Field* pane on the right-hand side, drag the variables of your choice into *Values* or *Axis*.

Values let you visualize aggregate measures (e.g. Total Revenue) Axis let you visualize categories (e.g. Sales Person)

# Aggregating data

>

Power BI sums numerical fields when visualizing them under *Values*. However, you can choose different aggregations

- Select the visualization you just created
- Go to the Visualizations section on the right-hand side
- Go to Values—the visualized column should be there
- On the selected column—click on the dropdown arrow 🗸 and change the aggregation (i.e., AVERAGE, MAX, COUNT, etc..)

# Data Visualizations in Power Bl

Power BI provides a wide range of data visualizations. Here is a list of the most useful visualizations you have in Power BI

- Bar Charts: Horizontal bars used for comparing specific values across categories (e.g. sales by region)
- Column Charts: Vertical columns for comparing specific values across categories
- Line Charts: Used for looking at a numeric value over time (e.g. revenue over time)
- Area Chart: Based on the line chart with the difference that the area between the axis and line is filled in (e.g. sales by month)
- E Scatter: Displays one set of numerical data along the horizontal axis and another set along the vertical axis (e.g. relation age and loan)
- Combo Chart: Combines a column chart and a line chart (e.g. actual sales performance vs target)
- Treemaps: Used to visualize categories with colored rectangles, sized with respect to their value (e.g. product category based on sales)
- Pie Chart: Circle divided into slices representing a category's proportion of the whole (e.g. market share)
- O Donut Chart: Similar to pie charts; used to show the proportion of sectors to a whole (e.g. market share)
- (C) Maps: Used to map categorical and quantitative information to spatial locations (e.g. sales per state)
- Cards: Used for displaying a single fact or single data point (e.g. total sales)
- Table: Grid used to display data in a logical series of rows and columns (e.g. all products with sold items)

# Power Query Editor in Power Bl >

Power Query is Microsoft's data transformation and data preparation engine. It is part of Power BI Desktop, and lets you connect to one or many data sources, shape and transform data to meet your needs, and load it into Power Bl.

# **Open the Power Query Editor**

## While loading data

- Underneath the *Home* tab, click on *Get Data*
- Choose any of your datasets and double click
- Click on *Transform* Data

## When data is already loaded

- Go to the 🖽 Data View
- Under Queries in the Home tab of the ribbon, click on Transform Data drop-down, then on the Transform Data button

# Using the Power Query Editor

# Removing rows

- You can remove rows dependent on their location, and properties
- Click on the *Home* tab in the *Query* ribbon
- Click on *Remove Rows* in the *Reduce Rows* group
- Choose which option to remove, whether Remove Top Rows, Remove Bottom Rows, etc..
- Choose the number of rows to remove
- You can undo your action by removing it from the Applied Steps list on the right-hand side

# Adding a new column

You can create new columns based on existing or new data

- Click on the Add Column tab in the Query ribbon
- Click on *Custom Column* in the *General* group
- Name your new column by using the New Column Name option
- Define the new column formula under the custom column formula using the available data

# **Replace values**

# You can replace one value with another value wherever that value is found in a column

- In the Power Query Editor, select the cell or column you want to replace
- Click on the column or value, and click on *Replace Values* under the *Home* tab under the *Transform* group
- Fill the Value to Find and Replace With fields to complete your operation

Appending datasets

- - Merge Queries

  - Select the first table and the second table you would like to merge
  - the second dataset



# Data profiling

- Click on the *View* tab in the *Query* ribbon
- In the Data Preview tab—tick the options you want to visualize
- Tick Column Quality to see the amount of missing data • Tick Column Distribution to see the statistical distribution under every column



Data Analysis Expressions (DAX) is a calculation language used in Power BI that lets you create calculations and perform data analysis. It is used to create calculated columns, measures, and custom tables. DAX functions are predefined formulas that perform calculations on specific values called arguments.

# Sample data

deal_size	sales_person	date	customer_name
1,000	Maria Shuttleworth	30-03-2022	Acme Inc.
3,000	Nuno Rocha	29-03-2022	Spotflix
2,300	Terence Mickey	13-04-2022	DataChamp

# Simple aggregations

# EXAMPLES

# Logical functions

create conditional results

# EXAMPLES

# **Text Functions**

- different text string.

# EXAMPLES

• Change column customer\_name be only lower case customer\_name = LOWER('sales\_data'[customer\_name])

# Date and time functions

- EXAMPLES

# You can append one dataset to another

• Click on Append Queries under the Home tab under the Combine group Select to append either Two tables or Three or more tables • Add tables to append under the provided section in the same window

You can use merge tables based on a related column

Click on Merge Queries under the Home tab under the Combine group

• Select the columns you would like to join the tables on by clicking on the column from the first dataset, and from

- Select the Join Kind that suits your operation:



• Click on Ok-new columns will be added to your current table

Data Profiling is a feature in Power Query that provides intuitive information about your data

• Tick Column Profile to see summary statistics and more detailed frequency information of columns

# DAX Expressions

Throughout this section, we'll use the columns listed in this sample table of `sales\_data`

- SUM(<column>) adds all the numbers in a column
- AVERAGE(<column>) returns the average (arithmetic mean) of all numbers in a column
- MEDIAN(<column>) returns the median of numbers in a column
- MIN/MAX(<column>) returns the smallest/biggest value in a column
- COUNT(<column>) counts the number of cells in a column that contain non-blank values
- DISTINCTCOUNT(<column>) counts the number of distinct values in a column.

• Sum of all deals — SUM('sales\_data'[deal\_size])

 Average deal size — AVERAGE('sales\_data'[deal\_size]) • Distinct number of customers — DISTINCTCOUNT('sales\_data'[customer\_name])

• IF(<logical\_test>, <value\_if\_true>[, <value\_if\_false>]) check the result of an expression and

• Create a column called large\_deal that returns "Yes" if deal\_size is bigger than 2,000 and "No" otherwise large\_deal = IF( 'sales\_data'[deal\_size] > 2000, "Yes", "No")

• LEFT(<text>, <num\_chars>) returns the specified number of characters from the start of a text • LOWER(<text>) converts a text string to all lowercase letters • UPPER (<text>) converts a text string to all uppercase letters • REPLACE(<old\_text>, <start\_num>, <num\_chars>, <new\_text>) replaces part of a text string with a

• CALENDAR(<start date>, <end date>) generates a column of continuous sets of dates • DATE(<year>, <month>, <day>) returns the specified date in the datetime format • WEEKDAY(<date>, <return\_type>) returns 1-7 corresponding to the day of the week of a date (return\_type) indicates week start and end (1: Sunday-Saturday, 2: Monday-Sunday)

• Return the day of week of each deal week\_day = WEEKDAY('sales\_data'[date], 2)

# Learn Data Skills Online at <u>www.DataCamp.com</u>



#### What is **Power BI**?

"It is Microsoft's Self-Service Business Intelligence tool for processing and analyzing data."

#### Components

> Power BI Desktop—Desktop application

 Report—Multi-page carvas visible to end users. It serves for the placement of visuals, buttons, images, slicers, etc.
 > Data—Preview pane for data loaded into a model.
 > Model—Editable scheme of relationships between tables in a model. Pages can be used in a model for easier navigation.
 > Power Query—A tool for connecting, transforming, and combining data.

"Apart from the standard version, there is also a version for Report Server."

 > Power BI Service—A cloud service enabling access to, and sharing and administration of, output data.
 > Workspace—There are three types of workspaces:
 > Personal, Team, and Develop a template app. They serve as storage and enable controlled access to output data.
 > Dashboard—A space consisting of tiles in which visuals and report pages are stored.\*

> Report—A report of pages containing visuals.\*
 > Worksheet—A published Excel worksheet. Can be used as a tile on a dashboard.

- > Dataset—A published sequence for fetching and transforming data from Power BI Desktop.
- > Dataflow—Online Power Query representing a special dataset outside of Power BI Desktop.\*
- > Application—A single location combining one or more reports or dashboards \*
- Admin portal—Administration portal that lets you configure capacities, permissions, and capabilities for individual users and workspaces.
- \*Can be created and edited in the Power BI Service environment.

> Data Gateway—On-premises data gateway that lets you transport data from an internal network or a custom device to the Power BI Service.

> Power BI Mobile—Mobile app for viewing reports. Mobile view is applied, if it exists, otherwise the desktop view is used.
 > Report Server—On-premises version of Power BI Service.
 > Report Builder—A tool for creating page reports.

#### Built-in and additional languages

#### **Built-in languages**

> M/Query Language—Lets you transform data in Power Query.

- > DAX (Data Analysis Expressions)—Lets you define custom calculated tables, columns, and measures in Power BI Desktop.
- "Both languages are natively available in Power BI, which eliminates the need to install anything."

#### Additional languages

Data

Brothers

> Python—Lets you fetch data and create visuals.
 Requires installation of the Python language on your computer and enabling Python scripting.
 > R—Lets you fetch and transform data and create visuals.
 Requires installation of the R language on your computer and enabling R scripting.

#### Power Query

Works with data fetched from data sources using connectors. This data is then processed at the Power BI app level and stored to an **in-memory** database in the program background. This means that data is not processed at the source level. The basic unit in Power Query is **query**, which means one sequence consisting of **steps**. A step is a data command that dictates what should happen to the data when it is loaded into Power BI. The basic definition of each step is based on its use:

Connecting data—Each query begins with a function that provides data for the subsequent steps. E.g., data can be loaded from Excel, SQL database, SharePoint etc. Connection steps can also be used later.

→ Transforming data—Steps that modify the structure of the data. These steps include features such as Pivot Column, converting columns to rows, grouping data, splitting columns, removing columns, etc. Transformation steps are necessary in order to clean data from not entirely clean data sources.

 Combining data—Data split into multiple source files needs to be combined so that it can be analyzed in bulk. Functions include merging queries and appending queries.

> Merge queries—This function merges queries based on the selected key. The primary query then contains a column which can be used to extract data from a secondary query. Supports typical join types:

0		0	٢	$\bigcirc$	0
Left outer	Right outer	Full outer	Inner	Left anti	Right anti

> Custom function—A query intended to apply a pre-defined sequence of steps so that the author does not need to create them repeatedly. The custom function can also accept input data (values, sheets, etc.) to be used in the sequence.

> Parameter—Values independent of datasets. These values can then be used in queries. Values enable the quick editing of a model because they can be changed in the Power BI Service environment.

#### Dataflow

The basic unit is a table or **Entity** consisting of columns or **Fields**. Just like Queries in Power Query, Entities in Dataflows consist of sequences of steps. The result of such steps is stored in native Azure Data Lake Gen 2.

"You can connect a custom Data Lake where the data will be stored."

There are three types of entities:

- Standard entity—It only works with data fetched directly from a data source or with data from non-stored entities within the same dataflow.
- **Computed entity\***—It uses data from another stored entity within the same dataflow.
- Linked entity\*—Uses data from an entity located in another dataflow. If data in the original entity is updated, the new data is directly passed to all linked entities.
- \*Can only be used in a dedicated Power BI Premium workspace.

"It supports custom functions as well as parameters."

#### DAX

Language developed for data analysis. It enables the creation of the following objects using expressions:

- > Calculated Columns
- > Calculated Tables

Each expression starts with the = sign, followed by links to tables/columns/functions/measures and

operators. The following operators are supported:
>Arithmetic {+, -, /, \*, ^}

- > Comparison { = , == , > , < , >= , <= , <> }
- > Text concatenation { & , && , II , IN }
- > Precedence { ( , ) }

Operators and functions require that all values/columns used are of the same data type or of a type that can be freely converted; such as a date or a number.

#### Visualization

Visualizations or visuals let you present data in various graphical forms, from graphs to tables, maps, and values. Some visuals are linked to other services outside Power BI, such as Power Apps.

#### 

In addition to basic visuals, Power BI supports creating custom visuals. Custom visuals can be added using a file import or from a free Marketplace offering certified and non-certified visuals. Certification is optional, but it verifies whether, among other things, a visual accesses external services and resources.

#### Themes

Serves as a single location for configuring all native graphical settings for visuals and pages.



By default, you can choose from 19 predefined themes. Custom themes can be added.

A custom theme can be applied in two different ways: > Modification of an existing theme—A native window that lets you modify a theme directly in the Power Bl environment. > Importing a JSON file—Any file you create only defines the formatting that should change. Everything else remains the same. The advantage of this approach is that you can customize any single visual.

"The resulting theme **can be exported** in the JSON format and **used** in any report without the need to create a theme from scratch."

JAK NA POWER POWER BI CHEATSHEET

#### Drill Down

**The Visual** that supports the embedding of hierarchies enables drilling down to the embedded hierarchy's individual levels using the following symbols:

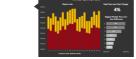
- Drill up to a higher-level hierarchy
  - Drill down to a specific field
- Drill down to the next level in the hierarchy
- Expand next-level hierarchy

#### Tooltip/Custom Tooltip

**> Tooltip** — A default detail preview pane which appears above a visual when you hover over its values.



 Custom Tooltip —A custom tooltip is a customdesigned report page identified as descriptive.
 When you hover over visual, a page appears with content filtered based on criteria specified by the value in the visual.



#### Drill-through

**Drill-through** lets you pass from a data overview visual to a page with specific details. The target page is displayed with all the applied filters affecting the value from which the drill-through originated.

257.589	Display as a table		
237,305			607
100,385	Include Exclude		12
16	Drill Through		Decomposition
134	Group		418
358,125	Сору	•	940

#### Bookmarks

**Bookmarks** capture the currently configured view or a report page visual. Later, you can go back to that state by selecting the saved bookmark. Setting options:

> Data—Stores filters, applied sort order in visuals and slicers.
 By selecting the bookmark, you can re-apply the corresponding settings.

 Display—Stores the state of the display for visuals and report elements (buttons, images, etc.). By selecting the bookmark, you can go back to the previously stored state of the display.

**Current page**—Stores the currently displayed page. By selecting the bookmark, you can go back the to stored page.

#### License

#### Per-user License

• Free—Can be obtained for any Microsoft work or school email account. Intended

for personal use. Users with this license can only use the personal workspace. They cannot share or consume shared content.

or consume snared content. "If it is not available in Premium workspace"

If it's not available in Premium Workspace<sup>2</sup>
 Pro—It is associated with a work/school account priced at €8.40 per month or it is included in the E5 license. Intended for team collaboration. Let's users access team workspaces, consume shared content, and use apps.
 Premium per User – Includes all Power BI Pro license capabilities, and adds features such as paginated reports, AI, greater frequency for refresh rate, XILA endpoint and other capabilities that are only available to Premium subscribers.

#### Per-tenant License

> Premium—Premium is set

up for individual workspaces. 0 to N workspaces can be used with a single version of this license. It provides dedicated server computing power based on license type: P1, P2, P3, P4\*, P5\*. It offers more space for datasets, extended metrics for individual workspaces, managed consumption of dedicated capacity, linking of Azure AI features with datasets, and access for users with **Free** licenses to shared content. Prices start at €4,212.30.

\*Only available upon special request. Intended for models larger than 100GB.

 Embedded—Supports embedding dashboards and reports in custom apps.
 Report Server—Included in Premium or SQL Server Enterprise licenses.

Report Server—Included in Premium or SQL Server Enterprise licenses

## Administration

on the Workspaces tab.

whole organization.

protect your content.

**External Tools** 

Featured section

external tools:

> Tabular Editor 😰

> DAX studio 🏙

> ALM Toolkit 🍞

› VertiPaq Analyzer

 Use metrics—Usage metrics let you monitor Power BI usage for your organization.

Users—The Users tab provides a link to the Microsoft 365 admin center
 Audit logs—The Audit logs tab provides a link to the Security & Compliance center.

 > Tenant settings—Tenant settings enable fine-grained control over features made available to your organization. It controls which features will be enabled or disabled and for which users and groups.

- Capacity settings—The Power BI Premium tab enables you to manage any Power BI Premium and Embedded capacities.
- Embed codes—You can view the embed codes that are generated for your tenant to share reports publicly. You can also revoke or delete codes.
   Organization visuals—You can control which type of Power BI visuals
- users can access across the organization. > Azure connections—You can control workspace-level storage permissions for Azure Data Lake Gen 2. > Workspaces—You can view the workspaces that exist in your tenant

> Custom branding-You can customize the look of Power BI for your

> Featured content—You can manage all the content promoted in the

> Protection metrics-The report shows how sensitivity labels help

They simplify the use of Power BI and extend the

capabilities offered in Power BI. These tools are

mostly developed by the community. Recommended



#### What is **DAX**?

" Data Analysis Expressions (DAX) is a library of functions and operators combined to create formulas and expressions "

#### Introduction to DAX

#### > Where to find

> Power BI, Power Pivot for Excel, Microsoft Analysis Services > Purpose

> DAX was created to enumerate formulas across the data model, where the data is stored in the form of tables, which can be linked together through the sessions. They may have a cardinality of either 1: 1, 1: N, or M: N and your direction, which decides which table filters which. These sessions are either active or inactive. The active session is automatically and participates in the calculation. The inactive is involved in this when it is activated, for example, by a function USERELATIONSHIP()

#### **Basic concepts**

> Constructs and their notation > Table – 'Table' Column – [Column] -> 'Table'[Column] > Measure - [NameOfMeasure]

Comments Single-line (CTRL + ') - // or --

> Data types > INTEGER > DECIMAI > CURRENCY > DATETIME > BOOLEAN > STRING > VARIANT (not implemented in Power BI)

> Multi-line - /\* \*/

> BINARY

> DAX can work very well with some types as well combined as if it were the same type. If so, for example, the DATETIME and INTEGER data types are supported operator "+" then it is possible to use them together.

Example: DATETIME ([Date]) + INTEGER (1) = DATETIME ([Date] + 1)

#### Operators

> Arithmetic { + , - , / , \* , ^ } > Comparative { = , == , > , < , >= , <= , <> } > Joining text { & } > Logic { && . II . IN. NOT } > Prioritization {(,)}

#### Calculated Columns

> They behave like any other column in the table. Instead of coming from a data source, they are created through a DAX expression evaluated based on the current context line, and we cannot get values of another row directly.

> Import mode. Their evaluation and storage is in progress when processing the model. > DirectQuery mode. They are evaluated at runtime, which may

slow down the model Profit = Trades[Quantity]\*Trades[UnitPrice]

#### Measures

> They do not compare row-based calculations, but they perform aggregation of row-based values input contexts that the environment passes to the calculation. Because of this, there can be no pre-counting result. It must be evaluated only at the moment when Measure is called. > The condition is that they must always be linked to the table to store their code, which is possible at any time alter. Because their calculation is no longer directly dependent, it is common practice to have one separate Measure Table, which groups all Measures into myself. For clarity, they are

✓ ■ MeasureTable

> 🗖 Average

> 🗀 Count

Example of Measure:

SalesVolume = SUM (Trades[Quantity])

therefore further divided into folders

#### Variables

> Variables in DAX calculations allow avoiding repeated recalculations of the same procedure. Which might look like this:

NumberSort = VAR selectedNumber = SELECTEDVALUE( Table[Number] ) RETURN

IF( selectedNumber < 4, selectedNumber, 5)

Their declaration uses the word VAR after followed by the name "=" and the expression. The first using the word VAR creates a section for DAX where possible declare such variables 1 to X. Individual variables always require a comment for their declaration VAR before setting the name. To end this section, the word RETURN that it defines is a necessary return point for calculations. > Variables are local only.

- > If there is a variable in the formula that is not used to get the result, this variable does not evaluate. (Lazy Evaluation)
- > Evaluation of variables is performed based on evaluated context instead of the context in which the variable is used directly. Within one. The expression can be multiple VAR / RETURN sections that always serve to evaluate the currently evaluated context.

They can store both the value and the whole table

#### Calculation contexts

All calculations are evaluated on a base basis some context that the environment brings to the

calculation. (Evaluation context)		
› Context Filter -	Country	Revenue
The following calculation calculates the profit forindividual sales.	Australia	2,838,077.18
the profit formatividual sales.	Canada	73,959.95
Revenue = SUMX( Trades,	Germany	340,392.76
Trades[Quantity]* Trades[UnitPrice]	Japan	1,833,026.16
	Mexico	320,833.27
)	Nigeria	378,202.44
If I place this calculation in a table	Total	5,784,491.77

without a Country column, then the

result will be 5,784,491.77. With this column, we get "Total" the same as the previous calculation. Still, the individual records provide us with a FILTER context that filters in calculating the input the SUMX function's input. They behave the same way, for example, AXES in the chart.

> The filter context is can be adjusted with various functions, such as FILTER.ALL. ALLSELECTED

> Row context - Unlike the previous one, this context does not filter the table. It is used to iterate over tables and evaluate values columns. They are typical, but at the same time. specific example calculated columns that are calculated from data that are valid for the table row being evaluated. In particular that, manual creation is not required when creating the line context because DAX makes it. Above the mentioned example with the use of **SUMX** also hides in itself line context. Because SUMX is the function for that specified, the table in the first argument performs an iterative pass and evaluates the calculation line by line. The line context is possible to use even nested. Or, for each row of the table, evaluates each row of a different table

#### Calculate type function

> CALCULATE, and CALCULATETABLE are functions that can programmatically set the context filter. In addition to this feature converts any existing line context to a context filter. > Calculate and Calculatetable syntax: CALCULATE / CALCULATETABLE ( <expression> [, <filter1> [, ... ]] > The section filter within the Calculate expression is NOT of type boolean but Table type. Nevertheless, boolean can be used as an argument.

> Example of using the calculate function in a cumulative calculation the sum of sales for the last 12 months: CALCULATE (

JAK NA POWER BI CHEATSHEET

SUM (Trades[Quantity]).

DATESINPERIOD( DateKey[Date],

MAX (DateKey[Date])

))

YEAR > Syntax Sugar:

> [TradeVolume](Trades[Dealer] = 1)

CALCULATE ( [TradeVolume], Trades[Dealer] = 1)

CALCULATE ( [TradeVolume], FILTER ( ALL (Trades[Dealer]) Trades[Dealer] = 1) )

#### Calcuation Groups

Example:

> They are very similar to Calculated members from MDX. In Power BI, it is not possible to create them directly in the Desktop application environment, but an External Tool Tabular Editor is required.

This is a set of Calculation Items grouped according to their purpose and whose purpose is to prepare an expression, which can be used for different input measures, so it doesn't have to write the same expression multiple times. To where she would be, but the input measure is placed SELECTEDMEASURE()

	Name 💌 O
	Hight&Low
( SELECTEDMEASURE(),	Labels
Trades[Dealer] = 1)	First&Last

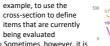
> From a visual point of view, the Calculation Group looks like a table with just two columns, "Name," "Ordinal," and rows that indicate the individual Calculation Items. In addition to facilitating the reusability of the prepared expressions also provide the ability to modify the output format of individual calculations. Within this section, "Format String Expression "often uses the DAX function SELECTEDMEASUREFORMATSTRING(), which returns a format string associated with the Measures being evaluated.



RETURN

SELECTEDMEASUREFORMATSTRING() & " " & \_selectedCurrency

> In Power BI, they can all be evaluated pre-prepared items. or it is possible, for First&Last



necessary to enable the evaluation of Calculation Items only for Specific Measures. In that case, it is possible to use the ISSELECTEDMEASURE() function, whose output is a value of type boolean or the SELECTEDMEASURENAME() function that returns the name of the currently inserted measure as a string

#### Conditions

> Like most languages, DAX uses the IF function. Within this language, it is defined by syntax:

IF ( <logical\_test>, <value\_if\_true>[, <value\_if\_false>]) Where false, the branch is optional. The IF function explicitly evaluates only a branch that is based on the result of a logical test relevant.

> If both branches need to be evaluated, then there is a function IF.EAGER() whose syntax is the same as IF itself but

evaluates as:

VAR value if true = <value if true> VAR \_value\_if\_false = <value\_if\_false>

RETURN

IF (<logical\_test>, \_value\_if\_true, \_value\_if\_false) > IF has an alternative as IFERROR. Evaluates the expression and return the output from the <value\_if\_error> branch only if the expression returns an error. Otherwise, it returns the value of the expression itself.

> DAX supports concatenation of conditions, both using submerged ones IF, so thanks to the SWITCH function. It evaluates the expression against the list values and returns one of several possible result expressions

#### Hierarchy

> DAX itself has no capability within the hierarchy to automatically convert your calculations to parent or child levels. Therefore, each level must Prepare Your Measures, which are then displayed based on the ISINSCOPE function. She tests which level to go just evaluating. Evaluation takes place from the bottom to the top level.

> The native data model used by DAX does not directly support its parent/child hierarchy. On the other hand, DAX contains functions that can convert this hierarchy to separate columns. > PATH - It accepts two parameters, where the first parameter is the key ID

column tables. The second parameter is the column that holds the parent ID of the row. The result of this function then looks like this: 1/2/3/4 Syntax: PATH( <ID\_columnName> sparent\_columnName>) > PATHITEM - Returns a specific item based on the specified position

from the string, resulting from the PATH function. Positions are counted from left to right. The inverted view uses the PATHITEMREVERSE function. Syntax: PATHITEM( <path>, <position>[, <type>] ) > PATHILENGTH – Returns the number of parent elements to the specified

item in given the PATH result, including itself. Syntax: PATHLENGTH( <path>)

> PATHCONTAINS - Returns true if the specified item is specified exists in the specified PATH path. Syntax: PATHCONTAINS( <nath>, <item>)

#### **DAX** Oueries

> The basic building block of DAX queries is the expression EVALUATE followed by any expression whose output is a table Example

**ΕVALUATE** 

ALL (Trades[Dealer])

> The EVALUATE statement can be divided into three primary sections. Each section has its specific purpose and its introductory word

> Definition - It always starts with the word DEFINE. This section defines local entities such as tables, columns, variables, and measures. There can be one section definition for an entire query, although a query can contain multiple EVALUATEs

> Query - It always starts with the word EVALUATE. This section contains the table expression to evaluate and return as a result.

> Result - This is a section that is optional and starts with the word ORDER BY. It contains the possibility to sort the result based on the inserted inputs.

Example:

DEFINE

VAR \_tax = 0.79 EVALUATE

ADDCOLUMNS

Trades.

"AdjustedpProfit" (Trades[Quantity] \* Trades[UnitPrice]) \* \_tax

#### ORDER BY [AdjustednProfit]

> This type of notation is used, for example, in DAX Studio (daxstudio.org). It is a publicly available tool that provides free access to query validation, code debugging, and query O DA performance measurement. > DAX studio has the ability to connect directly to

DAX Analysis Services, Power BI a Power Pivot for Excel STUDIO

#### Recommended sources

> Marco Russo & Alberto Ferrari > Daynatterns com > dax.guide > The Definitive Guide to DAX







#### What is **Power Query**?

#### "An IDE for M development"

#### Components

 > Ribbon – A ribbon containing settings and pre-built features by Power Query itself rewrites in M language for user convenience.
 > Queries – simply a named M expression. Queries can be moved into

groups **> Primitive** – A primitive value is a single-part value, such as a number, logical, date, text, or null. A **null** value can be used to indicate the absence of any data.

> List – The list is an ordered sequence of values. M supports endless lists. Lists define the characters "{" and "}" indicate the beginning and the end of the list.

> Record – A record is a set of fields, where the field is a pair of which form the name and value. The name is a text value that is in the field record unique.

> Table – A table is a set of values arranged in named columns and rows. Table can be operated on as if it is a list of records, or as if it is a record of lists. Table[Field]' (field reference syntax for records) returns a list of values in that field. Table[i] (list index access syntax) returns a record representing a row of the table.

> Function – A function is a value that when called using arguments creates a new value. Functions are written by listing the function argumets in parentheses, followed by the transition symbol "->-" and the expression defining the function. This expression usually refers to argumets by name. There are also functions without argumets. Parameter – The parameter stores a value that can be used for.

Parameter - The parameter stores a value that can be used for transformations. In addition to the name of the parameter and the value it stores, it also has other properties that provide metadata. The undeniable advantage of the parameter is that it can be changed from the **Power BI Service** environment without the need for direct Intervention in the data set. Syntax of parameter is as regular query only thing that is special is that the metadata follows a specific format.

 Formula Bar – Displays the currently loaded step and allows you to edit it. To be able to see formula bar, It has to be enabled in the ribbon menu inside View category.

 Query settings — Settings that include the ability to edit the name and description of the query. It also contains an overview of all currently applied steps. Applied Steps are the variables defined in a let expression and they are represented by variables names.

 Data preview — A component that displays a preview of the data in the currently selected transformation step.

> Status bar — This is the bar located at the bottom of the screen. The row contains information about the approximate state of the rows, columns, and time the data was last reviewed. In addition to this information, there is profiling source information for the columns. Here it is possible to switch the profiling from 1000 rows to the entire data set.

#### Functions in Power Query

Knowledge of functions is your best helper when working with a functional language such as  ${\bf M}.$  Functions are called with parentheses.

> Shared – Is a keyword that loads all functions (including help and example) and enumerators in result set. The call of function is made inside empty query using by = # shared

= #shared

#### Functions can be divided into two categories:

Data

Brothers

#### Data values

Each value type is associated with a literal syntax, a set of values of that type, a set of operators defined above that set of values, and an internal type attributed to the newly created values.

- Null null
   Logical true, false
- Number 1, 2, 3, ..
- > Time #time(HH,MM,SS)
- > Date #date(yyyy,mm,ss)
- > DateTime #datetime(yyyy,mm,dd,HH,MM,SS) > DateTimeZone -

#datetimezone(yyyy,mm,dd,HH,MM,SS,9,00)

> Text - "text"

> Binary - #binary("link")

- > List {1, 2, 3}
  > Record [ A = 1, B = 2 ]
- > Record [A = 1, B = 2]
  > Table #table({columns}, {{first row contenct}, {}...})\*
- > Function (x) => x + 1
- > **Type** type { number }, type table [ A = any, B = text ] \* The index of the first row of the table is the same as for the records in sheet 0
- \* The index of the first row of the table is the same as for the records in shee

#### Operators

There are several operators within the M language, but not every operator can be used for all types of values.

- > (x) Parenthesized expression
- $\mathbf{x}[\mathbf{i}]$  Field Reference. Return value from record, list of values from table.
- \*X{i} Item access. Return value from list, record from table. "Placing the "?" Character after the operator returns null if the index is not in the list "
- x(...) Function invocation
- > {1..10} Automatic list creation from 1 to 10
- ... Not implemented

#### Mathematical operators - +, -, \*, / Comparative operators

- >>, >= Greater than, greater than or equal to
- < , <= Less than, less than or equal to
- > = , <> is equal, is not equal. Equal returns true even for null = null

#### > Logical operators

- > and short-circuiting conjunction
- or short-circuiting disjunction
   not logical negation

#### > Type operators

> as – Is compatible nullable-primitive type or error

- is Test if compatible nullable-primitive type
- > Metadata The word meta assigns metadata to a value. Example of assigning metadata to variable x:
- "x meta y" or "x meta [name = x, value = 123,...]"
- Within Power Query, the priority of the operators applies, so for example "X + Y \* Z" will be evaluated as "X + (Y \* Z)"

#### Comments

- M language supports **two** versions of comments: > Single-line comments – can be created by // before code
- > Shortcut: CTRL + '
- > Multi-line comments can be created by /\* before code and \*/ after code
- > Shortcut: ALT + SHIFT + A

#### let expression

The expression let is used to capture the value from an intermediate calculation in a named variable. These named variables are local in scope to the 'let' expression. The construction of the term let looks like this:

#### 

#### in returnVariable

When it is evaluated, the following always applies:

Expressions in variables define a new range containing identifiers from the production of the list of variables and must be present when evaluating terms within a list variables. The expressions in the list of variables are they can refer to each other

- All variables must be evaluated before the term let is evaluated.
   If expressions in variables are not available, let will not be evaluated
- Errors that occur during query evaluation propagate as an error to other linked queries.

#### Conditions

Even in Power Query, there is an "if" expression, which, based on the inserted condition, decides whether the result will be a true-expression or a false-expression.

- Syntactic form of If expression:
- Condition entry: If x > 2 then 1 else 0
- If [Month] > [Fiscal\_Month] then true else false
- If expression is the only conditional in M. If you have multiple predicates to test, you must chain together like: if <predicate>
- then < true-expression >
- else if <predicate>
- then < false-true-expression >
- else < false-false-expression > When evaluating the conditions, the following applies:
- If the value created by evaluating the if a condition is not a logical value, then an error with the reason code
- "Expression.Error,, is raised
- A true-expression is evaluated only if the if condition evaluates to true. Otherwise, false-expression is evaluated.
   If expressions in variables are not available, they must not be
- evaluated > The error that occurred during the evaluation of the condition

will spread further either in the form of a failure of the entire query or "**Error**" value in the record.

#### The expression try... otherwise

Capturing errors is possible, for example, using the **try** expression. An attempt is made to evaluate the expression after the word **try**. If an error occurs during the evaluation, the expression after the word **otherwise** is applied

JAK NA POWER BI CHEATSHEET

Syntax example: try Date.From([textDate]) otherwise null

#### Custom function

Example of custom function entries: (x, y) => Number.From(x) + Number.From(y)

(x) =>

- out = Number.From(x) +
- Number.From(Date.From(DateTime.LocalNow())) in
- out

The input argumets to the functions are of two types: > Required – All commonly written argumets in (). Without these argumets, the function cannot be called. > Optional – Such a parameter may or may not be to function to enter. Mark the parameter as optional by placing text before the argument name "Optional". For example (optional x). If it does not happen fulfillment of an optional argument, so be the same for for calculation purposes, but its value will be null. Optional arguments. Syntax Sugar

add1ToField1 = ( ) => [field1] + 1,

sugar for field access of a Record named `

**Query Folding** 

addColumn(Source, "NewName", add1ToField1)

equivalent:

Source =

addColumn

Source = ....

this feature.

Valid functions

> Row filtering

> Invalid functions

DEMO

List.Transform/

Keywords

#sections, #shared, #table, #time

#table(

))

> Remove, Rename columns

> Adding columns with Index

LastStep[Year]{[ID]}

\*This means that you can get the

> Change the data type of a column

> Grouping, summarizing, pivot and unpivot

> Connect queries based on the same data source

> Merge queries based on different data sources

Operators can be combined. For example, as follows:

> Production of a DateKey dimension goes like this:

value from another step based on the index of the column

type table [Date=date, Day=Int64,Type, Month=Int64,Type,

MonthName=text, Year=Int64.Type,Quarter=Int64.Type],

Date.MonthName(\_), Date.Year(\_), Date.QuarterOfYear(\_)}

and, as, each, else, error, false, if, in, is, let, meta, not,

otherwise, or, section, shared, then, true, try, type, #binary,

#date, #datetime, #datetimezone, #duration, #infinity, #nan,

List.Dates(start date, (start date-endd ate),

#duration(1, 0, 0, 0)), each { , Date.Day( ), Date.Month( ), 1899

> Merge and extract data from queries

> Add custom columns with simple logic

let

let

> Each is essentially a syntactic abbreviation for declaring non-

addColumn = Table.AddColumn(Source, "NewName", each [field1] + 1)

The second piece of syntax sugar is that bare square brackets are syntax

As the name implies, it is about composing. Specifically, the

is then implemented against the data source. Data sources

that supports Query folding are resources that support the

concept of query languages as relational database sources.

This means that, for example, a CSV or XML file as a flat file

with data will definitely not be supported by Query Folding.

until after the data is loaded, but it is possible to get the data

ready immediately. Unfortunately, not every source supports

Therefore, the transformation does not have to take place

steps in Power Query are composed into a single query, which

type functions, using a single formal parameter named.

Therefore, the following notations are semantically

Arguments can be annotated with 'as <type>' to indicate required type of the argument. The function will throw a type error if called with arguments of the wrong type. Functions can also have annotated return of them. This annotation is provided as:

(x as number, y as text) as logical => <expression>

let first = (x) => () => let out =  $\{1, x\}$  in out in first

> Errors caused by evaluating expressions in a list of

further either as a failure or as an "Error" value

**Recursive functions** 

if x = 0 then 1 else x \* @Factorial(x - 1).

expressions or in a function expression will propagate

The number of arguments created from the argument

the function, otherwise an error will occur with reason

For recursive functions is necessary to use the character "@"

recursive function is the factorial. The function for the factorial

Functions can be called against specific arguments. However, if

sheet, or an entire column in a table, it is necessary to append

context record, it applies the procedure behind it. Each is never

required! It simply makes it easier to define a function in-line

the function needs to be executed for each record, an entire

the word **each** to the code. As the name implies, for each

for functions which require a function as their argument.

which refers to the function within its calculation. A typical

list must be compatible with the formal argumets of

When evaluating functions, it holds that:

code "Expression Error"

can be written as follows:

Factorial = (x) =>

Result // = 6

Each

Result = Factorial(3)

let

in

The return of the functions is very different. The output can be a sheet, a table, one value but also other functions. This means that one function can produce another function. Such a function is written as follows:



# 7.3 Data Visualization

# datacamp

# **The Data Visualization Cheat Sheet**

Learn Data Visualization online at <a href="http://www.DataCamp.com">www.DataCamp.com</a>

# How to use this cheat sheet

Use this cheat sheet for inspiration when making your next data visualizations. For more data visualization cheat sheets, check out our cheat sheets repository here.

# Capture a trend

# Line chart

The most straightforward way to capture how a numeric variable is changing over time

## USE CASES

- 1. Revenue in \$ over time
- 2. Energy consumption in kWh over time
- 3. Google searches over time

# Multi-line chart



Captures multiple numeric variables over time. It can include multiple axes allowing comparison of different units and scale ranges

# USE CASES

- 1. Apple vs Amazon stocks over time
- 2. Lebron vs Steph Curry
- searches over time 3. Bitcoin vs Ethereum price over time





Shows how a numeric value progresses by shading the area between line and the x-axis

## USE CASES

1. Total sales over time 2. Active users over time

# Stacked area chart



Most commonly used variation of area charts, the best use is to track the breakdown of a numeric value by subgroups

## USE CASES

- 1. Active users over time by segment
- 2. Total revenue over time by country

# Visualize relationships

## Bar chart



One of the easiest charts to read which helps in quick comparison of categorical data. One axis contains categories and the other axis represents values

## USE CASES

- 1. Volume of google searches by region
- 2. Market share in revenue by product

## Column chart



Also known as a vertical bar chart, where the categories are placed on the x-axis. These are preferred over bar charts for short labels, date ranges, or negatives in values

## USE CASES

1. Brand market share 2. Profit Analysis by region

# Scatter plot

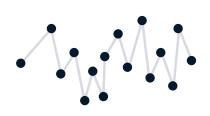


Most commonly used chart when observing the relationship between two variables. It is especially useful for quickly surfacing potential correlations between data points

## USE CASES

- 1. Display the relationship between time-on-platform and churn 2. Display the relationship
- between salary and years spent at company

# **Connected scatterplot**



A hybrid between a scatter plot and a line plot, the scatter dots are connected with a line

# USE CASES

- 1. Cryptocurrency price index
- 2. Visualizing timelines and events when analyzing two variables

# Part-to-whole charts

# **Pie chart**

One of the most common ways to show part to whole data. It is also commonly used with percentages

## USE CASES

1. Voting preference by age group 2. Market share of cloud providers

# Donut pie chart



The donut pie chart is a variant of the pie chart, the difference being it has a hole in the center for readability

## USE CASES

- 1. Android OS market share
- 2. Monthly sales by channel

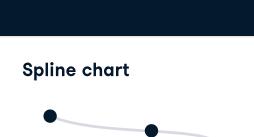
# **Heat maps**

that use color shading to represent data trends.

#### USE CASES

- 1. Average monthly temperatures
- across the year

2. Departments with the highest amount of attrition over time





It differs in that data points are connected with smoothed curves to account for missing values, as opposed to straight lines

## USE CASES

- 1. Electricity consumption over time
- 2. CO2 emissions over time

# Visualize a single value

## Card



Cards are great for showing and tracking KPIs in dashboards or presentations

#### USE CASES

- 1. Revenue to date on a sales dashboard
- 2. Total sign-ups after a promotion



 $\sim$  $\sim$ 

Best to be used on small datasets, it displays tabular data in a table

# USE CASES

- 1. Account executive
- leaderboards
- 2. Registrations per webinar

# Gauge chart



This chart is often used in executive dashboard reports to show relevant KPIs

## USE CASES

- 1. NPS scores
- 2. Revenue to target

# **Bubble chart**



Often used to visualize data points with 3 dimensions, namely visualized on the xaxis, y-axis, and with the size of the bubble. It tries to show relations between data points using location and size

## USE CASES

- 1. Adwords analysis: CPC vs Conversions vs Share of total conversions 2. Relationship between life
- expectancy, GDP per capita, & population size

## Word cloud chart

Data XXX Analyst Science Engineer

A convenient visualization for visualizing the most prevalent

## USE CASES

1. Top 100 used words by customers in customer service tickets

# Visualize a flow

# Sankey chart



Useful for representing flows in systems. This flow can be any measurable quantity

## USE CASES

1. Energy flow between countries

# Chord chart



Useful for presenting weighted relationships or flows between nodes. Especially useful for highlighting the dominant or important flows

## USE CASES

- 1. Export between countries to showcase biggest export partners 2. Supply chain volumes
- between the largest warehouses

# **Network chart**

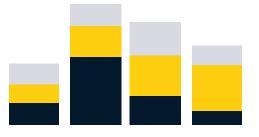
Similar to a graph, it consists of nodes and interconnected edges. It illustrates how different items have relationships with each other

- 2. Supply chain volumes between warehouses

words that appear in a text

Heatmaps are two-dimensional charts

# Stacked column chart



Best to compare subcategories within categorical data. Can also be used to compare percentages

# USE CASES

- 1. Quarterly sales per region
- 2. Total car sales by producer

## **Treemap charts**

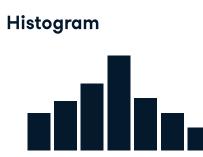


2D rectangles whose size is proportional to the value being measured and can be used to display hierarchically structured data

## USE CASES

- 1. Grocery sales count with categories
- 2. Stock price comparison by industry and company

# Capture distributions



Shows the distribution of a variable. It converts numerical data into bins as columns. The x-axis shows the range, and the y-axis represents the frequency

# USE CASES

1. Distribution of salaries in an organization 2. Distribution of height in one cohort

# Box plot

Shows the distribution of a variable using 5 key summary statistics minimum, first quartile, median, third quartile, and maximum

# USE CASES

1. Gas efficiency of vehicles 2. Time spent reading across readers

# **Violin plot**

A variation of the box plot.

1. Time spent in restaurants

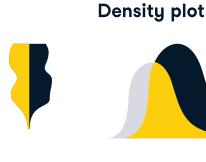
2. Length of pill effects by

across age groups

It also shows the full

USE CASES

distribution of the data



Visualizes a distribution by using smoothing to allow smoother distributions and alongside summary statistics better capture the distribution shape of the data

## USE CASES

- 1. Distribution of price of hotel listings
- 2. Comparing NPS scores by customer segment



## USE CASES

1. How different airports are connected worldwide 2. Social media friend group analysis

# Learn Data Skills Online at www.DataCamp.com



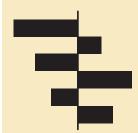
# Deviation

Emphasise variations (+/-) from a fixed reference point. Typically the reference point is zero but it can also be a target or a long-term average. Can also be used to show sentiment (positive/neutral/negative).

Example FT uses

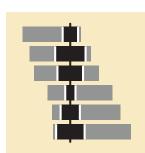
Trade surplus/deficit, climate change

# **Diverging bar**



A simple standard bar chart that can handle both negative and positive magnitude values.

# Diverging stacked bar



Perfect for presenting survey results which involve sentiment (eg disagree/neutral/ agree).

# Spine chart

Splits a single value into 2 contrasting components (eg Male/Female). 

# Surplus/deficit filled line



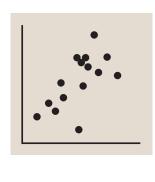
The shaded area of these charts allows a balance to be shown – either against a baseline or between two series.

# Correlation

Show the relationship between two or more variables. Be mindful that, unless you tell them otherwise, many readers will assume the relationships you show them to be causal (i.e. one causes the other).

Example FT uses Inflation & unemployment, income & life expectancy

# Scatterplot



The standard way to show the relationship between two continuous variables, each of which has its own axis.

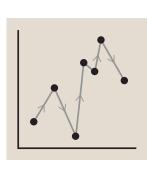
# Line + Column



A good way of showing the relationship between an amount (columns)

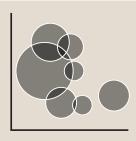
and a rate (line).

# **Connected scatterplot**

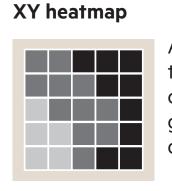


Usually used to show how the relationship between 2 variables has changed over time.

# Bubble



Like a scatterplot, but adds additional detail by sizing the circles according to a third variable.



A good way of showing the patterns between 2 categories of data, less good at showing fine differences in amounts.

# **VISUal** vocabulary

# Designing with data

There are so many ways to visualise data - how do we know which one to pick? Use the categories across the top to decide which data relationship is most important in your story, then look at the different types of chart within the category to form some initial ideas about what might work best. This list is not meant to be exhaustive, nor a wizard, but is a useful starting point for making informative and meaningful data visualisations.

FT graphic: Alan Smith; Chris Campbell; Ian Bott; Liz Faunce; Graham Parrish; Billy Ehrenberg; Paul McCallum; Martin Stabe Inspired by the Graphic Continuum by Jon Schwabish and Severino Ribecca



# Ranking

Use where an item's position in an ordered list is more important than its absolute or relative value. Don't be afraid to highlight the points of interest.

Example FT uses Wealth, deprivation, league tables, constituency election results

# Ordered bar



Standard bar charts display the ranks of values much more easily when sorted into order.

Ordered column



Dot strip plot

.....

Lollipop chart

**\_\_\_\_** 

Slope

•••••

Ordered proportional symbol Use when there are big variations between

See above.

values and/or seeing fine differences between data is not so important.

Dots placed in order on a strip are a space-efficient method of laying out ranks across multiple categories.

> Perfect for showing how ranks have changed over time or vary between **•** categories.

> > Lollipops draw more attention to the data value than standard bar/column and can also show rank and value effectively.

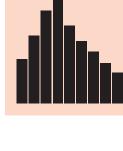
# Distribution

Show values in a dataset and how often they occur. The shape (or 'skew') of a distribution can be a memorable way of highlighting the lack of uniformity or equality in the data.

# Example FT uses

Income distribution, population (age/sex) distribution

# Histogram



Boxplot

\_\_\_\_

The standard way to show a statistical distribution - keep the gaps between columns small to highlight the 'shape' of the data.

Summarise multiple distributions by showing the median (centre) and range of the data

Similar to a box plot

complex distributions

(data that cannot be

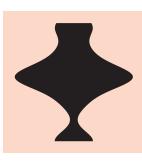
summarised with

simple average).

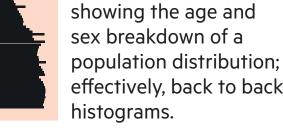
but more effective with

# Violin plot

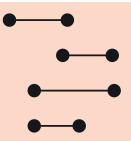
÷



# Population pyramid A standard way for



Good for showing many dots have the



A simple way of or range (min/max) of data across

# Barcode plot



Like dot strip plots, good for displaying all the data in a table, they work best when highlighting individual values.

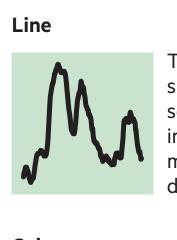
Cumulative curve

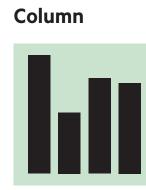
A good way of showing how unequal a distribution is: y axis is always cumulative frequency, x axis is always a measure.

# Change over Time

Give emphasis to changing trends. These can be short (intra-day) movements or extended series traversing decades or centuries: Choosing the correct time period is important to provide suitable context for the reader.

Example FT uses Share price movements, economic time series

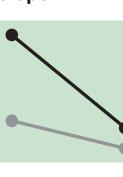








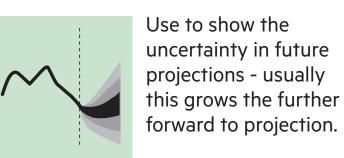




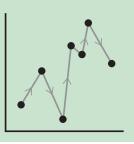
Area chart



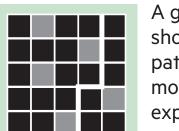
# Fan chart (projections)



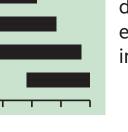
Connected scatterplot



Calendar heatmap



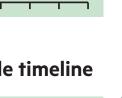
**Priestley timeline** 

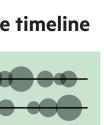


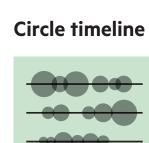
Seismogram

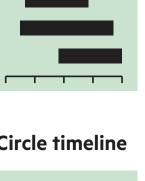
Another alternative to the circle timeline for showing series where there are big variations in the data.





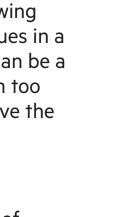




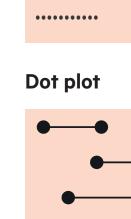












•••••••••

same value.

Dot strip plot

individual values in a distribution, can be a ••• ••• problem when too

showing the change multiple categories.

The standard way to show a changing time series. If data are irregular, consider markers to represent data points.

Columns work well for showing change over time - but usually best with only one series of data at

A good way of showing the relationship over time between an amount (columns) and a rate (line).

Usually focused on day-to-day activity, these charts show 
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I</t

> Good for showing changing data as long as the data can be simplified into 2 or 3 points without missing a key part of story.

Use with care – these are good at showing changes to total, but seeing change in components can be very difficult.

Use to show the uncertainty in future projections - usually this grows the further

A good way of showing changing data for two variables whenever there is a relatively clear pattern of progression.

> A great way of showing temporal patterns (daily, weekly monthly) – at the expense of showing precision in quantity

Great when date and duration are key elements of the story in the data.

Good for showing discrete values of varying size across — multiple categories (eg earthquakes by contintent).



Show how a single entity can be broken down into its component elements. If the reader's interest is solely in the size of the components consider a magnitude-type chart instead.

Example FT uses

Fiscal budgets, company structures, national election results

# Stacked column



showing part-to-whole relationships but can be difficult to read with more than a few components

A simple way of

# Proportional stacked bar



A good way of showing the size and proportion of data at the same time – as long as the data are not too complicated.

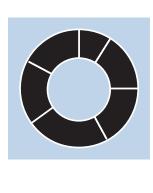


A common way of showing part-to-whole data – but be aware that it's difficult to

accurately compare the

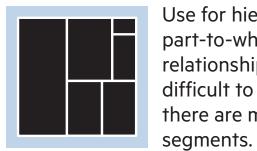
size of the segments.

# Donut



Similar to a pie chart – but the centre can be a good way of making space to include more information about the data (eg. total).

# Treemap



\_\_\_\_\_ Use for hierarchical part-to-whole relationships; can be difficult to read when there are many small

# Voronoi

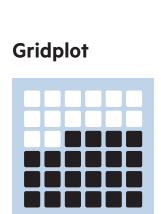


A way of turning points into areas – any point within each area is closer to the central point than any other centroid.



visualisaing hierarchical part-to-whole relationships. Use sparingly (if at all) for obvious reasons.

A hemicycle, often used for visualising political results in parliaments.



multiple layout form



Generally only used for schematic



# Magnitude

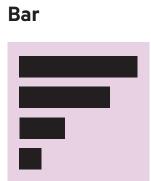
Show size comparisons. These can be relative (just being able to see larger/bigger) or absolute (need to see fine differences). Usually these show a 'counted' number (for example, barrels, dollars or people) rather than a calculated rate or per cent.

Example FT uses Commodity production, market capitalisation





The standard way to compare the size of things. Must always start at 0 on the axis.

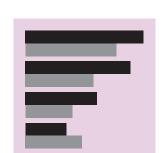


See above. Good when the data are not time series and labels have long category names.

# Paired column

As per standard column but allows for multiple series. Can become tricky to read with more than 2 series.

Paired bar



# See above.

# Proportional stacked bar



A good way of showing the size and proportion of data at

important.

Proportional symbol

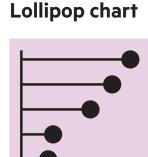
Use when there are big variations between values and/or seeing fine differences

between data is not so

not too complicated.

lsotype (pictogram)

Excellent solution in **Some instances – use** only with whole numbers (do not slice off an arm to represent a decimal).



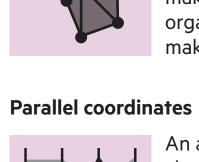
Lollipop charts draw more attention to the data value than standard bar/column does not HAVE to start

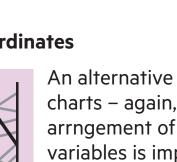
Radar chart



charts – again, the arrngement of the

at zero (but preferable) A space-efficient way of showing value pf multiple variables– but make sure they are

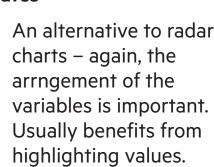


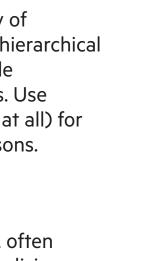


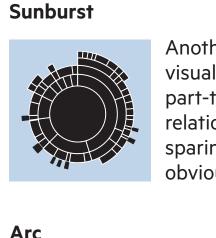




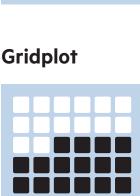
organised in a way that makes sense to reader.



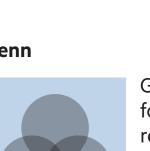


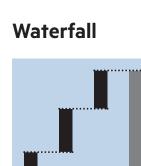


Arc



Good for showing % information, they work best when used **On whole numbers** and work well in

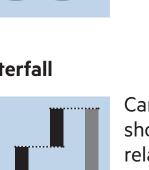




showing part-to-whole relationships where some of the components are negative.

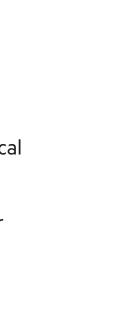
representation.





Can be useful for







# Spatial

Used only when precise locations or geographical patterns in data are more important to the reader than anything else.

Example FT uses Locator maps, population density, natural resource locations, natural disaster risk/impact, catchment areas, variation in election results

# Basic choropleth (rate/ratio)

The standard approach for putting data on a map – should always be rates rather than totals and use a sensible base geography.

Use for totals rather

than rates – be wary

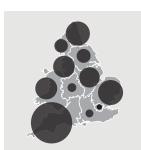
that small differences

in data will be hard to

Proportional symbol (count/magnitde)

see

For showing



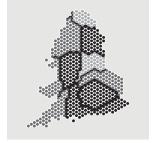
Flow map



Contour map

For showing areas of equal value on a map. Can use deviation colour schemes for showing +/- values

# Equalised cartogram



Converting each unit on a map to a regular and equally-sized shape – good for representing voting regions with equal value.

Scaled cartogram (value)

Stretching and shrinking a map so that each area is sized according to a particular value.

# Dot density



Heat map

Used to show the location of individual events/locations – make sure to annotate any patterns the reader should see.

Grid-based data values mapped with an

intensity colour scale. As choropleth map – but not snapped to an admin/political unit.



Show the reader volumes or intensity of movement between two or more states or conditions. These might be logical sequences or geographical locations.

# Example FT uses

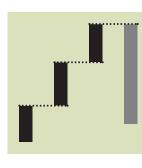
Movement of funds, trade, migrants, lawsuits, information; relationship graphs.

# Sankey



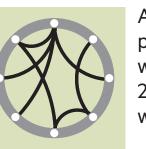
Shows changes in flows from one condition to at least one other; good for tracing the eventual outcome of a complex process.

# Waterfall



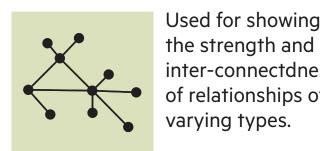
Designed to show the sequencing of data through a flow process, typically budgets. Can include +/- components.

# Chord



A complex but A complex but powerful diagram which can illustrate 2-way flows (and net winner) in a matrix.

# Network



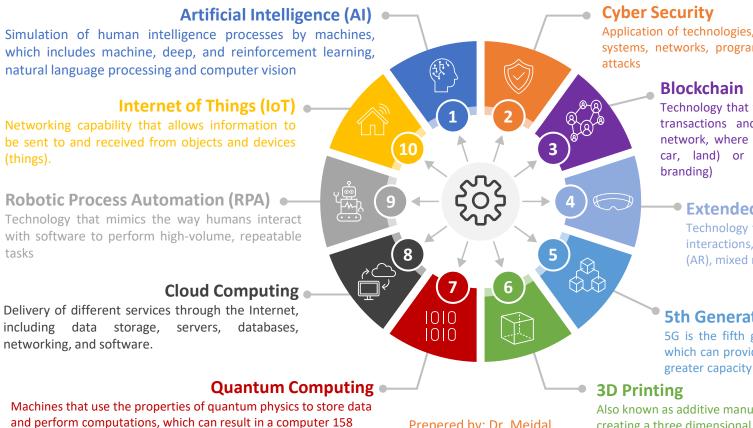
Used for showing the strength and inter-connectdness of relationships of





# 8. Emerging Technologies

# Top 10 Technology Trends



million times faster than the most sophisticated supercomputer we have in the world today

Prepered by: Dr. Mejdal

Application of technologies, processes and controls to protect systems, networks, programs, devices and data from cyber attacks

Technology that facilitates the process of recording transactions and tracking assets in a business network, where an asset can be tangible (house, car, land) or intangible (patents, copyrights, branding)

## • Extended Reality (XR)

Technology that represents all human-machine interactions, which includes augmented reality (AR), mixed reality (MR) and virtual reality (VR)

## 5th Generation Network (5G)

5G is the fifth generation of wireless technology, which can provide higher speed, lower latency and greater capacity than 4G LTE networks

Also known as additive manufacturing, is a method of creating a three dimensional object layer-by-layer using a computer created design.

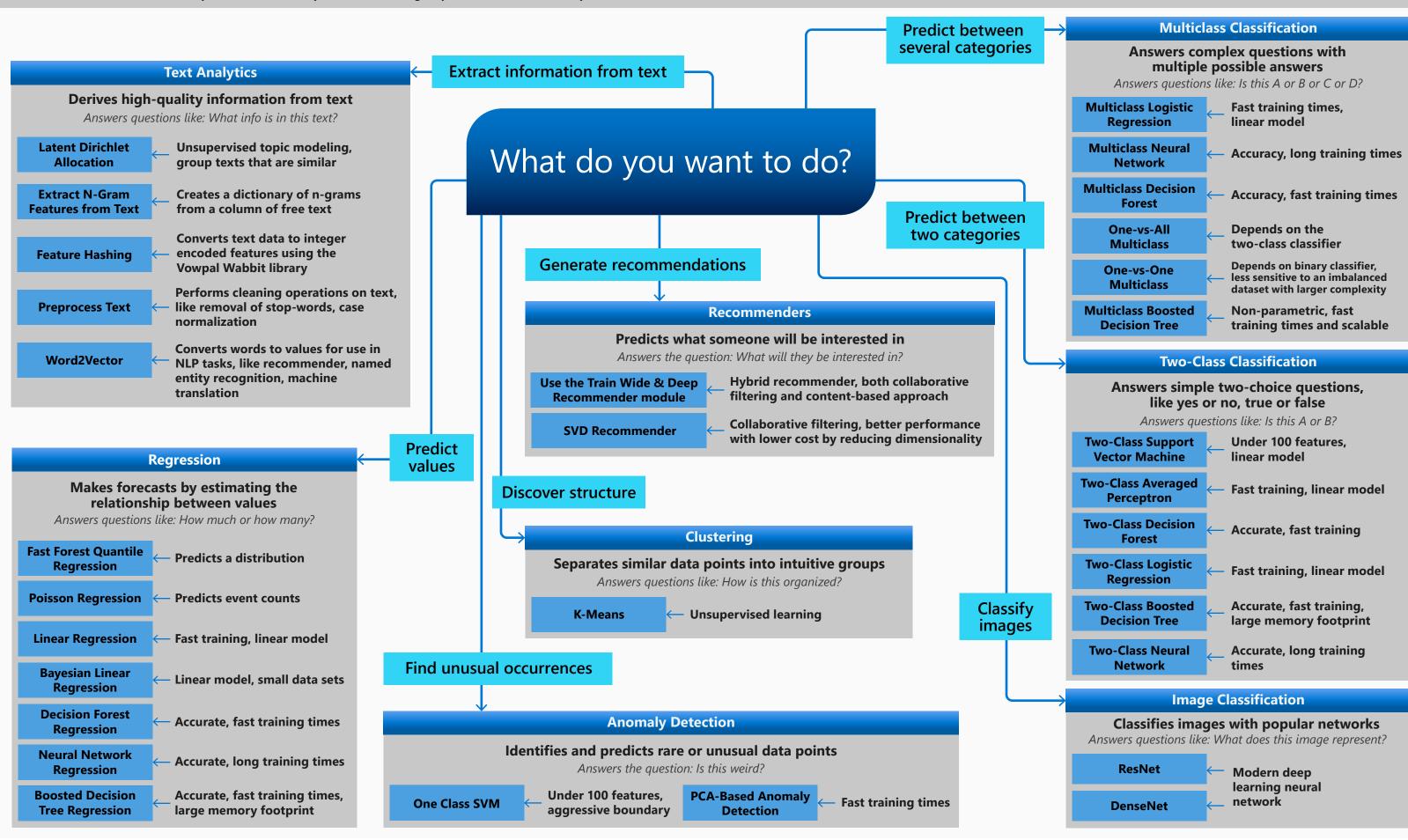


# 8.1 Artificial Intelligence (AI)



# **Machine Learning Algorithm Cheat Sheet**

This cheat sheet helps you choose the best machine learning algorithm for your predictive analytics solution. Your decision is driven by both the nature of your data and the goal you want to achieve with your data.



© 2021 Microsoft Corporation. All rights reserved. Share this poster: aka.ms/mlcheatsheet





# 8.2 Internet of Things (IoT)



Cheax shear Internet of Things

The Internet of Things (IoT) refers to a system of interrelated, internet-connected objects that can collect and transfer data over a wireless network without human intervention. IoT devices are basically "smart" devices that make "smart" systems (for e.g., smart home, smart factory, smart farms, smart cities, etc.).

Overview	Challenge	Solution	Results
If machines could talk - and tell you when they were going to break down.	When a machine is not functioning properly, productivity drops. How can you tell, quickly, that something is wrong?	loT sensors on agricultural machinery were connected to the cloud. Managers & workers received immediate notifications when the machine was not functioning properly.	This also resulted in predictive data analysis from the information provided. Engineers started to build a picture and predict when the machine would break and plan maintenance. This increased productivity to 100%.
Remote Monitoring - Improving work performance.	Staff shortages for a civil engineering company reduced the number of onsite visits that could be carried out.	Rather than manually taking readings and producing reports, data was collected on loT sensors, uploaded to the cloud and this generated a report.	Improved efficiencies, enabling the company to take on more projects and increase their customer base.
Blinding the competition - how to get a USP with loT.	Consumers want smarter, simpler everyday devices.	Blinds that are IoT enabled - features include a camera, thermometer, lighting and sensors, allowing for seamless integration into the home environment.	"Work still on-going on loT products at Bloc Blinds as it is an exciting space with loads of potential!" <u>Read how</u> <u>Magherafelt's Bloc</u> <u>Blinds has</u> incorporated IoT technologies





European Union European Regional Development Fund

Investment for Growth and Jobs



# 8.3 Robotic Process Automation (RPA)

# ROBOTIC PROCESS AUTOMATION CHEAT SHEET

# **RPA Basics**

## RPA

RPA is a Robotic Process Automation which is used for automating the current workflows with the help of robots to reduce human intervention at every point Robotic: Machines that mimics the human activities and actions are called as **robots** Process: Sequence of steps which is used to perform a particular task Automation: Any process which is done by robot without any human intervention and provides high degree of accuracy

VENDOR	Free Version	Pricing	Usability	Global Coverage
Another Monday				Europe
AntWorks				
Arago				
Automation Anywhere		Per Process	Drag & Drop, Macro Recording	Global
BluePrism		Per bot	Drag & Drop	
Contextor				EMEA & North America
Jidoka				
Kofax				Global
Kryon Systems				
Nice systems				Global
Pega				Global
Redwood software				
UiPath	UiPath Community Edition	Per bot	Drag & Drop, Macro Recording	Global
Visual Cron	45 day free trial	Per server	Drag & Drop	
WorkFusion	WorkFusion RPA Express	Per process	Drag & Drop, Macro Recording	Global

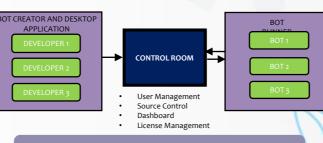
## Blue Prism

It is built with a simple concept by replicating user activity. **SQL Server** is the tool where data is stored

#### BUSINESS LOCIC +User Interface VIRTUAL BUSINESS LOCIC +User Interface VIRTUAL BUSINESS LOCIC +User Interface VIRTUAL BUSINESS LOCIC +User Interface BUSINESS BUSINESS LOCIC +User Interface BUSINESS BUSI

## Automation Anywhere

The architecture consists of **client**, **bots** and **control room**. Automation Anywhere is a software designed to virtually automate any computer process. **SQL express** is a tool where data is stored.



## Ui Path

It is a simple software automation and application integration expert. **UiPath** consists of three parts

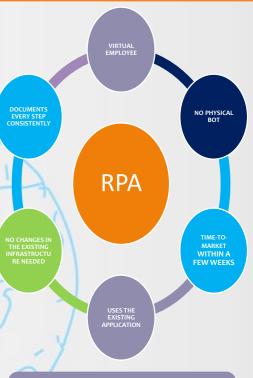
- UiPath studio
- UiPath robots
- UiPath orchestrator
   Browser (desktop) Mobile devices
   Information tier
   PORTAL
   Logic tier
   Access Manager
   Data tier
   External Users Db

## **RPA Tool Factors**

- Technology: Many organizations perform the tasks outside the local system using Virtual Machines or Citrix. Hence, the tool must be platform independent and support any type of application
- Scalability: How quickly and easily the tool responds to the business requirements
- Security: Implementations of security controls must be measured
- Total Cost: The initial set-up cost, vendor license fee, maintenance cost must be taken into consideration while selecting the tool
- Ease of use and control: The tool must be user friendly to increase the efficiency and employee satisfaction
- Vendor Experience: It improves the speed of implementation by reducing the work required to implement RPA software
- Maintenance and support: To make sure that the required service level agreements are met
   Quick Deployment: The tool must be able to work as a real end-user by interacting with applications.

## Benefits of Using RPA

- Reduces burden on IT: It does not disturb underlying legacy systems
- Reliability: As the bots can work 24\*7 effectively Cost cutting technology: It reduces the costs by reducing the size of manual workforce
- No coding required: To use RPA tools, a person need not have the programming skills
- Accuracy: It functions with accuracy and is less prone to errors
- Productivity rate: Execution time is much faster than the manual process approach
- Compliance: It follows the rules to provide audit free trail
- Consistency: Repetitive tasks are performed in the same way
- Increase employee engagement: It lets the employee to focus on value-added activities



## Applications of RPA

- It is used in customer service, to automate service order management and quality reporting
- Travel and logistics: for ticket booking, passenger details and accounting
- Human Resource: New employees joining formalities, payroll process and hiring shortlisted candidates
- Health care: In patient registrations and billing
- Banking and Financial services: It can be used for card activation and fraud claims and discovery
- Government: Change of address and license renewal



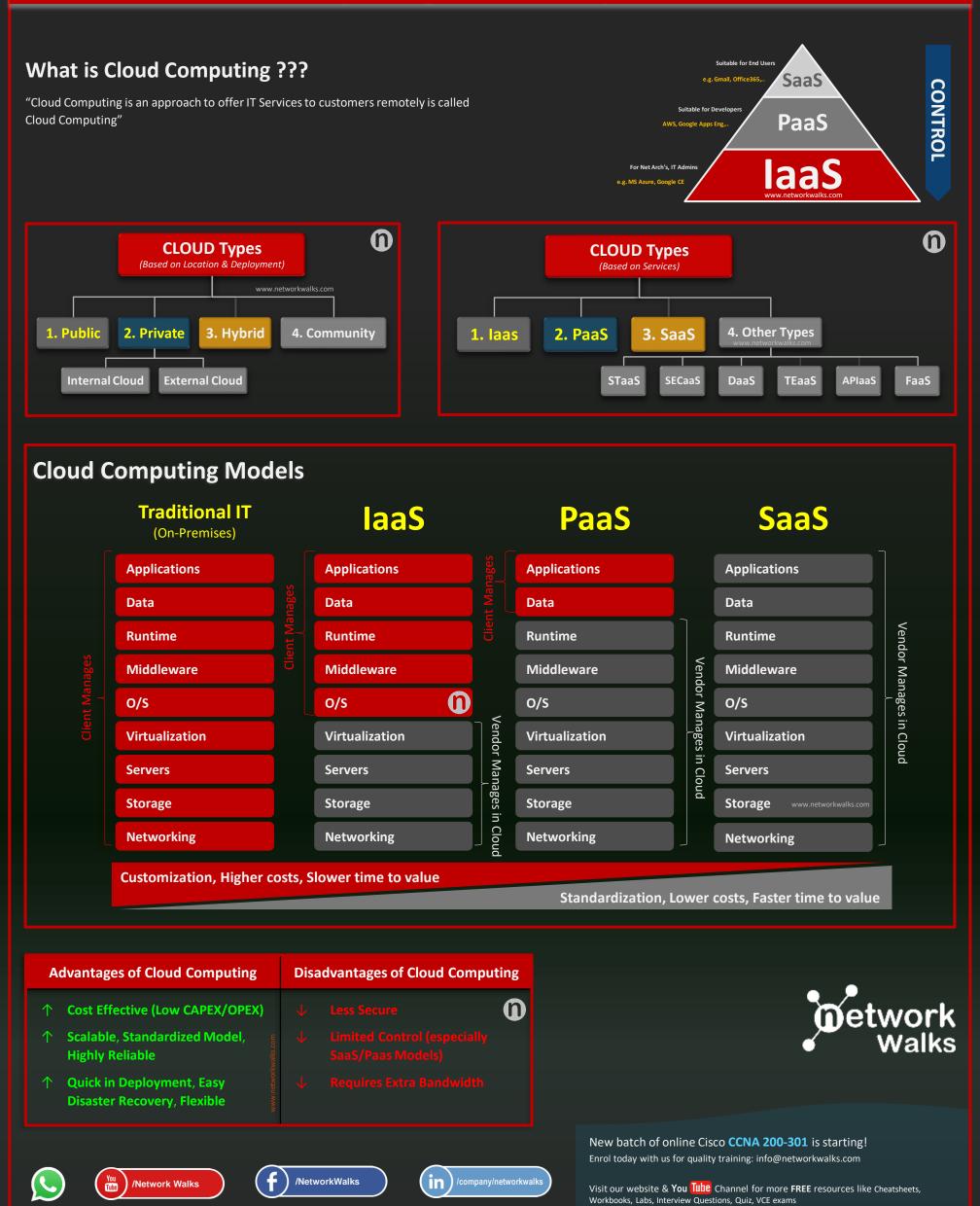
RPA Training using UiPath



# 8.4 Cloud Computing

# **Cloud Computing – Summary Cheat sheet**

Version 2



Your Feedback, Comments are always Welcomed: info@networkwalks.com

n

Network Walks Training Academy www.networkwalks.com

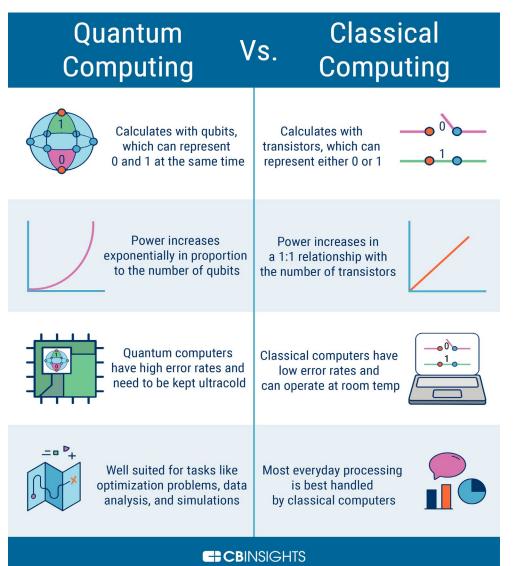


# 8.5 Quantum Computing

#### • History of Quantum computing

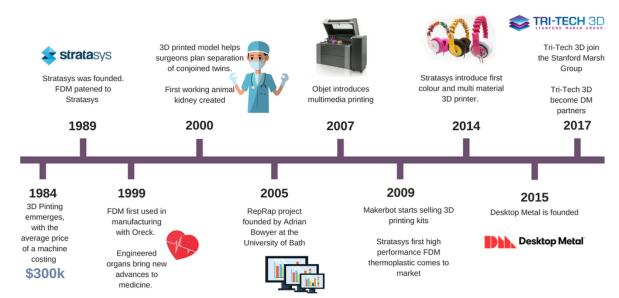


Classic computing vs Quantum computing



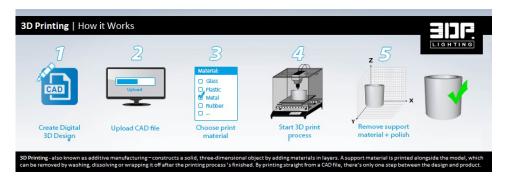


# 8.6 3D Printing

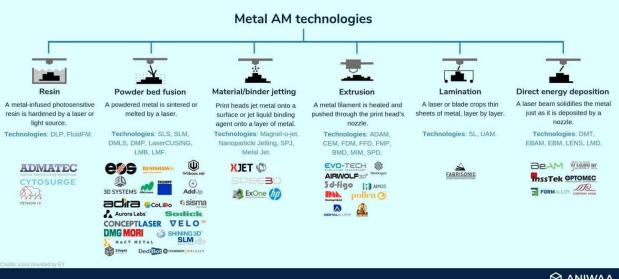


#### History of 3D printing Technology •

How 3D printing works? •



Type of 3D Printing Technologies •



# Cheatography

## 3D Printing Options Cheat Sheet by [deleted] via cheatography.com/2754/cs/16505/

#### Introduction

Our manufacturing industry generally – and 3D printing specifically – is driven by innovation. Indeed, key technological developments and new applications in industrial-grade 3D printing, or additive manufacturing, continue to advance this technology, which has only been around for a little more than 30 years.

Designers and engineers can now choose from several distinct classes of 3D printing technologies. Your choice of "tool" just depends on what it is you're designing and what its final application is. Here's a brief roundup of some of the main industrial-grade 3D printing options:

Source: https://www.medicaldesignandoutsourcing.com/3d-printingoptions-medical-device-development/

#### Stereolithography (SL)

Stereolithography (SL) uses an ultraviolet laser that draws on the surface of a liquid thermoset resin to create thousands of thin layers until final parts are formed. SL is used to create concept models, cosmetic prototypes, and complex parts with intricate geometries.

#### Selective laser sintering (SLS)

Selective laser sintering (SLS) uses a CO2 laser that lightly fuses nylon-based powder, layer by layer, until final thermoplastic parts are created. SLS produces accurate prototypes and functional production parts.

#### Direct metal laser sintering (DMLS)

Direct metal laser sintering (DMLS) uses a fiber laser system that draws onto a surface of atomized metal powder, welding the powder into fully dense metal parts. DMLS builds fully functional metal prototypes and production parts and works well to reduce metal components in multipart assemblies.

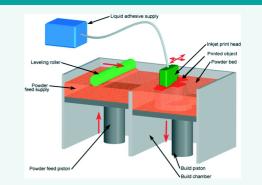
#### PolyJet

PolyJet uses a jetting process in which small droplets of liquid photopolymer are sprayed from multiple jets onto a build platform and cured in layers that form elastomeric parts. PolyJet produces multi-material prototypes with flexible features at varying durometers and is often used to concept overmolding designs.

#### By [deleted] cheatography.com/deleted-2754/

Published 23rd December, 2018. Last updated 23rd December, 2018. Page 1 of 1. Sponsored by **Readable.com** Measure your website readability! https://readable.com

#### **3D Printing**



#### **Fused Deposition Modeling (FDM)**

Fused deposition modeling (FDM) works by feeding a filament or spool of plastic into a heated nozzle, which then extrudes successive layers of thermoplastics onto the workpiece. FDM offers a wide thermoplastic material selection and is leveraged for iterative prototyping.

#### Continuous Liquid Interface Production (CLIP)

Carbon is the name of the company that is using a process called CLIP, Continuous Liquid Interface Production, which builds parts from the top down, unlike other additive technologies that work from the bottom up. Final plastic parts exhibit excellent mechanical properties and surface finishes.

#### Multi Jet Fusion

Multi Jet Fusion process selectively applies fusing and detailing agents across a bed of nylon powder, which are fused in thousands of layers by heating elements into a solid functional component. Final parts exhibit improved surface roughness, fine feature resolution, and more isotropic mechanical properties when compared to processes like SLS.

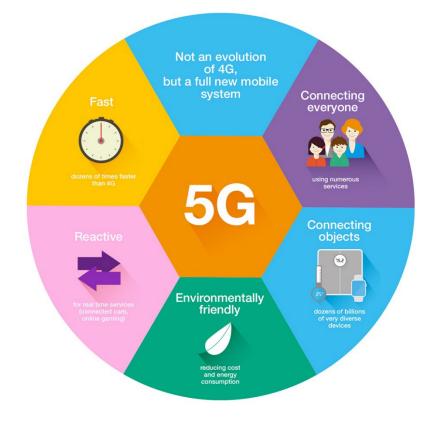


# 8.7 5th Generation Network (5G)

## • History of 5G network



## • Advantages of 5G network





# 8.8 Extended Reality (XR)

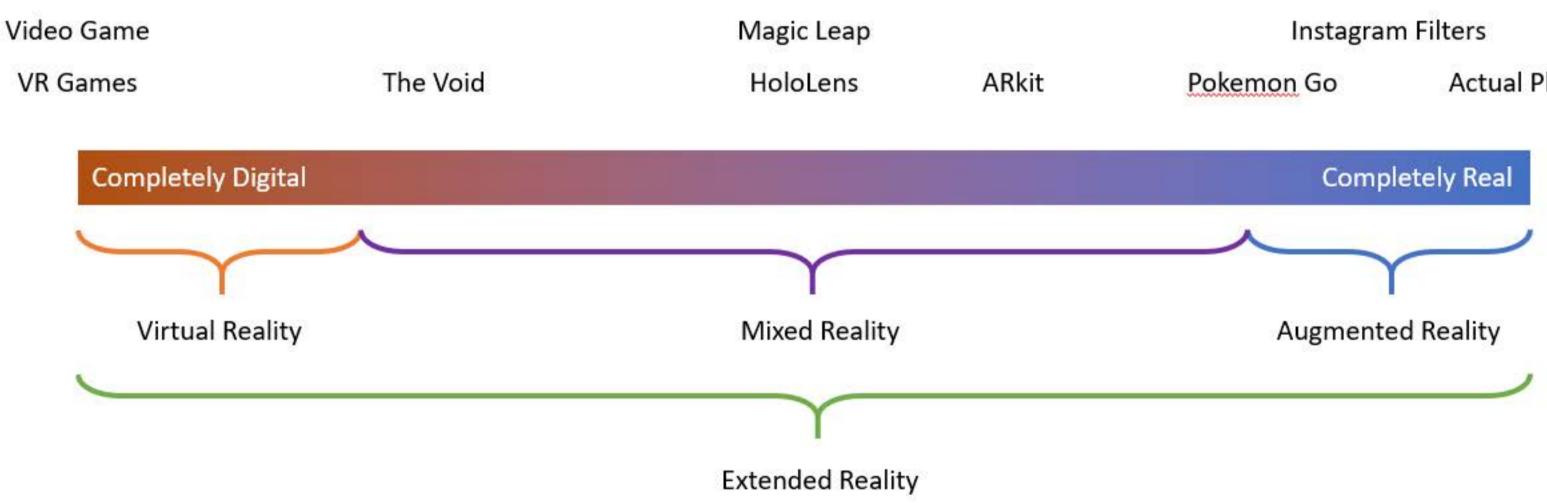


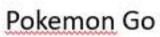
Chears Sheer **Immersive** Tech

Immersive technologies create experiences by merging the physical world with a digital / simulated reality. Augmented reality (AR) and virtual reality (VR) are the two main types of immersive tech. AR blends computer-generated information onto the user's real environment. VR uses computergenerated information to provide a full sense of immersion.

Overview	Challenge	Solution	Results
VR versus Covid-19 - Continuing business remotely.	Due to Covid-19, restrictions a company who delivered forklift driver training was unable to provide practical instruction sessions.	Use of virtual reality to train construction workers to use a forklift remotely. Also being used to deliver Health & Safety training to oil rig workers before they are deployed. VR is being used to show potential apprentices what experience they would receive without distracting workers from their activities.	Training occurred remotely using VR for the first two days, then onsite with the forklift on the 3rd. This improved productivity, as the forklift became available for work purposes for two days, when it was previously required for training.
Adding value through AR – Business growth in a post Covid-19 environment.	An interior design consultancy was struggling with a lack of business due to Covid-19 restrictions. Designers were unable to visit customers' homes.	Customers used AR to take measurements, and designers were able to create and place items within their 3D world.	They were able to survive Covid–19 by bringing their business online. Also, they can now access a wider client base and scale their business, as with AR their market is no longer restricted by distance.
Creating Tourist Destinations using AR – increasing visitor numbers & revenue.	Tourists, day trippers, and locals are all looking for new experiences.	An interactive AR app to create a huge scale solar system trail. There is also an app for those following the trail giving information on the planets and solar system.	Will reach millions, bring people together and showcase UK creativity globally. <u>Read more about</u> <u>Northern Ireland's</u> <u>Our Place in Space –</u> <u>a 3 dimensional</u> <u>sculpture trail,</u> <u>interactive AR app!</u>
Invest Northern Ireland	European Union European Regional Development Fund Investment for Growth and Jobs	Development Fund under the Investment for Growth & Jobs Northern Ireland (2014–2020) Programme.	

# Reality – Virtuality Spectrum

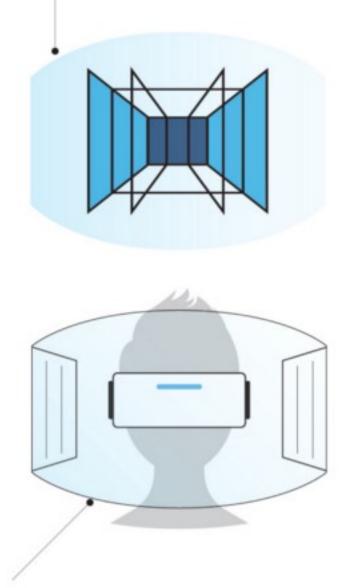




**Actual Photos** 

# VIRTUAL REALITY (VR)

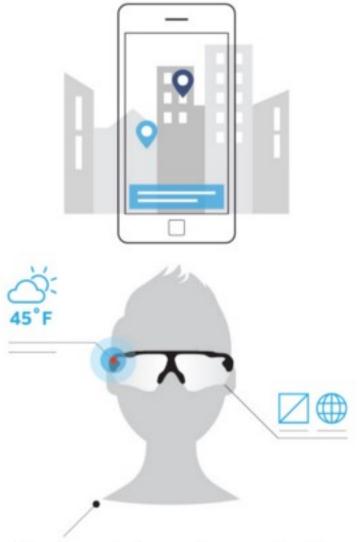
Completely digital environment



Fully enclosed, synthetic experience with no sense of the real world.

# **AUGMENTED REALITY (AR)**

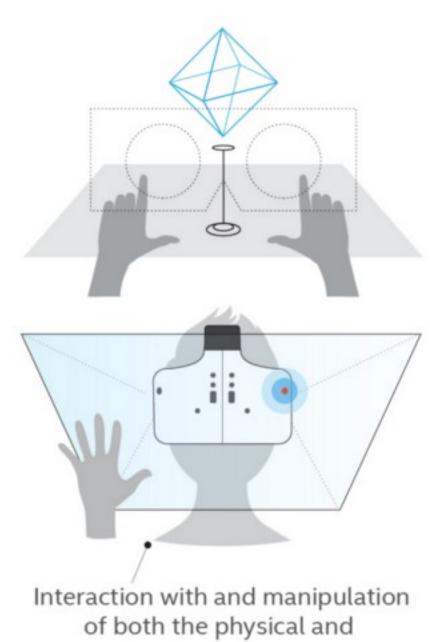
Real world with digital information overlay



Real world remains central to the experience, enhanced by virtual details.

# MERGED REALITY (MR)

Real and the virtual are intertwined



virtual environment.



# 8.9 Blockchain





Blockchain is a shared, immutable ledger for recording transactions and tracking assets in a business network. With blockchain technology, companies can establish a transparent supply chain, verify the authenticity of product claims and gain consumer trust.

Overview	Challenge	Solution	Results
Connecting through Blockchain – Increasing productivity and profitability.	How to track devices through a supply chain in the medical devices sector?	Each supplier was invited to the blockchain, added data about the parts they were making, and added status.	Able to view the origin of raw materials, have an estimated delivery of parts from 30 different suppliers, and guarantee provenance of medical grade products. This led to a significant increase in productivity, approximately 5-10% increase in turnover due to increased capacity, and 70% reduction in manual report creation.
Food and Beverage Traceability – Showcasing the truth of your product.	How to create a safer and more transparent supply chain?	Used blockchain technology and included a QR code on food and beverage products to allow the customer to see exactly how the product was made.	Using blockchain enhanced the transparency of product information, reduced food fraud, and waste. It increased food safety and profitability among food suppliers. And it improved consumer information, boosting trust & sales. See how a County Down enterprise uses blockchain to reveal everything that a consumer would want to know about their beer!
Invest <b>Northern</b>	European Union European Regional Development Fund		

Investment for Growth and Jobs

Ireland



# 8.10 Cyber Security





Visit <u>ref.customguide.com</u>

#### Security Risks

Businesses worldwide are at risk for security breaches. While large, well-known companies seem like a likely target, small and mediumsized organizations and individuals are also at risk. There are many ways data can be compromised, including viruses, phishing scams, hardware and software vulnerabilities, and network security holes.

#### Did you know?





**11%** of U.S. adults have had personal information stolen

1 in 5 people have had an email or social media

account hacked



applications

are vulnerable



Only **40%** of adults know how to protect themselves online

#### **Confidential Information**

When dealing with security, confidentiality means private information is never viewed by unauthorized parties. Confidential information must only be accessible to those authorized to view the sensitive data. Confidential information includes:

#### **Personal Information**

- Social Security Number
- Home Address
- Salary History
- Performance Issues
- Credit Card Numbers

#### **Corporate Information**

- Processes
- Customer Lists
- Research and Development
- Business Strategies
- Objectives and Projections

## Firewalls

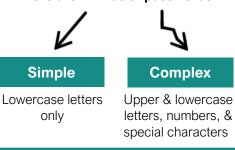
A firewall acts like a security guard and prevents unauthorized people or programs from accessing a network or computer from the Internet. There are hardwarebased firewalls, which create a protective barrier between internal networks and the outside world, and software firewalls, which are often part of your operating system.

#### Passwords

The first line of defense in maintaining system security is using complex passwords. Use passwords that are at least 8 characters long and include a combination of numbers, upper and lowercase letters, and special characters. Hackers have tools that can break easy passwords in just a few minutes.

#### How Long Does it Take to Crack a Password?

#### There are 2 kinds of passwords:



#### Complex passwords are **EXPONENTIALLY** More difficult to crack

Use them!

Here's how long it takes to crack a password when it's **simple** vs. **complex** 

Characters	Password	Time to crack
8[	ghiouhel ghiouH3l	<mark>4</mark> hours, 7 min 6 months
9 —[	houtheouh Houtheo!2	<mark>4</mark> days, 11 hours 1060 years
10 —[	ghotuhilhg gh34uhilh!	112 days 1500 years
11 —[	wopthiendhf w3pthi7ndh!	8 years, 3 months 232,800 years
12 —[	whithgildnzq @hi3hg5ldnq!	210 years 15,368,300 years
		Source: mywot.com

© 2021 CustomGuide, Inc.

#### Malware

Malware is short for "malicious software." It is written to infect the host computer. Common types of malware include:



Replicating computer program that infects computers



Hijacks your computer or browser and displays annoying advertisements



Secretly tracks your internet activities and information

Malicious program that tries to trick you into running it

TROJAN

#### Online Browsing

Browsers communicate to websites with a protocol called HTTP, which stands for Hyper Text Transfer Protocol. HTTPS is the secure version of HTTP. Websites that use HTTPS encrypt all communication between your browser and the site.



https://www.website.com

like a padlock, in the address bar

to show the site is secure. You

should always ensure security

when logging in or transferring

confidential information.



Sites without HTTPS are not secure and should never be used when dealing with person

http://www.website.com

used when dealing with personal data. If you are simply reading an article or checking the weather, HTTP is acceptable.

## Email and Phishing

#### Network Security

- Use Wi-Fi password security and change the default password
- Set permissions for shared files
- Only connect to known, secure public Wi-Fi and ensure HTTPS-enabled sites are used for sensitive data
- Keep your operating system updated
- Perform regular security checks
- Browse smart!

A phishing email tries to trick consumers into providing confidential data to steal money or information. These emails appear to be from a credible source, such as a bank, government entity, or service provider. Here are some things to look for in a phishing email:

