A mostly complete chart of

# Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

**Legend:**
- Backfed Input Cell
- Input Cell
- Noisy Input Cell
- Hidden Cell
- Probablistic Hidden Cell
- Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- Different Memory Cell
- Kernel
- Convolution or Pool

Perceptron (P)

Feed Forward (FF)

Radial Basis Network (RBF)

Deep Feed Forward (DFF)

Recurrent Neural Network (RNN)

Long / Short Term Memory (LSTM)

Gated Recurrent Unit (GRU)

Auto Encoder (AE)

Variational AE (VAE)

Denoising AE (DAE)

Sparse AE (SAE)

Markov Chain (MC)

Hopfield Network (HN)

Boltzmann Machine (BM)

Restricted BM (RBM)

Deep Belief Network (DBN)

Deep Convolutional Network (DCN)

Deconvolutional Network (DN)

Deep Convolutional Inverse Graphics Network (DCIGN)

Generative Adversarial Network (GAN)

Liquid State Machine (LSM)

Extreme Learning Machine (ELM)

Echo State Network (ESN)

Deep Residual Network (DRN)

Kohonen Network (KN)

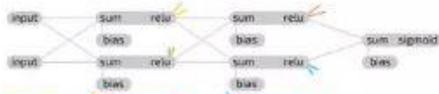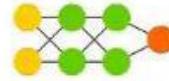Support Vector Machine (SVM)

Neural Turing Machine (NTM)

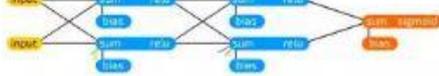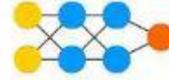An informative chart to build

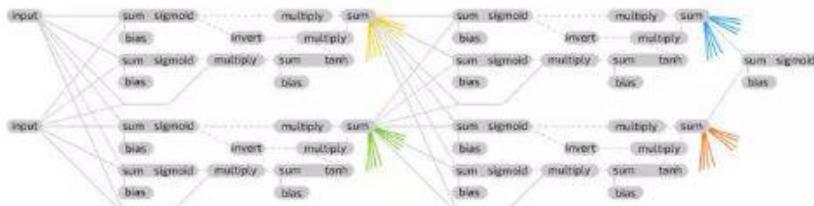# Neural Network Graphs

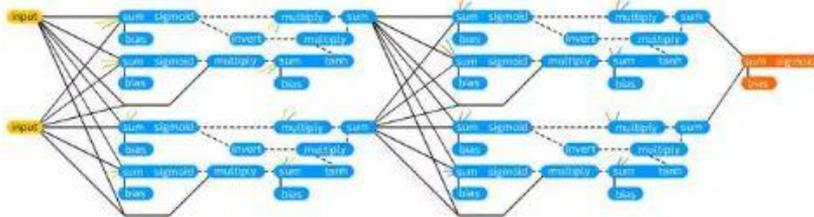©2016 Fjodor van Veen - asimovinstitute.org
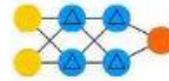


Deep Feed Forward Example
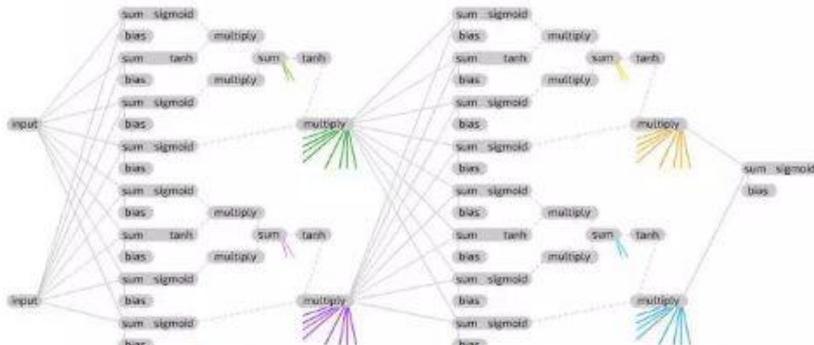
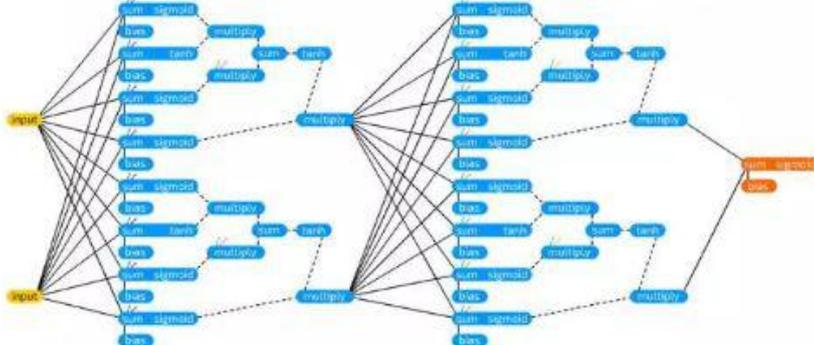Deep Recurrent Example
(previous iteration)

Deep Recurrent Example

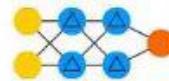Deep GRU Example
(previous iteration)

Deep GRU Example

Deep LSTM Example
(previous iteration)

Deep LSTM Example

## Linear Vector Spaces:

**Definition:** A linear vector space, $X$, is a set of elements (vectors) defined over a scalar field, $F$, that satisfies the following conditions:

1) if $x \in X$ and $y \in X$ then $x+y \in X$. 2) $x+y = y+x$. 3) $(x+y)+z = x+(y+z)$
4) There is a unique vector $0 \in X$, such that $x+0 = x$ for all $x \in X$.
5) For each vector $x \in X$ there is a unique vector in X, to be called $(-x)$, such that $x+(-x) = 0$. 6) multiplication, for all scalars $a \in F$, and all vectors $x \in X$,
7) For any $x \in X$, $1 \cdot x = x$ (for scalar 1).
8) For any two scalars $a \in F$ and $b \in F$ and any $x \in X$, $a(bx) = (ab)x$.
9) $(a+b)x = ax + bx$. 10) $a(x+y) = ax + ay$.

## Linear Independence:
Consider n vectors $\{x_1, x_2, \ldots, x_n\}$. If there exists n scalars $a_1, a_2, \ldots, a_n$, at least one of which is nonzero, such that $a_1 x_1 + a_2 x_2 + \ldots + a_n x_n = 0$, then the $\{x_i\}$ are linearly dependent.

## Spanning a Space:
Let $X$ be a linear vector space and let $\{u_1, u_2, \ldots, u_n\}$ be a subset of vectors in $X$. This subset spans X if and only if for every vector $x \in X$ there exist scalars $x_1, x_2, \ldots, x_n$ such that $x = x_1 u_1 + x_2 u_2 + \ldots + x_n u_n$.

## Inner Product:
$(x,y)$ for any scalar function of x and y.
1. $(x,y) = (y,x)$ 2. $(x, ay_1 + by_2) = a(x,y_1) + b(x,y_2)$
3. $(x,x) \geq 0$, where equality holds iff x is the zero vector.

## Norm:
A scalar function $\|x\|$ is called a norm if it satisfies:
1. $\|x\| \geq 0$ 2. $\|x\| = 0$ if and only if $x = 0$.
3. $\|ax\| = |a|\|x\|$ 4. $\|x + y\| \leq \|x\| + \|y\|$

## Angle:
The angle $\theta$ bet. 2 vectors x and y is defined by $\cos\theta = \frac{(x,y)}{\|x\| \|y\|}$

## Orthogonality:
2 vectors $x, y \in X$ are said to be orthogonal if $(x,y) = 0$.

## Gram Schmidt Orthogonalization:
Assume that we have $n$ independent vectors $y_1, y_2, \ldots, y_n$. From these vectors we will obtain $n$ orthogonal vectors $v_1, v_2, \ldots, v_n$.

$$v_1 = y_1, \quad v_k = y_k - \sum_{i=1}^{k-1} \frac{(v_i, y_k)}{(v_i, v_i)} v_i,$$

where $\frac{(v_i, y_k)}{(v_i, v_i)} v_i$ is the projection of $y_k$ on $v_i$

## Vector Expansions:

$$x = \sum_{i=1}^{n} x_i v_i = x_1 v_1 + x_2 v_2 + \cdots + x_n v_n,$$

$$\text{for orthogonal vectors, } x_j = \frac{(v_j, x)}{(v_j, v_j)}$$

## Reciprocal Basis Vectors:

$$(r_i, v_j) = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases}, \quad x_j = (r_j, x)$$

To compute the reciprocal basis vectors: set $\mathbf{B} = [v_1 \ v_2 \ \ldots \ v_n]$,
$\mathbf{R} = [r_1 \ r_2 \ldots \ r_n]$, $\mathbf{R}^T = \mathbf{B}^{-1}$ In matrix form: $\mathbf{x}^v = \mathbf{B}^{-1} \mathbf{x}^s$

## Transformations:
A *transformation* consists of three parts:
domain: $X = \{x_i\}$, range: $Y = \{y_i\}$, and a rule relating each $x_i \in X$ to an element $y_i \in Y$.

**Linear Transformations:** transformation $A$ is *linear* if:
1. for all $x_1, x_2 \in X$, $A(x_1 + x_2) = A(x_1) + A(x_2)$
2. for all $x \in X$, $a \in R$, $A(ax) = aA(x)$

**Matrix Representations:**
Let $\{v_1, v_2, \ldots, v_n\}$ be a basis for vector space $X$, and let $\{u_1, u_2, \ldots, u_n\}$ be a basis for vector space $Y$. Let $A$ be a linear transformation with domain $X$ and range $Y$: $A(x) = y$
The coefficients of the matrix representation are obtained from

$$A(v_j) = \sum_{i=1}^{m} a_{ij} u_i$$

**Change of Basis:** $B_t = [t_1 \ t_2 \ldots t_n]$, $B_w = [w_1 \ w_2 \ldots w_n]$
$$A' = [B_w^{-1} A B_t]$$

**Eigenvalues & Eigenvectors:** $Az = \lambda z$, $|[A - \lambda I]| = 0$
**Diagonalization:** $B = [z_1 \ z_2 \ldots z_n]$,
where $\{z_1, z_2, \ldots, z_n\}$ are the eigenvectors of a square matrix A,
$[B^{-1}AB] = diag([\lambda_1 \ \lambda_2 \ldots \lambda_n])$

## Perceptron Architecture:
$\mathbf{a} = hardlim(\mathbf{Wp} + \mathbf{b})$, $\mathbf{W} = [\ _1\mathbf{w}^T \ _2\mathbf{w}^T \ \ldots \ _S\mathbf{w}^T]^T$,
$$a_i = hardlim(n_i) = hardlim(\ _i\mathbf{w}^T \mathbf{p} + b_i)$$

**Decision Boundary:** $\ _i\mathbf{w}^T \mathbf{p} + b_i = 0$
The decision boundary is always orthogonal to the weight vector.
Single-layer perceptrons can only classify linearly separable vectors.

## Perceptron Learning Rule
$$\mathbf{W}^{new} = \mathbf{W}^{old} + \mathbf{e}\mathbf{p}^T, \mathbf{b}^{new} = \mathbf{b}^{old} + \mathbf{e},$$
$$\text{where} \quad \mathbf{e} = \mathbf{t} - \mathbf{a}$$

## Hebb's Postulate:
"When an axon of cell A is near enough to excite a cell B and repeatedly or persistently takes part in firing it, some growth process or metabolic change takes place in one or both cells such that A's efficiency, as one of the cells firing B, is increased."

## Linear Associator: $\mathbf{a} = purelin(\mathbf{Wp})$

## The Hebb Rule:
*Supervised Form:* $w_{ij}^{new} = w_{ij}^{old} + t_{qi} P_{qi}$
$$\mathbf{W} = \mathbf{t}_1 \mathbf{P}_1^T + \mathbf{t}_2 \mathbf{P}_2^T + \cdots + \mathbf{t}_Q \mathbf{P}_Q^T$$

$$\mathbf{W} = [\mathbf{t}_1 \ \mathbf{t}_2 \ldots \mathbf{t}_Q] \begin{bmatrix} \mathbf{p}_1^T \\ \mathbf{p}_2^T \\ \vdots \\ \mathbf{p}_Q^T \end{bmatrix} = \mathbf{TP}^T$$

## Pseudoinverse Rule: $\mathbf{W} = \mathbf{TP}^+$
When the number, $R$, of rows of $\mathbf{P}$ is greater than the number of columns, $Q$, of $\mathbf{P}$ and the columns of $\mathbf{P}$ are independent, then the pseudoinverse can be computed by $\mathbf{P}^+ = (\mathbf{P}^T \mathbf{P})^{-1} \mathbf{P}^T$

## Variations of Hebbian Learning:
**Filtered Learning** (Ch.16): $\mathbf{W}^{new} = (1 - \gamma)\mathbf{W}^{old} + \alpha \mathbf{t}_q \mathbf{p}_q^T$

**Delta Rule** (Ch.10): $\mathbf{W}^{new} = \mathbf{W}^{old} + \alpha(\mathbf{t}_q - \mathbf{a}_q)\mathbf{p}_q^T$

**Unsupervised Hebb** (Ch.13): $\mathbf{W}^{new} = \mathbf{W}^{old} + \alpha \mathbf{a}_q \mathbf{p}_q^T$

## Taylor:
$F(\mathbf{x}) = F(\mathbf{x}^*) + \nabla F(\mathbf{x})^T|_{\mathbf{x}=\mathbf{x}^*} (\mathbf{x} - \mathbf{x}^*) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)\nabla^2 F(\mathbf{x})^T|_{\mathbf{x}=\mathbf{x}^*}(\mathbf{x} - \mathbf{x}^*) + \cdots$

## Grad
$\nabla F(\mathbf{x}) = \left[\frac{\partial}{\partial x_1} F(\mathbf{x}) \ \frac{\partial}{\partial x_2} F(\mathbf{x}) \ \ldots \ \frac{\partial}{\partial x_n} F(\mathbf{x})\right]^T$

## Hessian:
$$\nabla^2 F(\mathbf{x}) = \begin{bmatrix} \frac{\partial}{\partial x_1^2} F(\mathbf{x}) & \frac{\partial}{\partial x_1 \partial x_2} F(\mathbf{x}) & \ldots & \frac{\partial}{\partial x_1 \partial x_n} F(\mathbf{x}) \\ \frac{\partial}{\partial x_2 \partial x_1} F(\mathbf{x}) & \frac{\partial}{\partial x_2^2} F(\mathbf{x}) & \ldots & \frac{\partial}{\partial x_2 \partial x_n} F(\mathbf{x}) \\ \vdots & \vdots & & \vdots \\ \frac{\partial}{\partial x_n \partial x_1} F(\mathbf{x}) & \frac{\partial}{\partial x_n \partial x_2} F(\mathbf{x}) & \ldots & \frac{\partial}{\partial x_n^2} F(\mathbf{x}) \end{bmatrix}$$

## Directional Derivatives:
**1st Dir.Der.:** $\frac{\mathbf{p}^T \nabla F(\mathbf{x})}{\|\mathbf{p}\|}$, **2nd Dir.Der.:** $\frac{\mathbf{p}^T \nabla^2 F(\mathbf{x})\mathbf{p}}{\|\mathbf{p}\|^2}$

## Minima:
*Strong Minimum:* if a scalar $\delta > 0$ exists, such that $F(x) < F(x + \Delta x)$ for all $\Delta x$ such that $\delta > \|\Delta x\| > 0$.
*Global Minimum:* if $F(x) < F(x + \Delta x)$ for all $\Delta x \neq 0$
*Weak Minimum:* if it is not a strong minimum, and a scalar $\delta > 0$ exists, such that $F(x) \leq F(x + \Delta x)$ for all $\Delta x$ such that $\delta > \|\Delta x\| > 0$.

## Necessary Conditions for Optimality:
*1st-Order Condition:* $\nabla F(x)|_{x=x^*} = 0$ (Stationary Points)
*2nd-Order Condition:* $\nabla^2 F(x)|_{x=x^*} \geq 0$ (Positive Semi-definite Hessian Matrix).

## Quadratic fn.: $F(x) = \frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{d}^T \mathbf{x} + c$
$\nabla F(x) = \mathbf{Ax} + \mathbf{d}$, $\nabla^2 F(x) = \mathbf{A}$, $\lambda_{min} \leq \frac{\mathbf{p}^T \mathbf{A} \mathbf{p}}{\|\mathbf{p}\|^2} \leq \lambda_{max}$

**General Minimization Algorithm:**
$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k \quad or \quad \Delta\mathbf{x}_k = (\mathbf{x}_{k+1} - \mathbf{x}_k) = \alpha_k \mathbf{p}_k$$
**Steepest Descent Algorithm:**
$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{g}_k \quad where, \quad \mathbf{g}_k = \nabla F(\mathbf{x})|_{\mathbf{x}=\mathbf{x}_k}$$
**Stable Learning Rate:** $(\alpha_k = \alpha, \text{ constant}) \ \alpha < \frac{2}{\lambda_{max}}$

$\{\lambda_1 \ \lambda_2, \dots, \lambda_n\}$ Eigenvalues of Hessian matrix A
**Learning Rate to Minimize Along the Line:**
$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k \overset{is}{\Rightarrow} \alpha_k = -\frac{\mathbf{g}_k^T \mathbf{p}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k} \text{ (For quadratic fn.)}$$
**After Minimization Along the Line:**
$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k \ \Rightarrow \ \mathbf{g}_{k+1}^T \mathbf{p}_k = 0$$

**ADALINE:** $\mathbf{a} = purelin(\mathbf{Wp} + \mathbf{b})$
**Mean Square Error:** *(for ADALINE it is a quadratic fn.)*
$$F(\mathbf{x}) = E[e^2] = E[(t-a)^2] = E[(t - \mathbf{x}^T \mathbf{z})^2]$$
$$F(\mathbf{x}) = c - 2\mathbf{x}^T \mathbf{h} + \mathbf{x}^T \mathbf{R} \mathbf{x},$$
$c = E[t^2], \ \mathbf{h} = E[t\mathbf{z}] \ and \ \mathbf{R} = E[\mathbf{zz}^T] \Rightarrow \mathbf{A} = 2\mathbf{R}, \mathbf{d} = -2\mathbf{h}$
Unique minimum, if it exists, is $\mathbf{x}^* = \mathbf{R}^{-1}\mathbf{h}$,
where $\mathbf{x} = \begin{bmatrix} _1\mathbf{w} \\ b \end{bmatrix}$ and $\mathbf{z} = \begin{bmatrix} \mathbf{p} \\ 1 \end{bmatrix}$
**LMS Algorithm:** $\mathbf{W}(k+1) = \mathbf{W}(k) + 2\alpha\, \mathbf{e}(k)\, \mathbf{p}^T(k)$
$$\mathbf{b}(k+1) = \mathbf{b}(k) + 2\alpha\, \mathbf{e}(k)$$
**Convergence Point:** $\mathbf{x}^* = \mathbf{R}^{-1}\mathbf{h}$
**Stable Learning Rate:** $0 < \alpha < 1/\lambda_{max}$ where
$\lambda_{max}$ is the maximum eigenvalue of R
**Adaptive Filter ADALINE:**
$$a(k) = purelin(\mathbf{W}\mathbf{p}(k) + b) = \sum_{i=1}^{R} w_{1,i}\, y(k-i+1) + b$$

**Backpropagation Algorithm:**
**Performance Index:**
Mean Square error: $F(\mathbf{x}) = E[\mathbf{e}^T\mathbf{e}] = E[(\mathbf{t}-\mathbf{a})^T(\mathbf{t}-\mathbf{a})]$
**Approximate Performance Index:** (single sample)
$$\hat{F}(\mathbf{x}) = \mathbf{e}^T(k)\mathbf{e}(k) = (\mathbf{t}(k) - \mathbf{a}(k))^T(\mathbf{t}(k) - \mathbf{a}(k))$$
**Sensitivity:** $\mathbf{s}^m = \frac{\partial \hat{F}}{\partial \mathbf{n}^m} = \begin{bmatrix} \frac{\partial \hat{F}}{\partial n_1^m} & \frac{\partial \hat{F}}{\partial n_2^m} & \cdots & \frac{\partial \hat{F}}{\partial n_{s^m}^m} \end{bmatrix}^T$

**Forward Propagation:** $\mathbf{a}^0 = \mathbf{p}$,
$$\mathbf{a}^{m+1} = \mathbf{f}^{m+1}(\mathbf{W}^{m+1}\mathbf{a}^m + \mathbf{b}^{m+1}) \ for \ m = 0,1,\dots,M-1$$
$$\mathbf{a} = \mathbf{a}^M$$
**Backward Propagation:** $\mathbf{s}^M = -2\dot{\mathbf{F}}^M(\mathbf{n}^M)(\mathbf{t}-\mathbf{a})$,
$$\mathbf{s}^m = \dot{\mathbf{F}}^m(\mathbf{n}^m)(\mathbf{W}^{m+1})^T \mathbf{s}^{m+1} \ for \ m = M-1, \dots, 2, 1, \text{where}$$
$$\dot{\mathbf{F}}^m(\mathbf{n}^m) = diag([\dot{f}^m(n_1^m) \ \dot{f}^m(n_2^m) \ \dots \ \dot{f}^m(n_{s^m}^m)])$$
$$\dot{f}^m(n_j^m) = \frac{\partial f^m(n_j^m)}{\partial n_j^m}$$
**Weight Update (Approximate Steepest Descent):**
$$\mathbf{W}^m(k+1) = \mathbf{W}^m(k) - \alpha\mathbf{s}^m(\mathbf{a}^{m-1})^T$$
$$\mathbf{b}^m(k+1) = \mathbf{b}^m(k) - \alpha\mathbf{s}^m$$

**\*Heuristic Variations of Backpropagation:**
**Batching:** The parameters are updated only after the entire training set has been presented. The gradients calculated for each training example are averaged together to produce a more accurate estimate of the gradient.(If the training set is complete, i.e., covers all possible input/output pairs, then the gradient estimate will be exact.)
**Backpropagation with Momentum (MOBP):**
$$\Delta\mathbf{W}^m(k) = \gamma\Delta\mathbf{W}^m(k-1) - (1-\gamma)\alpha\, \mathbf{s}^m(\mathbf{a}^{m-1})^T$$
$$\Delta\mathbf{b}^m(k) = \gamma\Delta\mathbf{b}^m(k-1) - (1-\gamma)\alpha\, \mathbf{s}^m$$
**Variable Learning Rate Backpropagation (VLBP)**
1. If the squared error (over the entire training set) increases by more than some set percentage $\zeta$ (typically one to five percent) after a weight update, then the weight update is discarded, the learning rate is multiplied by some factor $\rho < 1$, and the momentum coefficient $\gamma$ (if it is used) is set to zero.
2. If the squared error decreases after a weight update, then the weight update is accepted and the learning rate is multiplied by some factor $\eta > 1$. If $\gamma$ has been previously set to zero, it is reset to its original value.
3. If the squared error increases by less than $\zeta$, then the weight update is accepted but the learning rate and the momentum coefficient are unchanged.

**Association:** $\mathbf{a} = hardlim(\mathbf{W}^0\mathbf{p}^0 + \mathbf{Wp} + b)$
An association is a link between the inputs and outputs of a network so that when a stimulus A is presented to the network, it will output a response B.
**Associative Learning Rules:**
**Unsupervised Hebb Rule:**
$$\mathbf{W}(q) = \mathbf{W}(q-1) + \alpha\, \mathbf{a}(q)\mathbf{p}^T(q)$$
**Hebb with Decay:**
$$\mathbf{W}(q) = (1-\gamma)\mathbf{W}(q-1) + \alpha\, \mathbf{a}(q)\mathbf{p}^T(q)$$
**Instar:** $a = hardlim(\mathbf{Wp} + b), \quad a = hardlim(_1\mathbf{w}^T\mathbf{p} + b)$
The instar is activated for $_1\mathbf{w}^T\mathbf{p} = \|_1\mathbf{w}\|\|\mathbf{p}\|cos\theta \geq -b$
where $\theta$ is the angle between $\mathbf{p}$ and $_1\mathbf{w}$.
**Instar Rule:**
$$_i\mathbf{w}(q) = {}_i\mathbf{w}(q-1) + \alpha\, a_i(q)(\mathbf{p}(q) - {}_i\mathbf{w}(q-1))$$
$$_i\mathbf{w}(q) = (1-\alpha)\, {}_i\mathbf{w}(q-1) + \alpha\, \mathbf{p}(q), if \ (a_i(q) = 1)$$
**Kohonen Rule:**
$$_i\mathbf{w}(q) = {}_i\mathbf{w}(q-1) + \alpha(\mathbf{p}(q) - {}_i\mathbf{w}(q-1)) \ \ for \ i \in X(q)$$
**Outstar Rule:** $\mathbf{a} = satlins(\mathbf{Wp})$
$$\mathbf{w}_j(q) = \mathbf{w}_j(q-1) + \alpha\left(\mathbf{a}(q) - \mathbf{w}_j(q-1)\right)p_j(q)$$

**Competitive Layer:** $\mathbf{a} = compet(\mathbf{Wp}) = compet(\mathbf{n})$
**Competitive Learning with the Kohonen Rule:**
$$_{i^*}\mathbf{w}(q) = {}_{i^*}\mathbf{w}(q-1) + \alpha\left(\mathbf{p}(q) - {}_{i^*}\mathbf{w}(q-1)\right)$$
$$= (1-\alpha)\, {}_{i^*}\mathbf{w}(q-1) + \alpha\, \mathbf{p}(q)$$
$_i\mathbf{w}(q) = {}_i\mathbf{w}(q-1), \ i \neq i^*$ where $i^*$ is the winning neuron.
**Self-Organizing with the Kohonen Rule:**
$$_i\mathbf{w}(q) = {}_i\mathbf{w}(q-1) + \alpha\left(\mathbf{p}(q) - {}_i\mathbf{w}(q-1)\right)$$
$$= (1-\alpha)\, {}_i\mathbf{w}(q-1) + \alpha\, \mathbf{p}(q), \ i \in N_{i^*}(d)$$
$$N_i(d) = \{j, d_{i,j} \leq d\}$$
**LVQ Network:** $(w_{k,i}^2 = 1) \Rightarrow$ subclass $i$ is a part of class $k$
$$n_i^1 = -\|_i\mathbf{w}^1 - \mathbf{p}\|, \mathbf{a}^1 = compet(\mathbf{n}^1), \ \mathbf{a}^2 = \mathbf{W}^2\mathbf{a}^1$$
**LVQ Network Learning with the Kohonen Rule:**
$$_{i^*}\mathbf{w}^1(q) = {}_{i^*}\mathbf{w}^1(q-1) + \alpha\left(\mathbf{p}(q) - {}_{i^*}\mathbf{w}^1(q-1)\right),$$
$$if \ a_{k^*}^2 = t_{k^*} = 1$$
$$_{i^*}\mathbf{w}^1(q) = {}_{i^*}\mathbf{w}^1(q-1) - \alpha\left(\mathbf{p}(q) - {}_{i^*}\mathbf{w}^1(q-1)\right),$$
$$if \ a_{k^*}^2 = 1 \neq t_{k^*} = 0$$

---

$hardlim: a = \begin{cases} 0 & n < 0 \\ 1 & n \geq 0 \end{cases}$, $hardlims: a = \begin{cases} -1 & n < 0 \\ +1 & n \geq 0 \end{cases}$, $purelin: a = n$, $Logsig: a = \frac{1}{1+e^{-n}}$, $tansig: a = \frac{e^n - e^{-n}}{e^n + e^{-n}}$, $poslin: a = \begin{cases} 0 & n < 0 \\ n & n \geq 0 \end{cases}$,

$compet: a = \begin{cases} 1 & \text{neuron with max } n \\ 0 & \text{all other neurons} \end{cases}$, $satlin: a = \begin{cases} 0 & n < 0 \\ n & -1 \leq n \leq 1, \\ 1 & n > 1 \end{cases}$, $satlins: a = \begin{cases} -1 & n < 0 \\ n & -1 \leq n \leq 1, \\ 1 & n > 1 \end{cases}$

**\*\*HINT:** $diag([1\ 2\ 3]) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix}$

$Delay: a(t) = u(t-1), Integrator: a(t) = \int_0^t u(\tau)d\tau + a(0)$

机器学习概览



# MACHINE LEARNING IN EMOJI

SUPERVISED · UNSUPERVISED · REINFORCEMENT

SUPERVISED — human builds model based on input / output
UNSUPERVISED — human input, machine output, human utilizes if satisfactory
REINFORCEMENT — human input, machine output, human reward/punish, cycle continues

## BASIC REGRESSION

LINEAR — linear_model.LinearRegression()
Lots of numerical data

LOGISTIC — linear_model.LogisticRegression()
Target variable is categorical

## CLASSIFICATION

NEURAL NET — neural_network.MLPClassifier()
Complex relationships. Prone to overfitting
Basically magic.

K-NN — neighbors.KNeighborsClassifier()
Group membership based on proximity

DECISION TREE — tree.DecisionTreeClassifier()
If/then/else. Non-contiguous data
Can also be regression

RANDOM FOREST — ensemble.RandomForestClassifier()
Find best split randomly
Can also be regression

SVM — svm.SVC()  svm.LinearSVC()
Maximum margin classifier. Fundamental
Data Science algorithm

NAIVE BAYES — GaussianNB() MultinomialNB() BernoulliNB()
Updating knowledge step by step with new info

## CLUSTER ANALYSIS

K-MEANS — cluster.KMeans()
Similar datum into groups based on centroids

ANOMALY DETECTION — covariance.EllipticalEnvelope()
Finding outliers through grouping

## FEATURE REDUCTION

T-DISTRIB STOCHASTIC NEIB EMBEDDING — manifold.TSNE()
Visualize high dimensional data. Convert similarity to joint probabilities

PRINCIPLE COMPONENT ANALYSIS — decomposition.PCA()
Distill feature space into components that describe greatest variance

CANONICAL CORRELATION ANALYSIS — decomposition.CCA()
Making sense of cross-correlation matrices

LINEAR DISCRIMINANT ANALYSIS — lda.LDA()
Linear combination of features that separates classes

## OTHER IMPORTANT CONCEPTS

BIAS VARIANCE TRADEOFF

UNDERFITTING / OVERFITTING

INERTIA

ACCURACY FUNCTION — $(TP + TN) / (P + N)$

PRECISION FUNCTION — $TP / (TP + FP)$

SPECIFICITY FUNCTION — $TN / (FP + TN)$

SENSITIVITY FUNCTION — $TP / (TP + FN)$

@emilyinamillion made this

机器学习：Scikit-learn 算法



scikit-learn algorithm cheat-sheet

## Scikit-learn

Scikit-learn 是 Python 实现的一个免费机器学习算法库，基本涵盖了机器学习涉及到的各个方面，包括数据预处理、特征提取、模型构建、模型训练、模型验证以及模型评价等等。Scikit-learn 还可以和 Python 的 NumPy 和 SciPy 库协同工作。

机器学习之算法

这一来自 Microsoft Azure 的 cheatsheet 能够帮助你为你的预测分析选择合适的机器学习算法，根据不同的数据性质为你推荐最合适的算法。



用 Python 实现的数据科学

## Tensor Flow

TensorFlow™ 是一个采用数据流图（data flow graphs），用于数值计算的开源软件库。节点（Nodes）在图中表示数学操作，图中的线（edges）则表示在节点间相互联系的多维数据数组，即张量（tensor）。它灵活的架构让你可以在多种平台上展开计算，例如台式计算机中的一个或多个 CPU（或 GPU），服务器，移动设备等等。TensorFlow 最初由 Google 大脑小组（隶属于 Google 机器智能研究机构）的研究员和工程师们开发出来，用于机器学习和深度神经网络方面的研究，但这个系统的通用性使其也可广泛用于其他计算领域。

## About

### TensorFlow

TensorFlow™ is an open source software library for numerical computation using data flow graphs. TensorFlow was originally developed for the purposes of conducting machine learning and deep neural networks research, but the system is general enough to be applicable in a wide variety of other domains as well.

### Skflow

Scikit Flow provides a set of high level model classes that you can use to easily integrate with your existing Scikit-learn pipeline code. Scikit Flow is a simplified interface for TensorFlow, to get people started on predictive analytics and data mining. Scikit Flow has been merged into TensorFlow since version 0.8 and now called TensorFlow Learn.

### Keras

Keras is a minimalist, highly modular neural networks library, written in Python and capable of running on top of either TensorFlow or Theano

## Installation

### How to install new package in Python:
```
pip install <package-name>
```
Example: `pip install requests`

**How to install tensorflow?**
device = **cpu/gpu**
python_version = **cp27/cp34**
```
sudo pip install
https://storage.googleapis.com/
tensorflow/linux/$device/tensorflow-
0.8.0-$python_version-none-linux_x86
_64.whl
```

**How to install Skflow**
```
pip install sklearn
```
**How to install Keras**
```
pip install keras
```
update ~/.keras/keras.json - replace "theano" by "tensorflow"

## Helpers

### Python helper
### Important functions

`type(object)`
Get object type

`help(object)`
Get help for object (list of available methods, attributes, signatures and so on)

`dir(object)`
Get list of object attributes (fields, functions)

`str(object)`
Transform an object to string

`object?`
Shows documentations about the object

`globals()`
Return the dictionary containing the current scope's global variables.

`locals()`
Update and return a dictionary containing the current scope's local variables.

`id(object)`
Return the identity of an object. This is guaranteed to be unique among simultaneously existing objects.
```
import __builtin__
dir(__builtin__)
```
Other built-in functions

## TensorFlow

### Main classes
```
tf.Graph()
tf.Operation()
tf.Tensor()
tf.Session()
```
### Some useful functions
```
tf.get_default_session()
tf.get_default_graph()
tf.reset_default_graph()
ops.reset_default_graph()
tf.device("/cpu:0")
tf.name_scope(value)
tf.convert_to_tensor(value)
```
### TensorFlow Optimizers
```
GradientDescentOptimizer
AdadeltaOptimizer
AdagradOptimizer
MomentumOptimizer
AdamOptimizer
FtrlOptimizer
RMSPropOptimizer
```
### Reduction
```
reduce_sum
reduce_prod
reduce_min
reduce_max
reduce_mean
reduce_all
reduce_any
accumulate_n
```
### Activation functions
```
tf.nn?
relu
relu6
elu
softplus
softsign
dropout
bias_add
sigmoid
tanh
sigmoid_cross_entropy_with_logits
softmax
log_softmax
softmax_cross_entropy_with_logits
sparse_softmax_cross_entropy_with_logits
weighted_cross_entropy_with_logits
etc.
```

## Skflow

### Main classes
```
TensorFlowClassifier
TensorFlowRegressor
TensorFlowDNNClassifier
TensorFlowDNNRegressor
TensorFlowLinearClassifier
TensorFlowLinearRegressor
TensorFlowRNNClassifier
TensorFlowRNNRegressor
```

TensorFlowEstimator

**Each classifier and regressor have following fields**
n_classes=0 (Regressor), n_classes are expected to be input (Classifiers)
batch_size=32,
steps=200, // except
TensorFlowRNNClassifier - there is 50
optimizer='Adagrad',
learning_rate=0.1,

Keras

2017 年，Google 的 TensorFlow 团队决定在 TensorFlow 的核心库中支持 Keras。 Chollet 解释说，Keras 被认为是一个接口（interface）而不是一个端到端（end-to-end）的机器学习框架。 它提供了一个更高级，更直观的抽象集，使得配置神经网络变得更容易，无论后端科学计算库如何。

## Numpy

NumPy 的目标是 Python 的 CPython 参考实现，这是一个非优化的字节码解释器（non-optimizingbytecodeinterpreter）。 为这个版本的 Python 编写的数学算法通常运行速度要比其他版本的慢。NumPy 通过提供在数组上高效运行的多维数组和函数和运算符来解决缓慢问题，这需要重写一些代码，主要用到 NumPy 的内部循环。

## Pandas

"Pandas"这个名字来源于"面板数据（panel data）"一词，是一个多维结构化数据集的计量经济学术语。



## Data Wrangling

"Data wrangler"一词开始渗透到流行文化之中。 在 2017 年的电影"金刚：骷髅岛"中，演员马克·埃文·杰克逊（Marc Evan Jackson）饰演的角色就被介绍为"史蒂夫·伍德沃德（SteveWoodward），我们的 data

wrangler"。

# Data Wrangling
## with pandas
## Cheat Sheet
http://pandas.pydata.org

## Tidy Data – A foundation for wrangling in pandas

In a tidy data set:

Each **variable** is saved in its own **column**

Each **observation** is saved in its own **row**

Tidy data complements pandas's vectorized operations. pandas will automatically preserve observations as you manipulate variables. No other format works as intuitively with pandas.

M * A

## Syntax – Creating DataFrames

```
df = pd.DataFrame(
        {"a" : [4 ,5, 6],
         "b" : [7, 8, 9],
         "c" : [10, 11, 12]},
         index = [1, 2, 3])
```
Specify values for each column.

```
df = pd.DataFrame(
        [[4, 7, 10],
         [5, 8, 11],
         [6, 9, 12]],
         index=[1, 2, 3],
         columns=['a', 'b', 'c'])
```
Specify values for each row.

```
df = pd.DataFrame(
        {"a" : [4 ,5, 6],
         "b" : [7, 8, 9],
         "c" : [10, 11, 12]},
    index = pd.MultiIndex.from_tuples(
        [('d',1),('d',2),('e',2)],
            names=['n','v'])))
```
Create DataFrame with a MultiIndex

## Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.
```
df = (pd.melt(df)
        .rename(columns={
            'variable' : 'var',
            'value' : 'val'})
        .query('val >= 200')
    )
```

## Reshaping Data – Change the layout of a data set

**pd.melt(df)**
Gather columns into rows.

**df.pivot(columns='var', values='val')**
Spread rows into columns.

**pd.concat([df1,df2])**
Append rows of DataFrames

**pd.concat([df1,df2], axis=1)**
Append columns of DataFrames

**df.sort_values('mpg')**
Order rows by values of a column (low to high).

**df.sort_values('mpg',ascending=False)**
Order rows by values of a column (high to low).

**df.rename(columns = {'y':'year'})**
Rename the columns of a DataFrame

**df.sort_index()**
Sort the index of a DataFrame

**df.reset_index()**
Reset index of DataFrame to row numbers, moving index to columns.

**df.drop(['Length','Height'], axis=1)**
Drop columns from DataFrame

## Subset Observations (Rows)

**df[df.Length > 7]**
Extract rows that meet logical criteria.

**df.drop_duplicates()**
Remove duplicate rows (only considers columns).

**df.head(n)**
Select first n rows.

**df.tail(n)**
Select last n rows.

**df.sample(frac=0.5)**
Randomly select fraction of rows.

**df.sample(n=10)**
Randomly select n rows.

**df.iloc[10:20]**
Select rows by position.

**df.nlargest(n, 'value')**
Select and order top n entries.

**df.nsmallest(n, 'value')**
Select and order bottom n entries.

### Logic in Python (and pandas)

| | | | | | |
|---|---|---|---|---|---|
| < | Less than | != | | Not equal to | |
| > | Greater than | df.column.isin(values) | | Group membership | |
| == | Equals | pd.isnull(obj) | | Is NaN | |
| <= | Less than or equals | pd.notnull(obj) | | Is not NaN | |
| >= | Greater than or equals | &,|,~,^,df.any(),df.all() | | Logical and, or, not, xor, any, all | |

## Subset Variables (Columns)

**df[['width','length','species']]**
Select multiple columns with specific names.

**df['width']** or **df.width**
Select single column with specific name.

**df.filter(regex='regex')**
Select columns whose name matches regular expression regex.

### regex (Regular Expressions) Examples

| | |
|---|---|
| '\.' | Matches strings containing a period '.' |
| 'Length$' | Matches strings ending with word 'Length' |
| '^Sepal' | Matches strings beginning with the word 'Sepal' |
| '^x[1-5]$' | Matches strings beginning with 'x' and ending with 1,2,3,4,5 |
| '^(?!Species$).*' | Matches strings except the string 'Species' |

**df.loc[:, 'x2':'x4']**
Select all columns between x2 and x4 (inclusive).

**df.iloc[:,[1,2,5]]**
Select columns in positions 1, 2 and 5 (first column is 0).

**df.loc[df['a'] > 10, ['a','c']]**
Select rows meeting logical condition, and only the specific columns .

## Summarize Data

```
df['w'].value_counts()
    Count number of rows with each unique value of variable
len(df)
    # of rows in DataFrame.
df['w'].nunique()
    # of distinct values in a column.
df.describe()
    Basic descriptive statistics for each column (or GroupBy)
```

pandas provides a large set of **summary functions** that operate on different kinds of pandas objects (DataFrame columns, Series, GroupBy, Expanding and Rolling (see below)) and produce single values for each of the groups. When applied to a DataFrame, the result is returned as a pandas Series for each column. Examples:

```
sum()                          min()
    Sum values of each object.     Minimum value in each object.
count()                        max()
    Count non-NA/null values of     Maximum value in each object.
    each object.               mean()
median()                           Mean value of each object.
    Median value of each object.  var()
quantile([0.25,0.75])              Variance of each object.
    Quantiles of each object.  std()
apply(function)                    Standard deviation of each
    Apply function to each object.    object.
```

## Group Data

```
df.groupby(by="col")
    Return a GroupBy object,
    grouped by values in column
    named "col".

df.groupby(level="ind")
    Return a GroupBy object,
    grouped by values in index
    level named "ind".
```

All of the summary functions listed above can be applied to a group. Additional GroupBy functions:

```
size()                         agg(function)
    Size of each group.            Aggregate group using function.
```

## Windows

```
df.expanding()
    Return an Expanding object allowing summary functions to be
    applied cumulatively.
df.rolling(n)
    Return a Rolling object allowing summary functions to be
    applied to windows of length n.
```

## Handling Missing Data

```
df.dropna()
    Drop rows with any column having NA/null data.
df.fillna(value)
    Replace all NA/null data with value.
```

## Make New Columns

```
df.assign(Area=lambda df: df.Length*df.Height)
    Compute and append one or more new columns.
df['Volume'] = df.Length*df.Height*df.Depth
    Add single column.
pd.qcut(df.col, n, labels=False)
    Bin column into n buckets.
```

pandas provides a large set of **vector functions** that operate on all columns of a DataFrame or a single selected column (a pandas Series). These functions produce vectors of values for each of the columns, or a single Series for the individual Series. Examples:

```
max(axis=1)                    min(axis=1)
    Element-wise max.              Element-wise min.
clip(lower=-10,upper=10)       abs()
    Trim values at input thresholds  Absolute value.
```

The examples below can also be applied to groups. In this case, the function is applied on a per-group basis, and the returned vectors are of the length of the original DataFrame.

```
shift(1)                       shift(-1)
    Copy with values shifted by 1.   Copy with values lagged by 1.
rank(method='dense')           cumsum()
    Ranks with no gaps.            Cumulative sum.
rank(method='min')             cummax()
    Ranks. Ties get min rank.      Cumulative max.
rank(pct=True)                 cummin()
    Ranks rescaled to interval [0, 1]. Cumulative min.
rank(method='first')           cumprod()
    Ranks. Ties go to first value.  Cumulative product.
```

## Plotting

```
df.plot.hist()                 df.plot.scatter(x='w',y='h')
    Histogram for each column       Scatter chart using pairs of points
```

## Combine Data Sets

```
Standard Joins

pd.merge(adf, bdf,
    how='left', on='x1')
    Join matching rows from bdf to adf.

pd.merge(adf, bdf,
    how='right', on='x1')
    Join matching rows from adf to bdf.

pd.merge(adf, bdf,
    how='inner', on='x1')
    Join data. Retain only rows in both sets.

pd.merge(adf, bdf,
    how='outer', on='x1')
    Join data. Retain all values, all rows.

Filtering Joins

adf[adf.x1.isin(bdf.x1)]
    All rows in adf that have a match in bdf.

adf[~adf.x1.isin(bdf.x1)]
    All rows in adf that do not have a match in bdf.

Set-like Operations

pd.merge(ydf, zdf)
    Rows that appear in both ydf and zdf
    (Intersection).

pd.merge(ydf, zdf, how='outer')
    Rows that appear in either or both ydf and zdf
    (Union).

pd.merge(ydf, zdf, how='outer',
    indicator=True)
    .query('_merge == "left_only"')
    .drop(['_merge'],axis=1)
    Rows that appear in ydf but not zdf (Setdiff).
```

---

## Data Wrangling with dplyr and tidyr

# Data Wrangling
### with dplyr and tidyr
Cheat Sheet
RStudio

## Tidy Data - A foundation for wrangling in R

In a tidy data set:
Each **variable** is saved in its own **column**
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.

## 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: local data frame [150 x 5]

  Sepal.Length Sepal.Width Petal.Length
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 not shown: Petal.Width (dbl),
   Species (fctr)
```

**dplyr::glimpse(iris)**
Information dense summary of tbl data.

**utils::View(iris)**
View data set in spreadsheet-like display (note capital V).

**dplyr::%>%**
Passes object on left hand side as first argument (or . argument) of function on righthand side.

```
x %>% 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)
```

## Reshaping Data - Change the layout of a data set

**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.

**tidyr::gather(cases, "year", "n", 2:4)**
Gather columns into rows.

**tidyr::spread(pollution, size, amount)**
Spread rows into columns.

**tidyr::separate(storms, date, c("y", "m", "d"))**
Separate one column into several.

**tidyr::unite(data, col, ..., sep)**
Unite several columns into one.

## 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).

## 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.
```

## 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 |

## Summarise Data

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).

Summarise uses **summary functions**, functions that take a vector of values and return a single value, such as:

dplyr::first
First value of a vector.
dplyr::last
Last value of a vector.
dplyr::nth
Nth value of a vector.
dplyr::n
# of values in a vector.
dplyr::n_distinct
# of distinct values in a vector.
IQR
IQR of a vector.

min
Minimum value in a vector.
max
Maximum value in a vector.
mean
Mean value of a vector.
median
Median value of a vector.
var
Variance of a vector.
sd
Standard deviation of a 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::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::lead
Copy with values shifted by 1.
dplyr::lag
Copy with values lagged by 1.
dplyr::dense_rank
Ranks with no gaps.
dplyr::min_rank
Ranks. Ties get min rank.
dplyr::percent_rank
Ranks rescaled to [0, 1].
dplyr::row_number
Ranks. Ties got to first value.
dplyr::ntile
Bin vector into n buckets.
dplyr::between
Are values between a and b?
dplyr::cume_dist
Cumulative distribution.

dplyr::cumall
Cumulative all
dplyr::cumany
Cumulative any
dplyr::cummean
Cumulative mean
cumsum
Cumulative sum
cummax
Cumulative max
cummin
Cumulative min
cumprod
Cumulative prod
pmax
Element-wise max
pmin
Element-wise min

iris %>% group_by(Species) %>% mutate(...)
Compute new variables by group.

## Combine Data Sets

### Mutating Joins

dplyr::left_join(a, b, by = "x1")
Join matching rows from b to a.
dplyr::right_join(a, b, by = "x1")
Join matching rows from a to b.
dplyr::inner_join(a, b, by = "x1")
Join data. Retain only rows in both sets.
dplyr::full_join(a, b, by = "x1")
Join data. Retain all values, all rows.

### Filtering Joins

dplyr::semi_join(a, b, by = "x1")
All rows in a that have a match in b.
dplyr::anti_join(a, b, by = "x1")
All rows in a that do not have a match in b.

### Set Operations

dplyr::intersect(y, z)
Rows that appear in both y and z.
dplyr::union(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

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

## SciPy

SciPy 构建在 NumPy 数组对象（array object）上，是 NumPy 堆栈（stack）的一部分。NumPy 堆栈包括 Matplotlib，pandas 和 SymPy 等工具，以及一套扩展的科学计算库。这个 NumPy 堆栈与其他应用程序（如 MATLAB，GNU Octave 和 Scilab）具有相似的用户。NumPy 堆栈有时也被称为 SciPy 堆栈。

# Matplotlib

matplotlib 是 Python 编程语言的一个绘图库及其数值数学扩展 NumPy。它为利用通用的图形用户界面工具包，如 Tkinter, wxPython, Qt 或 GTK+向应用程序嵌入式绘图提供了面向对象的应用程序接口（API）。还有一个基于状态机（如开放图形库 OpenGL）的程序 pylab 接口，设计成与 MATLAB 非常类似--尽管使用起来有些不堪。SciPy 就利用了 matplotlib。

pyplot 是 matplotlib 的一个模块，它提供了一个类似 MATLAB 的接口。matplotlib 被设计得用起来像 MATLAB，具有使用 Python 的能力。免费是其优点。

数据可视化

# Data Visualization
## with ggplot2
### Cheat Sheet

RStudio

## Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same few components: a **data** set, a set of **geoms**—visual marks that represent data points, and a **coordinate system**.



To display data values, map variables in the data set to aesthetic properties of the geom like **size, color**, and **x** and **y** locations.



Build a graph with **qplot()** or **ggplot()**

```
qplot(x = cty, y = hwy, color = cyl, data = mpg, geom = "point")
```
Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

```
ggplot(data = mpg, aes(x = cty, y = hwy))
```
Begins a plot that you finish by adding layers to. No defaults, but provides more control than qplot().

```
ggplot(mpg, aes(hwy, cty)) +
  geom_point(aes(color = cyl)) +
  geom_smooth(method ="lm") +
  coord_cartesian() +
  scale_color_gradient() +
  theme_bw()
```

Add a new layer to a plot with a **geom_*()** or **stat_*()** function. Each provides a geom, a set of aesthetic mappings, and a default stat and position adjustment.

**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.

## Geoms
Use a geom to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

### One Variable

**Continuous**
a <- ggplot(mpg, aes(hwy))

a + **geom_area**(stat = "bin")
x, y, alpha, color, fill, linetype, size
b + geom_area(aes(y = ..density..), stat = "bin")

a + **geom_density**(kernel = "gaussian")
x, y, alpha, color, fill, linetype, size, weight
b + geom_density(aes(y = ..county..))

a + **geom_dotplot**()
x, y, alpha, color, fill

a + **geom_freqpoly**()
x, y, alpha, color, linetype, size
b + geom_freqpoly(aes(y = ..density..))

a + **geom_histogram**(binwidth = 5)
x, y, alpha, color, fill, linetype, size, weight
b + geom_histogram(aes(y = ..density..))

**Discrete**
b <- ggplot(mpg, aes(fl))

b + **geom_bar**()
x, alpha, color, fill, linetype, size, weight

### Graphical Primitives

c <- ggplot(map, aes(long, lat))

c + **geom_polygon**(aes(group = group))
x, y, alpha, color, fill, linetype, size

d <- ggplot(economics, aes(date, unemploy))

d + **geom_path**(lineend="butt",
linejoin="round', linemitre=1)
x, y, alpha, color, linetype, size

d + **geom_ribbon**(aes(ymin=unemploy - 900,
ymax=unemploy + 900))
x, ymax, ymin, alpha, color, fill, linetype, size

e <- ggplot(seals, aes(x = long, y = lat))

e + **geom_segment**(aes(
xend = long + delta_long,
yend = lat + delta_lat))
x, xend, y, yend, alpha, color, linetype, size

e + **geom_rect**(aes(xmin = long, ymin = lat,
xmax= long + delta_long,
ymax = lat + delta_lat))
xmax, xmin, ymax, ymin, alpha, color, fill,
linetype, size

### Two Variables

**Continuous X, Continuous Y**
f <- ggplot(mpg, aes(cty, hwy))

f + **geom_blank**()

f + **geom_jitter**()
x, y, alpha, color, fill, shape, size

f + **geom_point**()
x, y, alpha, color, fill, shape, size

f + **geom_quantile**()
x, y, alpha, color, linetype, size, weight

f + **geom_rug**(sides = "bl")
alpha, color, linetype, size

f + **geom_smooth**(model = lm)
x, y, alpha, color, fill, linetype, size, weight

f + **geom_text**(aes(label = cty))
x, y, label, alpha, angle, color, family, fontface,
hjust, lineheight, size, vjust

**Discrete X, Continuous Y**
g <- ggplot(mpg, aes(class, hwy))

g + **geom_bar**(stat = "identity")
x, y, alpha, color, fill, linetype, size, weight

g + **geom_boxplot**()
lower, middle, upper, x, ymax, ymin, alpha,
color, fill, linetype, shape, size, weight

g + **geom_dotplot**(binaxis = "y",
stackdir = "center")
x, y, alpha, color, fill

g + **geom_violin**(scale = "area")
x, y, alpha, color, fill, linetype, size, weight

**Discrete X, Discrete Y**
h <- ggplot(diamonds, aes(cut, color))

h + **geom_jitter**()
x, y, alpha, color, fill, shape, size

### Continuous Bivariate Distribution
i <- ggplot(movies, aes(year, rating))

i + **geom_bin2d**(binwidth = c(5, 0.5))
xmax, xmin, ymax, ymin, alpha, color, fill,
linetype, size, weight

i + **geom_density2d**()
x, y, alpha, colour, linetype, size

i + **geom_hex**()
x, y, alpha, colour, fill size

### Continuous Function
j <- ggplot(economics, aes(date, unemploy))

j + **geom_area**()
x, y, alpha, color, fill, linetype, size

j + **geom_line**()
x, y, alpha, color, linetype, size

j + **geom_step**(direction = "hv")
x, y, alpha, color, linetype, size

### Visualizing error
df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)
k <- ggplot(df, aes(grp, fit, ymin = fit-se, ymax = fit+se))

k + **geom_crossbar**(fatten = 2)
x, y, ymax, ymin, alpha, color, fill, linetype,
size

k + **geom_errorbar**()
x, ymax, ymin, alpha, color, linetype, size,
width (also **geom_errorbarh**())

k + **geom_linerange**()
x, ymin, ymax, alpha, color, linetype, size

k + **geom_pointrange**()
x, y, ymin, ymax, alpha, color, fill, linetype,
size

### Maps
data <- data.frame(murder = USArrests$Murder,
state = tolower(rownames(USArrests)))
map <- map_data("state")
l <- ggplot(data, aes(fill = murder))

l + **geom_map**(aes(map_id = state), map = map) +
**expand_limits**(x = map$long, y = map$lat)
map_id, alpha, color, fill, linetype, size

### Three Variables
seals$z <- with(seals, sqrt(delta_long^2 + delta_lat^2))
m <- ggplot(seals, aes(long, lat))

m + **geom_contour**(aes(z = z))
x, y, z, alpha, colour, linetype, size, weight

m + **geom_raster**(aes(fill = z), hjust=0.5,
vjust=0.5, interpolate=FALSE)
x, y, alpha, fill

m + **geom_tile**(aes(fill = z))
x, y, alpha, color, fill, linetype, size

RStudio® is a trademark of RStudio, Inc. • CC BY RStudio • info@rstudio.com • 844-448-1212 • rstudio.com

Learn more at docs.ggplot2.org • ggplot2 0.9.3.1 • Updated: 3/15

---

## Stats
An alternative way to build a layer

Some plots visualize a **transformation** of the original data set. Use a **stat** to choose a common transformation to visualize, e.g. a + **geom_bar**(stat = "bin")



Each stat creates additional variables to map aesthetics to. These variables use a common **..name..** syntax.

stat functions and geom functions both combine a stat with a geom to make a layer, i.e. stat_bin(geom="bar") does the same as geom_bar(stat="bin")

i + **stat_density2d**(aes(fill = ..level..),
geom = "polygon", n = 100)

**geom for layer** | **parameters for stat**

a + **stat_bin**(binwidth = 1, origin = 10)
x, y | ..count.., ..ncount.., ..density.., ..ndensity..

a + **stat_bindot**(binwidth = 1, binaxis = "x")
x, y | ..count.., ..ncount..

a + **stat_density**(adjust = 1, kernel = "gaussian")
x, y | ..count.., ..density.., ..scaled..

f + **stat_bin2d**(bins = 30, drop = TRUE)
x, y, fill | ..count.., ..density..

f + **stat_binhex**(bins = 30)
x, y, fill | ..count.., ..density..

f + **stat_density2d**(contour = TRUE, n = 100)
x, y, color, size | ..level..

m + **stat_contour**(aes(z = z))
x, y, z, order | ..level..

m + **stat_spoke**(aes(radius = z, angle = z))
angle, radius, x, xend, y, yend | ..x.., ..xend.., ..y.., ..yend..

m + **stat_summary_hex**(aes(z = z), bins = 30, fun = mean)
x, y, z, fill | ..value..

m + **stat_summary2d**(aes(z = z), bins = 30, fun = mean)
x, y, z, fill | ..value..

g + **stat_boxplot**(coef = 1.5)
x, y | ..lower.., ..middle.., ..upper.., ..outliers..

g + **stat_ydensity**(adjust = 1, kernel = "gaussian", scale = "area")
x, y | ..density.., ..scaled.., ..count.., ..n.., ..violinwidth.., ..width..

f + **stat_ecdf**(n = 40)
x, y | ..x.., ..y..

f + **stat_quantile**(quantiles = c(0.25, 0.5, 0.75), formula = y ~ log(x),
method = "rq")
x, y | ..quantile.., ..x.., ..y..

f + **stat_smooth**(method = "auto", formula = y ~ x, se = TRUE, n = 80,
fullrange = FALSE, level = 0.95)
x, y | ..se.., ..x.., ..y.., ..ymin.., ..ymax..

ggplot() + **stat_function**(aes(x = -3:3),
fun = dnorm, n = 101, args = list(sd=0.5))
x | ..y..

f + **stat_identity**()

ggplot() + **stat_qq**(aes(sample=1:100), distribution = qt,
dparams = list(df=5))
sample, x, y | ..x.., ..y..

f + **stat_sum**()
x, y, size | ..size..

f + **stat_summary**(fun.data = "mean_cl_boot")

f + **stat_unique**()

## Scales

**Scales** control how a plot maps data values to the visual values of an aesthetic. To change the mapping, add a custom scale.

n <- b + **geom_bar**(aes(fill = fl))

n + **scale_fill_manual**(
values = c("skyblue", "royalblue", "blue", "navy"),
limits = c("d", "e", "p", "r"), breaks =c("d", "e", "p", "r"),
name = "fuel", labels = c("D", "E", "P", "R"))

### General Purpose scales
Use with any aesthetic:
alpha, color, fill, linetype, shape, size

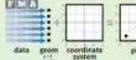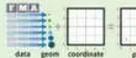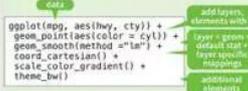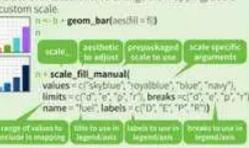**scale_*_continuous**() - map cont' values to visual values
**scale_*_discrete**() - map discrete values to visual values
**scale_*_identity**() - use data values as visual values
**scale_*_manual**(values = c()) - map discrete values to
manually chosen visual values

### X and Y location scales
Use with x or y aesthetics (x shown here)

**scale_x_date**(labels = date_format("%m/%d"),
breaks = date_breaks("2 weeks")) - treat x
values as dates. See ?strptime for label formats.

**scale_x_datetime**() - treat x values as date times. Use
same arguments as scale_x_date().

**scale_x_log10**() - Plot x on log10 scale

**scale_x_reverse**() - Reverse direction of x axis

**scale_x_sqrt**() - Plot x on square root scale

### Color and fill scales

**Discrete**

n <- b + **geom_bar**(
aes(fill = fl))

n + **scale_fill_brewer**(
palette = "Blues")
For palette choices:
library(RColorBrewer)
display.brewer.all()

n + **scale_fill_grey**(
start = 0.2, end = 0.8,
na.value = "red")

**Continuous**

o <- a + **geom_dotplot**(
aes(fill = ..x..))

o + **scale_fill_gradient**(
low = "red",
high = "yellow")

o + **scale_fill_gradient2**(
low = "red", high = "blue",
mid = "white", midpoint = 25)

o + **scale_fill_gradientn**(
colours = terrain.colors(6))
Also: rainbow(), heat.colors(),
topo.colors(), cm.colors(),
RColorBrewer::brewer.pal()

### Shape scales

p <- f + **geom_point**(
aes(shape = fl))

p + **scale_shape**(
solid = FALSE)

p + **scale_shape_manual**(
values = c(3:7))
Shape values shown in
chart on right

**Manual shape values**

### Size scales

q <- f + **geom_point**(
aes(size = cyl))

q + **scale_size_area**(max = 6)
Value mapped to area of circle
(not radius)

## Coordinate Systems

r <- b + **geom_bar**()

r + **coord_cartesian**(xlim = c(0, 5))
xlim, ylim
The default cartesian coordinate system

r + **coord_fixed**(ratio = 1/2)
ratio, xlim, ylim
Cartesian coordinates with fixed aspect
ratio between x and y units

r + **coord_flip**()
xlim, ylim
Flipped Cartesian coordinates

r + **coord_polar**(theta = "x", direction=1 )
theta, start, direction
Polar coordinates

r + **coord_trans**(ytrans = "sqrt")
xtrans, ytrans, limx, limy
Transformed cartesian coordinates. Set
extras and strains to the name
of a window function.

z + **coord_map**(projection = "ortho",
orientation=c(41, -74, 0))
projection, orientation, xlim, ylim
Map projections from the mapproj package
(mercator (default), azequalarea, 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))

s + **geom_bar**(position = "dodge")
Arrange elements side by side

s + **geom_bar**(position = "fill")
Stack elements on top of one another,
normalize height

s + **geom_bar**(position = "stack")
Stack elements on top of one another

f + **geom_point**(position = "jitter")
Add random noise to X and Y position
of each element to avoid overplotting

Each position adjustment can be recast as a function
with manual **width** and **height** arguments

s + **geom_bar**(position = position_dodge(width = 1))

## 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()

t + **facet_grid**(. ~ fl)
facet into columns based on fl

t + **facet_grid**(year ~ .)
facet into rows based on year

t + **facet_grid**(year ~ fl)
facet into both rows and columns

t + **facet_wrap**(~ fl)
wrap facets into a rectangular layout

Set **scales** to let axis limits vary across facets

t + **facet_grid**(y ~ x, 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 labels

t + facet_grid(. ~ fl, labeller = label_both)
t + facet_grid(. ~ fl, labeller = label_bquote(alpha ^ .(x)))
t + facet_grid(. ~ fl, labeller = label_parsed)

## Labels

t + **ggtitle**("New Plot Title")
Add a main title above the plot

t + **xlab**("New X label")
Change the label on the X axis

t + **ylab**("New Y label")
Change the label on the Y axis

t + **labs**(title = "New title", x = "New x", y = "New y")
All of the above

Use scale functions
to update legend
labels

## Legends

t + **theme**(legend.position = "bottom")
Place legend at "bottom", "top", "left", or "right"

t + **guides**(color = "none")
Set legend type for each aesthetic: colorbar, legend,
or none (no legend)

t + **scale_fill_discrete**(name = "Title",
labels = c("A", "B", "C"))
Set legend title and labels with a scale function.

## Themes

t + **theme_bw**()
White background
with grid lines

t + **theme_grey**()
Grey background
(default theme)

t + **theme_classic**()
White background
no gridlines

t + **theme_minimal**()
Minimal theme

ggthemes - Package with additional ggplot2 themes

## Zooming

**Without clipping** (preferred)

t + **coord_cartesian**(
xlim = c(0, 100), ylim = c(10, 20))

**With clipping** (removes unseen data points)

t + **xlim**(0, 100) + **ylim**(10, 20)

t + **scale_x_continuous**(limits = c(0, 100)) +
**scale_y_continuous**(limits = c(0, 100))

RStudio® is a trademark of RStudio, Inc. • CC BY RStudio • info@rstudio.com • 844-448-1212 • rstudio.com

Learn more at docs.ggplot2.org • ggplot2 0.9.3.1 • Updated: 3/15

## PySpark

### Python For Data Science Cheat Sheet
#### PySpark Basics
Learn Python for data science interactively at www.DataCamp.com

#### Spark
PySpark is the Spark Python API that exposes the Spark programming model to Python

#### Initializing Spark

**SparkContext**
```
>>> from pyspark import SparkContext
>>> sc = SparkContext(master = 'local[2]')
```

**Inspect SparkContext**
```
>>> sc.version                    Retrieve SparkContext version
>>> sc.pythonVer                  Retrieve Python version
>>> sc.master                     Master URL to connect to
>>> str(sc.sparkHome)             Path where Spark is installed on worker nodes
>>> str(sc.sparkUser())           Retrieve name of the Spark User running SparkContext
>>> sc.appName                    Return application name
>>> sc.applicationId              Retrieve application ID
>>> sc.defaultParallelism         Return default level of parallelism
>>> sc.defaultMinPartitions       Default minimum number of partitions for RDDs
```

**Configuration**
```
>>> from pyspark import SparkConf, SparkContext
>>> conf = (SparkConf()
         .setMaster("local")
         .setAppName("My app")
         .set("spark.executor.memory", "1g"))
>>> sc = SparkContext(conf = conf)
```

**Using The Shell**
In the PySpark shell, a special interpreter-aware SparkContext is already created in the variable called sc.
```
$ ./bin/spark-shell --master local[2]
$ ./bin/pyspark --master local[4] --py-files code.py
```
Set which master the context connects to with the --master argument, and add Python .zip, .egg or .py files to the runtime path by passing a comma-separated list to --py-files.

#### Loading Data

**Parallelized Collections**
```
>>> rdd = sc.parallelize([('a',7),('a',2),('b',2)])
>>> rdd2 = sc.parallelize([('a',2),('d',1),('b',1)])
>>> rdd3 = sc.parallelize(range(100))
>>> rdd4 = sc.parallelize([("a","b","c","d"),
                           ("b","p","r")])
```

**External Data**
Read either one text file from HDFS, a local file system or any Hadoop-supported file system URI with textFile(), or read in a directory of text files with wholeTextFiles().
```
>>> textFile = sc.textFile("/my/directory/*.txt")
>>> textFile2 = sc.wholeTextFiles("/my/directory/")
```

#### Retrieving RDD Information

**Basic Information**
```
>>> rdd.getNumPartitions()        List the number of partitions
>>> rdd.count()                   Count RDD instances
>>> rdd.countByKey()              Count RDD instances by key
>>> rdd.countByValue()            Count RDD instances by value
>>> rdd.collectAsMap()            Return (key,value) pairs as a dictionary
>>> rdd3.sum()                    Sum of RDD elements
>>> sc.parallelize([]).isEmpty()  Check whether RDD is empty
```

**Summary**
```
>>> rdd3.max()                    Maximum value of RDD elements
>>> rdd3.min()                    Minimum value of RDD elements
>>> rdd3.mean()                   Mean value of RDD elements
>>> rdd3.stdev()                  Standard deviation of RDD elements
>>> rdd3.variance()               Compute variance of RDD elements
>>> rdd3.histogram(3)             Compute histogram by bins
>>> rdd3.stats()                  Summary statistics (count, mean, stdev, max & min)
```

**Applying Functions**
```
>>> rdd.map(lambda x: x+(x[1],x[0])).collect()    Apply a function to each RDD element
>>> rdd5 = rdd.flatMap(lambda x: x+(x[1],x[0]))   Apply a function to each RDD element and flatten the result
>>> rdd5.collect()
>>> rdd4.flatMapValues(lambda x: x).collect()     Apply a flatMap function to each (key,value) pair of rdd4 without changing the keys
```

**Selecting Data**

Getting
```
>>> rdd.collect()                 Return a list with all RDD elements
>>> rdd.take(2)                   Take first 2 RDD elements
>>> rdd.first()                   Take first RDD element
>>> rdd.top(2)                    Take top 2 RDD elements
```

Sampling
```
>>> rdd3.sample(False, 0.15, 81).collect()    Return sampled subset of rdd3
```

Filtering
```
>>> rdd.filter(lambda x: "a" in x).collect()    Filter the RDD
>>> rdd5.distinct().collect()     Return distinct RDD values
>>> rdd.keys().collect()          Return (key,value) RDD's keys
```

**Iterating**
```
>>> def g(x): print(x)
>>> rdd.foreach(g)                Apply a function to all RDD elements
```

#### Reshaping Data

**Reducing**
```
>>> rdd.reduceByKey(lambda x,y : x+y).collect()    Merge the rdd values for each key
>>> rdd.reduce(lambda a, b: a + b)                 Merge the rdd values
```

**Grouping by**
```
>>> rdd3.groupBy(lambda x: x % 2).mapValues(list).collect()    Return RDD of grouped values
>>> rdd.groupByKey().mapValues(list).collect()                 Group rdd by key
```

**Aggregating**
```
>>> seqOp = (lambda x,y: (x[0]+y,x[1]+1))
>>> combOp = (lambda x,y:(x[0]+y[0],x[1]+y[1]))
>>> rdd3.aggregate((0,0),seqOp,combOp)             Aggregate RDD elements of each partition and then the results
>>> rdd.aggregateByKey((0,0),seqOp,combOp).collect()    Aggregate values of each RDD key
>>> rdd3.fold(0,add)                               Aggregate the elements of each partition, and then the results
>>> rdd.foldByKey(0, add).collect()                Merge the values for each key
>>> rdd3.keyBy(lambda x: x*x).collect()            Create tuples of RDD elements by applying a function
```

**Mathematical Operations**
```
>>> rdd.subtract(rdd2).collect()                   Return each rdd value not contained in rdd2
>>> rdd2.subtractByKey(rdd).collect()              Return each (key,value) pair of rdd2 with no matching key in rdd
>>> rdd.cartesian(rdd2).collect()                  Return the Cartesian product of rdd and rdd2
```

**Sort**
```
>>> rdd2.sortBy(lambda x: x[1]).collect()          Sort RDD by given function
>>> rdd2.sortByKey().collect()                     Sort (key, value) RDD by key
```

**Repartitioning**
```
>>> rdd.repartition(4)            New RDD with 4 partitions
>>> rdd.coalesce(1)               Decrease the number of partitions in the RDD to 1
```

**Saving**
```
>>> rdd.saveAsTextFile("rdd.txt")
>>> rdd.saveAsHadoopFile("hdfs://namenodehost/parent/child", 'org.apache.hadoop.mapred.TextOutputFormat')
```

**Stopping SparkContext**
```
>>> sc.stop()
```

**Execution**
```
$ ./bin/spark-submit examples/src/main/python/pi.py
```
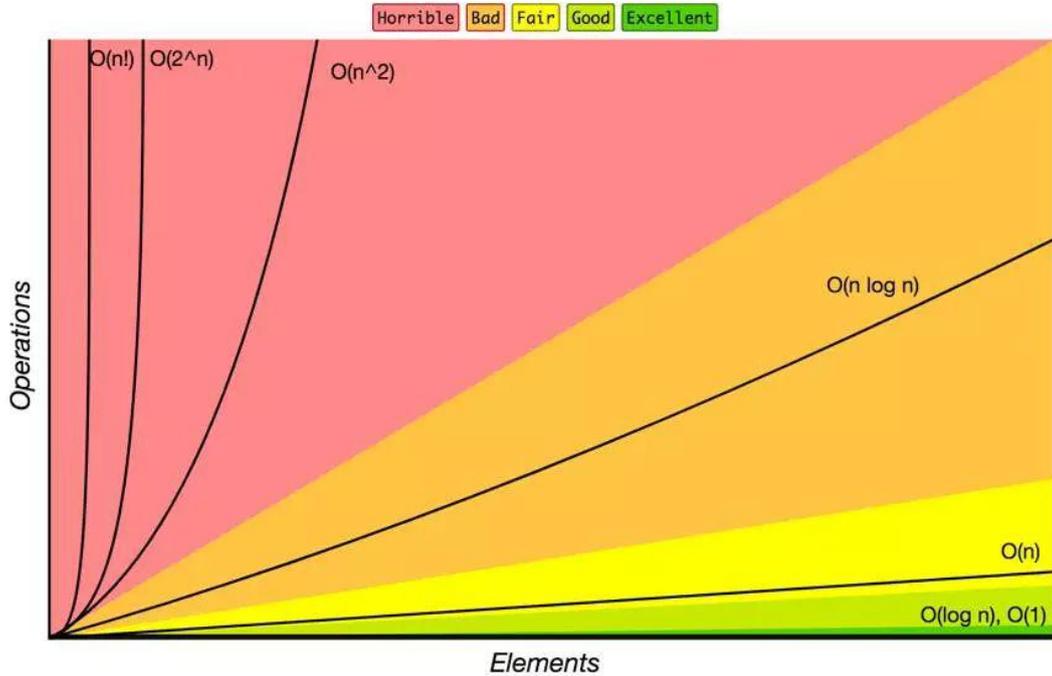
**DataCamp**
Learn Python for Data Science Interactively

## Big-O

# Big-O Complexity Chart

Horrible | Bad | Fair | Good | Excellent



Operations

O(n!) | O(2^n) | O(n^2)

O(n log n)

O(n)

O(log n), O(1)

Elements

## Common Data Structure Operations

| Data Structure | Time Complexity | | | | | | | | Space Complexity |
|---|---|---|---|---|---|---|---|---|---|
| | Average | | | | Worst | | | | Worst |
| | Access | Search | Insertion | Deletion | Access | Search | Insertion | Deletion | |
| Array | $\Theta(1)$ | $\Theta(n)$ | $\Theta(n)$ | $\Theta(n)$ | $O(1)$ | $O(n)$ | $O(n)$ | $O(n)$ | $O(n)$ |
| Stack | $\Theta(n)$ | $\Theta(n)$ | $\Theta(1)$ | $\Theta(1)$ | $O(n)$ | $O(n)$ | $O(1)$ | $O(1)$ | $O(n)$ |
| Queue | $\Theta(n)$ | $\Theta(n)$ | $\Theta(1)$ | $\Theta(1)$ | $O(n)$ | $O(n)$ | $O(1)$ | $O(1)$ | $O(n)$ |
| Singly-Linked List | $\Theta(n)$ | $\Theta(n)$ | $\Theta(1)$ | $\Theta(1)$ | $O(n)$ | $O(n)$ | $O(1)$ | $O(1)$ | $O(n)$ |
| Doubly-Linked List | $\Theta(n)$ | $\Theta(n)$ | $\Theta(1)$ | $\Theta(1)$ | $O(n)$ | $O(n)$ | $O(1)$ | $O(1)$ | $O(n)$ |
| Skip List | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $O(n)$ | $O(n)$ | $O(n)$ | $O(n)$ | $O(n \log(n))$ |
| Hash Table | N/A | $\Theta(1)$ | $\Theta(1)$ | $\Theta(1)$ | N/A | $O(n)$ | $O(n)$ | $O(n)$ | $O(n)$ |
| Binary Search Tree | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $O(n)$ | $O(n)$ | $O(n)$ | $O(n)$ | $O(n)$ |
| Cartesian Tree | N/A | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $\Theta(\log(n))$ | N/A | $O(n)$ | $O(n)$ | $O(n)$ | $O(n)$ |
| B-Tree | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $O(\log(n))$ | $O(\log(n))$ | $O(\log(n))$ | $O(\log(n))$ | $O(n)$ |
| Red-Black Tree | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $O(\log(n))$ | $O(\log(n))$ | $O(\log(n))$ | $O(\log(n))$ | $O(n)$ |
| Splay Tree | N/A | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $\Theta(\log(n))$ | N/A | $O(\log(n))$ | $O(\log(n))$ | $O(\log(n))$ | $O(n)$ |
| AVL Tree | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $O(\log(n))$ | $O(\log(n))$ | $O(\log(n))$ | $O(\log(n))$ | $O(n)$ |
| KD Tree | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $\Theta(\log(n))$ | $O(n)$ | $O(n)$ | $O(n)$ | $O(n)$ | $O(n)$ |

# Array Sorting Algorithms

| Algorithm | Time Complexity | | | Space Complexity |
|---|---|---|---|---|
| | Best | Average | Worst | Worst |
| Quicksort | Ω(n log(n)) | θ(n log(n)) | O(n^2) | O(log(n)) |
| Mergesort | Ω(n log(n)) | θ(n log(n)) | O(n log(n)) | O(n) |
| Timsort | Ω(n) | θ(n log(n)) | O(n log(n)) | O(n) |
| Heapsort | Ω(n log(n)) | θ(n log(n)) | O(n log(n)) | O(1) |
| Bubble Sort | Ω(n) | θ(n^2) | O(n^2) | O(1) |
| Insertion Sort | Ω(n) | θ(n^2) | O(n^2) | O(1) |
| Selection Sort | Ω(n^2) | θ(n^2) | O(n^2) | O(1) |
| Tree Sort | Ω(n log(n)) | θ(n log(n)) | O(n^2) | O(n) |
| Shell Sort | Ω(n log(n)) | θ(n(log(n))^2) | O(n(log(n))^2) | O(1) |
| Bucket Sort | Ω(n+k) | θ(n+k) | O(n^2) | O(n) |
| Radix Sort | Ω(nk) | θ(nk) | O(nk) | O(n+k) |
| Counting Sort | Ω(n+k) | θ(n+k) | O(n+k) | O(k) |
| Cubesort | Ω(n) | θ(n log(n)) | O(n log(n)) | O(n) |

参考资料：

Big-O Algorithm Cheat Sheet: http://bigocheatsheet.com/
Bokeh Cheat Sheet:https://s3.amazonaws.com/assets.datacamp.com/blog_assets/Python_Bokeh_Cheat_Sheet.pdf
Data Science Cheat Sheet: https://www.datacamp.com/community/tutorials/python-data-science-cheat-sheet-basics
Data Wrangling Cheat Sheet:https://www.rstudio.com/wp-content/uploads/2015/02/data-wrangling-cheatsheet.pdf
Data Wrangling: https://en.wikipedia.org/wiki/Data_wrangling
Ggplot Cheat Sheet: https://www.rstudio.com/wp-content/uploads/2015/03/ggplot2-cheatsheet.pdf
Keras Cheat Sheet: https://www.datacamp.com/community/blog/keras-cheat-sheet#gs.DRKeNMs
Keras: https://en.wikipedia.org/wiki/Keras
Machine Learning Cheat Sheet: https://ai.icymi.email/new-machinelearning-cheat-sheet-by-emily-barry-abdsc/
Machine Learning Cheat Sheet: https://docs.microsoft.com/en-in/azure/machine-learning/machine-learning-algorithm-cheat-sheet
ML Cheat Sheet:: http://peekaboo-vision.blogspot.com/2013/01/machine-learning-cheat-sheet-for-scikit.html
Matplotlib Cheat Sheet: https://www.datacamp.com/community/blog/py

thon-matplotlib-cheat-sheet#gs.uEKySpY

Matpotlib: https://en.wikipedia.org/wiki/Matplotlib

Neural Networks Cheat Sheet: http://www.asimovinstitute.org/neural-network-zoo/

Neural Networks Graph Cheat Sheet: http://www.asimovinstitute.org/blog/

Neural Networks: https://www.quora.com/Where-can-find-a-cheat-sheet-for-neural-network

Numpy Cheat Sheet: https://www.datacamp.com/community/blog/python-numpy-cheat-sheet#gs.AK5ZBgE

NumPy: https://en.wikipedia.org/wiki/NumPy

Pandas Cheat Sheet: https://www.datacamp.com/community/blog/python-pandas-cheat-sheet#gs.oundfxM

Pandas: https://en.wikipedia.org/wiki/Pandas_(software)

Pandas Cheat Sheet: https://www.datacamp.com/community/blog/pandas-cheat-sheet-python#gs.HPFoRIc

Pyspark Cheat Sheet: https://www.datacamp.com/community/blog/pyspark-cheat-sheet-python#gs.L=J1zxQ

Scikit Cheat Sheet: https://www.datacamp.com/community/blog/scikit-learn-cheat-sheet

Scikit-learn: https://en.wikipedia.org/wiki/Scikit-learn

Scikit-learn Cheat Sheet: http://peekaboo-vision.blogspot.com/2013/01/machine-learning-cheat-sheet-for-scikit.html

Scipy Cheat Sheet: https://www.datacamp.com/community/blog/python-scipy-cheat-sheet#gs.JDSg3OI

SciPy: https://en.wikipedia.org/wiki/SciPy

TesorFlow Cheat Sheet: https://www.altoros.com/tensorflow-cheat-sheet.html

Tensor Flow: https://en.wikipedia.org/wiki/TensorFlow