





# A mostly complete chart of Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

-  Backfed Input Cell
-  Input Cell
-  Noisy Input Cell
-  Hidden Cell
-  Probabilistic Hidden Cell
-  Spiking Hidden Cell
-  Output Cell
-  Match Input Output Cell
-  Recurrent Cell
-  Memory Cell
-  Different Memory Cell
-  Kernel
-  Convolution or Pool

Deep Feed Forward (DFF)



Perceptron (P)



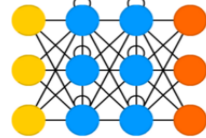
Feed Forward (FF)



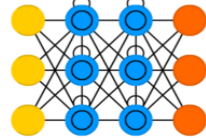
Radial Basis Network (RBF)



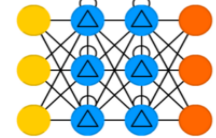
Recurrent Neural Network (RNN)



Long / Short Term Memory (LSTM)



Gated Recurrent Unit (GRU)



Auto Encoder (AE)



Variational AE (VAE)



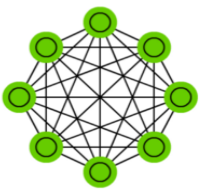
Denosing AE (DAE)



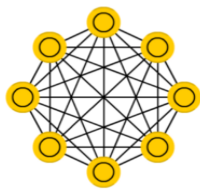
Sparse AE (SAE)



Markov Chain (MC)



Hopfield Network (HN)



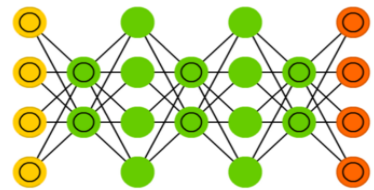
Boltzmann Machine (BM)



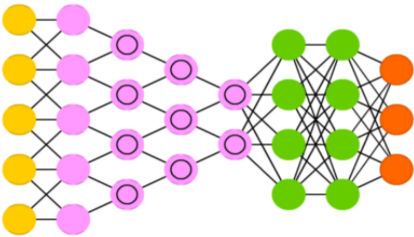
Restricted BM (RBM)



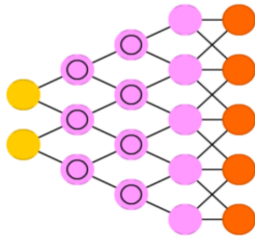
Deep Belief Network (DBN)



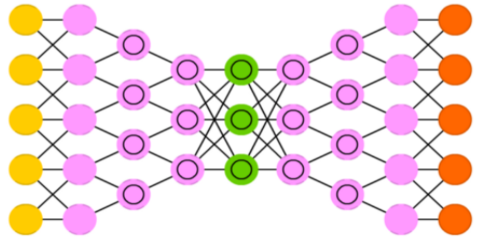
Deep Convolutional Network (DCN)



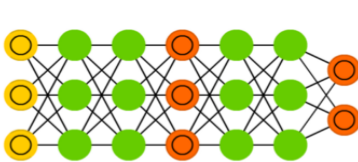
Deconvolutional Network (DN)



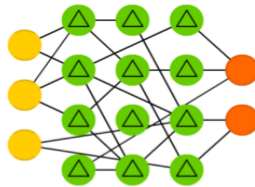
Deep Convolutional Inverse Graphics Network (DCIGN)



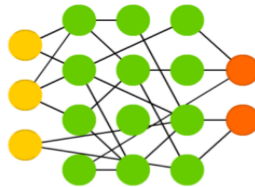
Generative Adversarial Network (GAN)



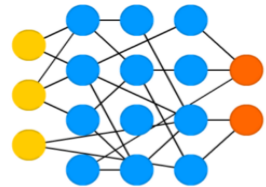
Liquid State Machine (LSM)



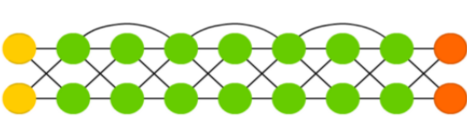
Extreme Learning Machine (ELM)



Echo State Network (ESN)



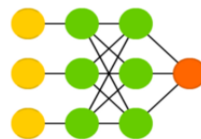
Deep Residual Network (DRN)



Kohonen Network (KN)



Support Vector Machine (SVM)

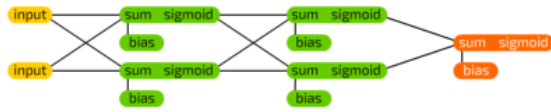


Neural Turing Machine (NTM)

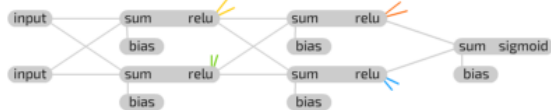
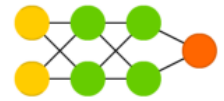


# An informative chart to build Neural Network Graphs

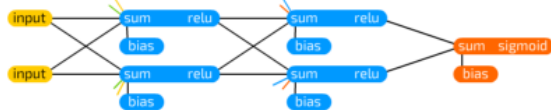
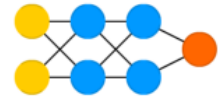
©2016 Fjodor van Veen - asimovinstitute.org



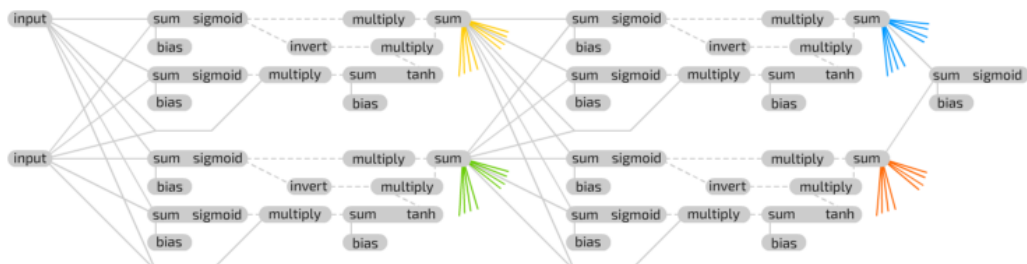
Deep Feed Forward Example



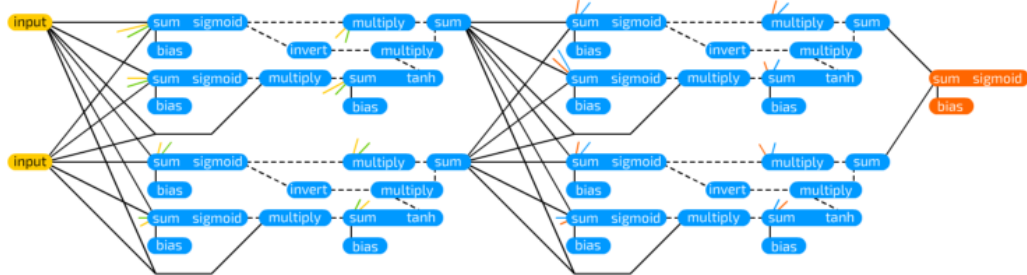
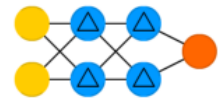
Deep Recurrent Example  
(previous iteration)



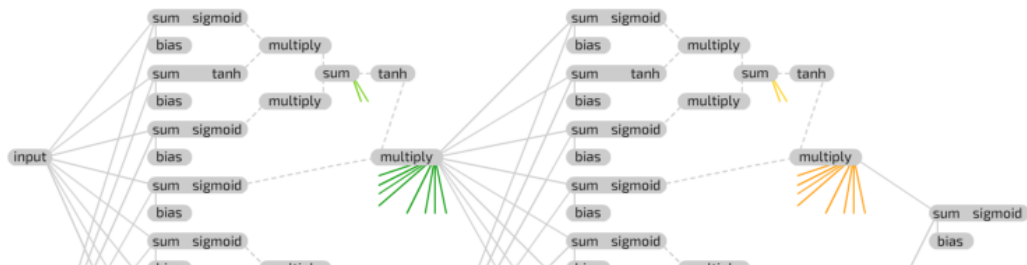
Deep Recurrent Example



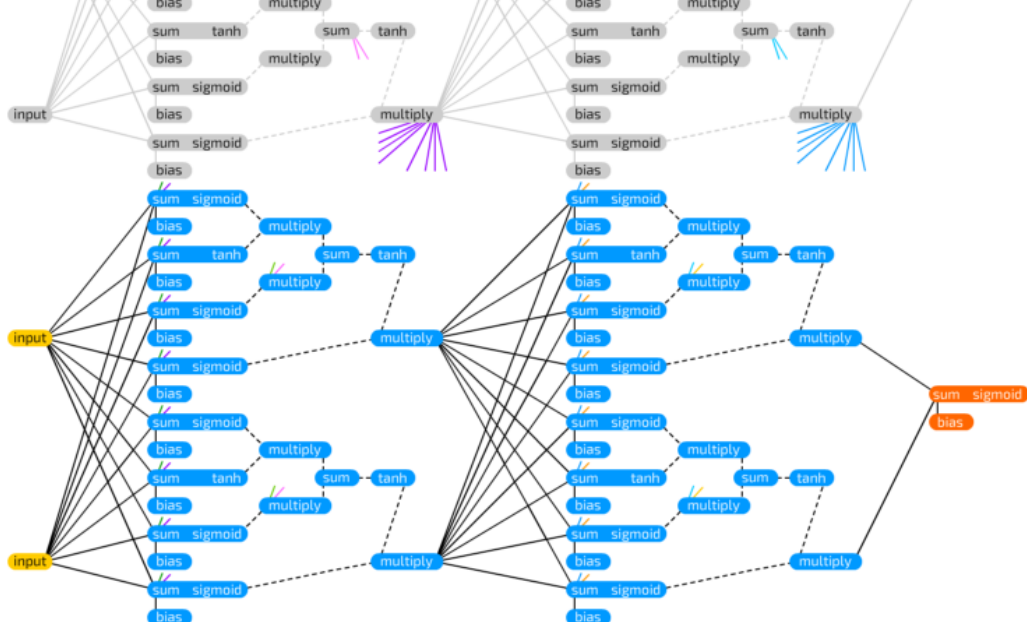
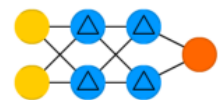
Deep GRU Example  
(previous iteration)



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  $\theta \in X$ , such that  $x+\theta = x$  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  $x+(-x) = \theta$ .
- 6) multiplication, for all scalars  $a \in F$ , and all vectors  $x \in X$ ,
- 7) For any  $x \in X$ ,  $1x = x$  (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  $\{x_1, x_2, \dots, x_n\}$ . If there exists  $n$  scalars  $a_1, a_2, \dots, a_n$ , at least one of which is nonzero, such that  $a_1x_1 + a_2x_2 + \dots + a_nx_n = \theta$ , then the  $\{x_i\}$  are linearly dependent.

**Spanning a Space:**

Let  $X$  be a linear vector space and let  $\{u_1, u_2, \dots, u_n\}$  be a subset of vectors in  $X$ . This subset spans  $X$  if and only if for every vector  $x \in X$  there exist scalars  $x_1, x_2, \dots, x_n$  such that  $x = x_1u_1 + x_2u_2 + \dots + x_nu_n$ .

**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  $x = \theta$ .
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  $y_1, y_2, \dots, y_n$ . From these vectors we will obtain  $n$  orthogonal vectors  $v_1, v_2, \dots, v_n$ .

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

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

**Vector Expansions:**

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

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  $B = [v_1 \ v_2 \ \dots \ v_n]$ ,

$R = [r_1 \ r_2 \ \dots \ r_n]$ ,  $R^T = B^{-1}$  In matrix form:  $x^v = B^{-1} x^s$

**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  $\{v_1, v_2, \dots, v_n\}$  be a basis for vector space  $X$ , and let  $\{u_1, u_2, \dots, u_n\}$  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:**  $B_t = [t_1 \ t_2 \ \dots \ t_n]$ ,  $B_w = [w_1 \ w_2 \ \dots \ w_n]$   
 $A' = [B_w^{-1} A B_t]$

**Eigenvalues & Eigenvectors:**  $Az = \lambda z$ ,  $|[A - \lambda I]| = 0$

**Diagonalization:**  $B = [z_1 \ z_2 \ \dots \ z_n]$ ,

where  $\{z_1, z_2, \dots, z_n\}$  are the eigenvectors of a square matrix  $A$ ,  
 $[B^{-1} A B] = \text{diag}([\lambda_1 \ \lambda_2 \ \dots \ \lambda_n])$

**Perceptron Architecture:**

$$a = \text{hardlim}(Wp + b), \quad W = [{}_1w^T \ {}_2w^T \ \dots \ {}_sw^T]^T, \\ a_i = \text{hardlim}(n_i) = \text{hardlim}({}_i w^T p + b_i)$$

**Decision Boundary:**  ${}_i w^T p + b_i = 0$

The decision boundary is always orthogonal to the weight vector. Single-layer perceptrons can only classify linearly separable vectors.

**Perceptron Learning Rule**

$$W^{new} = W^{old} + ep^T, \quad b^{new} = b^{old} + e, \\ \text{where } e = t - a$$

**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:**  $a = \text{purelin}(Wp)$

**The Hebb Rule:** Supervised Form:  $w_{ij}^{new} = w_{ij}^{old} + t_{qi} p_{qi}$

$$W = t_1 P_1^T + t_2 P_2^T + \dots + t_Q P_Q^T$$

$$W = [t_1 \ t_2 \ \dots \ t_Q] \begin{bmatrix} P_1^T \\ P_2^T \\ \vdots \\ P_Q^T \end{bmatrix} = TP^T$$

**Pseudoinverse Rule:**  $W = TP^+$

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

**Variations of Hebbian Learning:**

**Filtered Learning (Ch.14):**  $W^{new} = (1 - \gamma)W^{old} + \alpha t_q p_q^T$

**Delta Rule (Ch.10):**  $W^{new} = W^{old} + \alpha(t_q - a_q)p_q^T$

**Unsupervised Hebb (Ch.13):**  $W^{new} = W^{old} + \alpha a_q p_q^T$

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

**Grad**  $\nabla F(x) = \left[ \frac{\partial}{\partial x_1} F(x) \quad \frac{\partial}{\partial x_2} F(x) \quad \dots \quad \frac{\partial}{\partial x_n} F(x) \right]^T$

**Hessian:**  $\nabla^2 F(x) =$

$$\begin{bmatrix} \frac{\partial^2}{\partial x_1^2} F(x) & \frac{\partial^2}{\partial x_1 \partial x_2} F(x) & \dots & \frac{\partial^2}{\partial x_1 \partial x_n} F(x) \\ \frac{\partial^2}{\partial x_2 \partial x_1} F(x) & \frac{\partial^2}{\partial x_2^2} F(x) & \dots & \frac{\partial^2}{\partial x_2 \partial x_n} F(x) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial x_n \partial x_1} F(x) & \frac{\partial^2}{\partial x_n \partial x_2} F(x) & \dots & \frac{\partial^2}{\partial x_n^2} F(x) \end{bmatrix}$$

**Directional Derivatives:**

$$1^{st} \text{ Dir. Der.} = \frac{p^T \nabla F(x)}{\|p\|}, \quad 2^{nd} \text{ Dir. Der.} = \frac{p^T \nabla^2 F(x) p}{\|p\|^2}$$

**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:**

**1<sup>st</sup>-Order Condition:**  $\nabla F(x)|_{x=x^*} = 0$  (Stationary Points)

**2<sup>nd</sup>-Order Condition:**  $\nabla^2 F(x)|_{x=x^*} \geq 0$  (Positive Semi-definite Hessian Matrix).

**Quadratic fn.:**  $F(x) = \frac{1}{2} x^T A x + d^T x + c$

$$\nabla F(x) = Ax + d, \quad \nabla^2 F(x) = A, \quad \lambda_{min} \leq \frac{p^T A p}{\|p\|^2} \leq \lambda_{max}$$

<p><b>General Minimization Algorithm:</b>  <math>\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k</math> or <math>\Delta \mathbf{x}_k = (\mathbf{x}_{k+1} - \mathbf{x}_k) = \alpha_k \mathbf{p}_k</math></p> <p><b>Steepest Descent Algorithm:</b>  <math>\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{g}_k</math> where, <math>\mathbf{g}_k = \nabla F(\mathbf{x}) _{\mathbf{x}=\mathbf{x}_k}</math></p> <p><b>Stable Learning Rate:</b> (<math>\alpha_k = \alpha</math>, constant) <math>\alpha &lt; \frac{2}{\lambda_{max}}</math>  <math>\{\lambda_1, \lambda_2, \dots, \lambda_n\}</math> Eigenvalues of Hessian matrix A</p> <p><b>Learning Rate to Minimize Along the Line:</b>  <math>\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k \Rightarrow \alpha_k = -\frac{\mathbf{g}_k^T \mathbf{p}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}</math> (For quadratic fn.)</p> <p><b>After Minimization Along the Line:</b>  <math>\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k \Rightarrow \mathbf{g}_{k+1}^T \mathbf{p}_k = 0</math></p>	<p><b>*Heuristic Variations of Backpropagation:</b></p> <p><b>Batching:</b> 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.)</p> <p><b>Backpropagation with Momentum (MOBP):</b>  <math>\Delta \mathbf{W}^m(k) = \gamma \Delta \mathbf{W}^m(k-1) - (1-\gamma) \alpha \mathbf{s}^m (\mathbf{a}^{m-1})^T</math>  <math>\Delta \mathbf{b}^m(k) = \gamma \Delta \mathbf{b}^m(k-1) - (1-\gamma) \alpha \mathbf{s}^m</math></p> <p><b>Variable Learning Rate Backpropagation (VLBP)</b>  1. If the squared error (over the entire training set) increases by more than some set percentage <math>\zeta</math> (typically one to five percent) after a weight update, then the weight update is discarded, the learning rate is multiplied by some factor <math>\rho &lt; 1</math>, and the momentum coefficient <math>\gamma</math> (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 <math>\eta &gt; 1</math>. If <math>\gamma</math> has been previously set to zero, it is reset to its original value.  3. If the squared error increases by less than <math>\zeta</math>, then the weight update is accepted but the learning rate and the momentum coefficient are unchanged.</p>
<p><b>ADALINE: a = purelin(Wp + b)</b></p> <p><b>Mean Square Error:</b> (for ADALINE it is a quadratic fn.)  <math>F(\mathbf{x}) = E[e^2] = E[(t - a)^2] = E[(t - \mathbf{x}^T \mathbf{z})^2]</math>  <math>F(\mathbf{x}) = c - 2\mathbf{x}^T \mathbf{h} + \mathbf{x}^T \mathbf{R} \mathbf{x}</math>,  <math>c = E[t^2]</math>, <math>\mathbf{h} = E[t\mathbf{z}]</math> and <math>\mathbf{R} = E[\mathbf{z}\mathbf{z}^T] \Rightarrow \Lambda = 2\mathbf{R}</math>, <math>\mathbf{d} = -2\mathbf{h}</math>  Unique minimum, if it exists, is <math>\mathbf{x}^* = \mathbf{R}^{-1} \mathbf{h}</math>,  where <math>\mathbf{x} = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}</math> and <math>\mathbf{z} = \begin{bmatrix} \mathbf{p} \\ 1 \end{bmatrix}</math></p> <p><b>LMS Algorithm:</b> <math>\mathbf{W}(k+1) = \mathbf{W}(k) + 2\alpha \mathbf{e}(k) \mathbf{p}^T(k)</math>  <math>\mathbf{b}(k+1) = \mathbf{b}(k) + 2\alpha \mathbf{e}(k)</math></p> <p><b>Convergence Point:</b> <math>\mathbf{x}^* = \mathbf{R}^{-1} \mathbf{h}</math></p> <p><b>Stable Learning Rate:</b> <math>0 &lt; \alpha &lt; 1/\lambda_{max}</math> where <math>\lambda_{max}</math> is the maximum eigenvalue of R</p> <p><b>Adaptive Filter ADALINE:</b>  <math>a(k) = \text{purelin}(\mathbf{W}\mathbf{p}(k) + b) = \sum_{i=1}^R \mathbf{w}_{i,i} y(k-i+1) + b</math></p>	<p><b>Association: a = hardlim(W<sup>0</sup>p<sup>0</sup> + Wp + b)</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.</p> <p><b>Associative Learning Rules:</b></p> <p><b>Unsupervised Hebb Rule:</b>  <math>\mathbf{W}(q) = \mathbf{W}(q-1) + \alpha \mathbf{a}(q) \mathbf{p}^T(q)</math></p> <p><b>Hebb with Decay:</b>  <math>\mathbf{W}(q) = (1-\gamma) \mathbf{W}(q-1) + \alpha \mathbf{a}(q) \mathbf{p}^T(q)</math></p> <p><b>Instar: a = hardlim(Wp + b), a = hardlim(<sub>1</sub>w<sup>T</sup>p + b)</b>  The instar is activated for <math>{}_1\mathbf{w}^T \mathbf{p} = \ \mathbf{w}\  \ \mathbf{p}\  \cos \theta \geq -b</math>  where <math>\theta</math> is the angle between <math>\mathbf{p}</math> and <math>\mathbf{w}</math>.</p> <p><b>Instar Rule:</b>  <math>{}_i\mathbf{w}(q) = {}_i\mathbf{w}(q-1) + \alpha a_i(q) (\mathbf{p}(q) - {}_i\mathbf{w}(q-1))</math>  <math>{}_i\mathbf{w}(q) = (1-\alpha) {}_i\mathbf{w}(q-1) + \alpha \mathbf{p}(q)</math>, if <math>(a_i(q) = 1)</math></p> <p><b>Kohonen Rule:</b>  <math>{}_i\mathbf{w}(q) = {}_i\mathbf{w}(q-1) + \alpha (\mathbf{p}(q) - {}_i\mathbf{w}(q-1))</math> for <math>i \in X(q)</math></p> <p><b>Outstar Rule: a = satlins(Wp)</b>  <math>\mathbf{w}_j(q) = \mathbf{w}_j(q-1) + \alpha (\mathbf{a}(q) - \mathbf{w}_j(q-1)) p_j(q)</math></p>
<p><b>Backpropagation Algorithm:</b></p> <p><b>Performance Index:</b>  Mean Square error: <math>F(\mathbf{x}) = E[\mathbf{e}^T \mathbf{e}] = E[(\mathbf{t} - \mathbf{a})^T (\mathbf{t} - \mathbf{a})]</math></p> <p><b>Approximate Performance Index: (single sample)</b>  <math>\hat{F}(\mathbf{x}) = \mathbf{e}^T(k) \mathbf{e}(k) = (\mathbf{t}(k) - \mathbf{a}(k))^T (\mathbf{t}(k) - \mathbf{a}(k))</math></p> <p><b>Sensitivity: <math>\mathbf{s}^m = \frac{\partial \hat{F}}{\partial \mathbf{n}^m} = \begin{bmatrix} \frac{\partial \hat{F}}{\partial n_1^m} &amp; \frac{\partial \hat{F}}{\partial n_2^m} &amp; \dots &amp; \frac{\partial \hat{F}}{\partial n_{s^m}^m} \end{bmatrix}^T</math></b></p> <p><b>Forward Propagation: <math>\mathbf{a}^0 = \mathbf{p}</math>,</b>  <math>\mathbf{a}^{m+1} = \mathbf{f}^{m+1}(\mathbf{W}^{m+1} \mathbf{a}^m + \mathbf{b}^{m+1})</math> for <math>m = 0, 1, \dots, M-1</math>  <math>\mathbf{a} = \mathbf{a}^M</math></p> <p><b>Backward Propagation: <math>\mathbf{s}^M = -2\hat{\mathbf{F}}^M(\mathbf{n}^M)(\mathbf{t} - \mathbf{a})</math>,</b>  <math>\mathbf{s}^m = \hat{\mathbf{F}}^m(\mathbf{n}^m)(\mathbf{W}^{m+1})^T \mathbf{s}^{m+1}</math> for <math>m = M-1, \dots, 2, 1</math>, where  <math>\hat{\mathbf{F}}^m(\mathbf{n}^m) = \text{diag}([\hat{f}^m(n_1^m) \quad \hat{f}^m(n_2^m) \quad \dots \quad \hat{f}^m(n_{s^m}^m)])</math>  <math>\hat{f}^m(n_j^m) = \frac{\partial f^m(n_j^m)}{\partial n_j^m}</math></p> <p><b>Weight Update (Approximate Steepest Descent):</b>  <math>\mathbf{W}^m(k+1) = \mathbf{W}^m(k) - \alpha \mathbf{s}^m (\mathbf{a}^{m-1})^T</math>  <math>\mathbf{b}^m(k+1) = \mathbf{b}^m(k) - \alpha \mathbf{s}^m</math></p>	<p><b>Competitive Layer: a = compet(Wp) = compet(n)</b></p> <p><b>Competitive Learning with the Kohonen Rule:</b>  <math>{}_i\mathbf{w}(q) = {}_i\mathbf{w}(q-1) + \alpha (\mathbf{p}(q) - {}_i\mathbf{w}(q-1))</math>  <math>= (1-\alpha) {}_i\mathbf{w}(q-1) + \alpha \mathbf{p}(q)</math>  <math>{}_i\mathbf{w}(q) = {}_i\mathbf{w}(q-1)</math>, <math>i \neq i^*</math> where <math>i^*</math> is the winning neuron.</p> <p><b>Self-Organizing with the Kohonen Rule:</b>  <math>{}_i\mathbf{w}(q) = {}_i\mathbf{w}(q-1) + \alpha (\mathbf{p}(q) - {}_i\mathbf{w}(q-1))</math>  <math>= (1-\alpha) {}_i\mathbf{w}(q-1) + \alpha \mathbf{p}(q)</math>, <math>i \in N_{i^*}(d)</math>  <math>N_{i^*}(d) = \{j, d_{i,j} \leq d\}</math></p> <p><b>LVO Network: (<math>w_{k,i}^2 = 1</math>) <math>\Rightarrow</math> subclass <math>i</math> is a part of class <math>k</math></b>  <math>n_i^1 = -\ \mathbf{w}^1 - \mathbf{p}\ </math>, <math>\mathbf{a}^1 = \text{compet}(n^1)</math>, <math>\mathbf{a}^2 = \mathbf{W}^2 \mathbf{a}^1</math></p> <p><b>LVQ Network Learning with the Kohonen Rule:</b>  <math>{}_i\mathbf{w}^1(q) = {}_i\mathbf{w}^1(q-1) + \alpha (\mathbf{p}(q) - {}_i\mathbf{w}^1(q-1))</math>,  if <math>a_k^2 = t_k^* = 1</math>  <math>{}_i\mathbf{w}^1(q) = {}_i\mathbf{w}^1(q-1) - \alpha (\mathbf{p}(q) - {}_i\mathbf{w}^1(q-1))</math>,  if <math>a_k^2 = 1 \neq t_k^* = 0</math></p>
<p><b>hardlim: a = <math>\begin{cases} 0 &amp; n &lt; 0 \\ 1 &amp; n \geq 0 \end{cases}</math>, hardlims: a = <math>\begin{cases} -1 &amp; n &lt; 0 \\ +1 &amp; n \geq 0 \end{cases}</math>, purelin: a = n, Logsig: a = <math>\frac{1}{1+e^{-n}}</math>, tanstig: a = <math>\frac{e^n - e^{-n}}{e^n + e^{-n}}</math>, postlin: a = <math>\begin{cases} 0 &amp; n &lt; 0 \\ n &amp; n \geq 0 \end{cases}</math></b></p> <p><b>compet: a = <math>\begin{cases} 1 &amp; \text{neuron with max } n \\ 0 &amp; \text{all other neurons} \end{cases}</math>, satlin: a = <math>\begin{cases} 0 &amp; n &lt; 0 \\ n &amp; -1 \leq n \leq 1 \\ 1 &amp; n &gt; 1 \end{cases}</math>, satlins: a = <math>\begin{cases} -1 &amp; n &lt; 0 \\ -1 \leq n \leq 1 \\ 1 &amp; n &gt; 1 \end{cases}</math></b></p> <p><b>Delay: a(t) = u(t-1), Integrator: a(t) = <math>\int_0^t u(\tau) d\tau + a(0)</math></b></p>	<p><b>**HINT:</b></p> $\text{diag}([1 \ 2 \ 3]) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix}$

# MACHINE LEARNING IN EMOJI


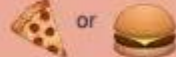
SUPERVISED

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 

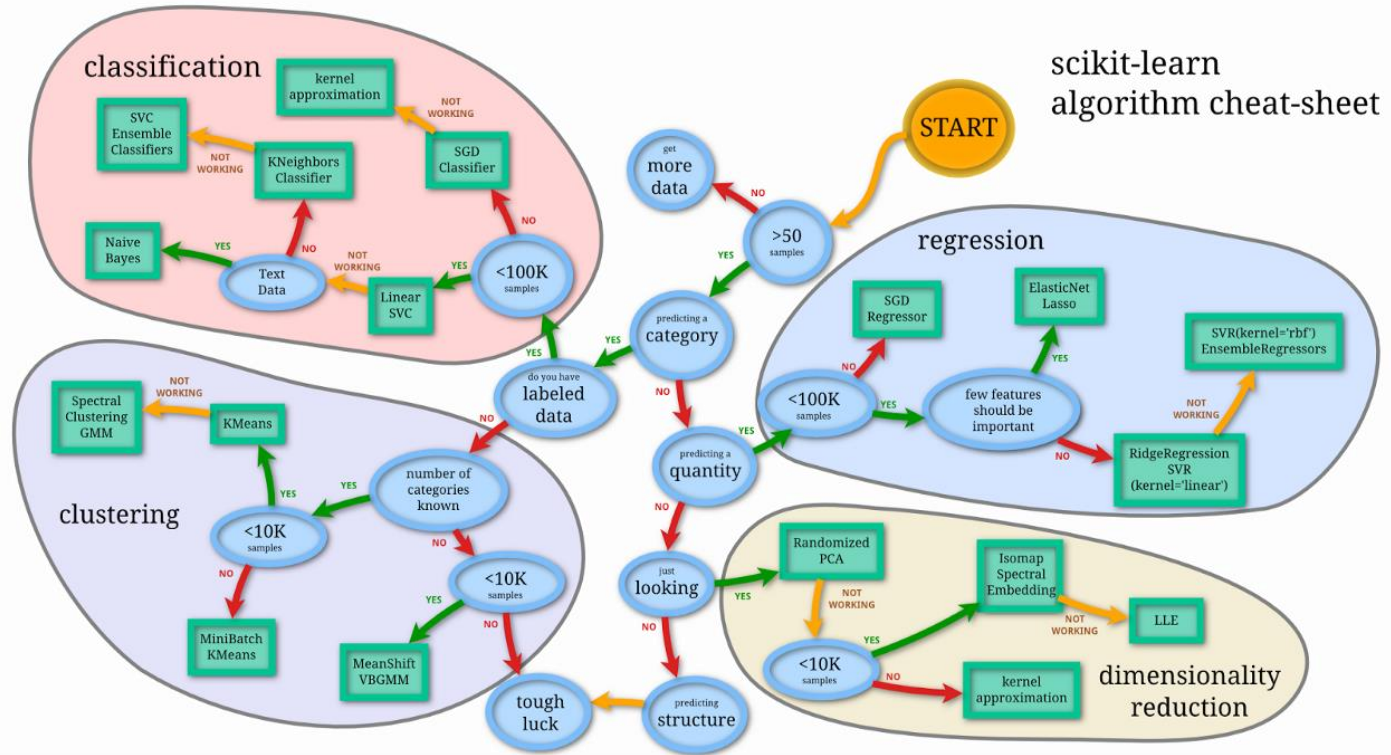
## FEATURE REDUCTION

- T-DISTRIBUTION 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**  $TP / (TP + FN)$  

@emilynamillion made this



## Python For Data Science Cheat Sheet

### Scikit-Learn

Learn Python for data science interactively at [www.DataCamp.com](http://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.cross_validation 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)
>>> X_test = scaler.transform(X_test)
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
>>> knn.fit(X_train, y_train)
>>> 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
>>> X = np.random.random((10,5))
>>> y = np.array(['M', 'M', 'F', 'F', 'M', 'F', 'M', 'F', 'F', 'F'])
>>> X[X < 0.7] = 0
```

#### Training And Test Data

```
>>> from sklearn.cross_validation import train_test_split
>>> X_train, X_test, y_train, y_test = train_test_split(X,
                                                    y,
                                                    Random_state=0)
```

#### 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)
```

##### Normalization

```
>>> from sklearn.preprocessing import Normalizer
>>> scaler = Normalizer().fit(X_train)
>>> normalized_X = scaler.transform(X_train)
>>> normalized_X_test = scaler.transform(X_test)
```

##### Binarization

```
>>> from sklearn.preprocessing import Binarizer
>>> binarizer = Binarizer(threshold=0.0).fit(X)
>>> binary_X = binarizer.transform(X)
```

##### 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)
```

#### 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.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)
```

Fit the model to the data

##### Unsupervised Learning

```
>>> k_means.fit(X_train)
>>> pca_model = pca.fit_transform(X_train)
```

Fit the model to the data  
Fit to data, then transform it

#### Prediction

```
Supervised Estimators
>>> y_pred = svc.predict(np.random.random((2,5)))
>>> y_pred = lr.predict(X_test)
>>> y_pred = knn.predict_probs(X_test)

Unsupervised Estimators
>>> y_pred = k_means.predict(X_test)
```

Predict labels  
Predict labels  
Estimate probability of a label

Predict labels in clustering algos

#### Evaluate Your Model's Performance

##### Classification Metrics

```
Accuracy Score
>>> knn.score(X_test, y_test)
>>> from sklearn.metrics import accuracy_score
>>> accuracy_score(y_test, y_pred)

Classification Report
>>> from sklearn.metrics import classification_report
>>> print(classification_report(y_test, y_pred))

Confusion Matrix
>>> from sklearn.metrics import confusion_matrix
>>> print(confusion_matrix(y_test, y_pred))
```

Estimator score method  
Metric scoring functions  
Precision, recall, f1-score and support

##### 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_true, y_pred)

R2 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

```
>>> from sklearn.grid_search import GridSearchCV
>>> params = {"n_neighbors": np.arange(1,3),
            "metric": ["euclidean", "cityblock"]}
>>> grid = GridSearchCV(estimator=knn,
                    param_grid=params)
>>> grid.fit(X_train, y_train)
>>> print(grid.best_score_)
>>> print(grid.best_estimator_.n_neighbors)
```

##### Randomized Parameter Optimization

```
>>> from sklearn.grid_search import RandomizedSearchCV
>>> params = {"n_neighbors": range(1,5),
            "weights": ["uniform", "distance"]}
>>> rsearch = RandomizedSearchCV(estimator=knn,
                             param_distributions=params,
                             cv=4,
                             n_iter=8,
                             random_state=5)
>>> rsearch.fit(X_train, y_train)
>>> print(rsearch.best_score_)
```

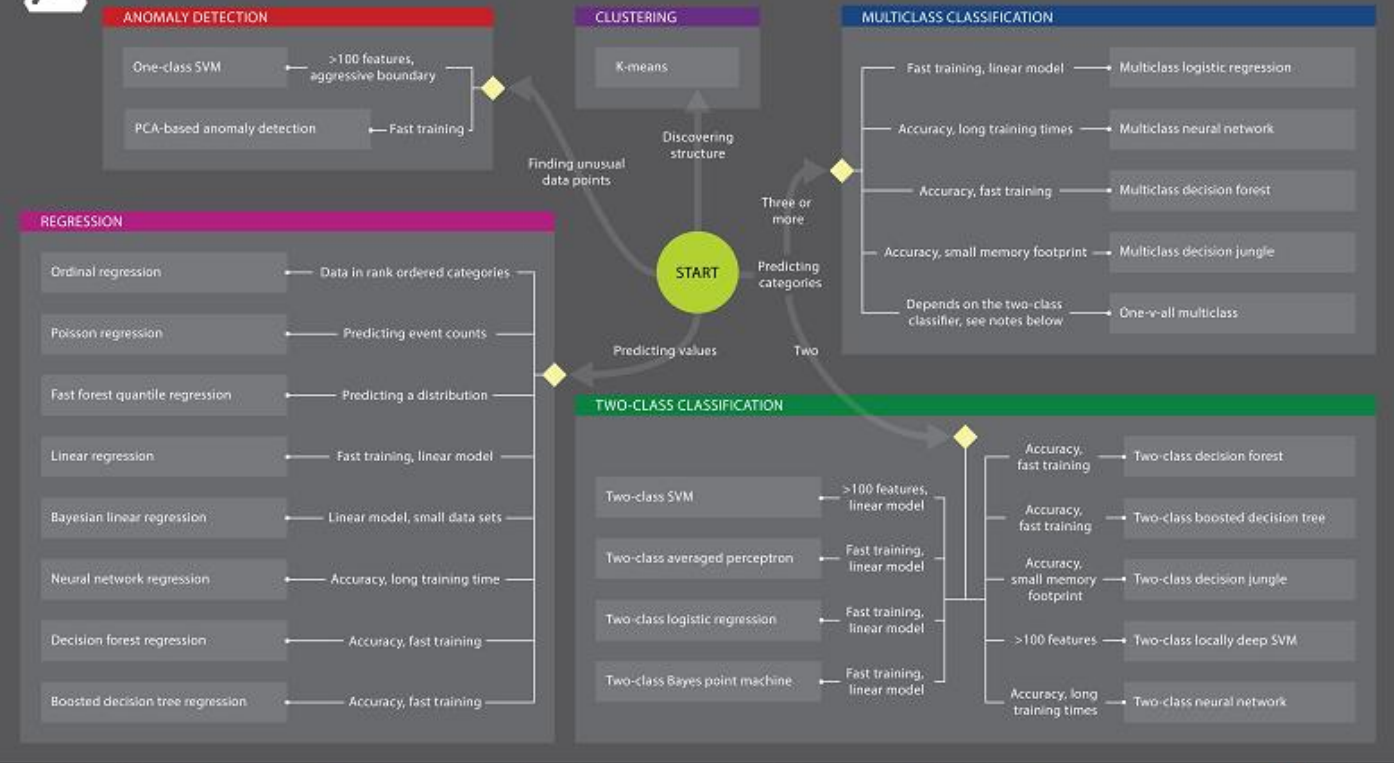
DataCamp

Learn Python for Data Science Interactively



# 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.



© 2015 Microsoft Corporation. All rights reserved. Created by the Azure Machine Learning Team. Email: AzurePoster@microsoft.com Download this poster: <http://aka.ms/MLCheatSheet> Microsoft

## Python For Data Science Cheat Sheet

### Python Basics

Learn More Python for Data Science Interactively at [www.datacamp.com](http://www.datacamp.com)

#### Variables and Data Types

##### Variable Assignment

```
>>> x=5
>>> x
5
```

##### Calculations With Variables

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

##### 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
True
```

#### Lists

Also see NumPy Arrays

```
>>> 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 0

<b>Subset</b>	Select item at index 1 Select 3rd last item
>>> my_list[1]	
>>> my_list[-3]	
<b>Slice</b>	Select items at index 1 and 2 Select items after index 0 Select items before index 3 Copy my_list
>>> my_list[1:3]	
>>> my_list[1:]	
>>> my_list[:3]	
>>> my_list[:]	
<b>Subset Lists of Lists</b>	my_list[listOfItemOfList]
>>> my_list2[1][0]	
>>> my_list2[1][:2]	

##### 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
True
```

##### List Methods

>>> my_list.index(a)	Get the index of an item
>>> my_list.count(a)	Count an item
>>> my_list.append('!')	Append an item at a time
>>> my_list.remove('!')	Remove an item
>>> del(my_list[0:1])	Remove an item
>>> my_list.reverse()	Reverse the list
>>> my_list.extend('!')	Append an item
>>> my_list.pop(-1)	Remove an item
>>> my_list.insert(0, '!')	Insert an item
>>> my_list.sort()	Sort the list

##### String Operations

Index starts at 0

```
>>> my_string[3]
'i'
>>> my_string[4:9]
'ngIsAwesome'
```

##### String Methods

>>> my_string.upper()	String to uppercase
>>> my_string.lower()	String to lowercase
>>> my_string.count('w')	Count String elements
>>> my_string.replace('e', 'i')	Replace String elements
>>> my_string.strip()	Strip whitespace from ends

#### Libraries

##### Import libraries

```
>>> import numpy
>>> import numpy as np
Selective import
>>> from math import pi
```

pandas  
Data analysis

Machine learning

NumPy  
Scientific computing

matplotlib  
2D plotting

#### Install Python

ANACONDA  
Leading open data science platform powered by Python

spyder  
Free IDE that is included with Anaconda

jupyter  
Create and share documents with live code, visualizations, text, ...

#### NumPy Arrays

Also see Lists

```
>>> my_list = [1, 2, 3, 4]
>>> my_array = np.array(my_list)
>>> my_2darray = np.array([[1,2,3], [4,5,6]])
```

##### Selecting NumPy Array Elements

Index starts at 0

<b>Subset</b>	Select item at index 1
>>> my_array[1]	2
<b>Slice</b>	Select items at index 0 and 1
>>> my_array[0:2]	array([1, 2])
<b>Subset 2D NumPy arrays</b>	my_2darray[rows, columns]
>>> my_2darray[:,0]	array([1, 4])

##### NumPy Array Operations

```
>>> my_array > 3
array([False, False, False,  True], dtype=bool)
>>> my_array * 2
array([2, 4, 6, 8])
>>> my_array + np.array([5, 6, 7, 8])
array([16, 8, 10, 12])
```

##### NumPy Array Functions

>>> my_array.shape	Get the dimensions of the array
>>> np.append(other_array)	Append items to an array
>>> np.insert(my_array, 1, 5)	Insert items in an array
>>> np.delete(my_array, [1])	Delete items in an array
>>> np.mean(my_array)	Mean of the array
>>> np.median(my_array)	Median of the array
>>> my_array.corrcoef()	Correlation coefficient
>>> np.std(my_array)	Standard deviation

DataCamp  
Learn Python for Data Science Interactively

Bokeh

Learn Bokeh Interactively at [www.DataCamp.com](http://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

```
>>> from bokeh.plotting import figure
>>> from bokeh.io import output_file, show
>>> x = [1, 2, 3, 4, 5]
>>> y = [6, 7, 2, 4, 5]
>>> p = figure(title="simple line example",
>>>            x_axis_label='x',
>>>            y_axis_label='y')
>>> p.line(x, y, legend="Temp.", line_width=2)
>>> output_file("lines.html")
>>> show(p)
```

1 Data

Also see Lists, NumPy & Pandas

Under the hood, your data is converted to Column Data Sources. You can also do this manually:

```
>>> import numpy as np
>>> import pandas as pd
>>> df = pd.DataFrame(np.array([[33.9, 4, 65, 'US'],
>>>                             [32.4, 4, 66, 'Asia'],
>>>                             [21.4, 4, 109, 'Europe']])),
>>>                   columns=['mpg', 'cyl', 'hp', 'origin'],
>>>                   index=['Toyota', 'Fiat', 'Volvo'])
>>> from bokeh.models import ColumnDataSource
>>> cds_df = ColumnDataSource(df)
```

2 Plotting

```
>>> from bokeh.plotting import figure
>>> p1 = figure(plot_width=300, tools='pan,box_zoom')
>>> p2 = figure(plot_width=300, plot_height=300,
>>>            x_range=(0, 8), y_range=(0, 8))
>>> p3 = figure()
```

3 Renderers & Visual Customizations

Glyphs

```
Scatter Markers
>>> p1.circle(np.array([1,2,3]), np.array([3,2,1]),
>>>           fill_color="white")
>>> p2.square(np.array([1.5,3.5,5.5]), [1,4,3],
>>>           color="blue", size=1)

Line Glyphs
>>> p1.line([1,2,3,4], [3,4,5,6], line_width=2)
>>> p2.multi_line(pd.DataFrame([[1,2,3],[5,6,7]]),
>>>               pd.DataFrame([[3,4,5],[3,2,1]]),
>>>               color="blue")
```

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])
```

Legends

Legend Location

```
Inside Plot Area
>>> p.legend.location = 'bottom_left'

Outside Plot Area
>>> r1 = p2.asterisk(np.array([1,2,3]), np.array([3,2,1])
>>> r2 = p2.line([1,2,3,4], [3,4,5,6])
>>> legend = Legend(items=[("One", [p1, r1]), ("Two", [r2])], location=(0, -30))
>>> p.add_layout(legend, 'right')
```

Customized Glyphs

Also see Data

```
Selection and Non-Selection Glyphs
>>> p.circle('mpg', 'cyl', source=cds_df,
>>>         selection_color='red',
>>>         nonselection_alpha=0.1)

Hover Glyphs
>>> hover = HoverTool(tooltips=None, mode='vline')
>>> p.add_tools(hover)

Colormapping
>>> color_mapper = CategoricalColorMapper(
>>>     factors=['Europe', 'Asia', 'US'],
>>>     palette=['red', 'green', 'blue'])
>>> p.circle('mpg', 'cyl', source=cds_df,
>>>         color=dict(field='origin',
>>>                   transform=color_mapper),
>>>         legend='Origin')
```

Linked Plots

Linked Axes

```
>>> p2.x_range = p1.x_range
>>> p2.y_range = p1.y_range
```

Linked Brushing

```
>>> p4 = figure(plot_width = 100, tools='box_select,lasso_select')
>>> p4.circle('mpg', 'cyl', source=cds_df)
>>> p5 = figure(plot_width = 200, tools='box_select,lasso_select')
>>> p5.circle('mpg', 'hp', source=cds_df)
>>> layout = row(p4,p5)
```

Also see Data

4 Output

Output to HTML File

```
>>> from bokeh.io import output_file, show
>>> output_file('my_bar_chart.html', mode='cdn')
```

Notebook Output

```
>>> from bokeh.io import output_notebook, show
>>> output_notebook()
```

Embedding

Standalone HTML

```
>>> from bokeh.embed import file_html
>>> html = file_html(p, CDN, "my_plot")
```

Components

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

5 Show or Save Your Plots

```
>>> show(p1)
>>> show(layout)

>>> save(p1)
>>> save(layout)
```

Statistical Charts With Bokeh

Also see Data

Bokeh's high-level `bokeh.charts` interface is ideal for quickly creating statistical charts

Bar Chart

```
>>> from bokeh.charts import Bar
>>> p = Bar(df, stacked=True, palette=['red','blue'])
```

Box Plot

```
>>> from bokeh.charts import BoxPlot
>>> p = BoxPlot(df, values='vals', label='cyl',
>>>             legend='bottom_right')
```

Histogram

```
>>> from bokeh.charts import Histogram
>>> p = Histogram(df, title='Histogram')
```

Scatter Plot

```
>>> from bokeh.charts import Scatter
>>> p = Scatter(df, x='mpg', y='hp', marker='square',
>>>            xlabel='Miles Per Gallon',
>>>            ylabel='Horsepower')
```

DataCamp

Learn Python For Data Science Interactively





## 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,
```

# Python For Data Science Cheat Sheet

## Keras

Learn Python for data science interactively at [www.DataCamp.com](http://www.DataCamp.com)



### 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 = Sequential()
>>> 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.layers 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'))
```

### 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=['mse'])
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(lr=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])
```

## DataCamp

Learn Python for Data Science Interactively



# Python For Data Science Cheat Sheet

## Pandas Basics

Learn Python for Data Science interactively at [www.DataCamp.com](http://www.DataCamp.com)



### Pandas

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



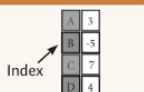
Use the following import convention:

```
>>> import pandas as pd
```

### Pandas Data Structures

#### Series

A one-dimensional labeled array capable of holding any data type



```
>>> s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])
```

#### DataFrame

A two-dimensional labeled data structure with columns of potentially different types

```
Columns
Country Capital Population
1 India New Delhi 1303171035
2 Brazil Brasilia 207847528

Index
0 Belgium Brussels 11190846
1 India New Delhi 1303171035
2 Brazil Brasilia 207847528

>>> data = {'Country': ['Belgium', 'India', 'Brazil'],
'Capital': ['Brussels', 'New Delhi', 'Brasilia'],
'Population': [11190846, 1303171035, 207847528]}
>>> df = pd.DataFrame(data,
columns=['Country', 'Capital', 'Population'])
```

### I/O

#### Read and Write to CSV

```
>>> pd.read_csv('file.csv', header=None, nrows=5)
>>> pd.to_csv('myDataFrame.csv')
```

#### Read and Write to Excel

```
>>> pd.read_excel('file.xlsx')
>>> pd.to_excel('dir/myDataFrame.xlsx', sheet_name='Sheet1')
Read multiple sheets from the same file
xlsx = pd.ExcelFile('file.xls')
>>> df = pd.read_excel(xlsx, 'Sheet1')
```

### Asking For Help

```
>>> help(pd.Series.loc)
```

### Selection

Also see NumPy Arrays

#### Getting

```
>>> s['b']
-5
Get one element
>>> df[1]
Country Capital Population
1 India New Delhi 1303171035
2 Brazil Brasilia 207847528
Get subset of a DataFrame
```

#### Selecting, Boolean Indexing & Setting

```
By Position
>>> df.iloc[0], (0)
'Belgium'
Select single value by row & column
>>> df.iat[0], (0)
'Belgium'
By Label
>>> df.loc[0], ['Country']
'Belgium'
Select single value by row & column labels
>>> df.at[0], ['Country']
'Belgium'
By Label/Position
>>> df.ix[2]
Country Brazil
Capital Brasilia
Population 207847528
Select single row of subset of rows
>>> df.ix[:, 'Capital']
0 Brussels
1 New Delhi
2 Brasilia
Select a single column of subset of columns
>>> df.ix[1, 'Capital']
'New Delhi'
Select rows and columns
Boolean Indexing
>>> s[~(s > 1)]
>>> s[(s < -1) | (s > 2)]
>>> df[df['Population'] > 1200000000]
Use filter to adjust DataFrame
Setting
>>> s['a'] = 6
Set index a of Series s to 6
```

### Dropping

```
>>> s.drop(['a', 'c'])
Drop values from rows (axis=0)
>>> df.drop('Country', axis=1)
Drop values from columns (axis=1)
```

### Sort & Rank

```
>>> df.sort_index(by='Country')
Sort by row or column index
>>> s.order()
Sort a series by its values
>>> df.rank()
Assign ranks to entries
```

### Retrieving Series/DataFrame Information

```
Basic Information
>>> df.shape (rows, columns)
>>> df.index Describe index
>>> df.columns Describe DataFrame columns
>>> df.info() Info on DataFrame
>>> df.count() Number of non-NA values
Summary
>>> df.sum() Sum of values
>>> df.cumsum() Cumulative sum of values
>>> df.min()/df.max() Minimum/Maximum values
>>> df.idxmin()/df.idxmax() Minimum/Maximum index value
>>> df.describe() Summary statistics
>>> df.mean() Mean of values
>>> df.median() Median of values
```

### Applying Functions

```
>>> f = lambda x: x*2
>>> df.apply(f)
Apply function
>>> df.applymap(f)
Apply function element-wise
```

### Data Alignment

```
Internal Data Alignment
NA values are introduced in the indices that don't overlap:
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])
>>> s + s3
a 10.0
b NaN
c 5.0
d 7.0
```

### Arithmetic Operations with Fill Methods

```
You can also do the internal data alignment yourself with the help of the fill methods:
>>> s.add(s3, fill_value=0)
a 10.0
b -5.0
c 5.0
d 7.0
>>> s.sub(s3, fill_value=2)
>>> s.div(s3, fill_value=4)
>>> s.mul(s3, fill_value=3)
```

## DataCamp

Learn Python for Data Science Interactively



# Python For Data Science Cheat Sheet

## NumPy Basics

Learn Python for Data Science Interactively at [www.datacamp.com](https://www.datacamp.com)



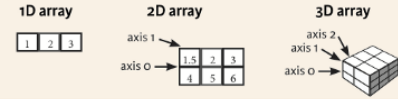
### NumPy

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.



Use the following import convention:  
>>> import numpy as np

#### NumPy Arrays



### Creating Arrays

```
>>> a = np.array([1,2,3])
>>> b = np.array([[1.5,2,3], (4,5,6)], dtype = float)
>>> c = np.array([[1.5,2,3], (4,5,6)], [(3,2,1), (4,5,6)],
                dtype = float)
```

#### Initial Placeholders

```
>>> np.zeros((3,4))           Create an array of zeros
>>> np.ones((2,3,4),dtype=np.int16) Create an array of ones
>>> d = np.arange(10,25,5)    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
```

### I/O

#### Saving & Loading On Disk

```
>>> np.save('my_array', a)
>>> np.savetxt('array.npy', 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           Signed 64-bit integer types
>>> np.float32        Standard double-precision floating point
>>> np.complex         Complex numbers represented by 128 floats
>>> np.bool            Boolean type storing TRUE and FALSE values
>>> np.object          Python object type
>>> np.string_         Fixed-length string type
>>> np.unicode_        Fixed-length unicode type
```

### Inspecting Your Array

```
>>> a.shape           Array dimensions
>>> len(a)            Length of array
>>> b.ndim            Number of array dimensions
>>> e.size            Number of array elements
>>> b.dtype           Data type of array elements
>>> b.dtype.name      Name of data type
>>> b.astype(int)     Convert an array to a different type
```

### Asking For Help

```
>>> np.info(np.ndarray.dtype)
```

### Array Mathematics

#### Arithmetic Operations

```
>>> g = a - b           Subtraction
>>> array([[ -0.5,  0. ,  0. ],
          [-3. , -3. , -3. ]])
>>> np.subtract(a,b)   Subtraction
>>> b + a               Addition
>>> array([[ 2.5,  4. ,  6. ],
          [ 5. ,  7. ,  9. ]])
>>> np.add(b,a)        Addition
>>> a / b               Division
>>> array([[ 0.66666667,  1. ,  1. ],
          [ 0.25 ,  0.4 ,  0.5 ]])
>>> np.divide(a,b)     Division
>>> a * b               Multiplication
>>> array([[ 1.5,  4. ,  9. ],
          [ 4. , 10. , 18. ]])
>>> np.multiply(a,b)   Multiplication
>>> np.exp(b)           Exponentiation
>>> np.sqrt(b)          Square root
>>> np.sin(a)           Print sines of an array
>>> np.cos(b)           Element-wise cosine
>>> np.log(a)           Element-wise natural logarithm
>>> np.dot(E)           Dot product
>>> array([[ 7. ,  7. ],
          [ 7. ,  7.]])
```

#### Comparison

```
>>> a == b             Element-wise comparison
>>> array([[False,  True,  True],
          [False, False, False]], dtype=bool)
>>> a < 2              Element-wise comparison
>>> array([[ True, False, False],
          [ True,  True,  True]], dtype=bool)
>>> np.array_equal(a, b) Array-wise comparison
```

#### Aggregate Functions

```
>>> a.sum()            Array-wise sum
>>> a.min()            Array-wise minimum value
>>> b.max(axis=0)       Maximum value of an array row
>>> b.cumsum(axis=1)   Cumulative sum of the elements
>>> a.mean()           Mean
>>> b.median()         Median
>>> a.corrcoef()        Correlation coefficient
>>> np.std(b)           Standard deviation
```

### Copying Arrays

```
>>> h = a.view()       Create a view of the array with the same data
>>> np.copy(a)         Create a copy of the array
>>> h = a.copy()       Create a deep copy of the array
```

### Sorting Arrays

```
>>> a.sort()           Sort an array
>>> c.sort(axis=0)     Sort the elements of an array's axis
```

### Subsetting, Slicing, Indexing

Also see Lists

```
Subsetting
>>> a[2]               Select the element at the 2nd index
>>> b[1,2]             Select the element at row 0 column 2
                       (equivalent to b[1][2])
>>> c[0]               Select items at index 0 and 1
>>> b[0:2,1]           Select items at rows 0 and 1 in column 1
                       (equivalent to b[0:1, :])
>>> a[0]               Select all items at row 0
                       (equivalent to b[0:1, :])
>>> c[1,...]           Same as [1, :, :]
>>> a[ : :-1]          Reversed array a
>>> array([[3, 2, 1],
          [4, 5, 6]])

Boolean Indexing
>>> a[a<2]            Select elements from a less than 2
>>> array([[1],
          [2, 3]])

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]])
>>> b[[1, 0, 1, 0]][[1, 0, 1, 2, 0]] Select a subset of the matrix's rows
>>> array([[4, 5, 6, 1.5],
          [2, 5, 2, 3, 1.5]])
          [1, 5, 2, 3, 1.5]])
```

### Array Manipulation

```
Transposing Array
>>> i = np.transpose(b) Permute array dimensions
>>> i.T                 Permute array dimensions

Changing Array Shape
>>> b.ravel()           Flatten the array
>>> g.reshape(3,-2)     Reshape, but don't change data

Adding/Removing Elements
>>> h.resize((2,6))     Return a new array with shape (2,6)
>>> np.append(i, g)     Append items to an array
>>> np.insert(a, 1, 5)  Insert items in an array
>>> np.delete(a, [1])   Delete items from an array

Combining Arrays
>>> np.concatenate((a,d),axis=0) Concatenate arrays
>>> array([ 1,  2,  3, 10, 15, 20])
>>> np.vstack((a,b))    Stack arrays vertically (row-wise)
>>> array([[ 1,  2,  3,  ],
          [ 1.5,  2,  3,  ],
          [ 4,  5,  6,  ]])
>>> np.r_*(e,E)         Stack arrays vertically (row-wise)
>>> np.hstack((e,f))    Stack arrays horizontally (column-wise)
>>> array([[ 7,  7,  0,  1]])
>>> np.column_stack((a,d)) Create stacked column-wise arrays
>>> array([[ 1, 10],
          [ 2, 20],
          [ 3, 20]])
>>> np.c_[a,d]          Create stacked column-wise arrays

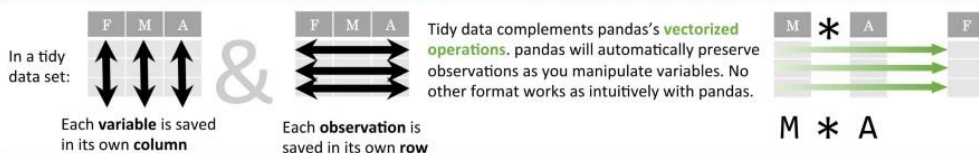
Splitting Arrays
>>> np.hsplit(a,3)      Split the array horizontally at the 3rd
>>> [array([1]),array([2]),array([3])] index
>>> np.vsplit(c,2)     Split the array vertically at the 2nd index
>>> (array([[1.5, 2, 1. ]],
          [ 4,  5,  6, ]]),
   array([[ 3,  2,  3],
          [ 4,  5,  6, ]]))
```

DataCamp

Learn Python for Data Science Interactively



## Tidy Data – A foundation for wrangling in pandas



## Syntax – Creating DataFrames

	a	b	c
1	4	7	10
2	5	8	11
3	6	9	12

```
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.

n	v	a	b	c
d	1	4	7	10
e	2	5	8	11
e	2	6	9	12

```
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).

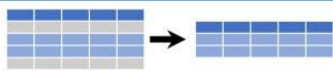
df.rename(columns = {'y': 'year'})
Rename the columns of a DataFrame

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.
```

## 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.
```

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

<http://pandas.pydata.org/> This cheat sheet inspired by RStudio Data Wrangling Cheatsheet (<https://www.rstudio.com/wp-content/uploads/2015/02/data-wrangling-cheatsheet.pdf>) Written by Irv Lustin, Princeton Consultants

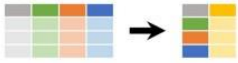
## Summarize Data

**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.

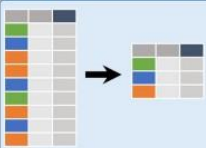
**df.describe()**  
Basic descriptive statistics for each column (or GroupBy)



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:

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

**agg(function)**  
Aggregate group using function.

## Handling Missing Data

**df.dropna()**  
Drop rows with any column having NA/null data.

**df.fillna(value)**  
Replace all NA/null data with value.

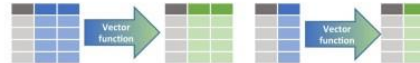
## Make New Columns



**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.



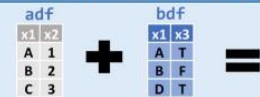
pandas provides a large set of **vector functions** that operate on all columns of a DataFrame or a single selected column (a pandas Series). These functions produce vectors of values for each of the columns, or a single Series for the individual Series. Examples:

<b>max(axis=1)</b> Element-wise max.	<b>min(axis=1)</b> Element-wise min.
<b>clip(lower=-10, upper=10)</b> Trim values at input thresholds	<b>abs()</b> 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.

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

## 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).

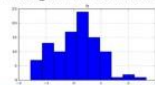
## Windows

**df.expanding()**  
Return an Expanding object allowing summary functions to be applied cumulatively.

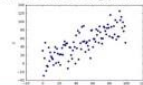
**df.rolling(n)**  
Return a Rolling object allowing summary functions to be applied to windows of length n.

## Plotting

**df.plot.hist()**  
Histogram for each column



**df.plot.scatter(x='w', y='h')**  
Scatter chart using pairs of points



# Data Wrangling with dplyr and tidyr

## Cheat Sheet



### 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:

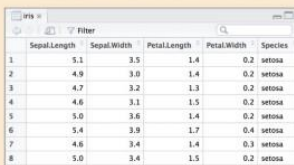
```
Source: local data frame [150 x 5]
  Sepal.Length Sepal.Width Petal.Length
1           5.1           3.5           1.4
2           4.9           3.0           1.4
3           4.7           3.2           1.3
4           4.6           3.1           1.5
5           5.0           3.6           1.4
..          ...           ...           ...
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.

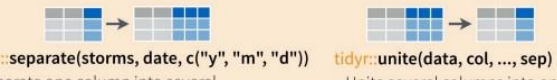
```
iris
  %>%
  group_by(Species) %>%
  summarise(avg = mean(Sepal.Width)) %>%
  arrange(avg)
```

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



**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.

**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).

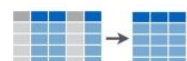
**dplyr::rename(tb, y = year)**  
Rename the columns of a data frame.

### Subset 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

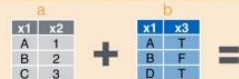
**pmax**

Element-wise max

**pmin**

Element-wise min

## Combine Data Sets



**Mutating Joins**

x1	x2	x3
A	1	T
B	2	F
C	3	NA

**dplyr::left\_join(a, b, by = "x1")**  
Join matching rows from b to a.

x1	x2	x3
A	T	1
B	F	2
D	T	NA

**dplyr::right\_join(a, b, by = "x1")**  
Join matching rows from a to b.

x1	x2	x3
A	1	T
B	2	F

**dplyr::inner\_join(a, b, by = "x1")**  
Join data. Retain only rows in both sets.

x1	x2	x3
A	1	T
C	3	NA
D	NA	T

**dplyr::full\_join(a, b, by = "x1")**  
Join data. Retain all values, all rows.

**Filtering Joins**

x1	x2
A	1
B	2

**dplyr::semi\_join(a, b, by = "x1")**  
All rows in a that have a match in b.

x1	x2
C	3

**dplyr::anti\_join(a, b, by = "x1")**  
All rows in a that do not have a match in b.



**Set Operations**

x1	x2
B	2
C	3

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

x1	x2
A	1
B	2
C	3
D	4

**dplyr::union(y, z)**  
Rows that appear in either or both y and z.

x1	x2
A	1

**dplyr::setdiff(y, z)**  
Rows that appear in y but not z.

**Binding**

x1	x2
A	1
B	2
C	3
D	4

**dplyr::bind\_rows(y, z)**  
Append z to y as new rows.

x1	x2	x1	x2
A	1	B	2
B	2	C	3
C	3	D	4

**dplyr::bind\_cols(y, z)**  
Append z to y as new columns.  
Caution: matches rows by position.

# Python For Data Science Cheat Sheet

## SciPy - Linear Algebra

Learn More Python for Data Science Interactively at [www.datacamp.com](http://www.datacamp.com)



### SciPy

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



### Interacting With NumPy

Also see 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] Create a dense meshgrid
>>> np.ogrid[0:2,0:2] Create an open meshgrid
>>> np.r_[3,0]+5,-1:1:10j Stack arrays vertically (row-wise)
>>> np.c_[b,c] Create stacked column-wise arrays
```

### Shape Manipulation

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

### Polynomials

```
>>> from numpy import poly1d
>>> p = poly1d([3,4,5]) Create a polynomial object
```

### Vectorizing Functions

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

### Type Handling

```
>>> np.real(b) Return the real part of the array elements
>>> np.imag(b) Return the imaginary part of the array elements
>>> np.real_if_close(c,tol=1000) Return a real array if complex parts close to 0
>>> np.cast['*f'](np.pi) Cast object to a data type
```

### Other Useful Functions

```
>>> np.angle(b,deg=True) Return the angle of the complex argument
>>> g = np.linspace(0,np.pi,num=5) 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
>>> misc.comb(10,3,exact=True) Combine N things taken at k time
>>> misc.central_diff_weights(3) 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`.

Also see NumPy

```
>>> 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)))
>>> D = np.mat([[13,4], [5,6]])
```

### Basic Matrix Routines

```
Inverse
>>> linalg.inv(A) Inverse Inverse
>>> A.T Transpose matrix Transpose matrix
>>> A.H Conjugate transpose Conjugate transpose

Trace
>>> np.trace(A) Trace Trace

Norm
>>> linalg.norm(A) Frobenius norm Frobenius norm
>>> linalg.norm(A,1) L1 norm (max column sum) L1 norm (max column sum)
>>> linalg.norm(A,np.inf) Linf norm (max row sum) Linf norm (max row sum)

Rank
>>> np.linalg.matrix_rank(C) Matrix rank Matrix rank

Determinant
>>> linalg.det(A) Determinant Determinant

Solving linear problems
>>> linalg.solve(A,b) Solver for dense matrices Solver for dense matrices
>>> E = np.mat(A.T) Solver for dense matrices Solver for dense matrices
>>> linalg.lstsq(F,E) Least-squares solution to linear matrix equation Least-squares solution to linear matrix equation

Generalized inverse
>>> linalg.pinv(C) Compute the pseudo-inverse of a matrix (least-squares solver) Compute the pseudo-inverse of a matrix (SVD)
>>> linalg.pinv2(C) Compute the pseudo-inverse of a matrix (SVD) Compute the pseudo-inverse of a matrix (SVD)
```

### Creating Sparse Matrices

```
>>> F = np.eye(3, k=1) Create a 2x2 identity matrix Create a 2x2 identity matrix
>>> G = np.mat(np.identity(2)) Create a 2x2 identity matrix Create a 2x2 identity matrix
>>> C[C > 0.5] = 0
>>> H = sparse.csr_matrix(C) Compressed Sparse Row matrix Compressed Sparse Row matrix
>>> T = sparse.csc_matrix(D) Compressed Sparse Column matrix Compressed Sparse Column matrix
>>> J = sparse.dok_matrix(A) Dictionary Of Keys matrix Dictionary Of Keys matrix
>>> E.todense() Sparse matrix to full matrix Sparse matrix to full matrix
>>> sparse.isspmatrix_csc(A) Identify sparse matrix Identify sparse matrix
```

### Sparse Matrix Routines

```
Inverse
>>> sparse.linalg.inv(I) Inverse Inverse

Norm
>>> sparse.linalg.norm(I) Norm Norm

Solving linear problems
>>> sparse.linalg.spsolve(H,I) Solver for sparse matrices Solver for sparse matrices
```

### Sparse Matrix Functions

```
>>> sparse.linalg.expm(I) Sparse matrix exponential Sparse matrix exponential
```

### Asking For Help

```
>>> help(scipy.linalg.diagsvd)
>>> np.info(np.matrix)
```

### Matrix Functions

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

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

Division
>>> np.divide(A,D) Division Division

Multiplication
>>> A @ B Multiplication operator (Python 3) Multiplication operator
>>> np.multiply(D,A) Multiplication Multiplication
>>> np.dot(A,D) Vector dot product Vector dot product
>>> np.vdot(A,D) Inner product Inner product
>>> np.outer(A,D) Outer product Outer product
>>> np.tensordot(A,D) Tensor dot product Tensor dot product
>>> np.kron(A,D) Kronecker product Kronecker product

Exponential Functions
>>> linalg.expm(A) Matrix exponential Matrix exponential
>>> linalg.expm2(A) Matrix exponential (Taylor Series) Matrix exponential (Taylor Series)
>>> linalg.expm3(D) Matrix exponential (eigenvalue decomposition) Matrix exponential (eigenvalue decomposition)

Logarithm Function
>>> linalg.logm(A) Matrix logarithm Matrix logarithm

Trigonometric Functions
>>> linalg.sinm(D) Matrix sine Matrix sine
>>> linalg.cosm(D) Matrix cosine Matrix cosine
>>> linalg.tanm(A) Matrix tangent Matrix tangent

Hyperbolic Trigonometric Functions
>>> linalg.sinhm(D) Hyperbolic matrix sine Hyperbolic matrix sine
>>> linalg.coshm(D) Hyperbolic matrix cosine Hyperbolic matrix cosine
>>> linalg.tanhm(A) Hyperbolic matrix tangent Hyperbolic matrix tangent

Matrix Sign Function
>>> np.signm(A) Matrix sign function Matrix sign function

Matrix Square Root
>>> linalg.sqrtm(A) Matrix square root Matrix square root

Arbitrary Functions
>>> linalg.funm(A, lambda x: x*x) Evaluate matrix function Evaluate matrix function
```

### Decompositions

```
Eigenvalues and Eigenvectors
>>> la, v = linalg.eig(A) Solve ordinary or generalized eigenvalue problem for square matrix Solve ordinary or generalized eigenvalue problem for square matrix
>>> l1, l2 = la Unpack eigenvalues Unpack eigenvalues
>>> v[:,0] First eigenvector First eigenvector
>>> v[:,1] Second eigenvector Second eigenvector
>>> linalg.eigvals(A) Unpack eigenvalues Unpack eigenvalues

Singular Value Decomposition
>>> U, s, Vh = linalg.svd(B) Singular Value Decomposition (SVD) Singular Value Decomposition (SVD)
>>> M, N = B.shape Construct sigma matrix in SVD Construct sigma matrix in SVD
>>> Sig = linalg.diagsvd(s,M,N) Construct sigma matrix in SVD Construct sigma matrix in SVD

LU Decomposition
>>> P, L, U = linalg.lu(C) LU Decomposition LU Decomposition
```

### Sparse Matrix Decompositions

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

DataCamp

Learn Python for Data Science Interactively



# Python For Data Science Cheat Sheet

## Matplotlib

Learn Python Interactively at [www.datacamp.com](http://www.datacamp.com)



### Matplotlib

Matplotlib is a Python 2D plotting library which produces publication-quality figures in a variety of hardcopy formats and interactive environments across platforms.



### 1 Prepare The Data

Also see Lists & NumPy

```
1D Data
>>> import numpy as np
>>> x = np.linspace(0,100, 100)
>>> y = np.cos(x)
>>> z = np.sin(x)

2D Data or Images
>>> data = 2 * np.random.random((10, 10))
>>> data2 = 3 * np.random.random((10, 10))
>>> Y, X = np.mgrid[-3:3:100j, -3:3:100j]
>>> U = -1 * X**2 + Y
>>> from matplotlib.cbook import get_sample_data
>>> img = np.load(get_sample_data('axes_grid/bivariate_normal.npy'))
```

### 2 Create Plot

```
>>> import matplotlib.pyplot as plt

Figure
>>> fig = plt.figure()
>>> fig2 = plt.figure(figsize=plt.figaspect(2.0))

Axes
All plotting is done with respect to an Axes. In most cases, a subplot will fit your needs. A subplot is an axes on a grid system.
>>> fig.add_axes()
>>> ax1 = fig.add_subplot(221) # row-col-num
>>> ax3 = fig.add_subplot(212)
>>> fig3, axes = plt.subplots(nrows=2,ncols=2)
>>> fig4, axes2 = plt.subplots(ncols=3)
```

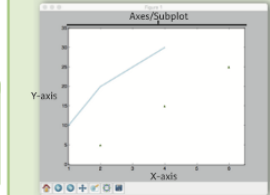
### 3 Plotting Routines

```
1D Data
>>> fig, ax = plt.subplots()
>>> lines = ax.plot(x,y) Draw points with lines or markers connecting them
>>> ax.scatter(x,y) Draw unconnected points, scaled or colored
>>> axes[0,0].bar([1,2,3],[3,4,5]) Plot vertical rectangles (constant width)
>>> axes[1,0].barh([0.5,1,2.5],[0,1,2]) Plot horizontal rectangles (constant height)
>>> axes[1,1].axhline(0.45) Draw a horizontal line across axes
>>> axes[0,1].axvline(0.65) Draw a vertical line across axes
>>> ax.fill(x,y,color='blue') Draw filled polygons
>>> ax.fill_between(x,y,color='yellow') Fill between y-values and 0

2D Data or Images
>>> fig, ax = plt.subplots()
>>> im = ax.imshow(img, cmap='gist_earth', interpolation='nearest', vmin=-4, vmax=2) Colormapped or RGB arrays
```

## Plot Anatomy & Workflow

### Plot Anatomy



### Workflow

The basic steps to creating plots with matplotlib are:

- 1 Prepare data
- 2 Create plot
- 3 Plot
- 4 Customize plot
- 5 Save plot
- 6 Show plot

### 4 Customize Plot

```
Colors, Color Bars & Color Maps
>>> plt.plot(x, x, x, x**2, x, x**3)
>>> ax.plot(x, y, alpha=0.4)
>>> ax.plot(x, y, c='k')
>>> fig.colorbar(img, orientation='horizontal')
>>> im = ax.imshow(img, cmap='seismic')

Markers
>>> fig, ax = plt.subplots()
>>> ax.scatter(x,y,marker=".")
>>> ax.plot(x,y,marker="o")

Linestyles
>>> plt.plot(x,y,linewidth=4.0)
>>> plt.plot(x,y,ls='solid')
>>> plt.plot(x,y,ls='--')
>>> plt.plot(x,y,'--',x**2,y**2,'-.')
>>> plt.setp(lines,color='r',linewidth=4.0)

Text & Annotations
>>> ax.text(1, 2, 'Example Graph', style='italic')
>>> ax.annotate("Sine", xy=(8, 0), xycoords='data', xytext=(10.5, 0), textcoords='data', arrowprops=dict(arrowstyle="→", connectionstyle="arc3"),)
```

```
Mathtext
>>> plt.title(r'$\sigma_i=15\$', fontsize=20)

Limits, Legends & Layouts
Limits & Autoscaling
>>> ax.margins(x=0,y=0.1)
>>> ax.axis('equal')
>>> ax.set_xlim([0,10.5],ylim=[-1.5,1.5])
>>> ax.set_ylim(0,10.5)

Legends
>>> ax.set(title='An Example Axes', ylabel='Y-Axis', xlabel='X-Axis')

Ticks
>>> ax.xaxis.set(ticks=range(1,5), ticklabel=[3,100,-12,"foo"])
>>> ax.tick_params(axis='y', direction='inout', length=10)

Subplot Spacing
>>> fig3.subplots_adjust(wspace=0.5, hspace=0.3, left=0.125, right=0.9, top=0.9, bottom=0.1)

Axis Spines
>>> ax1.spines['top'].set_visible(False)
>>> ax1.spines['bottom'].set_position(('outward',10))
```

### 5 Save Plot

```
Save figures
>>> plt.savefig('foo.png')
Save transparent figures
>>> plt.savefig('foo.png', transparent=True)
```

### 6 Show Plot

```
>>> plt.show()
```

### Close & Clear

```
>>> plt.cla() Clear an axis
>>> plt.clf() Clear the entire figure
>>> plt.close() Close a window
```

DataCamp

Learn Python for Data Science Interactively





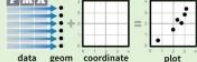
# Data Visualization with ggplot2

Cheat Sheet



## 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()`

aesthetic mappings    data    geom

`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()`.

data    add layers, elements with +

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

layer = geom + default stat + layer specific mappings

additional elements

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(fill))`

`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, fill, 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



`f + 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, 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, color, fill, linetype, size



`i + geom_hex()`  
x, y, alpha, color, 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, shape, 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_raster(aes(fill = z), hjust=0.5, vjust=0.5, interpolate=FALSE)`  
x, y, alpha, fill



`m + geom_contour(aes(z = z))`  
x, y, z, alpha, color, linetype, size, weight



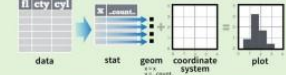
`m + geom_tile(aes(fill = z))`  
x, y, alpha, color, fill, linetype, size



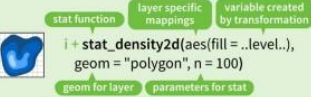
`m + geom_tile(aes(fill = z))`  
x, y, alpha, color, fill, linetype, size

## 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)
  
```

```

a + stat_bin(binwidth = 1, origin = 1)
x, y | ..count.., ..ncount.., ..density..
a + stat_bin2d(binwidth = 1, binaxis = "x")
x, y | ..count.., ..ncount..
a + stat_bin2d(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, 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 | ..quartile.., ..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_uniq()
  
```

## 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 = fill))
  
```

Scale: `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 = fill))
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 = 2.5)
o + scale_fill_gradientn(colors = terrain.colors(6))
Also: rainbow(), heat.colors(), topo.colors(), cm.colors(), RColorBrewer::brewer.pal()
  
```

```

Shape scales
Manual shape values
p <- f + geom_point(aes(shape = fill))
p + scale_shape(solid = FALSE)
p + scale_shape_manual(values = c(3, 7))
Shape values shown in chart on right
  
```

```

Size scales
q <- f + geom_point(aes(size = cyl))
o + scale_size_area(max = 6)
Value mapped to area of circle (not radius)
  
```

## Coordinate Systems

```

r <- b + geom_bar()
xlim, ylim

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))
  
```

## Themes

```

r + theme_bw()
White background with grid lines
r + theme_classic()
White background no gridlines
r + theme_grey()
Grey background (default theme)
r + theme_minimal()
Minimal theme
ggthemes - Package with additional ggplot2 themes
  
```

## 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)
fl: c fl: d fl: e fl: p fl: r
t + facet_grid(~ fl, labeller = label_quote(alpha ^ (.)))
alpha alpha alpha alpha alpha
t + facet_grid(~ fl, labeller = label_parsed)
c d e p r
  
```

## 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
  
```

## 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.
  
```

## 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))
  
```

# Python For Data Science Cheat Sheet

## PySpark Basics

Learn Python for data science interactively at [www.DataCamp.com](http://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
>>> sc.pythonVer
>>> sc.master
>>> str(sc.sparkHome)
>>> str(sc.sparkUser())
>>> sc.appName
>>> sc.applicationId
>>> sc.defaultParallelism
>>> sc.defaultMinPartitions
```

Retrieve SparkContext version  
Retrieve Python version  
Master URL to connect to  
Path where Spark is installed on worker nodes  
Retrieve name of the Spark User running SparkContext  
Return application name  
Retrieve application ID  
Return default level of parallelism  
Default minimum number of partitions for RDDs

### Configuration

```
>>> from pyspark import SparkConf, SparkContext
>>> conf = SparkConf()
>>> conf.setMaster("local")
>>> conf.setAppName("My app")
>>> conf.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', 'm', 'y', 'm', 'e'], ('b', 'm', 'e'))
```

#### 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()
>>> rdd.count()
>>> rdd.countByKey()
>>> rdd.countByValue()
>>> rdd.collectAsMap()
>>> rdd3.sum()
>>> sc.parallelize().isEmpty()
```

List the number of partitions  
Count RDD instances  
Count RDD instances by key  
Count RDD instances by value  
Return (key,value) pairs as a dictionary  
Sum of RDD elements  
Check whether RDD is empty

#### Summary

```
>>> rdd3.max()
>>> rdd3.min()
>>> rdd3.mean()
>>> rdd3.stdev()
>>> rdd3.variance()
>>> rdd3.histogram()
>>> rdd3.stats()
```

Maximum value of RDD elements  
Minimum value of RDD elements  
Mean value of RDD elements  
Standard deviation of RDD elements  
Compute variance of RDD elements  
Compute histogram by bins  
Summary statistics (count, mean, stdev, max & min)

### Applying Functions

```
>>> rdd.map(lambda x: x*(x[1],x[0]))
>>> rdd5 = rdd.flatMap(lambda x: x*(x[1],x[0]))
>>> rdd5.collect()
>>> rdd4.flatMapValues(lambda x: x)
>>> rdd4.collect()
```

Apply a function to each RDD element  
Apply a function to each RDD element and flatten the result  
Apply a flatMap function to each (key,value) pair of rdd4 without changing the keys

### Selecting Data

#### Getting

```
>>> rdd.collect()
>>> rdd.take(2)
>>> rdd.first()
>>> rdd.top(2)
>>> rdd3.sample(False, 0.15, 81).collect()
```

Return a list with all RDD elements  
Take first 2 RDD elements  
Take first RDD element  
Take top 2 RDD elements  
Return sampled subset of rdd3

#### Filtering

```
>>> rdd.filter(lambda x: "a" in x)
>>> rdd5.distinct().collect()
>>> rdd.keys().collect()
```

Filter the RDD  
Return distinct RDD values  
Return (key,value) RDD's keys

### Iterating

```
>>> def f(x): print(x)
>>> rdd.foreach(f)
>>> ('b', 2)
>>> ('a', 2)
```

Apply a function to all RDD elements

### Reshaping Data

#### Reducing

```
>>> rdd.reduceByKey(lambda x,y: x+y)
>>> rdd.reduce(lambda a, b: a + b)
```

Merge the rdd values for each key  
Merge the rdd values

#### Grouping by

```
>>> rdd3.groupBy(lambda x: x % 2)
>>> rdd.groupByKey()
>>> rdd.foldByKey(0, add)
>>> rdd.fold(0, add)
```

Return RDD of grouped values  
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)
>>> rdd3.aggregateByKey((0,0),seqOp,combOp)
>>> rdd3.fold(0,add)
>>> rdd.foldByKey(0, add)
>>> rdd3.keyBy(lambda x: x*x)
>>> rdd.cartesian(rdd2).collect()
```

Aggregate RDD elements of each partition and then the results  
Aggregate values of each RDD key  
Aggregate the elements of each partition, and then the results  
Merge the values for each key  
Create tuples of RDD elements by applying a function

### Mathematical Operations

```
>>> rdd.subtract(rdd2)
>>> rdd2.subtractByKey(rdd)
>>> rdd.cartesian(rdd2).collect()
```

Return each rdd value not contained in rdd2  
Return each (key,value) pair of rdd2 with no matching key in rdd  
Return the Cartesian product of rdd and rdd2

### Sort

```
>>> rdd2.sortBy(lambda x: x[1])
>>> rdd2.sortBy(lambda x: x[1], ('a', 2))
>>> rdd2.sortBy(lambda x: x[1], ('a', 2))
>>> rdd.cartesian(rdd2).collect()
```

Sort RDD by given function  
Sort (key, value) RDD by key

### Repartitioning

```
>>> rdd.repartition(4)
>>> rdd.coalesce(1)
```

New RDD with 4 partitions  
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



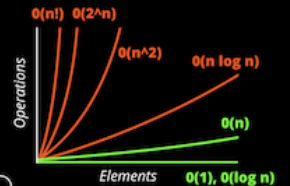
### LEGEND

TIME Complexity VS. SPACE Complexity



## <BIG-O-CHEATSHEET>

www.bigocheatsheet.com



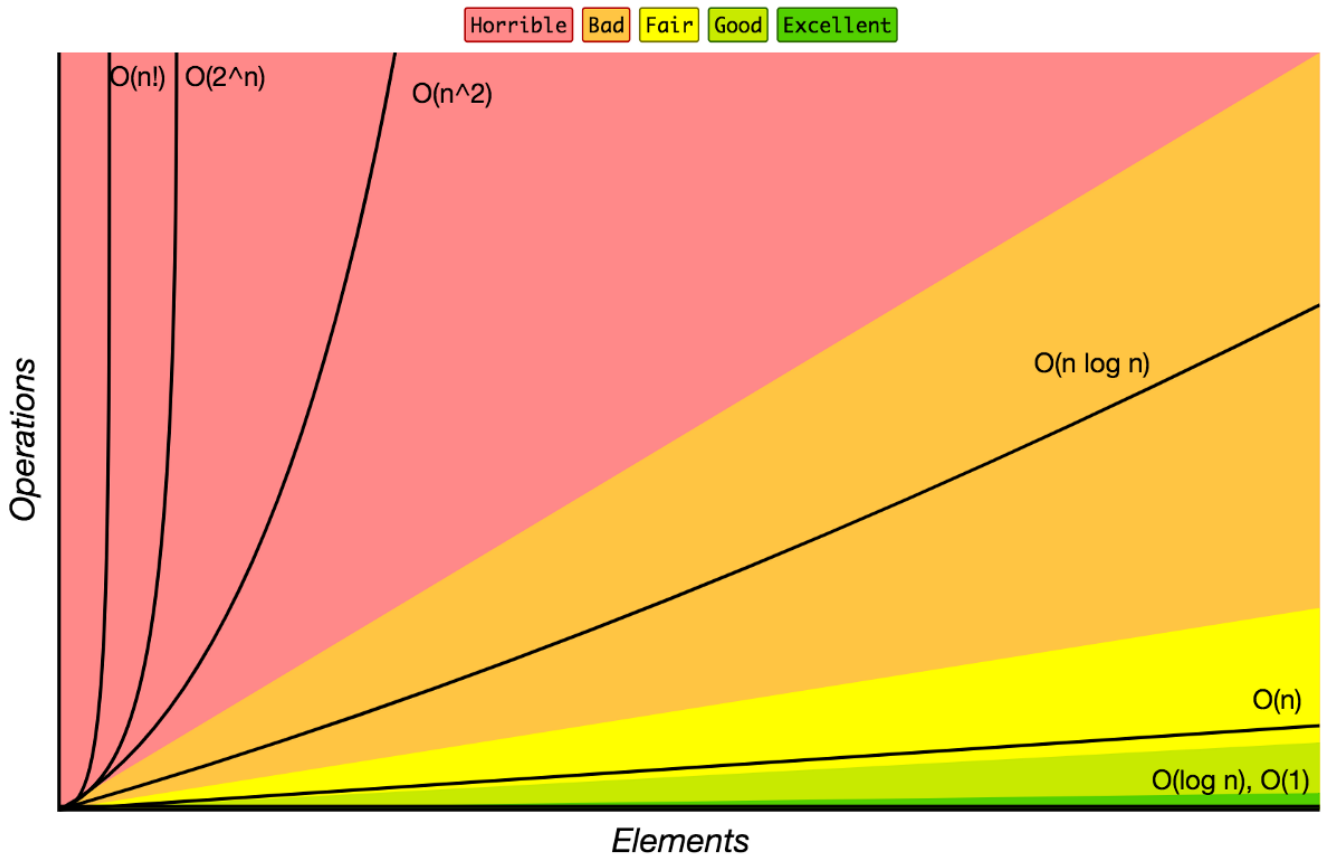
### DATA STRUCTURE Operations

DATA Structure	TIME Complexity				SPACE Complexity			
	Access	Search	Insertion	Deletion	Access	Search	Insertion	Deletion
Array	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$
Stack	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$
Queue	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$
Singly-Linked List	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$
Doubly-Linked List	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$
Skip List	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n \log(n))$
Hash Table	N/A	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$	N/A	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$
Binary Search Tree	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$
Cartesian Tree	N/A	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	N/A	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$
B-Tree	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$
Red-Black Tree	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$
Splay Tree	N/A	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	N/A	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$
AVL Tree	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$
KD Tree	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$

### ARRAY SORTING Algorithms

ARRAY Algorithms	TIME Complexity			SPACE Complexity
	Best	Average	Worst	Worst
Quicksort	$\Omega(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n^2)$	$\Theta(\log(n))$
Mergesort	$\Omega(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n)$
Timsort	$\Omega(n)$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(1)$
Heapsort	$\Omega(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(1)$
Bubble Sort	$\Omega(n)$	$\Theta(n^2)$	$\Theta(n^2)$	$\Theta(1)$
Insertion Sort	$\Omega(n)$	$\Theta(n^2)$	$\Theta(n^2)$	$\Theta(1)$
Selection Sort	$\Omega(n^2)$	$\Theta(n^2)$	$\Theta(n^2)$	$\Theta(1)$
Tree Sort	$\Omega(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n^2)$	$\Theta(n)$
Shell Sort	$\Omega(n \log(n))$	$\Theta(n(\log(n))^2)$	$\Theta(n(\log(n))^2)$	$\Theta(1)$
Bucket Sort	$\Omega(n+k)$	$\Theta(n+k)$	$\Theta(n^2)$	$\Theta(n)$
Radix Sort	$\Omega(nk)$	$\Theta(nk)$	$\Theta(nk)$	$\Theta(n+k)$
Counting Sort	$\Omega(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(k)$
Cubesort	$\Omega(n)$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n)$

# Big-O Complexity Chart



## Common Data Structure Operations

Data Structure	Time Complexity								Space Complexity
	Average				Worst				Worst
	Access	Search	Insertion	Deletion	Access	Search	Insertion	Deletion	
Array	$O(1)$	$\theta(n)$	$\theta(n)$	$\theta(n)$	$O(1)$	$O(n)$	$O(n)$	$O(n)$	$O(n)$
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
<u>Quicksort</u>	$\Omega(n \log(n))$	$\theta(n \log(n))$	$O(n^2)$	$O(\log(n))$
<u>Mergesort</u>	$\Omega(n \log(n))$	$\theta(n \log(n))$	$O(n \log(n))$	$O(n)$
<u>Timsort</u>	$\Omega(n)$	$\theta(n \log(n))$	$O(n \log(n))$	$O(n)$
<u>Heapsort</u>	$\Omega(n \log(n))$	$\theta(n \log(n))$	$O(n \log(n))$	$O(1)$
<u>Bubble Sort</u>	$\Omega(n)$	$\theta(n^2)$	$O(n^2)$	$O(1)$
<u>Insertion Sort</u>	$\Omega(n)$	$\theta(n^2)$	$O(n^2)$	$O(1)$
<u>Selection Sort</u>	$\Omega(n^2)$	$\theta(n^2)$	$O(n^2)$	$O(1)$
<u>Tree Sort</u>	$\Omega(n \log(n))$	$\theta(n \log(n))$	$O(n^2)$	$O(n)$
<u>Shell Sort</u>	$\Omega(n \log(n))$	$\theta(n(\log(n))^2)$	$O(n(\log(n))^2)$	$O(1)$
<u>Bucket Sort</u>	$\Omega(n+k)$	$\theta(n+k)$	$O(n^2)$	$O(n)$
<u>Radix Sort</u>	$\Omega(nk)$	$\theta(nk)$	$O(nk)$	$O(n+k)$
<u>Counting Sort</u>	$\Omega(n+k)$	$\theta(n+k)$	$O(n+k)$	$O(k)$
<u>Cubesort</u>	$\Omega(n)$	$\theta(n \log(n))$	$O(n \log(n))$	$O(n)$