# Cheat Sheet of Machine Learning and Python (and Math) Cheat Sheets

Robbie Allen   Follow

Jun 1, 2017 · 4 min read

*If you like this article, check out another by Robbie:*
*My Curated List of AI and Machine Learning Resources*



There are many facets to Machine Learning. As I started brushing up on the

subject, I came across various "cheat sheets" that compactly listed all the key points I needed to know for a given topic. Eventually, I compiled over 20 Machine Learning-related cheat sheets. Some I reference frequently and thought others may benefit from them too. This post contains 27 of the better cheat sheets I've found on the web. Let me know if I'm missing any you like.

Given how rapidly the Machine Learning space is evolving, I imagine these will go out of date quickly, but at least as of June 1, 2017, they are pretty current.

If you want all of the cheat sheets without having to download them individually like I did, I created a zip file containing all 27. Enjoy!
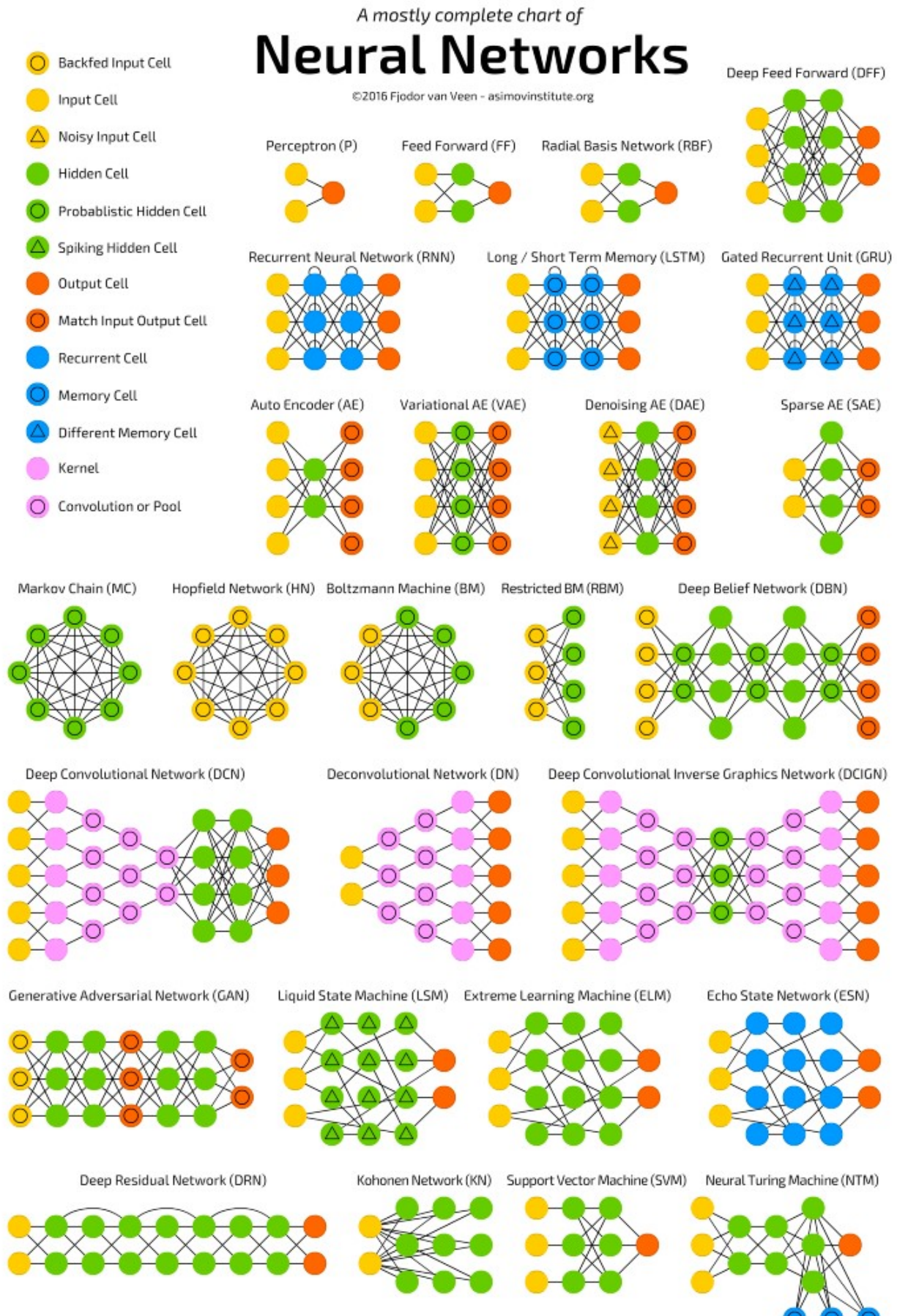
If you like this post, give it a ❤ below.

# Machine Learning

There are a handful of helpful flowcharts and tables of Machine Learning algorithms. I've included only the most comprehensive ones I've found.
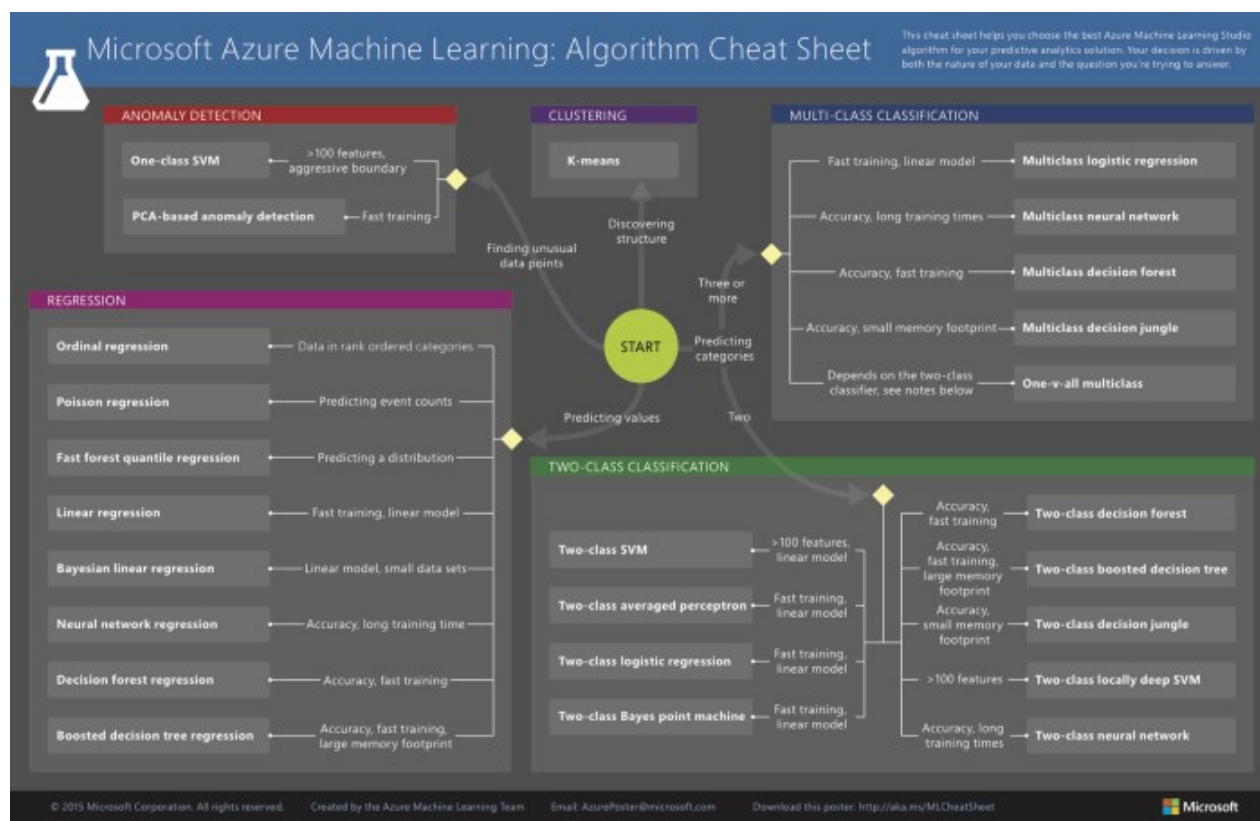
### Neural Network Architectures

Source: http://www.asimovinstitute.org/neural-network-zoo/

# A mostly complete chart of
# Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

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

**Perceptron (P)**

**Feed Forward (FF)**

**Radial Basis Network (RBF)**

**Deep Feed Forward (DFF)**

**Recurrent Neural Network (RNN)**

**Long / Short Term Memory (LSTM)**

**Gated Recurrent Unit (GRU)**

**Auto Encoder (AE)**

**Variational AE (VAE)**

**Denoising AE (DAE)**

**Sparse AE (SAE)**

**Markov Chain (MC)**

**Hopfield Network (HN)**

**Boltzmann Machine (BM)**

**Restricted BM (RBM)**

**Deep Belief Network (DBN)**

**Deep Convolutional Network (DCN)**

**Deconvolutional Network (DN)**

**Deep Convolutional Inverse Graphics Network (DCIGN)**

**Generative Adversarial Network (GAN)**

**Liquid State Machine (LSM)**

**Extreme Learning Machine (ELM)**

**Echo State Network (ESN)**

**Deep Residual Network (DRN)**

**Kohonen Network (KN)**

**Support Vector Machine (SVM)**

**Neural Turing Machine (NTM)**

## Microsoft Azure Algorithm Flowchar

Source: https://docs.microsoft.com/en-us/azure/machine-learning/machine-learning-algorithm-cheat-sheet



Machine learning algorithm cheat sheet for Microsoft Azure Machine Learning Studio

## SAS Algorithm Flowchart

Source: http://blogs.sas.com/content/subconsciousmusings/2017/04/12/machine-learning-algorithm-use/

SAS: Which machine learning algorithm should I use?

## Algorithm Summary

Source: http://machinelearningmastery.com/a-tour-of-machine-learning-algorithms/

A Tour of Machine Learning Algorithms

Source: http://thinkbigdata.in/best-known-machine-learning-algorithms-infographic/

# the world of machine learning algorithms – a summary

## think big data

### regression

Ordinary Least Squares Regression (OLSR)
Linear Regression
Logistic Regression
Stepwise Regression
Multivariate Adaptive Regression Splines (MARS)
Locally Estimated Scatterplot Smoothing (LOESS)
Jackknife Regression

### regularization

Ridge Regression
Least Absolute Shrinkage and Selection Operator (LASSO)
Elastic Net
Least-Angle Regression (LARS))

### instance based

also called cake-based, memory-based

k-Nearest Neighbour (kNN)
Learning Vector Quantization (LVQ)
Self-Organizing Map (SOM)
Locally Weighted Learning (LWL)

### dimesionality reduction

Principal Component Analysis (PCA)
Principal Component Regression (PCR)
Partial Least Squares Regression (PLSR)
Sammon Mapping
Multidimensional Scaling (MDS)
Projection Pursuit
Discriminant Analysis (LDA, MDA, QDA, FDA)

### deep learning

Deep Boltzmann Machine (DBM)
Deep Belief Networks (DBN)
Convolutional Neural Network (CNN)
Stacked Auto-Encoders

### associated rule

Apriori
Eclat
FP-Growth

### ensemble

Logit Boost (Boosting)
Bootstrapped Aggregation (Bagging)
AdaBoost
Stacked Generalization (blending)
Gradient Boosting Machines (GBM)
Gradient Boosted Regression Trees (GBRT)
Random Forest

### bayesian

Naive Bayes
Gaussian Naive Bayes
Multinomial Naive Bayes
Averaged One-Dependence Estimators (AODE)
Bayesian Belief Network (BBN)
Bayesian Network (BN)
Hidden Markov Models
Conditional random fields (CRFs)

### decision tree

Classification and Regression Tree (CART)
Iterative Dichotomiser 3 (ID3)
C4.5 and C5.0 (different versions of a powerful approach)
Chi-squared Automatic Interaction Detection (CHAID)
Decision Stump
M5
Random Forests
Conditional Decision Trees

### clustering

Single-linkage clustering
k-Means
k-Medians
Expectation Maximisation (EM)
Hierarchical Clustering
Fuzzy clustering
DBSCAN
OPTICS algorithm
Non Negative Matrix Factorization
Latent Dirichlet allocation (LDA)

### neural networks

Self Organizing Map
Perceptron
Back-Propagation
Hopfield Network
Radial Basis Function Network (RBFN)
Backpropagation
Autoencoders
Hopfield networks
Boltzmann machines
Restricted Boltzmann Machines
Spiking Neural Networks
Learning Vector quantization (LVQ)

### ...and others

Support Vector Machines (SVM)
Evolutionary Algorithms
Inductive Logic Programming (ILP)
Reinforcement Learning (Q-Learning, Temporal Difference,
State-Action-Reward-State-Action (SARSA))
ANOVA
Information Fuzzy Network (IFN)
Page Rank
Conditional Random Fields (CRF)

## Algorithm Pro/Con

Source: https://blog.dataiku.com/machine-learning-explained-algorithms-are-your-friend

# data iku

## TOP PREDICTION ALGORITHMS

| TYPE | NAME | DESCRIPTION | ADVANTAGES | DISADVANTAGES |
|---|---|---|---|---|
| Linear | Linear regression | The "best fit" line through all data points. Predictions are numerical. | Easy to understand -- you clearly see what the biggest drivers of the model are. | X Sometimes too simple to capture complex relationships between variables.<br>X Tendency for the model to "overfit". |
| Linear | Logistic regression | The adaptation of **linear regression** to problems of classification (e.g., yes/no questions, groups, etc.) | Also easy to understand. | X Sometimes too simple to capture complex relationships between variables.<br>X Tendency for the model to "overfit". |
| Tree-based | Decision tree | A graph that uses a **branching method** to match all possible outcomes of a decision. | Easy to understand and implement. | X Not often used on its own for prediction because it's also often too simple and not powerful enough for complex data. |
| Tree-based | Random Forest | Takes the average of many decision trees, each of which is made with a sample of the data. Each tree is weaker than a full decision tree, but **by combining them we get better overall performance**. | A sort of "wisdom of the crowd". Tends to result in very high quality models. Fast to train. | X Can be slow to output predictions relative to other algorithms.<br>X Not easy to understand predictions. |
| Tree-based | Gradient Boosting | Uses even weaker decision trees, that are increasingly **focused on "hard"** examples. | High-performing. | X A small change in the feature set or training set can create radical changes in the model.<br>X Not easy to understand predictions. |
| Neural networks | Neural networks | Mimics the behavior of the brain. Neural networks are interconnected neurons that pass messages to each other. Deep learning uses several **layers of neural networks put one after the other.** | Can handle extremely complex tasks - no other algorithm comes close in image recognition. | X Very, very slow to train, because they have so many layers. Require a lot of power.<br>X Almost impossible to understand predictions. |

©2017 Dataiku, Inc. | www.dataiku.com | contact@dataiku.com | @dataiku

# Python

Unsurprisingly, there are a lot of online resources available for Python. For this section, I've only included the best cheat sheets I've come across.

## Algorithms

Source: [https://www.analyticsvidhya.com/blog/2015/09/full-cheatsheet-machine-learning-algorithms/](https://www.analyticsvidhya.com/blog/2015/09/full-cheatsheet-machine-learning-algorithms/)

## CHEATSHEET

# Machine Learning Algorithms

( Python and R Codes)

# Types

**Supervised Learning**
- Decision Tree · Random Forest
- kNN · Logistic Regression

**Unsupervised Learning**
- Apriori algorithm · k-means
- Hierarchical Clustering

**Reinforcement Learning**
- Markov Decision Process
- Q Learning

**Python Code**

**R Code**

**Linear Regression**

```
#Import Library
#Import other necessary libraries like pandas,
#numpy...
from sklearn import linear_model
#Load Train and Test datasets
#Identify feature and response variable(s) and
#values must be numeric and numpy arrays
x_train=input_variables_values_training_datasets
y_train=target_variables_values_training_datasets
x_test=input_variables_values_test_datasets
#Create linear regression object
linear = linear_model.LinearRegression()
#Train the model using the training sets and
#check score
linear.fit(x_train, y_train)
linear.score(x_train, y_train)
#Equation coefficient and Intercept
print('Coefficient: \n', linear.coef_)
print('Intercept: \n', linear.intercept_)
#Predict Output
predicted= linear.predict(x_test)
```

```
#Load Train and Test datasets
#Identify feature and response variable(s) and
#values must be numeric and numpy arrays
x_train <- input_variables_values_training_datasets
y_train <- target_variables_values_training_datasets
x_test <- input_variables_values_test_datasets
x <- cbind(x_train,y_train)
#Train the model using the training sets and
#check score
linear <- lm(y_train ~ ., data = x)
summary(linear)
#Predict Output
predicted= predict(linear,x_test)
```

## Python Basics

Source: http://datasciencefree.com/python.pdf



Source: https://www.datacamp.com/community/tutorials/python-data-science-cheat-sheet-basics#gs.0x1rxEA

## Numpy

Source: https://www.dataquest.io/blog/numpy-cheat-sheet/

# Data Science Cheat Sheet
## NumPy

**KEY**
We'll use shorthand in this cheat sheet
arr - A numpy Array object

**IMPORTS**
Import these to start
import numpy as np

### IMPORTING/EXPORTING
np.loadtxt('file.txt') - From a text file
np.genfromtxt('file.csv',delimiter=',') - From a CSV file
np.savetxt('file.txt',arr,delimiter=' ') - Writes to a text file
np.savetxt('file.csv',arr,delimiter=',') - Writes to a CSV file

### CREATING ARRAYS
np.array([1,2,3]) - One dimensional array
np.array([(1,2,3),(4,5,6)]) - Two dimensional array
np.zeros(3) - 1D array of length 3 all values 0
np.ones((3,4)) - 3x4 array with all values 1
np.eye(5) - 5x5 array of 0 with 1 on diagonal (Identity matrix)
np.linspace(0,100,6) - Array of 6 evenly divided values from 0 to 100
np.arange(0,10,3) - Array of values from 0 to less than 10 with step 3 (eg [0,3,6,9])
np.full((2,3),8) - 2x3 array with all values 8
np.random.rand(4,5) - 4x5 array of random floats between 0-1
np.random.rand(6,7)*100 - 6x7 array of random floats between 0-100
np.random.randint(5,size=(2,3)) - 2x3 array with random ints between 0-4

### INSPECTING PROPERTIES
arr.size - Returns number of elements in arr
arr.shape - Returns dimensions of arr (rows, columns)
arr.dtype - Returns type of elements in arr
arr.astype(dtype) - Convert arr elements to type dtype
arr.tolist() - Convert arr to a Python list
np.info(np.eye) - View documentation for np.eye

### COPYING/SORTING/RESHAPING
np.copy(arr) - Copies arr to new memory
arr.view(dtype) - Creates view of arr elements with type dtype
arr.sort() - Sorts arr
arr.sort(axis=0) - Sorts specific axis of arr
two_d_arr.flatten() - Flattens 2D array two_d_arr to 1D

arr.T - Transposes arr (rows become columns and vice versa)
arr.reshape(3,4) - Reshapes arr to 3 rows, 4 columns without changing data
arr.resize((5,6)) - Changes arr shape to 5x6 and fills new values with 0

### ADDING/REMOVING ELEMENTS
np.append(arr,values) - Appends values to end of arr
np.insert(arr,2,values) - Inserts values into arr before index 2
np.delete(arr,3,axis=0) - Deletes row on index 3 of arr
np.delete(arr,4,axis=1) - Deletes column on index 4 of arr

### COMBINING/SPLITTING
np.concatenate((arr1,arr2),axis=0) - Adds arr2 as rows to the end of arr1
np.concatenate((arr1,arr2),axis=1) - Adds arr2 as columns to end of arr1
np.split(arr,3) - Splits arr into 3 sub-arrays
np.hsplit(arr,5) - Splits arr horizontally on the 5th index

### INDEXING/SLICING/SUBSETTING
arr[5] - Returns the element at index 5
arr[2,5] - Returns the 2D array element on index [2][5]
arr[1]=4 - Assigns array element on index 1 the value 4
arr[1,3]=10 - Assigns array element on index [1][3] the value 10
arr[0:3] - Returns the elements at indices 0,1,2 (On a 2D array: returns rows 0,1,2)
arr[0:3,4] - Returns the elements on rows 0,1,2 at column 4
arr[:2] - Returns the elements at indices 0,1 (On a 2D array: returns rows 0,1)
arr[:,1] - Returns the elements at index 1 on all rows
arr<5 - Returns an array with boolean values
(arr1<3) & (arr2>5) - Returns an array with boolean values
~arr - Inverts a boolean array
arr[arr<5] - Returns array elements smaller than 5

### SCALAR MATH
np.add(arr,1) - Add 1 to each array element
np.subtract(arr,2) - Subtract 2 from each array element
np.multiply(arr,3) - Multiply each array element by 3
np.divide(arr,4) - Divide each array element by 4 (returns np.nan for division by zero)
np.power(arr,5) - Raise each array element to the 5th power

### VECTOR MATH
np.add(arr1,arr2) - Elementwise add arr2 to arr1
np.subtract(arr1,arr2) - Elementwise subtract arr2 from arr1
np.multiply(arr1,arr2) - Elementwise multiply arr1 by arr2
np.divide(arr1,arr2) - Elementwise divide arr1 by arr2
np.power(arr1,arr2) - Elementwise raise arr1 raised to the power of arr2
np.array_equal(arr1,arr2) - Returns True if the arrays have the same elements and shape
np.sqrt(arr) - Square root of each element in the array
np.sin(arr) - Sine of each element in the array
np.log(arr) - Natural log of each element in the array
np.abs(arr) - Absolute value of each element in the array
np.ceil(arr) - Rounds up to the nearest int
np.floor(arr) - Rounds down to the nearest int
np.round(arr) - Rounds to the nearest int

### STATISTICS
np.mean(arr,axis=0) - Returns mean along specific axis
arr.sum() - Returns sum of arr
arr.min() - Returns minimum value of arr
arr.max(axis=0) - Returns maximum value of specific axis
np.var(arr) - Returns the variance of array
np.std(arr,axis=1) - Returns the standard deviation of specific axis
arr.corrcoef() - Returns correlation coefficient of array

Source: http://datasciencefree.com/numpy.pdf

Source: https://www.datacamp.com/community/blog/python-numpy-cheat-sheet#gs.Nw3V6CE

Source: https://github.com/donnemartin/data-science-ipython-notebooks
/blob/master/numpy/numpy.ipynb

## NumPy

Credits: Forked from [Parallel Machine Learning with scikit-learn and IPython](#) by Olivier Grisel

- NumPy Arrays, dtype, and shape
- Common Array Operations
- Reshape and Update In-Place
- Combine Arrays
- Create Sample Data

```
In [1]: import numpy as np
```

### NumPy Arrays, dtypes, and shapes

```
In [2]: a = np.array([1, 2, 3])
        print(a)
        print(a.shape)
        print(a.dtype)

        [1 2 3]
        (3,)
        int64
```

```
In [3]: b = np.array([[0, 2, 4], [1, 3, 5]])
        print(b)
        print(b.shape)
        print(b.dtype)

        [[0 2 4]
         [1 3 5]]
        (2, 3)
        int64
```

# Pandas

Source: [http://datasciencefree.com/pandas.pdf](http://datasciencefree.com/pandas.pdf)

Source: https://www.datacamp.com/community/blog/python-pandas-cheat-sheet#gs.S4P4T=U

Source: https://github.com/donnemartin/data-science-ipython-notebooks /blob/master/pandas/pandas.ipynb

## Pandas

Credits: The following are notes taken while working through Python for Data Analysis by Wes McKinney

- Series
- DataFrame
- Reindexing
- Dropping Entries
- Indexing, Selecting, Filtering
- Arithmetic and Data Alignment
- Function Application and Mapping
- Sorting and Ranking
- Axis Indices with Duplicate Values
- Summarizing and Computing Descriptive Statistics
- Cleaning Data (Under Construction)
- Input and Output (Under Construction)

```
In [1]: from pandas import Series, DataFrame
        import pandas as pd
        import numpy as np
```

### Series

A Series is a one-dimensional array-like object containing an array of data and an associated array of data labels. The data can be any NumPy data type and the labels are the Series' index.

Create a Series:

```
In [2]: ser_1 = Series([1, 1, 2, -3, -5, 8, 13])
        ser_1
```

```
Out[2]: 0    1
        1    1
        2    2
```

## Matplotlib

Source: https://www.datacamp.com/community/blog/python-matplotlib-cheat-sheet

Source: https://github.com/donnemartin/data-science-ipython-notebooks /blob/master/matplotlib/matplotlib.ipynb

## matplotlib

Credits: Content forked from Parallel Machine Learning with scikit-learn and IPython by Olivier Grisel

- Setting Global Parameters
- Basic Plots
- Histograms
- Two Histograms on the Same Plot
- Scatter Plots

```
In [1]:  %matplotlib inline
         import pandas as pd
         import numpy as np
         import pylab as plt
         import seaborn
```

### Setting Global Parameters

```
In [2]:  # Set the global default size of matplotlib figures
         plt.rc('figure', figsize=(10, 5))

         # Set seaborn aesthetic parameters to defaults
         seaborn.set()
```

### Basic Plots

```
In [3]:  x = np.linspace(0, 2, 10)

         plt.plot(x, x, 'o-', label='linear')
         plt.plot(x, x ** 2, 'x-', label='quadratic')

         plt.legend(loc='best')
         plt.title('Linear vs Quadratic progression')
```

## Scikit Learn

About  Write  Help  Legal

Source: https://www.datacamp.com/community/blog/scikit-learn-cheat-sheet#gs.fZ2A1Jk