



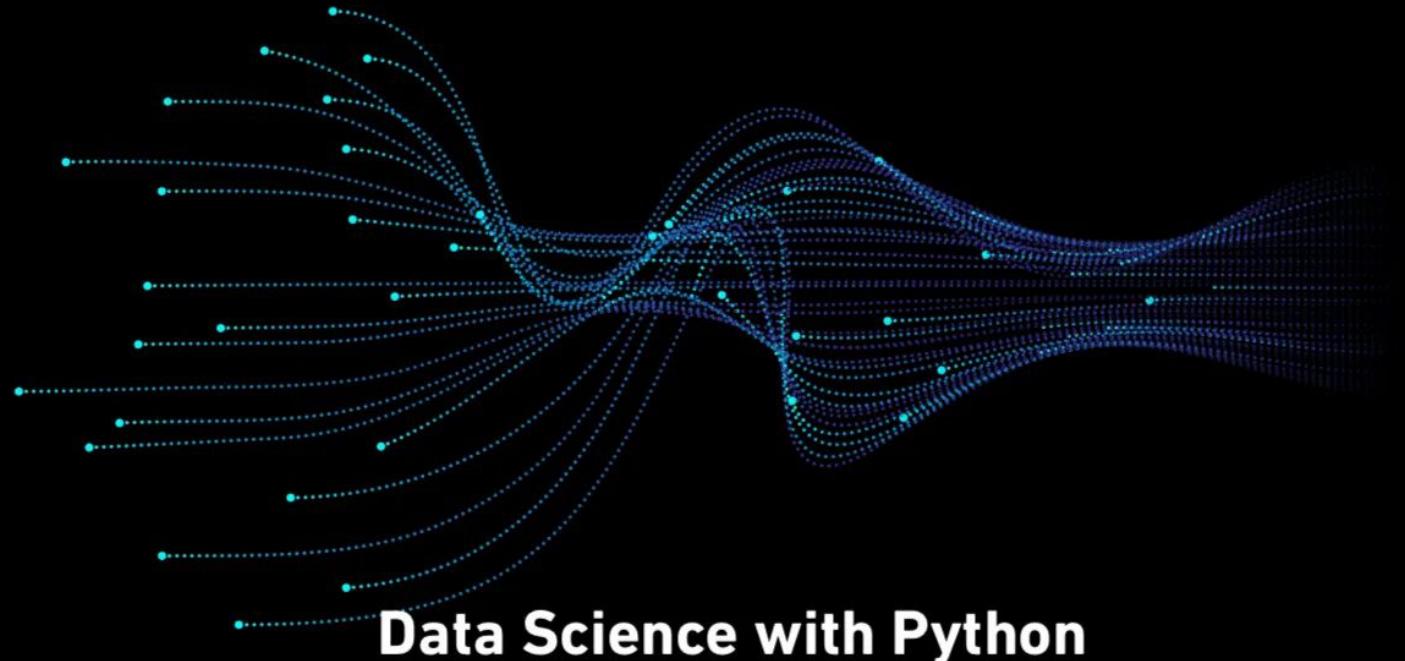
Cheat Sheets for AI

Neural Networks,
Machine Learning,
DeepLearning &
Big Data

**The Most Complete List
of Best AI Cheat Sheets**

BecomingHuman.AI

Table of Content



Data Science with Python

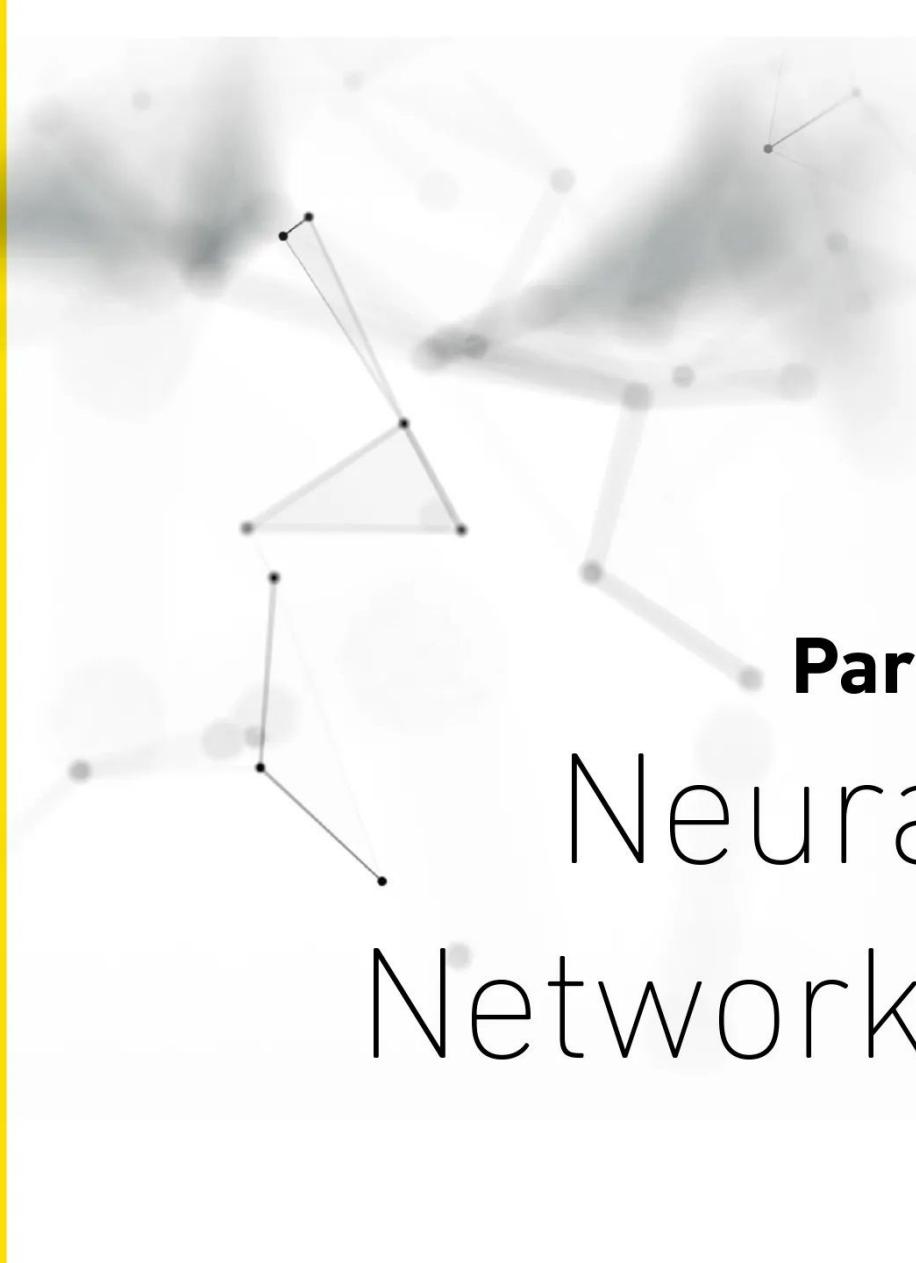
Neural Networks

- 03 Neural Networks Basics
- 04 Neural Network Graphs

Machine Learning

- 06 Machine Learning Basics
- 07 Scikit Learn with Python
- 08 Scikit Learn Algorithm
- 09 Choosing ML Algorithm

- | | | | |
|----|----------------|----|-----------------------------------|
| 11 | Tensor Flow | 17 | Pandas |
| 12 | Python Basics | 18 | Data Wrangling with Pandas |
| 13 | PySpark Basics | 19 | Data Wrangling with dplyr & tidyr |
| 14 | Numpy Basics | 20 | SciPi |
| 15 | Bokeh | 21 | MatPlotLib |
| 16 | Karas | 22 | Data Visualization with ggplot |
| 23 | Big-O | | |



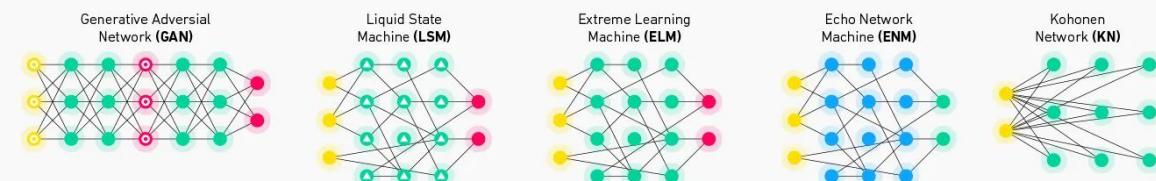
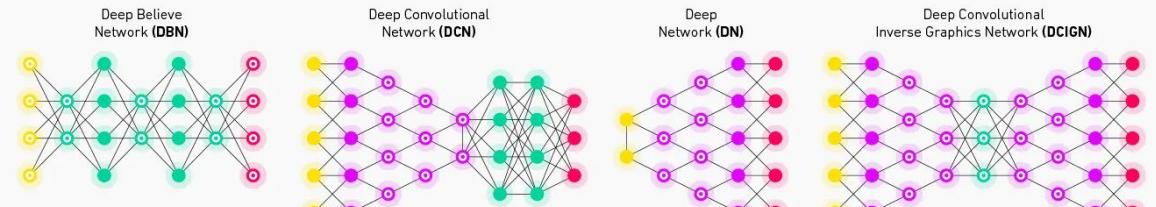
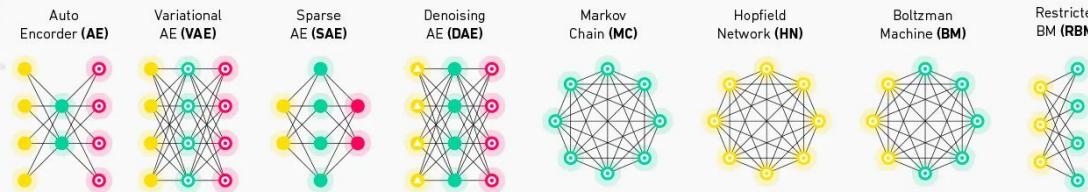
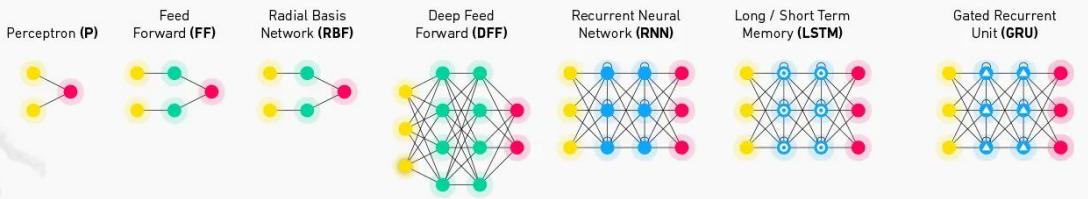
Part 1
Neural
Networks

Neural Networks Basic Cheat Sheet

BecomingHuman.AI

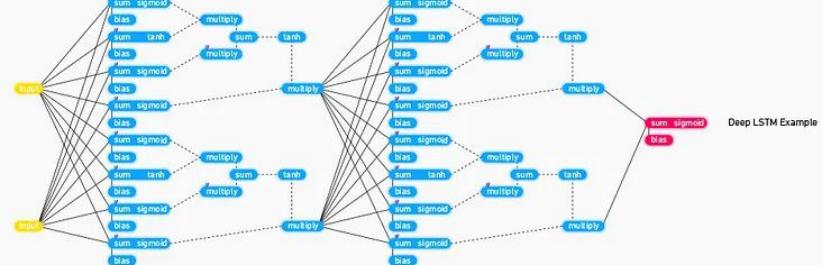
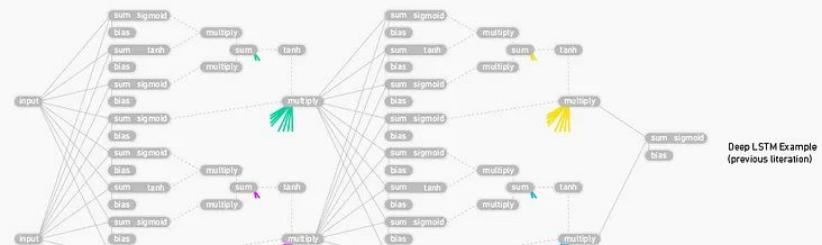
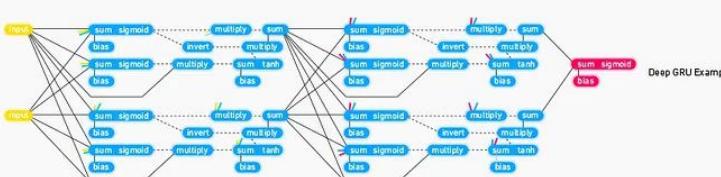
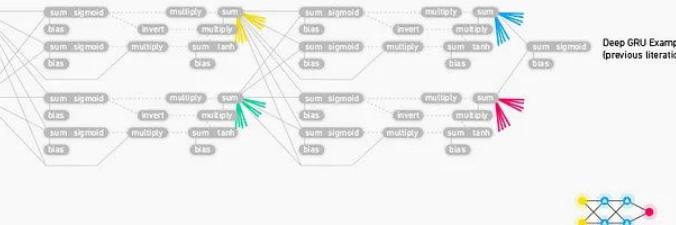
Index

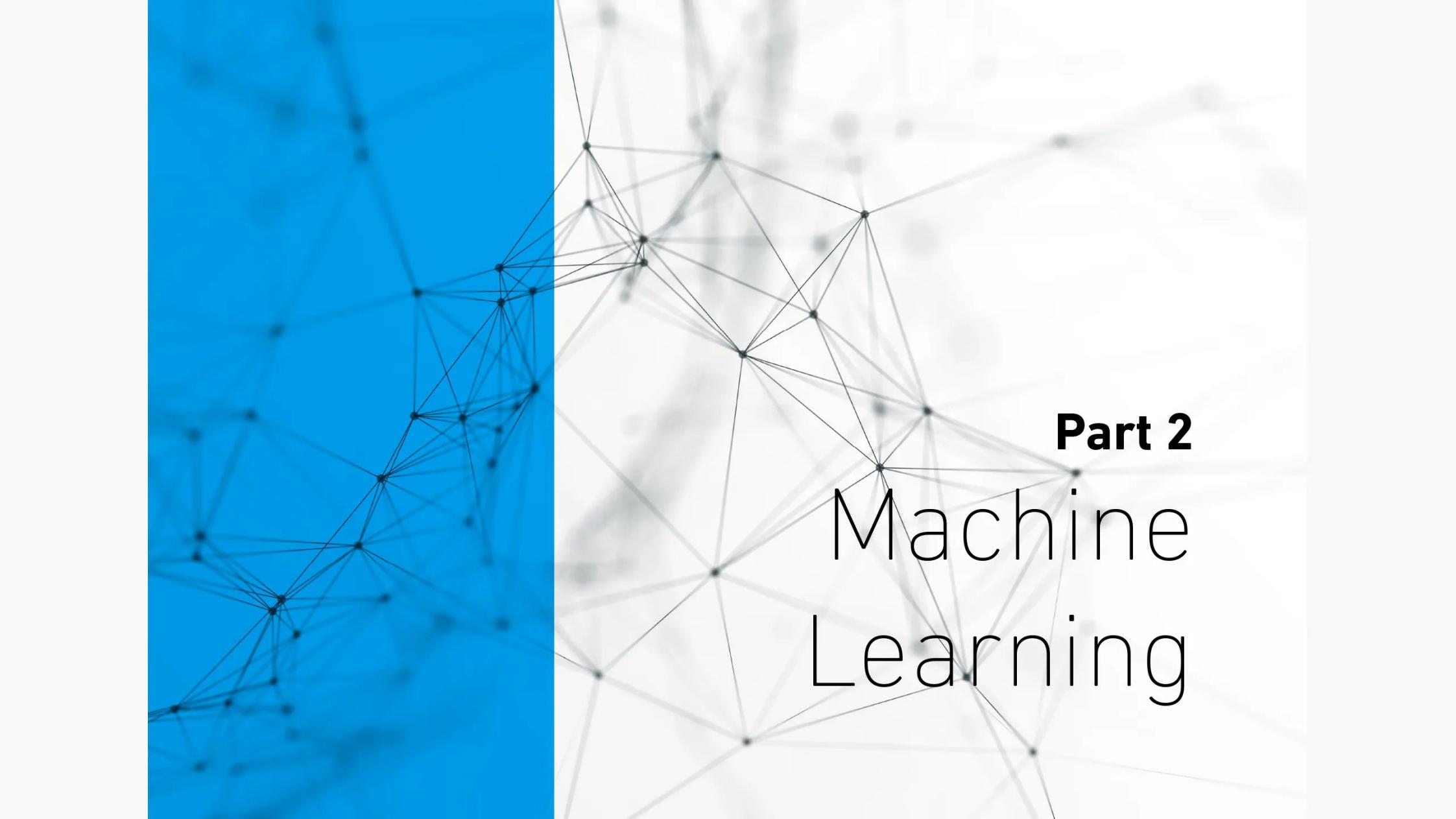
- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Different Memory Cell
- Kernel
- Convolutional or Pool



Neural Networks Graphs Cheat Sheet

BecomingHuman.AI



The background of the slide features a complex network graph composed of numerous small, semi-transparent black dots connected by thin gray lines, creating a sense of data points and connections. The left side of the background is a solid blue color, while the right side is white.

Part 2

Machine Learning

MachineLearning Overview

MACHINE LEARNING IN EMOJI

BecomingHuman.AI

SUPERVISED

UNSUPERVISED

REINFORCEMENT

BASIC REGRESSION

LINEAR

`linear_model.LinearRegression()`
Lots of numerical data



LOGISTIC

`linear_model.LogisticRegression()`
Target variable is categorical



CLUSTER ANALYSIS

K-MEANS

`cluster.KMeans()`
Similar datum into groups based on centroids



ANOMALY DETECTION

`covariance.EllipticalEnvelope()`
Finding outliers through grouping

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



FEATURE REDUCTION

T-DISTRIB STOCHASTIC NEIB EMBEDDING

`manifold.TSNE()`



Visual 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

`manifold.TSNE()`

SPECIFICITY FUNCTION

$TN / (FP+TN)$

SENSITIVITY FUNCTION

$TP / (TP+FN)$

Cheat-Sheet Skicit learn

Phyton For Data Science

BecomingHuman.AI



Skicit Learn

Skicit Learn is an open source Python library that implements a range of machine learning, processing, cross validation and visualization algorithm using a unified

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,test,y_train = train_test_split(X,y,random_state=33)
>>> scalar = preprocessing.StandardScaler().fit(X_train)
>>> X_train = scalar.transform(X_train)
>>> X_test = scalar.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)
```

Prediction

Supervised Estimators
>>> y_pred = svc.predict(random.random(2,5))
>>> y_pred = lr.predict(X)
>>> y_pred = knn.predict_proba(X_test)

Predict labels
Estimate probability of a label

Unsupervised Estimators

```
>>> y_pred = kmeans.predict(X_test)
```

Predict labels in clustering algos

Loading the Data

Your data needs to be numeric and stored as NumPy arrays or SciPy sparse matrix. other types that they 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([P,H,M,F,F,M,F,N,I,V,F,F])
>>> X[X < 0.7] = 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
>>> scalar = Normalizer().fit(X_train)
>>> normalized_X = scalar.transform(X_train)
>>> normalized_X_test = scalar.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 Imputer
>>> imp = Imputer(missing_values=0, strategy='mean', axis=0)
>>> imp.fit_transform(X_train)
```

Imputing Missing Values

```
>>> from sklearn.preprocessing import Imputer
>>> imp = Imputer(missing_values=0, strategy='mean', axis=0)
>>> imp.fit_transform(X_train)
```

Generating Polynomial Features

```
>>> from sklearn.preprocessing import PolynomialFeatures
>>> poly = PolynomialFeatures(5)
>>> poly.fit_transform(X)
```

Evaluate Your Model's Performance

Classification Metrics

Accuracy Score
>>> from sklearn.metrics import accuracy_score
>>> accuracy_score(y_test,y_pred)

Estimator score method
Metric scoring functions

Classification Report

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

Precision, recall, f1-score
and support

Confusion Matrix

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

Regression Metrics

Mean Absolute Error
>>> from sklearn.metrics import mean_absolute_error
>>> y_true = [3,-0.5,2]
>>> mean_absolute_error(y_true,y_pred)

Mean Squared Error

```
>>> from sklearn.metrics import mean_squared_error
>>> mean_squared_error(y_test,y_pred)
```

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))
```

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
>>> kmeans.fit(X_train)
>>> pca_model = pca.fit_transform(X_train)

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

Create Your Model

Supervised Learning Estimators

Linear Regression
>>> from sklearn.linear_model import LinearRegression
>>> lr = LinearRegression(normalize=True)

Support Vector Machines (SVM)
>>> from sklearn.svm import SVC
>>> svc = SVC(kernel='linear')

Naive Bayes
>>> from sklearn.naive_bayes import GaussianNB
>>> gnb = GaussianNB()

KNN
>>> from sklearn import neighbors
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)

Unsupervised Learning Estimators

Principal Component Analysis (PCA)
>>> from sklearn.decomposition import PCA
>>> pca = PCA(n_components=0.95)

K Means
>>> from sklearn.cluster import KMeans
>>> k_means = KMeans(n_clusters=3, random_state=0)

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)
```

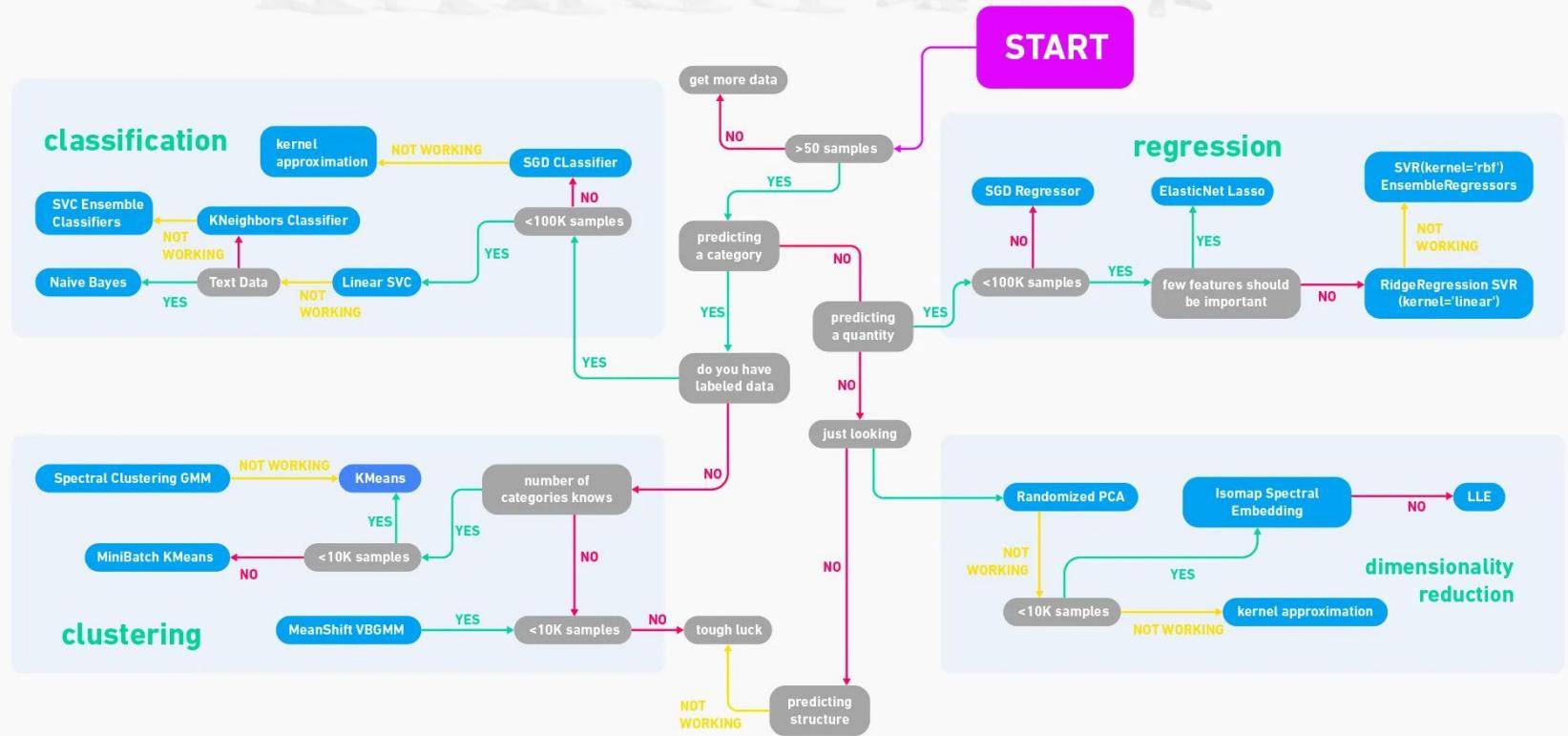
Tune Your Model

Grid Search
>>> from sklearn.grid_search import GridSearchCV
>>> params = {'n_neighbors': np.arange(1,5),
... '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(knn,
... param_distributions=params,
... cv=4,
... n_iter=8,
... random_state=5)
>>> rsearch.fit(X_train,y_train)
>>> print(rsearch.best_score_)

Skicit-learn Algorithm

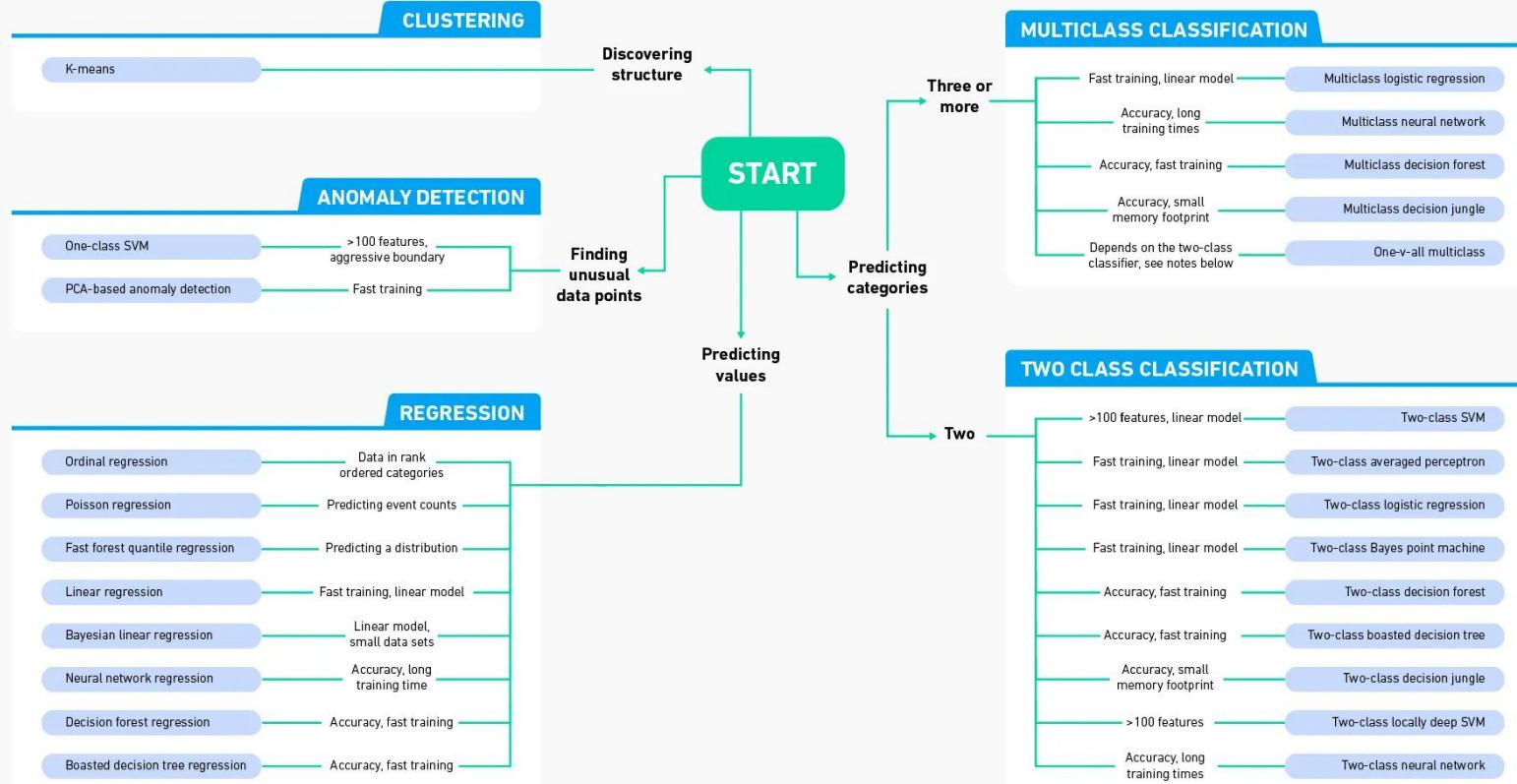
BecomingHuman.AI



Algorithm Cheat Sheet

BecomingHuman.AI

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.



The background of the slide features a complex network graph composed of numerous small, semi-transparent grey dots connected by thin grey lines, forming a web-like structure. This pattern is divided vertically down the center. The left half of the graph is set against a solid orange background, while the right half is set against a white background.

Part 3

Data Science with Python

Tensor Flow Cheat Sheet

BecomingHuman.AI



In May 2017 Google announced the second-generation of the TPU, as well as the availability of the TPUs in Google Compute Engine.[12] The second-generation TPUs deliver up to 180 teraflops of performance, and when organized into clusters of 64 TPUs provide up to 11.5 petaflops.

Info

TensorFlow

TensorFlow™ is an open source software library created by Google for numerical computation and large scale computation. Tensorflow bundles together Machine Learning, Deep learning models and frameworks and makes them useful by way of common metaphor.

Keras

Keras is an open sourced neural networks library, written in Python and is built for fast experimentation via deep neural networks and modular design. It is capable of running on top of TensorFlow, Theano, Microsoft Cognitive Toolkit, or PlaidML.

Skflow

Scikit Flow is a high level interface base on tensorflow which can be used like sklearn. You can build your own model on your own data quickly without rewriting extra code.provides a set of high level model classes that you can use to easily integrate with your existing Scikit-learn pipeline code.

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
sudo pip install
```

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
```

Tensor Flow

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 (Classifier)
batch_size=32,
steps=200, // except
TensorFlowRNNClassifier - there is 50
optimizer='Adagrad',
learning_rate=0.1,

Each class has a method fit

```
fit(X, y, monitor=None, logdir=None)
X: matrix or tensor of shape [n_samples, n_features...]. Can be iterator that returns arrays of features. The training input samples for fitting the model.
Y: vector or matrix [n_samples] or [n_samples, n_outputs]. Can be iterator that returns array of targets. The training target values (class labels in classification, real numbers in regression).
monitor: Monitor object to print training progress and invoke early stopping
logdir: the directory to save the log file that can be used for optional visualization.
predict (X, axis=1, batch_size=None)
Args:
X: array-like matrix, [n_samples, n_features...] or iterator.
axis: Which axis to argmax for classification.
By default axis 1 (next after batch) is used. Use 2 for sequence predictions.
batch_size: If test set is too big, use batch size to split it into mini batches. By default the batch_size member variable is used.
Returns:
y: array of shape [n_samples]. The predicted classes or predicted value.
```

Phyton For Data Science

Cheat-Sheet Phyton Basic

BecomingHuman.AI



Variables and Data Types

Variable Assignment

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

Calculations With Variables

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

Calculations With Variables

str0	'5', '3.45', 'True'	Variables to strings
int0	5, 3, 1	Variables to integers
float0	5.0, 1.0	Variables to floats
bool0	True, True, True	Variables to booleans

Asking For Help

```
>>> help(str)
```

Lists

Also see NumPy Arrays

Selecting List Elements

Index starts at 0

Subset

```
>>> my_list[1]
```

```
>>> my_list[-3]
```

Slice

```
>>> my_list[1:3]
```

```
>>> my_list[1:]
```

```
>>> my_list[:3]
```

```
>>> my_list[::]
```

Subset Lists of Lists

```
>>> my_list[2][1:0]
```

```
>>> my_list[2][1:2]
```

Select item at index 1

Select 3rd last item

Select items at index 1 and 2

Select items after index 0

Select items before index 3

Copy my_list

my_list[1][1:0]

my_list[1][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('t')
```

Append an Item at a time

```
>>> my_list.remove('t')
```

Remove an Item

```
>>> del(my_list[0:1])
```

Remove an Item

```
>>> my_list.reverse()
```

Reverse the list

```
>>> my_list.extend('t')
```

Append an Item

```
>>> my_list.pop(-1)
```

Remove an Item

```
>>> my_list.insert(0,'t')
```

Insert an Item

```
>>> my_list.sort()
```

Sort the list

Numpy Arrays

Also see Lists

Selecting Numpy Array Elements

Index starts at 0

Subset

```
>>> my_list = [1, 2, 3, 4]
```

```
>>> my_array = np.array(my_list)
```

```
>>> my_2darray =
```

```
np.array([[1,2,3],[4,5,6]])
```

Select item at index 1

Select items at index 0 and 1

Select items before index 3

Copy my_list

```
my_list[1][1:0]
```

```
my_list[1][1:2]
```

Numpy Array Operations

```
>>> my_array = 3
```

```
array([False, False, False], dtype=bool)
```

```
>>> my_array * 2
```

```
array([0, 0, 0], dtype=bool)
```

```
>>> my_array + np.array([5, 6, 7, 8])
```

```
array([6, 10, 12])
```

Numpy Array Operations

```
>>> 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.correlate()
```

Correlation coefficient

```
>>> np.std(my_array)
```

Standard deviation

Strings

Also see NumPy Arrays

String Operations

```
>>> my_string = 'thisStringIsAwesome'
```

```
>>> my_string
```

```
'thisStringIsAwesome'
```

String Operations

```
>>> my_string * 2
```

```
'thisStringIsAwesomethisStringIsAwesome'
```

```
>>> my_string * 1000
```

```
'thisStringIsAwesomethisStringIsAwesome...'
```

String Methods

String to uppercase

String to lowercase

Count String elements

Replace String elements

Strip whitespaces

Libraries

Import libraries

```
>>> import numpy
```

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

```
Selective Import
```

```
>>> from math import pi
```

Install Python



Leading open data science platform
powered by Python



Free IDE that is included
with Anaconda



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

Python For Data Science Cheat Sheet

PySpark - RDD Basics

BecomingHuman.AI



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]")
```

Calculations With Variables

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

Configuration

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

Configuration

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),(b,1),(b,1)])  
>>> rdd3 = sc.parallelize(range(100))  
>>> rdd4 = sc.parallelize([(a,"x"),(c,"y"),  
                         (b,"p"),("r")])
```

External Data

Read either one text file from HDFS, a local file system or any Hadoop-supported file system URI with `textFile()`, or read in a directory of text files with `wholeTextFiles()`.

```
>>> textFile = sc.textFile("/mydirectory/*txt")  
>>> textFile2 = sc.wholeTextFiles("/mydirectory/")
```

Selecting Data

Getting

```
>>> rdd.collect()      Return a list with all RDD elements  
[(a, 7), (a, 2), (b, 2)]  
>>> rdd.take(2)       Take first 2 RDD elements  
[(a, 7), (a, 2)]  
>>> rdd.first()       Take first RDD element  
(a, 7)  
>>> rdd.top(2)       Take top 2 RDD elements  
(b, 2), (a, 7)
```

Sampling

```
>>> rdd3.sample(False, 0.15, 81).collect()  Return sampled subset of rdd3  
[3,4,27,31,40,41,42,43,60,76,79,80,86,97]
```

Filtering

```
>>> rdd.filter(lambda x: "a" in x).collect()  Filter the RDD  
[(a,7),(a,2)]  
>>> rdd5.distinct().collect()  Return distinct RDD values  
[a,2,b]  
>>> rdd.keys().collect()  Return (key,value) RDD's keys  
[a, a, b]
```

Iterating

Getting

```
>>> def g(x): print(x)  
>>> rdd.foreach(g)  
[a, 7]  
(b, 2)  
(a, 2)
```

Retrieving RDD Information

Basic Information

```
>>> rdd.getNumPartitions()  List the number of partitions  
3  
>>> rdd.count()           Count RDD instances  
3  
>>> rdd.collectByKey()   Count RDD instances by key  
defaultdict<type 'int'>[(a,2),(b,1)]  
>>> rdd.countByValue()   Count RDD instances by value  
defaultdict<type 'int'>[(c,1),(a,2),(b,1),(a,7)]  
>>> rdd.collectAsMap()  Return (key,value) pairs as a dictionary  
{a: 2, b: 1}  
>>> rdd3.sum()            Sum of RDD elements  
4950  
>>> sc.parallelize([]).isEmpty()  Sum of RDD elements  
True  
>>> sc.parallelize([]).isNotEmpty()  Check whether RDD is empty  
False
```

Summary

```
>>> rdd3.max()           Maximum value of RDD elements  
99  
>>> rdd3.min()           Minimum value of RDD elements  
0  
>>> rdd3.mean()          Mean value of RDD elements  
49.5  
>>> rdd3.stdev()         Standard deviation of RDD elements  
28.8660047722118  
>>> rdd3.variance()      Compute variance of RDD elements  
833.25  
>>> rdd3.histogram(3)    Compute histogram by bins  
((0,33.66,99],[33,33,34])  
>>> rdd3.stats()         Summary statistics (count, mean, stdev, max & min)
```

Applying Functions

```
>>> rdd.map(lambda x: x+(x[1]*x[0]))  Apply a function to each RDD element  
[(a,7),(a,2),(a,b)]  
>>> rdd5 = rdd.flatMap(lambda x: x*(x[1]*x[0]))  Apply a function to each RDD element and flatten the result  
[(a,7),(a,a,2,a,b,2,b)]  
>>> rdd4 = rdd.flatMap(lambda x: x)  Apply a flatMap function to each (key,value) pair of rdd4 without changing the keys  
[(a),(a),(a),(a),(b),(b),(b)]
```

Mathematical Operations

```
>>> rdd.subtract(rdd2)  Return each rdd value not contained  
[(b,2),(a,7)]  
>>> rdd2.subtractByKey(rdd)  Return each (key,value) pair of rdd2 with no matching key in rdd  
[(a,1)]  
>>> rdd.cartesian(rdd2).collect()  Return the Cartesian product of rdd and rdd2
```

Sort

```
>>> rdd2.sortBy(lambda x: x[1]).collect()  Sort RDD by given function  
[(d,1),(b,1),(c,2)]  
>>> rdd2.partitionBy(2).sortByKey().collect()  RDD by key  
[(a,2),(b,1),(c,1)]
```

Reshaping Data

Reducing

```
>>> rdd.reduceByKey(lambda xy : x+y)  Merge the rdd values for each key  
[(a,9),(b,2)]  
>>> rdd.reduceByKey(lambda a, b: a + b)  Merge the rdd values  
(a,7),(b,2),(b,2)
```

Grouping by

```
>>> rdd3.groupByKey().mapValues(list).collect()  Return RDD of grouped values  
[(a,[7,2]),(b,[2])]  
>>> rdd3.groupByKey().mapValues(list).collect()  Group rdd by key  
[(a,[7,2]),(b,[2])]
```

Aggregating

```
>>> seqOp = (lambda x: (x[0]+x[1])*1)  Aggregate RDD elements of each partition and then the results  
>>> combOp = (lambda x,y:(x[0]+y[0],x[1]+y[1]))  
>>> rdd3.aggregate((0,0),seqOp,combOp)  Aggregate values of each RDD key  
(4950,100)  
>>> rdd3.aggregateByKey((0,0),seqOp,combOp)  Aggregate the elements of each partition, and then the results  
[(a,(9,2)),(b,(2,1))]  Merge the values for each key  
4950  
>>> rdd3.foldByKey(0,add)  Create tuples of RDD elements by applying a function  
[(a,9),(b,2)]  
>>> rdd3.keyBy(lambda x: x*x).collect()  Create tuples of RDD elements by applying a function
```

Reshaping Data

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

Saving

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

Stopping SparkContext

```
>>> sc.stop()
```

Execution

```
$ ./bin/spark-submit examples/src/main/python/pi.py
```




Bokeh Cheat Sheet

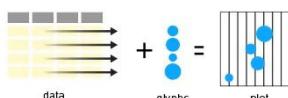
BecomingHuman.AI



Data Types

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] step 1
>>> y = [6, 7, 8, 4, 5] step 2
>>> p = figure(title='simple line example', step 3
              x_axis_label='x',
              y_axis_label='y')
>>> p.line(x, y, legend='Temp', line_width=2) step 4
>>> output_file('lines.html') step 5
>>> show(p) step 6
```

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, 65, 'US'],
                               [32.4, 66, 'Asia'],
                               [31.4, 109, 'Europe']]),
                               columns=['mpg','cyl','origin'],
                               index=['Toyota', 'Fiat', 'Volvo'])

>>> from bokeh.models import ColumnDataSource
>>> cds_df = ColumnDataSource(df)
```

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()
```

Show or Save Your Plots

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

Renderers & Visual Customizations

Glyphs



Scatter Markers
>>> p1.circle(array([1,2,3]), np.array([3,2,1]),
 fill_color='white')



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]])
```

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')
```

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()
```

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)
```

Customized Glyphs



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



Hover Glyphs
>>> hover = HoverTool(tooltips=None, mode='vline')
>>> p3.add_tools(hover)



Colormapping
>>> color_mapper = CategoricalColorMapper(
 factors=['US', 'Asia', 'Europe'],
 palette=['blue', 'red', 'green'])
>>> p3.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')
```

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])
```

Legend Orientation

```
>>> p.legend.orientation = 'horizontal'
>>> p.legend.orientation = 'vertical'
```

Legend Background & Border

```
>>> p.legend.border_line_color = "navy"
>>> p.legend.background_fill_color = "white"
```

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’])

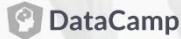
Box Plot
>>> from bokeh.charts import BoxPlot
>>> p = BoxPlot(df, values=‘vals’, label=‘cyl’, legend=‘bottom_right’)

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’,

Keras Cheat Sheet

BecomingHuman.AI



K 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])
```

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
>>> model.fit(data,labels,epochs=10,batch_size=32)
>>> predictions = model.predict(data)
```

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.csv'),delimiter=',')
>>> X = data[:,0:8]
>>> y = data[:,8]
```

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,
>                  activation='relu',
>                  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(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),padding='same'))
>>> 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 output shape
>>> model.summary()
Model summary representation
>>> model.get_config()
Model configuration
>>> model.get_weights()
List all weight tensors in the model
```

Prediction

```
>>> model3.predict(x_test4,batch_size=32)
>>> model3.predict_classes(x_test4,batch_size=32)
```

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])
```

Compile Model

MLP: Binary Classification

```
>>> model.compile(optimizer='adam',
>                  loss='binary_crossentropy',
>                  metrics=[accuracy])
```

MLP: Multi-Class Classification

```
>>> model.compile(optimizer='rmsprop',
>                  loss='categorical_crossentropy',
>                  metrics=[accuracy])
```

MLP: Regression

```
>>> model.compile(optimizer='rmsprop',
>                  loss='mse',
>                  metrics=[mae])
```

Recurrent Neural Network

```
>>> model3.compile(loss='binary_crossentropy',
>                  optimizer='adam',
>                  metrics=[accuracy])
```

Save/ Reload Models

```
>>> from keras.models import load_model
>>> model3.save('model_file.h5')
>>> my_model = load_model('my_model.h5')
```

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)
```

Preprocessing

Sequence Padding

```
>>> from keras.preprocessing import sequence
>>> x_train = sequence.pad_sequences(x_train4,maxlen=80)
>>> x_test = 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)
```

Train and Test Sets

```
>>> from sklearn.model_selection import train_test_split
>>> X_train,X_test,y_train,y_test = train_test_split(
>             ...,
>             test_size=0.33,
>             random_state=42)
```

Standardization/Normalization

```
>>> from sklearn.preprocessing import StandardScaler
>>> scaler = StandardScaler().fit(x_train)
>>> standardized_X = scaler.transform(x_train2)
>>> standardized_X_test = scaler.transform(x_test2)
```

Pandas Basics

Cheat Sheet

BecomingHuman.AI



Use the following import convention: >>> import pandas as pd

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

Pandas Data Structures

Series

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

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

Data Frame

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

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

Dropping

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

Sort & Rank

```
>>> df.sort_index()             Sort by labels along an axis
>>> df.sort_values(by='Country') Sort by the values along an axis
>>> df.rank()                  Assign ranks to entries
```

Retrieving Series/ DataFrame 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
```

Selection

Also see NumPy Arrays

Getting

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

Selecting, Boolean Indexing & Setting

By Position

```
>>> df.loc[0,0]               Select single value by row & column
'Belgium'
>>> df.loc[0, ['Country']]    Select single value by row & column labels
'Belgium'
>>> df.at[0, 'Country']      Belgium
```

By Label/Position

```
>>> df.ix[2]                 Select single row of subset of rows
2   Country      Brazil   Population
   India      New Delhi  1303171035
   Brazil     Brasilia  207847528
>>> df.ix[1, 'Capital']     Select a single column of subset of columns
0 Brussels
1 New Delhi
2 Brasilia
>>> df.ix[1, 'Capital']     Select rows and columns
'New Delhi'
```

Boolean Indexing

```
>>> s[~(s > 1)]            Series s where value is not > 1
>>> s[(s < 1) | (s > 2)]  s where value is < 1 or > 2
>>> df[df['Population']>1200000000] Use filter to adjust DataFrame
```

Setting

```
>>> s['a'] = 6               Set Index a of Series s to 6
```

I/O

Read and Write to CSV

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

Read and Write to Excel

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

Asking For Help

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

Applying Functions

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

Data Alignment

Internal Data Alignment

NaN values are introduced in the indices that don't overlap:

```
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])
>>> s3 + 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)
b 5.0
c 5.0
d 7.0
>>> s.sub(s3, fill_value=2)
>>> s.div(s3, fill_value=4)
```

Pandas

Cheat Sheet

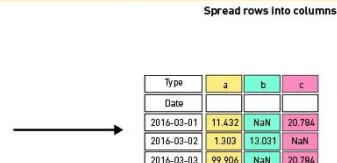
BecomingHuman.AI

Pandas Data Structures

Pivot

```
>>> df3 = df2.pivot(index='Date',
                   columns='Type',
                   values='Value')
```

	Date	Type	Value
0	2016-03-01	a	11.432
1	2016-03-02	b	12.031
2	2016-03-01	c	20.784
3	2016-03-03	a	99.906
4	2016-03-02	a	1.203
5	2016-03-03	c	20.784

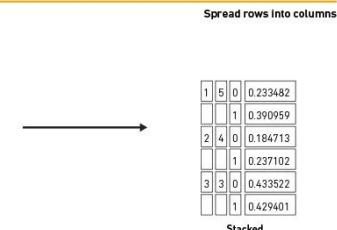


Pivot Table

```
>>> df4 = pd.pivot_table(df2,
                        values='Value',
                        index='Date',
                        columns='Type')
```

	0	1
1	5	0.233482
2	4	0.184713
3	3	0.433522

Unstacked

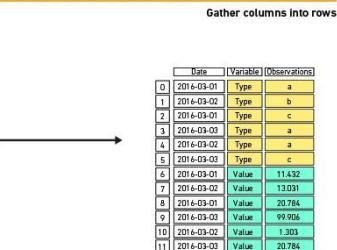


Stacked

Melt

```
>>> pd.melt(df2,
            id_vars=['Date'],
            value_vars=['Type', 'Value'],
            value_name='Observations')
```

	Date	Type	Value
0	2016-03-01	a	11.432
1	2016-03-02	b	13.031
2	2016-03-01	c	20.784
3	2016-03-03	a	99.906
4	2016-03-02	a	1.203
5	2016-03-03	c	20.784



Advanced Indexing

Also see NumPy Arrays

Selecting

```
>>> df3.loc[:,(df3>1).any()]
>>> df3.loc[:,(df3>1).all()]
>>> df3.loc[:,df3.isnull().any()]
>>> df3.loc[:,df3.notnull().all()]
Indexing With isin
>>> df[(df['Country'].isin(df['Type']))]
>>> df.filter(items=['a','b'])
>>> df.select(lambda x: not x%5)
```

Where

```
>>> s.where(s > 0)
Query
>>> df6.query('second > first')
```

Setting/Resetting Index

```
>>> df.set_index('Country')
>>> df4 = df.reset_index()
>>> df = df.rename(index=str,
                  columns={'Country':'cntry',
                           'Capital':'cptl',
                           'Population':'ppln'})
```

Select cols with any vals >1
Select cols with val > 1
Select cols with NaN
Select cols without NaN

Find same elements
Filter on values
Select specific elements

Subset the data

Query DataFrame

Set the index
Reset the index
Rename DataFrame

Reindexing

```
>>> s2 = s.reindex(['a','d','b','c'])
```

Forward Filling

```
>>> df.reindex(range(4),
               method='ffill')
Country Capital Population
0 Belgium Brussels 11190846
1 India New Delhi 1303171035
2 Brazil Brasilia 207847528
3 Brazil Brasilia 207847528
```

Forward Filling

```
>>> s3 = s.reindex(range(5),
                   method='ffill')
0 3
1 3
2 3
3 3
4 3
```

Multilabeling

```
>>> arrays = [np.array([1,2,3]),
             np.array([4,5,3])]
>>> df5 = pd.DataFrame(np.random.rand(3, 2), index=arrays)
>>> tuples = list(zip(*arrays))
>>> index = pd.MultiIndex.from_tuples(tuples,
                                      names=['first', 'second'])
>>> df6 = pd.DataFrame(np.random.rand(3, 2), index=index)
>>> df2.set_index(['Date', 'Type'])
```

Duplicate Data

```
>>> s3.unique()
>>> df2.duplicated(['Type'])
>>> df2.drop_duplicates(['Type'], keep='last')
>>> df2.index.duplicated()
```

Return unique values
Check duplicates
Drop duplicates
Drop duplicates

Grouping Data

Aggregation

```
>>> df2.groupby(by=['Date','Type']).mean()
>>> df4.groupby(level=0).sum()
>>> df4.groupby(level=0).agg('a':lambda x: x.sum() / len(x), 'b': np.sum())
```

Transformation

```
>>> customSum = lambda x: (x*x)%2
>>> df4.groupby(level=0).transform(customSum)
```

Check duplicates

Drop duplicates

Drop duplicates

Missing Data

```
>>> df.dropna()
>>> df3.fillna(df3.mean())
>>> df2.replace('a', 'f')
```

Drop NaN values
Fill NaN values with a predetermined value
Replace values with others

Combining Data

data1	data2
X1	X2
a	11.432
b	1.303
c	99.906
	NaN

Pivot
>>> pd.merge(data1, data2, how='left', on='X1')
a 11.432 20.784
b 1.303 NaN
c 99.906 NaN

Pivot
>>> pd.merge(data1, data2, how='right', on='X1')
a 11.432 20.784
b 1.303 NaN
d NaN 20.784

Pivot
>>> pd.merge(data1, data2, how='inner', on='X1')
a 11.432 20.784
b 1.303 NaN

Join
>>> data1.join(data2, how='right')

Concatenate

Vertical

```
>>> s.append(s2)
```

Horizontal/Vertical

```
>>> pd.concat([s2], axis=1, keys=[0|1|2|3])
>>> pd.concat([data1, data2], axis=1, join='inner')
```

Dates

```
>>> df2['Date'] = pd.to_datetime(df2['Date'])
>>> df2['Date'] = pd.date_range('2000-1-1', periods=6,
                                freq='M')
>>> dates = [datetime(2012,5,1), datetime(2012,5,2)]
>>> index = pd.DatetimeIndex(dates)
>>> index = pd.date_range(datetime(2012,2,1), end, freq='BM')
```

Visualization

```
>>> import matplotlib.pyplot as plt
>>> s.plot()
>>> df2.plot()
>>> plt.show()
```

Data Wrangling with pandas Cheat Sheet

BecomingHuman.AI

Syntax Creating DataFrames

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

```
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=['b', 'c'])
```

Specify values for each row.

n	v	a	b	c
d	1	4	7	10
d	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')
     )
```

Windows

`df.exploding()`
Return an Exploding 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.

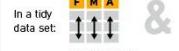
Windows

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

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



Tidy Data A foundation for wrangling in pandas

In a tidy data set:

 Each variable is saved in its own column
 Each observation is saved in its own row

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

 M * A

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.


`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(columns=['Length','Height'])`
Drop columns from DataFrame

Subset Observations (Rows)

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

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

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

`df.loc[10:20]`
Select rows by position.

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

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

`df.sample(frac=0.5)`
Select multiple columns with specific names.

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

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

`Logic in Python (and pandas)`

`~` Matches strings containing a period

`Length` Matches strings ending with word Length

`^Sep$` Matches strings beginning with the word Sep

`^([1-5]{1})` Matches strings beginning with 'x' and ending with 1,2,3,4

`^(?P<spec>){n}` Matches strings except the string Spec

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

`df.loc[:,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.

`df.rolling(3).sum()` Return rolling window sum over axis=0.

`df.rolling(3).mean()` Return rolling window mean over axis=0.

`df.rolling(3).std()` Return rolling window standard deviation over axis=0.

`df.rolling(3).min()` Return rolling window minimum over axis=0.

`df.rolling(3).max()` Return rolling window maximum over axis=0.

`df.rolling(3).count()` Return rolling window count over axis=0.

`df.rolling(3).sum().shift(1)` Shifted rolling window sum over axis=0.

`df.rolling(3).mean().shift(1)` Shifted rolling window mean over axis=0.

`df.rolling(3).std().shift(1)` Shifted rolling window standard deviation over axis=0.

`df.rolling(3).min().shift(1)` Shifted rolling window minimum over axis=0.

`df.rolling(3).max().shift(1)` Shifted rolling window maximum over axis=0.

`df.rolling(3).count().shift(1)` Shifted rolling window count over axis=0.

`df.rolling(3).sum().shift(-1)` Shifted rolling window sum over axis=0.

`df.rolling(3).mean().shift(-1)` Shifted rolling window mean over axis=0.

`df.rolling(3).std().shift(-1)` Shifted rolling window standard deviation over axis=0.

`df.rolling(3).min().shift(-1)` Shifted rolling window minimum over axis=0.

`df.rolling(3).max().shift(-1)` Shifted rolling window maximum over axis=0.

`df.rolling(3).count().shift(-1)` Shifted rolling window count over axis=0.

`df.rolling(3).sum().dropna()` Drop missing values from sum.

`df.rolling(3).mean().dropna()` Drop missing values from mean.

`df.rolling(3).std().dropna()` Drop missing values from standard deviation.

`df.rolling(3).min().dropna()` Drop missing values from minimum.

`df.rolling(3).max().dropna()` Drop missing values from maximum.

`df.rolling(3).count().dropna()` Drop missing values from count.

`df.rolling(3).sum().dropna(how='any')` Drop missing values from sum.

`df.rolling(3).mean().dropna(how='any')` Drop missing values from mean.

`df.rolling(3).std().dropna(how='any')` Drop missing values from standard deviation.

`df.rolling(3).min().dropna(how='any')` Drop missing values from minimum.

`df.rolling(3).max().dropna(how='any')` Drop missing values from maximum.

`df.rolling(3).count().dropna(how='any')` Drop missing values from count.

`df.rolling(3).sum().dropna(how='all')` Drop missing values from sum.

`df.rolling(3).mean().dropna(how='all')` Drop missing values from mean.

`df.rolling(3).std().dropna(how='all')` Drop missing values from standard deviation.

`df.rolling(3).min().dropna(how='all')` Drop missing values from minimum.

`df.rolling(3).max().dropna(how='all')` Drop missing values from maximum.

`df.rolling(3).count().dropna(how='all')` Drop missing values from count.

`df.rolling(3).sum().dropna(how='ffill')` Fill missing values with previous value.

`df.rolling(3).mean().dropna(how='ffill')` Fill missing values with previous mean.

`df.rolling(3).std().dropna(how='ffill')` Fill missing values with previous standard deviation.

`df.rolling(3).min().dropna(how='ffill')` Fill missing values with previous minimum.

`df.rolling(3).max().dropna(how='ffill')` Fill missing values with previous maximum.

`df.rolling(3).count().dropna(how='ffill')` Fill missing values with previous count.

`df.rolling(3).sum().dropna(how='bfill')` Fill missing values with next value.

`df.rolling(3).mean().dropna(how='bfill')` Fill missing values with next mean.

`df.rolling(3).std().dropna(how='bfill')` Fill missing values with next standard deviation.

`df.rolling(3).min().dropna(how='bfill')` Fill missing values with next minimum.

`df.rolling(3).max().dropna(how='bfill')` Fill missing values with next maximum.

`df.rolling(3).count().dropna(how='bfill')` Fill missing values with next count.

`df.rolling(3).sum().dropna(how='ffill', limit=1)` Fill missing values with previous value, limit 1.

`df.rolling(3).mean().dropna(how='ffill', limit=1)` Fill missing values with previous mean, limit 1.

`df.rolling(3).std().dropna(how='ffill', limit=1)` Fill missing values with previous standard deviation, limit 1.

`df.rolling(3).min().dropna(how='ffill', limit=1)` Fill missing values with previous minimum, limit 1.

`df.rolling(3).max().dropna(how='ffill', limit=1)` Fill missing values with previous maximum, limit 1.

`df.rolling(3).count().dropna(how='ffill', limit=1)` Fill missing values with previous count, limit 1.

`df.rolling(3).sum().dropna(how='bfill', limit=1)` Fill missing values with next value, limit 1.

`df.rolling(3).mean().dropna(how='bfill', limit=1)` Fill missing values with next mean, limit 1.

`df.rolling(3).std().dropna(how='bfill', limit=1)` Fill missing values with next standard deviation, limit 1.

`df.rolling(3).min().dropna(how='bfill', limit=1)` Fill missing values with next minimum, limit 1.

`df.rolling(3).max().dropna(how='bfill', limit=1)` Fill missing values with next maximum, limit 1.

`df.rolling(3).count().dropna(how='bfill', limit=1)` Fill missing values with next count, limit 1.

`df.rolling(3).sum().dropna(how='ffill', limit=2)` Fill missing values with previous value, limit 2.

`df.rolling(3).mean().dropna(how='ffill', limit=2)` Fill missing values with previous mean, limit 2.

`df.rolling(3).std().dropna(how='ffill', limit=2)` Fill missing values with previous standard deviation, limit 2.

`df.rolling(3).min().dropna(how='ffill', limit=2)` Fill missing values with previous minimum, limit 2.

`df.rolling(3).max().dropna(how='ffill', limit=2)` Fill missing values with previous maximum, limit 2.

`df.rolling(3).count().dropna(how='ffill', limit=2)` Fill missing values with previous count, limit 2.

`df.rolling(3).sum().dropna(how='bfill', limit=2)` Fill missing values with next value, limit 2.

`df.rolling(3).mean().dropna(how='bfill', limit=2)` Fill missing values with next mean, limit 2.

`df.rolling(3).std().dropna(how='bfill', limit=2)` Fill missing values with next standard deviation, limit 2.

`df.rolling(3).min().dropna(how='bfill', limit=2)` Fill missing values with next minimum, limit 2.

`df.rolling(3).max().dropna(how='bfill', limit=2)` Fill missing values with next maximum, limit 2.

`df.rolling(3).count().dropna(how='bfill', limit=2)` Fill missing values with next count, limit 2.

`df.rolling(3).sum().dropna(how='ffill', limit=3)` Fill missing values with previous value, limit 3.

`df.rolling(3).mean().dropna(how='ffill', limit=3)` Fill missing values with previous mean, limit 3.

`df.rolling(3).std().dropna(how='ffill', limit=3)` Fill missing values with previous standard deviation, limit 3.

`df.rolling(3).min().dropna(how='ffill', limit=3)` Fill missing values with previous minimum, limit 3.

`df.rolling(3).max().dropna(how='ffill', limit=3)` Fill missing values with previous maximum, limit 3.

`df.rolling(3).count().dropna(how='ffill', limit=3)` Fill missing values with previous count, limit 3.

`df.rolling(3).sum().dropna(how='bfill', limit=3)` Fill missing values with next value, limit 3.

`df.rolling(3).mean().dropna(how='bfill', limit=3)` Fill missing values with next mean, limit 3.

`df.rolling(3).std().dropna(how='bfill', limit=3)` Fill missing values with next standard deviation, limit 3.

`df.rolling(3).min().dropna(how='bfill', limit=3)` Fill missing values with next minimum, limit 3.

`df.rolling(3).max().dropna(how='bfill', limit=3)` Fill missing values with next maximum, limit 3.

`df.rolling(3).count().dropna(how='bfill', limit=3)` Fill missing values with next count, limit 3.

`df.rolling(3).sum().dropna(how='ffill', limit=4)` Fill missing values with previous value, limit 4.

`df.rolling(3).mean().dropna(how='ffill', limit=4)` Fill missing values with previous mean, limit 4.

`df.rolling(3).std().dropna(how='ffill', limit=4)` Fill missing values with previous standard deviation, limit 4.

`df.rolling(3).min().dropna(how='ffill', limit=4)` Fill missing values with previous minimum, limit 4.

`df.rolling(3).max().dropna(how='ffill', limit=4)` Fill missing values with previous maximum, limit 4.

`df.rolling(3).count().dropna(how='ffill', limit=4)` Fill missing values with previous count, limit 4.

`df.rolling(3).sum().dropna(how='bfill', limit=4)` Fill missing values with next value, limit 4.

`df.rolling(3).mean().dropna(how='bfill', limit=4)` Fill missing values with next mean, limit 4.

`df.rolling(3).std().dropna(how='bfill', limit=4)` Fill missing values with next standard deviation, limit 4.

`df.rolling(3).min().dropna(how='bfill', limit=4)` Fill missing values with next minimum, limit 4.

`df.rolling(3).max().dropna(how='bfill', limit=4)` Fill missing values with next maximum, limit 4.

`df.rolling(3).count().dropna(how='bfill', limit=4)` Fill missing values with next count, limit 4.

`df.rolling(3).sum().dropna(how='ffill', limit=5)` Fill missing values with previous value, limit 5.

`df.rolling(3).mean().dropna(how='ffill', limit=5)` Fill missing values with previous mean, limit 5.

`df.rolling(3).std().dropna(how='ffill', limit=5)` Fill missing values with previous standard deviation, limit 5.

`df.rolling(3).min().dropna(how='ffill', limit=5)` Fill missing values with previous minimum, limit 5.

`df.rolling(3).max().dropna(how='ffill', limit=5)` Fill missing values with previous maximum, limit 5.

`df.rolling(3).count().dropna(how='ffill', limit=5)` Fill missing values with previous count, limit 5.

`df.rolling(3).sum().dropna(how='bfill', limit=5)` Fill missing values with next value, limit 5.

`df.rolling(3).mean().dropna(how='bfill', limit=5)` Fill missing values with next mean, limit 5.

`df.rolling(3).std().dropna(how='bfill', limit=5)` Fill missing values with next standard deviation, limit 5.

`df.rolling(3).min().dropna(how='bfill', limit=5)` Fill missing values with next minimum, limit 5.

`df.rolling(3).max().dropna(how='bfill', limit=5)` Fill missing values with next maximum, limit 5.

`df.rolling(3).count().dropna(how='bfill', limit=5)` Fill missing values with next count, limit 5.

`df.rolling(3).sum().dropna(how='ffill', limit=6)` Fill missing values with previous value, limit 6.

`df.rolling(3).mean().dropna(how='ffill', limit=6)` Fill missing values with previous mean, limit 6.

`df.rolling(3).std().dropna(how='ffill', limit=6)` Fill missing values with previous standard deviation, limit 6.

`df.rolling(3).min().dropna(how='ffill', limit=6)` Fill missing values with previous minimum, limit 6.

`df.rolling(3).max().dropna(how='ffill', limit=6)` Fill missing values with previous maximum, limit 6.

`df.rolling(3).count().dropna(how='ffill', limit=6)` Fill missing values with previous count, limit 6.

`df.rolling(3).sum().dropna(how='bfill', limit=6)` Fill missing values with next value, limit 6.

`df.rolling(3).mean().dropna(how='bfill', limit=6)` Fill missing values with next mean, limit 6.

`df.rolling(3).std().dropna(how='bfill', limit=6)` Fill missing values with next standard deviation, limit 6.

`df.rolling(3).min().dropna(how='bfill', limit=6)` Fill missing values with next minimum, limit 6.

`df.rolling(3).max().dropna(how='bfill', limit=6)` Fill missing values with next maximum, limit 6.

`df.rolling(3).count().dropna(how='bfill', limit=6)` Fill missing values with next count, limit 6.

`df.rolling(3).sum().dropna(how='ffill', limit=7)` Fill missing values with previous value, limit 7.

`df.rolling(3).mean().dropna(how='ffill', limit=7)` Fill missing values with previous mean, limit 7.

`df.rolling(3).std().dropna(how='ffill', limit=7)` Fill missing values with previous standard deviation, limit 7.

`df.rolling(3).min().dropna(how='ffill', limit=7)` Fill missing values with previous minimum, limit 7.

`df.rolling(3).max().dropna(how='ffill', limit=7)` Fill missing values with previous maximum, limit 7.

`df.rolling(3).count().dropna(how='ffill', limit=7)` Fill missing values with previous count, limit 7.

`df.rolling(3).sum().dropna(how='bfill', limit=7)` Fill missing values with next value, limit 7.

`df.rolling(3).mean().dropna(how='bfill', limit=7)` Fill missing values with next mean, limit 7.

`df.rolling(3).std().dropna(how='bfill', limit=7)` Fill missing values with next standard deviation, limit 7.

`df.rolling(3).min().dropna(how='bfill', limit=7)` Fill missing values with next minimum, limit 7.

Data Wrangling with dplyr and tidyr

Cheat Sheet

BecomingHuman.AI

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)

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
4	4.6	3.1	1.5	0.2	setosa
5	5.0	3.6	1.4	0.2	setosa
6	5.4	3.9	1.7	0.4	setosa
7	4.6	3.4	1.4	0.3	setosa
8	5.0	3.4	1.5	0.2	setosa

```
dplyr::%>%
```

Passes object on left hand side as first argument (or .argument) of function on right-hand 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)
```

Tidy Data A foundation for wrangling in R

In a tidy data set:

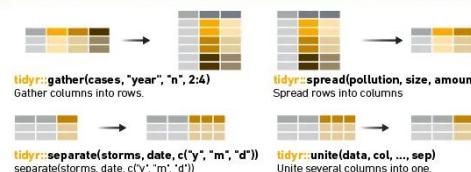


Each variable is saved in its own column
Each observation is saved in its own row

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



Reshaping Data Change the layout of a data set



```
dplyr::data_frame(a = 1:3, b = 4:6)
```

Combine vectors into data frame (optimized).

```
dplyr::arrange(mtcars, mpg)
```

Order rows by values of a column (low to high).

```
dplyr::arrange(mtcars, desc(mpg))
```

Order rows by values of a column (high to low).

```
dplyr::rename(tb, y = year)
```

Rename the columns of a data frame.

Subset Observations (Rows)



```
dplyr::filter(iris, Sepal.Length > 7)
```

Extract rows that meet logical criteria.

```
dplyr::distinct(iris)
```

Remove duplicate rows.

```
dplyr::sample_frac(iris, 0.5, replace = TRUE)
```

Randomly select fraction of rows.

```
dplyr::sample_n(iris, 10, replace = TRUE)
```

Randomly select n rows.

```
dplyr::slice(iris, 10:15)
```

Select rows by position.

```
dplyr::top_n(storms, 2, date)
```

Select and order top n entries (by group if grouped data).

Logic in R - ? Comparison_These
< Less than 1x
> Greater than >x
== Equal to ==x
!= Less than or equal to !=x
>= Greater than or equal to >=x

Boolean operators
& &
| |
! !

group_by(Species) %>%

summarise(avg = mean(Sepal.Width)) %>%

arrange(avg)

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(...)
```

```
iris %>% group_by(Species) %>% mutate(...)
```

Compute separate summary row for each group.

Compute new variables by group.

```
dplyr::group_modify(iris, funs(...))
```

Compute new variables by group.

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).

summary
function

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.

```
dplyr::var
```

Variance of a vector.

```
dplyr::sd
```

Standard deviation of a vector.

```
dplyr::IQR
```

IQR of a vector.

```
dplyr::min
```

Minimum value in a vector.

```
dplyr::max
```

Maximum value in a vector.

```
dplyr::mean
```

Mean value of a vector.

```
dplyr::median
```

Median value of a vector.

```
dplyr::range
```

Range of a vector.

```
dplyr::cumsum
```

Cumulative sum.

```
dplyr::cummax
```

Cumulative max.

```
dplyr::cummin
```

Cumulative min.

```
dplyr::cumprod
```

Cumulative prod.

```
dplyr::nrow
```

Number of rows in a data frame.

```
dplyr::between
```

Are values between a and b?

```
dplyr::cume_dist
```

Cumulative distribution.

```
dplyr::pmin
```

Element-wise min.

Make New Variables



```
dplyr::mutate(iris, sepal = Sepal.Length + Sepal.Width)
```

Compute and append one or more new columns.

```
dplyr::mutate_each(iris, funs(rank))
```

Apply window function to each column.

```
dplyr::transmute(iris, sepal = Sepal.Length + Sepal.Width)
```

Compute one or more new columns. Drop original columns

window
function

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::pmax
```

Element-wise max.

```
dplyr::pmin
```

Element-wise min.

Combine Data Sets

A + B = C

Y + Z = W

Set Operations

X - Y = Z

X ∩ Y = Z

X ∪ Y = Z

X ⊕ Y = Z

Filtering Joins

X ⊂ Y = Z

X ⊃ Y = Z

Binding

X ⊕ Y = Z

Dplyr Bindings

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.

Scipy Linear Algebra Cheat Sheet

BecomingHuman.AI



Interacting With NumPy

Also see NumPy

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

Index Tricks

>>> np.mgrid[0:5,0:5]	Create a dense meshgrid
>>> np.mgrid[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 (number of samples)
>>> g[3] += np.pi	
>>> 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 N-point central derivative
>>> misc.derivative(myfunc,1.0)	Find the n-th derivative of a function at a point

Linear Algebra

Also see NumPy

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

Creating Matrices

```
>>> A = np.matrix(np.random.random((2,2)))
>>> B = np.asmatrix(b)
>>> C = np.matrix(np.random.random((10,5)))
>>> D = np.mat([[3,4], [5,6]])
```

Basic Matrix Routines

Inverse	
>>> A.I	Inverse
>>> linalg.inv(A)	Inverse

Transposition

>>> A.T	Transpose matrix
>>> A.H	Conjugate transposition

Trace

>>> np.trace(A)	Trace
>>> linalg.norm(A)	Frobenius norm
>>> linalg.norm(A)	L1 norm (max column sum)
>>> linalg.norm(A,np.inf)	L inf norm (max row sum)

Rank

>>> np.linalg.matrix_rank(C)	Matrix rank
>>> linalg.det(A)	Determinant

Solving linear problems

>>> linalg.solve(a,b)	Solver for dense matrices
>>> np.linalg.lstsq(F,E)	Solver for dense matrices
>>> linalg.lsqr(F,E)	Least-squares solution to linear matrix

Generalized inverse

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

Creating Matrices

>>> F = np.eye(3, k=1)	Create a 2X2 Identity matrix
>>> G = np.mat(np.identity(2))	Create a 2x2 Identity matrix
>>> C[> 0.5] = 0	
>>> H = sparse.csr_matrix(C)	Compressed Sparse Row matrix
>>> I = sparse.csc_matrix(D)	Compressed Sparse Column matrix
>>> J = sparse.dok_matrix(A)	Dictionary Of Keys matrix
>>> E.todense()	Sparse matrix to full matrix
>>> sparse.isspmatrix_csc(A)	Identify sparse matrix

Matrix Functions

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

Subtraction

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

Division

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

Multiplication

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

Exponential Functions

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

Logarithm Function

>>> linalg.logm(A)	Matrix logarithm
--------------------	------------------

Trigonometric Functions

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

Hyperbolic Trigonometric Functions

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

Matrix Sign Function

>>> np.signm(A)	Matrix sign function
-----------------	----------------------

Matrix Square Root

>>> linalg.sqrtm(A)	Matrix square root
---------------------	--------------------

Arbitrary Functions

>>> linalg.funm(A, lambda x: x**k)	Evaluate matrix function
------------------------------------	--------------------------

Sparse Matrix Routines

Inverse	Inverse
>>> sparse.linalg.inv(H)	

Norm

>>> sparse.linalg.norm(H)	Norm
---------------------------	------

Solving linear problems

>>> sparse.linalg.spsolve(H,I)	Solver for sparse matrices
--------------------------------	----------------------------

Sparse Matrix Functions

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

Decompositions

Eigenvalues and Eigenvectors	
-------------------------------------	--

>>> la, lv = linalg.eig(A)	Solve ordinary or generalized eigenvalue problem for square matrix
>>> l, l2 = la	
>>> V[0]	First eigenvector
>>> V[1]	Second eigenvector
>>> linalg.eigvals(A)	Unpack eigenvalues

Singular Value Decomposition

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

LU Decomposition

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

Sparse Matrix Decompositions

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

Asking For Help

>>> help(scipy.linalg.diagsvd)

>>> np.info(np.matrix)

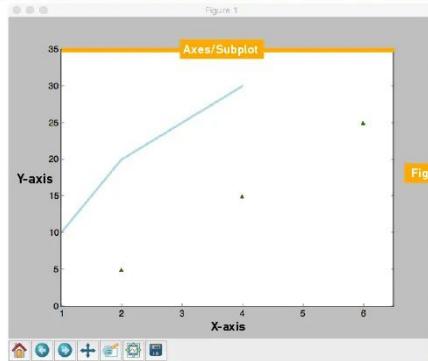
Matplotlib is a Python 2D plotting library which produces publication-quality figures in a variety of hardcopy formats and interactive environments across platforms.

Matplotlib Cheat Sheet

BecomingHuman.AI

Anatomy & Workflow

Plot Anatomy



Workflow

- 01** Prepare data
- 02** Create plot
- 03** Plot
- 04** Customize plot
- 05** Save plot
- 06** Show plot

```
>>> import matplotlib.pyplot as plt
step 1: >>> x = [1,2,3,4]
>>> y = [10,20,25,30]
step 2: >>> fig = plt.figure()
step 3: >>> ax = fig.add_subplot(111)
step 4: >>> ax.plot(x, y, color='lightblue', linewidth=3)
>>> ax.scatter([2,4,6], [5,15,25], color='darkgreen', marker='^')
step 5: >>> ax.set_xlim(1, 6)
>>> plt.savefig('foo.png')
>>> plt.show()
```

Prepare The Data

Also see [Lists & NumPy](#)

Index Tricks

```
>>> import numpy as np
>>> x = np.linspace(0, 10, 100)
>>> y = np.cos(x)
>>> z = np.sin(x)

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

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)
>>> fig1, axes = plt.subplots(nrows=2, ncols=2)
>>> fig4, axes2 = plt.subplots(ncols=3)
```

Plotting Routines

1D Data

>>> 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].barh([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

```
>>> fig, ax = plt.subplots()
>>> im = ax.imshow(img,
                  arrays cmap='gist_earth',
                  interpolation='nearest',
                  vmin=-2,
                  vmax=2)
```

Customize Plot

Colors, Color Bars & Color Maps

```
>>> plt.plot(x, y, x**2, x**3)
>>> ax.plot(x, y, alpha = 0.4)
>>> ax.plot(x, c=k)
>>> fig.colorbar(im, orientation='horizontal')
>>> im = plt.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=6.0)
>>> 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, -2.1, 'Example Graph',
           style='italic',
           size=16)
>>> ax.annotate(some_text, (8, 0),
               xycoords='data',
               xytext=(10.5, 0),
               textcoords='data',
               arrowprops=dict(arrowstyle='<-',
                               connectionstyle='arc3'),)
```

MathText

```
>>> plt.title(r'$\sigma_{\mathrm{I}}$=15$', fontsize=20)
```

Vector Fields

```
>>> axes[0,1].arrow(0,0,0.5,0.5)
>>> axes[1,1].quiver(y,z)
>>> axes[0,1].streamplot(X,Y,U,V)
```

Data Distributions

>>> ax1.hist(y)	Plot a histogram
>>> ax3.boxplot(y)	Make a box and whisker plot
>>> ax3.violinplot(z)	Make a violin plot

>>> axes[2,0].pcolor(data2)	Pseudocolor plot of 2D array
>>> axes[2,0].pcolormesh(data)	Pseudocolor plot of 2D array
>>> CS = plt.contour(Y,X,U)	Plot contours
>>> axes[2,1].contour(data1)	Plot filled contours
>>> axes[2,2]= ax.clabel(CS)	Label a contour plot

Limits, Legends & Layouts

>>> ax.margins(x=0.0,y=0.1)	Add padding to a plot
>>> ax.axis('equal')	Set the aspect ratio of the plot to 1
>>> ax.set_xlim(0,10.5),ylim=[-1.5,1.5])	Set limits for x- and y-axis
>>> ax.set_xlim(0,10.5)	Set limits for x-axis

>>> ax.set_title('An Example Axes', y=1.05, xlabel='X-Axis', ylabel='Y-Axis')	Set a title and x- and y-axis labels
>>> ax.legend(loc='best')	No overlapping plot elements

>>> ax.xaxis.set(ticks=range(1,5), ticklabels=[3,100,-12,'foo'])	Manually set x-ticks
>>> ax.yaxis.set(ticks=range(1,5), ticklabels=[3,100,-12,'foo'])	Make y-ticks longer and go in and out

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

>>> ax1.spines['top'].set_visible(False)	Axis Spines
>>> ax1.spines['bottom'].set_position(('outward',10))	Move the top axis line for a plot invisible
>>> ax1.spines['bottom'].set_position(('outward',10))	Move the bottom axis line outward

Save Plot

Save figures

```
>>> plt.savefig('foo.png')
```

Save transparent figures

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

Show Plot

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

Close & Clear

```
>>> plt.cla()
```

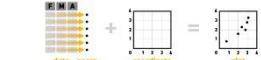
```
>>> plt.clf()
```

```
>>> plt.close()
```

Data Visualisation 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/y locations



Build a graph with `plot()` or `ggplot()`

`ggplot(mapping, geom)`
Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

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

Data
`ggplot(mapping, aes(hwy, cyl)) + layer + geom + default stat + geom_point(mapping, color = cyl) + layer + geom + default stat + geom_smooth(mapping, method = "lm") + layer + geom + mapping + scale_color_gradient() + theme_bw() + 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

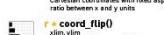
`ggsave("plot.png", width = 5, height = 5)`
Saves last plot as 5" x 5" file named "plot.png" in working directory. Matches file type to file extension.

Coordinate Systems

`r <- b + geom_bar()`



`r + coord_fixed(ratio = 1/2)`
Cartesian coordinates with fixed aspect ratio, x and y axis



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

`r + coord_trans(trans = "sqrt")`
Square root transformation. Set trans to "sqrt" for square root.

`r + coord_map(projection = "ortho", orientation = c(0, 0, 0))`
Map projections from the `grid` package
Mercator (default), stereographic, lagrange, etc.

RStudio® is a trademark of RStudio, Inc. • CC BY SA BecomingHuman • Team@BecomingHuman.AI • 844-448-1212 • BecomingHuman.AI • Updated: 2018-03

Creative Commons Data Visualisation with ggplot2 by RStudio is licensed under CC BY SA 4.0

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_rect(stat = "bin")
    a + geom_density(kernel = "density", stat = "bin")
  a + geom_hex()
    a + geom_hex(stat = "density", ..count..)
  a + geom_hex(stat = "frequency")
    a + geom_hex(stat = "density", ..density..)
```

Discrete

```
b <- ggplot(mpg, aes(fct))
  b + geom_bar()
```

Geographic Primitives

```
c <- ggplot(mtcars, aes(carb))
  c + geom_polygon(mapping = group)
  c + geom_path(mapping = group)
  c + geom_point(mapping = group)
  c + geom_rect(mapping = group, ymin = -900, ymax = +900, fill = NA)
  c + geom_violin(mapping = group)
```

Continuous

```
d <- ggplot(mtcars, aes(cyl))
  d + geom_boxplot()
  d + geom_dotplot(binaxis = "y", stackdir = "center")
  d + geom_hex(stat = "hex", binwidth = 1, binheight = 1, fill = "#f0f0f0")
```

Discrete X, Continuous Y

```
e <- ggplot(mtcars, aes(cyl))
  e + geom_barstat(stat = "identity")
  e + geom_crossbar(stat = "identity")
  e + geom_hex(stat = "identity")
```

Discrete X, Discrete Y

```
f <- ggplot(diamonds, aes(cut, color))
  f + geom_jitter()
  f + geom_rect(mapping = list(ymin = lat, ymax = lat + delta_y),
    xmin = xmn, xmax = xmx, alpha = 0.5, fill = "black")
```

Maps

Three Variables

`sealSize <- whitebeals.squidata; sealSize$z <- 20`
`m <- ggplot(sealSize, aes(z = z))`

`m + geom_contour(aes(z = z))`
Contour lines

`m + geom_raster(aes(z = z), hijust = 0.5, interpolation = "bilinear")`
Raster

`m + geom_hex(aes(z = z))`
Hex grid

Faceting

Facets divide a plot into subplots based on the values of one or more discrete variables

`t <- ggplot(mtcars, aes(wt)) + geom_point()`

`t + facet_grid(~ am)`
Facet into columns based on am

`t + facet_grid(~ J)`
Facet into rows based on J

`t + facet_grid(~ J | am)`
Facet into both rows and columns

`t + facet_wrap(~ J)`
Wrap facets into a rectangular layout of one or more rows and columns

`t + facet_wrap(~ am)`
Wrap facets into a rectangular layout of one or more rows and columns

`t + facet_grid(~ am) + scales = "free_x"`
x and y axis limits adjust to individual facets

`*free_x*, *x* x axis limits adjust`

`*free_y*, *y* y axis limits adjust`

Set scales to let ggplot2 know which facets

`t + facet_grid(~ am) + scales = "fixed"`

`x + y` x and y axis limits

`*free_x*, *x* x axis limits adjust`

`*free_y*, *y* y axis limits adjust`

Set labeller to adjust facet labels

`t + facet_grid(~ am, labeller = label_both)`

`t + facet_grid(~ am, labeller = label_bquote)`

`t + facet_grid(~ am, labeller = label_c)`

`t + facet_grid(~ am, labeller = label_d)`

`t + facet_grid(~ am, labeller = label_g)`

`t + facet_grid(~ am, labeller = label_l)`

`t + facet_grid(~ am, labeller = label_n)`

`t + facet_grid(~ am, labeller = label_s)`

`t + facet_grid(~ am, labeller = label_t)`

`t + facet_grid(~ am, labeller = label_v)`

`t + facet_grid(~ am, labeller = label_w)`

`t + facet_grid(~ am, labeller = label_x)`

`t + facet_grid(~ am, labeller = label_y)`

`t + facet_grid(~ am, labeller = label_z)`

`t + facet_grid(~ am, labeller = label_a)`

`t + facet_grid(~ am, labeller = label_b)`

`t + facet_grid(~ am, labeller = label_c)`

`t + facet_grid(~ am, labeller = label_d)`

`t + facet_grid(~ am, labeller = label_e)`

`t + facet_grid(~ am, labeller = label_f)`

`t + facet_grid(~ am, labeller = label_g)`

`t + facet_grid(~ am, labeller = label_h)`

`t + facet_grid(~ am, labeller = label_i)`

`t + facet_grid(~ am, labeller = label_j)`

`t + facet_grid(~ am, labeller = label_k)`

`t + facet_grid(~ am, labeller = label_l)`

`t + facet_grid(~ am, labeller = label_m)`

`t + facet_grid(~ am, labeller = label_n)`

`t + facet_grid(~ am, labeller = label_o)`

`t + facet_grid(~ am, labeller = label_p)`

`t + facet_grid(~ am, labeller = label_q)`

`t + facet_grid(~ am, labeller = label_r)`

`t + facet_grid(~ am, labeller = label_s)`

`t + facet_grid(~ am, labeller = label_t)`

`t + facet_grid(~ am, labeller = label_u)`

`t + facet_grid(~ am, labeller = label_v)`

`t + facet_grid(~ am, labeller = label_w)`

`t + facet_grid(~ am, labeller = label_x)`

`t + facet_grid(~ am, labeller = label_y)`

`t + facet_grid(~ am, labeller = label_z)`

`t + facet_grid(~ am, labeller = label_a)`

`t + facet_grid(~ am, labeller = label_b)`

`t + facet_grid(~ am, labeller = label_c)`

`t + facet_grid(~ am, labeller = label_d)`

`t + facet_grid(~ am, labeller = label_e)`

`t + facet_grid(~ am, labeller = label_f)`

`t + facet_grid(~ am, labeller = label_g)`

`t + facet_grid(~ am, labeller = label_h)`

`t + facet_grid(~ am, labeller = label_i)`

`t + facet_grid(~ am, labeller = label_j)`

`t + facet_grid(~ am, labeller = label_k)`

`t + facet_grid(~ am, labeller = label_l)`

`t + facet_grid(~ am, labeller = label_m)`

`t + facet_grid(~ am, labeller = label_n)`

`t + facet_grid(~ am, labeller = label_o)`

`t + facet_grid(~ am, labeller = label_p)`

`t + facet_grid(~ am, labeller = label_q)`

`t + facet_grid(~ am, labeller = label_r)`

`t + facet_grid(~ am, labeller = label_s)`

`t + facet_grid(~ am, labeller = label_t)`

`t + facet_grid(~ am, labeller = label_u)`

`t + facet_grid(~ am, labeller = label_v)`

`t + facet_grid(~ am, labeller = label_w)`

`t + facet_grid(~ am, labeller = label_x)`

`t + facet_grid(~ am, labeller = label_y)`

`t + facet_grid(~ am, labeller = label_z)`

`t + facet_grid(~ am, labeller = label_a)`

`t + facet_grid(~ am, labeller = label_b)`

`t + facet_grid(~ am, labeller = label_c)`

`t + facet_grid(~ am, labeller = label_d)`

`t + facet_grid(~ am, labeller = label_e)`

`t + facet_grid(~ am, labeller = label_f)`

`t + facet_grid(~ am, labeller = label_g)`

`t + facet_grid(~ am, labeller = label_h)`

`t + facet_grid(~ am, labeller = label_i)`

`t + facet_grid(~ am, labeller = label_j)`

`t + facet_grid(~ am, labeller = label_k)`

`t + facet_grid(~ am, labeller = label_l)`

`t + facet_grid(~ am, labeller = label_m)`

`t + facet_grid(~ am, labeller = label_n)`

`t + facet_grid(~ am, labeller = label_o)`

`t + facet_grid(~ am, labeller = label_p)`

`t + facet_grid(~ am, labeller = label_q)`

`t + facet_grid(~ am, labeller = label_r)`

`t + facet_grid(~ am, labeller = label_s)`

`t + facet_grid(~ am, labeller = label_t)`

`t + facet_grid(~ am, labeller = label_u)`

`t + facet_grid(~ am, labeller = label_v)`

`t + facet_grid(~ am, labeller = label_w)`

`t + facet_grid(~ am, labeller = label_x)`

`t + facet_grid(~ am, labeller = label_y)`

`t + facet_grid(~ am, labeller = label_z)`

`t + facet_grid(~ am, labeller = label_a)`

`t + facet_grid(~ am, labeller = label_b)`

`t + facet_grid(~ am, labeller = label_c)`

`t + facet_grid(~ am, labeller = label_d)`

`t + facet_grid(~ am, labeller = label_e)`

`t + facet_grid(~ am, labeller = label_f)`

`t + facet_grid(~ am, labeller = label_g)`

`t + facet_grid(~ am, labeller = label_h)`

`t + facet_grid(~ am, labeller = label_i)`

`t + facet_grid(~ am, labeller = label_j)`

`t + facet_grid(~ am, labeller = label_k)`

`t + facet_grid(~ am, labeller = label_l)`

`t + facet_grid(~ am, labeller = label_m)`

`t + facet_grid(~ am, labeller = label_n)`

`t + facet_grid(~ am, labeller = label_o)`

`t + facet_grid(~ am, labeller = label_p)`

`t + facet_grid(~ am, labeller = label_q)`

`t + facet_grid(~ am, labeller = label_r)`

`t + facet_grid(~ am, labeller = label_s)`

`t + facet_grid(~ am, labeller = label_t)`

`t + facet_grid(~ am, labeller = label_u)`

`t + facet_grid(~ am, labeller = label_v)`

`t + facet_grid(~ am, labeller = label_w)`

`t + facet_grid(~ am, labeller = label_x)`

`t + facet_grid(~ am, labeller = label_y)`

`t + facet_grid(~ am, labeller = label_z)`

`t + facet_grid(~ am, labeller = label_a)`

`t + facet_grid(~ am, labeller = label_b)`

`t + facet_grid(~ am, labeller = label_c)`

`t + facet_grid(~ am, labeller = label_d)`

`t + facet_grid(~ am, labeller = label_e)`

`t + facet_grid(~ am, labeller = label_f)`

`t + facet_grid(~ am, labeller = label_g)`

`t + facet_grid(~ am, labeller = label_h)`

`t + facet_grid(~ am, labeller = label_i)`

`t + facet_grid(~ am, labeller = label_j)`

`t + facet_grid(~ am, labeller = label_k)`

`t + facet_grid(~ am, labeller = label_l)`

`t + facet_grid(~ am, labeller = label_m)`

`t + facet_grid(~ am, labeller = label_n)`

`t + facet_grid(~ am, labeller = label_o)`

`t + facet_grid(~ am, labeller = label_p)`

`t + facet_grid(~ am, labeller = label_q)`

`t + facet_grid(~ am, labeller = label_r)`

`t + facet_grid(~ am, labeller = label_s)`

`t + facet_grid(~ am, labeller = label_t)`

`t + facet_grid(~ am, labeller = label_u)`

`t + facet_grid(~ am, labeller = label_v)`

`t + facet_grid(~ am, labeller = label_w)`

`t + facet_grid(~ am, labeller = label_x)`

`t + facet_grid(~ am, labeller = label_y)`

`t + facet_grid(~ am, labeller = label_z)`

`t + facet_grid(~ am, labeller = label_a)`

`t + facet_grid(~ am, labeller = label_b)`

`t + facet_grid(~ am, labeller = label_c)`

`t + facet_grid(~ am, labeller = label_d)`

`t + facet_grid(~ am, labeller = label_e)`

`t + facet_grid(~ am, labeller = label_f)`

`t + facet_grid(~ am, labeller = label_g)`

`t + facet_grid(~ am, labeller = label_h)`

`t + facet_grid(~ am, labeller = label_i)`

`t + facet_grid(~ am, labeller = label_j)`

`t + facet_grid(~ am, labeller = label_k)`

`t + facet_grid(~ am, labeller = label_l)`

`t + facet_grid(~ am, labeller = label_m)`

`t + facet_grid(~ am, labeller = label_n)`

`t + facet_grid(~ am, labeller = label_o)`

`t + facet_grid(~ am, labeller = label_p)`

`t + facet_grid(~ am, labeller = label_q)`

`t + facet_grid(~ am, labeller = label_r)`

`t + facet_grid(~ am, labeller = label_s)`

`t + facet_grid(~ am, labeller = label_t)`

`t + facet_grid(~ am, labeller = label_u)`

`t + facet_grid(~ am, labeller = label_v)`

`t + facet_grid(~ am, labeller = label_w)`

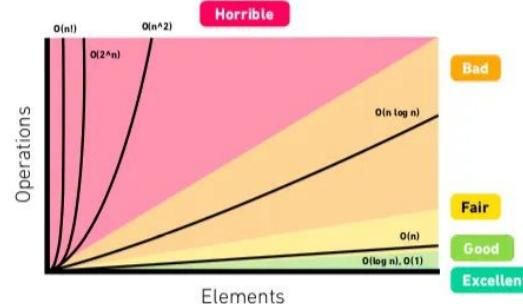
`t + facet_grid(~ am, labeller = label_x)`

`t + facet_grid(~ am, labeller = label_y)`

`t + facet_grid(~ am, labeller = label_z)`

`t + facet_grid(~ am, labeller = label_a)`

Big-O Complexity Chart



Big-O Cheat Sheet

BecomingHuman.AI

Data Structure Operation

	Time Complexity				Space Complexity			
	Average				Worst			
	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	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(n)$
Queue	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(n)$
Singly-Linked List	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(n)$
Doubly-Linked List	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(n)$
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	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	N/A	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$
B-Tree	N/A	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(n)$
Red-Black Tree	$\Theta(\log(n))$	$\Theta(n)$						
Splay Tree	N/A	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(\log(n))$	N/A	$\Theta(\log(n))$	$\Theta(\log(n))$	$\Theta(n)$
AVL Tree	$\Theta(\log(n))$	$\Theta(n)$						
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

	Time Complexity				Space Complexity			
	Best				Worst			
	Access	Search	Insertion	Deletion	Access	Search	Insertion	Deletion
Quicksort	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n^2)$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$
Mergesort	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$
Timsort	$\Theta(n)$	$\Theta(n)$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n)$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n)$
Heapsort	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$
Bubble Sort	$\Theta(n)$	$\Theta(n)$	$\Theta(n^2)$	$\Theta(n^2)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n^2)$	$\Theta(n)$
Insertion Sort	$\Theta(n)$	$\Theta(n)$	$\Theta(n^2)$	$\Theta(n^2)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n^2)$	$\Theta(n)$
Selection Sort	$\Theta(n^2)$	$\Theta(n^2)$	$\Theta(n^2)$	$\Theta(n^2)$	$\Theta(n^2)$	$\Theta(n^2)$	$\Theta(n^2)$	$\Theta(n^2)$
Tree Sort	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n^2)$	$\Theta(n^2)$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n^2)$	$\Theta(n \log(n))$
Shell Sort	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))^2$	$\Theta(n \log(n))^2$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$
Bucket Sort	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$
Radix Sort	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$
Counting Sort	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$	$\Theta(n+k)$
Cubesort	$\Theta(n)$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n)$	$\Theta(n \log(n))$	$\Theta(n \log(n))$	$\Theta(n \log(n))$