

A mostly complete chart of
Neural Networks

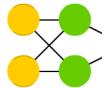
©2016 Fjodor van Veen - asimovinstitute.org

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

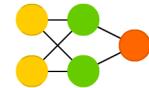
Perceptron (P)



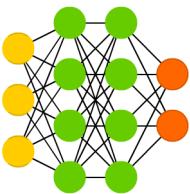
Feed Forward (FF)



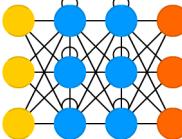
Radial Basis Network (RBF)



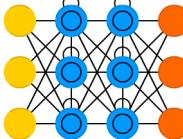
Deep Feed Forward (DFF)



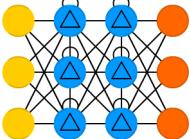
Recurrent Neural Network (RNN)



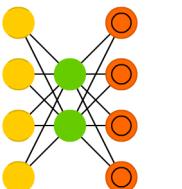
Long / Short Term Memory (LSTM)



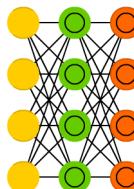
Gated Recurrent Unit (GRU)



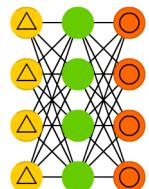
Auto Encoder (AE)



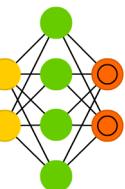
Variational AE (VAE)



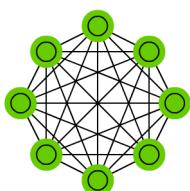
Denoising AE (DAE)



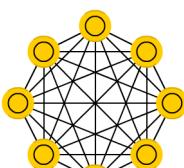
Sparse AE (SAE)



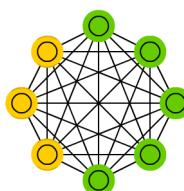
Markov Chain (MC)



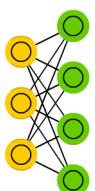
Hopfield Network (HN)



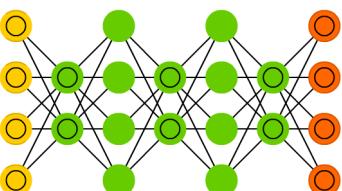
Boltzmann Machine (BM)



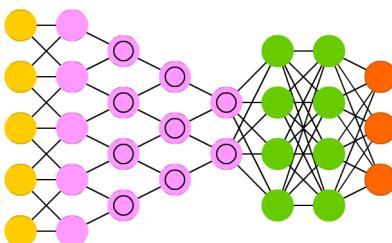
Restricted BM (RBM)



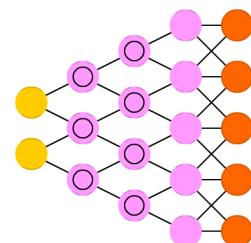
Deep Belief Network (DBN)



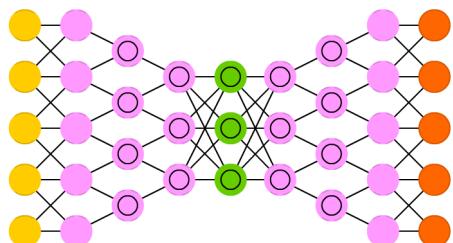
Deep Convolutional Network (DCN)



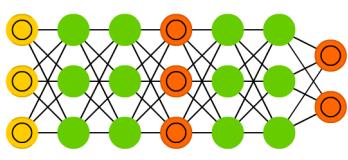
Deconvolutional Network (DN)



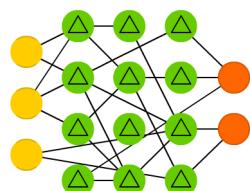
Deep Convolutional Inverse Graphics Network (DCIGN)



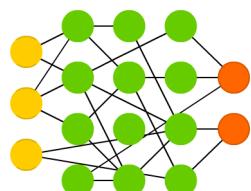
Generative Adversarial Network (GAN)



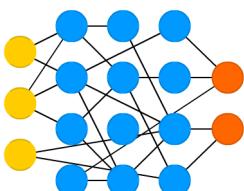
Liquid State Machine (LSM)



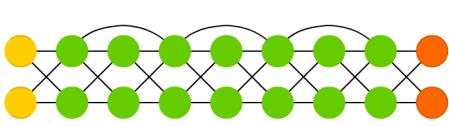
Extreme Learning Machine (ELM)



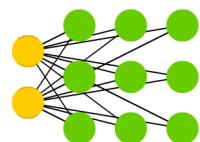
Echo State Network (ESN)



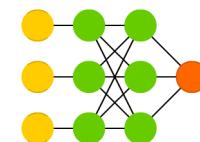
Deep Residual Network (DRN)



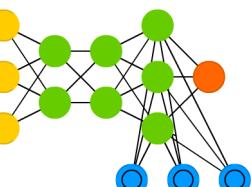
Kohonen Network (KN)



Support Vector Machine (SVM)



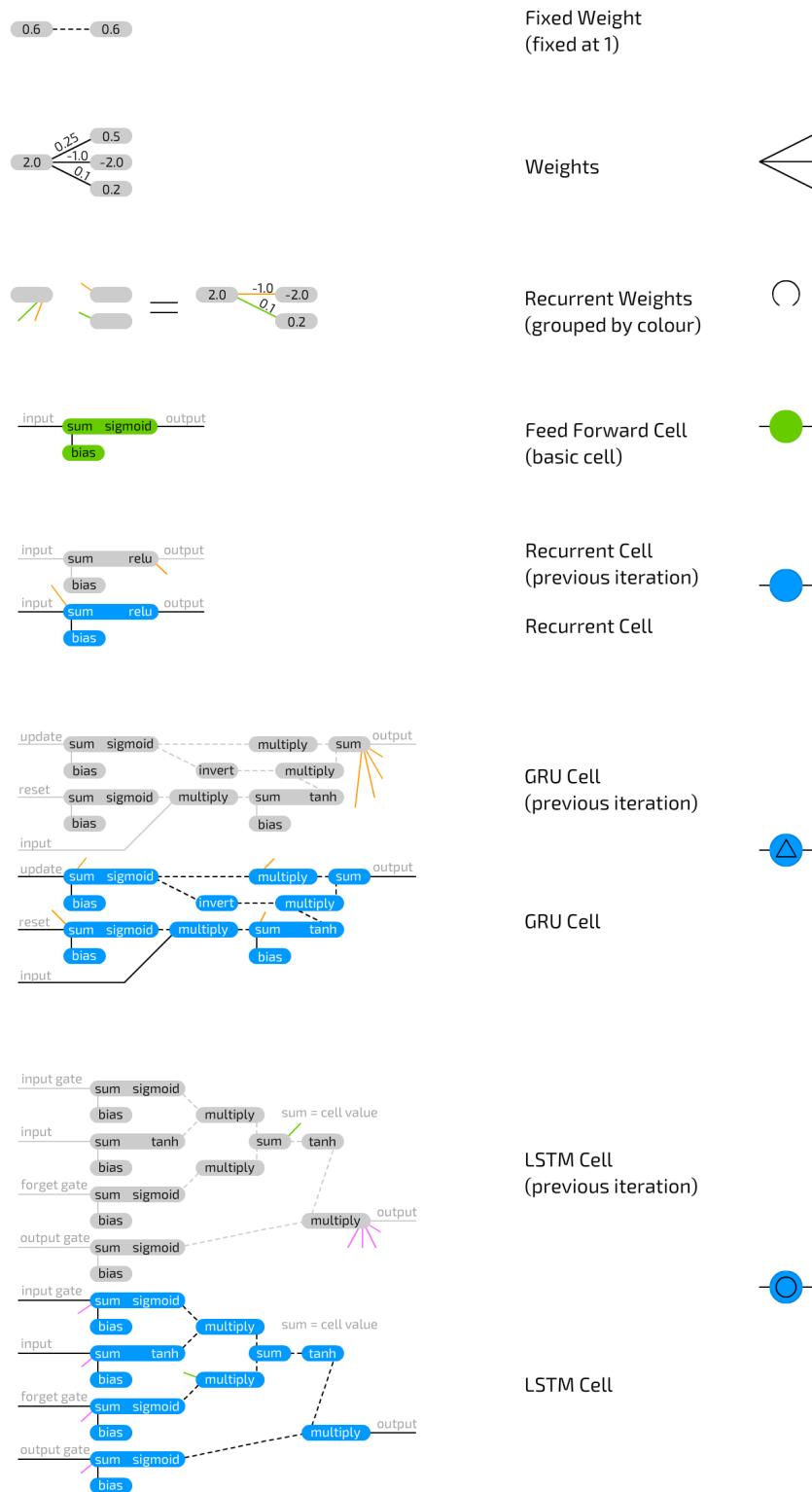
Neural Turing Machine (NTM)



An informative chart to build

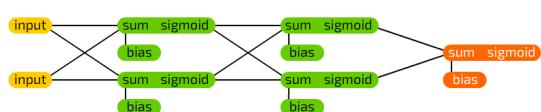
Neural Network Cells

©2016 Fjodor van Veen - asimovinstitute.org

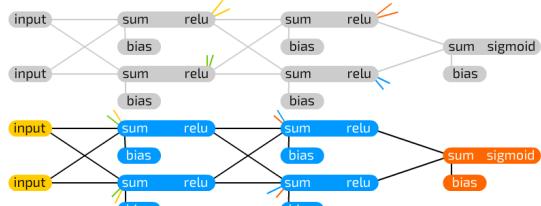
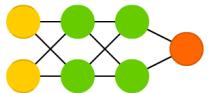


An informative chart to build
Neural Network Graphs

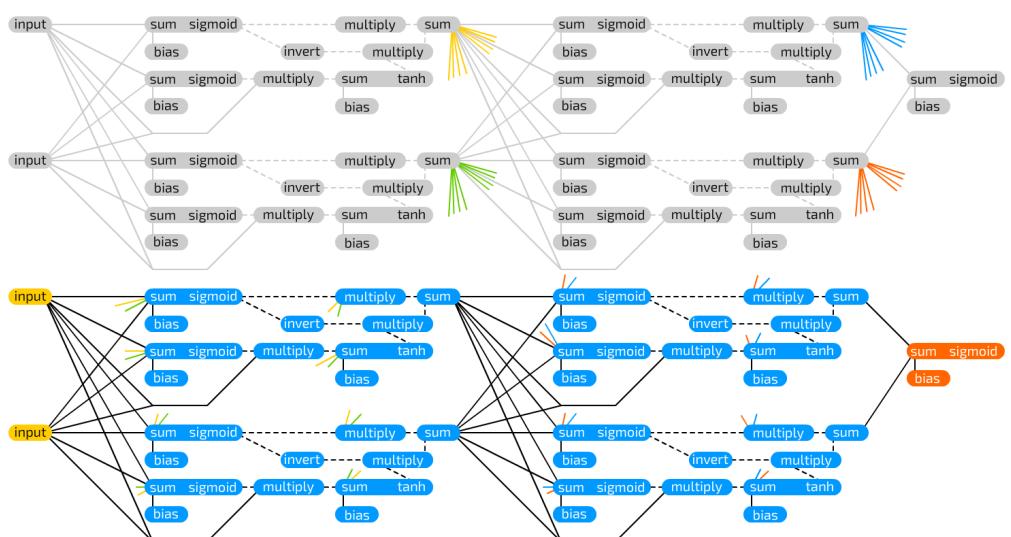
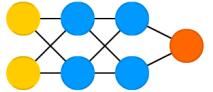
©2016 Fjodor van Veen - asimovinstitute.org



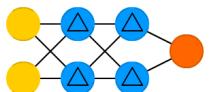
Deep Feed Forward Example



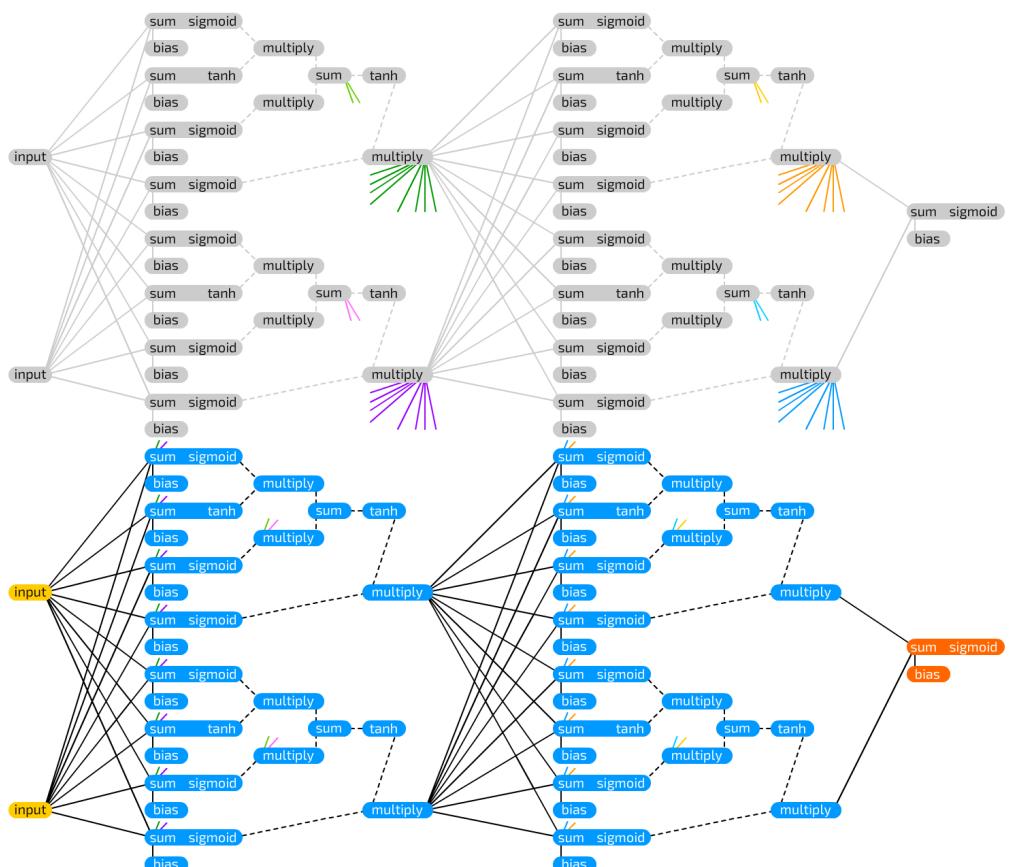
Deep Recurrent Example
(previous iteration)



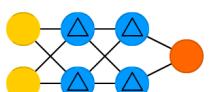
Deep GRU Example
(previous iteration)



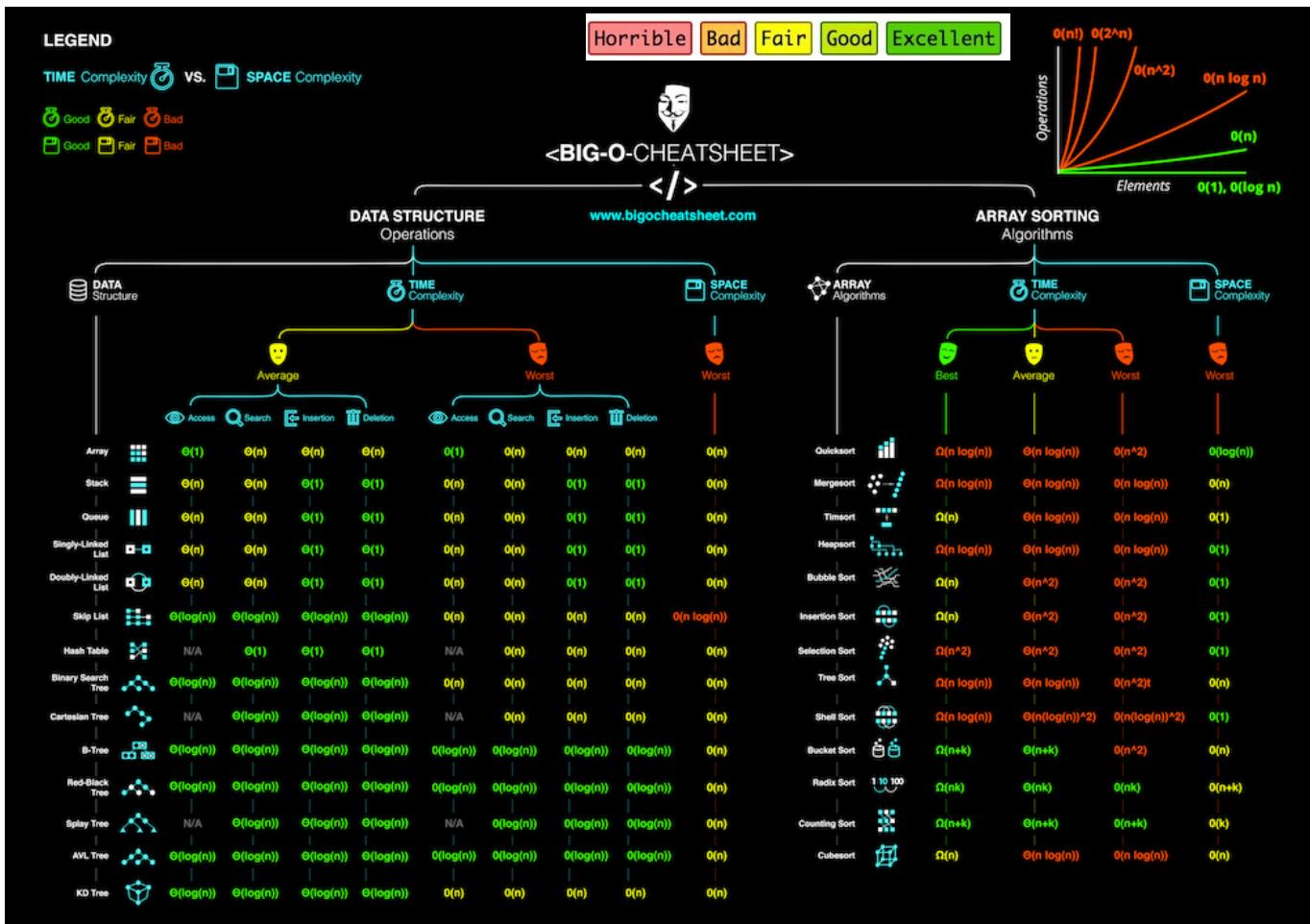
Deep GRU Example



Deep LSTM Example
(previous iteration)



Deep LSTM Example



Common Data Structure Operations

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

MACHINE LEARNING IN EMOJI

SUPERVISED

UNSUPERVISED

REINFORCEMENT

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

BASIC REGRESSION

- LINEAR** linear_model.LinearRegression()
Lots of numerical data   
- LOGISTIC** linear_model.LogisticRegression()
Target variable is categorical  or 

CLASSIFICATION

- NEURAL NET** neural_network.MLPClassifier()
Complex relationships. Prone to overfitting
Basically magic. 
- K-NN** neighbors.KNeighborsClassifier()
Group membership based on proximity 
- DECISION TREE** tree.DecisionTreeClassifier()
If/then/else. Non-contiguous data
Can also be regression  
- RANDOM FOREST** ensemble.RandomForestClassifier()
Find best split randomly
Can also be regression    

- SVM** svm.SVC() svm.LinearSVC()
Maximum margin classifier. Fundamental Data Science algorithm 
- NAIVE BAYES** GaussianNB() MultinomialNB() BernoulliNB()
Updating knowledge step by step with new info 

CLUSTER ANALYSIS

- K-MEANS** cluster.KMeans()
Similar datum into groups based on centroids 
- ANOMALY DETECTION** covariance.EllipticalEnvelope()
Finding outliers through grouping    

FEATURE REDUCTION

- T-DISTRIB STOCHASTIC NEIB EMBEDDING** manifold.TSNE()
Visualize high dimensional data. Convert similarity to joint probabilities 
- PRINCIPLE COMPONENT ANALYSIS** decomposition.PCA()
Distill feature space into components that describe greatest variance 
- CANONICAL CORRELATION ANALYSIS** decomposition.CCA()
Making sense of cross-correlation matrices 
- LINEAR DISCRIMINANT ANALYSIS** Ida.LDA()
Linear combination of features that separates classes  

OTHER IMPORTANT CONCEPTS

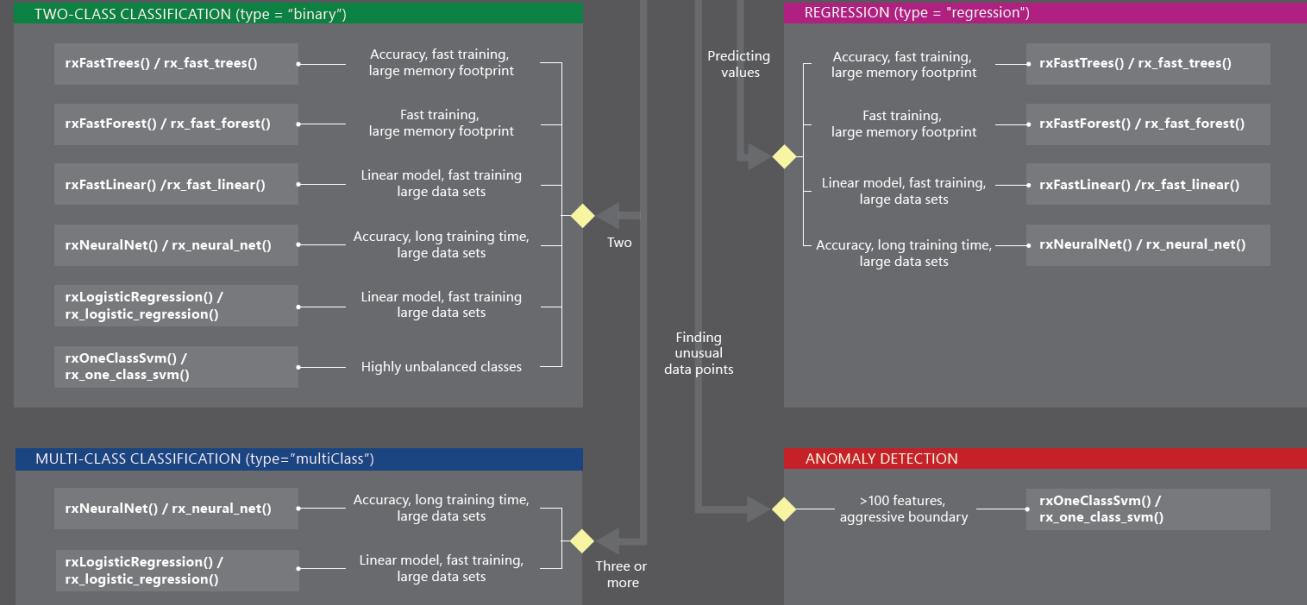
- BIAS VARIANCE TRADEOFF**
- UNDERFITTING / OVERFITTING**
- INERTIA**
- ACCURACY FUNCTION** $(TP + TN) / (P + N)$
- Precision Function** $TP / (TP + FP)$
- Specificity Function** $TN / (FP + TN)$
- Sensitivity Function** $TP / (TP + FN)$





MicrosoftML: Algorithm Cheat Sheet for R / Python

This cheat sheet helps you choose the best MicrosoftML 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.



© 2017 Microsoft Corporation. All rights reserved.



Python For Data Science Cheat Sheet

NumPy Basics

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



NumPy

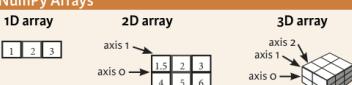
The NumPy library is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays.

Use the following import convention:

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



NumPy Arrays



Creating Arrays

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

Create an array of zeros
Create an array of ones
Create an array of evenly spaced values (step value)
Create an array of evenly spaced values (number of samples)
Create a constant array
Create a 2x2 identity matrix
Create an array with random values
Create an empty array

Initial Placeholders

```
>>> np.zeros((3,4))
>>> np.ones((2,3,4),dtype=np.int16)
>>> d = np.arange(10,25,5)
>>> np.linspace(0,2,9)
>>> e = np.full((2,2),7)
>>> f = np.eye(2)
>>> np.random.random((2,2))
>>> np.empty((3,2))
```

Create an array of zeros
Create an array of ones
Create an array of evenly spaced values (step value)
Create an array of evenly spaced values (number of samples)
Create a constant array
Create a 2x2 identity matrix
Create an array with random values
Create an empty array

I/O

Saving & Loading On Disk

```
>>> np.save('my_array', a)
>>> np.savetxt('array.npz', a, b)
>>> np.load('my_array.npy')
```

Saving & Loading Text Files

```
>>> np.loadtxt("myfile.txt")
>>> np.genfromtxt("my_file.csv", delimiter=',')
>>> np.savetxt("myarray.txt", a, delimiter=" ")
```

Data Types

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

Inspecting Your Array

```
>>> a.shape
>>> len(a)
>>> b.ndim
>>> a.size
>>> b.dtype
>>> b.dtype.name
>>> b.astype(int)
```

Array dimensions
Length of array
Number of array dimensions
Number of array elements
Data type of array elements
Name of data type
Convert an array to a different type

Asking For Help

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

Array Mathematics

Arithmetic Operations

```
>>> q = a - b
>>> array([[-0.5, 0., 0.],
           [-3., -3., -3.]])
>>> np.subtract(a,b)
>>> b + a
>>> array([[ 2.5,  4.,  6.],
           [ 5.,  7.,  9.]])
>>> np.add(b,a)
>>> a / b
>>> array([[ 0.66666667,  1.        ,  1.        ],
           [ 0.25,  0.4,  0.5        ]])
>>> np.divide(a,b)
>>> a * b
>>> array([[ 1.5,  4.,  9.],
           [ 10.,  10.,  18.]])
>>> np.multiply(a,b)
>>> np.exp(b)
>>> np.sqrt(b)
>>> np.sin(a)
>>> np.cos(b)
>>> np.log(a)
>>> a.dot(f)
>>> array([[ 1.,  2.,  3.],
           [ 7.,  7.,  7.]])
```

Subtraction
Addition
Multiplication
Division
Element-wise comparison
Comparison
Aggregate Functions

Comparison

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

Element-wise comparison
Element-wise comparison
Array-wise comparison

Copying Arrays

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

Create a view of the array with the same data
Create a copy of the array
Create a deep copy of the array

Sorting Arrays

```
>>> a.argsort()
>>> c.argsort(axis=0)
```

Sort an array
Sort the elements of an array's axis

Subsetting, Slicing, Indexing

Also see Lists

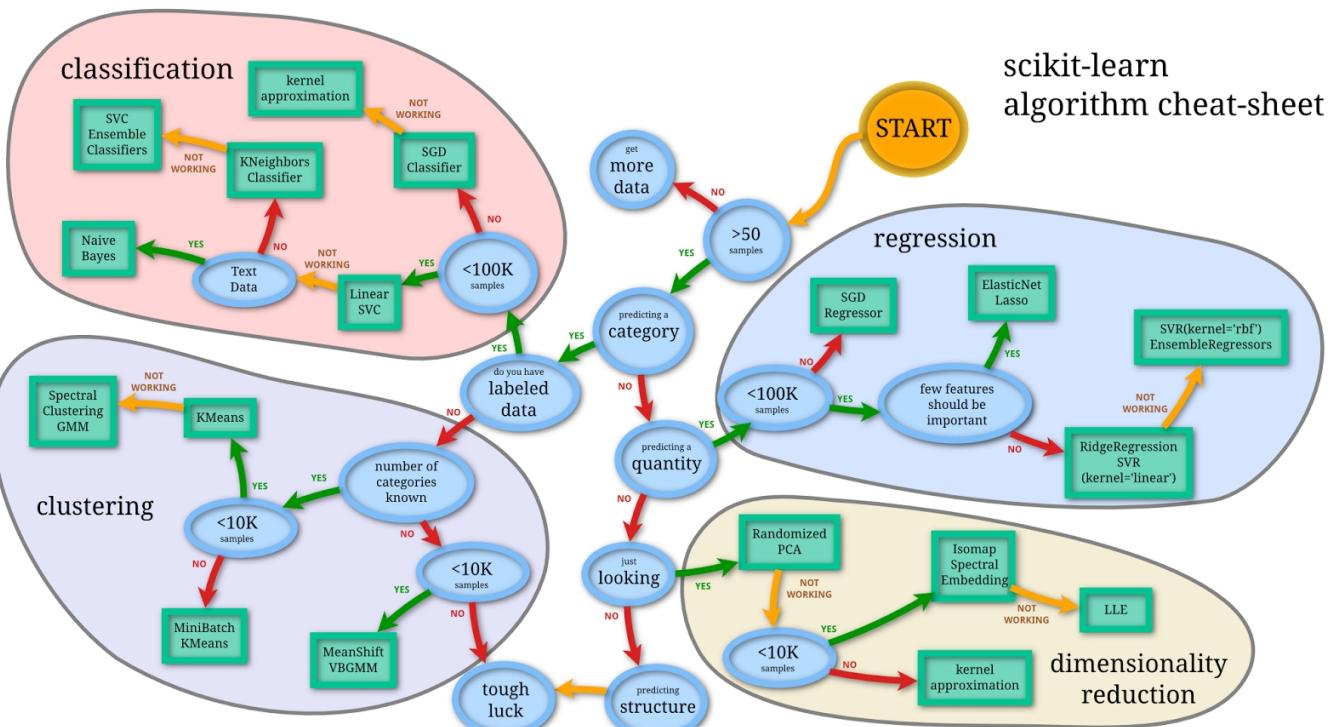
Subsetting	<code>a[2]</code> <code>3</code> <code>b[1,2]</code> <code>6,0</code>	Select the element at the 2nd index Select the element at row 0 column 2 (equivalent to <code>b[1][2]</code>) Select items at index 0 and 1 Select all items at row 0 (equivalent to <code>b[0,:]</code>) Same as <code>[1,:,:]</code>
Slicing	<code>a[0:2]</code> <code>array([1, 2])</code> <code>b[0:2,1]</code> <code>array([2, 5, 8])</code>	Select items at index 0 and 1 Select items at rows 0 and 1 in column 1 Select all items at row 0 (equivalent to <code>b[0,:,1]</code>) Same as <code>[1,:,1]</code>
Boolean Indexing	<code>a[[1, 0, 2, 1]]</code> <code>array([1, 0, 2, 1])</code>	Select elements from a less than 2
Fancy Indexing	<code>b[[1, 0, 1, 2, 0]]</code> <code>array([1, 0, 1, 2, 0])</code> <code>b[[1, 0, 1, 0, 1, 0, 1, 2, 0]]</code> <code>array([[1, 0, 2, 1, 0, 1, 2, 0]])</code>	Select elements <code>(1,0),(0,1),(1,2)</code> and <code>(0,0)</code> Select a subset of the matrix's rows and columns

Array Manipulation

Transposing Array	<code>a.T</code>	Permute array dimensions Permute array dimensions
Changing Array Shape	<code>a.ravel()</code> <code>g.reshape(3,-2)</code>	Flatten the array Reshape, but don't change data
Adding/Removing Elements	<code>a.resize((2,6))</code> <code>a[[1,2,3,4,5,6]]</code> <code>np.insert(a, 1, 5)</code> <code>np.delete(a, [1])</code>	Return a new array with shape <code>(2,6)</code> Append items to an array Insert items in an array Delete items from an array
Combining Arrays	<code>np.concatenate((a,d),axis=0)</code> <code>array([1, 2, 3, 10, 15, 20])</code> <code>np.vstack((a,b))</code> <code>array([[1, 2, 3, 4, 5, 6],</code> <code>[7, 8, 9, 10, 11, 12]])</code> <code>np.hstack((e,f))</code> <code>array([[1., 2., 3., 4., 5., 6.],</code> <code>[7., 8., 9., 10., 11., 12.]])</code> <code>np.column_stack((a,d))</code> <code>array([[1, 2, 3, 4, 5, 6],</code> <code>[7, 8, 9, 10, 11, 12],</code> <code>[1, 2, 3, 4, 5, 6],</code> <code>[7, 8, 9, 10, 11, 12]])</code> <code>np.e_(a,d)</code>	Concatenate arrays Stack arrays vertically (row-wise) Stack arrays vertically (row-wise) Stack arrays horizontally (column-wise) Create stacked column-wise arrays Create stacked column-wise arrays
Splitting Arrays	<code>np.hsplit(a,3)</code> <code>[array([1]),array([2]),array([3])]</code> <code>np.vsplit(e,2)</code> <code>[array([[1., 2., 3., 4., 5., 6.]]),</code> <code>[array([[7., 8., 9., 10., 11., 12.]])]</code> <code>np.dsplit(e,2)</code> <code>[array([[1., 2., 3., 4., 5., 6.],</code> <code>[7., 8., 9., 10., 11., 12.]])]</code>	Split the array horizontally at the 3rd index Split the array vertically at the 2nd index

DataCamp
Learn Python for Data Science interactively





Python For Data Science Cheat Sheet

Keras

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



Keras

Keras is a powerful and easy-to-use deep learning library for Theano and TensorFlow that provides a high-level neural networks API to develop and evaluate deep learning models.

A Basic Example

```
>>> import numpy as np
>>> from keras.models import Sequential
>>> from keras.layers import Dense
>>> data = np.random.random((1000,100))
>>> labels = np.random.randint(2,size=(1000,1))
>>> model = Sequential()
>>> model.add(Dense(32,
      activation='relu',
      input_dim=100))
>>> model.add(Dense(1, activation='sigmoid'))
>>> model.compile(optimizer='rmsprop',
      loss='binary_crossentropy',
      metrics=['accuracy'])
>>> model.fit(data,labels,epochs=10,batch_size=32)
>>> predictions = model.predict(data)
```

Data

Also see NumPy, Pandas & Scikit-Learn

Your data needs to be stored as NumPy arrays or as a list of NumPy arrays. Ideally, you split the data in training and test sets, for which you can also resort to the `train_test_split` module of `sklearn.cross_validation`.

Keras Data Sets

```
>>> from keras.datasets import boston_housing,
      mnist,
      cifar10,
      imdb
>>> (x_train,y_train),(x_test,y_test) = mnist.load_data()
>>> (x_train2,y_train2),(x_test2,y_test2) = boston_housing.load_data()
>>> (x_train3,y_train3),(x_test3,y_test3) = cifar10.load_data()
>>> (x_train4,y_train4),(x_test4,y_test4) = imdb.load_data(num_words=20000)
>>> num_classes = 10
```

Other

```
>>> from urllib.request import urlopen
>>> data = np.loadtxt(urlopen("http://archive.ics.uci.edu/ml/machine-learning-databases/pima-indians-diabetes/pima-indians-diabetes.data"),delimiter=",")
>>> data = data[:,0:-1]
>>> y = data[:,-1]
```

Preprocessing

Sequence Padding

```
>>> from keras.preprocessing import sequence
>>> x_train4 = sequence.pad_sequences(x_train4,maxlen=80)
>>> x_test4 = sequence.pad_sequences(x_test4,maxlen=80)
```

One-Hot Encoding

```
>>> from keras.utils import to_categorical
>>> y_train1 = to_categorical(y_train, num_classes)
>>> Y_test1 = to_categorical(y_test, num_classes)
>>> Y_train3 = to_categorical(y_train3, num_classes)
>>> Y_test3 = to_categorical(y_test3, num_classes)
```

Model Architecture

Sequential Model

```
>>> from keras.models import Sequential
>>> model = Sequential()
>>> model2 = Sequential()
>>> model3 = Sequential()
```

Multilayer Perceptron (MLP)

Binary Classification

```
>>> from keras.layers import Dense
>>> model.add(Dense(12,
      input_dim=8,
      kernel_initializer='uniform',
      activation='relu'))
>>> model.add(Dense(8,kernel_initializer='uniform',activation='relu'))
>>> model.add(Dense(1,kernel_initializer='uniform',activation='sigmoid'))
```

Multi-Class Classification

```
>>> from keras.layers import Dropout
>>> model.add(Dense(512,activation='relu',input_shape=(784,)))
>>> model.add(Dropout(0.2))
>>> model.add(Dense(512,activation='relu'))
>>> model.add(Dropout(0.2))
>>> model.add(Dense(10,activation='softmax'))
```

Regression

```
>>> model.add(Dense(64,activation='relu',input_dim=train_data.shape[1]))
>>> model.add(Dense(1))
```

Convolutional Neural Network (CNN)

```
>>> from keras.layers import Activation,Conv2D,MaxPooling2D,Flatten
>>> model2.add(Conv2D(32,(3,3),padding='same',input_shape=x_train.shape[1:]))
>>> model2.add(Activation('relu'))
>>> model2.add(MaxPooling2D(pool_size=(2,2)))
>>> model2.add(Flatten())
>>> model2.add(Dense(128,activation='relu'))
>>> model2.add(Dropout(0.25))
>>> model2.add(Dense(64,(3,3),padding='same'))
>>> model2.add(Activation('relu'))
>>> model2.add(Dense(64,(3,3)))
>>> model2.add(Activation('relu'))
>>> model2.add(Dense(1,activation='softmax'))
>>> 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'))
```

Also see NumPy & Scikit-Learn

Train and Test Sets

```
>>> from sklearn.model_selection import train_test_split
>>> X_train5,X_test5,y_train5,y_test5 = train_test_split(X,
      y,
      test_size=0.33,
      random_state=42)
```

Standardization/Normalization

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

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

Compile Model

>>> model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])	MLP: Binary Classification
>>> model.compile(optimizer='rmsprop', loss='categorical_crossentropy', metrics=['accuracy'])	MLP: Multi-Class Classification
>>> model.compile(optimizer='rmsprop', loss='mse', metrics=['mae'])	MLP: Regression
>>> model3.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])	Recurrent Neural Network

Model Training

```
>>> model3.fit(x_train4,
      y_train4,
      batch_size=32,
      epochs=15,
      verbose=1,
      validation_data=(x_test4,y_test4))
```

Evaluate Your Model's Performance

```
>>> score = model3.evaluate(x_test,
      y_test,
      batch_size=32)
```

Prediction

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

Save/Reload Models

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

Model Fine-tuning

Optimization Parameters

```
>>> from keras.optimizers import RMSprop
>>> opt = RMSprop(lr=0.0001, decay=1e-6)
>>> model2.compile(loss='categorical_crossentropy',
      optimizer=opt,
      metrics=['accuracy'])
```

Early Stopping

```
>>> from keras.callbacks import EarlyStopping
>>> early_stopping_monitor = EarlyStopping(patience=2)
>>> model3.fit(x_train4,
      y_train4,
      batch_size=32,
      epochs=15,
      validation_data=(x_test4,y_test4),
      callbacks=[early_stopping_monitor])
```



TensorFlow Cheat Sheet



About

TensorFlow

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

Skflow

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

Keras

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

Installation

How to install new package in Python:

```
pip install <package-name>
```

Example: `pip install requests`

How to install tensorflow?

```
device = cpu/gpu
```

```
python_version = cp27/cp34
```

```
sudo pip install
```

```
https://storage.googleapis.com/tensorflow/linux/$device/tensorflow-0.8.0-$python_version-none-linux_x86_64.whl
```

How to install Skflow

```
pip install sklearn
```

How to install Keras

```
pip install keras
```

update ~/.keras/keras.json - replace "theano" by "tensorflow"

Helpers

Python helper

Important functions

`type(object)`

Get object type

`help(object)`

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

`dir(object)`

Get list of object attributes (fields, functions)

`str(object)`

Transform an object to string

`object?`

Shows documentations about the object

`globals()`

Return the dictionary containing the current scope's global variables.

`locals()`

Update and return a dictionary containing the current scope's local variables.

`id(object)`

Return the identity of an object. This is guaranteed to be unique among simultaneously existing objects.

```
import __builtin__  
dir(__builtin__)
```

Other built-in functions

TensorFlow

Main classes

```
tf.Graph()  
tf.Operation()  
tf.Tensor()  
tf.Session()
```

Some useful functions

```
tf.get_default_session()  
tf.get_default_graph()  
tf.reset_default_graph()  
ops.reset_default_graph()  
tf.device("/cpu:0")  
tf.name_scope(value)  
tf.convert_to_tensor(value)
```

TensorFlow Optimizers

```
GradientDescentOptimizer  
AdadeltaOptimizer  
AdagradOptimizer  
MomentumOptimizer  
AdamOptimizer  
FtrlOptimizer  
RMSPropOptimizer
```

Reduction

```
reduce_sum  
reduce_prod  
reduce_min  
reduce_max  
reduce_mean  
reduce_all  
reduce_any  
accumulate_n
```

Activation functions

```
tf.nn?  
relu  
relu6  
elu  
softplus  
softsign  
dropout  
bias_add  
sigmoid  
tanh  
sigmoid_cross_entropy_with_logits  
softmax  
log_softmax  
softmax_cross_entropy_with_logits  
sparse_softmax_cross_entropy_with_logits  
weighted_cross_entropy_with_logits  
etc.
```

Skflow

Main classes

```
TensorFlowClassifier  
TensorFlowRegressor  
TensorFlowDNNClassifier  
TensorFlowDNNRegressor  
TensorFlowLinearClassifier  
TensorFlowLinearRegressor  
TensorFlowRNNClassifier  
TensorFlowRNNRegressor
```

TensorFlowEstimator

Each classifier and regressor have following fields

`n_classes=0` (Regressor), `n_classes` are expected to be input (Classifiers)

`batch_size=32`,

`steps=200`, // except

`TensorFlowRNNClassifier` - there is 50

`optimizer='Adagrad'`,

`learning_rate=0.1`,

`class_weight=None`,

`clip_gradients=5.0`,

`continue_training=False`,

`config=None`,

`verbose=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.

Useful links

TensorFlow API

https://www.tensorflow.org/versions/r0.9/api_docs/index.html

TensorFlow How-Tos

https://www.tensorflow.org/versions/r0.9/how_tos/index.html

Skflow github

<https://github.com/tensorflow/skflow>

Skflow API

<http://scikit-learn.org/stable/modules/classes.html#module-sklearn.cluster>

Keras

<http://keras.io/>

Our github

<https://github.com/Altoras/TensorFlow-training>

TensorFlow github

<https://github.com/tensorflow/tensorflow>

Python For Data Science Cheat Sheet

Pandas Basics

Learn Python for Data Science [Interactively](#) at [www.DataCamp.com](#)



Pandas

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



Use the following import convention:

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

Pandas Data Structures

Series

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



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

DataFrame

Columns

	Country	Capital	Population
1	Belgium	Brussels	11190846
2	India	New Delhi	1303171035
3	Brazil	Brasilia	207847528

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'])
```

I/O

Read and Write to CSV

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

Read and Write to Excel

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

Asking For Help

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

Selection

Also see NumPy Arrays

Getting

```
>>> s['b']
->
>>> df[1:]
Country Capital Population
1 India New Delhi 1303171035
2 Brazil Brasilia 207847528
```

Get one element

Get subset of a DataFrame

Selecting, Boolean Indexing & Setting

By Position

```
>>> df.iloc[[0], [0]]
Belgium
>>> df.iat[[0], [0]]
Belgium
```

Select single value by row & column

By Label

```
>>> df.loc[[0], ['Country']]
Belgium
>>> df.at[[0], ['Country']]
Belgium
```

Select single value by row & column labels

By Label/Position

```
>>> df.ix[2]
Country Brazil
Capital Brasilia
Population 207847528
>>> df.ix[:, 'Capital']
0 Brussels
1 New Delhi
2 Brasilia
```

Select single row of subset of rows

Setting

```
>>> df.ix[1, 'Capital'] = 'New Delhi'
```

Select a single column of subset of columns

Boolean Indexing

```
>>> s[s > 1]
s[(s < -1) | (s > 2)]
>>> df[df['Population'] > 1200000000]
```

Series s where value is not > 1

s where value is <-1 or > 2

Use filter to adjust DataFrame

Setting

```
>>> s['a'] = 6
```

Set index a of Series s to 6

Dropping

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

Drop values from rows (axis=0)

Drop values from columns (axis=1)

Sort & Rank

```
>>> df.sort_index(by='Country')
>>> s.order()
>>> df.rank()
```

Sort by row or column index

Sort a series by its values

Assign ranks to entries

Retrieving Series/DataFrame Information

Basic Information

>>> df.shape	(rows,columns)
>>> df.index	Describe index
>>> df.columns	Describe DataFrame columns
>>> df.info()	Info on DataFrame
>>> df.count()	Number of non-NA values

Summary

>>> df.sum()	Sum of values
>>> df.cumsum()	Cumulative sum of values
>>> df.min() / df.max()	Minimum/maximum values
>>> df.idmin() / df.idmax()	Minimum/Maximum index value
>>> df.describe()	Summary statistics
>>> df.mean()	Mean of values
>>> df.median()	Median of values

Applying Functions

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

Data Alignment

Internal Data Alignment

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

```
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])
>>> s3
a    10.0
b    NaN
c    5.0
d    7.0
```

Arithmetic Operations with Fill Methods

You can also do the internal data alignment yourself with the help of the fill methods:

```
>>> s.add(s3, fill_value=0)
a    10.0
b    -5.0
c     5.0
d     7.0
>>> s.sub(s3, fill_value=2)
s3
a    10.0
b    2.0
c    5.0
d    7.0
>>> s.div(s3, fill_value=4)
s3
a    10.0
b    2.0
c    5.0
d    7.0
>>> s.mul(s3, fill_value=3)
s3
a    10.0
b    2.0
c    5.0
d    7.0
```

DataCamp

Learn Python for Data Science [Interactively](#)

Also see NumPy

Linear Algebra

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

```
>>> from scipy import linalg, sparse
```

Creating Matrices

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

Basic Matrix Routines

Inverse	Inverse
>>> A.I	Inverse
>>> linalg.inv(A)	Inverse
Transposition	Transpose matrix
>>> A.T	Conjugate transpose
>>> A.H	Trace
Trace	Frobenius norm
>>> np.trace(A)	L1 norm (max column sum)
Norm	L1 norm (max row sum)
>>> linalg.norm(A)	Linf norm (max row sum)
Rank	Matrix rank
>>> np.linalg.matrix_rank(C)	Determinant
Inverse	Solver for dense matrices
>>> linalg.det(A)	Solver for dense matrices
Transposition	Least-squares solution to linear matrix equation
>>> A.T	Compute the pseudo-inverse of a matrix (least-squares solver)
Conjugate transpose	Compute the pseudo-inverse of a matrix (SVD)

Creating Sparse Matrices

>>> F = np.eye(3, k=1)	Create a 2x2 identity matrix
>>> G = np.mat(np.identity(2))	Create a 2x2 identity matrix
>>> C1 = np.zeros((3,3)) >>> H = sparse.csr_matrix(C1)	Compressed Sparse Row matrix
>>> I = sparse.csc_matrix(D) >>> J = sparse.dok_matrix(A) >>> E = np.mat(np.zeros((3,3))) >>> E.todense()	Compressed Sparse Column matrix
>>> K = sparse.bsr_matrix(F) >>> linalg.lsqsq(F, E)	Dictionary Of Keys matrix
>>> L = sparse.dok_matrix(A) >>> M = sparse.csc_matrix(I) >>> N = sparse.csr_matrix(J) >>> O = sparse.csc_matrix(K) >>> P = sparse.csr_matrix(L) >>> Q = sparse.csc_matrix(M) >>> R = sparse.bsr_matrix(N) >>> S = sparse.csr_matrix(O) >>> T = sparse.bsr_matrix(P) >>> U = sparse.csc_matrix(Q) >>> V = sparse.bsr_matrix(R) >>> W = sparse.csc_matrix(S) >>> X = sparse.bsr_matrix(T) >>> Y = sparse.csc_matrix(U) >>> Z = sparse.bsr_matrix(V) >>> A = sparse.csr_matrix(F) >>> B = sparse.csc_matrix(G) >>> C = sparse.bsr_matrix(H) >>> D = sparse.csc_matrix(I) >>> E = sparse.bsr_matrix(J) >>> F = sparse.csc_matrix(K) >>> G = sparse.bsr_matrix(L) >>> H = sparse.csc_matrix(M) >>> I = sparse.bsr_matrix(N) >>> J = sparse.csc_matrix(O) >>> K = sparse.bsr_matrix(P) >>> L = sparse.csc_matrix(Q) >>> M = sparse.bsr_matrix(R) >>> N = sparse.csc_matrix(S) >>> O = sparse.bsr_matrix(T) >>> P = sparse.csc_matrix(U) >>> Q = sparse.bsr_matrix(V) >>> R = sparse.csc_matrix(W) >>> S = sparse.bsr_matrix(X) >>> T = sparse.csc_matrix(Y) >>> U = sparse.bsr_matrix(Z)	Sparse matrix to full matrix
>>> sparse.isspmatrix_csc(A)	Identify sparse matrix

Sparse Matrix Routines

Inverse	Inverse
>>> sparse.linalg.inv(I)	Norm
Norm	Solver for sparse matrices
>>> sparse.linalg.norm(I)	Solver for sparse matrices
Solving linear problems	Solver for sparse matrices
>>> sparse.linalg.solve(H, I)	Solver for sparse matrices
>>> E = np.mat(np.zeros((3,3))) >>> F = np.mat(np.zeros((3,3))) >>> G = np.mat(np.zeros((3,3))) >>> H = np.mat(np.zeros((3,3))) >>> I = np.mat(np.zeros((3,3))) >>> J = np.mat(np.zeros((3,3))) >>> K = np.mat(np.zeros((3,3))) >>> L = np.mat(np.zeros((3,3))) >>> M = np.mat(np.zeros((3,3))) >>> N = np.mat(np.zeros((3,3))) >>> O = np.mat(np.zeros((3,3))) >>> P = np.mat(np.zeros((3,3))) >>> Q = np.mat(np.zeros((3,3))) >>> R = np.mat(np.zeros((3,3))) >>> S = np.mat(np.zeros((3,3))) >>> T = np.mat(np.zeros((3,3))) >>> U = np.mat(np.zeros((3,3))) >>> V = np.mat(np.zeros((3,3))) >>> W = np.mat(np.zeros((3,3))) >>> X = np.mat(np.zeros((3,3))) >>> Y = np.mat(np.zeros((3,3))) >>> Z = np.mat(np.zeros((3,3))) >>> A = np.mat(np.zeros((3,3))) >>> B = np.mat(np.zeros((3,3))) >>> C = np.mat(np.zeros((3,3))) >>> D = np.mat(np.zeros((3,3))) >>> E = np.mat(np.zeros((3,3))) >>> F = np.mat(np.zeros((3,3))) >>> G = np.mat(np.zeros((3,3))) >>> H = np.mat(np.zeros((3,3))) >>> I = np.mat(np.zeros((3,3))) >>> J = np.mat(np.zeros((3,3))) >>> K = np.mat(np.zeros((3,3))) >>> L = np.mat(np.zeros((3,3))) >>> M = np.mat(np.zeros((3,3))) >>> N = np.mat(np.zeros((3,3))) >>> O = np.mat(np.zeros((3,3))) >>> P = np.mat(np.zeros((3,3))) >>> Q = np.mat(np.zeros((3,3))) >>> R = np.mat(np.zeros((3,3))) >>> S = np.mat(np.zeros((3,3))) >>> T = np.mat(np.zeros((3,3))) >>> U = np.mat(np.zeros((3,3))) >>> V = np.mat(np.zeros((3,3))) >>> W = np.mat(np.zeros((3,3))) >>> X = np.mat(np.zeros((3,3))) >>> Y = np.mat(np.zeros((3,3))) >>> Z = np.mat(np.zeros((3,3))) >>> A = np.mat(np.zeros((3,3))) >>> B = np.mat(np.zeros((3,3))) >>> C = np.mat(np.zeros((3,3))) >>> D = np.mat(np.zeros((3,3))) >>> E = np.mat(np.zeros((3,3))) >>> F = np.mat(np.zeros((3,3))) >>> G = np.mat(np.zeros((3,3))) >>> H = np.mat(np.zeros((3,3))) >>> I = np.mat(np.zeros((3,3))) >>> J = np.mat(np.zeros((3,3))) >>> K = np.mat(np.zeros((3,3))) >>> L = np.mat(np.zeros((3,3))) >>> M = np.mat(np.zeros((3,3))) >>> N = np.mat(np.zeros((3,3))) >>> O = np.mat(np.zeros((3,3))) >>> P = np.mat(np.zeros((3,3))) >>> Q = np.mat(np.zeros((3,3))) >>> R = np.mat(np.zeros((3,3))) >>> S = np.mat(np.zeros((3,3))) >>> T = np.mat(np.zeros((3,3))) >>> U = np.mat(np.zeros((3,3))) >>> V = np.mat(np.zeros((3,3))) >>> W = np.mat(np.zeros((3,3))) >>> X = np.mat(np.zeros((3,3))) >>> Y = np.mat(np.zeros((3,3))) >>> Z = np.mat(np.zeros((3,3))) >>> A = np.mat(np.zeros((3,3))) >>> B = np.mat(np.zeros((3,3))) >>> C = np.mat(np.zeros((3,3))) >>> D = np.mat(np.zeros((3,3))) >>> E = np.mat(np.zeros((3,3))) >>> F = np.mat(np.zeros((3,3))) >>> G = np.mat(np.zeros((3,3))) >>> H = np.mat(np.zeros((3,3))) >>> I = np.mat(np.zeros((3,3))) >>> J = np.mat(np.zeros((3,3))) >>> K = np.mat(np.zeros((3,3))) >>> L = np.mat(np.zeros((3,3))) >>> M = np.mat(np.zeros((3,3))) >>> N = np.mat(np.zeros((3,3))) >>> O = np.mat(np.zeros((3,3))) >>> P = np.mat(np.zeros((3,3))) >>> Q = np.mat(np.zeros((3,3))) >>> R = np.mat(np.zeros((3,3))) >>> S = np.mat(np.zeros((3,3))) >>> T = np.mat(np.zeros((3,3))) >>> U = np.mat(np.zeros((3,3))) >>> V = np.mat(np.zeros((3,3))) >>> W = np.mat(np.zeros((3,3))) >>> X = np.mat(np.zeros((3,3))) >>> Y = np.mat(np.zeros((3,3))) >>> Z = np.mat(np.zeros((3,3))) >>> A = np.mat(np.zeros((3,3))) >>> B = np.mat(np.zeros((3,3))) >>> C = np.mat(np.zeros((3,3))) >>> D = np.mat(np.zeros((3,3))) >>> E = np.mat(np.zeros((3,3))) >>> F = np.mat(np.zeros((3,3))) >>> G = np.mat(np.zeros((3,3))) >>> H = np.mat(np.zeros((3,3))) >>> I = np.mat(np.zeros((3,3))) >>> J = np.mat(np.zeros((3,3))) >>> K = np.mat(np.zeros((3,3))) >>> L = np.mat(np.zeros((3,3))) >>> M = np.mat(np.zeros((3,3))) >>> N = np.mat(np.zeros((3,3))) >>> O = np.mat(np.zeros((3,3))) >>> P = np.mat(np.zeros((3,3))) >>> Q = np.mat(np.zeros((3,3))) >>> R = np.mat(np.zeros((3,3))) >>> S = np.mat(np.zeros((3,3))) >>> T = np.mat(np.zeros((3,3))) >>> U = np.mat(np.zeros((3,3))) >>> V = np.mat(np.zeros((3,3))) >>> W = np.mat(np.zeros((3,3))) >>> X = np.mat(np.zeros((3,3))) >>> Y = np.mat(np.zeros((3,3))) >>> Z = np.mat(np.zeros((3,3))) >>> A = np.mat(np.zeros((3,3))) >>> B = np.mat(np.zeros((3,3))) >>> C = np.mat(np.zeros((3,3))) >>> D = np.mat(np.zeros((3,3))) >>> E = np.mat(np.zeros((3,3))) >>> F = np.mat(np.zeros((3,3))) >>> G = np.mat(np.zeros((3,3))) >>> H = np.mat(np.zeros((3,3))) >>> I = np.mat(np.zeros((3,3))) >>> J = np.mat(np.zeros((3,3))) >>> K = np.mat(np.zeros((3,3))) >>> L = np.mat(np.zeros((3,3))) >>> M = np.mat(np.zeros((3,3))) >>> N = np.mat(np.zeros((3,3))) >>> O = np.mat(np.zeros((3,3))) >>> P = np.mat(np.zeros((3,3))) >>> Q = np.mat(np.zeros((3,3))) >>> R = np.mat(np.zeros((3,3))) >>> S = np.mat(np.zeros((3,3))) >>> T = np.mat(np.zeros((3,3))) >>> U = np.mat(np.zeros((3,3))) >>> V = np.mat(np.zeros((3,3))) >>> W = np.mat(np.zeros((3,3))) >>> X = np.mat(np.zeros((3,3))) >>> Y = np.mat(np.zeros((3,3))) >>> Z = np.mat(np.zeros((3,3))) >>> A = np.mat(np.zeros((3,3))) >>> B = np.mat(np.zeros((3,3))) >>> C = np.mat(np.zeros((3,3))) >>> D = np.mat(np.zeros((3,3))) >>> E = np.mat(np.zeros((3,3))) >>> F = np.mat(np.zeros((3,3))) >>> G = np.mat(np.zeros((3,3))) >>> H = np.mat(np.zeros((3,3))) >>> I = np.mat(np.zeros((3,3))) >>> J = np.mat(np.zeros((3,3))) >>> K = np.mat(np.zeros((3,3))) >>> L = np.mat(np.zeros((3,3))) >>> M = np.mat(np.zeros((3,3))) >>> N = np.mat(np.zeros((3,3))) >>> O = np.mat(np.zeros((3,3))) >>> P = np.mat(np.zeros((3,3))) >>> Q = np.mat(np.zeros((3,3))) >>> R = np.mat(np.zeros((3,3))) >>> S = np.mat(np.zeros((3,3))) >>> T = np.mat(np.zeros((3,3))) >>> U = np.mat(np.zeros((3,3))) >>> V = np.mat(np.zeros((3,3))) >>> W = np.mat(np.zeros((3,3))) >>> X = np.mat(np.zeros((3,3))) >>> Y = np.mat(np.zeros((3,3))) >>> Z = np.mat(np.zeros((3,3))) >>> A = np.mat(np.zeros((3,3))) >>> B = np.mat(np.zeros((3,3))) >>> C = np.mat(np.zeros((3,3))) >>> D = np.mat(np.zeros((3,3))) >>> E = np.mat(np.zeros((3,3))) >>> F = np.mat(np.zeros((3,3))) >>> G = np.mat(np.zeros((3,3))) >>> H = np.mat(np.zeros((3,3))) >>> I = np.mat(np.zeros((3,3))) >>> J = np.mat(np.zeros((3,3))) >>> K = np.mat(np.zeros((3,3))) >>> L = np.mat(np.zeros((3,3))) >>> M = np.mat(np.zeros((3,3))) >>> N = np.mat(np.zeros((3,3))) >>> O = np.mat(np.zeros((3,3))) >>> P = np.mat(np.zeros((3,3))) >>> Q = np.mat(np.zeros((3,3))) >>> R = np.mat(np.zeros((3,3))) >>> S = np.mat(np.zeros((3,3))) >>> T = np.mat(np.zeros((3,3))) >>> U = np.mat(np.zeros((3,3))) >>> V = np.mat(np.zeros((3,3))) >>> W = np.mat(np.zeros((3,3))) >>> X = np.mat(np.zeros((3,3))) >>> Y = np.mat(np.zeros((3,3))) >>> Z = np.mat(np.zeros((3,3))) >>> A = np.mat(np.zeros((3,3))) >>> B = np.mat(np.zeros((3,3))) >>> C = np.mat(np.zeros((3,3))) >>> D = np.mat(np.zeros((3,3))) >>> E = np.mat(np.zeros((3,3))) >>> F = np.mat(np.zeros((3,3))) >>> G = np.mat(np.zeros((3,3))) >>> H = np.mat(np.zeros((3,3))) >>> I = np.mat(np.zeros((3,3))) >>> J = np.mat(np.zeros((3,3))) >>> K = np.mat(np.zeros((3,3))) >>> L = np.mat(np.zeros((3,3))) >>> M = np.mat(np.zeros((3,3))) >>> N = np.mat(np.zeros((3,3))) >>> O = np.mat(np.zeros((3,3))) >>> P = np.mat(np.zeros((3,3))) >>> Q = np.mat(np.zeros((3,3))) >>> R = np.mat(np.zeros((3,3))) >>> S = np.mat(np.zeros((3,3))) >>> T = np.mat(np.zeros((3,3))) >>> U = np.mat(np.zeros((3,3))) >>> V = np.mat(np.zeros((3,3))) >>> W = np.mat(np.zeros((3,3))) >>> X = np.mat(np.zeros((3,3))) >>> Y = np.mat(np.zeros((3,3))) >>> Z = np.mat(np.zeros((3,3))) >>> A = np.mat(np.zeros((3,3))) >>> B = np.mat(np.zeros((3,3))) >>> C = np.mat(np.zeros((3,3))) >>> D = np.mat(np.zeros((3,3))) >>> E = np.mat(np.zeros((3,3))) >>> F = np.mat(np.zeros((3,3))) >>> G = np.mat(np.zeros((3,3))) >>> H = np.mat(np.zeros((3,3))) >>> I = np.mat(np.zeros((3,3))) >>> J = np.mat(np.zeros((3,3))) >>> K = np.mat(np.zeros((3,3))) >>> L = np.mat(np.zeros((3,3))) >>> M = np.mat(np.zeros((3,3))) >>> N = np.mat(np.zeros((3,3))) >>> O = np.mat(np.zeros((3,3))) >>> P = np.mat(np.zeros((3,3))) >>> Q = np.mat(np.zeros((3,3))) >>> R = np.mat(np.zeros((3,3))) >>> S = np.mat(np.zeros((3,3))) >>> T = np.mat(np.zeros((3,3))) >>> U = np.mat(np.zeros((3,3))) >>> V = np.mat(np.zeros((3,3))) >>> W = np.mat(np.zeros((3,3))) >>> X = np.mat(np.zeros((3,3))) >>> Y = np.mat(np.zeros((3,3))) >>> Z = np.mat(np.zeros((3,3))) >>> A = np.mat(np.zeros((3,3))) >>> B = np.mat(np.zeros((3,3))) >>> C = np.mat(np.zeros((3,3))) >>> D = np.mat(np.zeros((3,3))) >>> E = np.mat(np.zeros((3,3))) >>> F = np.mat(np.zeros((3,3))) >>> G = np	

Data Wrangling

with pandas

Cheat Sheet

<http://pandas.pydata.org>

Syntax – Creating DataFrames

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

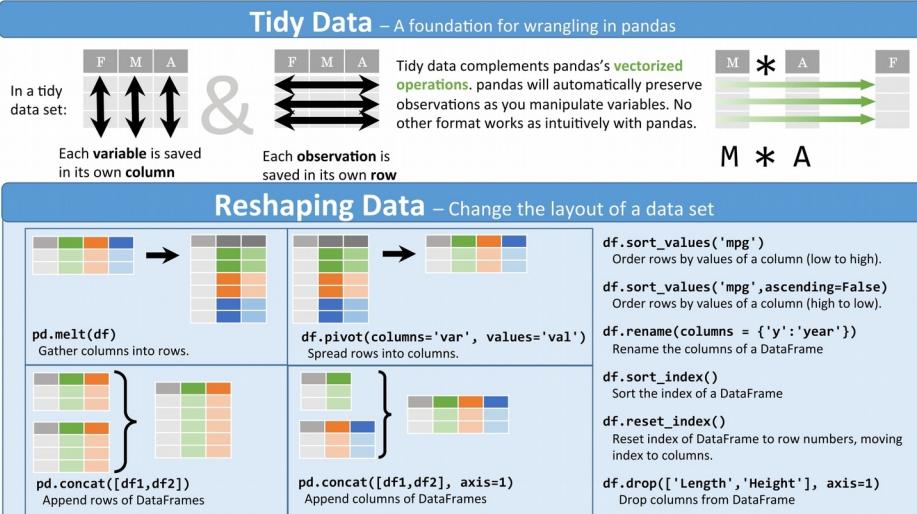
df = pd.DataFrame([
    [4, 7, 10],
    [5, 8, 11],
    [6, 9, 12]],
index=[1, 2, 3],
columns=['a', 'b', 'c'])
Specify values for each row.

df = pd.DataFrame([
    {"a": 4, "b": 7, "c": 10},
    {"b": 5, "c": 8, "a": 11},
    {"c": 6, "a": 9, "b": 12}],
index=[1, 2, 3],
columns=['a', 'b', 'c'])
Create DataFrame with a MultiIndex
```

Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

```
df = (pd.melt(df)
      .rename(columns={
          'variable' : 'var',
          'value' : 'val'})
      .query('val > 200')
     )
```



Reshaping Data – Change the layout of a data set

pd.melt(df) Gather columns into rows.

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

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

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

Subset Observations (Rows)

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

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

df.head(n) Select first n rows.

df.tail(n) Select last n rows.

Subset Variables (Columns)

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

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

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

regex (Regular Expressions) Examples	
'\.'	Matches strings containing a period ''
'Length\$'	Matches strings ending with word 'Length'
'Sepal'	Matches strings beginning with the word 'Sepal'
'x[1-5]\$'	Matches strings beginning with 'x' and ending with 1,2,3,4,5
'^!(?i:Species)\$'."	Matches strings except the string 'Species'

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

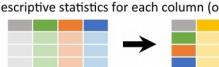
df.iloc[:, [1,2,5]] Select columns in positions 1, 2 and 5 (first column is 0).

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

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

Summarize Data

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



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

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

Group Data

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

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

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

size() Size of each group.

agg(function) Aggregate group using function.

Handling Missing Data

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

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

Make New Columns

df.assign(Area=lambda df: df.Length*df.Height) Compute and append one or more new columns.

df['Volume'] = df.Length*df.Height*df.Depth Add single column.

pd.qcut(df.col, n, labels=False) Bin column into n buckets.

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

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

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

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

Windows

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

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

df.plot.hist() Histogram for each column

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

Combine Data Sets

adf + bdf = Standard Joins

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

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

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

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

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

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

ydf + zdf = Set-like Operations

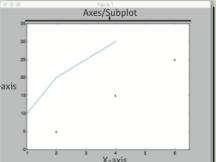
pd.merge(ydf, zdf) Rows that appear in both ydf and zdf (Intersection).

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

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

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

Python For Data Science Cheat Sheet					
<h2>Scikit-Learn</h2> <p>Learn Python for data science interactively at www.DataCamp.com</p> 	<h2>Create Your Model</h2> <h3>Supervised Learning Estimators</h3> <pre> 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) </pre> <h3>Unsupervised Learning Estimators</h3> <pre> 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) </pre>				
<h2>Scikit-learn</h2> <p>Scikit-learn is an open source Python library that implements a range of machine learning, preprocessing, cross-validation and visualization algorithms using a unified interface.</p>	<h2>Evaluate Your Model's Performance</h2> <h3>Classification Metrics</h3> <p>Accuracy Score</p> <pre>>>> knn.score(X_test, y_test) >>> from sklearn.metrics import accuracy_score >>> accuracy_score(y_test, y_pred)</pre> <p>Estimator score method Metric scoring functions</p> <p>Classification Report</p> <pre>>>> from sklearn.metrics import classification_report >>> print(classification_report(y_test, y_pred))</pre> <p>Precision, recall, f-score and support</p> <p>Confusion Matrix</p> <pre>>>> from sklearn.metrics import confusion_matrix >>> print(confusion_matrix(y_test, y_pred))</pre>				
<h2>A Basic Example</h2> <pre> >>> from sklearn import neighbors, datasets, preprocessing >>> from sklearn.cross_validation import train_test_split >>> from sklearn import metrics, linear_model, accuracy_score >>> iris = datasets.load_iris() >>> X, y = iris.data[:, :2], iris.target >>> X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=33) >>> 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) </pre>	<h3>Regression Metrics</h3> <p>Mean Absolute Error</p> <pre>>>> from sklearn.metrics import mean_absolute_error >>> y_true = [3, -0.5, 2] >>> mean_absolute_error(y_true, y_pred)</pre> <p>Mean Squared Error</p> <pre>>>> from sklearn.metrics import mean_squared_error >>> mean_squared_error(y_test, y_pred)</pre> <p>R² Score</p> <pre>>>> from sklearn.metrics import r2_score >>> r2_score(y_true, y_pred)</pre>				
<h2>Loading The Data</h2> <p>Also see NumPy & Pandas</p> <p>Your data needs to be numeric and stored as NumPy arrays or SciPy sparse matrices. Other types that are convertible to numeric arrays, such as Pandas DataFrame, are also acceptable.</p> <pre> >>> import numpy as np >>> np.__version__ >>> y = np.array(['M', 'M', 'F', 'F', 'M', 'F', 'M', 'M', 'F', 'F', 'F']) >>> X[(X < 0.7) == 0] </pre>	<h2>Model Fitting</h2> <table border="1"> <tr> <td> Supervised learning <pre> >>> lr.fit(X, y) >>> knn.fit(X_train, y_train) >>> svc.fit(X_train, y_train) </pre> </td><td> Fit the model to the data </td></tr> <tr> <td> Unsupervised Learning <pre> >>> k_means.fit(X_train) >>> pca_model = pca.fit_transform(X_train) </pre> </td><td> Fit the model to the data Fit to data, then transform it </td></tr> </table>	Supervised learning <pre> >>> lr.fit(X, y) >>> knn.fit(X_train, y_train) >>> svc.fit(X_train, y_train) </pre>	Fit the model to the data	Unsupervised Learning <pre> >>> k_means.fit(X_train) >>> pca_model = pca.fit_transform(X_train) </pre>	Fit the model to the data Fit to data, then transform it
Supervised learning <pre> >>> lr.fit(X, y) >>> knn.fit(X_train, y_train) >>> svc.fit(X_train, y_train) </pre>	Fit the model to the data				
Unsupervised Learning <pre> >>> k_means.fit(X_train) >>> pca_model = pca.fit_transform(X_train) </pre>	Fit the model to the data Fit to data, then transform it				
<h2>Training And Test Data</h2> <pre> >>> from sklearn.cross_validation import train_test_split >>> X_train, X_test, y_train, y_test = train_test_split(X, ... y, ... random_state=0) </pre>	<h2>Prediction</h2> <table border="1"> <tr> <td> Supervised Estimators <pre> >>> y_pred = svc.predict(np.random((2, 5))) >>> y_pred = lr.predict(X_test) >>> y_pred = knn.predict_proba(X_test) </pre> </td><td> Predict labels Predict labels Estimate probability of a label </td></tr> <tr> <td> Unsupervised Estimators <pre> >>> y_pred = k_means.predict(X_test) </pre> </td><td> Predict labels in clustering algos </td></tr> </table>	Supervised Estimators <pre> >>> y_pred = svc.predict(np.random((2, 5))) >>> y_pred = lr.predict(X_test) >>> y_pred = knn.predict_proba(X_test) </pre>	Predict labels Predict labels Estimate probability of a label	Unsupervised Estimators <pre> >>> y_pred = k_means.predict(X_test) </pre>	Predict labels in clustering algos
Supervised Estimators <pre> >>> y_pred = svc.predict(np.random((2, 5))) >>> y_pred = lr.predict(X_test) >>> y_pred = knn.predict_proba(X_test) </pre>	Predict labels Predict labels Estimate probability of a label				
Unsupervised Estimators <pre> >>> y_pred = k_means.predict(X_test) </pre>	Predict labels in clustering algos				
<h2>Preprocessing The Data</h2> <h3>Standardization</h3> <pre> >>> from sklearn.preprocessing import StandardScaler >>> scalar = StandardScaler().fit(X_train) >>> standardized_X = scalar.transform(X_train) >>> standardized_X_test = scalar.transform(X_test) </pre>	<h3>Encoding Categorical Features</h3> <pre> >>> from sklearn.preprocessing import LabelEncoder >>> enc = LabelEncoder() >>> y = enc.fit_transform(y) </pre>				
<h3>Normalization</h3> <pre> >>> from sklearn.preprocessing import Normalizer >>> scaler = Normalizer().fit(X_train) >>> normalized_X = scaler.transform(X_train) >>> normalized_X_test = scaler.transform(X_test) </pre>	<h3>Imputing Missing Values</h3> <pre> >>> from sklearn.preprocessing import Imputer >>> imp = Imputer(missing_values=0, strategy='mean', axis=0) >>> imp.fit_transform(X_train) </pre>				
<h3>Binarization</h3> <pre> >>> from sklearn.preprocessing import Binarizer >>> binarizer = Binarizer(threshold=0.0).fit(X) >>> binary_X = binarizer.transform(X) </pre>	<h3>Generating Polynomial Features</h3> <pre> >>> from sklearn.preprocessing import PolynomialFeatures >>> poly = PolynomialFeatures(5) >>> poly.fit_transform(X) </pre>				
<h2>Grid Search</h2> <p>Grid Search</p> <pre> >>> from sklearn.grid_search import GridSearchCV >>> params = {'n_neighbors': np.arange(1, 3), ... 'metric': ['euclidean', 'cityblock']} >>> grid = GridSearchCV(knn, params) >>> grid.fit(X_train, y_train) >>> print(grid.best_score_) >>> print(grid.best_estimator_.n_neighbors) </pre>					
<h2>Randomized Parameter Optimization</h2> <p>Randomized Parameter Optimization</p> <pre> >>> from sklearn.grid_search import RandomizedSearchCV >>> params = {'n_neighbors': range(1, 5), ... 'weights': ['uniform', 'distance']} >>> search = RandomizedSearchCV(knn, ... param_distributions=params, ... n_iter=8, ... random_state=5) >>> search.fit(X_train, y_train) >>> print(search.best_score_) </pre>					
<p>Learn Python for Data Science Interactively</p> 					

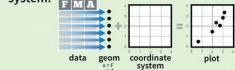
Python For Data Science Cheat Sheet	
Matplotlib	
<p>Learn Python Interactively at www.DataCamp.com</p> 	
<h2>Matplotlib</h2> <p>Matplotlib is a Python 2D plotting library which produces publication-quality figures in a variety of hardcopy formats and interactive environments across platforms.</p> 	<h2>Plot Anatomy & Workflow</h2>
<h3>Plot Anatomy</h3>  <p>Figure</p>	<h3>Workflow</h3> <p>The basic steps to creating plots with matplotlib are:</p> <ol style="list-style-type: none"> 1 Prepare data 2 Create plot 3 Plot 4 Customize plot 5 Save plot 6 Show plot <pre>>>> import matplotlib.pyplot as plt >>> x = [1,2,3,4] >>> y = [10,20,25,30] >>> fig = plt.figure() >>> ax = fig.add_subplot(111) >>> ax.plot(x, y, color='lightblue', linewidth=3) # Step 3.d >>> ax.scatter([2,4,6], [15,15,25], color='darkgreen', s=[100,200,50]) >>> ax.set_xlim(1, 6.5) >>> plt.savefig('foo.png') # Step 5 >>> plt.show() # Step 6</pre>
<h3>1 Prepare The Data</h3> <p>Also see Lists & NumPy</p> <p>1D Data</p> <pre>>>> import numpy as np >>> x = np.linspace(0, 10, 100) >>> y = np.sin(x) >>> z = np.sin(x)</pre> <p>2D Data or Images</p> <pre>>>> data = 2 * np.random.random((10, 10)) >>> data2 = 2 * np.random.random((10, 10)) >>> X, Y = np.meshgrid(np.linspace(-3, 100), np.linspace(-3, 100)) >>> U, V = -Y * X**2, -X**2 >>> from matplotlib.cbook import get_sample_data >>> img = np.load(get_sample_data('axes_grid/bivariate_normal.npy'))</pre>	<h3>4 Customize Plot</h3> <p>Colors, Color Bars & Color Maps</p> <pre>>>> plt.plot(x, x, x**2, x, x**3) >>> plt.plot(x, y, alpha=0.4) >>> ax.plot(x, y, c='k') >>> fig.colorbar(im, orientation='horizontal') >>> im = ax.imshow(img, cmap='seismic')</pre> <p>Markers</p> <pre>>>> fig, ax = plt.subplots() >>> ax.plot(x,y,marker="*") >>> ax.plot(x,y,marker="o")</pre> <p>Linestyles</p> <pre>>>> plt.plot(x,y,linewidth=4.0) >>> plt.plot(x,y,ls="solid") >>> plt.plot(x,y,ls="--") >>> plt.plot(x,y,ls="x*x2,y**2,-") >>> plt.setp(lines,color='r',linewidth=4.0)</pre> <p>Text & Annotations</p> <pre>>>> ax.text(1, -2.1, 'Example Graph', style='italic') >>> ax.annotate("S1", xy=(8, 0), xycoords='data', xytext=(10.5, 0), textcoords="data", arrowprops=dict(arrowsize=->, connectionstyle="arc3"),)</pre>
<h3>2 Create Plot</h3> <p>Also see Matplotlib</p> <p>Figure</p> <pre>>>> fig = plt.figure() >>> fig2 = plt.figure(figsize=plt.figaspect(2.0))</pre> <p>Axes</p> <p>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.</p> <pre>>>> fig, add_axes() >>> ax1 = fig.add_subplot(221) # row-col-num >>> ax2 = fig.add_subplot(212) >>> fig3, axes = plt.subplots(mrows=2, ncols=2) >>> fig4, axes2 = plt.subplots(ncols=3)</pre>	<p>Mathtext</p> <pre>>>> plt.title(r'\$\sigma_i=15\$', fontsize=20)</pre> <p>Limits, Legends & Layouts</p> <pre>>>> ax.margins(x=0.0, y=0.1) >>> ax.xaxis.set('equal') >>> ax.set(xlim=[0,10.5], ylim=[-1.5,1.5]) >>> ax.set_xlim(0,10.5)</pre> <p>Legend</p> <pre>>>> ax.set(title='An Example Axes', xlabel='X-Axis', ylabel='Y-Axis') >>> ax.legend(loc='best')</pre> <p>Ticks</p> <pre>>>> ax.xaxis.set(ticks=range(1,5), ticklabels=[3,100,-12,"foo"]) >>> ax.tick_params(labelsize='y', direction='inout', length=10)</pre> <p>Subplot Spacing</p> <pre>>>> fig.subplots_adjust(wspace=0.5, hspace=0.3, left=0.125, right=0.9, top=0.5, bottom=0.1)</pre> <p>Axis Spines</p> <pre>>>> ax1.spines['top'].set_visible(False) >>> ax1.spines['bottom'].set_position('outward',10)</pre>
<h3>3 Plotting Routines</h3> <p>1D Data</p> <pre>>>> fig, ax = plt.subplots() >>> lines = ax.plot(x,y) >>> ax.scatter(x,y) >>> axes[0,0].bar([1,2,3],[3, 4, 5]) >>> axes[1,0].bar([0,1,2,3],[0, 1, 2, 1]) >>> axes[0,1].bar([0,1,2,3],[0, 1, 2, 1]) >>> axes[0,1].axline(0,0.65) >>> ax.fill(x,y,color='blue') >>> ax.fill_between(x,y,color='yellow')</pre> <p>2D Data or Images</p> <pre>>>> fig, ax = plt.subplots() >>> im = ax.imshow(img, cmap='gist_earth', interpolation='nearest', vmin=-2, vmax=2)</pre>	<p>Vector Fields</p> <pre>>>> axes[0,1].arrow(0,0,0.5,0.5) >>> axes[1,1].quiver(y,z) >>> axes[0,1].streamplot(X,Y,u,v)</pre> <p>Data Structures</p> <pre>>>> ax1.hist(y) >>> ax3.boxplot(y) >>> ax3.violinplot(z)</pre> <p>Pseudocolor</p> <pre>>>> axes[0,0].pcolor(data2) >>> axes[0,0].pcolormesh(data2) >>> CS = plt.contour(Y,X,U) >>> axes[2,0].contourf(data1) >>> axes[2,0].axlabel(CS) >>> axes[2,0].clabel(CS)</pre>
<p>Colormapped or RGB arrays</p>	<p>Save figures</p> <pre>>>> plt.savefig('foo.png')</pre> <p>Save transparent figures</p> <pre>>>> plt.savefig('foo.png', transparent=True)</pre>
<p>Close & Clear</p>	<p>Clear an axis</p> <pre>>>> plt.cla()</pre> <p>Clear the entire figure</p> <pre>>>> plt.clf()</pre> <p>Close a window</p> <pre>>>> plt.close()</pre>

Data Visualization with ggplot2 Cheat Sheet

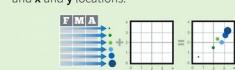


Basics

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



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



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

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

ggplot(data = mpg, aes(x = cyl, y = hwy))

Begins a plot that you finish by adding layers to. No defaults, but provides more control than **plot()**.

```
data
ggplot(mpg, aes(hvy, cty)) +
  geom_point(aes(color = cyl)) +
  geom_smooth(method = "lm") +
  coord_cartesian() +
  scale_color_gradient() +
  theme_bw()
  add layers elements with +
  layer specific mappings
  default stat
  layer specific
  mappings
  additional elements
```

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

last_plot()

Returns the last plot

ggsave("plot.png", width = 5, height = 5)

Saves last plot as 5 x 5" file named "plot.png" in working directory. Matches file type to file extension.

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

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

One Variable

Continuous

```
a + ggplot(mpg, aes(hwy))
  a + geom_area(stat = "bin")
    x, y, alpha, color, fill, linetype, size
    b + geom_area(aes(y = ..density..), stat = "bin")
      a + geom_density(kern = "gaussian")
        x, y, alpha, color, fill, linetype, size, weight
        b + geom_density(aes(y = ..count..))
          a + geom_dotplot()
            x, y, alpha, color, fill
            b + geom_freqpoly()
              x, y, alpha, color, linetype, size
              b + geom_freqpoly(aes(y = ..density..))
                a + geom_histogram(binwidth = 5)
                  x, y, alpha, color, fill, linetype, size, weight
                  b + geom_histogram(aes(y = ..density..))
```

Discrete

```
b + ggplot(mpg, aes(fl))
  b + geom_bar()
    x, alpha, color, fill, linetype, size, weight
```

Graphical Primitives

```
c < ggplot(map, aes(long, lat))
  c + geom_polygon(aes(group = group))
    x, y, alpha, color, fill, linetype, size
```

```
d < ggplot(economics, aes(date, unemploy))
  d + geom_path(lineend = "butt",
    linejoin = "round", linemitre = 1)
    x, y, alpha, color, linetype, size
  d + geom_ribbon(aes(ymin = unemploy - 900,
    ymax = unemploy + 900))
    x, ymax, ymin, alpha, color, fill, linetype, size
```

```
e < ggplot(seals, aes(x = long, y = lat))
  e + geom_segment(aes(
    xend = long + delta_long,
    yend = lat + delta_lat))
    x, end, y, end, alpha, color, linetype, size
  e + geom_rect(aes(xmin = long, ymin = lat,
    xmax = long + delta_long,
    ymax = lat + delta_lat))
    xmax, xmin, ymax, ymin, alpha, color, fill,
    linetype, size
```

```
f + geom_pointrange()
  f + geom_jitter()
    x, y, alpha, color, fill, shape, size
  f + geom_text(aes(label = cty))
    x, y, label, alpha, angle, color, family, fontface,
    hjust, lineheight, size, vjust
```

Two Variables

Continuous X, Continuous Y

```
f < ggplot(mpg, aes(cty, hwy))
  f + geom_blank()
  f + geom_jitter()
    x, y, alpha, color, fill, shape, size
  f + geom_point()
    x, y, alpha, color, fill, shape, size
  f + geom_quantile()
    x, y, alpha, color, linetype, size, weight
  f + geom_rug(sides = "bl")
    alpha, color, linetype, size
  f + geom_smooth(model = lm)
    x, y, alpha, color, fill, linetype, size, weight
  C f + geom_text(aes(label = cty))
    x, y, label, alpha, angle, color, family, fontface,
    hjust, lineheight, size, vjust
```

Discrete X, Continuous Y

```
g < ggplot(mpg, aes(class, hwy))
  g + geom_bar(stat = "identity")
    x, y, alpha, color, fill, linetype, size, weight
  g + geom_boxplot()
    lower, middle, upper, x, ymin, alpha, color, fill,
    linetype, shape, size, weight
  g + geom_dotplot(binxaxis = "y",
    stackdir = "center")
    x, y, alpha, color, fill
  g + geom_violin(scale = "area")
    x, y, alpha, color, fill, linetype, size, weight
```

Discrete X, Discrete Y

```
h < ggplot(diamonds, aes(cut, color))
  h + geom_jitter()
    x, y, alpha, color, fill, shape, size
```

Three Variables

```
sealsSz <- with(seals, sqrt(delta_long^2 + delta_lat^2))
m < ggplot(seals, aes(long, lat))
  m + geom_raster(aes(fill = z), hjust = 0.5,
    vjust = 0.5, interpolate = FALSE)
    x, y, alpha, fill
  m + geom_contour(aes(z = z))
    x, y, z, alpha, colour, linetype, size, weight
```

Visualizing error

```
df < data.frame(grp = c("A", "B"))
f < ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
  k + geom_crossbar(fatten = 2)
    x, y, ymax, ymin, alpha, color, fill, linetype,
    size
```

Maps

```
data < data.frame(murder = USArrests$Murder,
  state = tolower(rownames(USArrests)))
map <- map_data("state")
l < ggplot(data, aes(fill = murder))
  l + geom_map(aes(map_id = state), map = map) +
  expand_limits(x = map$long, y = map$lat)
  map_id, alpha, color, fill, linetype, size
```

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

Stats

Stats - An alternative way to build a layer

Some plots visualize a **transformation** of the original data set.

Use a **stat** to choose a common transformation to visualize, e.g. **a + geom_bar(stat = "bin")**

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

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

1D distributions

```
a + stat_bin(binwidth = 1, origin = 10)
x, y, fill | count, ..count.., density, ..density..
a + stat_bin(binwidth = 1, bins = 20)
x, y, fill | count, density..
```

```
a + stat_density(adjust = 1, kernel = "gaussian")
x, y, ..scaled.., ..scaled..,
```

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

```
f + stat_hex(bins = 30)
x, y, fill | count, density..
```

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

```
m + stat_contour(aes(z = z))
x, y, z, fill | ..level..
```

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

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

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

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

```
g + stat_ydensity(aes(y = y, ..density.. = ..density..),
  ..density.., ..count.., ..n.., ..violinwidth.., ..width..)
```

```
f + stat_ecdf(p = 40)
Functions
x, y | ..x.., ..y..
```

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

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

```
ggplot() + stat_function(aes(x = -3, y = 0))
General Purpose
fun = dnorm, n = 101, args = list(d = 0.5)
x | ..x.., ..y..
```

```
f + stat_identity()
ggplot() + stat_qq(aes(sample = 1:100), distribution = qt,
  dparams = list(fit = 0))
sample, x, y | ..x.., ..y..
```

```
f + stat_kde()
x, y | ..size..
```

```
f + stat_summary(fun.data = "mean_cl_boot")
x, y | ..size..
```

```
f + stat_uniques()

```

Scales

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

range of values to include in mapping
title to use in legend/axis
labels to use in legend/axis
breaks to use in legend/axis

General Purpose scales

Use with any aesthetic: alpha, color, fill, linetype, shape, size

```
scale_*_continuous() - map cont values to visual values
```

```
scale_*_discrete() - map discrete values to visual values
```

```
scale_*_identity() - use data values as visual values
```

```
scale_*_manual() - map values to visual values
```

X and Y location scales

Use with x or y aesthetics (x shown here)

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

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

```
scale_x_log10() - Plot on log10 scale
```

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

```
scale_x_sqrt() - Plot on square root scale
```

Color and fill scales

Discrete

```
a < b + geom_bar(aes(fill = fl))
  n <= b + geom_bar(aes(fill = fl))
    n + scale_fill_manual(values = c("skyblue", "royalblue", "blue", "navy"),
      limits = c("D", "E", "P", "R"), breaks = c("D", "E", "P", "R"),
      name = "fuel", labels = c("D", "E", "P", "R"))
```

range of values to include in mapping
title to use in legend/axis
labels to use in legend/axis
breaks to use in legend/axis

Continuous

```
a < b + geom_dotplot(aes(fill = ..count..))
  n <= a + geom_dotplot(aes(fill = ..count..))
    n + scale_fill_continuous(low = "red", high = "yellow",
      midpoint = 0, trans = "log")
      n + scale_fill_grey(midpoint = 0.5)
        n + scale_fill_hexbin(colors = hexbin.colors(), topo.colors(),
          RColorBrewer::brewer.pal)
```

range of values to include in mapping
title to use in legend/axis
labels to use in legend/axis
breaks to use in legend/axis

Shape scales

Manual shape values

```
p < f + geom_point(aes(shape = fl))
  p + scale_shape_manual(values = 1:24)
    p + scale_shape_discrete()
```

range of values to include in mapping
title to use in legend/axis
labels to use in legend/axis
breaks to use in legend/axis

Size scales

Manual size values

```
q < f + geom_point(aes(size = cyl))
  q + scale_size_area(max = 6, aes(size = cyl))
    q + scale_size_area(max = 6, limit = 6)
```

range of values to include in mapping
title to use in legend/axis
labels to use in legend/axis
breaks to use in legend/axis

Coordinate Systems

```
r < b + geom_bar()
  r + coord_cartesian(xlim = c(0, 5))
    xlim, ylim
```

The default cartesian coordinate system

```
r + coord_fixed(ratio = 1/2)
  ratio, xlim, ylim
```

Cartesian coordinates with fixed aspect ratio between x and y units

```
r + coord_flip()
  xlim, ylim
```

Flipped Cartesian coordinates

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

Polar coordinates

```
r + coord_trans(ytrans = "sqrt")
  xtrans, ytrans, xlim, ylim
```

Transformed cartesian coordinates. Set extras and strains to the name of a window function.

```
z + coord_map(projection = "ortho",
  orientation = c(41, -74, 0))
  projection, orientation, xlim, ylim
```

Map projections from the mapproj package (mercator (default), azequalarea, lagrange, etc.)

Set scales to let axis limits vary across facets

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

- "free_x" - x axis limits adjust
- "free_y" - y axis limits adjust

Set labeller to adjust facet labels

```
t + facet_grid(~ fl, labeller = label_both)
  fl: c | fl: d | fl: e | fl: f | fl: r
```

```
t + facet_grid(~ fl, labeller = label_bquote(alpha ^ ..x..))
  alpha^c | alpha^d | alpha^e | alpha^f | alpha^r
```

```
t + facet_grid(~ fl, labeller = label_parsed)
  c | d | e | f | p | r
```

Position Adjustments

Position adjustments determine how to arrange geoms that would otherwise occupy the same space.

```
s < ggplot(mpg, aes(lt))
  s + geom_bar(position = "dodge")
```

Arrange elements side by side

```
s + geom_bar(position = "stack")
```

Stack elements on top of one another, normalize height

```
s + geom_bar(position = "stack")
```

Stack elements on top of one another

```
f + geom_point(position = "jitter")
```

Add random noise to X and Y position of each element to avoid overplotting

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

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

Themes

```
r + theme_bw()
  White background with grid lines
```

```
r + theme_classic()
  White background no gridlines
```

```
r + theme_grey()
  Grey background (default theme)
```

```
r + theme_minimal()
  Minimal theme
```

ggthemes - Package with additional ggplot2 themes

Faceting

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

```
t < ggplot(mpg, aes(cty, hwy)) + geom_point()
```

Set scales to let axis limits vary across facets

```
t + facet_grid(~ x, scales = "free")
```

- "free_x" - x axis limits adjust
- "free_y" - y axis limits adjust

All of the above

Labels

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

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

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

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

Use scale functions to update legend labels

Legends

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

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

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

Zooming

Without clipping (preferred)

```
t + coord_cartesian(xlim = c(0, 100), ylim = c(10, 20))
  White background with grid lines
```

With clipping (removes unseen data points)

```
t + xlim(0, 100) + ylim(10, 20)
  White background with grid lines
```

```
t + scale_x_continuous(limits = c(0, 100)) +
  scale_y_continuous(limits = c(0, 100))
  White background with grid lines
```

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

Python For Data Science Cheat Sheet

PySpark - SQL Basics

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



PySpark & Spark SQL

Spark SQL is Apache Spark's module for working with structured data.



Initializing SparkSession

A SparkSession can be used create DataFrame, register DataFrame as tables, execute SQL over tables, cache tables, and read parquet files.

```
>>> from pyspark.sql import SparkSession
>>> spark = SparkSession \
...     .builder \
...     .appName("Python Spark SQL basic example") \
...     .config("spark.some.config.option", "some-value") \
...     .getOrCreate()
```

Creating DataFrames

From RDDs

```
>>> from pyspark.sql.types import *
Infer Schema
>>> sc = spark.sparkContext
>>> lines = sc.textFile("people.txt")
>>> parts = lines.map(lambda l: l.split(","))
>>> people = parts.map(lambda p: Row(name=p[0], age=int(p[1])))
>>> peopleDF = spark.createDataFrame(people)
Specify Schema
>>> people = parts.map(lambda p: Row(name=p[0], age=pint(p[1]).strip())))
>>> schemaString = "name age"
>>> fields = [StructField(field_name, StringType(), True) for field_name in schemaString.split()]
>>> schema = StructType(fields)
>>> spark.createDataFrame(people, schema).show()
+-----+
| name | age |
+-----+
|  Phillip | 29 |
| Jonathan | 30 |
+-----+
```

From Spark Data Sources

```
JSON
>>> df = spark.read.json("customer.json")
>>> df.show()
+-----+
| address.firstName | lastName | phoneNumber |
+-----+
| New York, 10021, N.Y. | John | 333 888-1234, ho... |
| New York, 10021, N.Y. | John | 333 888-1234, ho... |
+-----+
>>> df2 = spark.read.load("people.json", format="json")
Parquet files
>>> df3 = spark.read.load("users.parquet")
TXT files
>>> df4 = spark.read.text("people.txt")
```

Inspect Data

>>> df.dtypes	Return df column names and data types
>>> df.show()	Display the content of df
>>> df.head(1)	Return first n rows
>>> df.take(2)	Return the first n rows
>>> df.schema	Return the schema of df

>>> df.describe().show()	Compute summary statistics
>>> df.columns	Return the columns of df
>>> df.count()	Count the number of rows in df
>>> df.distinct().count()	Count the number of distinct rows in df
>>> df.printSchema()	Print the schema of df
>>> df.explain()	Print the (logical and physical) plans

Python For Data Science Cheat Sheet

PySpark - RDD Basics

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



Spark

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



Initializing Spark

SparkContext

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

Inspect SparkContext

```
>>> sc.version
>>> sc.pythonVer
>>> sc.master
>>> sc.parallelize(["local"])
>>> sc.sparkUser()
>>> sc.appName
>>> sc.applicationId
>>> sc.defaultParallelism
>>> sc.defaultMinPartitions
```

Configuration

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

Using The Shell

In the PySpark shell, a special interpreter-aware SparkContext is already created in the variable called `sc`.

```
$ ./bin/spark-shell --master local[2]
$ ./bin/pyspark --master local[4] --py-files code.py
```

Set which master the context connects to with the `--master` argument, and add Python zip, egg or py files to the runtime path by passing a comma-separated list to `--py-files`.

Loading Data

Parallelized Collections

```
>>> rdd = sc.parallelize([(1,2),(1,2),(1,2)])
>>> rdd2 = sc.parallelize([(1,2),(1,1),(1,1)])
>>> rdd3 = sc.parallelize(range(100))
>>> rdd4 = sc.parallelize([(1,"x","y","z"), ("b","p","r")])
```

External Data

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

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

Duplicate Values

```
>>> df = df.dropDuplicates()
```

Queries

```
>>> from pyspark.sql import functions as F
>>> df.select("firstName").show()
+-----+
| firstName |
+-----+
| Phillip |
| Jonathan |
+-----+
>>> df.select("firstName", "age")
+-----+
| firstName | age |
+-----+
| Phillip | 29 |
| Jonathan | 30 |
+-----+
>>> df.withColumn("age_sq", F.col("age") * F.col("age")).show()
+-----+
| firstName | age | age_sq |
+-----+
| Phillip | 29 | 841 |
| Jonathan | 30 | 900 |
+-----+
>>> df.select("firstName", "age").where(F.col("age") > 24).show()
+-----+
| firstName | age |
+-----+
| Jonathan | 30 |
+-----+
>>> df.select("firstName", "age").where(F.when(df.age > 30, 1).otherwise(0)).show()
+-----+
| firstName | age |
+-----+
| Phillip | 29 |
+-----+
>>> df.filter(df.firstName.isin("Jane", "Boris")).collect()
[Row(firstName='Jane', age=30), Row(firstName='Boris', age=30)]
>>> df.filter(df.firstName.isin("Jane", "Boris")).collect()
[Row(firstName='Jane', age=30), Row(firstName='Boris', age=30)]
>>> df.filter(df.lastName.like("Smith")).collect()
[Row(firstName='Phillip', lastName='Smith'), Row(firstName='Jonathan', lastName='Smith')]
>>> df.filter(df.lastName.startswith("Sm")).collect()
[Row(firstName='Phillip', lastName='Smith'), Row(firstName='Jonathan', lastName='Smith')]
>>> df.filter(df.lastName.endswith("th")).collect()
[Row(firstName='Phillip', lastName='Smith'), Row(firstName='Jonathan', lastName='Smith')]
>>> df.select(df.firstName.substr(1, 3).alias("name")).collect()
[Row(name='Phi'), Row(name='Jon')]
>>> df.filter(df.age.between(22, 24)).collect()
[Row(firstName='Phillip', age=29), Row(firstName='Jonathan', age=30)]
```

Add, Update & Remove Columns

Adding Columns

```
>>> df = df.withColumn("city", df.address.city) \
...     .withColumn("postalCode", df.address.postalCode) \
...     .withColumn("state", df.address.state) \
...     .withColumn("streetAddress", df.address.streetAddress) \
...     .withColumn("phoneNumbers", explode(df.phoneNumber.number)) \
...     .withColumn("phoneType", explode(df.phoneNumber.type))
```

Updating Columns

```
>>> df = df.withColumnRenamed('telephoneNumber', 'phoneNumber')
```

Removing Columns

```
>>> df = df.drop("address", "phoneNumber")
>>> df = df.drop(df.address).drop(df.phoneNumber)
```

GroupBy

```
>>> df.groupby("age") \
...     .count() \
...     .show()
```

Filter

```
>>> df.filter(df["age"]>24).show()
```

SORT

```
>>> peopleDF.sort(peopleDF.age.desc()).collect()
>>> df.sort("age", ascending=False).collect()
>>> df.orderBy(["age", "city"], ascending=[0,1]).collect()
```

Missing & Replacing Values

```
>>> df.na.fill(50).show()
>>> df.na.drop().show()
>>> df.na.replace(10, 20).show()
```

Repartitioning

```
>>> df.repartition(10) \
...     .getNumPartitions()
df with 10 partitions
>>> df.coalesce(1).getNumPartitions()
df with 1 partition
```

Running SQL Queries Programmatically

Registering DataFrames as Views

```
>>> peopleDF.createGlobalTempView("people")
>>> df.createTempView("customers")
>>> df.createOrReplaceTempView("customer")
```

Query Views

```
>>> df5 = spark.sql("SELECT * FROM customer").show()
>>> peopleDF2 = spark.sql("SELECT * FROM global_temp.people") \
...     .show()
```

Output

Data Structures

```
>>> rdd1 = df.rdd
Convert df into an RDD
>>> df.toJSON().first()
Convert df into a RDD of string
>>> df.toPandas()
Return the contents of df as Pandas DataFrame
```

Write & Save to Files

```
>>> df.select("firstName", "city") \
...     .write \
...     .save("nameAndCity.parquet")
>>> df.select("firstName", "age") \
...     .write \
...     .save("namesAndAges.json", format="json")
```

Stopping SparkSession

```
>>> spark.stop()
```

DataCamp
Learn Python for Data Science interactively



Python For Data Science Cheat Sheet

PySpark - RDD Basics

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



Spark

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



Initializing Spark

SparkContext

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

Inspect SparkContext

```
>>> sc.version
>>> sc.pythonVer
>>> sc.master
>>> sc.parallelize(["local"])
>>> sc.sparkUser()
>>> sc.appName
>>> sc.applicationId
>>> sc.defaultParallelism
>>> sc.defaultMinPartitions
```

Configuration

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

Using The Shell

In the PySpark shell, a special interpreter-aware SparkContext is already created in the variable called `sc`.

```
$ ./bin/spark-shell --master local[2]
$ ./bin/pyspark --master local[4] --py-files code.py
```

Set which master the context connects to with the `--master` argument, and add Python zip, egg or py files to the runtime path by passing a comma-separated list to `--py-files`.

Retrieving RDD Information

Basic Information

```
>>> rdd.getNumPartitions()
3
>>> rdd.count()
3
>>> rdd.countByKey()
defaultdict(<type 'int'>, {'a':2, 'b':1})
>>> rdd.countByValue()
{'a':2}
>>> rdd.collectAsMap()
{'a': 2, 'b': 1}
>>> rdd.sum()
495
>>> sc.parallelize([]).isEmpty()
True
```

Summary

>>> rdd3.max() 99	Maximum value of RDD elements
>>> rdd3.min() 0	Minimum value of RDD elements
>>> rdd3.mean() 49.5	Mean value of RDD elements
>>> rdd3.stdev() 28.86607004772218	Standard deviation of RDD elements
>>> rdd3.variance() 833.0	Compute variance of RDD elements
>>> rdd3.histogram(3) ([0,33,66,99],[33,33,34]) >>> rdd3.stats()	Compute histogram by bins
	Summary statistics (count, mean, stdev, max & min)

Applying Functions

>>> rdd.map(lambda x: x+(x[1],x[0])) .collect() [('a',7,7), ('a',2,2), ('b',2,2), ('b',1,1)] >>> rdd2 = rdd.flatMap(lambda x: x+(x[1],x[0])) .collect() [('a',7,7), ('a',2,2), ('b',2,2), ('b',1,1)] >>> rdd4 = rdd.flatMapValues(lambda x: x) .collect() [('a','x'), ('a','y'), ('a','z'), ('b','p'), ('b','r')]	Apply a function to each RDD element
>>> rdd3.map(lambda x: x+(x[1],x[0])) .collect() [('a',7,7), ('a',2,2), ('b',2,2), ('b',1,1)] >>> rdd3.flatMap(lambda x: x+(x[1],x[0])) .collect() [('a',7,7), ('a',2,2), ('b',2,2), ('b',1,1)]	Apply a function to each RDD element and flatten the result
>>> rdd3.map(lambda x: (x[0]+x[1],x[0])) .collect() [('a',9,9), ('b',2,2)] >>> rdd3.aggregateByKey(0,0,seqOp,combOp) (4950,100) >>> rdd3.aggregateByKey(0,0,seqOp,combOp) (4950,100) >>> rdd3.foldByKey(0, add) 4950 >>> rdd3.foldByKey(0, add) 4950 [('a',9), ('b',2)] >>> rdd3.foldByKey(0, add) 4950 [('a',9), ('b',2)]	Aggregate RDD elements of each partition and then the results
>>> rdd3.aggregateByKey(0,0,seqOp,combOp) (4950,100) >>> rdd3.foldByKey(0, add) 4950 [('a',9), ('b',2)] >>> rdd3.foldByKey(0, add) 4950 [('a',9), ('b',2)]	Aggregate values of each RDD key
>>> rdd3.map(lambda x: x+x) .collect() [('a',2), ('b',1), ('c',0)]	Create tuples of RDD elements by applying a function

Reshaping Data

Reducing

```
>>> rdd.reduceByKey(lambda x,y : x+y)
.collect()
[('a',9), ('b',2)]
>>> rdd.reduce(lambda a, b: a + b)
('a',7,7,2,2)
```

Merge the RDD values for each key

Merge the RDD values

Grouping by

```
>>> rdd3.groupByKey(x: x % 2)
.collect()
[('a',7), ('a',2,2,2,2)]
>>> rdd3.groupByKey(x: x % 2)
.collect()
[('a',7,7,2,2,2,2)]
```

Return RDD of grouped values

Group RDD by key

Aggregating

```
>>> rdd.map(lambda x,y: (x[0]+y,x[1]+1))
.collect()
[('a',9,2), ('b',2,2)]
>>> rdd3.aggregateByKey(0,0,seqOp,combOp)
(4950,100)
>>> rdd3.aggregateByKey(0,0,seqOp,combOp)
(4950,100)
>>> rdd3.foldByKey(0, add)
4950
>>> rdd3.foldByKey(0, add)
4950
[('a',9), ('b',2)]
>>> rdd3.foldByKey(0, add)
4950
[('a',9), ('b',2)]

```

Aggregate RDD elements of each partition and then the results

Aggregate values of each RDD key

Mathematical Operations

```
>>> rdd.subtract(rdd2)
.collect()
[('b',2), ('a',7)]
>>> rdd2.subtractByKey(rdd)
.collect()
[('d',1)]
>>> rdd.cartesian(rdd2).collect()
[('a',2), ('b',1), ('c',1), ('d',1)]
```

Return each RDD value not contained in rdd2

Return each (key,value) pair of rdd2 with no matching key in rdd

Return the Cartesian product of rdd and rdd2

SORT

```
>>> rdd2.sortBy(lambda x: x[1])
.collect()
[('a',2), ('b',1)]
>>> rdd2.orderBy(lambda x: x[1])
.collect()
[('a',2), ('b',1), ('c',1)]
```

Sort RDD by given function

Sort (key, value) RDD by key

Repartitioning

```
>>> rdd2.repartition(4)
.collect()
[('a',2), ('b',1), ('c',1), ('d',1)]
```

New RDD with 4 partitions

Decrease the number of partitions in the RDD to 1

Saving

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

Stopping SparkContext

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

Execution

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

DataCamp
Learn Python for Data Science interactively

