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

Pandas API Reference   Pandas User Guide

## Tidy Data – A foundation for wrangling in pandas

In a tidy data set:



&



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

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

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



## Creating DataFrames



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

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



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

## Method Chaining

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

## Reshaping Data – Change layout, sorting, reindexing, renaming



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



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



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



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

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

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

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

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

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

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

## Subset Observations - rows



```
df[df.Length > 7]
```
Extract rows that meet logical criteria.
```
df.drop_duplicates()
```
Remove duplicate rows (only considers columns).
```
df.sample(frac=0.5)
```
Randomly select fraction of rows.
```
df.sample(n=10)   Randomly select n rows.
df.nlargest(n, 'value')
```
Select and order top n entries.
```
df.nsmallest(n, 'value')
```
Select and order bottom n entries.
```
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*.

## Using query

query() allows Boolean expressions for filtering rows.
```
df.query('Length > 7')
df.query('Length > 7 and Width < 8')
df.query('Name.str.startswith("abc")',
         engine="python")
```

## Subsets - rows and columns

Use **df.loc[]** and **df.iloc[]** to select only rows, only columns or both.
Use **df.at[]** and **df.iat[]** to access a single value by row and column.
First index selects rows, second index columns.

```
df.iloc[10:20]
```
Select rows 10-20.
```
df.iloc[:, [1, 2, 5]]
```
Select columns in positions 1, 2 and 5 (first column is 0).
```
df.loc[:, 'x2':'x4']
```
Select all columns between x2 and x4 (inclusive).
```
df.loc[df['a'] > 10, ['a', 'c']]
```
Select rows meeting logical condition, and only the specific columns .
```
df.iat[1, 2]   Access single value by index
df.at[4, 'A']   Access single value by label
```

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

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

Cheatsheet for pandas (http://pandas.pydata.org/ originally written by Irv Lustig, Princeton Consultants, inspired by Rstudio Data Wrangling Cheatsheet

# Summarize Data

`df['w'].value_counts()`
    Count number of rows with each unique value of variable

`len(df)`
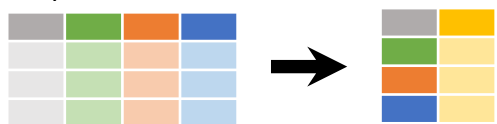    # of rows in DataFrame.

`df.shape`
    Tuple of # of rows, # of columns in DataFrame.

`df['w'].nunique()`
    # of distinct values in a column.

`df.describe()`
    Basic descriptive and 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 pandas Series for each column. Examples:

`sum()`
    Sum values of each object.

`count()`
    Count non-NA/null values of each object.

`median()`
    Median value of each object.

`quantile([0.25,0.75])`
    Quantiles of each object.

`apply(function)`
    Apply function to each object.

`min()`
    Minimum value in each object.

`max()`
    Maximum value in each object.

`mean()`
    Mean value of each object.

`var()`
    Variance of each object.

`std()`
    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.

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)`
    Copy with values shifted by 1.

`rank(method='dense')`
    Ranks with no gaps.

`rank(method='min')`
    Ranks. Ties get min rank.

`rank(pct=True)`
    Ranks rescaled to interval [0, 1].

`rank(method='first')`
    Ranks. Ties go to first value.

`shift(-1)`
    Copy with values lagged by 1.

`cumsum()`
    Cumulative sum.

`cummax()`
    Cumulative max.

`cummin()`
    Cumulative min.

`cumprod()`
    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.
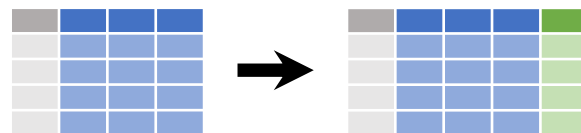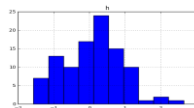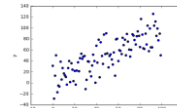
# Handling Missing Data

`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)`
    Element-wise max.

`clip(lower=-10,upper=10)`
    Trim values at input thresholds

`min(axis=1)`
    Element-wise min.

`abs()`
    Absolute value.

# Plotting

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

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

# Combine Data Sets

**adf**

| x1 | x2 |
|----|----|
| A | 1 |
| B | 2 |
| C | 3 |

**+**

**bdf**

| x1 | x3 |
|----|----|
| A | T |
| B | F |
| D | T |

**=**

### Standard Joins

| x1 | x2 | x3 |
|----|----|----|
| A | 1 | T |
| B | 2 | F |
| C | 3 | NaN |

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

| x1 | x2 | x3 |
|----|----|----|
| A | 1.0 | T |
| B | 2.0 | F |
| D | NaN | T |

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

| x1 | x2 | x3 |
|----|----|----|
| A | 1 | T |
| B | 2 | F |

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

| x1 | x2 | x3 |
|----|----|----|
| A | 1 | T |
| B | 2 | F |
| C | 3 | NaN |
| D | NaN | T |

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

### Filtering Joins

| x1 | x2 |
|----|----|
| A | 1 |
| B | 2 |

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

| x1 | x2 |
|----|----|
| C | 3 |

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

**ydf**

| x1 | x2 |
|----|----|
| A | 1 |
| B | 2 |
| C | 3 |

**+**

**zdf**

| x1 | x2 |
|----|----|
| B | 2 |
| C | 3 |
| D | 4 |

**=**

### Set-like Operations

| x1 | x2 |
|----|----|
| B | 2 |
| C | 3 |

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

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

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

| x1 | x2 |
|----|----|
| A | 1 |

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

# Data Science Cheat Sheet
## NumPy

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

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

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

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

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

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

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

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

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

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

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

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

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

# 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

**1D array**

| 1 | 2 | 3 |
|---|---|---|

**2D array**

axis 1
axis 0

| 1.5 | 2 | 3 |
|-----|---|---|
| 4   | 5 | 6 |

**3D array**

axis 2
axis 1
axis 0

## Creating Arrays

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

### Initial Placeholders

```
>>> np.zeros((3,4))                 Create an array of zeros
>>> np.ones((2,3,4),dtype=np.int16) Create an array of ones
>>> d = np.arange(10,25,5)          Create an array of evenly
                                    spaced values (step value)
>>> np.linspace(0,2,9)              Create an array of evenly
                                    spaced values (number of samples)
>>> e = np.full((2,2),7)            Create a constant array
>>> f = np.eye(2)                   Create a 2X2 identity matrix
>>> np.random.random((2,2))         Create an array with random values
>>> np.empty((3,2))                 Create an empty array
```

## I/O

### Saving & Loading On Disk

```
>>> np.save('my_array', a)
>>> np.savez('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

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

## Inspecting Your Array

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

## Asking For Help

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

## Array Mathematics

### Arithmetic Operations

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

### Comparison

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

### Aggregate Functions

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

## Copying Arrays

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

## Sorting Arrays

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

## Subsetting, Slicing, Indexing

### Subsetting

```
>>> a[2]       Select the element at the 2nd index
  3
>>> b[1,2]     Select the element at row 1 column 2
  6.0          (equivalent to b[1][2])
```

### Slicing

```
>>> a[0:2]                Select items at index 0 and 1
  array([1, 2])
>>> b[0:2,1]              Select items at rows 0 and 1 in column 1
  array([ 2.,  5.])
>>> b[:1]                 Select all items at row 0
  array([[1.5, 2., 3.]])  (equivalent to b[0:1, :])
>>> c[1,...]              Same as [1,:,:]
  array([[[ 3.,  2.,  1.],
          [ 4.,  5.,  6.]]])
>>> a[ : :-1]             Reversed array a
  array([3, 2, 1])
```

### Boolean Indexing

```
>>> a[a<2]     Select elements from a less than 2
  array([1])
```

### Fancy Indexing

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

## Array Manipulation

### Transposing Array

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

### Changing Array Shape

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

### Adding/Removing Elements

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

### Combining Arrays

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

### Splitting Arrays

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

# Python For Data Science *Cheat Sheet*
## Matplotlib

Learn Python **Interactively** at www.DataCamp.com

## Matplotlib

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

## Plot Anatomy & Workflow

### Plot Anatomy



### Workflow

The basic steps to creating plots with matplotlib are:

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

```
>>> import matplotlib.pyplot as plt
>>> x = [1,2,3,4]              Step 1
>>> y = [10,20,25,30]
>>> fig = plt.figure()        Step 2
>>> ax = fig.add_subplot(111)  Step 3
>>> ax.plot(x, y, color='lightblue', linewidth=3)   Step 3, 4
>>> ax.scatter([2,4,6],
               [5,15,25],
               color='darkgreen',
               marker='^')
>>> ax.set_xlim(1, 6.5)
>>> plt.savefig('foo.png')
>>> plt.show()                Step 6
```

## 1 Prepare The Data

**Also see Lists & NumPy**

### 1D Data

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

### 2D Data or Images

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

## 2 Create Plot

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

### Figure

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

### Axes

All plotting is done with respect to an `Axes`. In most cases, a subplot will fit your needs. A subplot is an axes on a grid system.

```
>>> fig.add_axes()
>>> ax1 = fig.add_subplot(221) # row-col-num
>>> ax3 = fig.add_subplot(212)
>>> fig3, axes = plt.subplots(nrows=2,ncols=2)
>>> fig4, axes2 = plt.subplots(ncols=3)
```

## 4 Customize Plot

### Colors, Color Bars & Color Maps

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

### Markers

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

### Linestyles

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

### Text & Annotations

```
>>> ax.text(1,
            -2.1,
            'Example Graph',
            style='italic')
>>> ax.annotate("Sine",
                xy=(8, 0),
                xycoords='data',
                xytext=(10.5, 0),
                textcoords='data',
                arrowprops=dict(arrowstyle="->",
                                connectionstyle="arc3"),)
```

### Mathtext

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

### Limits, Legends & Layouts

**Limits & Autoscaling**
```
>>> ax.margins(x=0.0,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
```

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

**Ticks**
```
>>> ax.xaxis.set(ticks=range(1,5),    Manually set x-ticks
                 ticklabels=[3,100,-12,"foo"])
>>> ax.tick_params(axis='y',          Make y-ticks longer and go in and out
                   direction='inout',
                   length=10)
```

**Subplot Spacing**
```
>>> fig3.subplots_adjust(wspace=0.5,   Adjust the spacing between subplots
                         hspace=0.3,
                         left=0.125,
                         right=0.9,
                         top=0.9,
                         bottom=0.1)
>>> fig.tight_layout()                 Fit subplot(s) in to the figure area
```

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

## 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].bar([1,2,3],[3,4,5])  Plot vertical rectangles (constant width)
>>> axes[1,0].barh([0.5,1,2.5],[0,1,2])  Plot horiontal rectangles (constant height)
>>> axes[1,1].axhline(0.45)         Draw a horizontal line across axes
>>> axes[0,1].axvline(0.65)         Draw a vertical line across axes
>>> ax.fill(x,y,color='blue')       Draw filled polygons
>>> ax.fill_between(x,y,color='yellow')  Fill between y-values and 0
```

### 2D Data or Images

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

### Vector Fields

```
>>> axes[0,1].arrow(0,0,0.5,0.5)     Add an arrow to the axes
>>> axes[1,1].quiver(y,z)            Plot a 2D field of arrows
>>> axes[0,1].streamplot(X,Y,U,V)    Plot 2D vector fields
```

### Data Distributions

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

```
>>> axes2[0].pcolor(data2)      Pseudocolor plot of 2D array
>>> axes2[0].pcolormesh(data)   Pseudocolor plot of 2D array
>>> CS = plt.contour(Y,X,U)     Plot contours
>>> axes2[2].contourf(data1)    Plot filled contours
>>> axes2[2]= ax.clabel(CS)     Label a contour plot
```

## 5 Save Plot

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

## 6 Show Plot

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

## Close & Clear

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

# Python For Data Science *Cheat Sheet*
## Scikit-Learn

Learn Python for data science **Interactively** at www.DataCamp.com

## Scikit-learn

**Scikit-learn** is an open source Python library that implements a range of machine learning, preprocessing, cross-validation and visualization algorithms using a unified interface.

### A Basic Example

```
>>> from sklearn import neighbors, datasets, preprocessing
>>> from sklearn.model_selection import train_test_split
>>> from sklearn.metrics import accuracy_score
>>> iris = datasets.load_iris()
>>> X, y = iris.data[:, :2], iris.target
>>> X_train, X_test, y_train, y_test = train_test_split(X,y,random_state=33)
>>> scaler = preprocessing.StandardScaler().fit(X_train)
>>> X_train = scaler.transform(X_train)
>>> X_test = scaler.transform(X_test)
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
>>> knn.fit(X_train, y_train)
>>> y_pred = knn.predict(X_test)
>>> accuracy_score(y_test, y_pred)
```

## Loading The Data                    Also see **NumPy & Pandas**

Your data needs to be numeric and stored as NumPy arrays or SciPy sparse matrices. Other types that are convertible to numeric arrays, such as Pandas DataFrame, are also acceptable.

```
>>> import numpy as np
>>> X = np.random.random((10,5))
>>> y = np.array(['M','M','F','F','M','F','M','M','F','F','F'])
>>> X[X < 0.7] = 0
```

## Training And Test Data

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

## Preprocessing The Data

### Standardization

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

### Normalization

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

### Binarization

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

### Encoding Categorical Features

```
>>> from sklearn.preprocessing import LabelEncoder
>>> enc = LabelEncoder()
>>> y = enc.fit_transform(y)
```

### Imputing Missing Values

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

### Generating Polynomial Features

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

## Create Your Model

### Supervised Learning Estimators

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

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

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

#### KNN
```
>>> from sklearn 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)
```

## Model Fitting

### Supervised learning
```
>>> lr.fit(X, y)
>>> knn.fit(X_train, y_train)
>>> svc.fit(X_train, y_train)
```
Fit the model to the data

### Unsupervised Learning
```
>>> k_means.fit(X_train)
>>> pca_model = pca.fit_transform(X_train)
```
Fit the model to the data
Fit to data, then transform it

## Prediction

### Supervised Estimators
```
>>> y_pred = svc.predict(np.random.random((2,5)))
>>> y_pred = lr.predict(X_test)
>>> y_pred = knn.predict_proba(X_test)
```
Predict labels
Predict labels
Estimate probability of a label

### Unsupervised Estimators
```
>>> y_pred = k_means.predict(X_test)
```
Predict labels in clustering algos

## Evaluate Your Model's Performance

### Classification Metrics

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

## Tune Your Model

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

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

scikit-learn algorithm cheat-sheet