



HACKWAGON

• ACADEMY •

DS102

Data Analysis with Python

The Cheat Sheet

All the imports - Data Manipulation

pandas	<code>import pandas as pd</code>
numpy	<code>import numpy as np</code>
matplotlib	<code>import matplotlib.pyplot as plt</code>
seaborn	<code>import seaborn as sns</code>
sklearn	<code>import sklearn</code>
statsmodels	<code>import statsmodels</code>
scipy	<code>import scipy</code>



All the imports - Text Mining, Web Scraping

wordcloud	<code>import wordcloud</code>
requests	<code>import requests</code>
bs4	<code>from bs4 import BeautifulSoup</code>
nltk	<code>import nltk</code> <code>nltk.download('corpus')</code>



Instantiating a DataFrame

Instantiate a df from a CSV file	<pre>df = pd.read_csv('x.csv')</pre> <p>Variations:</p> <pre>df = pd.read_csv('arrivals.csv', sep=' ', index_col=0) df = pd.read_csv('arrivals.csv', index_col=0)</pre>									
Instantiate a df from a list of dict	<pre>students = ['name': 'Ben', 'age': 22, 'gender': 'M'], {'name': 'John', 'age': 32, 'gender': 'M'}, {'name': 'Charlotte', 'age': 21, 'gender': 'F'}]</pre> <pre>df = pd.DataFrame(students)</pre>									
Instantiate a df from a dict	<pre>student = {'name': ['Ben','John'], 'age': [22,32], 'gender': ['M','M']}</pre> <pre>df = pd.DataFrame.from_dict(student)</pre> <table><tr><th>name</th><th>age</th><th>gender</th></tr><tr><td>Ben</td><td>22</td><td>M</td></tr><tr><td>John</td><td>32</td><td>M</td></tr></table>	name	age	gender	Ben	22	M	John	32	M
name	age	gender								
Ben	22	M								
John	32	M								



Properties of a DataFrame

Gets the number of columns and rows in the DataFrame and returns it in a tuple	<code>df.shape</code>
Gets all the data types of each column	<code>df.dtypes</code>
Gets all the column labels of the DataFrame (header names)	<code>df.columns</code>
Gets the column count and the data type for each column	<code>df.info()</code>
Gets statistical data of the DataFrame such as count, mean, standard deviation, min value, max value	<code>df.describe()</code> / <code>s.describe()</code>



Peeking into a DataFrame

Obtains the first 5 (default) data in the DataFrame unless specified. Eg. <code>df.head(2)</code> obtains the first 2 data	<code>df.head()</code> / <code>s.head()</code> Variations: <code>df.head(2)</code> , <code>s.head(2)</code>
Obtains 1 (default) random data from the DataFrame without replacement unless specified Eg. <code>df.sample(2, replace=True)</code> obtains 2 random data with replacement	<code>df.sample()</code> / <code>s.sample()</code> Variations: <code>df.sample(3)</code> , <code>s.sample(3)</code> , <code>df.sample(4, replace=True)</code>
Gets an array / list of unique values	<code>df.unique()</code> / <code>s.unique()</code>
Gets the number of unique values	<code>df.nunique()</code> / <code>s.nunique()</code>
Gets the count of non-NA values for each column	<code>df.count()</code> / <code>s.count()</code>



Key

df: pandas DataFrame

s : pandas Series

nsc: NOT syntactically correct

Peeking into a DataFrame

Get the minimum value for each column of the df	<code>df.min()</code> / <code>s.min()</code>
Get the maximum value for each column of the df	<code>df.max()</code> / <code>s.max()</code>
Get the summation of all the values for each column	<code>df.sum()</code> / <code>s.sum()</code>
Get the cumulative sum of all the values for each column	<code>df.cumsum()</code> / <code>s.cumsum()</code>



Key

df: pandas DataFrame

s : pandas Series

nsc: NOT syntactically correct

Sorting

Sorts the Dataframe by the column name specified	<code>df.sort_values('col_x')</code>
Sorts the Dataframe by the column name specified	<code>df.sort_values(['col_x'])</code>
Sorts the Dataframe by the column name specified in the ascending order	<code>df.sort_values('col_x', ascending=True)</code> Variations: <code>df.sort_values('col_x', ascending=False)</code>
Sorts the dataframe by column: 'col_x' followed by column: 'col_y'	<code>df.sort_values(['col_x', 'col_y'])</code> Variations: <code>df.sort_values(['col_x', 'col_y', ascending=[True, False])</code>



Key

df: pandas DataFrame

s : pandas Series

nsc: NOT syntactically correct

Retrieval from a DataFrame

Locates a selection based on the position index	<code>df.iloc[0]</code> Variations: <code>df.iloc[[1, 2, 3]]</code> , <code>df.iloc[1:10]</code>
Access a group of row and column	<code>df.loc[0]</code> Variations: <code>df.loc[[1, 2, 3]]</code> , <code>df.loc[range(0, 10)]</code>
Retrieves all values for the specified column in the dataframe	<code>df['col_x']</code>
Retrieves all values for the specified columns in the dataframe	<code>df[['col_x', 'col_y']]</code>
The value for the specified column in the series	<code>s['col_x']</code>



Simple Statistics & Moments of a Distribution

Returns the mean of each numeric column in the dataframe / series	<code>df.mean()</code> / <code>s.mean()</code>
Returns the standard deviation of each numeric column in the dataframe / series	<code>df.std()</code> / <code>s.std()</code> Variations: <code>df.sample(t)</code> , <code>s.sample(t)</code> , <code>df.sample(t, replace=True)</code>
Returns the variance of each numeric column in the dataframe / series	<code>df.var()</code> / <code>s.var()</code>
Returns the median of each numeric column in the dataframe / series	<code>df.median()</code> / <code>s.median()</code>



Simple Statistics & Moments of a Distribution

Returns the mode of each numeric column in the dataframe / series	<code>df.mode()</code> / <code>s.mode()</code>
Returns a unbiased skew of the dataframe / series	<code>df.skew()</code> / <code>s.skew()</code>
Returns a unbiased kurtosis of the dataframe / series	<code>df.kurt()</code> / <code>s.kurt()</code>



Filtering

Filter for values greater than <i>i</i>	<pre>df[df['col_1'] > i]</pre> <p>Variations: <code>df[df['col_1'] >= i],</code> <code>df[df['col_1'] < i],</code> <code>df[df['col_1'] <= i]</code></p>
Filter for values equal to <i>i</i> (Applies to strings as well)	<pre>df[df['col_1'] == i]</pre> <p>Variations: <code>df[df['col_1'] == 'string_i']</code></p>
Filter for having a list of values (Applies to strings as well)	<pre>df[df['col_1'].isin([i, j, k])]</pre> <p>Variations: <code>df[df['col_1'].isin(['string_i', 'string_j', 'string_k'])]</code></p>



Filtering (AND, OR, NOT)

Using NOT on a condition	<code>df[~(df['col_1'] > i)]</code>
Using AND on two conditions (record must fulfill condition_1 AND condition_2)	<code>df[(condition_1) & (condition_2)]</code> <code>df[(df['col_1'] > lo) & (df['col_1'] < hi)]</code> Variations: <code>df[(condition_1) & (condition_2) & ...] (nsc)</code>
Using OR on two conditions (record can fulfill condition_1 OR condition_2)	<code>df[(condition_1) (condition_2)]</code> <code>df[(df['col_1'] < lo) (df['col_1'] > hi)]</code> Variations: <code>df[(condition_1) (condition_2) ...] (nsc)</code>



Handling Empty Values

Replaces all NA / NaN values with the specified value	<code>df.fillna()</code>
Drops rows with any NA / NaN values	<code>df.dropna()</code> Variations: <code>df[(condition_1) & (condition_2) & ...] (nsc)</code>
Returns a same sized object that shows if the value is a NA / NaN. If value is null, it will be shown as True else, False	<code>df.isnull()</code> Variations: <code>df[(condition_1) (condition_2) ...] (nsc)</code>



Substitutions

Renames the column's name or index name	<code>df.rename()</code> Variations: <code>df.rename(columns={})</code> , <code>df.rename({}, axis='index')</code>
Replaces the specified values with another value and returns the Dataframe.	<code>df.replace()</code> Variations: <code>df[(condition_1) & (condition_2) & ...]</code> (nsc)



Removing Items

Drops the row by index	<code>df.drop(axis=0)</code> Variations: <code>df.drop('column name', axis=0)</code>
Drops the column of the specified value	<code>df.drop(axis=1)</code> Variations: <code>df.drop('column name', axis=1, inplace=True)</code>



Creating New Columns

Using simple arithmetic	<code>df['x'] = df['a'] + df['b'] + df['c'] + ...</code>
Using apply and a function (1 column)	<pre>def minus_5(x): return x-5 df['col_name'].apply(minus_5)</pre>
Using apply and a function (2 columns)	<pre>cols_name = ['col1', 'col2'] for col in cols_name: df[col].apply(minus_5)</pre>
Using apply and a function (named params)	<pre>def add(x, y=4): return x+y df['col_name'].apply(add, y=5)</pre>



Creating New Columns

Using apply and a function (lambda)	<pre>D = {True: 'Yes', False: 'No'} df['yn'] = df['Legendary'].apply(lambda x : return D[x])</pre>
Using apply and a function (numpy functions)	<pre>df['yn'] = df['Attack'].apply(np.int)</pre>



Aggregation

groupby() on one column	<code>df.groupby('column1')</code>
groupby() on two or more columns	<code>df.groupby(['column1','column2'])</code>



Aggregation Functions

Gives the maximum value of the first numerical valued column in the group	<code>groupby().max()</code>
Gives the minimum value of the first numerical valued column in the group	<code>groupby().min()</code>
Shows the number of occurrence that falls in the group	<code>groupby().size()</code>
Aggregates using one or more functions over the groups	<code>groupby().agg()</code> Variations: <code>df.groupby().agg(['min','max'])</code>
Computes the mean of the groups	<code>groupby().mean()</code>
Transforms means combining based on the idmax	<code>groupby().transform(idmax)</code>



Pivot

df

	date	game	no_players	score
0	14/12/2018	basketball	14	87
1	15/12/2018	basketball	13	76
2	16/12/2018	basketball	15	54
3	14/12/2018	soccer	14	2
4	15/12/2018	soccer	15	4
5	16/12/2018	soccer	14	1
6	14/12/2018	tennis	4	24
7	15/12/2018	tennis	4	23
8	16/12/2018	tennis	2	17

Key

df: pandas DataFrame

s : pandas Series

nsc: NOT syntactically correct ²¹



Pivot

```
df.pivot(index = 'game',  
columns = 'date')
```

date	no_players			score		
	14/12/2018	15/12/2018	16/12/2018	14/12/2018	15/12/2018	16/12/2018
game						
basketball	14	13	15	87	76	54
soccer	14	15	14	2	4	1
tennis	4	4	2	24	23	17



Pivot

```
df.pivot(index='game',columns='date',value  
s='no_players')
```

date	14/12/2018	15/12/2018	16/12/2018
game			
basketball	14	13	15
soccer	14	15	14
tennis	4	4	2



Pivot Tables

Returns a pivot table with the resultant index and columns which can aggregate the data based on what you specify in the index.

Values are used to specify which data we want to grab from the pivot table which return them in a column. If there are multiple indexes with the same key, it will take the average. However, if you were to specify aggfunc to be np.sum, it will return the sum of those with the same index.

Aggfunc can take a list of functions to apply to the index. E.g. len, np.mean()

```
pd.pivot_table(df,  
index=[],  
columns=[],  
values=[],  
aggfunc=len)
```



Pivot Tables

```
df.pivot_table(index='game',columns='no_players')
```

	score				
	2	4	13	14	15
game					
basketball	NaN	NaN	76.0	87.0	54.0
soccer	NaN	NaN	NaN	1.5	4.0
tennis	17.0	23.5	NaN	NaN	NaN



Pivot Table (Margins)

```
import numpy as np
```

```
df.pivot_table(index='game', columns='date',  
               margins=True, aggfunc=np.sum)
```

date	no_players				score			
	14/12/2018	15/12/2018	16/12/2018	All	14/12/2018	15/12/2018	16/12/2018	All
game								
basketball	14	13	15	42	87	76	54	217
soccer	14	15	14	43	2	4	1	7
tennis	4	4	2	10	24	23	17	64
All	32	32	31	95	113	103	72	288



Key

df: pandas DataFrame

s : pandas Series

nsc: NOT syntactically correct

Merge & Join

Returns a merged Dataframe of a and b with the name of 'c' as the header	<code>pd.merge(a, b, on = 'c')</code>
Returns a joined dataframe with the left dataframe as a and the right as b. The rows which are not filled will be replaced by 'NaN'	<code>a.join(lsuffix = 'a', rsuffix = 'b')</code>
Returns a transposed a with the rows as columns and vice versa	<code>a.transpose()</code> <code>a.T()</code>



More EDA (Quartiles, IQR, Correlation)

Returns the 25th percentile of s	<code>np.percentile(s, 25)</code>
Returns the 25th, 50th, and 75th percentile of s	<code>np.percentile(s, [25, 50, 75])</code>
Returns the correlation between s1 and s2	<code>s1.corr(s2)</code>



Binning

<p>Sort and separate values into bins. Bins are ranges such as '50,100' or '0,50', evenly spreaded out.</p> <p>A bin represents a range of that variable. if an observation falls within that range then it belongs to the bin</p>	<pre>s.cut(bins=4, right= False)</pre>
<p>Sort and separate values into bins but qcut will choose bins such that there will be the same number of records in each bin. q is used to specify the number of bins. Qcut is frequency based.</p>	<pre>s.qcut(q = 5)</pre>



Dropping

Drops all duplicate, only keeping the last occurrence	<code>df.drop_duplicates(keep = 'last')</code>
Drops the column 'A'	<code>df.drop(columns = ['A'])</code>



Plotting I

Creates a line chart of the values in the dataframe	<code>df.plot(kind='line', x='', y='')</code>
Creates a bar chart of the values in the dataframe	<code>df.plot(kind='bar', x='', y='')</code>
Creates a scatter plot of the values in the dataframe	<code>df.plot(kind='scatter', x='', y='')</code>
Creates a boxplot of the values in the series	<code>s.plot(kind='box', x='', y='')</code>



Key

df: pandas DataFrame

s : pandas Series

nsc: NOT syntactically correct ³¹

Plotting II

Creates a distribution plot. Mainly used to look at univariate distribution.	<code>sns.distplot(kind='line', x='', y='')</code>
Creates a kdeplot. Mainly used to estimate the probability density function of a variable.	<code>sns.kdeplot(kind='bar', x='', y='')</code>
Creates a plot of two variables with bivariate and univariate graphs	<code>sns.jointplot(kind='scatter', x='', y='')</code>
Creates a violin plot (combination of boxplot and kernel density estimate). Shows the distribution of quantitative data	<code>sns.violinplot(kind='box', x='', y='')</code>



Plotting III

Creates a heatmap using the Seaborn data visualization library	<code>sns.heatmap()</code>
Creates a swarmplot using the Seaborn data visualization library. Swarmplot is a form Categorical scatterplot	<code>sns.swarmplot()</code>



fig, ax and Multiple Subplots

Creates a figure and one subplot	<code>fig, ax = plt.subplots()</code>
Creates a new figure with another figure number	<code>fig = plt.figure()</code> Variation: <code>fig = plt.figure(figsize=(16, 6))</code>
Adds an Axes to the figure. <code>add_subplot(number of rows, number of columns, index position on the grid)</code>	<code>ax1 = fig.add_subplot(1,2,1)</code> Variation: <code>ax1 = fig.add_subplot(1,2,1)</code>



ax Settings

Sets the title of the graph with a given string input	<code>ax.set_title(t)</code>
Sets the label text of the x and y axis	<code>ax.set_xlabel(y)</code> <code>ax.set_ylabel(y)</code>
Sets the x-tick labels with a list of labels	<code>ax.set_xticks(np.arange(3), ['label1', 'label2', 'label3'])</code>



Extracting from a website

Sends a request to the server to get the contents of it	<code>resp = requests.get(<Website URL>)</code>
Converts the HTML format data to a BeautifulSoup data type for extraction	<code>soup = BeautifulSoup(resp.text, 'html.parser')</code>
Allows us to see the HTML data in a readable format with indentations.	<code>soup.prettify()</code>



Extracting from a website

<p>Finds a specified tag in the soup data Returns a result and if the tag is not found, returns a None</p> <p>Attrs → if you want to find a specified tag to find</p> <p>Eg: div tag with class attribute 'about-widget'</p>	<pre>Sentence = soup.find()</pre> <p>Variation: <code>soup.find(<tagname>, attrs={'class' : '<classname>'})</code></p>
<p>Similar to <code>.find()</code>, instead this will return a list containing the results and if the tag is not found, it returns an empty list</p>	<pre>soup.find_all()</pre>
<p>Gets the text of sentence</p>	<pre>sentence.text</pre>



WordCloud

Generates a word cloud based on the string that you input	<code>WordCloud(width=800, height=400, background_color='black', max_words=100).generate(<variable name>)</code>
Shows the generated WordCloud	<code>plt.imshow()</code>

