



Pandas Cheatsheet

KEY

We'll use shorthand in this cheat sheet

`df` - A pandas DataFrame object

`s` - A pandas Series object

IMPORTS

Import these to start

```
import pandas as pd  
import numpy as np
```

IMPORTING DATA

- If file you are importing is in different directory so in place of filename, write path of your file.

CODE	WORKING
<code>pd.read_csv(filename)</code>	From a CSV file
<code>pd.read_table(filename)</code>	From a delimited text file (like TSV)
<code>pd.read_excel(filename)</code>	From an Excel file
<code>pd.read_sql(query, connection_object)</code>	Reads from a SQL table/database
<code>pd.read_json(json_string)</code>	Reads from a JSON formatted string, URL or file.

CODE	WORKING
<code>pd.DataFrame(dict)</code>	From a dict, # keys for columns names, # values for data as lists.

EXPORTING DATA

CODE	WORKING
<code>df.to_csv(filename)</code>	Writes to a CSV file
<code>df.to_excel(filename)</code>	Writes to an Excel file
<code>df.to_sql(table_name, connection_object)</code>	Writes to a SQL table
<code>df.to_json(filename)</code>	Writes to a file in JSON format

VIEWING/INSPECTING DATA

CODE	WORKING
<code>df.head(n)</code>	First n rows of the DataFrame
<code>df.tail(n)</code>	Last n rows of the DataFrame
<code>df.shape</code>	Number of rows and columns
<code>df.info()</code>	Index, Datatype and Memory information
<code>df.describe()</code>	Summary statistics for numerical columns
<code>s.value_counts(dropna=False)</code>	Views unique values and counts
<code>df.apply(pd.Series.value_counts)</code>	Unique values and counts for all columns

SELECTION

CODE	WORKING
<code>df[col]</code>	Returns column with label col as series
<code>df[[col1, col2]]</code>	Returns Columns as a new DataFrame
<code>s.iloc[0]</code>	Selection by position
<code>s.loc[0]</code>	Selection by index
<code>df.iloc[0, :]</code>	First row
<code>df.iloc[0, 0]</code>	First element of first column

DATA CLEANING

CODE	WORKING
<code>df.columns = ['a', 'b', 'c']</code>	Renames columns
<code>pd.isnull()</code>	Checks for null Values, Returns Boolean Array
<code>pd.notnull()</code>	Opposite of <code>s.isnull()</code>
<code>pd.dropna()</code>	Drops all rows that contain null values
<code>df.dropna(axis=1)</code>	Drops all columns that contain null values.
<code>df.dropna(thresh=n)</code>	Drop all columns that have fewer than <code>n</code> non-NaN values
<code>df.fillna(x)</code>	Replaces all null values with <code>x</code>
<code>s.fillna(s.mean())</code>	Replaces all null values with the mean(mean can be replaced with almost any function from the statistics section)
<code>s.astype(float)</code>	Converts the datatype of the series to float
<code>s.replace(1, 'one')</code>	Replaces all values equal to <code>1</code> with ' <code>one</code> '
<code>s.replace([1, 3], ['one', 'three'])</code>	Replaces all <code>1</code> with ' <code>one</code> 'and <code>3</code> with ' <code>three</code> '
<code>df.rename(columns=lambda x: x + 1)</code>	Mass renaming of columns
<code>df.rename(columns={'old_name': 'new_name'})</code>	Selective renaming
<code>df.set_index('column_one')</code>	Changes the index
<code>df.rename(index = lambda x: x + 1)</code>	Mass renaming of index

FILTER, SORT, & GROUPBY

CODE	WORKING
<code>df[df[col] > 5]</code>	Rows where the <code>col</code> column is greater than <code>5</code>
<code>df[(df[col] > 5) & (df[col] < 7)]</code>	Rows where <code>7 > col > 5</code>
<code>df.sort_values(col1)</code>	Sorts values by <code>col1</code> in ascending order
<code>df.sort_values(col2, ascending = False)</code>	Sorts values by <code>col2</code> in descending order
<code>df.sort_values([col1, col2], ascending = [True, False])</code>	Sorts values by <code>col1</code> in ascending order then <code>col2</code> in descending order.
<code>df.groupby(col)</code>	Returns a groupby object for values from one columns

CODE	WORKING
<code>df.groupby([col1, col2])</code>	Returns a groupby object values from multiple columns
<code>df.groupby(col1)[col2].mean()</code>	Returns the mean of the values in col2 , grouped by the values in col1 (mean can be replaced with almost any function from the statistics section)
<code>df.groupby(col1).agg(np.mean)</code>	Finds the average across all columns for every unique column 1 group
<code>df.apply(np.mean)</code>	Applies a function across each column
<code>df.apply(np.max, axis = 1)</code>	Applies a function across each row.

JOIN/COMBINE

CODE	WORKING
<code>df1.append(df2)</code>	Adds the rows in df1 to the end of df2 (columns should be identical)
<code>pd.concat([df1, df2], axis=1)</code>	Adds the columns in df1 to the end of df2 (rows should be identical)

STATISTICS

CODE	WORKING
<code>df.describe()</code>	Summary statistics for numerical columns
<code>df.mean()</code>	Returns the mean of all columns
<code>df.corr()</code>	Returns the correlation between columns in a DataFrame
<code>df.count()</code>	Returns the number of non-null values in each DataFrame column
<code>df.max()</code>	Return the highest value in each column
<code>df.min()</code>	Returns the lowest value in each column
<code>df.median()</code>	Returns the median of each column
<code>df.std()</code>	Returns the standard deviation of each column