

# Pandas Basics Cheatsheet 2.0

## Importing Data

From a CSV	<code>pd.read_csv(filename)</code>
From a delimited text file (like TSV)	<code>pd.read_table(filename)</code>
From an Excel	<code>pd.read_excel(filename)</code>
Read from a SQL	<code>pd.read_sql(query, connection_object)</code>
Read from a JSON formatted string, URL or file	<code>pd.read_json(json_string)</code>
Parses an html URL, string or file and extracts tables to a list of dataframes	<code>pd.read_html(url)</code>
Takes the contents of your clipboard and passes it to <code>read_table()</code>	<code>pd.read_clipboard()</code>
From a dict, keys for columns names, values for data as lists	<code>pd.DataFrame(dict)</code>

## Exporting Data

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

## Create Test Objects

5 columns and 20 rows of random floats	<code>pd.DataFrame(np.random.rand(20,5))</code>
Create a series from an iterable <code>my_list</code>	<code>pd.Series(my_list)</code>
Add a date index	<code>df.index = pd.date_range('1900/1/30', periods=df.shape[0])</code>

## Viewing/Inspecting Data

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

## Selection

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

## Statistics

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

## Data Cleaning

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

## Filter, Sort, and Groupby

Rows where the column <code>col</code> is greater than 0.5	<code>df[df[col] &gt; 0.5]</code>
Rows where $0.7 > col > 0.5$	<code>df[(df[col] &gt; 0.5) &amp; (df[col] &lt; 0.7)]</code>
Sort values by <code>col1</code> in ascending order	<code>df.sort_values(col1)</code>
Sort values by <code>col2</code> in descending order	<code>df.sort_values(col2,ascending=False)</code>
Sort values by <code>col1</code> in ascending order then <code>col2</code> in descending order	<code>df.sort_values([col1,col2],ascending=[True,False])</code>
Returns a groupby object for values from one column	<code>df.groupby(col)</code>
Returns groupby object for values from multiple columns	<code>df.groupby([col1,col2])</code>
Returns the mean of the values in <code>col2</code> , grouped by the values in <code>col1</code> (mean can be replaced with almost any function)	<code>df.groupby(col1)[col2]</code>
Apply the function <code>np.mean()</code> across each column	<code>df.apply(np.mean)</code>