

Importing Data	Exploring Data
<p>pd.read_csv(filename) # From a CSV file pd.read_table(filename) # From a delimited TSV file</p> <p>pd.read_excel(filename) # From an Excel file</p> <p>pd.read_sql(query, connection_object) # Reads from a SQL table/database</p> <p>pd.read_json(json_string) # Reads from a JSON formatted string, URL or file.</p> <p>pd.read_html(url) # Parses an html URL to a list of dataframes</p> <p>pd.read_clipboard() # Takes the contents of clipboard and passes it to read_table()</p> <p>pd.DataFrame(dict) # From a dict, keys for columns names, values for data as lists</p>	<p>df.shape() # Prints number of rows and columns in dataframe</p> <p>df.head(n) # Prints first n rows of the DataFrame</p> <p>df.tail(n) # Prints last n rows of the DataFrame</p> <p>df.info() # Index, Datatype and Memory information</p> <p>df.describe() # Summary statistics for numerical columns</p> <p>s.value_counts(dropna=False) # Views unique values and counts</p> <p>df.apply(pd.Series.value_counts) # Unique values and counts for all columns</p> <p>df.describe() # Summary statistics for numerical columns</p> <p>df.mean() # Returns the mean of all columns</p> <p>df.corr() # Returns the correlation between columns in a DataFrame</p>
Selecting Data	

<p><code>df[col]</code> # Returns column with label col as Series</p> <p><code>df[[col1, col2]]</code> # Returns Columns as a new DataFrame</p> <p><code>s.iloc[0]</code> # Selection by position (selects first element)</p> <p><code>s.loc[0]</code> # Selection by index (selects element at index 0)</p> <p><code>df.iloc[0,:]</code> # First row <code>df.iloc[0,0]</code> # First element of first column</p> <p><code>df[df[col] > 0.5]</code> # Rows where the col column is greater than 0.5</p> <p><code>df[(df[col] > 0.5) & (df[col] < 0.7)]</code> # Rows where $0.5 < \text{col} < 0.7$</p>	<p><code>df.count()</code> # Returns the number of non-null values in each DataFrame column</p> <p><code>df.max()</code> # Returns the highest value in each column</p> <p><code>df.min()</code> # Returns the lowest value in each column</p> <p><code>df.median()</code> # Returns the median of each column</p> <p><code>df.std()</code> # Returns the standard deviation of each column</p>
Filter, Sort, and Group By	Data Cleaning
<p><code>df.sort_values(col1)</code> # Sorts values by col1 in ascending order</p> <p><code>df.sort_values(col2,ascending=False)</code> # Sorts values by col2 in descending order</p> <p><code>df.sort_values([col1,col2], ascending=[True,False])</code> # Sorts values by col1 ascending col2 descending</p> <p><code>df.groupby(col)</code> # Returns a groupby object for values from one column</p> <p><code>df.groupby([col1,col2])</code> # Returns a groupby object values from multiple columns</p> <p><code>df.groupby(col1)[col2].mean()</code> # Returns the mean of the values in col2, grouped by the values in col1</p> <p><code>df.pivot_table(index=col1, values= col2,col3], aggfunc=mean)</code> # Creates a pivot table</p>	<p><code>df.columns = ['a','b','c']</code> # Renames columns</p> <p><code>pd.isnull()</code> # Checks for null Values, Returns Boolean Array</p> <p><code>pd.notnull()</code> # Opposite of <code>s.isnull()</code></p> <p><code>df.dropna()</code> # Drops all rows that contain null values</p> <p><code>df.dropna(axis=1)</code> # Drops all columns that contain null values</p> <p><code>df.dropna(axis=1,thresh=n)</code> # Drops all rows have have less than n non null values</p> <p><code>df.fillna(x)</code> # Replaces all null values with x</p> <p><code>s.fillna(s.mean())</code> # Replaces all null values with the mean</p> <p><code>s.astype(float)</code> # Converts the datatype of the series to float</p>

<pre>df.groupby(col1).agg(np.mean) # Finds the average across all columns for every unique column 1</pre> <pre>df.apply(np.mean) # Applies a function across each column</pre> <pre>df.apply(np.max, axis=1) # Applies a function across each row</pre>	<pre>s.replace(1,'one') # Replaces all values equal to 1 with 'one'</pre>
---	---