

Exploratory Data Analysis

Importing Data	
Function	Description
<code>pd.read_csv(file_name)</code>	Read from a csv file
<code>pd.read_csv(file_name, sep='\t')</code>	Read from a csv file separated by tabs
<code>pd.read_excel(file_name)</code>	Read from excel file
<code>pd.read_table(file_name)</code>	Read from a delimited text file
<code>pd.read_sql(sql_query, connection_object)</code>	Read from a database
<code>pd.read_json("string, url or file")</code>	Read from a json string, url or a file
<code>pd.read_html(URL)</code>	Read from a url or a file
Data Exploration	
Function	Description
<code>df.info()</code>	Provides information like datatype, shape of the dataset and memory usage
<code>df.describe()</code>	Provides information like count, mean, min, max, standard deviation and quantiles
<code>df.shape</code>	Returns the shape of the dataset
<code>df.head()</code>	Prints top 5 rows of the dataset
<code>df.tail()</code>	Prints last 5 rows of the dataset
<code>df.column_name.value_counts()</code>	Returns count of the unique classes in a column
<code>df.count()</code>	Returns total number of observations in each column
<code>df.column_name.unique()</code>	Returns unique classes in the column
Filter data	
Function	Description
<code>df.loc[condition]</code>	Returns the rows based on one condition
<code>df[(condition) & (condition)]</code>	Returns the rows based on two conditions (& operator)
<code>df[(condition) (condition)]</code>	Returns the rows based on two conditions (operator)
<code>df.loc[(condition) & (condition)]</code>	Returns the rows based on two conditions (& operator) using loc
<code>df.loc[(condition) (condition)]</code>	Returns the rows based on two conditions (operator) using loc
Renaming Columns and Indices	
Function	Description
<code>df.columns = ['Column 1', 'Column 2', ...]</code>	Rename the columns by passing a list
<code>df.rename(columns={'old_name': 'new_name'})</code>	Rename the columns using rename function
<code>df.rename(index={'old_name': 'new_name'})</code>	Rename the indices using rename function
<code>df.set_index("Column_name")</code>	Set the column as indices
Statistical Functions	
Function	Description
<code>df.mean()</code>	Finds the mean of every column
<code>df.median()</code>	Finds the median of every column
<code>df.column_name.mode()</code>	Finds the mode of a column
<code>df.corr()</code>	Creates a correlation table
<code>df.max()</code>	Finds the max value from a column
<code>df.min()</code>	Finds the min value from a column
<code>df.std()</code>	Finds the standard deviation of each column
<code>df.cov()</code>	Creates a covariance matrix
Sort and Group By	
Function	Description
<code>df.sort_values(col, ascending)</code>	Sorts the dataframe on the basis of a column
<code>df.sort_values([col1, col2, ...], ascending)</code>	Sorts the dataframe on the basis of multiple columns
<code>df.groupby(column_name)</code>	Groups a dataframe by the column name
<code>df.groupby([column_1, column_2, ...])</code>	Groups a dataframe by multiple column names
<code>df.groupby(column_1)[column_2].mean()</code>	Finds the mean of the column from the group
<code>df.groupby(column_1).agg(np.mean())</code>	Finds the mean of all the columns from the group
<code>df.apply(function, axis)</code>	Applies a function on all the columns (axis=1) or rows (axis=0) of a dataframe
Append, Concat, Join, Merge	
Function	Description
<code>df1.append(df2)</code>	Appends a dataframe df2 to df1
<code>pd.concat([df1, df2], axis)</code>	Concates multiple dataframes based on axis value
<code>df1.join(df2, on=col1, how='inner')</code>	Joins a dataframe df2 with df1 on some column
<code>pd.merge(left, right, on, how)</code>	Merge two columns on a column

EDA Cheat Sheet

Null Value Analysis and Data Cleaning

Function	Description
<code>df.isnull()</code>	Returns True where the value is null
<code>df.isnull().sum()</code>	Returns the count of null values in each column
<code>df.isnull().sum().sum()</code>	Returns the count of all the null values from a dataframe
<code>df.notnull()</code>	Returns True where the value is not null
<code>df.dropna(axis, thresh)</code>	Drops the columns (axis=1) or rows (axis=0) having null values based on threshold
<code>df.fillna(value)</code>	Fills the cells having null values with the passed value
<code>df.replace('old_value', 'new_value')</code>	Replace a value by a new value
<code>df.replace([old_1, old_2], [new_1, new_2])</code>	Replace multiple values with multiple new values
<code>df.column_name.astype('data_type')</code>	Change the data type of the column

Selecting rows and columns

Function	Description
<code>df.column_name</code>	Select the column using. Note: a column having white spaces cannot be selected by this method
<code>df["column_name"]</code>	Select a column
<code>df[["column_name_1", "column_name_2", ...]]</code>	Select multiple columns
<code>df.iloc[: , :]</code>	Pass the row and column start and end indices to extract selected rows and columns
<code>df.iloc[index_position]</code>	Pass the index position to extract rows
<code>df.loc[index_value]</code>	Pass the index value to extract rows

Write Data

Function	Description
<code>df.to_csv(file_name)</code>	Write the data from df to a csv file
<code>df.to_excel(file_name)</code>	Write the data from df to an excel file
<code>df.to_html(file_name)</code>	Write the data from df to a html file
<code>df.to_sql(table_name, connection_object)</code>	Write the data from df to a table in a database
<code>df.to_json(file_name)</code>	Write the data from df to a json file

Duplicates

Function	Description
<code>df.duplicated(keep='first')</code>	Find the first occurring duplicates.
<code>df.drop_duplicates(keep, inplace)</code>	Drop the duplicate rows