

# Pandas Cheatsheet

## Dataframe

Load dataset

```
pd.read_csv(csv_file)  
pd.read_excel(excel_file)
```

Python ▾

Set index

```
df.set_index(keys=col_name, inplace=True)
```

Python ▾

Get basic information about the dataframe

```
df.head()  
df.describe()  
df.info()  
df.columns
```

Python ▾

Get size of dataset

```
df.shape
```

Python ▾

## Drop rows

Drop the rows where at least one of the elements is missing

```
df.dropna()
```

Python ▾

Drop the rows where all of the elements are missing

```
df.dropna(how='all')
```

Python ▾

Drop the rows where columns in subset are missing

```
df.dropna(subset=[col_name])
```

Python ▾



# Pandas Cheatsheet

## Drop rows

Drop the rows at index\_number

```
df.drop(index_number)
```

Python ▾

Drop the rows that do not meet condition

```
df.drop(df[<condition>].index)  
ex: df.drop(df[df['High Income'] == False].index)
```

Python ▾

## Delete a column

Delete a column

```
del df[col_name]
```

Python ▾

## Replace NaN values

Replace any NaN values with new\_value in the Dataframe

```
df.fillna(new_value, inplace=True)
```

Python ▾

Replace any NaN values with new\_value in col\_name

```
df[col_name].fillna(new_value, inplace=True)
```

Python ▾

## Sort dataset

Sort dataset by col\_name

```
df.sort_values(by=col_name)
```

Python ▾

Sort dataset by col\_name in a descending order

```
df.sort_values(by=col_name, ascending=False)
```

Python ▾

# Pandas Cheatsheet

## Unique values

Unique values in col\_name

```
df[[col_name]].drop_duplicates()  
df[[col_name_1, col_name_2]].drop_duplicates()
```

Python ▾

Number of unique values in col\_name

```
df[col_name].nunique()
```

Python ▾

## Select columns

Select specific columns using index location

```
df.iloc[rows_index, columns_index]
```

Python ▾

Select specific columns by column names

```
df[[col_name]]
```

Python ▾

## Filter dataset

Filter dataset based on criteria

```
df.loc[criteria]  
ex: df.loc[df['Age'] > 30]
```

Python ▾

Filter dataset based on multiple criteria

```
df.loc[(condition1) & (condition2)]  
df.loc[(condition1) | (condition2)]
```

Python ▾



# Pandas Cheatsheet

## Filter dataset

Filter dataset where col\_name is in list\_of\_col\_values

```
df.loc[df[col_name].isin([list_of_col_values])]  
ex: df.loc[df['State'].isin(['California', 'Texas', 'Florida'])]
```

Python ▾

Filter dataset by index value

```
df[df.index == 'SFO']
```

Python ▾

## Create new columns

String manipulation

```
df[new_col_name] = df[col_name_1] + ' ' + df[col_name_2]
```

Python ▾

Numerical column

```
df[new_col_name] = df[col_name_1] + df[col_name_2]
```

Python ▾

## Aggregate functions

Basic aggregate functions

```
df[col_name].sum()  
df[col_name].mean()  
df[col_name].min()  
df[col_name].max()  
df[col_name].count()
```

Python ▾

## Group by - you may put in multiple columns in groupby

- Group by col\_name
- Apply aggregate functions

```
df.groupby([col_name]).mean()  
df.groupby([col_name_1, col_name_2, col_name_3]).mean()
```

Python ▾

# Pandas Cheatsheet

## Group by - you may put in multiple columns in groupby

- Group by col\_name\_1
- Apply aggregate functions
- Sort by col\_name\_2

```
df.groupby([col_name_1]).mean().sort_values(by=col_name_2)
```

Python ▾

- Group by col\_name\_1
- Apply aggregate functions
- Display col\_name\_2

```
df.groupby([col_name_1])[col_name_2].mean()
```

Python ▾

- Group by col\_name\_1
- Apply aggregate functions
- Display col\_name\_2
- Sort by the result of aggregate functions

```
df.groupby([col_name_1])[col_name_2].mean().sort_values(ascending=False)
```

Python ▾

- Group by col\_name
- Apply multiple aggregate functions in agg()
- Display col\_name\_2

```
df.groupby([col_name_1])[col_name_2].agg(['sum', 'count'])
```

Python ▾

## Join dataframes

### Merge

```
```python
DataFrame.merge(right, how='inner', on=None, left_on=None, right_on=None,
left_index=False, right_index=False, sort=False, suffixes='_x', '_y', copy=True,
indicator=False, validate=None)
```

```

### Merge example

```
pd.merge(df1, df2, how='right', on=['First Name', 'Last Name'])
```

Python ▾

# Pandas Cheatsheet

## Join dataframes

### Join

```
DataFrame.join(other, on=None, how='left', lsuffix='', rsuffix='', sort=False)
```

Python ▾

### Join example

```
df1.join(df2, lsuffix='_1', rsuffix='_2')
```

Python ▾

## Series

### Create a series

#### Scalar

```
scalar_series = pd.Series(5, index =[0, 1, 2, 3, 4, 5])
```

Python ▾

#### List

```
array = np.array(['a','b','c','d'])
pd.Series(array)
```

Python ▾

#### List with indices

```
array = np.array(['a','b','c','d'])
pd.Series(array,index=[1,2,3,4])
```

Python ▾

```
names = np.array(['Daisy', 'Matt', 'Kelly', 'Mike'])
cities=["Atlanta", "San Francisco", "New York", "Seattle"]

pd.Series(names, index=cities)
```

Python ▾

# Pandas Cheatsheet

## Series

### Create a series

Dict

```
aDict = {'Apple':3, 'Banana':5, 'Cherry': 2, 'Mango': 13, 'Peach': 10}
pd.Series(aDict)
```

Python ▾

### Sort a Series

By index

```
s = pd.Series(aDict)
s.sort_index()
```

Python ▾

By value

```
s = pd.Series(aDict)
s.sort_values()
```

JavaScript ▾

Descending order

```
s.sort_index(ascending=False)
s.sort_values(ascending=False)
```

Python ▾

### Access an element

By index number

```
array = np.array(['a','b','c','d','e','f','g','h'])
s = pd.Series(array)

s[0]
s[:3]
```

Python ▾

# Pandas Cheatsheet

## Access an element

By index label

```
aDict = {'Apple':3, 'Banana':5, 'Cherry': 2, 'Peach': 10}  
s = pd.Series(aDict)  
  
s['Apple']  
s['Peach']
```

Python ▾

## Binary operations

Add

```
s1.add(s2)  
s1 + s2
```

Python ▾

Subtract

```
s1.sub(s2)  
s1 - s2
```

Python ▾

Multiply

```
s1.mul(s2)  
s1 * s2
```

Python ▾