

# Pandas Python Module

## Computing for Data Analytics (CPSC 4800)

Mourad Bouguerra  
mbouguerra@langara.ca

Langara College

June, 2022




## 1 Lesson's Learning Objectives

## 2 Pandas Overview

- Pandas Installation
- Pandas Series
- Pandas DataFrame
  - Reading/Writing Data
  - Basic Data Exploratory Analysis (DEA)
  - Data Filtering

## Learning Objectives

-  Upon **completion** of this lesson, you will be **able** to
- ☐ generate **Pandas Series** and **DataFrames**
  - ☐ use read **data** from **different sources** into a **Pandas DataFrame**
  - ☐ access and filter **Pandas DataFrame** rows and columns using different methods

## Pandas Module

❑ Pandas Python module provides

➡ fast, flexible, and efficient data structures that enables

✓ easy analysis & manipulation of relational/tabular<sup>a</sup> data

❑ Pandas Python module provides two data structures

➡ Series (1-dimensional)

➡ DataFrame (2-dimensional)

---

<sup>a</sup>Data in the format rows x columns.

## Pandas Installation

- ❑ **Pandas** can be installed from **terminal** using one of the following  
➡ **command-lines**

```
1 python -m pip install pandas
2 pip install pandas
3 conda install pandas
```

## Pandas Version

- ❑ To check **Pandas** version

```
pd.__version__
```

- ❑ To check **Pandas** dependencies **versions**

```
pd.show_versions()
```

## Pandas Series

❑ Pandas series is **one-dimensional** object that consists of

➡ **index**

➡ **values**

❑ The **index** identifies each **data point**

```
1 import pandas as pd
2 series = pd.Series(data=[1,2,3])
3 print(series)
4 # Output
5 # 0      1
6 # 1      2
7 # 2      3
8 # dtype: int64
```

## Pandas Series

❑ Pandas series is **one-dimensional** object that consists of

➡ **index**

➡ **values**

❑ The **index** identifies each **data point**

```
1 import pandas as pd
2 series = pd.Series(data=[1,2,3])
3 print(f'The index is {series.index}')
4 print(f'The values are {series.values}')
5 # Output
6 # The index is RangeIndex(start=0, stop=3, step=1)
7 # The values are [1 2 3]
```



## Pandas Series

- ❑ **Series index** identifies each **data point**
  - ➡ **row names**
- ❑ If **NOT specified**, a default **index** of integers
  - ➡ **0 to  $\text{len}(\text{Series}) - 1$**

```
1 import pandas as pd
2 series = pd.Series(data=[1,2,3], index=list('abc'), \
3     dtype='uint8')
4 print(series)
5 # Output
6 # a      1
7 # b      2
8 # c      3
9 # dtype: uint8
```

## Pandas Series

- ❑ **Series index** identifies each **data point**  
    ➡ **row names**
- ❑ Can be used to access **Series** elements

```
1 import pandas as pd
2 series = pd.Series(data=[1,2,3],index=list('abc'),\
3                     dtype='uint8')
4 print(f'The index is {series.index}')
5 print(f'The values are {series.values}')
6 # Output
7 # The index is Index(['a', 'b', 'c'], dtype='object')
8 # The values are [1 2 3]
```

## Pandas Series

- ❑ To access **Series** elements, we can
  - ➡ **Series index/index slicing**

```
1 print(f'series[1:2]={series[1:2]}')
2 # Output
3 # series[1:2]=b      2
4 # dtype: uint8
5 print(f"series['a':'b']={series['a':'b']}")
6 # Output
7 # a      1
8 # b      2
9 # dtype: uint8
```

# Class Activity

❑ Given the following series

- ✎ What is the Series index?
- ✎ What is the Series values?
- ✎ What is the data type of the series elements?

```
1 import pandas as pd
2 import random as r
3 series = pd.Series(data=[r.random() for i in range(10)])
```

Chinese  
Proverb

I Hear & I Forget, I See & I  
Remember, I Do & I Understand



# Class Activity

- ❑ Given the following series
- ✎ What is the Series index?
  - ✎ What is the Series values?
  - ✎ What is the data type of the series elements?

```
1 import pandas as pd
2 alphabet='abcdefghijklmnopqrstuvwxyz'
3 series = pd.Series([x for x in range(1,11)],\
4                     index=list(alphabet[0:10]),\
5                     dtype='uint8')
6 series['d':'i'] = 255
7 series[-1] = 100
```

Chinese  
Proverb

I Hear & I Forget, I See & I  
Remember, I Do & I Understand



# Class Activity

- ❑ Given the following series
- ✎ What is the Series index?
  - ✎ What is the Series values?
  - ✎ What is the data type of the series elements?

```
1 import pandas as pd
2 alphabet='abcdefghijklmnopqrstuvwxyz'
3 series = pd.Series([x for x in range(1,11)],\
4                     index=list(alphabet[0:10]),\
5                     dtype='uint8')
6 series[1:5] = 255
7 series[[-1,-3]] = 100
```

Chinese  
Proverb

I Hear & I Forget, I See & I  
Remember, I Do & I Understand



## Pandas Series

- ❑ **Series** can be **generated** from a **dictionary**
  - ➡ The series **index** is the dictionary **keys** set
  - ➡ The series **values** are the dictionary **values**

```
1 import pandas as pd
2 scores = {'quizzes':89.3,'assignments':97.7,\
3          'midterm':75.45,'final':99.87}
4 series = pd.Series(data=scores)
5 print(series)
6 # Output
7 # quizzes      89.30
8 # assignments  97.70
9 # midterm      75.45
10 # final        99.87
11 # dtype: float64
```

## Pandas Series Functions

❑ To get the **unique values** of a **series** use

➡ **unique()** method

```
1 import pandas as pd
2 colors = pd.Series(data=['red','red','green','blue',\
3     'blue','red','red','green','blue'])
4 print(f'{colors.unique()}')
5 # Output
6 # ['red' 'green' 'blue']
```



## Pandas Series Functions

❑ To get the **frequency** of each value in a **series** use

➡ **value\_counts()** method

```
1 print(f'Frequencies\n{colors.value_counts()}')
2 # Output
3 # Frequencies
4 # red      4
5 # blue     3
6 # green    2
7 # dtype: int64
```

## Pandas DataFrame

- ❑ Pandas DataFrame is **two-dimensional** object that consists of
  - ➡ **rows** and **columns**
  - ➡ Each **column** is a **Pandas** series
  - ➡ Each **column** elements must be of the **data type**
- ❑ Pandas DataFrame has both
  - ➡ **row** index
  - ➡ **column** index
- ❑ Pandas DataFrame is generally constructed by **reading** data

## Pandas DataFrame

- ❑ **Pandas** DataFrame can be constructed from a  
➡ **two-dimensional NumPy** object

```
1 import pandas as pd
2 import string
3 alphabet = string.ascii_uppercase
4 df = pd.DataFrame(np.random.randn(1000, 4), \
5                   columns=list('abcd'))
6 df = pd.DataFrame(np.random.randn(4, 10), \
7                   index=list(alphabet[:4], \
8                             columns=list(alphabet[:10])))
```

# Class Activity

- ❑ Given the following **DataFrames**
  - ✎ What is the **DataFrame rows** index?
  - ✎ What is the **DataFrame columns** index?
  - ✎ What is the content of the **DataFrame**?

```
1 import pandas as pd
2 import string
3 alphabet = string.ascii_uppercase
4 df_1 = pd.DataFrame(np.random.randn(10, 5))
5 df_2 = pd.DataFrame(np.random.randn(10, 5), \
6                     index=list(alphabet)[:10])
7 df_3 = pd.DataFrame(np.random.randn(5, 10), \
8                     index=list(alphabet)[:5], \
9                     columns=list(alphabet[:10]))
```

## Class Activity

❑ What is the output of the following **Python** script?

```
1 df_1 = pd.DataFrame([np.arange(5), np.arange(5, 10), np.arange(10, 15)])
2 df_2 = pd.DataFrame(np.arange(15).reshape(3, 5))
3 df_3 = pd.DataFrame(np.arange(15).reshape(5, 3))
4 print(f'df_1=\n{df_1}')
5 print(f'df_2=\n{df_2}')
6 print(f'df_3=\n{df_3}')
```

Chinese  
Proverb

I **Hear** & I **Forget**, I **See** & I  
**Remember**, I **Do** & I **Understand**



## Pandas DataFrame

❑ To add a **new column** to a **Pandas DataFrame** use

➡ **df['NewColumnName']=List**

```
1 df = pd.DataFrame(np.random.randn(4,10))  
2 df['NewColumn']=['One','Two','Three','Four']
```

## Pandas DataFrame

❑ To remove a **column** from **Pandas DataFrame** use

➡ `df.drop('ColumnName', inplace=True, axis=1)`

```
1 df = pd.DataFrame(np.random.randn(100, 4), \
2                     columns=['One', 'Two', 'Three', 'Four'])
3 # remove column Three
4 df.drop('Three', inplace=True, axis=1)
```

## Class Activity

❑ What is the output of the following **Python** script?

```
1 df = pd.DataFrame(np.arange(20).reshape(4, 5))
2 df['Key']=['One','Two','Three','Four']
3 print(f'df=\n{df}')
```

Chinese  
Proverb

Tell Me & I Forget,  
Teach Me & I Remember,  
Involve Me & I Learn





# Reading/Writing Data



**Figure:** Pandas Reading/Writing\*

---

\* [https://pandas.pydata.org/docs/getting\\_started/intro\\_tutorials/02\\_read\\_write.html](https://pandas.pydata.org/docs/getting_started/intro_tutorials/02_read_write.html)

## Class Activity

❑ Using **list comprehension**, display all the **Pandas read** functions

➡ All **Pandas read** functions start with **read**

❑ Using **list comprehension**, display all the **Pandas write** functions

➡ All **Pandas write** functions start with **to**

Chinese  
Proverb

I **Hear** & I **Forget**, I **See** & I  
**Remember**, I **Do** & I **Understand**



## Pandas DataFrame

❑ **Pandas DataFrame** can be constructed by reading different data format

➡ CSV

➡ Excel

➡ JSON

➡ HTML tables

```
1 import pandas as pd
2 df = pd.read_csv('data/arh.csv')
3 df = pd.read_table('data/arh.csv', sep=',')
4 df = pd.read_excel('data/arh.csv')
5 df = pd.read_clipboard()
```

## Class Activity

❑ Using **list comprehension**, display all the **DataFrame write** functions

➡ All **DataFrame write** functions start with **to**

Chinese  
Proverb

Tell Me & I Forget,  
Teach Me & I Remember,  
Involve Me & I Learn



## Pandas DataFrame

❑ To check the size of the file in **disk** and **memory**

```
1 import pandas as pd
2 import os
3 MB = 2**20
4 filename = 'data/arh.csv'
5 df = pd.read_csv(filename)
6 size_in_disk = os.stat(filename).st_size/MB
7 size_in_ram = df.memory_usage('deep').sum()/MB
8 print(f'The size of the file in disk is {size_in_disk:.2f} MB.')
9 print(f'The size of the file in ram is {size_in_ram:.2f} MB.')
```

## Pandas DataFrame

❑ To display **rows/columns** index of the **Pandas DataFrame** use

➡ **index** property

➡ **columns** property

❑ To get the numbers of **rows** and **column** of the **Pandas DataFrame** use

➡ **shape** property

```
1 df.index
2 df.columns
3 df.shape
```

## Pandas DataFrame

❑ To display the **first/last** rows of the **Pandas DataFrame** use

➡ **head()**

➡ **tail**

```
1 import pandas as pd
2 df = pd.read_csv('data/arh.csv')
3 df.head() # first five rows
4 df.tail() # last five rows
5 df.head(15) # first 15 rows
6 df.tail(20) # last 20 rows
```

## Pandas DataFrame

- ❑ To get detailed information about each **column** of the **Pandas DataFrame** use
  - ➡ **info()** function
- ❑ To get **RAM** usage of each **column** of the **Pandas DataFrame** use
  - ➡ **memory\_usage('deep')** function
- ❑ To get the **five-number** summary of the **numerical columns** of the **Pandas DataFrame** use
  - ➡ **describe()** function<sup>a</sup>

```
1 df.info()  
2 df.memory_usage('deep')  
3 df.describe()
```

---

<sup>a</sup>It also computes the **count**, **mean**, and **standard deviation**.



## Pandas DataFrame

- ❑ To check for **missing** values in **Pandas DataFrame** use  
    ➔ **isna()** function
- ❑ To display the number of **missing** values in each column use  
    ➔ **isna().sum()** function<sup>a</sup>

```
1 df.isna()  
2 df.isna().sum()
```

---

<sup>a</sup>It also computes the **count**, **mean**, and **standard deviation**.

## Pandas DataFrame

- ❑ To check for **duplicate** rows in **Pandas DataFrame** use  
    ➔ **uplicated()** function
- ❑ To remove **missing** values rows use  
    ➔ **dropna()** function
- ❑ To remove **duplicate** rows use  
    ➔ **drop\_duplicates()** function

```
1 df.duplicated()  
2 clean_df = df.dropna()  
3 clean_df = df.drop_duplicates()
```

## Pandas Data Filtering

❑ To select **rows** or **columns** from a **Pandas DataFrame** use the methods

➡ **loc()**

➡ **iloc()**

## Pandas DataFrame iloc() Method

- ❑ Pandas `loc()` is used to **select rows** and/or **columns** by  
    ➔ **integer position** of the **rows** and **columns** index
- ❑ In **slicing**, the **end** is **NOT** included

```
1 # To select the first 5 rows and the second & third
   columns
2 movies.iloc[0:5,1:3]
3 # To select last three rows with all the columns
4 movies.iloc[-3:,:]
5 movies.iloc[-3:]
6 # This will generate an error
7 movies.iloc[:,['title','star_rating']]
8 # .iloc requires numeric indexers
```

# Class Activity

❑ What does the following Python code accomplish?

```
1 movies.index
2 movies.columns
3 movies.iloc[-5:]
4 movies.iloc[:,2].unique()
5 movies.iloc[:,3].value_counts()
6 movies.iloc[[x for x in movies.index if x % 2 == 1],\
7             [1,0]]
```

Chinese  
Proverb

Tell Me & I Forget,  
Teach Me & I Remember,  
Involve Me & I Learn



## Pandas DataFrame loc() Method

- ❑ Pandas `loc()` is used to **select rows** and/or **columns** by  
    ➔ the **index** label
- ❑ Both sides of the **slicing** are **included** in `loc()`

```
1 # To select the first fourth rows with all the columns
2 movies.loc[0:3,:]
3 # To select all the rows for the title column
4 movies.loc[:, 'title']
5 # To select all the rows for the title and star_rating
   columns
6 movies.loc[:, ['title', 'star_rating']]
7 # This will generate an error
8 movies.loc[:, [1, 0]]
```

## Class Activity

❑ What does the following **Python** code accomplish?

```
1 movies.loc[975:978]
2 movies.loc[:, 'genre'].unique()
3 movies.loc[:, 'content_rating'].value_counts()
4 movies.loc[[x for x in movies.index if x % 2 == 1], \
5             ['title', 'star_rating']]
```

Chinese  
Proverb

Tell Me & I Forget,  
Teach Me & I Remember,  
Involve Me & I Learn



## Filtering Using Boolean Conditions

- ❑ To filter **Pandas DataFrame** by **column value** use
  - ➡ a **Boolean** condition on the **column value**
- ❑ **dot** operator can be only used if the **column name**
  - ➡ a **valid Python** variable name

```
1 # To filter only action movies
2 condition = movies.genre == 'Action'
3 # or
4 condition = movies['genre'] == 'Action'
5 movies[condition]
6 # or
7 movies.loc[condition,:]
```



## Filtering Using Boolean Conditions

- ❑ To filter **Pandas DataFrame** by **column value** use
  - ➡ a **Boolean** condition on the **column value**
- ❑ **dot** operator can be only used if the **column name**
  - ➡ a **valid Python** variable name

```
1 # To filter only movies of more than 2 hours
2 long_movies = movies.duration >= 120
3 # or
4 long_movies = movies['duration'] >= 120
5 movies[long_movies] # with all columns
6 movies.loc[long_movies, 'star_rating': 'duration']
```

## Filtering Using Boolean Conditions

- ❑ To combine multiple **Boolean** conditions
  - ➡ a **Boolean** condition on the **column value**
- ❑ **and** using the operator **&**
- ❑ **or** using the operator **|**
  - ➡ a **valid Python** variable name

```
1 # To filter only Drama genre movies
2 # with R content rating
3 condition = (movies.genre == 'Drama') &\
4             (movies.content_rating == 'R')
5 movies[condition]
6 movies.loc[:,condition]
```

## Class Activity

- ☐ Filter all the **drama** movies that have a **duration** more than **200** minutes
- ☐ Filter all the movies that have a **duration** more than **200** minutes or less than **90** minutes
- ☐ What are the top **10 rated** movies?
- ☐ What are the worst **10 rated** movies?

Chinese  
Proverb

**Tell Me & I Forget,  
Teach Me & I Remember,  
Involve Me & I Learn**

