

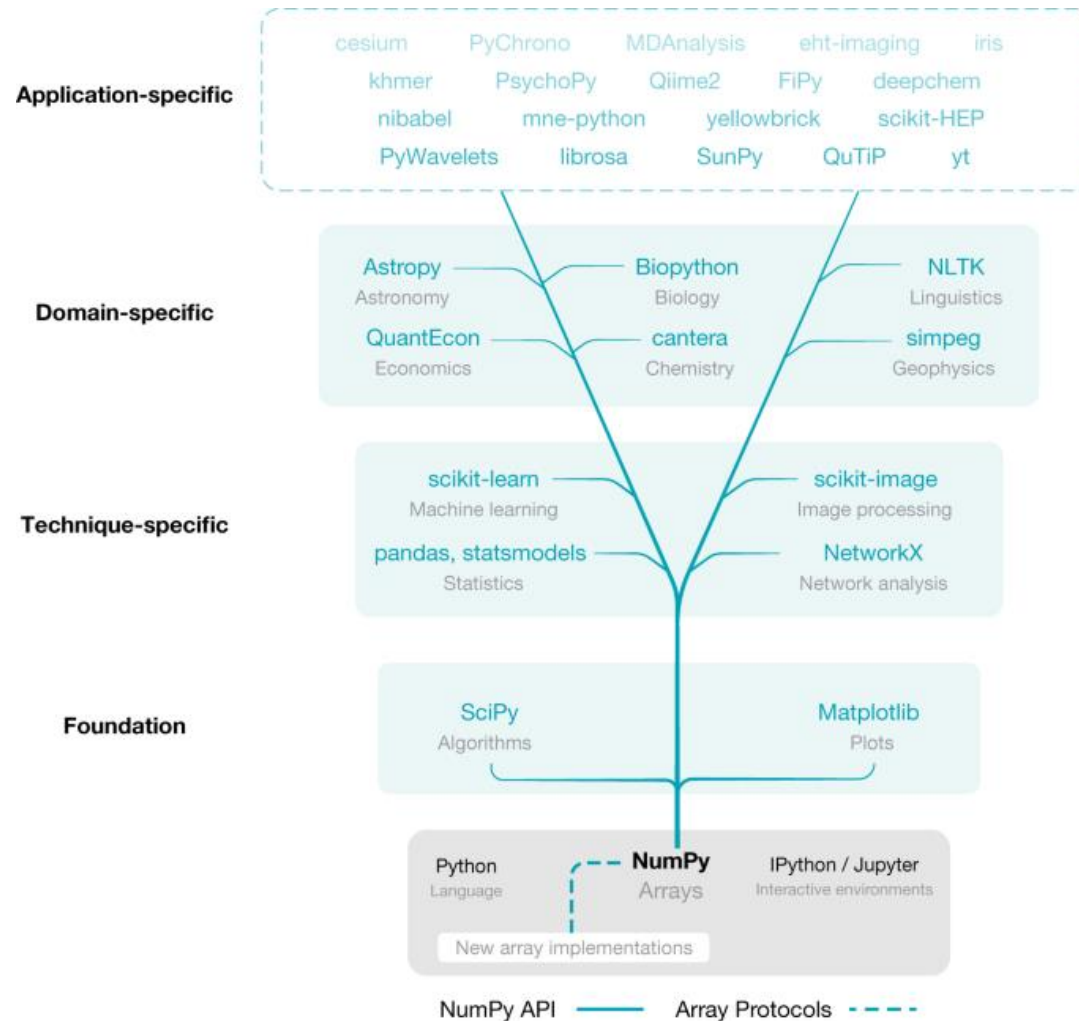
2022년 가을학기

Pandas



In Last Lecture

❖ Numpy 중요성



In Last Lecture

Python For Data Science Cheat Sheet

NumPy Basics

Learn Python for Data Science Interactively at www.DataCamp.com



NumPy

The NumPy library is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays.

Use the following import convention:

```
>>> import numpy as np
```



NumPy Arrays

1D array

```
[1 2 3]
```

2D array

```
axis 1  
axis 0  
[[1.5 2.3]  
 [4.5 6.]]
```

3D array

```
axis 2  
axis 1  
axis 0  
[[[1.5 2.3]  
 [4.5 6.]]  
 [[1.5 2.3]  
 [4.5 6.]]  
 [[1.5 2.3]  
 [4.5 6.]]]
```

Creating Arrays

```
>>> a = np.array([1,2,3])  
>>> b = np.array([(1.5,2,3), (4,5,6)], dtype=float)  
>>> c = np.array([(1.5,2,3), (4,5,6)], [(3,2,1), (4,5,6)],  
               dtype=float)
```

Initial Placeholders

```
>>> np.zeros((2,4))  
>>> np.ones((2,3,4),dtype=np.int16)  
>>> d = np.arange(10,25,5)  
  
>>> np.linspace(0,2,9)  
  
>>> e = np.full((2,2),7)  
>>> f = np.eye(2)  
>>> np.random.random((2,2))  
>>> np.empty((3,2))
```

Create an array of zeros
Create an array of ones
Create an array of evenly spaced values (step value)
Create an array of evenly spaced values (number of samples)
Create a constant array
Create a 2x2 Identity matrix
Create an array with random values
Create an empty array

I/O

Saving & Loading On Disk

```
>>> np.save('my_array', a)  
>>> np.savez('array.npz', a, b)  
>>> np.load('my_array.npy')
```

Saving & Loading Text Files

```
>>> np.loadtxt('myfile.txt')  
>>> np.genfromtxt('my_file.csv', delimiter=',')  
>>> np.savetxt('myarray.txt', a, delimiter=' ')
```

Data Types

```
>>> np.int64  
>>> np.float64  
>>> np.complex  
>>> np.bool  
>>> np.object  
>>> np.string_  
>>> np.unicode_
```

Signed 64-bit integer types
Standard double-precision floating point
Complex numbers represented by 128 floats
Boolean type storing TRUE and FALSE values
Python object type
Fixed-length string type
Fixed-length unicode type

Inspecting Your Array

```
>>> a.shape  
>>> len(a)  
>>> b.ndim  
>>> e.size  
>>> b.dtype  
>>> b.dtype.name  
>>> b.astype(int)
```

Array dimensions
Length of array
Number of array dimensions
Number of array elements
Data type of array elements
Name of data type
Convert an array to a different type

Asking For Help

```
>>> np.info(np.ndarray.dtype)
```

Array Mathematics

Arithmetic Operations

```
>>> g = a - b  
array([[ -0.5,  0. ,  0. ],  
       [ -3. , -3. , -3. ]])  
>>> np.subtract(a,b)  
  
>>> b = a  
array([[ 2.5,  4. ,  6. ],  
       [ 8. ,  7. ,  9. ]])  
>>> np.add(b,a)  
  
>>> a / b  
array([[ 0.408167,  1. ,  1. ],  
       [ 0.25 ,  0.4 ,  0.5 ]])  
>>> np.divide(a,b)  
  
>>> a * b  
array([[ 1.5,  4. ,  8. ],  
       [ 4. ,  10. ,  18. ]])  
>>> np.multiply(a,b)  
>>> np.exp(b)  
>>> np.sqrt(b)  
>>> np.sin(a)  
>>> np.cos(b)  
>>> np.log(a)  
>>> a.dot(f)  
array([[ 7. ,  7. ],  
       [ 7. ,  7.]])
```

Subtraction
Subtraction
Addition
Addition
Division
Division
Division
Multiplication
Multiplication
Exponentiation
Square root
Print sines of an array
Element-wise cosine
Element-wise natural logarithm
Dot product

Comparison

```
>>> a == b  
array([[False,  True,  True],  
       [False, False, False]], dtype=bool)  
  
>>> a < 2  
array([[ True,  False, False],  
       [ True,  True,  True]], dtype=bool)  
>>> np.array_equal(a, b)
```

Element-wise comparison
Element-wise comparison
Array-wise comparison

Aggregate Functions

```
>>> a.sum()  
>>> a.min()  
>>> b.max(axis=0)  
>>> b.cumsum(axis=1)  
>>> a.mean()  
>>> a.median()  
>>> a.corrcoef()  
>>> np.std(b)
```

Array-wise sum
Array-wise minimum value
Maximum value of an array row
Cumulative sum of the elements
Mean
Median
Correlation coefficient
Standard deviation

Copying Arrays

```
>>> h = a.view()  
>>> np.copy(a)  
>>> h = a.copy()
```

Create a view of the array with the same data
Create a copy of the array
Create a deep copy of the array

Sorting Arrays

```
>>> a.sort()  
>>> c.sort(axis=0)
```

Sort an array
Sort the elements of an array's axis

Subsetting, Slicing, Indexing

Also see Lists

```
>>> a[2]  
3  
>>> b[1,2]  
6.0  
  
>>> a[0:2]  
array([0, 2])  
>>> b[0:2,1]  
array([ 2. ,  2. ])  
  
>>> b[::1]  
array([1.5, 2. , 3. ])  
>>> a[1,...]  
array([[ 3. ,  2. ,  1. ],  
       [ 4. ,  5. ,  6. ]])  
  
>>> a[: :-1]  
array([3, 2, 1])  
  
>>> a[a<2]  
array([1])  
  
>>> b[[1, 0, 1, 0], [0, 1, 2, 0]]  
array([ 4. ,  2. ,  6. ,  1.5])  
>>> b[[1, 0, 1, 0]] [[0,2,2,0]]  
array([[ 1.5, 2. ,  3. ,  1.5],  
       [ 2. , 2. ,  2. ,  2.5]])
```

Select the element at the 2nd index
Select the element at row 1 column 2 (equivalent to b[1][2])
Select items at index 0 and 1
Select items at rows 0 and 1 in column 1
Select all items at row 0 (equivalent to a[0,:])
Same as b[:,1:]
Reversed array a
Select elements from a less than 2
Fancy Indexing
Select elements (1,0),(0,1),(1,2) and (0,0)
Select a subset of the matrix's rows and columns

Array Manipulation

```
>>> i = np.transpose(b)  
>>> i.T  
  
>>> g.reshape(3,-2)  
  
>>> h.resize((2,6))  
>>> np.append(h,g)  
>>> np.insert(a, 1, 5)  
>>> np.delete(a, (1))  
  
>>> np.concatenate((a,d),axis=0)  
array([ 1. ,  2. ,  3. , 10. , 18. , 20])  
>>> np.vstack((a,b))  
array([[ 1.5,  2. ,  3. ],  
       [ 4. ,  5. ,  6. ]])  
  
>>> np.d_[a,f]  
>>> np.hstack((e,f))  
array([[ 7. ,  7. ,  1. ,  0. ],  
       [ 7. ,  7. ,  0. ,  1.]])  
  
>>> np.column_stack((a,d))  
array([[ 1. ,  10. ],  
       [ 2. ,  18. ],  
       [ 3. ,  20.]])  
>>> np.d_[a,d]  
  
>>> np.hsplit(a,3)  
[array([1]), array([2]), array([3])]  
>>> np.vsplit(a,2)  
array([[ 1.5,  2. ,  3. ],  
       [ 4. ,  5. ,  6. ]])
```

Transpose Array
Permute array dimensions
Permute array dimensions
Flatten the array
Reshape, but don't change data
Adding/Removing Elements
Return a new array with shape (2,6)
Append items to an array
Insert items in an array
Delete items from an array
Combining Arrays
Concatenate arrays
Stack arrays vertically (row-wise)
Stack arrays vertically (row-wise)
Stack arrays horizontally (column-wise)
Create stacked column-wise arrays
Create stacked column-wise arrays
Splitting Arrays
Split the array horizontally at the 3rd index
Split the array vertically at the 2nd index



Day



Pandas

CONTENTS

- A. What is Pandas?
- B. Pandas Data Structure
- C. Data processing in Pandas



A

What is Pandas?

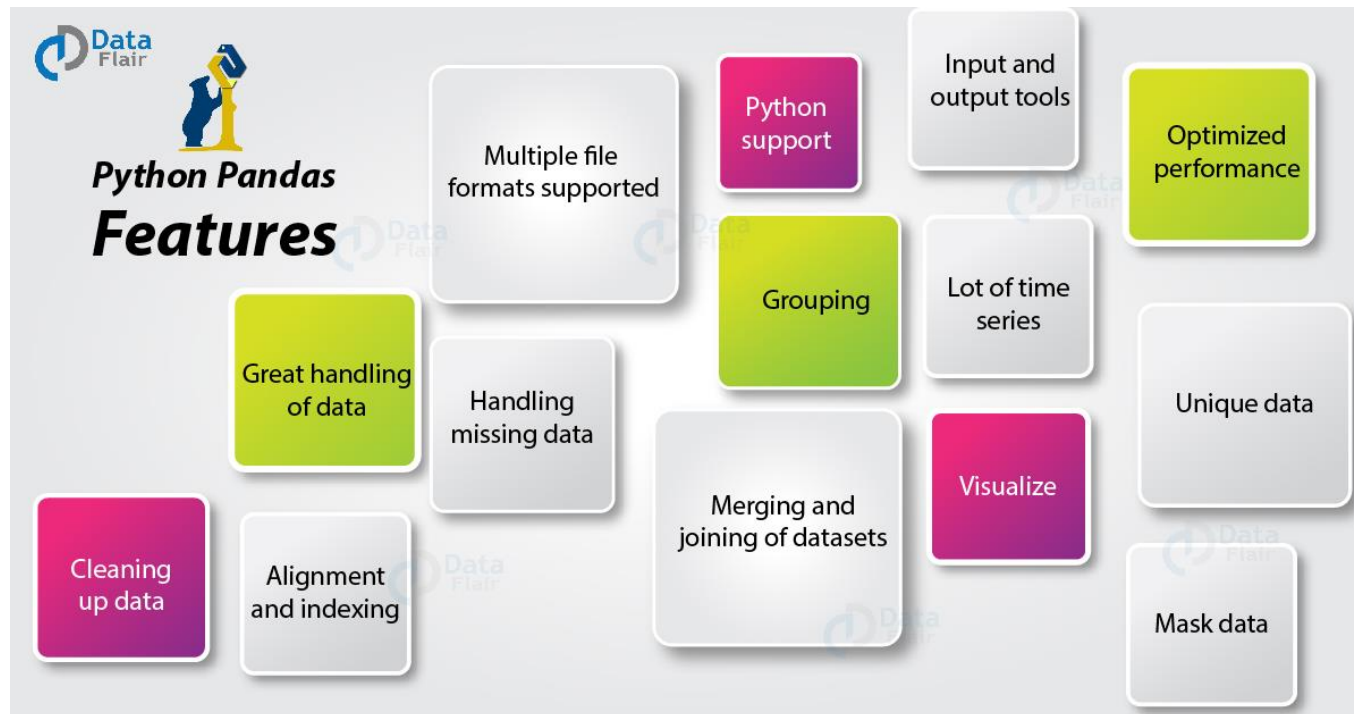
What is Pandas?

- ❖ An open-source library that is built on top of NumPy library
- ❖ One of the most popular Python libraries for data science
 - Data read, data cleaning, data transforming, and data analysis



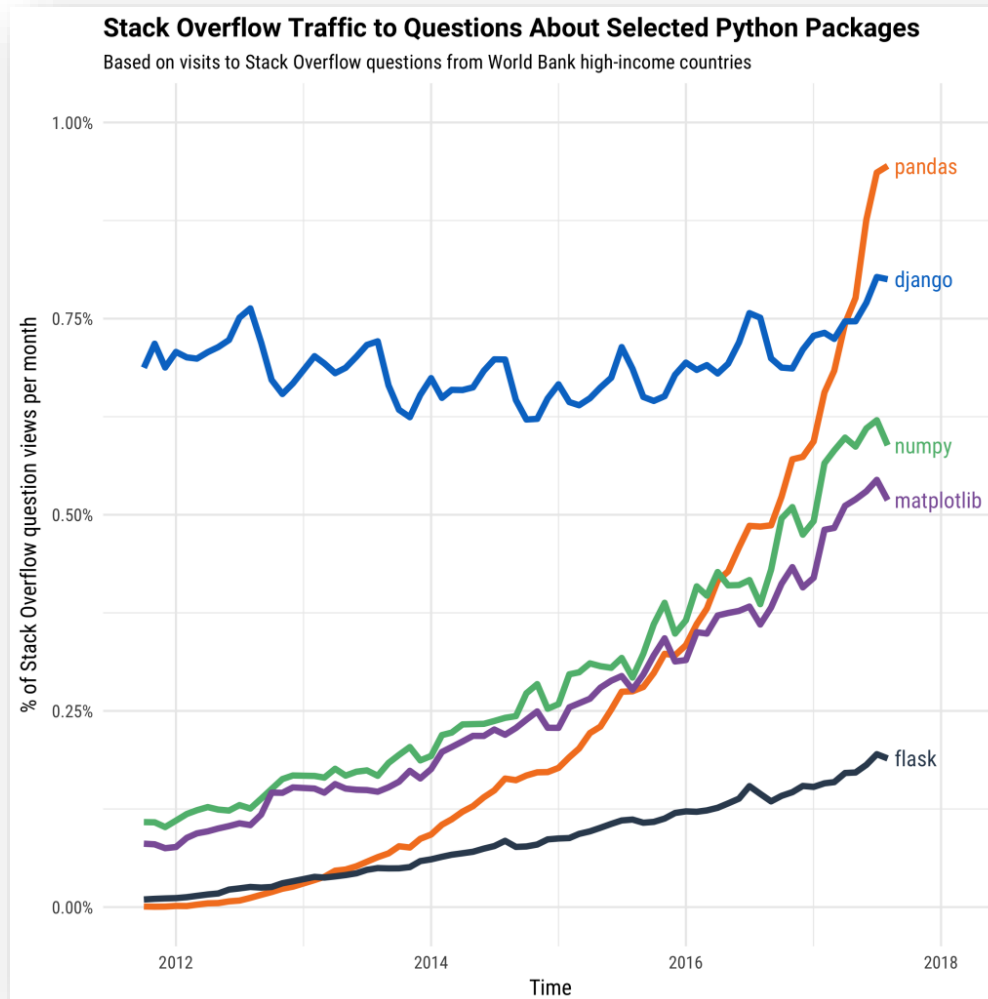
What is Pandas?

❖ Pandas features



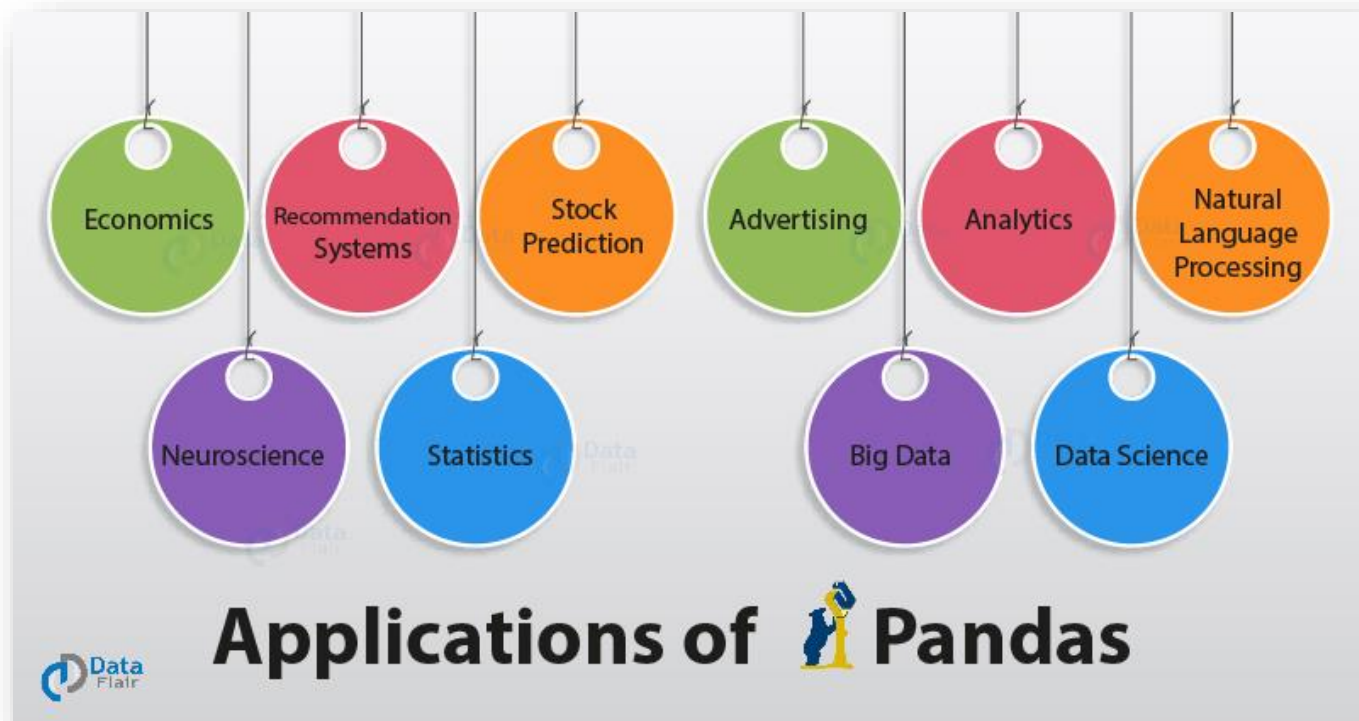
What is Pandas?

❖ Why Pandas?



What is Pandas?

❖ Pandas Application areas



What is Pandas?

- ❖ Jupyter Notebook을 실행
- ❖ 명령 프롬프트(CMD)를 실행하고 다음 명령어를 입력
 - `pip install pandas`
- ❖ 다른 라이브러리 또한 설치해야 함
 - Numpy
- ❖ pandas 라이브러리 버전 확인

```
import pandas          //Importing pandas library  
print(pandas.__version__) //Printing pandas library version
```



B

Pandas Data Structure

Pandas Data Structure

❖ Pandas의 두 가지 주요 구성 요소는 Series와 DataFrame이다

❖ Series

- A column

❖ Data frame

- Series 컬렉션으로 구성된 다차원 테이블

| Series | | | Series | | | DataFrame | | |
|--------|--------|---|--------|---------|---|-----------|--------|---------|
| | apples | | | oranges | | | apples | oranges |
| 0 | 3 | + | 0 | 0 | = | 0 | 3 | 0 |
| 1 | 2 | | 1 | 3 | | 1 | 2 | 3 |
| 2 | 0 | | 2 | 7 | | 2 | 0 | 7 |
| 3 | 1 | | 3 | 2 | | 3 | 1 | 2 |

Pandas Data Structure

❖ Series

- 레이블이 지정된 일차원 배열
- Syntax
 - `pandas.Series(data, index, dtype, copy)`
 - `data`
 - Input data in the form of ndarray, list, dict or scalar value
 - `index`
 - Index of the column
 - `dtype`
 - Data type
 - `copy`
 - Copy data

Pandas Data Structure

❖ Series

- Numpy를 이용한 Series 만드는 방법

```
import pandas as pd
import numpy as np

data = np.array(['a','b','c','d'])
series = pd.Series(data, index=[100,101,102,103])

print(series)
```

```
100 a
101 b
102 c
103 d
dtype: object
```

Pandas Data Structure

❖ Series

- dict를 이용한 Series 만드는 방법

```
import pandas as pd  
import numpy as np
```

```
data = {100 : 'a', 101 : 'b', 102 : 'c', 103 : 'd'}  
series = pd.Series(data)
```

```
print(series)
```

```
100 a  
101 b  
102 c  
103 d  
dtype: object
```


Pandas Data Structure

❖ Series

- Task: Create a series for the following column

| Index | Data |
|-------|------|
| 1 | 'A' |
| 2 | 'B' |
| 3 | 'C' |
| 4 | 'D' |
| 5 | 'E' |

Pandas Data Structure

❖ Series

- Task: Create a series for the following column

```
import pandas as pd
import numpy as np

data = np.array(['A','B','C','D', 'E'])
series = pd.Series(data, index=[1,2,3,4, 5])

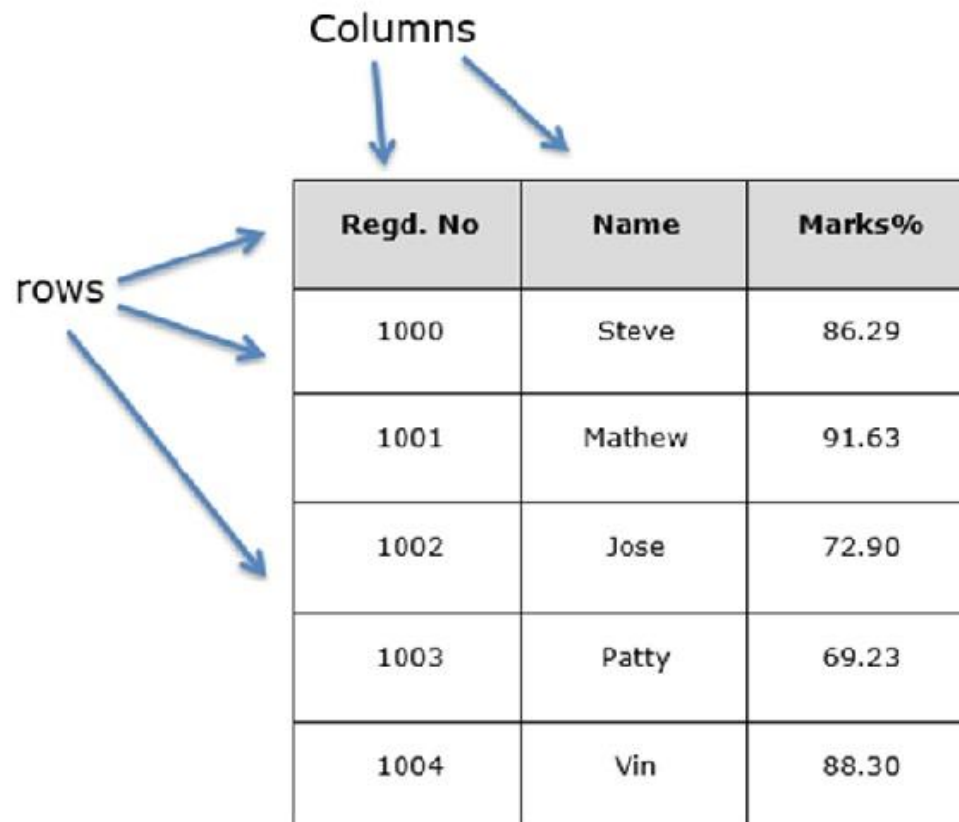
print(series)
```

```
1  A
2  B
3  C
4  D
5  E
dtype: object
```

Pandas Data Structure

❖ Data frame

- Multi-dimensional table with rows and columns



The diagram illustrates a Pandas DataFrame as a multi-dimensional table. It features a table with three columns and five rows. The columns are labeled 'Regd. No', 'Name', and 'Marks%'. The rows are labeled with registration numbers from 1000 to 1004. Blue arrows point from the labels 'Columns' and 'rows' to their respective parts of the table.

| Columns | | |
|----------|--------|--------|
| Regd. No | Name | Marks% |
| 1000 | Steve | 86.29 |
| 1001 | Mathew | 91.63 |
| 1002 | Jose | 72.90 |
| 1003 | Patty | 69.23 |
| 1004 | Vin | 88.30 |

rows

Pandas Data Structure

❖ Data frame

■ Syntax

- `pandas.DataFrame(data, index, columns, dtype, copy)`
- `data`
 - Input data in the form of series, ndarray, list, dict or scalar value
- `index`
 - Index of the data frame
- `columns`
 - Name/label of columns
- `dtype`
 - Data type
- `copy`
 - Copy data

Pandas Data Structure

❖ Data frame

- List를 이용한 data frame 만드는 방법

```
import pandas as pd

data = [['Tim',35],['Sonya',30],['Sunny',34]]

df = pd.DataFrame(data,columns=['Name','Age'],dtype=float)

print(df)
```

| | Name | Age |
|---|-------|------|
| 0 | Tim | 35.0 |
| 1 | Sonya | 30.0 |
| 2 | Sunny | 34.0 |

Pandas Data Structure

❖ Data frame

- Dict of series를 이용한 data frame 만드는 방법

```
import pandas as pd

data = {'one' : pd.Series([1, 2, 3], index=['a', 'b', 'c']),
        'two' : pd.Series([1, 2, 3, 4], index=['a', 'b', 'c', 'd'])}

df = pd.DataFrame(data)
print(df)
```

| | one | two |
|---|-----|-----|
| a | 1.0 | 1 |
| b | 2.0 | 2 |
| c | 3.0 | 3 |
| d | NaN | 4 |

Pandas Data Structure

❖ Data frame

- Dict of series를 사용하여 다음 테이블에 대한 data frame 만들어보기

| | Artist | Genre | Listeners | Plays |
|---|----------------|-------|-----------|------------|
| 0 | Billie Holiday | Jazz | 1,300,000 | 27,000,000 |
| 1 | Jimi Hendrix | Rock | 2,700,000 | 70,000,000 |
| 2 | Miles Davis | Jazz | 1,500,000 | 48,000,000 |
| 3 | SIA | Pop | 2,000,000 | 74,000,000 |

Pandas Data Structure

❖ Data frame

- dict of series를 사용하여 다음 테이블에 대한 data frame 만들어보기

```
import pandas as pd
import numpy as np

data = {'Artist' : pd.Series(['Billie Holiday', 'Jimi Hendrix', 'Miles Davis', 'SIA']),
        'Genre' : pd.Series(['Jazz', 'Rock', 'Jazz', 'Pop']),
        'Listeners' : pd.Series([1300000, 2700000, 1500000, 2000000]),
        'Plays' : pd.Series([27000000, 70000000, 48000000, 74000000])}

df = pd.DataFrame(data)
print(df)
```



C

Data Processing in Pandas

Data Processing in Pandas

❖ Pandas의 read_csv 함수를 이용하여 데이터 불러오기

- Save IMDB-Movie-Data.csv file in your local storage
- index_col 속성
 - CSVs don't have indexes like our Data Frames, so all we need to do is just designate the index_col when reading
- In our case, we select "Title" column as index

```
import pandas as pd  
  
movies_df = pd.read_csv("D:/IMDB-Movie-Data.csv",  
                        index_col="Title")  
print(movies_df)
```

Data Processing in Pandas

❖ head() 함수

- Print out a first five rows of your data frame
 - `movies_df.head()`
- We could also pass a number to head()
 - `movies_df.head(10)`

❖ tail() 함수

- Print out the last five rows
 - `movies_df.tail()`
- If you want to see the last two records, then pass a number
 - `movies_df.tail(2)`

Data Processing in Pandas

❖ `info()`

- Provides essential details about your dataset
- Features
 - Number of rows
 - Number of columns
 - Number of non-null values
 - Type of data in each column
 - How much memory the data frame is using

Data Processing in Pandas

❖ info()

- `movies_df.info()`

```
<class 'pandas.core.frame.DataFrame'>  
Index: 1000 entries, Guardians of the Galaxy to Nine Lives  
Data columns (total 11 columns):  
Rank                1000 non-null int64  
Genre               1000 non-null object  
Description         1000 non-null object  
Director           1000 non-null object  
Actors             1000 non-null object  
Year               1000 non-null int64  
Runtime (Minutes)  1000 non-null int64  
Rating             1000 non-null float64  
Votes             1000 non-null int64  
Revenue (Millions) 872 non-null float64  
Metascore          936 non-null float64  
dtypes: float64(3), int64(4), object(4)  
memory usage: 74.2+ KB
```

Data Processing in Pandas

❖ shape

- Outputs just a size of rows and columns
 - `movies_df.shape`
 - `(1000, 11)`
- Note that shape has no parentheses, where output is a simple tuple of format (rows, columns)
- Used frequently when cleaning and transforming data
 - For example, you might filter some rows based on some criteria and then want to know quickly how many rows were removed

Data Processing in Pandas

❖ describe()

- Outputs summary of the distribution of continuous variables
 - `movies_df.describe()`

| | Rank | Year | Runtime (Minutes) | Rating | Votes | Revenue (Millions) | Metascore |
|-------|-------------|-------------|-------------------|-------------|--------------|--------------------|------------|
| count | 1000.000000 | 1000.000000 | 1000.000000 | 1000.000000 | 1.000000e+03 | 872.000000 | 936.000000 |
| mean | 500.500000 | 2012.783000 | 113.172000 | 6.723200 | 1.698083e+05 | 82.956376 | 58.985043 |
| std | 288.819436 | 3.205962 | 18.810908 | 0.945429 | 1.887626e+05 | 103.253540 | 17.194757 |
| min | 1.000000 | 2006.000000 | 66.000000 | 1.900000 | 6.100000e+01 | 0.000000 | 11.000000 |
| 25% | 250.750000 | 2010.000000 | 100.000000 | 6.200000 | 3.630900e+04 | 13.270000 | 47.000000 |
| 50% | 500.500000 | 2014.000000 | 111.000000 | 6.800000 | 1.107990e+05 | 47.985000 | 59.500000 |
| 75% | 750.250000 | 2016.000000 | 123.000000 | 7.400000 | 2.399098e+05 | 113.715000 | 72.000000 |
| max | 1000.000000 | 2016.000000 | 191.000000 | 9.000000 | 1.791916e+06 | 936.630000 | 100.000000 |

Data Processing in Pandas

❖ describe()

- We can also use describe() function for categorical variables
 - `movies_df['Genre'].describe()`

```
count          1000
unique          207
top      Action,Adventure,Sci-Fi
freq           50
Name: Genre, dtype: object
```

- Features
 - count of rows
 - unique count of categories
 - top category
 - freq of top category

Data Processing in Pandas

❖ columns

- Print out column names of our dataset
 - `movies_df.columns`

```
Index(['Rank', 'Genre', 'Description', 'Director', 'Actors', 'Year',  
      'Runtime (Minutes)', 'Rating', 'Votes', 'Revenue (Millions)',  
      'Metascore'],  
      dtype='object')
```

- We can use `columns` to rename column names

Data Processing in Pandas

❖ `rename()`

- Rename certain or all columns via a dict

```
movies_df.rename(columns={  
    'Runtime (Minutes)': 'Runtime',  
}, inplace=True)
```

```
movies_df.columns
```

```
Index(['Rank', 'Genre', 'Description', 'Director', 'Actors', 'Year',  
      'Runtime', 'Rating', 'Votes', 'Revenue (Millions)', 'Metascore'],  
      dtype='object')
```

- Task
 - Change Revenue (Millions) -> 'Revenue_millions'

Data Processing in Pandas

❖ Data frame manipulation

- Output the following table
 - Note: Title is index variable

| | | Genre | Rating |
|-------------------------|--------------------------|-------|--------|
| Title | | | |
| Guardians of the Galaxy | Action,Adventure,Sci-Fi | | 8.1 |
| Prometheus | Adventure,Mystery,Sci-Fi | | 7.0 |
| Split | Horror,Thriller | | 7.3 |
| Sing | Animation,Comedy,Family | | 7.2 |
| Suicide Squad | Action,Adventure,Fantasy | | 6.2 |

- Source code

```
subset = movies_df[['Genre', 'Rating']]  
  
subset.head()
```

Data Processing in Pandas

❖ Data frame manipulation

■ Output movies taken in 2012

- `movies_df[movies_df['Year'] == 2012].head(5)`

| Title | Rank | Genre | Description | Director | Actors | Year | Runtime | Rating | Votes | Revenue_millions | Metascore |
|----------------------------|------|--------------------------|---|-------------------|---|------|---------|--------|---------|------------------|-----------|
| Prometheus | 2 | Adventure,Mystery,Sci-Fi | Following clues to the origin of mankind, a te... | Ridley Scott | Noomi Rapace, Logan Marshall-Green, Michael Fa... | 2012 | 124 | 7.0 | 485820 | 126.46 | 65.0 |
| The Avengers | 77 | Action,Sci-Fi | Earth's mightiest heroes must come together an... | Joss Whedon | Robert Downey Jr., Chris Evans, Scarlett Johan... | 2012 | 143 | 8.1 | 1045588 | 623.28 | 69.0 |
| The Dark Knight Rises | 125 | Action,Thriller | Eight years after the Joker's reign of anarchy... | Christopher Nolan | Christian Bale, Tom Hardy, Anne Hathaway,Gary ... | 2012 | 164 | 8.5 | 1222645 | 448.13 | 78.0 |
| The Place Beyond the Pines | 136 | Crime,Drama,Thriller | A motorcycle stunt rider turns to robbing bank... | Derek Cianfrance | Ryan Gosling, Bradley Cooper, Eva Mendes,Craig... | 2012 | 140 | 7.3 | 200090 | 21.38 | 68.0 |
| Django Unchained | 145 | Drama,Western | With the help of a German bounty hunter , a fr... | Quentin Tarantino | Jamie Foxx, Christoph Waltz, Leonardo DiCaprio... | 2012 | 165 | 8.4 | 1039115 | 162.80 | 81.0 |

Data Processing in Pandas

❖ Data frame manipulation

- Output the movies that have a rating of 8.6
 - `movies_df[movies_df['Rating'] >= 8.6].head(5)`

| Title | Rank | Genre | Description | Director | Actors | Year | Runtime | Rating | Votes | Revenue_millions | Metascore |
|-----------------|------|-------------------------|---|-------------------|---|------|---------|--------|---------|------------------|-----------|
| Interstellar | 37 | Adventure,Drama,Sci-Fi | A team of explorers travel through a wormhole ... | Christopher Nolan | Matthew McConaughey, Anne Hathaway, Jessica Ch... | 2014 | 169 | 8.6 | 1047747 | 187.99 | 74.0 |
| The Dark Knight | 55 | Action,Crime,Drama | When the menace known as the Joker wreaks havo... | Christopher Nolan | Christian Bale, Heath Ledger, Aaron Eckhart,Mi... | 2008 | 152 | 9.0 | 1791916 | 533.32 | 82.0 |
| Inception | 81 | Action,Adventure,Sci-Fi | A thief, who steals corporate secrets through ... | Christopher Nolan | Leonardo DiCaprio, Joseph Gordon-Levitt, Ellen... | 2010 | 148 | 8.8 | 1583625 | 292.57 | 74.0 |
| Kimi no na wa | 97 | Animation,Drama,Fantasy | Two strangers find themselves linked in a biza... | Makoto Shinkai | Ryūnosuke Kamiki, Mone Kamishiraishi, Ryō Nari... | 2016 | 106 | 8.6 | 34110 | 4.68 | 79.0 |
| Dangal | 118 | Action,Biography,Drama | Former wrestler Mahavir Singh Phogat and his t... | Nitesh Tiwari | Aamir Khan, Sakshi Tanwar, Fatima Sana Shaikh,... | 2016 | 161 | 8.8 | 48969 | 11.15 | NaN |

Data Processing in Pandas

❖ Data frame manipulation

- Output movies, which was directed by Christopher Nolan OR Ridley Scott
- Hint we can use OR (|) operator
- `movies_df[(movies_df['Director'] == 'Christopher Nolan') | (movies_df['Director'] == 'Ridley Scott')].head(5)`

| Title | Rank | Genre | Description | Director | Actors | Year | Runtime | Rating | Votes | Revenue_millions | Metascore |
|-----------------|------|--------------------------|---|-------------------|---|------|---------|--------|---------|------------------|-----------|
| Prometheus | 2 | Adventure,Mystery,Sci-Fi | Following clues to the origin of mankind, a te... | Ridley Scott | Noomi Rapace, Logan Marshall-Green, Michael Fa... | 2012 | 124 | 7.0 | 485820 | 126.46 | 65.0 |
| Interstellar | 37 | Adventure,Drama,Sci-Fi | A team of explorers travel through a wormhole ... | Christopher Nolan | Matthew McConaughey, Anne Hathaway, Jessica Ch... | 2014 | 169 | 8.6 | 1047747 | 187.99 | 74.0 |
| The Dark Knight | 55 | Action,Crime,Drama | When the menace known as the Joker wreaks havo... | Christopher Nolan | Christian Bale, Heath Ledger, Aaron Eckhart,Mi... | 2008 | 152 | 9.0 | 1791916 | 533.32 | 82.0 |
| The Prestige | 65 | Drama,Mystery,Sci-Fi | Two stage magicians engage in competitive one-... | Christopher Nolan | Christian Bale, Hugh Jackman, Scarlett Johanss... | 2006 | 130 | 8.5 | 913152 | 53.08 | 66.0 |
| Inception | 81 | Action,Adventure,Sci-Fi | A thief, who steals corporate secrets through ... | Christopher Nolan | Leonardo DiCaprio, Joseph Gordon-Levitt, Ellen... | 2010 | 148 | 8.8 | 1583625 | 292.57 | 74.0 |

Data Processing in Pandas

❖ Data frame manipulation

- `sort_value()` function for sorting by values
 - “by” argument that indicates the column name of data frame to be sorted
- Example of `sort_value()`
 - `movies_df.sort_values(by='Year', ascending=False).head(5)`

| | Rank | Genre | Description | Director | Actors | Year | Runtime | Rating | Votes | Revenue_millions | Metascore |
|------------------------------------|------|-------------------------|---|-----------------------|--|------|---------|--------|-------|------------------|-----------|
| Title | | | | | | | | | | | |
| Nine Lives | 1000 | Comedy,Family,Fantasy | A stuffy businessman finds himself trapped ins... | Barry Sonnenfeld | Kevin Spacey, Jennifer Garner, Robbie Amell, Ch... | 2016 | 87 | 5.3 | 12435 | 19.64 | 11.0 |
| Free Fire | 162 | Action,Comedy,Crime | Set in Boston in 1978, a meeting in a deserted... | Ben Wheatley | Sharlto Copley, Brie Larson, Armie Hammer, Cil... | 2016 | 90 | 7.0 | 6946 | 1.80 | 63.0 |
| Tall Men | 648 | Fantasy,Horror,Thriller | A challenged man is stalked by tall phantoms i... | Jonathan Holbrook | Dan Crisafulli, Kay Whitney, Richard Garcia, P... | 2016 | 133 | 3.2 | 173 | NaN | 57.0 |
| The Huntsman: Winter's War | 235 | Action,Adventure,Drama | Eric and fellow warrior Sara, raised as member... | Cedric Nicolas-Troyan | Chris Hemsworth, Jessica Chastain, Charlize Th... | 2016 | 114 | 6.1 | 66766 | 47.95 | 35.0 |
| Popstar: Never Stop Never Stopping | 654 | Comedy,Music | When it becomes clear that his solo album is a... | Akiva Schaffer | Andy Samberg, Jorma Taccone, Akiva Schaffer, Sa... | 2016 | 87 | 6.7 | 30875 | 9.39 | 68.0 |

Data Processing in Pandas

❖ Dealing with missing values

- Missing or null values
 - Indicate the quality of your dataset
- Check out where you have missing values using `isnull()` function
 - `movies_df.isnull().tail(5)`

| | Rank | Genre | Description | Director | Actors | Year | Runtime | Rating | Votes | Revenue_millions | Metascore |
|------------------------|-------|-------|-------------|----------|--------|-------|---------|--------|-------|------------------|-----------|
| Title | | | | | | | | | | | |
| Secret in Their Eyes | False | False | False | False | False | False | False | False | False | True | False |
| Hostel: Part II | False | False | False | False | False | False | False | False | False | False | False |
| Step Up 2: The Streets | False | False | False | False | False | False | False | False | False | False | False |
| Search Party | False | False | False | False | False | False | False | False | False | True | False |
| Nine Lives | False | False | False | False | False | False | False | False | False | False | False |

Data Processing in Pandas

❖ Dealing with missing values

- To see the summary of missing values
 - `movies_df.isnull().sum()`

```
Rank          0
Genre          0
Description    0
Director       0
Actors         0
Year           0
Runtime        0
Rating         0
Votes          0
Revenue_millions 128
Metascore      64
dtype: int64
```

Data Processing in Pandas

❖ Removing null values

- Use `dropna()` function to remove rows with null values
 - `movies_df.dropna(inplace=True)`
 - `movies_df.isnull().sum()`

```
Rank          0
Genre          0
Description    0
Director       0
Actors         0
Year           0
Runtime        0
Rating         0
Votes          0
Revenue_millions  0
Metascore      0
dtype: int64
```

Data Processing in Pandas

❖ Removing null values

- You can also remove columns using `dropna()`
 - `movies_df.dropna(axis=1)`
 - 0, or 'index' : Drop rows which contain missing values.
 - 1, or 'columns' : Drop columns which contain missing value.
- Hint
 - Removing null data is only suggested if you have a small amount of missing data

Data Processing in Pandas

❖ Data Imputation

- `fillna()` function to fill the nulls
 - `movies_df['Revenue_millions'].fillna(movies_df['Revenue_millions'].mean(), inplace=True)`
 - `movies_df.isnull().sum()`

```
Rank          0
Genre          0
Description    0
Director       0
Actors         0
Year           0
Runtime        0
Rating         0
Votes          0
Revenue_millions 0
Metascore      64
dtype: int64
```

Homework for Lecture 4

❖ Given the following task:

1. Read `usedcars.csv`
2. Show summary using `head()`, `tale()`, `info()` and `describe()`
3. Change name of columns
4. Perform at least three data manipulation
5. Check if we have null values
6. Replace null values with `mean()` values

❖ Submission: source code and result screenshots



감사합니다!