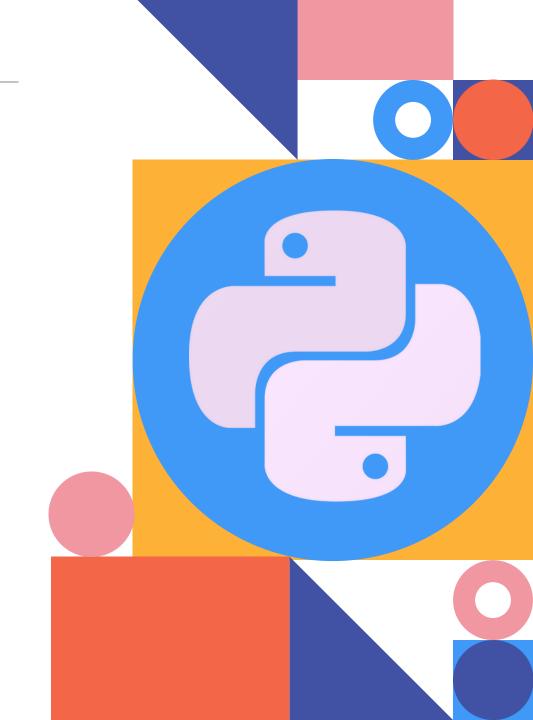# Data Analysis II

Yan-Fu Kuo

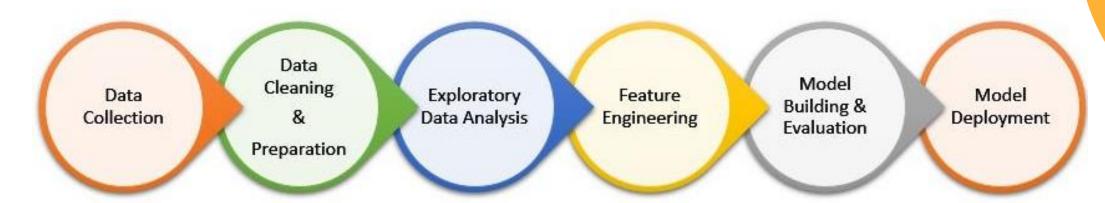Dept. of Biomechatronics Engineering
National Taiwan University

→

# Contents

# 01 Exploratory Data Analysis

# Exploratory Data Analysis (EDA)

An approach to analyzing data sets to summarize their main characteristics, often with visual methods

1. Understand your variables

2. Clean your dataset

3. Analyze relationships between variables

# COVID-19 in Taiwan

https://reurl.cc/QdjjaM

# COVID-19 in Taiwan

# 02 Pandas

# Pandas

A fast and efficient **DataFrame** object for data manipulation with integrated indexing

Tools for **reading and writing data** between in-memory data structures and different formats: CSV and text files, Microsoft Excel, SQL databases, and the fast HDF5 format

…

8

# Pandas Objects

## Series

One-dimensional array of indexed data, or a specialization of a Python dictionary

## Dataframe *

A generalization of a NumPy array, or an extended Python dictionary

## Index

An immutable array or as an ordered set

# Series

```
import pandas as pd

data = pd.Series([0.25, 0.5, 0.75, 1.0],
                 index=['a', 'b', 'c', 'd'])
data['b'] #0.5
```

```
data = pd.Series([0.25, 0.5, 0.75, 1.0],
                 index=[2, 5, 3, 7])
data[3]
```

| Index | Value |
|-------|-------|
| a | 0.25 |
| b | 0.50 |
| c | 0.75 |
| d | 1.00 |

What is the value of data[3] ?

# Series as Specialized Dictionary

```
population_dict = {'California': 38332521,
                   'Texas': 26448193,
                   'New York': 19651127,
                   'Florida': 19552860,
                   'Illinois': 12882135}

population = pd.Series(population_dict)


population_dict['California'] #38332521
population['California'] #38332521
```

```
population.sort_values(ascending=True)
```

| California | 38332521 |
|---|---|
| Texas | 16448193 |
| New York | 19651127 |
| Florida | 19552860 |
| Illinois | 12882135 |

| Illinois | 12882135 |
|---|---|
| Texas | 16448193 |
| Florida | 19552860 |
| New York | 19651127 |
| California | 38332521 |

# Creating a `DataFrame`

- From a single Series object

- From a list of dictionaries

- From clipboard

- From a two-dimensional NumPy array

- From a NumPy structured array

# DataFrame **from** Series

```
area_dict = {'California': 423967, 'Texas': 695662,
             'New York': 141297, 'Florida': 170312,
             'Illinois': 149995}
area = pd.Series(area_dict)
```

```
states = pd.DataFrame({'population': population,
                       'area': area})
```
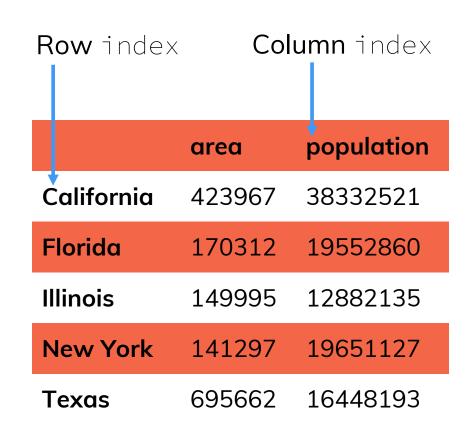
Series **2**    Series **1**

```
States['area']
```

Row index          Column index

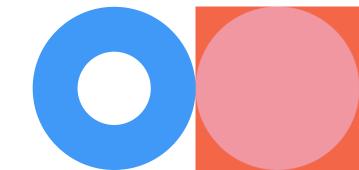| | area | population |
|---|---|---|
| **California** | 423967 | 38332521 |
| **Florida** | 170312 | 19552860 |
| **Illinois** | 149995 | 12882135 |
| **New York** | 141297 | 19651127 |
| **Texas** | 695662 | 16448193 |

# From a List of Dictionaries

```python
dict = {
    "col 1": [1, 2, 3],
    "col 2": [10, 20, 30],
    "col 3": list('xyz'),
    "col 4": ['a', 'b', 'c'],
    "col 5": pd.Series(range(3))
}

df = pd.DataFrame(dict)
```

|   | col 1 | col 2 | col 3 | col 4 | col 5 |
|---|-------|-------|-------|-------|-------|
| **0** | 1 | 10 | x | a | 0 |
| **1** | 2 | 20 | y | b | 1 |
| **2** | 3 | 30 | z | c | 2 |

# From Clipboard

Copy the text:

```
a    b      c      d
0    1      inf    1/1/00
2    71.38  N/A    5-Jan-13
4    54.59  nan    7/24/18
6    40.42  None   NaT
```

```
pd.read_clipboard()
```

# From CSV

Also check:
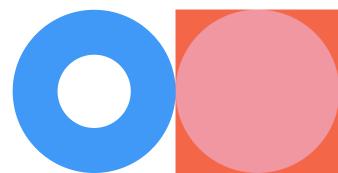
```
df = pd.read_csv('titanic.csv')
df.head(6)
```

```
df.tail(6)
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S |
| 1 | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| 2 | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S |
| 3 | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S |
| 4 | 5 | 0 | 3 | Allen, Mr. William Henry | male | 35.0 | 0 | 0 | 373450 | 8.0500 | NaN | S |
| 5 | 6 | 0 | 3 | Moran, Mr. James | male | NaN | 0 | 0 | 330877 | 8.4583 | NaN | Q |

# Titanic Machine Learning Competition

# Titanic Dataset

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | Braund, Mr. | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S |
| 1 | 2 | 1 | 1 | Cumings, Mrs. John... | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| 2 | 3 | 1 | 3 | Heikkinen, | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S |
| 3 | 4 | 1 | 1 | Futrelle, May Peel) | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S |
| 4 | 5 | 0 | 3 | Allen, Mr. | male | 35.0 | 0 | 0 | 373450 | 8.0500 | NaN | S |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 888 | 889 | 0 | 3 | Johnston, Miss. " | female | NaN | 1 | 2 | W./C. 6607 | 23.4500 | NaN | S |
| 889 | 890 | 1 | 1 | Behr, Mr. | male | 26.0 | 0 | 0 | 111369 | 30.0000 | C148 | C |
| 890 | 891 | 0 | 3 | Dooley, Mr. k | male | 32.0 | 0 | 0 | 370376 | 7.7500 | NaN | Q |

# df = read_csv()

*pandas.read_csv(filepath_or_buffer, sep=',', delimiter=None, header='infer',
names=None, index_col=None, usecols=None, squeeze=False, prefix=None,
mangle_dupe_cols=True, dtype=None, engine=None, converters=None,
true_values=None, false_values=None, skipinitialspace=False, skiprows=None,
skipfooter=0, nrows=None, na_values=None, keep_default_na=True,
na_filter=True, verbose=False, skip_blank_lines=True, parse_dates=False,
infer_datetime_format=False, keep_date_col=False, date_parser=None,
dayfirst=False, cache_dates=True, iterator=False, chunksize=None,
compression='infer', thousands=None, decimal='.', lineterminator=None,
quotechar='"', quoting=0, doublequote=True, escapechar=None, comment=None,
encoding=None, dialect=None, error_bad_lines=True, warn_bad_lines=True,
delim_whitespace=False, low_memory=True, memory_map=False,
float_precision=None)*

Memorize? Try `SHIFT` + `TAB` in jupyter

# df.info()

## Q: What are types of the objects?

df.info(memory_usage="deep")

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
 #   Column       Non-Null Count   Dtype
---  ------       --------------   -----
 0   PassengerId  891 non-null     int64
 1   Survived     891 non-null     int64
 2   Pclass       891 non-null     int64
 3   Name         891 non-null     object
 4   Sex          891 non-null     object
 5   Age          714 non-null     float64
 6   SibSp        891 non-null     int64
 7   Parch        891 non-null     int64
 8   Ticket       891 non-null     object
 9   Fare         891 non-null     float64
 10  Cabin        204 non-null     object
 11  Embarked     889 non-null     object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
 #   Column       Non-Null Count   Dtype
---  ------       --------------   -----
 0   PassengerId  891 non-null     int64
 1   Survived     891 non-null     int64
 2   Pclass       891 non-null     int64
 3   Name         891 non-null     object
 4   Sex          891 non-null     object
 5   Age          714 non-null     float64
 6   SibSp        891 non-null     int64
 7   Parch        891 non-null     int64
 8   Ticket       891 non-null     object
 9   Fare         891 non-null     float64
 10  Cabin        204 non-null     object
 11  Embarked     889 non-null     object
dtypes: float64(2), int64(5), object(5)
memory usage: 315.0 KB
```

# df.describe()

| | PassengerId | Survived | Pclass | Age | SibSp | Parch | Fare |
|---|---|---|---|---|---|---|---|
| count | 891.000000 | 891.000000 | 891.000000 | 714.000000 | 891.000000 | 891.000000 | 891.000000 |
| mean | 446.000000 | 0.383838 | 2.308642 | 29.699118 | 0.523008 | 0.381594 | 32.204208 |
| std | 257.353842 | 0.486592 | 0.836071 | 14.526497 | 1.102743 | 0.806057 | 49.693429 |
| min | 1.000000 | 0.000000 | 1.000000 | 0.420000 | 0.000000 | 0.000000 | 0.000000 |
| 25% | 223.500000 | 0.000000 | 2.000000 | 20.125000 | 0.000000 | 0.000000 | 7.910400 |
| 50% | 446.000000 | 0.000000 | 3.000000 | 28.000000 | 0.000000 | 0.000000 | 14.454200 |
| 75% | 668.500000 | 1.000000 | 3.000000 | 38.000000 | 1.000000 | 0.000000 | 31.000000 |
| max | 891.000000 | 1.000000 | 3.000000 | 80.000000 | 8.000000 | 6.000000 | 512.329200 |

# Dealing with Big Data

1PB Dataset              64GB ram       27" monitor

# Titanic Dataset (Review)

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | Braund, Mr. | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | **NaN** | S |
| **1** | 2 | 1 | 1 | Cumings, Mrs. John... | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| **2** | 3 | 1 | 3 | Heikkinen, | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S |
| **3** | 4 | 1 | 1 | Futrelle, May Peel) | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S |
| **4** | 5 | 0 | 3 | Allen, Mr. | male | 35.0 | 0 | 0 | 373450 | 8.0500 | **NaN** | S |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| **888** | 889 | 0 | 3 | Johnston, Miss. " | female | NaN | 1 | 2 | W./C. 6607 | 23.4500 | **NaN** | S |
| **889** | 890 | 1 | 1 | Behr, Mr. | male | 26.0 | 0 | 0 | 111369 | 30.0000 | C148 | C |
| **890** | 891 | 0 | 3 | Dooley, Mr. k | male | 32.0 | 0 | 0 | 370376 | 7.7500 | **NaN** | Q |

# Reducing Memory Usage – 1

Using the Pandas **Category** data type

```python
dtypes = {"Embarked": "category"}
cols = ['PassengerId', 'Name', 'Sex', 'Embarked']
df = pd.read_csv('titanic.csv',
                 dtype=dtypes, usecols=cols)

df.info(memory_usage="deep")
```

```
...
10     Cabin           204 non-null
11     Embarked        889 non-null
dtypes: float64(2), int64(5), object(5)
memory usage: 315.0 KB
```

```
...
10     Cabin              204 non-null
11     Embarked           889 non-null
dtypes: float64(2), int64(5), object(5)
memory usage: 135.0 KB
```

https://pandas.pydata.org/docs/reference/api/pandas.read_csv.html

# Reducing Memory Usage – 2

Batch processing

```python
reader = pd.read_csv('titanic.csv', chunksize=4)

for _, df_partial in zip(range(2), reader):
    display(df_partial)
```

**Q: How do you determine the batch size?**

https://docs.python.org/3/library/functions.html#zip

# Customizing DataFrame Display

```python
print("Default display.max_colwidth:",
      pd.get_option("display.max_colwidth"))
```

```
Default display.max_colwidth: 10
```

```python
pd.set_option("display.max_colwidth", 5)
```

# Customizing DataFrame `Display`

```python
pd.set_option("display.max_columns", 5)
pd.set_option("display.max_rows", 6)
```

| | PassengerId | Survived | ... | Cabin | Embarked |
|---|---|---|---|---|---|
| **0** | 1 | 0 | ... | **NaN** | S |
| **1** | 2 | 1 | ... | C85 | C |
| **2** | 3 | 1 | ... | NaN | S |
| **...** | ... | ... | ... | ... | ... |
| **888** | 889 | 0 | ... | **NaN** | S |
| **889** | 890 | 1 | ... | C148 | C |
| **890** | 891 | 0 | ... | **NaN** | Q |

# Exercise (5 mins)

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | Braund, Mr.... | male | 22.0 | 1 | 0 | A/5 21171 | 7.25 | NaN | S |
| 1 | 2 | 1 | 1 | Cumings, Mr... | female | 38.0 | 1 | 0 | PC 17599 | 71.28 | C85 | C |
| 2 | 3 | 1 | 3 | Heikkinen, ... | female | 26.0 | 0 | 0 | STON/O2. 31... | 7.92 | NaN | S |
| 3 | 4 | 1 | 1 | Futrelle, M... | female | 35.0 | 1 | 0 | 113803 | 53.10 | C123 | S |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 887 | 888 | 1 | 1 | Graham, Mis... | female | 19.0 | 0 | 0 | 112053 | 30.00 | B42 | S |
| 888 | 889 | 0 | 3 | Johnston, M... | female | NaN | 1 | 2 | W./C. 6607 | 23.45 | NaN | S |
| 889 | 890 | 1 | 1 | Behr, Mr. K... | male | 26.0 | 0 | 0 | 111369 | 30.00 | C148 | C |
| 890 | 891 | 0 | 3 | Dooley, Mr.... | male | 32.0 | 0 | 0 | 370376 | 7.75 | NaN | Q |

https://pandas.pydata.org/pandas-docs/stable/user_guide/options.html 🎤 🔍

# DataFrame Styling

```python
df_sample.style\
    .format('{:.1f}', subset='Fare')\
    .set_caption('Colorful Titanic Dataset')\
    .hide_index()\
    .bar('Age', vmin=0)\
    .highlight_max('Survived')\
    .background_gradient('Greens',
                         subset='Fare')\
    .highlight_null()
```

Colorful Titanic Dataset

| PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|
| 543 | 0 | 3 | female | 11 | 4 | 2 | 347082 | 31.3 | nan | S |
| 243 | 0 | 2 | male | 29 | 0 | 0 | W./C. 14263 | 10.5 | nan | S |
| 197 | 0 | 3 | male | nan | 0 | 0 | 368703 | 7.8 | nan | Q |
| 482 | 0 | 2 | male | nan | 0 | 0 | 239854 | 0.0 | nan | S |
| 683 | 0 | 3 | male | 20 | 0 | 0 | 6563 | 9.2 | nan | S |
| 693 | 1 | 3 | male | nan | 0 | 0 | 1601 | 56.5 | nan | S |
| 510 | 1 | 3 | male | 26 | 0 | 0 | 1601 | 56.5 | nan | S |
| 511 | 1 | 3 | male | 29 | 0 | 0 | 382651 | 7.8 | nan | Q |
| 568 | 0 | 3 | female | 29 | 0 | 4 | 349909 | 21.1 | nan | S |
| 227 | 1 | 2 | male | 19 | 0 | 0 | SW/PP 751 | 10.5 | nan | S |
| 722 | 0 | 3 | male | 17 | 1 | 0 | 350048 | 7.1 | nan | S |
| 645 | 1 | 3 | female | 1 | 2 | 1 | 2666 | 19.3 | nan | C |
| 681 | 0 | 3 | female | nan | 0 | 0 | 330935 | 8.1 | nan | Q |
| 114 | 0 | 3 | female | 20 | 1 | 0 | 4136 | 9.8 | nan | S |
| 630 | 0 | 3 | male | nan | 0 | 0 | 334912 | 7.7 | nan | Q |
| 411 | 0 | 3 | male | nan | 0 | 0 | 349222 | 7.9 | nan | S |
| 130 | 0 | 3 | male | 45 | 0 | 0 | 347061 | 7.0 | nan | S |

29

https://pandas.pydata.org/pandas-docs/stable/user_guide/style.html

# Exercise – Data Cleansing (5 mins)

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **542** | 543 | 0 | 3 | female | 11.0 | 4 | 2 | 347082 | 31.2750 | NaN | S |
| **242** | 243 | 0 | 2 | male | 29.0 | 0 | 0 | W./C. 14263 | 10.5000 | NaN | S |
| **196** | 197 | 0 | 3 | male | NaN | 0 | 0 | 368703 | 7.7500 | NaN | Q |
| **481** | 482 | 0 | 2 | male | NaN | 0 | 0 | 239854 | 0.0000 | NaN | S |
| **682** | 683 | 0 | 3 | male | 20.0 | 0 | 0 | 6563 | 9.2250 | NaN | S |

Try:

```
df.fillna(0)
```

```
df.fillna("Unknown")
```

```
df.Age.fillna(int(df.Age.mean()))
```

Q: Which on is better?
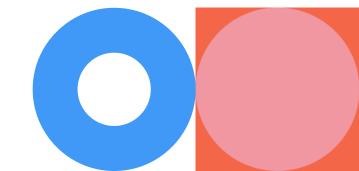
# Drop (Useless) Data

Drop column data

```
columns = ['Name', 'Ticket']
df.drop(columns, axis=1)
```

Drop row data

```
df.drop(2)
```

# DataFrame Indexing

`df.Pclass`   `df['Pclass']`   `df[0:4]`

| | |
|---|---|
| 0 | 3 |
| 1 | 1 |
| 2 | 3 |
| 3 | 1 |
| ... | ... |
| 887 | 1 |
| 888 | 3 |
| 889 | 1 |
| 890 | 3 |

| | PassengerId | Survived | ... | Embarked |
|---|---|---|---|---|
| 0 | 1 | 0 | ... | S |
| 1 | 2 | 1 | ... | C |
| 2 | 3 | 1 | ... | S |
| 3 | 4 | 1 | ... | S |

- Index refers to columns and slicing refers to rows
- <u>Cannot</u> do indexing on both columns and rows together!

# Indexer

Allow indexing DataFrame as if it is a simple NumPy array.

- `df.loc` is primarily label based, but may also be used with a boolean array.

- `df.iloc` is primarily integer position based (from 0 to length-1 of the axis), but may also be used with a boolean array.

# Data Selection with Indexer

```
df.loc[0:4, 'PassengerId': 'Age']
```

```
df.iloc[0:5, 0:5]
```

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S |
| **1** | 2 | 1 | 1 | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| **2** | 3 | 1 | 3 | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S |
| **3** | 4 | 1 | 1 | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S |
| **4** | 5 | 0 | 3 | male | 35.0 | 0 | 0 | 373450 | 8.0500 | NaN | S |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |

34

# Data Selection with Indexer

```
df.loc[df.Survived == 1, :]
```

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **1** | 2 | 1 | 1 | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| **2** | 3 | 1 | 3 | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S |
| **3** | 4 | 1 | 1 | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S |
| **8** | 9 | 1 | 3 | female | 27.0 | 0 | 2 | 347742 | 11.1333 | NaN | S |
| **9** | 10 | 1 | 2 | female | 14.0 | 1 | 0 | 237736 | 30.0708 | NaN | C |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| **875** | 876 | 1 | 3 | female | 15.0 | 0 | 0 | 2667 | 7.2250 | NaN | C |
| **879** | 880 | 1 | 1 | female | 56.0 | 0 | 1 | 11767 | 83.1583 | C50 | C |

# Data Selection – String and Num

## Masking

```python
female_and_age_under_20 = (df.Sex == 'female') & (df.Age < 20)
df[female_and_age_under_20]
```

## Query

```python
age = 70
df.query("Age > @age & Sex == 'male'")
```
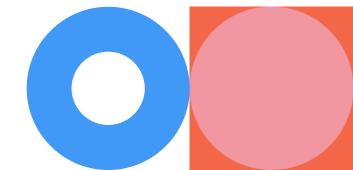
# Data Selection – String

```
df[df.Name.str.contains("Mr\.")]
```

```
tickets = ["SC/Paris 2123", "PC 17475"]
df[df.Ticket.isin(tickets)]
```

```
top_k = 3
top_tickets = df.Ticket.value_counts()[:top_k]
top_tickets.index
```

# Exercise

Show the passenger info of the most frequently purchased class of tickets

# GroupBy: Split, Apply, Combine

```
df.groupby("Pclass").Age.mean()
```

|   | Pclass | Age |
|---|--------|-----|
| **0** | 1 | 45 |
| **1** | 2 | 20 |
| **2** | 3 | 32 |
| **3** | 2 | 26 |
| **4** | 3 | 19 |
| **5** | 1 | 35 |

|   | Pclass | Age |
|---|--------|-----|
| **0** | 1 | 45 |
| **5** | 1 | 35 |

|   | Pclass | Age |
|---|--------|-----|
| **1** | 2 | 20 |
| **3** | 2 | 26 |

|   | Pclass | Age |
|---|--------|-----|
| **2** | 3 | 32 |
| **4** | 3 | 19 |

| Pclass | Age(mean) |
|--------|-----------|
| **1** | 38.23 |

| Pclass | Age(mean) |
|--------|-----------|
| **2** | 29.87 |

| Pclass | Age(mean) |
|--------|-----------|
| **3** | 25.14 |

| Pclass | Age(mean) |
|--------|-----------|
| 1 | 38.23 |
| 2 | 29.87 |
| 3 | 25.14 |

# Insights

## Does sex related to the survival rate in that accident?

```
df.groupby(['Sex','Survived']).Sex.count()
```

```
Sex     Survived
female  0            81
        1           233
male    0           468
        1           109
Name: Sex, dtype: int64
```



**female**     **male**

No Survived

Survived

# Exercise

Does ticket class (`pclass`) related to the survival rate in that accident?

thank you!