# Essential Python for Data Analyst

**Contents**

- review class (OOP)
- review loop + dictionary
- try except
- assert
- csv. json, context manager
- web scraping
- numpy pandas sqlite sklearn

```python
1 # super classs
2 class Person:
3     def __init__(self, name):
4         self.name = name
5     def hey (self):
6         print("hello world!")
7     def sayhi(self):
8         print(f"Hi my name is {self.name}!")
9
10 class Employee(Person):
11     def __init__(self, name, company):
12         super(). __init__(name)
13         self.company = company
14     def potato(self):
15         print("Potato is very good!")
16     def sayhi(self):
17         print(f"Hi my name is {self.name} and I work at {self.company}.")
18
```

```python
1 plume = Person("plume")
2 john = Employee("John", "Google")
3 plume.sayhi()
4 john.sayhi()
```

```
    Hi my name is plume!
    Hi my name is John and I work at Google.
```

```python
1 plume.hey()
```

```
    hello world!
```

```python
1 john.potato()
```

```
    Potato is very good!
```

```python
1 #try-except block
2 try:
3     1/0
4 except :
5     print("this is an error message")
```

```
    this is an error message
```

```python
1 try:
2     print("hello")
3 except ZeroDivisionError:
4     print("Cannot divide a number with zero")
5 except NameError:
6     print("Variable is not found")
7 except ValueError:
8     print("Please check your data type again!")
9 else:
10     print("Done!")
11 finally :
12     print("This is the last line")
```

```
    hello
    Done!
    This is the last line
```

```python
1 assert 1+1 ==5, "1+1 must be equal 2"
```

```
    --------------------------------------------------------------------------
    AssertionError                         Traceback (most recent call last)
    <ipython-input-31-044f29da9dba> in <cell line: 1>()
    ----> 1 assert 1+1 ==5, "1+1 must be equal 2"

    AssertionError: 1+1 must be equal 2
```

SEARCH STACK OVERFLOW

```python
1 #type hint
2 def check_this(a: int, b: int) -> int:
3     assert a+b > 10, "a+b should be more than 10"
4     return a+b
```

```python
1 check_this(10,4)
```

```
    14
```

```python
1 # lamda function(เหมาะกับ simple function)
```

```python
1 def greeting(name):
2     print("hello! "+ name)
```

```python
1 greeting = lambda name: print("hello! "+ name)
```

```python
1 greeting("pasit")
```

```
    hello! pasit
```

```python
1 # แบบมาตรฐาน
2 def odd_or_even(num):
3     if num% 2 == 0:
4         print("Even")
5     else:
6         print("Odd")
```

```python
1 # แบบสั้น (เหมาะกับ simple function)
2 odd_or_even = lambda num: "even" if num %2 == 0 else "odd"
```

```python
1 odd_or_even(4)
```

```
    'even'
```

```python
1 # review class
2 # try except
3 # assert 1+1==2, "error message"
4 # lambda function
```

```python
1 hello = lambda : print("hello!")
```

```python
1 hello()
```

```
    hello!
```

```python
1 # list of colored balls
2 balls = ["red", "red", "greed", "blue", "blue"]
```

```python
1 def summary(lst_items):
2     result = {}
3     for item in lst_items :
4         if item in result:
```

```
5              result[item] += 1
6         else:
7              result[item] = 1
8     return(result)
9
```

```
1 summary(balls)
```

```
    {'red': 2, 'greed': 1, 'blue': 2}
```

```
1 # comma separated values (CSV file)
2 import csv
3
4 try:
5     file = open("hotel.csv")
6     data = csv.reader(file)
7     file.close()
8 except FileNotFoundError:
9     print("File not found.")
10 else:
11     print("Load data successfully.")
```

```
    Load data successfully.
```

```
1 # context manager
2 import csv
3
4 try:
5     with open ("hotel.csv", "r") as file:
6         data = csv.reader(file)
7         for row in data:
8             print(row)
9 except:
10     print("Found some errors, please check again.")
```

```
    ['id', 'hotel ', 'location', 'price']
    ['1', 'Accor', 'Thailand', '50']
    ['2', 'Ibis', 'Italy', '60']
    ['3', 'London', 'UK', '35']
    ['4', 'Seoul', 'Korea', '30']
    ['5', 'le concord', 'USA', '25']
```

```
1 csv_data = []
2
3 with open ("hotel.csv", "r") as file:
4     data = csv.reader(file)
5     for row in data:
6         csv_data.append(row)
7
8 print(csv_data)
```

```
    [['id', 'hotel ', 'location', 'price'], ['1', 'Accor', 'Thailand', '50'], ['2', 'Ibis', 'Italy', '60'], ['3', 'London', 'UK', '35'], ['4
```

```
1 import pandas as pd
2 df= pd.read_csv("hotel.csv")
3
4 df["price_double"] = df["price"]*2
5
6 df.to_csv("hotel_modifiled.csv")
```

```
1 # write csv files
2 # "w" write mode
3 import csv
4
5 headers = ["id", "name", "age"]
6 data = [
7     [1, "pasit", 27],
8     [2, "jisoo", 26],
9     [3, "john", 32]
10 ]
11
12
13 with open("newfile.csv", "w") as file:
```

```
14     writer = csv.writer(file)
15     writer.writerow(headers)
16     writer.writerows(data)
```

```
1 !cat newfile.csv
```

```
id,name,age
1,pasit,27
2,jisoo,26
3,john,32
```

```
1 result=[]
2 with open("newfile.csv", "r") as file:
3     data = csv.reader(file)
4     for row in data:
5         result.append(row)
6 print(result)
```

```
[['id', 'name', 'age'], ['1', 'pasit', '27'], ['2', 'jisoo', '26'], ['3', 'john', '32']]
```

```
1 #regular expression
2
3 import re
4
5 #from re import search
6
7 text = "a tiger walks into bar"
8
9 result = re.search("duck", text)
10
11 if result:
12     print("Duck Found")
13 else:
14     print("Duck not found")
```

```
Duck not found
```

```
1 # list comprehension
2 numbers = [1, 2, 3, 4, 5]
```

```
1 result = [num*2 for num in numbers]
2 print(result)
```

```
[2, 4, 6, 8, 10]
```

```
1 oddNums = [num for num in numbers if num%2 == 1]
2 evenNums = [num for num in numbers if num%2 == 0]
3
4 print(oddNums, evenNums)
```

```
[1, 3, 5] [2, 4]
```

```
1 # json (JavaScript Object Notation)
2 import json
3
4 with open("pasit.json", "r") as file:
5     data = json.load(file)
```

```
1 print(data, type(data))
```

```
{'name': 'Pasit', 'age': 27, 'movie_lover': True, 'movies': ['batman', 'superman']} <class 'dict'>
```

```
1 data["name"]
```

```
'Pasit'
```

```
1 # create new .jason file
2 school = {
3     "name" : "bootcamp",
4     "batch": 7,
5     "duration": "4 months"
```

```
6 }
7
8 with open("newjsonfile.json", "w") as file:
9     json.dump(school, file)
```

```
1 !cat newjsonfile.json
```

```
{"name": "bootcamp", "batch": 7, "duration": "4 months"}
```

```
1 with open("newjsonfile.json", "r") as file:
2     data2 = json.load(file)
3
4 data2
```

```
{'name': 'bootcamp', 'batch': 7, 'duration': '4 months'}
```

## ▾ API

```
1 # API
2 # HW01 - find public API
3 from requests import get
4
5 nums = list(range(1, 6))
6 result = []
7
8 for i in nums:
9     url= f"https://swapi.dev/api/people/{i}"
10    response = get(url)
11    data = response.json()
12    row = [
13        data["name"],
14        data["height"],
15        data["mass"],
16        data["gender"]
17    ]
18    result.append(row)
19    print(result)
20
```

```
[['Luke Skywalker', '172', '77', 'male']]
[['Luke Skywalker', '172', '77', 'male'], ['C-3PO', '167', '75', 'n/a']]
[['Luke Skywalker', '172', '77', 'male'], ['C-3PO', '167', '75', 'n/a'], ['R2-D2', '96', '32', 'n/a']]
[['Luke Skywalker', '172', '77', 'male'], ['C-3PO', '167', '75', 'n/a'], ['R2-D2', '96', '32', 'n/a'], ['Darth Vader', '202', '136', 'ma
[['Luke Skywalker', '172', '77', 'male'], ['C-3PO', '167', '75', 'n/a'], ['R2-D2', '96', '32', 'n/a'], ['Darth Vader', '202', '136', 'ma
◂                                                                                                                                        ▸
```

```
1 import pandas as pd
2
3 df = pd.DataFrame(result, columns=["name", "height", "mass", "gender"])
4
5 df.to_csv("starwarapi.csv")
```

## ▾ Web scraping

```
1 #web scraping
2 #gazpacho
3
4 !pip install gazpacho
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Collecting gazpacho
  Downloading gazpacho-1.1.tar.gz (7.9 kB)
  Installing build dependencies ... done
  Getting requirements to build wheel ... done
  Preparing metadata (pyproject.toml) ... done
Building wheels for collected packages: gazpacho
  Building wheel for gazpacho (pyproject.toml) ... done
  Created wheel for gazpacho: filename=gazpacho-1.1-py3-none-any.whl size=7480 sha256=cae997e716e59e0a6c4f9eba3fccc6afd4f52f7df9f8bfe2b3
  Stored in directory: /root/.cache/pip/wheels/7b/82/22/0f73f5a5fa5abb29aac7a06941ff561e0bb415e558bb64d467
Successfully built gazpacho
Installing collected packages: gazpacho
Successfully installed gazpacho-1.1
```

```
1 from gazpacho import Soup
2 from requests import get
3
4 url = "https://www.imdb.com/search/title/?groups=top_100&sort=user_rating%2Cdesc"
5
6 response = get(url)
7
8 imdb = Soup(response.text)
```

```
1 imdb.find("h3")[9].strip()
```

    '10. The Lord of the Rings: The Fellowship of the Ring (2001)'

```
1 title_list = imdb.find("h3", {"class": "lister-item-header"})
```

```
1 titles = [title.strip() for title in title_list]
2 print(titles)
```

    ['1. The Shawshank Redemption (1994)', '2. The Godfather (1972)', '3. The Dark Knight (2008)', "4. Schindler's List (1993)", '5. The Lor

```
1 runtimes = imdb.find("span", {"class": "runtime"})
2 runtimes = [runtime.strip() for runtime in runtimes]
3 print(runtimes)
```

    ['142 min', '175 min', '152 min', '195 min', '201 min', '96 min', '202 min', '154 min', '148 min', '178 min', '139 min', '142 min', '179

```
1 import pandas as pd
2
3 df = pd.DataFrame({
4     "title" : titles,
5     "runtime": runtimes
6 })
7
8 df.head()
```

|   | title | runtime |
|---|---|---|
| 0 | 1. The Shawshank Redemption (1994) | 142 min |
| 1 | 2. The Godfather (1972) | 175 min |
| 2 | 3. The Dark Knight (2008) | 152 min |
| 3 | 4. Schindler's List (1993) | 195 min |
| 4 | 5. The Lord of the Rings: The Return of the Ki... | 201 min |

```
1 df.to_csv("imdb2.csv")
```

## Basic Data Tools

- numpy
- pandas
- sklearn

```
1 # numpy => numerical python
2 import numpy as np
3
4 num = [1, 2, 3, 4, 5]
5
6 num_arr = np.array(num)
7
8 print(num_arr*2, num_arr+2, num_arr-2, num_arr/5)
```

    [ 2  4  6  8 10] [3 4 5 6 7] [-1  0  1  2  3] [0.2 0.4 0.6 0.8 1. ]

```
1 # statistical basic calculation
2 np.sum(num_arr)
3 np.mean(num_arr)
4 np.std(num_arr)
5 np.min(num_arr)
6 np.max(num_arr)
7 len(num_arr)
```

    5

```
1 num_arr.mean()
```

    3.0

```
 1 # matrix 2D-array
 2 m1 = np.array([
 3     [1, 2],
 4     [3,4]
 5 ])
 6
 7 m2 = np.array([
 8     [5, 5],
 9     [4,2]
10 ])
```

```
1 print(m1*2, "\n", m2)
```

    [[2 4]
     [6 8]]
     [[5 5]
     [4 2]]

```
1 #intro to pandas
2 import pandas as pd
3 import numpy as np
```

```
1 hotel = pd.read_csv("hotel.csv")
```

```
1 hotel.shape
```

    (5, 4)

```
1 hotel.info()
```

    <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 5 entries, 0 to 4
    Data columns (total 4 columns):
     #   Column    Non-Null Count  Dtype
    ---  ------    --------------  -----
     0   id        5 non-null      int64
     1   hotel     5 non-null      object
     2   location  5 non-null      object
     3   price     5 non-null      int64
    dtypes: int64(2), object(2)
    memory usage: 288.0+ bytes

```
1 hotel.describe()
```

|       | id       | price    |
|-------|----------|----------|
| count | 5.000000 | 5.00000  |
| mean  | 3.000000 | 40.00000 |
| std   | 1.581139 | 14.57738 |
| min   | 1.000000 | 25.00000 |
| 25%   | 2.000000 | 30.00000 |
| 50%   | 3.000000 | 35.00000 |
| 75%   | 4.000000 | 50.00000 |
| max   | 5.000000 | 60.00000 |

```
1 hotel["price"]
```

```
0    50
1    60
2    35
3    30
4    25
Name: price, dtype: int64
```

```
1 hotel.price
```

```
0    50
1    60
2    35
3    30
4    25
Name: price, dtype: int64
```

```
1 hotel["vat"]= hotel["price"] *0.07
2
3 hotel
```

|   | id | hotel | location | price | vat |
|---|----|-------|----------|-------|-----|
| 0 | 1 | Accor | Thailand | 50 | 3.50 |
| 1 | 2 | Ibis | Italy | 60 | 4.20 |
| 2 | 3 | London | UK | 35 | 2.45 |
| 3 | 4 | Seoul | Korea | 30 | 2.10 |
| 4 | 5 | le concord | USA | 25 | 1.75 |

```
1 # query(filter rows)
2 value_hotel = hotel.query("price <= 40")
```

```
1 #drop column
2 new_hotel = hotel.drop(["id","vat"], axis=1)
```

```
1 hotel.query("price < 40 and location == 'UK' ")
```

|   | id | hotel | location | price | vat |
|---|----|-------|----------|-------|-----|
| 2 | 3 | London | UK | 35 | 2.45 |

```
1 hotel.query("price > 30 or location == 'UK' ")
```

|   | id | hotel | location | price | vat |
|---|----|-------|----------|-------|-----|
| 0 | 1 | Accor | Thailand | 50 | 3.50 |
| 1 | 2 | Ibis | Italy | 60 | 4.20 |
| 2 | 3 | London | UK | 35 | 2.45 |

```
1 #sklearn => machine learning most popular python
2 # template in sklearn to train_test_split
```

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.tree import DecisionTreeRegressor
3 from sklearn.ensemble import RandomForestRegressor
4 from sklearn.model_selection import train_test_split
5 import pandas as pd
6 import numpy as np
```

```
1 mtcars = pd.read_csv("https://gist.githubusercontent.com/seankross/a412dfbd88b3db70b74b/raw/5f23f993cd87c283ce766e7ac6b329ee7cc2e1d1/mtcar
```

```
1 mtcars.head()
```

| | model | mpg | cyl | disp | hp | drat | wt | qsec | vs | am | gear | carb |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Mazda RX4 | 21.0 | 6 | 160.0 | 110 | 3.90 | 2.620 | 16.46 | 0 | 1 | 4 | 4 |
| 1 | Mazda RX4 Wag | 21.0 | 6 | 160.0 | 110 | 3.90 | 2.875 | 17.02 | 0 | 1 | 4 | 4 |
| 2 | Datsun 710 | 22.8 | 4 | 108.0 | 93 | 3.85 | 2.320 | 18.61 | 1 | 1 | 4 | 1 |
| 3 | Hornet 4 Drive | 21.4 | 6 | 258.0 | 110 | 3.08 | 3.215 | 19.44 | 1 | 0 | 3 | 1 |

```
1 #prepare
2 X = mtcars[ ["hp", "wt", "am"] ]
3 y = mtcars["mpg"]
```

```
1 #split data
2
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42 )
```

```
1 #train model
2 #model = LinearRegression()
3 #model = DecisionTreeRegressor()
4 model = RandomForestRegressor()
5 model.fit(X_train, y_train)
6
7 #test model/ scoring
8 pred = model.predict(X_test)
9
10 # MAE mean absolute error
11 mae = np.mean(np.absolute((y_test - pred)))
12
13 #MSE
14 mse = np.mean((y_test - pred)**2)
15
16 print(mae, mse)
17 #evaluate model
18
```

```
2.0610000000000013 8.637949571428555
```

```
1
```