# 05 Dataset & Dataframe

## Telung Pan
## telung@mac.com

Datascientist

# 介紹新工具

- Python IDLE

- repl.it

- Jupyter Notebook 

  - conda  

  - miniconda

  - brew cask install miniconda

# Jupyter Notebook

- 網址 jupyter.org

- pip install - - upgrade pip
  pip install notebook
  pip install - - upgrade ipython jupyter

- conda install -c Conda-forge notebook

# return 語法 I

```
class DataShell:
    def __init__(self, x):
        return x
```

return int(self.value)

❖ 上面的語法我們以 return x 取代熟悉的 print()

❖ print() 直接輸出字串到 console，而 return 會離開所在的 function 或 method 然後將回傳值傳給呼叫他的人。

# return 語法 II

❖ 先建立一個 DataShell 類別

```
 1   # Create class: DataShell
 2 ▾ class DataShell:
 3       # Initialize class with self and dataList as arguments
 4 ▾   def __init__(self, dataList):
 5         # Set data as instance variable, and assign it the value of
    dataList
 6         self.data = dataList
 7     # Define method that returns data: show
 8 ▾   def show(self):
 9         return self.data
10     # Define method that prints average of data: avg
11 ▾   def avg(self):
12         # Declare avg and assign it the average of data
13         avg = sum(self.data)/float(len(self.data))
14         # Return avg
15         return avg
16 # Instantiate DataShell taking integer_list as argument: my_data_shell
17 my_data_shell = DataShell(integer_list)
```

❖ 呼叫 show 和 avg 方法，列印出來。

```
20   print(my_data_shell.show())
21   print(my_data_shell.avg())
```

# return 語法 III: 擴充更加強大的 DataShell 類別

❖ 匯入 (import) bumpy 和 pandas 套件 (packages)

**NumPy.org**

## NumPy

NumPy is the fundamental package for scientific computing with Python. It contains among other things:

- a powerful N-dimensional array object
- sophisticated (broadcasting) functions
- tools for integrating C/C++ and Fortran code
- useful linear algebra, Fourier transform, and random number capabilities

## pandas

**pandas** is a fast, powerful, flexible and easy to use open source data analysis and manipulation tool, built on top of the Python programming language.

# Panda 官方網站文件

https://pandas.pydata.org/pandas-docs/
stable/reference/api/pandas.read_table.html

```python
1  # Load numpy as np and pandas as pd
2  import numpy as np
3  import pandas as pd
4
5  # Create class: DataShell
6  class DataShell:
7
8      # Initialize class with self and inputFile as arguments
9      def __init__(self, inputFile):
10         self.file = inputFile
11
12     # Define generate_csv method, with self argument
13     def generate_csv(self):
14         self.data_as_csv = pd.read_csv(self.file)
15         return self.data_as_csv
```

```python
17  # Instantiate DataShell with us_life_expectancy as input argument
18  data_shell = DataShell(us_life_expectancy)
19
20  # Call data_shell's generate_csv method, assign it to df
21  df = data_shell.generate_csv()
22
23  # Print df
24  print(df)
```

❖ **Source code**
```
import numpy as np
import pandas as pd

class DataShell:
 def __init__(self, inputFile):
     self.file = inputFile
   def generate_csv(self):
     self.data_as_csv = pd.read_csv(self.file)
     return self.data_as_csv


us_life_expectancy = '/Users/telung/Documents/us_life_expectancy.csv'
data_shell = DataShell(us_life_expectancy)
df = data_shell.generate_csv()

print(df)
```

# 修改欄位名稱的程式

```python
1  # Create class DataShell
2  class DataShell:
3
4      # Define initialization method
5      def __init__(self, filepath):
6          self.filepath = filepath
7          self.data_as_csv = pd.read_csv(filepath)
8
9      # Define method rename_column, with arguments self, column_name, and
   new_column_name
10     def rename_column(self, column_name, new_column_name):
11         self.data_as_csv.columns = self.data_as_csv.columns.str.replace
   (column_name, new_column_name)
12
13 # Instantiate DataShell as us_data_shell with argument us_life_expectancy
14 us_data_shell = DataShell(us_life_expectancy)
15
```

```python
16  # Print the datatype of your object's data_as_csv attribute
17  print(us_data_shell.data_as_csv.dtypes)
18
19  # Rename your objects column 'code' to 'country_code'
20  us_data_shell.rename_column('code', 'country_code')
21
22  # Again, print the datatype of your object's data_as_csv attribute
23  print(us_data_shell.data_as_csv.dtypes)
```
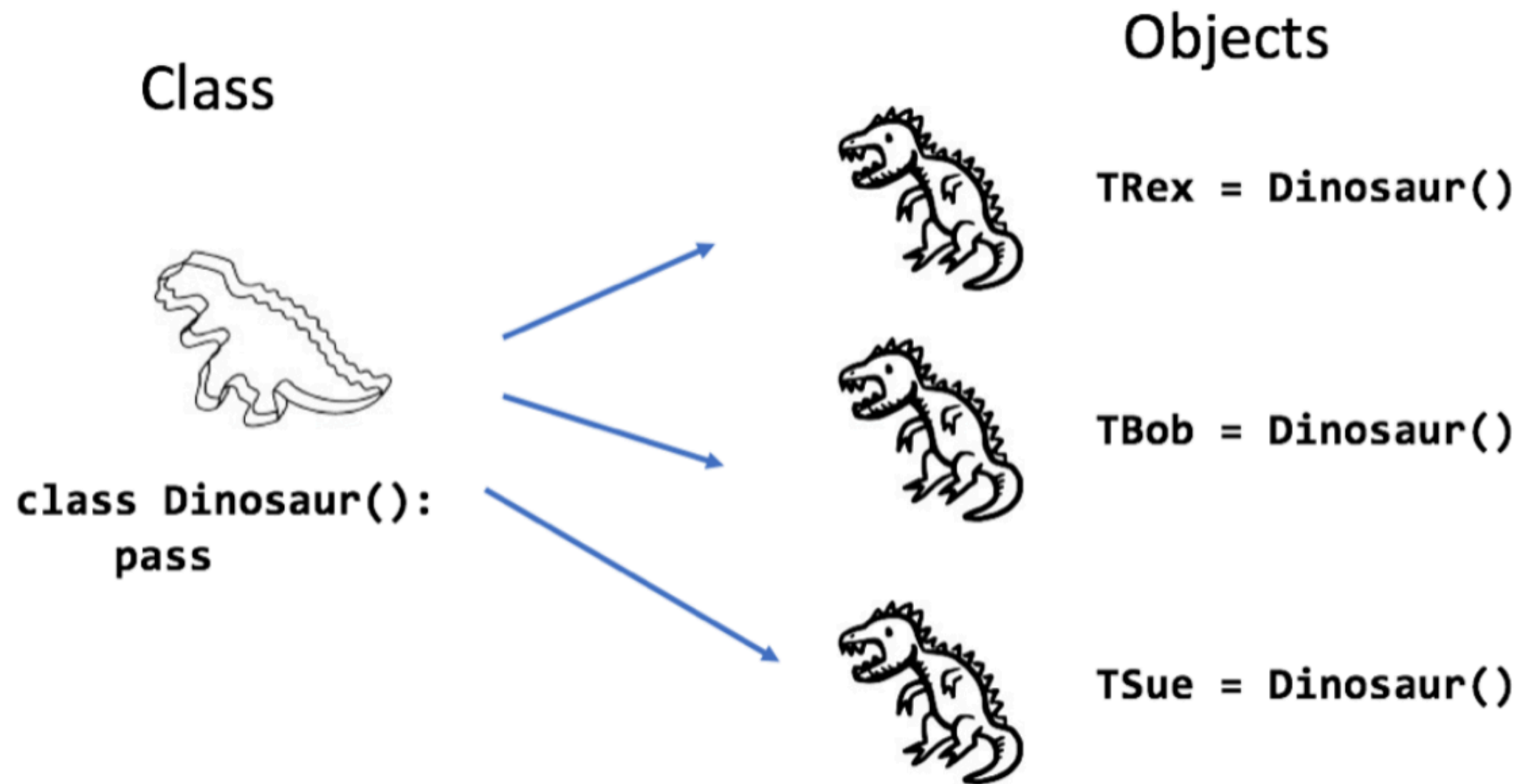
# 加入回傳自身的資訊

```python
class DataShell:
    def __init__(self, filepath):
        self.filepath = filepath
        self.data_as_csv = pd.read_csv(filepath)

    # 定義第一個方法 rename_column, with arguments self, column_name, and new_column_name
    def rename_column(self, column_name, new_column_name):
        self.data_as_csv.columns = self.data_as_csv.columns.str.replace(column_name, new_column_name)

    # 定義第二個 get_stats method, with argument self
    def get_stats(self):
        # Return a description data_as_csv
        return self.data_as_csv.describe()

us_data_shell = DataShell(us_life_expectancy)

# Print the output of your objects get_stats method
print(us_data_shell.get_stats())
```
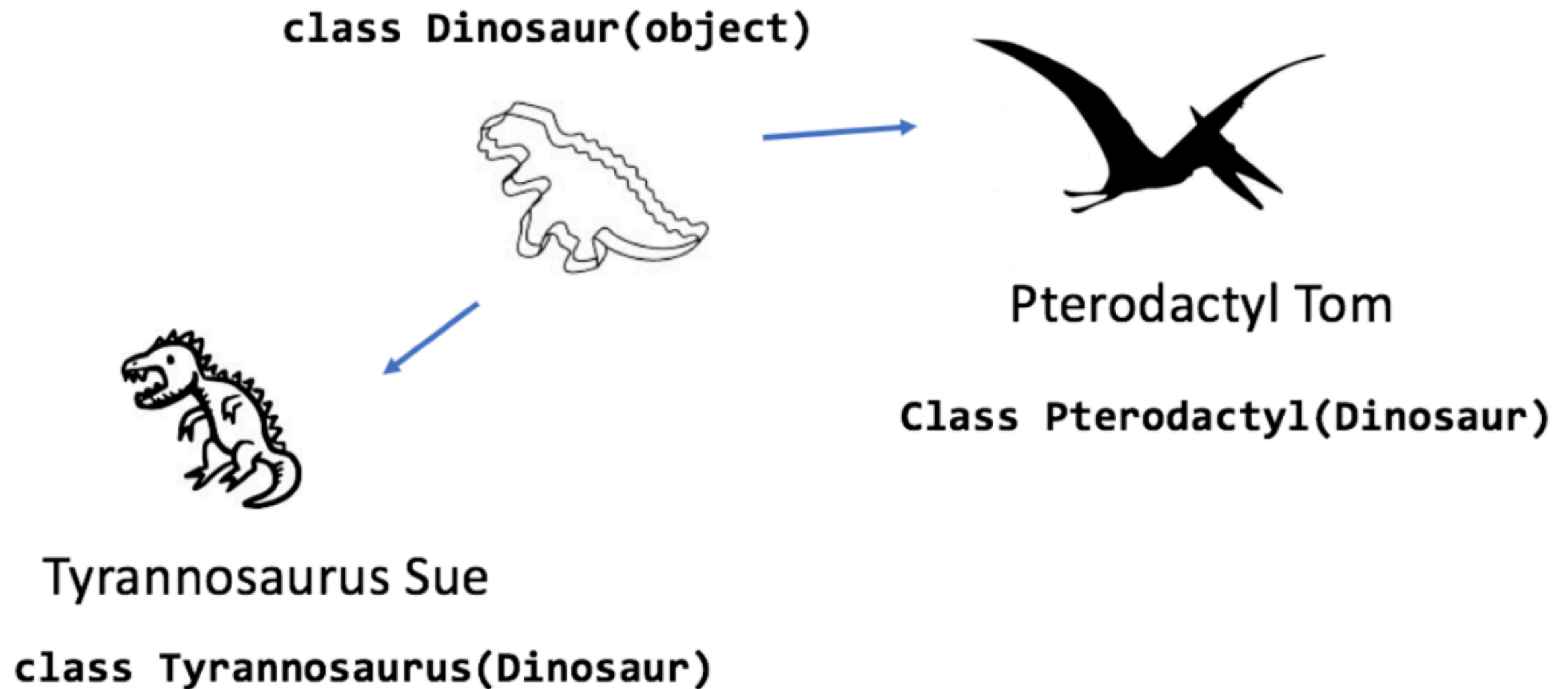
# 有關 Python 物件導向程式設計的一些建議指引

# 參考其他人的寫法

- ❖ GitHubCode

- ❖ Python 官方程式文件

- ❖ 參考 codebase:

  - ❖ Pandas  as Spark

  - ❖ PEP Style

# 繼承 (Inheritance)

# Is a and Has a Relationship

`class Dinosaur(object)`

Pterodactyl Tom

`Class Pterodactyl(Dinosaur)`

Tyrannosaurus Sue

`class Tyrannosaurus(Dinosaur)`

# Inheritance, simple one

```python
1   class Animal:
2     def __init__(self, name):
3       self.name = name
4
5   # Create a class Mammal, which inherits from Animal
6   class Mammal(Animal):
7     def __init__(self, name, animal_type):
8       self.animal_type = animal_type
9
10  # Create a class Reptile, which also inherits from Animal
11  class Reptile(Animal):
12    def __init__(self, name, animal_type):
13      self.animal_type = animal_type
14
15  # Instantiate a mammal with name 'Daisy' and animal_type 'dog': daisy
16  daisy = Mammal('Daisy', 'dog')
17
18  # Instantiate a reptile with name 'Stella' and animal_type 'alligator': stella
19  stella = Reptile('Stella', 'alligator')
20
21  # Print both daisy and stella
22  print(daisy)
23  print(stella)
```

# Inheritance, more complete

```python
1    # Create a class Vertebrate
2    class Vertebrate:
3        spinal_cord = True
4        def __init__(self, name):
5            self.name = name
6
7    # Create a class Mammal, which inherits from Vertebrate
8    class Mammal(Vertebrate):
9        def __init__(self, name, animal_type):
10            self.animal_type = animal_type
11            self.temperature_regulation = True
12
13    # Create a class Reptile, which also inherits from Vertebrate
14    class Reptile(Vertebrate):
15        def __init__(self, name, animal_type):
16            self.animal_type = animal_type
17            self.temperature_regulation = False
18
19    # Instantiate a mammal with name 'Daisy' and animal_type 'dog':
     daisy
20    daisy = Mammal('Daisy', 'dog')
```

```python
22  # Instantiate a reptile with name 'Stella' and animal_type
    'alligator': stella
23  stella = Reptile('Stella', 'alligator')
24
25  # Print stella's attributes spinal_cord and temperature_regulation
26  print("Stella Spinal cord: " + str(stella.spinal_cord))
27  print("Stella temperature regulation: " + str
    (stella.temperature_regulation))
28
29  # Print daisy's attributes spinal_cord and temperature_regulation
30  print("Daisy Spinal cord: " + str(daisy.spinal_cord))
31  print("Daisy temperature regulation: " + str
    (daisy.temperature_regulation))
```

# 實際案例練習：

```python
import pandas as pd

class DataShell:
    family = 'DataShell'
    def __init__(self, name, filepath):
        self.name = name
        self.filepath = filepath

class CsvDataShell(DataShell):
    def __init__(self, name, filepath):
        self.data = pd.read_csv(filepath)
        self.stats = self.data.describe()

class TsvDataShell(DataShell):
    # Initialization method with arguments self, name, filepath
    def __init__(self, name, filepath):
        # Instance variable data
        self.data = pd.read_table(filepath)
        # Instance variable stats
        self.stats = self.data.describe()
```

```python
23    # Instantiate CsvDataShell as us_data_shell, print
      us_data_shell.stats
24    us_data_shell = CsvDataShell("US",
      'us_life_expenctancy2.csv')
25    print(us_data_shell.stats)
26    print('-----------------------------\n')
27
28    # Instantiate TsvDataShell as france_data_shell, print
      france_data_shell.stats
29    france_data_shell = TsvDataShell("France",
      'france_life_expectancy3.csv')
30    print(france_data_shell.stats)
```