# Pandas 資料分析 (3)

張家瑋 副教授 國立臺中科技大學資訊工程系

#### 建立 DataFrame

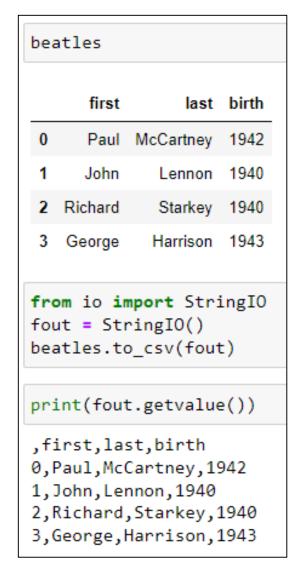
```
import pandas as pd
import numpy as np
pd.set option('max columns', 4, 'max rows', 10, 'max colwidth', 12)
fname = ['Paul', 'John', 'Richard', 'George']
lname = ['McCartney', 'Lennon', 'Starkey', 'Harrison']
birth = [1942, 1940, 1940, 1943]
people = {'first': fname, 'last': lname, 'birth': birth}
beatles = pd.DataFrame(people)
beatles
     first
               last birth
0
     Paul McCartney 1942
            Lennon 1940
     John
2 Richard
            Starkey 1940
 3 George
            Harrison 1943
```

```
beatles.index
RangeIndex(start=0, stop=4, step=1)
beatles.columns
Index(['first', 'last', 'birth'], dtype='object')
pd.DataFrame(people, index=['a', 'b', 'c', 'd'])
               last birth
      first
     Paul McCartney 1942
     John
             Lennon 1940
             Starkey 1940
c Richard
            Harrison 1943
d George
```

#### 建立 DataFrame

```
pd.DataFrame([{"first":"Paul","last":"McCartney", "birth":1942},
               {"first":"John","last":"Lennon", "birth":1940},
               {"first": "Richard", "last": "Starkey", "birth": 1940},
               {"first": "George", "last": "Harrison", "birth": 1943}])
                last birth
      first
     Paul McCartney 1942
      John
             Lennon 1940
2 Richard
             Starkey 1940
 3 George
            Harrison 1943
pd.DataFrame([{"first":"Paul","last":"McCartney", "birth":1942},
               {"first": "John", "last": "Lennon", "birth": 1940},
               {"first": "Richard", "last": "Starkey", "birth": 1940},
               {"first": "George", "last": "Harrison", "birth": 1943}],
               columns=['last', 'first', 'birth'])
               first birth
         last
               Paul 1942
0 McCartney
               John 1940
      Lennon
      Starkey Richard 1940
     Harrison George 1943
```

#### CSV 建立 DataFrame



```
fout.seek(0)
pd.read csv(fout)
   Unnamed: 0
                  first
                            last birth
                  Paul McCartney 1942
0
1
                 John
                          Lennon
                                1940
            2 Richard
                         Starkey 1940
2
 3
            3 George
                         Harrison 1943
fout.seek(0)
pd.read_csv(fout, index_col=0)
                 last birth
      first
      Paul McCartney 1942
      John
              Lennon 1940
2 Richard
              Starkey 1940
 3 George
             Harrison 1943
```

```
fout = StringIO()
beatles.to_csv(fout, index=False)
print(fout.getvalue())

first,last,birth
Paul,McCartney,1942
John,Lennon,1940
Richard,Starkey,1940
George,Harrison,1943
```

由於 csv 的索引無法直接當成 dataframe 的索引,因此使用 seek(0) 建立新索引。

原本的 csv 索引就變成了 Unnamed:0。

在 read\_csv() 中,直接指定 index\_col = 0,即特定(第一個)欄位作為 dataframe 索引。

```
diamonds = pd.read csv('data/diamonds.csv', nrows=1000)
diamonds
               cut ... y z
     carat
  0 0.23
              Ideal ... 3.98 2.43
      0.21 Premium ... 3.84 2.31
              Good ... 4.07 2.31
  2 0.23
      0.29 Premium ... 4.23 2.63
             Good ... 4.35 2.75
              Ideal ... 5.34 3.26
      0.54
              Ideal ... 5.74 3.57
      0.72
             Good ... 5.89 3.48
      0.74 Premium ... 5.77 3.58
     1.12 Premium ... 6.61 4.03
1000 rows x 10 columns
```

```
diamonds.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 10 columns):
    Column
            Non-Null Count Dtype
    carat
            1000 non-null float64
    cut
            1000 non-null object
    color
            1000 non-null object
    clarity 1000 non-null object
    depth
            1000 non-null float64
    table
            1000 non-null float64
    price
            1000 non-null int64
            1000 non-null float64
            1000 non-null float64
            1000 non-null float64
dtypes: float64(6), int64(1), object(3)
memory usage: 78.2+ KB
```

```
diamonds.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 10 columns):
    Column Non-Null Count Dtype
    carat 1000 non-null float64
    cut 1000 non-null object
    color
            1000 non-null object
            1000 non-null object
    clarity
    depth
            1000 non-null float64
    table
           1000 non-null float64
    price 1000 non-null int64
            1000 non-null float64
    Х
            1000 non-null float64
            1000 non-null float64
dtypes: float64(6), int64(1), object(3)
memory usage: 78.2+ KB
```

```
diamonds2 = pd.read csv('data/diamonds.csv', nrows=1000,
                      dtype={'carat': np.float32, 'depth': np.float32,
                            'table': np.float32, 'x': np.float32,
                            'y': np.float32, 'z': np.float32,
                            'price': np.int16})
diamonds2.info()
                           #改變成低精度型別來降低記憶體使用量
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 10 columns):
    Column Non-Null Count Dtype
    carat 1000 non-null float32
    cut 1000 non-null
                           object
    color 1000 non-null
                           object
    clarity 1000 non-null
                           object
    depth
           1000 non-null
                           float32
    table
                           float32
           1000 non-null
    price 1000 non-null
                           int16
            1000 non-null float32
            1000 non-null float32
            1000 non-null float32
dtypes: float32(6), int16(1), object(3)
memory usage: 49.0+ KB
```

	carat	depth	 у	Z
count	1000.000000	1000.000000	 1000.000000	1000.000000
mean	0.689280	61.722800	 5.599180	3.457530
std	0.195291	1.758879	 0.611974	0.389819
min	0.200000	53.000000	 3.750000	2.270000
25%	0.700000	60.900000	 5.630000	3.450000
50%	0.710000	61.800000	 5.760000	3.550000
75%	0.790000	62.600000	 5.910000	3.640000
max	1.270000	69.500000	 7.050000	4.330000

	carat	depth	у	Z
count	1000.000000	1000.000000	1000.000000	1000.000000
mean	0.689281	61.722824	5.599180	3.457533
std	0.195291	1.758878	0.611972	0.389819
min	0.200000	53.000000	3.750000	2.270000
25%	0.700000	60.900002	5.630000	3.450000
50%	0.710000	61.799999	5.760000	3.550000
75%	0.790000	62.599998	5.910000	3.640000
max	1.270000	69.500000	7.050000	4.330000

#改成低精度後,幾乎沒有差別,只有犧牲一點點精度,但省下了38%的記憶體。

記閱針億對 發 現 值都很單純  $\bigcirc$ 00 cla 等 category obje ry 存,欄 位 做 節

```
diamonds2.cut.value_counts()
Ideal
             333
Premium
             290
Very Good
             226
Good
              89
Fair
              62
Name: cut, dtype: int64
diamonds2.color.value_counts()
     240
     226
     139
     129
     125
      95
      46
Name: color, dtype: int64
diamonds2.clarity.value counts()
        306
SI1
VS2
        218
VS1
        159
SI2
        154
VVS2
         62
VVS1
         58
11
         29
         14
```

Name: clarity, dtype: int64

```
diamonds3 = pd.read csv('data/diamonds.csv', nrows=1000,
                       dtype={'carat': np.float32, 'depth': np.float32,
                              'table': np.float32, 'x': np.float32,
                              'v': np.float32, 'z': np.float32,
                              'price': np.int16,
                              'cut': 'category', 'color': 'category',
                              'clarity': 'category'})
diamonds3.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 10 columns):
             Non-Null Count Dtype
    Column
             1000 non-null float32
    carat
             1000 non-null category
    cut
             1000 non-null category
    color
    clarity 1000 non-null category
    depth
            1000 non-null float32
    table 1000 non-null float32
    price
             1000 non-null int16
    Х
             1000 non-null float32
             1000 non-null float32
    У
             1000 non-null float32
dtypes: category(3), float32(6), int16(1)
memory usage: 29.4 KB
```

#再將 object 改成 category 後,僅佔原本 78kb 的 37%。

```
cols = ['carat', 'cut', 'color', 'clarity', 'depth', 'table', 'price']
diamonds4 = pd.read csv('data/diamonds.csv', nrows=1000,
                      dtype={'carat': np.float32, 'depth': np.float32,
                             'table': np.float32, 'price': np.int16,
                             'cut': 'category', 'color': 'category',
                             'clarity': 'category'},
                      usecols=cols)
                                    #排除 x, y, Z 欄位。
diamonds4.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 7 columns):
    Column Non-Null Count Dtype
   carat 1000 non-null float32
  cut 1000 non-null category
   color 1000 non-null category
   clarity 1000 non-null category
   depth 1000 non-null float32
   table 1000 non-null float32
    price 1000 non-null int16
dtypes: category(3), float32(3), int16(1)
memory usage: 17.7 KB
```

```
cols = ['carat', 'cut', 'color', 'clarity', 'depth', 'table', 'price']
diamonds iter = pd.read csv('data/diamonds.csv', nrows=1000,
                            dtype={'carat': np.float32, 'depth': np.float32,
                                   'table': np.float32, 'price': np.int16,
                                   'cut': 'category', 'color': 'category',
                                   'clarity': 'category'},
                            usecols=cols.
                            chunksize=200)
def process(df):
    return f'processed {df.size} items'
for chunk in diamonds iter:
    print(process(chunk))
processed 1400 items
```

#### 查詢值域及使用的記憶體大小

```
np.iinfo(np.int8)
iinfo(min=-128, max=127, dtype=int8)

diamonds4['price'].min()
326

diamonds4['price'].max()
2898

np.finfo(np.float16)
finfo(resolution=0.001, min=-6.55040e+04, max=6.55040e+04, dtype=float16)
```

```
diamonds.price.memory_usage()

8128

diamonds.price.memory_usage(index=False)

8000

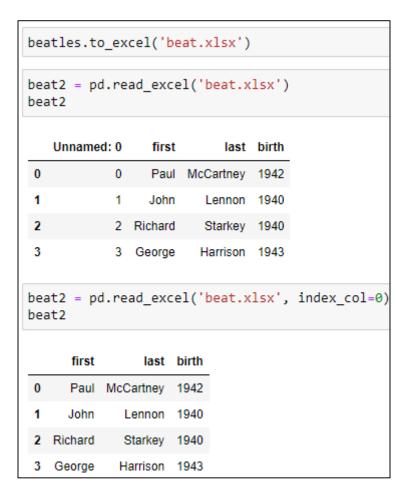
diamonds.cut.memory_usage(deep=True)

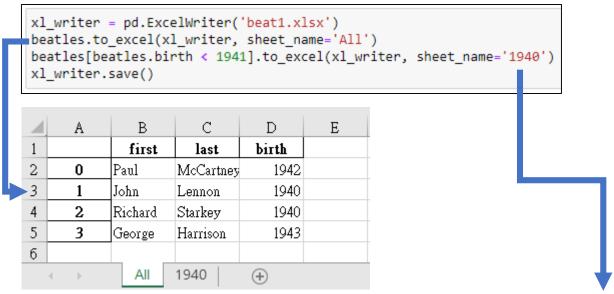
63461
```

#### 讀寫 EXCEL

```
!pip install xlwt
Collecting xlwt
 Downloading xlwt-1.3.0-py2.py3-none-any.whl (99 kB)
Installing collected packages: xlwt
Successfully installed xlwt-1.3.0
WARNING: You are using pip version 21.1.1; however, version 22.2.2 is available.
You should consider upgrading via the 'c:\users\test\appdata\local\programs\python\python38\python.exe -m pip install --upgrade
pip' command.
!pip install openpyxl
Collecting openpyxl
 Downloading openpyxl-3.0.10-py2.py3-none-any.whl (242 kB)
Collecting et-xmlfile
 Using cached et xmlfile-1.1.0-pv3-none-anv.whl (4.7 kB)
Installing collected packages: et-xmlfile, openpyxl
Successfully installed et-xmlfile-1.1.0 openpyxl-3.0.10
WARNING: You are using pip version 21.1.1; however, version 22.2.2 is available.
You should consider upgrading via the 'c:\users\test\appdata\local\programs\python\python38\python.exe -m pip install --upgrade
pip' command.
!pip install xlrd
Collecting xlrd
 Using cached xlrd-2.0.1-py2.py3-none-any.whl (96 kB)
Installing collected packages: xlrd
WARNING: You are using pip version 21.1.1; however, version 22.2.2 is available.
You should consider upgrading via the 'c:\users\test\appdata\local\programs\python\python38\python.exe -m pip install --upgrade
pip' command.
Successfully installed xlrd-2.0.1
```

#### 讀寫 EXCEL





$\Delta$	A	В	С	D	Е
1		first	last	birth	
2	1	John	Lennon	1940	
3	2	Richard	Starkey	1940	
4					
5					
б					
	<b>←</b> →	All	1940	<b>(+)</b>	

## 讀取 zip 中的 csv

#### #zip 中,只有一個 csv 的情況

autos = pd.read csv('data/vehicles.csv.zip') autos C:\Users\Admin\anaconda3\lib\site-packages\IPy have mixed types. Specify dtype option on impor has raised = await self.run ast nodes(code as barrels08 barrelsA08 ... phevHwy phevComb 0.0 ... 0 15.695714 0 0.0 ... 1 29.964545 2 12.207778 0.0 ... 0 0 3 29.964545 0.0 ... 0 4 17.347895 0.0 ... 0 39096 14.982273 0.0 ... 0 0.0 ... 39097 14.330870 0 0 0.0 ... 39098 15.695714 0 0 0.0 ... 39099 15.695714 0 39100 18.311667 0.0 ... 39101 rows x 83 columns

import zipfile # with zipfile.ZipFile('data/kaggle-survey-2018.zip') as z: print('\n'.join(z.namelist())) kag = pd.read csv(z.open('multipleChoiceResponses.csv')) kag questions = kag.iloc[0] survey = kag.iloc[1:] multipleChoiceResponses.csv p freeFormResponses.csv SurveySchema.csv **#** survey.head(2).T 有 1 2 Time from Start to Finish (seconds) 710 434 客 Female Q1 Male 個 Q1\_OTHER\_TEXT -1 -1 45-49 Q2 30-34 CQ3 United S... Indonesia 5 Q50\_Part\_5 NaN NaN Q50 Part 6 NaN NaN Q50\_Part\_7 NaN NaN 情 Q50\_Part\_8 NaN NaN 況 Q50\_OTHER\_TEXT -1 -1 395 rows × 2 columns

#### 存取資料庫

```
Iname birthyear
    fname
 id
     Paul McCartney
                         1942
     John
                         1940
             Lennon
sql = '''SELECT fname, birthyear from Band'
fnames = pd.read sql(sql, con)
fnames
   fname birthyear
    Paul
              1942
    John
              1940
```

```
import sqlalchemy as sa
engine = sa.create engine('sqlite:///data/beat.db', echo=True)
sa connection = engine.connect()
beat = pd.read sql('Band', sa connection, index col='id')
2021-11-10 17:44:47,791 INFO sqlalchemy.engine.base.Engine SELECT CAST('test plain returns' AS VARCHAR(60)) AS anon 1
2021-11-10 17:44:47,793 INFO sqlalchemv.engine.base.Engine ()
2021-11-10 17:44:47,795 INFO sqlalchemy.engine.base.Engine SELECT CAST('test unicode returns' AS VARCHAR(60)) AS anon 1
2021-11-10 17:44:47,797 INFO sqlalchemy.engine.base.Engine ()
2021-11-10 17:44:47,799 INFO sqlalchemy.engine.base.Engine PRAGMA main.table info("Band")
2021-11-10 17:44:47,801 INFO sqlalchemy.engine.base.Engine ()
2021-11-10 17:44:47,804 INFO sqlalchemy.engine.base.Engine SELECT name FROM sqlite master WHERE type='table' ORDER BY name
2021-11-10 17:44:47,805 INFO sqlalchemy.engine.base.Engine ()
2021-11-10 17:44:47,890 INFO sqlalchemy.engine.base.Engine PRAGMA main.table xinfo("Band")
2021-11-10 17:44:47.892 INFO sqlalchemy.engine.base.Engine ()
2021-11-10 17:44:47,894 INFO sqlalchemy.engine.base.Engine SELECT sql FROM (SELECT * FROM sqlite master UNION ALL SELECT * F
ROM sqlite temp master) WHERE name = 'Band' AND type = 'table'
2021-11-10 17:44:47,895 INFO sqlalchemy.engine.base.Engine ()
2021-11-10 17:44:47,962 INFO sqlalchemy.engine.base.Engine PRAGMA main.foreign key list("Band")
2021-11-10 17:44:47.963 INFO sqlalchemv.engine.base.Engine ()
2021-11-10 17:44:47,965 INFO sqlalchemy.engine.base.Engine PRAGMA temp.foreign_key_list("Band")
2021-11-10 17:44:47,965 INFO sqlalchemy.engine.base.Engine ()
2021-11-10 17:44:47.967 INFO sqlalchemy.engine.base.Engine SELECT sql FROM (SELECT * FROM sqlite master UNION ALL SELECT * F
ROM sqlite_temp_master) WHERE name = 'Band' AND type = 'table'
2021-11-10 17:44:47,969 INFO sqlalchemy.engine.base.Engine ()
2021-11-10 17:44:47,973 INFO sqlalchemy.engine.base.Engine PRAGMA main.index list("Band")
2021-11-10 17:44:47,975 INFO sqlalchemy.engine.base.Engine ()
2021-11-10 17:44:47,977 INFO sqlalchemy.engine.base.Engine PRAGMA temp.index list("Band")
2021-11-10 17:44:47,979 INFO sqlalchemy.engine.base.Engine ()
2021-11-10 17:44:47,982 INFO sqlalchemy.engine.base.Engine PRAGMA main.index list("Band")
2021-11-10 17:44:47,984 INFO sqlalchemy.engine.base.Engine ()
2021-11-10 17:44:47,985 INFO sqlalchemy.engine.base.Engine PRAGMA temp.index list("Band")
2021-11-10 17:44:47,987 INFO sqlalchemy.engine.base.Engine ()
2021-11-10 17:44:47,989 INFO sqlalchemy.engine.base.Engine SELECT sql FROM (SELECT * FROM sqlite master UNION ALL SELECT * F
ROM sqlite temp master) WHERE name = 'Band' AND type = 'table'
2021-11-10 17:44:47,990 INFO sqlalchemy.engine.base.Engine ()
2021-11-10 17:44:47,996 INFO sqlalchemy.engine.base.Engine SELECT "Band".id, "Band".fname, "Band".lname, "Band".birthyear
2021-11-10 17:44:47,997 INFO sqlalchemv.engine.base.Engine ()
```

```
import pandas as pd
import numpy as np
pd.set_option('max_columns', 4, 'max_rows', 10, 'max_colwidth', 12)
fname = ['Paul', 'John', 'Richard', 'George']
lname = ['McCartney', 'Lennon', 'Starkey', 'Harrison']
birth = [1942, 1940, 1940, 1943]
people = {'first': fname, 'last': lname, 'birth': birth}
import json
encoded = json.dumps(people)
encoded
'{"first": ["Paul", "John", "Richard", "George"], "last": ["McCartney", "Lennon", "Starkey", "Harrison"], "birth": [1942, 1940,
1940, 1943]}'
                                                                        beatles = pd.read json(encoded)
json.loads(encoded)
                                                                        beatles
{'first': ['Paul', 'John', 'Richard', 'George'],
                                                                              first
                                                                                        last birth
 'last': ['McCartney', 'Lennon', 'Starkey', 'Harrison'],
 'birth': [1942, 1940, 1940, 1943]}
                                                                              Paul McCartney 1942
                                                                         1
                                                                              John
                                                                                     Lennon 1940
                                                                         2 Richard
                                                                                     Starkey 1940
                                                                                                                              16
                                                                         3 George
                                                                                     Harrison 1943
```

```
records = beatles.to_json(orient='records')
records

'[{"first":"Paul","last":"McCartney","birth":1942},{"first":"John","last":"Lennon","birth":1940},{"first":"Richard","last":"Starkey","birth":1940},{"first":"George","last":"Harrison","birth":1943}]'

pd.read_json(records, orient='records')

first last birth

O Paul McCartney 1942

1 John Lennon 1940
2 Richard Starkey 1940
3 George Harrison 1943
```

```
split = beatles.to_json(orient='split')
split
'{"columns":["first","last","birth"],"index":[0,1,2,3],"data":[["Paul","McCartney",1942],["John","Lennon",1940],["Richard","Sta
rkey",1940],["George","Harrison",1943]]}'
pd.read json(split, orient='split')
     first
               last birth
     Paul McCartney 1942
     John
            Lennon 1940
2 Richard
            Starkey 1940
3 George
           Harrison 1943
index = beatles.to json(orient='index')
index
'{"0":{"first":"Paul","last":"McCartney","birth":1942},"1":{"first":"John","last":"Lennon","birth":1940},"2":{"first":"Richar
d","last":"Starkey","birth":1940},"3":{"first":"George","last":"Harrison","birth":1943}}'
pd.read json(index, orient='index')
     first
               last birth
     Paul McCartney 1942
            Lennon 1940
     John
2 Richard
            Starkey 1940
3 George
           Harrison 1943
```

```
values = beatles.to_json(orient='values')
values
'[["Paul", "McCartney", 1942], ["John", "Lennon", 1940], ["Richard", "Starkey", 1940], ["George", "Harrison", 1943]]
pd.read_json(values, orient='values')
        0
      Paul McCartney 1942
             Lennon 1940
1
      John
2 Richard
             Starkey 1940
 3 George
            Harrison 1943
(pd.read_json(values, orient='values')
   .rename(columns=dict(enumerate(['first', 'last', 'birth'])))
                last birth
      first
      Paul McCartney 1942
1
      John
             Lennon 1940
2 Richard
             Starkey 1940
 3 George
            Harrison 1943
```