

### Bài 10: Pandas 2

TS. Trịnh Tuấn Đạt Trường CNTT-TT, ĐHBK Hà Nội

1

1



# Nội dung

- 1. Thao tác đọc/ghi dữ liệu với file
- 2. Các thao tác cơ bản với Series và DataFrame
- 3. Đặc tả thống kê (Descriptive Statistics)
- 4. Sắp xếp trên DataFrame (Sorting)
- 5. Cắt nhỏ DataFrame (Slicing)
- 6. Loc DataFrame (Filtering)
- 7. Duyệt trên Series và DataFrame (Iteration)
- 8. Đánh lại nhãn cho DataFrame (Reindexing)
- 9. Gióng nhãn với reindex\_like

2

ว



## Nội dung

- 1. Thao tác đọc/ghi dữ liệu với file
- 2. Các thao tác cơ bản với Series và DataFrame
- 3. Đặc tả thống kê (Descriptive Statistics)
- 4. Sắp xếp trên DataFrame (Sorting)
- 5. Cắt nhỏ DataFrame (Slicing)
- 6. Loc DataFrame (Filtering)
- 7. Duyệt trên Series và DataFrame (Iteration)
- 8. Đánh lại nhãn cho DataFrame (Reindexing)
- 9. Gióng nhãn với reindex\_like

3

3



### 1. Thao tác đọc/ghi dữ liệu với file

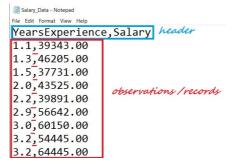
- Các định dạng file được hỗ trợ trong Pandas
  - Excel
  - CSV Comma Separated Values
  - JSON

4



### Đọc dữ liệu từ file CSV

- CSV là dạng thức đơn giản nhất để lưu trữ dữ liệu cấu trúc bảng dưới dạng thuần text
- Là định dạng dùng rất phổ biến trong khoa học dữ liệu
- Comma Separated Values: Các giá trị ngăn cách bằng dấu , nhưng có thể ngăn cách bằng dấu ;
- Cấu trúc file CSV



5



#### Đọc dữ liệu từ file CSV

Nếu dữ liệu chứa dấu, cần đặt trong cặp dấu "

```
id, name, organization, address, contact
    id, name, organization, address, contact
    1,Ranjith,tcs, "gachibowli, hyderabad",1234567890
    2,sudha, altimetrik, "electonic city, bangalore",78945632210
    4,Micheal, wipro, "dlf info city, chennai",2031456879
    4,kiran, infosys, "Rajiv Gandhi Infotech Park, pune",5520369741
    5,Ramesh, accenture, "financial district, hyderabad",5522001144
    6,mohammed, capegemini, "Rajiv Gandhi Infotech Park, pune",10023654989
```



### Đọc dữ liệu từ file CSV

Đọc dữ liệu file CSV vào DataFrame

```
Duration, Pulse, Maxpulse, Calories
                                          import pandas as pd
                                          df = pd.read csv('data.csv')
60.110.130.409.1
60,117,145,479.0
                                          print(df)
60,103,135,340.0
45,109,175,282.4
                                            Duration Pulse Maxpulse Calories
45,117,148,406.0
                                                 60
                                                      110
                                                             130
                                                                   409.1
60,102,127,300.5
                                                 60
                                                      117
                                                             145
                                                                   479.0
60,110,136,374.0
                                                 60
                                                      103
                                                             135
                                                                   340.0
                                                 45
45,104,134,253.3
                                                 45
                                                      117
                                                             148
30,109,133,195.1
60,98,124,269.0
                                          164
                                                 60
                                                      105
                                                             140
                                                                   290.8
                                                             145
145
                                                      110
                                          165
                                                 60
                                                                   300.4
                                                      115
                                                                   310.2
                                          [169 rows x 4 columns]
```

7



### Đọc dữ liệu từ file JSON

- JSON: JavaScript Object Notation
- Là đinh dang text, dùng để lưu trữ và truyền tải dữ liệu
- JSON có tính tự mô tả (self-describing), dễ hiểu
- Dùng rông rãi trong hầu hết các ngôn ngữ lập trình
- Ví du môt xâu dữ liêu lưu dưới dang JSON:

```
{"name":"John", "age":30, "car":null}
```

Các ngôn ngữ lập trình (Python, Java Script, Java, PHP, ...)
 có thể phân tích một xâu JSON thành các đối tượng JSON
 tương ứng

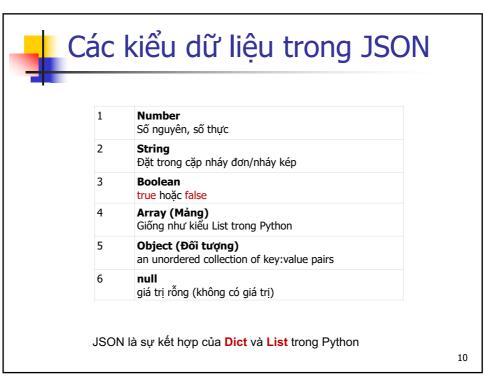
```
Ví dụ dữ liệu định dạng JSON

["Sunday", "Monday", "Tuesday",
"Wednesday", "Thursday", "Friday",
"Saturday"]

{
    "name":"John",
    "age":30,
    "cars":["Ford", "BMW", "Fiat"]
}

    "Museulse":{
    "0":110,
    "1":117,
    "2":103,
    "3":109,
    "4":117,
    "5":102
},

"Maxpulse":{
    "0":30,
    "1":145,
    "2":135,
    "3":175,
    "4":148,
    "5":127
}
```



```
Đọc dữ liệu từ file JSON
import pandas as pd
df = pd.read_json('data.json')
print(df)
 Duration Pulse Maxpulse Calories
      60
          110
                 130
                      409.1
                      479.0
      60
          117
                 145
2
                 135
                      340.0
      60
          103
      45
          109
                 175
                      282.4
      45
          117
                 148
                      406.0
164
      60
          105
                 140
                      290.8
165
      60
          110
                 145
                      300.4
                      310.2
166
          115
                 145
167
          120
                      320.4
168
          125
                      330.4
[169 rows x 4 columns]
                                                                        11
```



# Đọc dữ liệu từ file JSON

So sánh sự dễ dàng khi đọc dữ liệu từ file JSON, CSV, và file Excel?



# Ghi dữ liệu ra file CSV, JSON

 Khi ghi dữ liệu ra file CSV, có thể bỏ qua nhãn hàng (index=false), bỏ qua nhãn cột (header=false), chọn cột muốn ghi (tham số columns)



13



## Nội dung

- 1. Thao tác đọc/ghi dữ liệu với file
- 2. Các thao tác cơ bản với Series và DataFrame
- 3. Đặc tả thống kê (Descriptive Statistics)
- 4. Sắp xếp trên DataFrame (Sorting)
- 5. Cắt nhỏ DataFrame (Slicing)
- 6. Loc DataFrame (Filtering)
- 7. Duyệt trên Series và DataFrame (Iteration)
- 8. Đánh lại nhãn cho DataFrame (Reindexing)
- 9. Gióng nhãn với reindex\_like

```
Các thao tác cơ bản với Series
                                                                    [Index(['day1', 'day2',
'day3', 'day4', 'day5',
'day6'], dtype='object')]
import pandas as pd
                                                                    False
calories = {
     "day1": 420,
"day2": 380,
                                                                    [420 380 390 350 410 310]
     "day3": 390,
"day4": 350,
                                                                               420
                                                                    day1
                                                                    day2
                                                                               380
     "day5": 410,
"day6": 310}
                                                                    day3
                                                                    day4
                                                                               350
                                                                    day5
                                                                               410
myvar = pd.Series(calories)
                                                                    dtype: int64
                                                                    day1
                                                                               420
print(myvar.axes) # Lấy các nhãn (hoặc myvar.index)
                                                                   day2
day3
                                                                               380
print(myvar.empty) # Series có rong?
print(myvar.ndim) # Số chiều dữ liệu (là 1)
                                                                    dtype: int64
print(myvar.size) # Kích thước Series
print(myvar.values) # Series dưới dạng List
print(myvar.head()) # Tạo series con: 5 phần tử đầu
                                                                               390
                                                                    day4
                                                                               350
print(myvar.head(3)) # Tao series con: 3 phần tử dầu
print(myvar.tail()) # Tao series con: 5 phần tử cuối
                                                                    day5
                                                                               410
                                                                               310
                                                                    day6
                                                                   dtype: int64
print(myvar.tail(2)) # Tao series con: 2 phan tử cuối
                                                                    day5
                                                                    day6
                                                                               310
                                                                    dtype: int64
                                                                                                   15
```



Chuyển vị (Transpose: chuyển hàng → cột, cột → hàng)

```
import pandas as pd
data = {
    'Name':['Tom', 'Jack', 'Steve', 'Ricky'], 'Age':[28,34,29,42]}
df = pd.DataFrame(data, index=['rank1','rank2','rank3','rank4'])
print(df)
print(df.T)
         Name Age
         Tom
rank1
rank2
         Jack
                 34
rank3
        Steve
                 29
rank4 Ricky
     rank1 rank2 rank3 rank4
Name
      Tom Jack Steve Ricky
Age
         28
                34
                        29
                                42
                                                                               16
```



#### Các thao tác cơ bản với DataFrame

Lấy ra nhãn hàng, cột

```
import pandas as pd
data = {'Name':['Tom', 'Jack', 'Steve',
'Ricky'],'Age':[28,34,29,42]}
df = pd.DataFrame(data, index=['rank1','rank2','rank3','rank4'])
print(type(df.axes))
print(df.axes)
                                                      Name Age
                                             rank1 Tom
print(df.index)
                                             rank2
                                                      Jack
print(df.columns)
                                             rank3 Steve
                                                               29
                                             rank4 Ricky
<class 'list'>
[Index(['rank1', 'rank2', 'rank3', 'rank4'], dtype='object'),
Index(['Name', 'Age'], dtype='object')]
Index(['rank1', 'rank2', 'rank3', 'rank4'], dtype='object')
Index(['Name', 'Age'], dtype='object')
                                                                                                         17
```

17



#### Các thao tác cơ bản với DataFrame

Lấy ra kiểu dữ liệu cho các cột

```
import pandas as pd
data = {'Name':['Tom', 'Jack', 'Steve', 'Ricky'],'Age':[28,34,29,42]}
df = pd.DataFrame(data, index=['rank1','rank2','rank3','rank4'])
print(df.dtypes)
                                                                                  Name Age
print()
print(type(df.dtypes))
                                                                          rank1
                                                                                   Tom
                                                                                            28
                                                                          rank2
                                                                                   Jack
                                                                                            34
                                                                          rank3 Steve
                                                                                            29
Name
            object
                                                                          rank4 Ricky
             int64
Age
dtype: object
<class 'pandas.core.series.Series'>
```



#### Các thao tác cơ bản với DataFrame

Các thao tác khác

```
import pandas as pd
data = {'Name':['Tom', 'Jack', 'Steve', 'Ricky'],'Age':[28,34,29,42]}
df = pd.DataFrame(data, index=['rank1','rank2','rank3','rank4'])
print(df.empty) # Có rỗng không
                                                                                               Name Age
print(df.ndim) # Số chiều - luôn là 2
print(df.shape) # Số hàng, số cột
print(df.size) # Tổng số phần tử (hàng, cột)
                                                                                     rank1
                                                                                                Tom
print(df.values) # Trả về dưới dạng list
                                                                                     rank2
                                                                                                Jack
                                                                                     rank3 Steve
                                                                                                          29
False
                                                                                     rank4 Ricky
(4, 2)
[['Tom' 28]
 ['Jack' 34]
  ['Steve' 29]
['Ricky' 42]]
```

19



#### Các thao tác cơ bản với DataFrame

```
    Phương thức head, tail

import pandas as pd
data = {'Name':['Tom', 'Jack', 'Steve', 'Ricky'],'Age':[28,34,29,42]}
df = pd.DataFrame(data, index=['rank1','rank2','rank3','rank4'])
print(df.head())
print(df.head(1))
print(df.tail())
                                                      28
                                               Tom
                                    rank1
                                    rank2
                                               Jack
                                                        34
print(df.tail(2))
                                    rank3 Steve
                                     rank4 Ricky 42
                                            Name Age
                                     rank1 Tom 28
                                              Name Age
                                               Tom 28
                                     rank1
                                     rank2
                                               Jack
                                                        34
                                     rank3 Steve 29
                                     rank4 Ricky
                                                      42
                                              Name Age
                                     rank3 Steve
                                                      29
                                     rank4 Ricky 42
                                                                                       20
```

20



# Nội dung

- 1. Thao tác đọc/ghi dữ liệu với file
- 2. Các thao tác cơ bản với Series và DataFrame
- 3. Đặc tả thống kê (Descriptive Statistics)
- 4. Sắp xếp trên DataFrame (Sorting)
- 5. Cắt nhỏ DataFrame (Slicing)
- 6. Loc DataFrame (Filtering)
- 7. Duyệt trên Series và DataFrame (Iteration)
- 8. Đánh lại nhãn cho DataFrame (Reindexing)
- 9. Gióng nhãn với reindex\_like

21

21



# Đặc tả thống kê

Xét DataFrame sau:

	Name	Age	Rating
0	Tom	25	4.23
1	James	26	3.24
2	Ricky	25	3.98
3	Vin	23	2.56
4	Ricky	30	3.20
5	Smith	29	4.60
6	Jack	23	3.80
7	Lee	34	3.78
8	David	40	2.98
9	Gasper	30	4.80
10	Betina	51	4.10
11	Andres	46	3.65

```
Đặc tả thống kê - hàm sum
import pandas as pd
import numpy as np
  'Rating': [4.23,3.24,3.98,2.56,3.20,4.6,3.8,3.78,2.98,4.80,4.10,3.65]
df = pd.DataFrame(d)
sum = df.sum() # Tương đương với df.sum(axis = 0): cộng trên cột
print(type(sum))
print(sum)
                                                            3.24
<class 'pandas.core.series.Series'>
                                                    Ricky 25
                                                            3.98
        TomJamesRickyVinSteveSmithJackLeeDavidGasperBe..
                                                      Vin 23
                                                            2.56
                                       382
                                       44.92
Rating
dtype: object
                                                           23
```

Đặc tả thống kê – hàm sum

```
import pandas as pd
import numpy as np
  'Rating': [4.23,3.24,3.98,2.56,3.20,4.6,3.8,3.78,2.98,4.80,4.10,3.65]
                                          <class 'pandas.core.series.Series'>
df = pd.DataFrame(d)
                                          0 29.23
sum = df.sum(axis = 1) # Cong trên hàng
                                               29.24
print(type(sum))
                                               28.98
print(sum)
                          Name Age Rating
                                               25.56
                                               33.20
                                   4.23
                                               33.60
                                   3.24
                                               26.80
                                   3.98
                                               42.98
                                               34.80
                           Vin 23
                                   2.56
                                               55.10
                                         10
                                               49.65
                                         dtype: float64
                                                                      24
```

24

```
Đặc tả thống kê - hàm mean
import pandas as pd
import numpy as np
d = {
  'Rating': [4.23,3.24,3.98,2.56,3.20,4.6,3.8,3.78,2.98,4.80,4.10,3.65]
                                                  Name Age Rating
df = pd.DataFrame(d)
mean = df.mean(numeric_only=True)
                                                       25
                                                  Tom
                                                            4.23
print(type(meah))
                                                            3.24
                                                 James
print(mean)
                                                  Ricky
                                                       25
                                                            3.98
                                                   Vin
                                                       23
                                                            2.56
<class 'pandas.core.series.Series'>
                                                  Steve
                                                       30
                                                            3.20
         31.833333
         3.743333
                                                            4 60
Rating
dtype: float64
                                                             25
```

25



# Đặc tả thống kê

Sr.No.	Function	Description
1	count()	Number of non-null observations
2	sum()	Sum of values
3	mean()	Mean of Values
4	median()	Median of Values
5	mode()	Mode of values
6	std()	Standard Deviation of the Values
7	min()	Minimum Value
8	max()	Maximum Value
9	abs()	Absolute Value
10	prod()	Product of Values
11	cumsum()	Cumulative Sum
12	cumprod()	Cumulative Product

27

27



### Đặc tả thống kê – phương thức describe

Phương thức đặc tả

df.describe(include)

- include nhận 3 giá trị:
  - 'object': Thống kê với cột dữ liệu khác số
  - 'number': Thống kê với cột dữ liệu số
  - 'all': Thống kê tất cả các cột dữ liệu

```
include = 'object'
import pandas as pd
import numpy as np
d = {
  'Age': [25,26,25,23,30,29,23,34,40,30,51,46],
   'Rating': [4.23,3.24,3.98,2.56,3.20,4.6,3.8,3.78,2.98,4.80,4.10,3.65]
                                                 Name Age Rating
df = pd.DataFrame(d)
                                                  Tom
                                                      25
                                                           4.23
desc = df.describe(include='object')
print(type(desc))
                                                James
                                                           3.24
print(desc)
                                                 Ricky
                                                       25
                                                           3.98
<class 'pandas.core.frame.DataFrame'>
                                                  Vin
                                                           2.56
        Name
          12
count
                                                 Steve
                                                       30
                                                           3.20
          11
unique
                                                           4 AN
top
       Ricky
                                                              29
freq
            2
```

```
include = 'number'
import pandas as pd
import numpy as np
   'Name': ['Tom', 'James', 'Ricky', 'Vin', 'Ricky', 'Smith', 'Jack', 'Lee', 'David', 'Gasper', 'Betina', 'Andres'], 'Age': [25,26,25,23,30,29,23,34,40,30,51,46],
   'Rating': [4.23,3.24,3.98,2.56,3.20,4.6,3.8,3.78,2.98,4.80,4.10,3.65]
                                                                         Rating
                                                               Age
df = pd.DataFrame(d)
                                              count 12.000000 12.000000
desc = df.describe(include='number')
                                                       31.833333
                                                                      3.743333
                                              mean
print(desc)
                                              std
                                                        9.232682
                                                                      0.661628
                                              min
                                                       23.000000
                                                                      2.560000
                                              25%
                                                       25.000000
                                                                      3.230000
                                              50%
                                                       29.500000
                                                                      3.790000
                                                                      4.132500
                                              75%
                                                       35.500000
                                                       51.000000
                                                                       4.800000
                                              max
                                                                               30
```

```
include = 'all'
import pandas as pd
import numpy as np
   'Name': ['Tom','James','Ricky','Vin', 'Ricky','Smith','Jack',
'Lee','David','Gasper','Betina','Andres'],
'Age': [25,26,25,23,30,29,23,34,40,30,51,46],
    'Rating': [4.23,3.24,3.98,2.56,3.20,4.6,3.8,3.78,2.98,4.80,4.10,3.65]
                                                                           Rating
                                                      Name
                                                                   Age
df = pd.DataFrame(d)
                                                       12 12.000000 12.000000
                                            count
desc = df.describe(include='all')
                                            unique
                                                                  NaN
                                                                               NaN
print(desc)
                                                     Ricky
                                                                   NaN
                                                                               NaN
                                            top
                                            freq
                                                                   NaN
                                                                               NaN
                                                                        3.743333
                                                       NaN 31.833333
                                            mean
                                            std
                                                       NaN
                                                             9.232682
                                                                         0.661628
                                                       NaN 23.000000
                                                                         2.560000
                                            min
                                                       NaN 25.000000
                                                      NaN 29.500000
NaN 35.500000
                                            50%
                                                                         3.790000
                                                      NaN 51.000000
                                                                         4.800000
                                            max
                                                                                      31
```



# Nội dung

- 1. Thao tác đọc/ghi dữ liệu với file
- 2. Các thao tác cơ bản với Series và DataFrame
- 3. Đặc tả thống kê (Descriptive Statistics)
- 4. Sắp xếp trên DataFrame (Sorting)
- 5. Cắt nhỏ DataFrame (Slicing)
- 6. Loc DataFrame (Filtering)
- 7. Duyệt trên Series và DataFrame (Iteration)
- 8. Đánh lại nhãn cho DataFrame (Reindexing)
- 9. Gióng nhãn với reindex\_like



# Sắp xếp

- Có thể sắp xếp DataFrame theo nhãn các hàng/cột, hoặc theo giá trị
- Có thể sắp xếp DataFrame theo thứ tự tăng dần, giảm dần

33

34

33



# Sắp xếp nhãn của hàng/cột

- axis=1: sắp xếp nhãn của cột. axis=0: sắp xếp nhãn của hàng. Mặc đinh axis=0
- Ví dụ sắp xếp nhãn của cột:

```
import pandas as pd
                                        28
                                            15
                                                46
data = {
                                        34
                                                27
    'c': [5,9,1,10],
                                        29
                                            11
    'd': [28,34,29,42],
    'b': [15,16,11,24],
                                     а
    'a': [46,27,10,30]
                                 2 46 15
                                 0 27 16
df = pd.DataFrame(data,
                                 3 10
                                        11
    index = [2, 0, 3, 1])
                                        24
print(df)
df1 = df.sort_index(axis = 1)
print(df1)
```



# Sắp xếp nhãn của hàng/cột

- axis=1: sắp xếp theo nhãn của cột. axis=0: sắp xếp theo nhãn của hàng. Mặc định axis=0
- Ví du sắp xếp nhãn của hàng

```
import pandas as pd
                              2 5 28 15 46
data = {
                              0 9 34 16 27
   'c': [5,9,1,10],
                              3 1 29
                                        11
   'd': [28,34,29,42],
                              1 10 42
                                        24
   'b': [15,16,11,24],
                                  С
                                     d
                                        b
                                            а
                                 9
    'a': [46,27,10,30]
                                        16
                              1 10
                                     42
                                        24
df = pd.DataFrame(data,
                                 5 28
                                        15 46
   index = [2, 0, 3, 1])
                              3 1 29 11 10
print(df)
df1 = df.sort_index()
print(df1)
```

35



# Sắp xếp nhãn của hàng/cột

35

36

 Có thể sắp xếp nhãn theo thứ tự tăng dần (ascending=True) hoặc giảm dần (ascending=False). Mặc định là sắp xếp tăng dần

```
import pandas as pd
                                            b
data = {
                                       28
                                           15
                                               46
   'c': [5,9,1,10],
                                       34
                                           16
                                               27
    'd': [28,34,29,42],
                                3 1 29 11 10
    'b': [15,16,11,24],
    'a': [46,27,10,30]
                                2 28
df = pd.DataFrame(data,
                                       9 16 27
   index = [2, 0, 3, 1])
                                3 29
                                       1 11 10
print(df)
                                       10 24 30
df1 = df.sort_index(axis=1, ascending=False)
print(df1)
```



# Sắp xếp theo giá trị

- Có thể sắp xếp theo giá trị. Tham số axis chỉ ra sắp xếp các hàng/cột. Tham số by chỉ ra nhãn của hàng/cột được lấy làm cơ sở sắp xếp. Có thể là giá trị đơn hoặc một list.
- by đi theo axis. VD: nếu axis=1 là sắp xếp lại cột, by sẽ là nhãn của 1 hàng nào đó, và giá trị của hàng này được lấy làm cơ sở để sắp xếp các cột

```
import pandas as pd
data = {
    'a': [1,5,5,10],
    'b': [28,29,34,42],
    'c': [15,16,11,24],
    'd': [46,27,10,30]
}
df = pd.DataFrame(data)
print(df)

df1 = df.sort_values(
    axis=1, ascending=True,
    by=0)
print(df1)
```

```
d
         C
    28
        15
            46
   29
        16
            27
5
   34
        11
            10
10
   42
        24
            30
 a
    С
        b
            d
 1
   15
        28
            46
 5 16
        29
            27
   11
        34
            10
10 24
        42
            30
```

**Lưu ý**: trong ví dụ này, hàng không có nhãn, nên hàng sẽ nhận nhãn mặc định theo STT



# Sắp xếp theo giá trị

- by nhận giá trị List → có thể sắp xếp các giá trị trên nhiều hàng, nhiều cột
- VD: sắp xếp các hàng theo giá tri trên 2 côt: côt a trước, côt b sau

```
import pandas as pd
                                          28
                                              15
                                                   46
data = {
                                       5 29
                                               16
                                                   27
    'a': [1,5,5,10],
                                       5
                                          34
                                               11
                                                   10
    'b': [28,29,34,42],
                                      10
                                          42
                                                   30
    'c': [15,16,11,24],
                                          b
                                       а
                                               C
                                                   d
    'd': [46,27,10,30]
                                      10
                                          42
                                               24
                                                   30
                                       5
                                          34
                                               11
                                                   10
df = pd.DataFrame(data)
                                          29
                                       5
                                               16
                                                   27
print(df)
                                         28
                                              15
                                       1
                                                  46
df1 = df.sort_values(
    axis=0, ascending=False,
by=['a','b'])
```

print(df1)



### Nội dung

- 1. Thao tác đoc/ghi dữ liêu với file
- 2. Các thao tác cơ bản với Series và DataFrame
- 3. Đặc tả thống kê (Descriptive Statistics)
- 4. Sắp xếp trên DataFrame (Sorting)
- 5. Cắt nhỏ DataFrame (Slicing)
- 6. Loc DataFrame (Filtering)
- 7. Duyệt trên Series và DataFrame (Iteration)
- 8. Đánh lại nhãn cho DataFrame (Reindexing)
- 9. Gióng nhãn với reindex\_like

39

39



#### Cắt nhỏ dữ liệu (slicing) trong DataFrame

 Lấy ra tất cả các hàng trong 1 cột (lấy ra 1 cột) với thuộc tính loc

```
import pandas as pd

d = {'one':[1, 2, 3], 'two': [5, 6, 7]}
df = pd.DataFrame(d, index=['a', 'b', 'c'])

print(df, end="\n\n")

# Lấy ra tất cả các hàng trong 1 cột
print(df.loc[:,'two']) # Trả về 1 Series

one two
a 1 5
b 2 6
c 3 7
Name: two, dtype: int64
```

40



### Cắt nhỏ dữ liệu (slicing) trong DataFrame

Lấy ra tất cả các hàng trong nhiều cột với thuộc tính loc

```
import pandas as pd
                                                                          three
                                                                     5
                                                                               8
                                                        b
                                                               2
                                                                      6
                                                                               9
     'one':[1, 2, 3],
                                                                              10
'two': [5, 6, 7],
'three': [8, 9, 10]}
df = pd.DataFrame(d, index=['a', 'b', 'c'])
                                                                      2
                                                        b
print(df, end="\n\n")
# Lấy ra tất cả các hàng trong 2 cột
print(df.loc[:, ['two', 'one']]) # Tra ve 1 DataFrame
                                                                                   41
```

41



#### Cắt nhỏ dữ liệu (slicing) trong DataFrame

Lấy ra một vài hàng trong nhiều cột với thuộc tính loc

```
import pandas as pd
                                                      one
                                                            two
                                                              5
    one':[1, 2, 3],
                                                                     10
    'two': [5, 6, 7],
'three': [8, 9, 10]}
df = pd.DataFrame(d, index=['a', 'b', 'c'])
                                                       7
                                                  С
                                                            3
                                                        5
                                                              1
print(df, end="\n\n")
# Lấy ra tất cả dữ liệu trong 2 hàng, 2 cột
print(df.loc[['c', 'a'], ['two', 'one']]) # Tra ve 1 DataFrame
                                                                          42
```



### Cắt nhỏ dữ liệu (slicing) trong DataFrame

Lấy ra một dải các hàng trong tất cả các cột với thuộc tính

```
import pandas as pd
                                                                   two
                                                                          three
                                                                     5
                                                                               8
                                                        а
                                                        b
                                                               2
                                                                      6
                                                                               9
     'one':[1, 2, 3],
                                                                              10
'two': [5, 6, 7],
'three': [8, 9, 10]}
df = pd.DataFrame(d, index=['a', 'b', 'c'])
                                                             1
                                                                               8
print(df, end="\n\n")
# Lấy ra dải các hàng
print(df.loc['a':'b']) # Tra ve 1 DataFrame
                                                                                   43
```

43



### Cắt nhỏ dữ liệu (slicing) trong DataFrame

Lấy ra một dải các hàng và cột với thuộc tính loc

```
import pandas as pd
                                                      one
                                                            two
                                                             5
    'one':[1, 2, 3],
                                                                     10
    'two': [5, 6, 7],
'three': [8, 9, 10]}
df = pd.DataFrame(d, index=['a', 'b', 'c'])
                                                      5
                                                                 8
                                                  а
                                                  b
                                                        6
print(df, end="\n\n")
# Lấy ra dải các hàng và cột
print(df.loc['a':'b' , 'two':'three']) # Tra ve 1 DataFrame
                                                                          44
```



# Lấy ra DataFrame các giá trị Boolean

one

two

three

46

 Thuộc tính loc có thể dùng kết hợp với biểu thức

```
1
                                                                       5
                                                                                8
                                                         b
                                                                 2
import pandas as pd
d = {
    'one':[1, 2, 3],
                                                                   three
'two': [5, 6, 7],
'three': [8, 9, 10]}
df = pd.DataFrame(d, index=['a', 'b', 'c'])
                                                                5
                                                                          9
print(df, end="\n\n")
                                                                 two three
                                                             False
                                                                        True
                                                             False
# Lấy ra 1 dải các hàng và cột
print(df.loc['a':'b', 'two':'three']) # Trả về 1 DataFrame
# Lấy ra DataFrame các giá trị boolean theo điều kiện check
print(df.loc['a':'b', 'two':'three']>6) # Tra ve 1 DataFrame
                                                                                  45
```

45



#### Cắt nhỏ dữ liệu trong DataFrame với phương thức iloc

 iloc: tương tự như loc, nhưng không dùng nhãn, mà dùng STT, đánh số từ 0

```
import pandas as pd
d = \{ 'one' : [1, 2, 3], 'two' : [5, 6, 7], 'three' : [8, 9, 10] \}

df = pd.DataFrame(d, index=['a', 'b', 'c'])
                                                     one
                                                             two
                                                                    three
                                                                          8
                                                                6
                                                                          9
                                                 b
print(df.iloc[1:3, 0:1])
                                                                        10
                                                        3
print(df.iloc[1:2, :])
                                                     one
print(df.iloc[:1])
                                                 b
print(df.iloc[1:])
                                                     one
                                                             two
                                                                    three
                                                 b
```



## Nội dung

- 1. Thao tác đoc/ghi dữ liêu với file
- 2. Các thao tác cơ bản với Series và DataFrame
- 3. Đặc tả thống kê (Descriptive Statistics)
- 4. Sắp xếp trên DataFrame (Sorting)
- 5. Cắt nhỏ DataFrame (Slicing)
- 6. Loc DataFrame (Filtering)
- 7. Duyệt trên Series và DataFrame (Iteration)
- 8. Đánh lại nhãn cho DataFrame (Reindexing)
- 9. Gióng nhãn với reindex\_like

47

47



### Loc DataFrame

 Có thể dùng các toán tử logic trên các cột để lọc ra các hàng trong DataFrame

```
Name Age Rating
import pandas as pd
                                                Tom 25 4.23
import numpy as np
                                             1 James 26
                                                                3.24
                                             2 Ricky 25
                                                                3.98
 'Name': ['Tom','James','Ricky','Vin'],
'Age': [25,26,25,23],
'Rating': [4.23,3.24,3.98,2.56]
                                                 Name Age Rating
                                             0 Tom 25
                                                              4.23
                                             2 Ricky 25
                                                                3.98
df = pd.DataFrame(data)
print(df)
df1 = df[df.Age==25] # Hoặc dùng df["Age"]==25
print(df1)
                                                                     48
```



#### Lọc DataFrame – kết hợp điều kiện

 Có thể kết hợp các điều kiện. Lưu ý: Mỗi điều kiện phải đặt trong cặp dấu ngoặc tròn (). Toán tử and là &. Toán tử or là |. Toán tử not là ~

```
Name Age Rating
import pandas as pd
                                                 Tom 25
                                                            4.23
import numpy as np
                                               James
                                                        26
                                                              3.24
data = {
                                               Ricky
                                                        25
                                                              3.98
                                                      23
  'Name': ['Tom','James','Ricky','Vin'],
                                                              2.56
                                                Vin
  'Age': [25,26,25,23],
'Rating': [4.23,3.24,3.98,2.56]
                                                Name Age Rating
                                                      25
                                                            4.23
                                                Tom
                                                      26
                                                              3.24
                                              James
df = pd.DataFrame(data)
                                                      25
                                                              3.98
print(df)
df1 = df[(df.Name=="James") | (df["Age"]==25)]
print(df1)
                                                                   49
```

49



### Loc DataFrame – phương thức **isin**

Phương thức isin

```
Name Age Rating
import pandas as pd
                                                    Tom 25
                                                                 4.23
import numpy as np
                                                           26
                                                                    3.24
                                                   James
                                                   Ricky
                                                          25
                                                                    3.98
  'Name': ['Tom','James','Ricky','Vin'],
'Age': [25,26,25,23],
'Rating': [4.23,3.24,3.98,2.56]
                                                    Vin 23
                                                                    2.56
                                                    Name Age Rating
                                               2 Ricky 25
                                                                  3.98
                                                     Vin
                                                             23
                                                                    2.56
df = pd.DataFrame(data)
print(df)
names = ['Bob', 'David', 'Ricky', 'Vin']
df1 = df[df.Name.isin(names)]
print(df1)
                                                                        50
```



## Lọc DataFrame – Xử lý xâu

Xử lý xâu

```
Name Age Rating
import pandas as pd
                                                Tom
                                                      25
                                                            4.23
import numpy as np
                                                        26
                                                              3.24
                                               James
data = {
  'Name': ['Tom','James','Ricky','Vin'],
'Age': [25,26,25,23],
                                               Ricky
                                                        25
                                                              3.98
                                                       23
                                                              2.56
                                                Vin
 'Rating': [4.23,3.24,3.98,2.56]
                                                Name Age Rating
                                                             3.24
                                            1 James
                                                      26
df = pd.DataFrame(data)
print(df)
                                                Name Age Rating
                                                      25
                                                             3.98
df1 = df[df.Name.str.startswith("J")]
                                                 Vin
                                                              2.56
print(df1)
df2 = df[df.Name.str.contains("i")]
print(df2)
                                                                   51
```

51



#### Loc DataFrame – phương thức query

 Phương thức query: xử lý code ngắn và linh hoạt hơn. Nhưng dễ bị lỗi do câu truy vấn là dạng xâu

```
Name Age Rating
import pandas as pd
                                                      Tom
                                                             25
                                                                   4.23
import numpy as np
                                                             26
                                                                      3.24
                                                 1 James
                                                     Ricky
                                                             25
                                                                      3.98
  'Name': ['Tom','James','Ricky','Vin'],
'Age': [25,26,25,23],
'Rating': [4.23,3.24,3.98,2.56]
                                                     Vin 23
                                                                      2.56
                                                     Name Age Rating
                                                 1 James 26
2 Ricky 25
                                                                    3.24
df = pd.DataFrame(data)
print(df)
df1 = df.query('Name!="Tom" and Age>23')
print(df1)
                                                                           52
```



# Loc DataFrame – phương thức nlargest và nsmallest

 Lấy ra N hàng có giá trị lớn nhất/nhỏ nhất. Chỉ áp dụng được với cột có kiểu Number. Sẽ sắp xếp lại các hàng luôn.

```
Name Age Rating
import pandas as pd
                                             Tom 25
                                                         4.23
import numpy as np
                                           James
                                         1
                                                    26
                                                           3.24
data = {
                                            Ricky
                                                    25
                                                           3.98
                                                  23
 'Name': ['Tom','James','Ricky','Vin'],
                                                          2.56
                                             Vin
 'Age': [25,26,25,23],
'Rating': [4.23,3.24,3.98,2.56]
                                             Name Age Rating
                                                   26
                                                         3.24
                                         1 James
                                                           4.23
df = pd.DataFrame(data)
                                        2 Ricky
                                                          3.98
print(df)
df1 = df.nlargest(3, 'Age')
                                             Name Age Rating
print(df1)
                                             Vin 23
                                                        2.56
df2 = df.nsmallest(2, 'Rating')
                                         1 James
                                                    26
                                                           3.24
print(df2)
```

53



### Nội dung

- 1. Thao tác đọc/ghi dữ liệu với file
- 2. Các thao tác cơ bản với Series và DataFrame
- 3. Đặc tả thống kê (Descriptive Statistics)
- 4. Sắp xếp trên DataFrame (Sorting)
- 5. Cắt nhỏ DataFrame (Slicing)
- 6. Loc DataFrame (Filtering)
- 7. <u>Duyêt trên Series và DataFrame (Iteration)</u>
- 8. Đánh lại nhãn cho DataFrame (Reindexing)
- 9. Gióng nhãn với reindex\_like



### Duyệt trên DataFrame và Series

- Vòng lặp for i in Series: i lần lượt là các giá trị trong Series
- Vòng lặp for i in DataFrame: i lần lượt là các nhãn cột

```
import pandas as pd
data = {
     'Name':['Tom', 'Jack
'Age':[28,34,29,42],
                    'Jack', 'Nam', 'Ricky'],
    'Place': ['Paris', 'London', 'Hanoi', 'LA']
df = pd.DataFrame(data, index=['rank1', 'rank2', 'rank3', 'rank4'])
print(df)
                                               Name Age Place
for col in df:
                                               Tom
                                                      2.8
                                                            Paris
                                      rank1
                                                      34 London
                                               Jack
  print(col)
                                      rank2
                                               Nam 29 Hanoi
                                      rank3
for value in df['Name']:
   print(value, end=" ")
                                      rank4 Ricky 42
                                      Name
                                      Age
                                      Place
Lưu ý: df ['Name'] là một Series
                                                                            55
                                      Tom Jack Nam Ricky
```

55



### iteritems()

 DataFrame.iteritems(): trả về đối tượng dạng như dict, với key là nhãn của một cột và value là Series ứng với cột đó

```
Name <class
'pandas.core.series.Series'>
rank1
         Tom
rank3
          Nam
       Ricky
rank4
Name: Name, dtype: object
Age <class
'pandas.core.series.Series'>
rank1
        28
rank2
         34
rank3
        29
rank4
        42
Place <class
'pandas.core.series.Series'>
        Paris
rank2
        London
        Hanoi
rank3
rank4
Name: Place, dtype: object
```



 DataFrame.iterrows(): trả về đối tượng dạng như dict, với key là nhãn của một hàng và value là Series ứng với hàng đó

```
rank1
Name
           Tom
Age
            2.8
Place
        Paris
Name: rank1, dtype: object
rank2
Name
           Jack
Age
            34
Place
       London
Name: rank2, dtype: object
rank3
Name
           Nam
Age
            29
         Hanoi
Name: rank3, dtype: object
         Ricky
Name
Age
            42
Place
           LA
Name: rank4, dtype: object
                           57
```

57



• itertuples: trả về các dòng, mỗi dòng như một namedtuples

```
import pandas as pd
data = {
     a = {
'Name':['Tom', 'Jack', 'Nam', 'Ricky'],
'Age':[28,34,29,42],
'Place': ['Paris', 'London', 'Hanoi', 'LA']
                                                                          Name Age
df = pd.DataFrame(data,
                                                                          Tom 28 Paris
Jack 34 London
     index=['rank1', 'rank2', 'rank3', 'rank4'])
                                                                 rank2
                                                                           Nam 29
                                                                 rank3
                                                                                         Hanoi
                                                                 rank4
print(df)
for row in df.itertuples():
  print(type(row), end="\n\n") <class 'pandas.core.frame.Pandas'>
print(row, end="\n\n")
  print(list(row), end="\n\n") Pandas(Index='rank1', Name='Tom', Age=28, Place='Paris')
print(row.Place)
                                         ['rank1', 'Tom', 28, 'Paris']
  print(row[3])
  break;
                                         Paris
                                                                                                        58
```



# Duyệt trên DataFrame

- Lưu ý: khi duyệt DataFrame, chỉ nên đọc dữ liệu. Không nên thay đổi DataFrame, vì đây là 1 việc vô ích.
- Khi duyệt, các đối tượng được duyệt là bản sao dữ liệu từ DataFrame. Thay đổi trên các bản sao này sẽ không cập nhật được DataFrame

59

59



# Nội dung

- 1. Thao tác đọc/ghi dữ liệu với file
- 2. Các thao tác cơ bản với Series và DataFrame
- 3. Đặc tả thống kê (Descriptive Statistics)
- 4. Sắp xếp trên DataFrame (Sorting)
- 5. Cắt nhỏ DataFrame (Slicing)
- 6. Loc DataFrame (Filtering)
- 7. Duyệt trên Series và DataFrame (Iteration)
- 8. Đánh lại nhãn cho DataFrame (Reindexing)
- 9. Gióng nhãn với reindex\_like

60



# Đánh lại nhãn với reindex

Chọn và sắp xếp lại nhãn hàng, cột

```
import pandas as pd
data = {
     'Name':['Tom', 'Jack', 'Nam', 'Ricky'],
'Age':[28,34,29,42],
'Place': ['Paris', 'London', 'Hanoi', 'Newyork']
df = pd.DataFrame(data, index=['rank1','rank2','rank3','rank4'])
print(df)
df reindexed = df.reindex(
    index=['rank1', 'rank4'],
columns=['Age', 'Name', 'Language'])
                                                        Name Age
print(df_reindexed)
                                                         Tom 28
                                              rank1
                                                                         Paris
                                              rank2
                                                         Jack
                                                                 34
                                                                        London
                                                       Nam 29
                                              rank3
                                                                       Hanoi
                                              rank4 Ricky 42 Newyork
Lưu ý: Có thể bỏ cột cũ, sắp lại thứ tự
                                                              Name Language
cột, hoặc thêm cột mới, như cột
                                              rank1
                                                       28
                                                             Tom
                                                                        NaN
Language. Các ô trong cột mới sẽ có
                                              rank4
                                                        42 Ricky
                                                                             NaN
                                                                                       61
giá trị là NaN
```

61



# Đánh lại nhãn với reindex

Chọn và sắp xếp lại nhãn hàng, cột

```
import pandas as pd
data = {
    'Name':['Tom', 'Jack', 'Nam', 'Ricky'],
'Age':[28,34,29,42],
'Place': ['Paris', 'London', 'Hanoi', 'Newyork']
df = pd.DataFrame(data)
print(df)
df_reindexed = df.reindex(
                                                                          Place
                                                         Name
                                                                Age
    index=[0, 2, 3, 4], columns=['Age', 'Name', 'Language'])
                                                          Tom
                                                                  28
                                                                          Paris
                                                         Jack
                                                                   34
                                                                         London
                                                    1
print(df_reindexed)
                                                          Nam
                                                                   29
                                                     3 Ricky
                                                                   42 Newyork
                                                          Age
                                                                  Name
                                                                          Language
                                                     0 28.0
                                                                   Tom
                                                                                 NaN
                                                       29.0
                                                                   Nam
                                                                                 NaN
                                                     3 42.0 Ricky
                                                                                 NaN
                                                          NaN
                                                                   NaN
                                                                                 NaN
                                                                                          62
```



#### Đổi tên nhãn cho các cột, hàng

```
import pandas as pd
                                                                 Name Age
                                                       rank1
                                                                 Tom 28
                                                                              Paris
                                                       rank2
                                                                 Jack
                                                                          34 London
                                                                        29
                                                      rank3
                                                                 Nam
                                                                               Hanoi
     'Name':['Tom', 'Jack', 'Nam', 'Ricky'],
    'Age':[28,34,29,42],
'Place': ['Paris', 'London', 'Hanoi', 'LA']
                                                       Sau khi đổi tên nhãn:
df = pd.DataFrame(data,
    index=['rank1', 'rank2', 'rank3', 'rank4'])
                                                                                Place
                                                                   Tên Tuổi
                                                                         28
                                                       Hang 1
                                                                   Tom
                                                                                  Paris
                                                                           34 London
                                                       rank2
                                                                  Jack
print(df)
                                                       Hạng 3
                                                                         29 Hanoi
                                                       rank4
                                                                Ricky
df = df.rename(
    columns={'Name': 'Tên', 'Age': 'Tuổi', 'New': 'Mới'}, index={'rank1': 'Hạng 1', 'rank3': 'Hạng 3'}
print("\nSau khi đổi nhãn: \n")
print(df)
                                                                                      63
```

63



# Nội dung

- 1. Thao tác đọc/ghi dữ liệu với file
- 2. Các thao tác cơ bản với Series và DataFrame
- 3. Đặc tả thống kê (Descriptive Statistics)
- 4. Sắp xếp trên DataFrame (Sorting)
- 5. Cắt nhỏ DataFrame (Slicing)
- 6. Loc DataFrame (Filtering)
- 7. Duyệt trên Series và DataFrame (Iteration)
- 8. Đánh lại nhãn cho DataFrame (Reindexing)
- 9. Gióng nhãn với reindex like



#### Nâng cao: gióng nhãn với reindex\_like

65

65



# Gióng nhãn với reindex\_like

```
Lưu ý: df1. reindex_like(df2) sẽ lấy tất cả
                                     các nhãn của df2, nhưng lấy dữ liệu tương ứng
import pandas as pd
                                     với nhãn đó trong df1 để làm dữ liệu cho các cell
    'Name':['Tom', 'Jack', 'Nam', 'Ricky'],
'Age':[28,34,29,42],
    'Place': ['Paris', 'London', 'Hanoi', 'Newyork']
df1 = pd.DataFrame(data1,
index=['rank1','rank2','rank3','rank4'])
                                                         Name Age
                                                                       Place
                                               rank1
                                                         Tom
                                                                28
                                                                       Paris
                                               rank2
                                                         Jack
                                                                      London
                                               rank3
                                                         Nam
                                                                29
                                                                       Hanoi
    'Name':['David', 'Alice', 'Bob'],
                                                                42
    'Age':[29,24,20],
'Place': ['Rio', 'LA', 'Berlin']
                                               rank4
                                                       Ricky
                                                                     Newyork
                                                        Name Age
                                                                      Place
                                               rank2
                                                       David
                                                               29
                                                                        Rio
rank1
                                                       Alice
                                                                24
                                                                         LA
                                                               20
                                               rank5
                                                         Bob
                                                                     Berlin
                                                       Name
                                                               Age
                                                                      Place
                                               rank2
                                                       Jack 34.0
                                                                     London
print(df2)
                                               rank1
                                                        Tom 28.0
                                                                      Paris
df1 = df1.reindex_like(df2)
                                               rank5
                                                        NaN
                                                              NaN
                                                                        NaN
print(df1)
                                                                             66
```



### Điền bổ sung dữ liệu khi đánh lại nhãn

- Các phương thức reindex, reindex\_like có thể nhận tham số method để bổ sung dữ liêu ở các trường NaN
  - pad/ffill Điền giá trị tiến lên (giá trị với index phía trước dùng để điền cho giá trị NaN của index phía sau)
  - bfill/backfill Điền giá trị lùi lại (giá trị sau dùng để điền cho giá trị NaN phía trước)
  - nearest Tìm giá trị với index gần nhất để điền
    - Chỉ dùng được khi nhãn của hàng có kiểu Number
- Lưu ý: df1.reindex\_like(df2)sẽ lấy tất cả các nhãn của df2, nhưng lấy giá trị của df1 để điền bổ sung dữ liệu

67

67

#### Ví dụ: nhãn của hàng là kiểu Number Name Age Place import pandas as pd 28 Paris Tom London Jack Nam 29 Hanoi 4 Ricky 42 Newyork df1 = pd.DataFrame(data1, index=[1, 2, 3, 4])Name Age Place $data2 = {$ 'Name':['David', 'Alice', 'Bob', "Kelvin"], David Rio 'Age':[29,24,20, 22], 'Place':['Rio', 'LA', 'Berlin', "Rome"] LA Alice Bob 20 Berlin df2 = pd.DataFrame(data2, index=[4, 3, 10, -1]) 22 Kelvin Rome print(df1.reindex\_like(df2)) print(df1.reindex\_like(df2, method="ffill")) print(df1.reindex\_like(df2, method="bfill")) print(df1.reindex\_like(df2, method="nearest")) 68

```
Ví dụ: nhãn của hàng là kiểu Number
                                                                  Name
                                                                                    Place
import pandas as pd
                                                                 Ricky
                                                            4
                                                                          42.0
                                                                                 Newvork
                                                            3
                                                                   Nam
                                                                          29.0
                                                                                    Hanoi
data1 = {
    al = {
'Name':['Tom', 'Jack', 'Nam', 'Ricky'],
'Age':[<mark>28,34,29,42</mark>],
'Place': ['Paris', 'London', 'Hanoi', 'Newyork']
                                                            10
                                                                   NaN
                                                                           NaN
                                                                                      NaN
                                                                   NaN
                                                                           NaN
                                                                                      NaN
                                                                                    Place
                                                                           Age
                                                                 Ricky
                                                                          42.0
                                                                                 Newvork
df1 = pd.DataFrame(data1, index=[1, 2, 3, 4])
                                                            3
                                                                   Nam
                                                                          29.0
                                                                                    Hanoi
                                                            10
                                                                          42.0
                                                                 Rickv
                                                                                 Newyork
                                                                   NaN
                                                                           NaN
data2 = {
    'Name':['David', 'Alice', 'Bob', "Kelvin"],
                                                                  Name
                                                                           Age
                                                                                    Place
    'Age':[29,24,20, 22],
'Place': ['Rio', 'LA',
                                                                 Ricky
                                                                          42.0
                                                                                 Newyork
                             'Berlin', "Rome"]
                                                            3
                                                                   Nam
                                                                          29.0
                                                           10
                                                                   NaN
                                                                          NaN
                                                                                     NaN
df2 = pd.DataFrame(data2,
                            index=[4, 3, 10, -1])
                                                                   Tom
                                                                          28.0
                                                                                    Paris
                                                                  Name
                                                                                   Place
                                                                          Age
print(df1.reindex_like(df2))
print(df1.reindex_like(df2, method="ffill"))
                                                            4
                                                                 Ricky
                                                                           42
                                                                                Newyork
                                                            3
                                                                           29
print(df1.reindex_like(df2, method="bfill"))
                                                                   Nam
                                                                                   Hanoi
print(df1.reindex_like(df2, method="nearest"))
                                                           1.0
                                                                Ricky
                                                                           42
                                                                                Newyork
                                                                   Tom
                                                                           28
                                                                                   Paris
                                                                                        69
```

```
Lưu ý: dữ liệu được lấy từ df1 để điền
import pandas as pd
                                                                      Name Age
                                                                                    Place
                                                                              28
                                                                                      Paris
                                                                       Tom
data1 = {
    'Name':['Tom','Jack','Nam','Ricky',"Alice"],
                                                                       Jack
                                                                                    London
  'Age':[28,34,29,42,30],
'Place':['Paris','London','Hanoi','Newyork',"LA"]
                                                                       Nam
                                                                              29
                                                                                     Hanoi
                                                                     Ricky
                                                                               42 Newyork
df1 = pd.DataFrame(data1, index=[1, 2, 3, 4, 5])
                                                                       Alice
                                                                               30
data2 = {
    'Name':['David', 'Alice', 'Bob', "Kelvin"],
'Age':[29,24,20, 22],
'Place': ['Rio', 'LA', 'Berlin', "Rome"]
                                                                        Name
                                                                              Age Place
                                                                                       Rio
                                                                        David
                                                                                29
df2 = pd.DataFrame(data2, index=[4, 3, 10, -1])
                                                                         Alice
print(df1.reindex_like(df2))
                                                                    10
                                                                         Bob
                                                                                20
                                                                                     Berlin
print(df1.reindex_like(df2, method="ffill"))
print(df1.reindex_like(df2, method="bfill"))
                                                                    -1
                                                                       Kelvin
                                                                                22
                                                                                     Rome
print(df1.reindex_like(df2, method="nearest"))
                                                                                             70
```

```
Lưu ý: dữ liệu được lấy từ df1 để điền
                                                              Name
                                                                               Place
import pandas as pd
                                                             Ricky
                                                        4
                                                                     42.0
                                                                             Newvork
                                                        3
                                                               Nam
                                                                     29.0
                                                                               Hanoi
                                                        10
                                                               NaN
                                                                      NaN
                                                                                 NaN
  'Name':['Tom','Jack','Nam','Ricky',"Alice"],
                                                               NaN
                                                                      NaN
                                                                                 NaN
  'Age': [28,34,29,42,30],
                                                                               Place
                                                                      Age
  'Place':['Paris','London','Hanoi','Newyork',"LA"]
                                                             Ricky
                                                        4
                                                                     42.0
                                                                            Newvork
                                                        3
                                                               Nam
                                                                     29.0
                                                                               Hanoi
df1 = pd.DataFrame(data1, index=[1, 2, 3, 4, 5])
                                                        10
                                                                     30.0
                                                             Alice
                                                                                 LA
                                                               NaN
                                                                                 NaN
                                                              Name
                                                                      Age
                                                                               Place
data2 = {
    'Name':['David', 'Alice', 'Bob', "Kelvin"],
'Age':[29,24,20, 22],
'Place': ['Rio', 'LA', 'Berlin', "Rome"]
                                                             Ricky
                                                                     42.0
                                                                            Newyork
                                                        3
                                                                     29.0
                                                        1.0
                                                               NaN
                                                                      NaN
                                                                                 NaN
                                                               Tom
                                                                     28.0
                                                                               Paris
df2 = pd.DataFrame(data2, index=[4, 3, 10, -1])
                                                              Name
                                                                              Place
                                                                     Age
                                                        4
                                                             Ricky
                                                                      42
                                                                           Newyork
print(df1.reindex_like(df2))
                                                        3
                                                                       29
                                                               Nam
                                                                              Hanoi
print(df1.reindex_like(df2, method="ffill"))
print(df1.reindex_like(df2, method="bfill"))
                                                        1.0
                                                            Alice
                                                                      3.0
                                                                                 LA
print(df1.reindex_like(df2, method="nearest"))
                                                               Tom
                                                                       28
                                                                              Paris
                                                                                   71
```

```
Ví dụ: nhãn của hàng là kiểu String
                                                                                    Place
                                                                           Age
import pandas as pd
                                                                                     Paris
data1 = {
    'Name':['Tom', 'Jack', 'Nam', 'Ricky'],
'Age':[28,34,29,42],
'Place': ['Paris', 'London', 'Hanoi', 'Newyork']
                                                                       Jack
                                                                              34
                                                                                   London
                                                                       Nam
                                                                              29
                                                                                    Hanoi
df1 = pd.DataFrame(data1, index=['b', 'c', 'd', 'e'])
                                                                      Ricky
                                                                              42
                                                                                  Newyork
data2 = {
                                                                       Name
                                                                            Age Place
     'Name':['David', 'Alice', 'Bob', "Kelvin"],
    'Age':[29,24,20, 22],
'Place': ['Rio', 'LA'
                                                                       David
                                                                               29
                                                                                      Rio
                              'Berlin', "Rome"]
                                                                                      LA
df2 = pd.DataFrame(data2,
                             index=['e','c','f','a'])
                                                                        Bob
                                                                               20
                                                                                    Berlin
                                                                      Kelvin
                                                                                   Rome
print(df1.reindex_like(df2))
print(df1.reindex_like(df2, method="ffill"))
print(df1.reindex_like(df2, method="bfill"))
# print(df1.reindex_like(df2, method="nearest")) # Bi lõi néu dùng
                                                                                           72
```

```
Ví dụ: nhãn của hàng là kiểu String
import pandas as pd
                                                                    42.0
                                                           Ricky
                                                                           Newyork
                                                             Jack
                                                                    34.0
                                                                             London
data1 = {
    'Name':['Tom', 'Jack', 'Nam', 'Ricky'],
'Age':[28,34,29,42],
'Place': ['Paris', 'London', 'Hanoi', 'Newyork']
                                                              NaN
                                                                     NaN
                                                                                NaN
                                                             NaN
                                                                     NaN
                                                                                NaN
                                                                     Age
                                                                              Place
                                                           Ricky
                                                                    42.0
                                                                           Newvork
df1 = pd.DataFrame(data1, index=['b', 'c', 'd', 'e'])
                                                             Jack
                                                                    34.0
                                                                            London
                                                                    42.0
                                                           Rickv
                                                                           Newyork
                                                                     NaN
                                                             NaN
data2 = {
    'Name':['David', 'Alice', 'Bob', "Kelvin"],
                                                                     Age
                                                                              Place
    'Age':[29,24,20, 22],
'Place': ['Rio', 'LA',
                                                           Ricky
                                                                    42.0
                                                                           Newyork
                           'Berlin', "Rome"]
                                                            Jack 34.0
                                                                            London
                                                             NaN
                                                                    NaN
                                                                               NaN
df2 = pd.DataFrame(data2, index=['e','c','f','a'])
                                                                    28.0
                                                                              Paris
print(df1.reindex_like(df2))
print(df1.reindex_like(df2, method="ffill"))
print(df1.reindex_like(df2, method="bfill"))
# print(df1.reindex_like(df2, method="nearest")) # Bi loi néu dùng
                                                                                    73
```

```
Thiết lập giới hạn khi điền dữ liệu bổ sung
                                                                       Name Age
                                                                                    Place
import pandas as pd
                                                                        Tom
                                                                              28
                                                                                      Paris
    Jack
                                                                              34
                                                                                   London
                                                                       Nam
                                                                              29
                                                                                     Hanoi
                                                                    4 Ricky
                                                                               42 Newyork
df1 = pd.DataFrame(data1, index=[1, 2, 3, 4])
data2 = {
                                                                              Age
                                                                                   Place
                                                                        Name
    'Name':['David', 'Alice', 'Bob', "Kelvin"],
'Age':[29,24,20, 22],
    'Place': ['Rio', 'LA', 'Berlin', "Rome"]
                                                                         Alice
                                                                                24
                                                                                       LA
df2 = pd.DataFrame(data2, index=[-1, 4, 6, 7])
                                                                         Bob
                                                                                20
                                                                                     Berlin
print(df1.reindex_like(df2))
print(df1.reindex_like(df2, method="ffill", limit=1))
print(df1.reindex_like(df2, method="bfill", limit=1))
print(df1.reindex_like(df2, method="nearest", limit=1))
                                                                    7 Kelvin
                                                                                22
                                                                                     Rome
Lưu ý: Để dùng tham số limit, index của df2 phải đánh tăng dần.
                                                                                           74
Nếu index của df2 là [4, -1, 6, 7] sẽ bị lỗi
```

```
Thiết lập giới hạn khi điền dữ liệu bổ sung
                                                                   Name
                                                                                     Place
import pandas as pd
                                                                    NaN
                                                                            NaN
                                                                                      NaN
                                                                  Ricky
                                                                           42.0
                                                                                  Newyork
data1 = {
    'Name':['Tom', 'Jack', 'Nam', 'Ricky'],
'Age':[28,34,29,42],
'Place': ['Paris', 'London', 'Hanoi', 'Newyork']
                                                              6
                                                                    NaN
                                                                                       NaN
                                                                            NaN
                                                                    NaN
                                                                            NaN
                                                                                       NaN
                                                                            Age
                                                             -1
                                                                           NaN
                                                                                     NaN
                                                                    NaN
df1 = pd.DataFrame(data1, index=[1, 2, 3, 4])
                                                                  Ricky
                                                                           42.0
                                                                                  Newyork
                                                                           42.0
                                                                 Rickv
                                                                                  Newyork
data2 = {
                                                                    NaN
                                                                           NaN
                                                                                       NaN
     'Name':['David', 'Alice', 'Bob', "Kelvin"],
                                                                   Name
                                                                            Age
                                                                                     Place
     'Age':[29,24,20, 22],
    'Place': ['Rio', 'LA', 'Berlin', "Rome"]
                                                                    Tom
                                                                           28.0
                                                                                     Paris
                                                                 Ricky
                                                                           42.0
df2 = pd.DataFrame(data2, index=[-1, 4, 6, 7])
                                                              6
                                                                    NaN
                                                                           NaN
                                                                                       NaN
                                                                    NaN
                                                                            NaN
                                                                                       NaN
print(df1.reindex_like(df2))
                                                                   Name
                                                                                     Place
                                                                            Age
print(df1.reindex_like(df2, method="ffill", limit=1))
print(df1.reindex_like(df2, method="bfill", limit=1))
                                                                    Tom
                                                                           28.0
                                                                                     Paris
print(df1.reindex_like(df2, method="nearest", limit=1)) 4
                                                                  Ricky
                                                                           42.0
                                                                                  Newyork
                                                                 Ricky
                                                                           42.0
                                                                                  Newyork
                                                                    NaN
                                                                            NaN
                                                                                       NaN
Lưu ý: Đế dùng tham số limit, index của df2 phải đánh tăng dần.
                                                                                        75
```

Nếu index của df2 là [4, -1, 6, 7] sẽ bị lỗi

```
Thiết lập giới hạn khi điền dữ liệu bổ sung
                                                        Name
                                                                Age
                                                                       Place
import pandas as pd
                                                    -1
                                                         NaN
                                                                NaN
                                                                        NaN
                                                       Ricky
                                                               42.0
                                                                     Newyork
data1 = {
    'Name':['Tom', 'Jack
'Age':[28,34,29,42],
                'Jack', 'Nam', 'Ricky'],
                                                         NaN
                                                                NaN
                                                                         NaN
                                                         NaN
                                                                NaN
                                                                          NaN
    'Place': ['Paris', 'London', 'Hanoi', 'Newyork']
                                                        Name
                                                                        Place
                                                                Age
                                                   -1
                                                         NaN
                                                                        NaN
                                                                NaN
df1 = pd.DataFrame(data1, index=[1, 2, 3, 4])
                                                    4
                                                       Ricky
                                                               42.0
                                                                     Newyork
                                                    6
                                                       Ricky
                                                               42.0
                                                                     Newyork
data2 = {
                                                       Ricky
                                                               42.0
                                                                     Newyork
    'Name':['David', 'Alice', 'Bob', "Kelvin"],
'Age':[29,24,20, 22],
                                                                       Place
                                                        Name
                                                               Age
                                                               28.0
   'Place': ['Rio', 'LA', 'Berlin', "Rome"]
                                                   -1
                                                         Tom
                                                                       Paris
                                                       Ricky
                                                               42.0
                                                                     Newyork
df2 = pd.DataFrame(data2, index=[-1, 4, 6, 7])
                                                         NaN
                                                                NaN
                                                                         NaN
                                                         NaN
                                                                NaN
                                                                          NaN
print(df1.reindex_like(df2))
                                                        Name
                                                               Age
                                                                      Place
Tom
                                                                28
                                                                      Paris
                                                       Ricky
                                                                42
                                                                    Newyork
                                                       Ricky
                                                    6
                                                                    Newyork
                                                                42
                                                                42
                                                       Ricky
                                                                    Newyork
                                                                          76
```

```
Thiết lập giới hạn khi điền dữ liệu bổ sung
import pandas as pd
                                                                             Name
                                                                                   Age
                                                                                            Place
                                                                                             Paris
data1 = {
                                                                              Tom
                                                                                      28
     'Name':['Tom', 'Jack', 'Nam', 'Ricky'],
'Age':[28,34,29,42],
'Place': ['Paris', 'London', 'Hanoi', 'Newyork']
                                                                                            London
                                                                              Nam
                                                                                      29
                                                                                             Hanoi
df1 = pd.DataFrame(data1, index=['b', 'c', 'd', 'e'])
                                                                          e Ricky
                                                                                      42 Newyork
     'Name':['David', 'Alice', 'Bob', "Kelvin"],
                                                                              Name
                                                                                    Age Place
     'Age':[29,24,20, 22],
'Place': ['Rio', 'LA', 'Berlin', "Rome"]
                                                                          a David
                                                                                       29
                                                                                               Rio
df2 = pd.DataFrame(data2, index=['a', 'c', 'y', 'z'])
                                                                               Alice
                                                                                               LA
print(df1.reindex_like(df2))
                                                                                Bob
                                                                                       20
                                                                                            Berlin
print(df1.reindex_like(df2, method="ffill", limit=1))
print(df1.reindex_like(df2, method="bfill", limit=1))
                                                                          z Kelvin
                                                                                       22
                                                                                            Rome
                                                                                                    77
```

```
Thiết lập giới hạn khi điền dữ liệu bổ sung
                                                                Name
                                                                         Age
                                                                                 Place
import pandas as pd
                                                                 NaN
                                                                         NaN
                                                                                   NaN
                                                             С
                                                                 Jack
                                                                        34.0
                                                                                London
data1 = {
    1 = {
'Name':['Tom', 'Jack', 'Nam', 'Ricky'],
'Age':[28,34,29,42],
                                                                  NaN
                                                                         NaN
                                                                                   NaN
                                                                  NaN
                                                                         NaN
                                                                                    NaN
    'Place': ['Paris', 'London', 'Hanoi', 'Newyork']
                                                                  Name
                                                                                    Place
                                                                  NaN
                                                                                    NaN
                                                                          NaN
df1 = pd.DataFrame(data1, index=['b', 'c', 'd', 'e'])
                                                                  Jack
                                                                         34.0
                                                                                  London
                                                                Ricky
                                                                         42.0
                                                                                 Newyork
                                                             У
data2 = {
    'Name':['David', 'Alice', 'Bob', "Kelvin"],
'Age':[29,24,20, 22],
                                                                  NaN
                                                                          NaN
                                                                         Age
    'Place': ['Rio', 'LA', 'Berlin', "Rome"]
                                                                 Tom
                                                                        28.0
                                                                                 Paris
                                                                 Jack
                                                                        34.0
                                                                                London
df2 = pd.DataFrame(data2, index=['a', 'c', 'y', 'z'])
                                                                 NaN
                                                                         NaN
                                                                                   NaN
print(df1.reindex_like(df2))
print(df1.reindex_like(df2, method="ffill", limit=1))
print(df1.reindex_like(df2, method="bfill", limit=1))
                                                                                        78
```