Machine Translated by Google Bảng tóm tắt Python cho khoa học dữ liệu

Khái niêm cơ bản về gấu trúc

Tìm hiểu Python cho Khoa học dữ liệu một cách tương tác tại www.DataCamp.com



gấu trúc

Thư viên Pandas được xây dự ng trên NumPy và cung cấp các cấu trúc dữ liêu và công cụ phân tích dữ liêu dễ sử dụng cho ngôn ngữ lập trình Python.

Sử dụng các quy ước nhập sau: >>> nhập gấu trúc dưới dạng pd

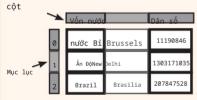
Cấu trúc dữ liêu Pandas

Mảng được gắn nhãn một chiều có khả năng chứa bất kỳ loại dữ liệu nào.



>>> s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])

Khuna dữ liêu



Cấu trúc dữ liêu được gắn nhãn hai chiều với các côt có thể có các kiểu khác nhau

```
>>> data = {'Quốc gia': ['Bỉ', 'Ấn Độ', 'Brazil'],
                  'Thủ đô': ['Brussels', 'New Delhi', 'Brasília'],
                  'Dân số': [11190846, 1303171035, 207847528]}
>>> df = pd.DataFrame(dí? liêu
                                 côt=['Quốc qia', 'Thủ đô', 'Dân số'])
```

Yêu cầu giúp đỡ

>>> trợ giúp(pd.Series.loc)

Lu a chon

```
>>> s['b']
                                                          Lấy một phần tử
                                                          Nhận tập hợp con của DataFrame
>>> df[1:1
      Ouốc nia
                       Dân số vốn
   1 Ấn Độ New Delhi 1303171035
   2 Brazil Brasilia 207847528
```

Chon, lập chỉ mục Boolean và thiết lập

```
theo vi trí
>>> df.iloc[[0],[0]]
                                                    Chọn một giá trị theo hàng
                                                    và côt
  'Nước Bỉ
 >> df.iat([0],[0])
  'Nurric Bi
 Theo nhãn
>>> df.loc[[0], ['Quốc gia']]
                                                    Chon một giá tri theo nhãn
                                                    hàng và cột
   'Nước Bỉ
```

Theo Nhãn/Vi trí

'Nước Bỉ

>>> df.at([0], ['Quốc gia'])

```
>>> df.ix[2]
                                                       Chon một hàng của tập
                                                       hợp con các hàng
                   Brazil
  Ouốc gia
                Brasilia
  Thủ đô
  Dân số 207847528
                                                       Chọn một cột trong tập
>>> df.ix[:,'Capital']
                                                       hợp con của các cột
          Bruxelles
        New Delhi
          Brasilia
```

>>> df.ix[1,'Vốn'] 'New Delhi

Lä	ĝр	chỉ	ı	nục	Boolea
>>>	s [~(s	>	1)]	>>>

s[(s < -1) | (s > 2)] >>>

df[df['Population']>1200000000] Sử dụng bộ lọc để điều chỉnh DataFrame

Cài đặt

>>> s['a'] = 6

Chọn hàng và cột

Chuỗi s có giá trị không >1

s trong đó giá tri là <-1 hoặc >2

Đặt chỉ mục a của Series s thành 6

Đọc và ghi vào CSV

>>> pd.read csv('file.csv', header=None, nrows=5)

Đọc và ghi vào Excel

>>> df.to csv('myDataFrame.csv')

>>> df = pd.read_excel(xlsx, 'Sheet1')

```
>>> pd.read_excel('file.xlsx')
>>> pd.to excel('dir/myDataFrame.xlsx', sheet name='Sheet1')
 Đọc nhiều trang từ cùng một tập tin
>>> xlsx = pd.ExcelFile('file.xls')
```

Đọc và ghi vào truy vấn SQL hoặc bảng cơ sở dữ liệu

```
>>> từ nhập sqlalchemy create_engine
>>> engine = create_engine('sqlite:///:memory:')
>>> pd.read_sql("CHON * Từ my_table;", công cụ)
>>> pd.read_sql_table('my_table', engine)
>>> pd.read_sql_query("CHQN * Từ my_table;", công cụ)
read_sql() là một trình bao bọc tiện lợi xung quanh read_sql_table() và
read_sql_query()
>>> pd.to_sql('myDf', động cơ)
```

Rơ i

```
Bổ giá trị từ các hàng (axis=0)
>>> s.drop(['a', 'c']) >>>
df.drop('Country', axis=1) Bổ giá trị khổi cột(axis=1)
```

Sắp xếp & Xếp hạng

```
>>> df.sort index() >>>
                                                     Sắp xếp theo nhãn dọc theo một trục
df.sort_values(by='Country') Sắp xếp theo các qiá trị dọc theo một trục
>>> df.rank()
                                                    Chỉ định thứ hang cho các mục
```

Truv xuất thông tin chuỗi/DataFrame

Thông tin cơ bản

```
>>> df.shape >>>
                                         (những hàng, những cột)
df.index >>>
                                        Mô tả chỉ số
df.columns >>>
                                        Mô tả các côt DataFrame
df.info() >>>
                                        Thông tin về DataFrame
df.count()
                                        Số lương giá tri không phải NA
```

Bản tóm tắt

```
Tổng các giá
tri >>> df.sum()
                                             Tổng tích lũ y của các giá
tri >>> df.cumsum() >>>
                                             Giá tri tối thiểu/tối đa
df.min()/df.max()
>>> df.idxmin()/df.idxmax() Giá tri chỉ
                                           muc tối thiểu/tối đa >>>
df.describe()
                                             thống kê tóm tắt
>>> df.mean() >>>
                                             Ý nghĩa của các giá trị
df.median()
                                             Giá tri trung bình
```

Áp dung hàm

```
>>> f = lambda x: x*2
>>> df.apply(f) >>>
                                      Áp dụng chức năng
df.applymap(f)
                                      Áp dụng từng phần tử hàm
```

Căn chỉnh dữ liêu

Căn chỉnh dữ liệu nội bộ

Giá tri NA được đưa vào trong các chỉ số không trùng nhau:

```
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])
>>>s+s3
         10.0
         NaN
         5.0
```

Bạn cũ ng có thể tự mình thực hiện việc căn chỉnh dữ liệu nội bộ với sự trợ giúp của các phươ ng thức điền:

```
>>> s.add(s3. fill value=0)
 b
         -5.0
  С
        5.0
  А
        7 0
>>> s.sub(s3, fill value=2)
>>> s.div(s3, fill_value=4)
>>> s.mul(s3, fill_value=3)
```

