Bài giảng:

Lập trình Python

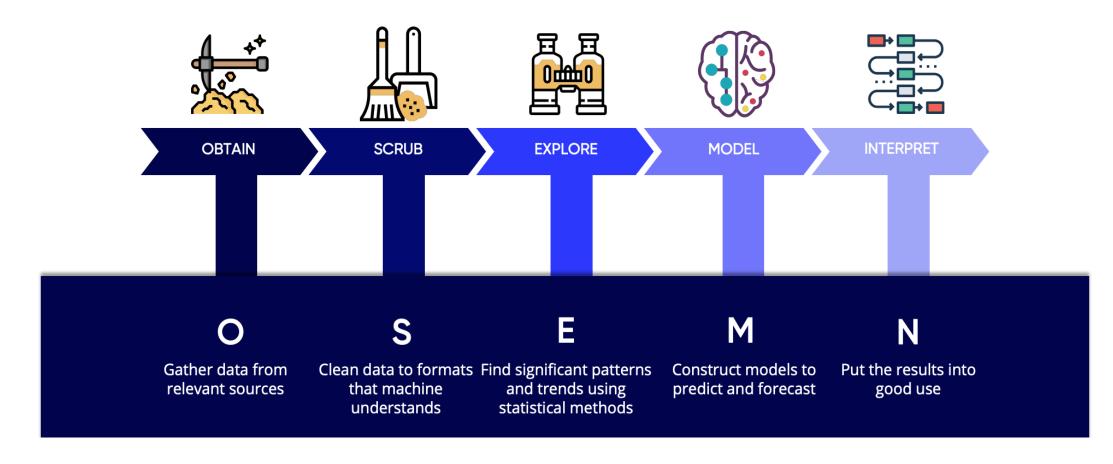
Bài 06: Làm việc với các nguồn dữ liệu khác nhau

Nội dung bài 6

- 1. Giới thiệu
- 2. Đọc các file dữ liệu với Pandas:
 - Giới thiệu Pandas
 - > Đọc dữ liệu từ file CSV, Excel, Json
- 3. Lấy dữ liệu từ API
- 4. Làm việc với dữ liệu từ MongoDB
- 5. Sử dụng sklearn đọc tập dữ liệu mẫu
- 6. Tạo Simulated Dataset với sklearn

1. Giới thiệu

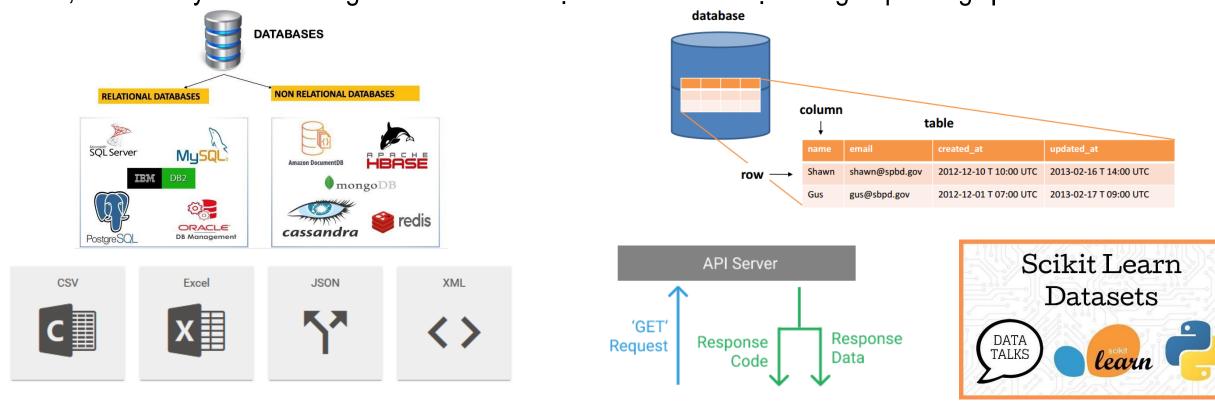
1. Giới thiệu



Để có thể phân tích dữ liệu cho bất kỳ mục đích gì, bước đầu tiên bao giờ cũngphải thu thập và đọc dữ liệu từ tất cả các nguồn, các dạng khác nhau

1. Giới thiệu

Dữ liệu tồn tại trong rất nhiều dạng khác nhau, từ dữ liệu trong các file cơ bản như Excel, CSV, text, Json, XML...hay lưu trữ trong các CSDL. Dữ liệu còn có thể được cung cấp thông qua các API.



Trong bài học này chúng ta sẽ tìm hiểu cách đọc dữ liệu từ một số dạng cơ bản và phổ

biến nhất.

2. Đọc dữ liệu với Pandas

2.1 Giới thiệu Pandas

- Pandas là một thư viện phần mềm viết cho ngôn ngữ lập trình Python. Nó được sử dụngcho thao tác và phân tích dữ liệu. Nó cung cấp cấu trúc dữ liệu đặc biệt và hoạt động chocác thao tác của các bảng số liệu và chuỗi thời gian. Pandas là thư viện miễn phí được phát hành theo giấy phép BSD.
- Python với Pandas được sử dụng trong nhiều lĩnh vực cả lĩnh vực học thuật và thương mại như: Tài chính, kinh tế, thống kê, phân tích, v.v.

Top 5 Python Libraries for Data Science puthon Report from Cloud Academy suggests that the top technical skill in demand for data engineers is python. 67 percent of job posts mentioned python. 1.) NUMPY Through NumPy, you can use it as an efficient multi-dimensional container of generic data. It also contains sophisticated (broadcasting) functions, tools for integrating C/C++ and Fortran code, useful linear algebra, Fourier transform, and random number capabilities

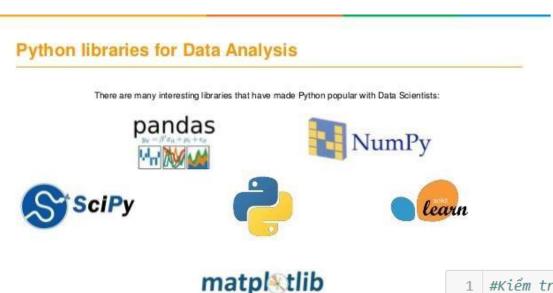
2.) PANDAS

and columns

Pandas provides high-performance, easy-to-use data structures and data

analysis tools for python. You can

store and manage data from tables by performing manipulation over rows



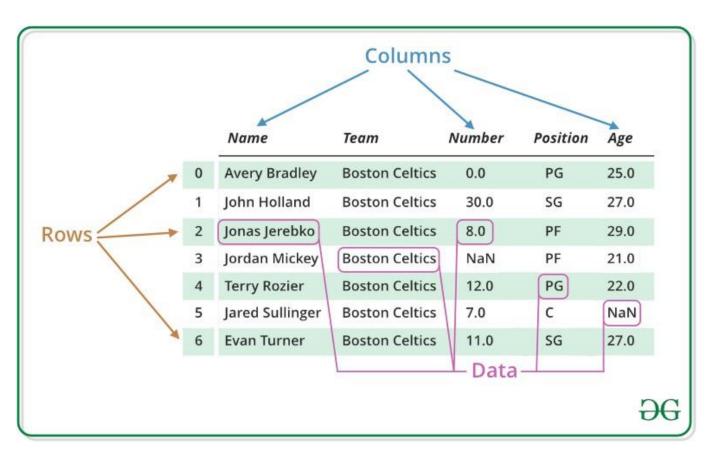
#Kiểm tra phiên bản của thư viện Pandas import pandas as pd

print('Version Pandas: ',pd. version

OSimplifearn. All rights reserved.

2.1 Giới thiệu Pandas

- Pandas có ba cấu trúc dữ liệu và nó được xây dựng dựa trên thư viện Numpy vậynên chúng hoạt động rất nhanh và hiệu quả: **Series, DataFrame, Panel**.
- Trong ba kiểu dữ liệu, DataFrame là kiểu dữ liệu được sử dụng rộng rãi nhất



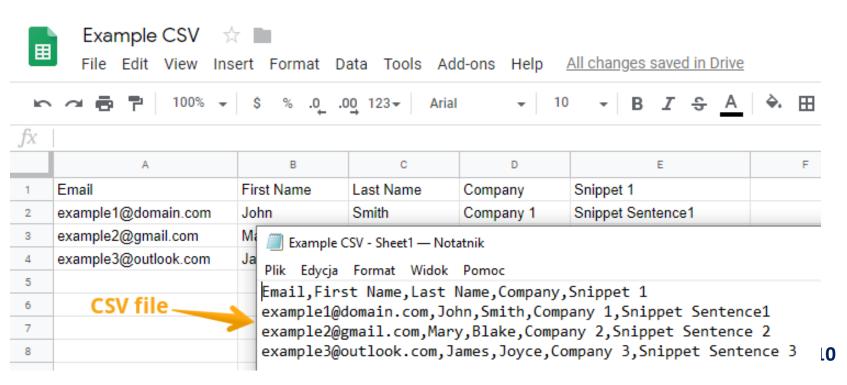
DataFrame: Cấu trúc dạng bảng 2D, kích thước có thể thay đổi được. Dữ liệu một cột là đồng nhất nhưng có thể không đồngnhất giữa các cột

2.2 Đọc dữ liệu từ file CSV

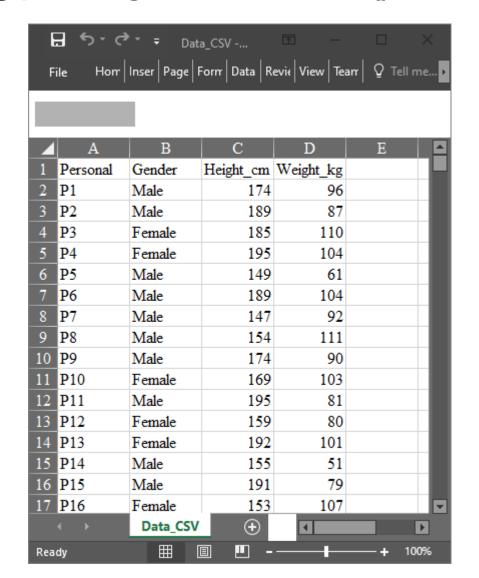
2.2 Doc file CSV

- CSV là một định dạng dữ liệu văn bản đơn giản có tên đầy đủ là Comma Separated Values. Với định dạng CSV này, các giá trị được chia tách với nhau bởi các dấu phẩy. Định dạng CSV phổ biến bởi vì chúng có tính tương thích cao, dễ dàng dichuyển từ phần mềm này sang phần mềm khác để sử dụng mà không lo gặp các xung đột.
- Tài liệu CSV cũng làm một trong những tài liệu phổ biến trên thế giới với khả nănglưu trữ nhỏ nhẹ.



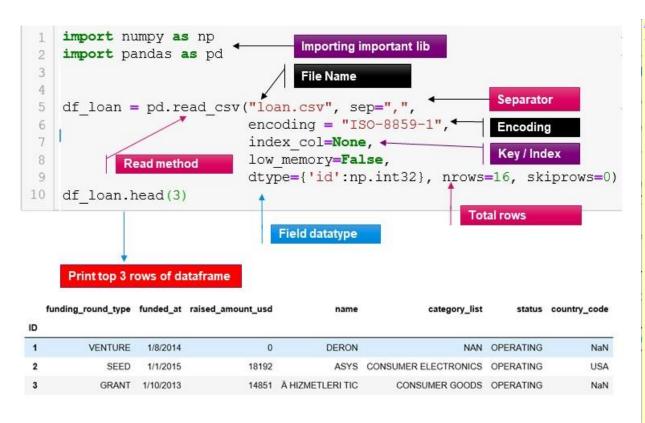


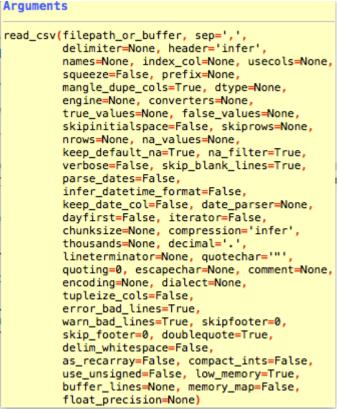
Sử dụng phương thức read_csv() đọc dữ liệu từ file .CSV



```
import pandas as pd
    path = 'Data Excersice\CSV\Data CSV.csv'
    #Sử dụng phương thức read csv
    data = pd.read csv(path)
    #Hiển thị thông tin biến Data
    data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 4 columns):
    Column
               Non-Null Count Dtype
    Personal
                                object
                500 non-null
                500 non-null
                                object
    Gender
    Height cm 500 non-null
                                int64
    Weight kg 500 non-null
                                int64
dtypes: int64(2), object(2)
memory usage: 15.8+ KB
```

Sử dụng phương thức read_csv() có rất nhiều tham số khác nhau để thiết lập cách thức đọc file .csv





Vd1: sử dụng tham số index_col để thiết lập cột index khi đọc file csv

```
1 #Sử dụng phương thức read csv()
    #Tham số: Thiết lập cột index là cột Personal
    data1 = pd.read_csv(path,
                       index col=0)
 4
    data1.info()
<class 'pandas.core.frame.DataFrame'>
Index: 500 entries, P1 to P500
Data columns (total 3 columns):
            Non-Null Count Dtype
    Column
#
   Gender 500 non-null object
    Height_cm 500 non-null
                               int64
    Weight kg 500 non-null
                               int64
dtypes: int64(2), object(1)
memory usage: 15.6+ KB
```

1	#Hiển thị dữ liệu	5	dòng	đầu	tiên
2	data1.head()				

	Gender	Height_cm	Weight_kg
Personal			
P1	Male	174	96
P2	Male	189	87
P3	Female	185	110
P4	Female	195	104
P5	Male	149	61

dtypes: int64(2)

memory usage: 1.7 KB

Vd2: Thiết lập tham số chỉ đọc 100 dòng đầu tiên và dữ liệu trong 2 cột Height_cm, Weight_kg

```
#Sử dụng phương thức read csv()
    #Thiết lập số hàng, cột muốn đọc dữ liệu
    data2 = pd.read csv(path,
                        nrows=100,
 4
                       usecols=['Height cm', 'Weight kg'])
    data2.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100 entries, 0 to 99
Data columns (total 2 columns):
               Non-Null Count Dtype
    Column
    Height cm 100 non-null
                               int64
    Weight kg 100 non-null
                               int64
```

```
1 #Hiển thị dữ liệu 5 dòng đầu tiên
2 data2.head()
```

	Height_cm	Weight_kg
0	174	96
1	189	87
2	185	110
3	195	104
4	149	61

2.2 Doc file CSV

Vd3: Thiết lập tham số đọc dữ liệu từ dòng thứ 5 trở đi, và đặt lại tên của từng cộtdữ liệu thành ['ID','Sex','H(cm)',W(kg)']

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 496 entries, 0 to 495
Data columns (total 4 columns):
     Column Non-Null Count Dtype
                          object
     ID
            496 non-null
            496 non-null
                            object
    Sex
            496 non-null
                            int64
    H(cm)
     W(kg)
            496 non-null
                            int64
dtypes: int64(2), object(2)
memory usage: 15.6+ KB
```

```
#Hiển thị 5 dòng dữ liệu đầu tiên
data3.head()
```

	ID	Sex	H(cm)	W(kg)
0	P5	Male	149	61
1	P6	Male	189	104
2	P7	Male	147	92
3	P8	Male	154	111
4	P9	Male	174	90



Yêu cầu 2.2.1: Sinh viên đọc dữ liệu dạng CSV lưu trong file csv_Data_Loan.csv với các tham số mặc định

A	В	С	D	Е	F	G	H	I	J	K	L	M	N O
loan_amnt		int_rate	emp_length	home_ownership	annual_inc	purpose	addr_state	đti	delinq_2yrs re	vol_util	total_acc	bad_loan	longest_cre verification_statu
	36 months	10.65	10	RENT	24000	credit_card	AZ	27.65	0	83.7	9	0	26 verified
2500	60 months	15.27		RENT	30000	car	GA	1	0	9.4		1	12 verified
55.0.00 St. 2000.00	36 months	15.96	10	RENT	12252	small_business	IL	8.72	0	98.5			10 not verified
10000	36 months	13.49	10	RENT	49200	other	CA	20	0	21	37		15 verified
5000	36 months	7.9	3	RENT	36000	wedding	AZ	11.2	0	28.3	12	0	7 verified
D-010101	36 months	18.64	9	RENT	48000	car	CA	5.35	0	87.5	4	0	4 verified
	60 months	21.28	4	OWN	40000	small_business	CA	5.55	0	32.6	13	1	7 verified
5375	60 months	12.69	0	RENT	15000	other	TX	18.08	0	36.5	3	1	7 verified
	60 months	14.65	5	OWN	72000	debt_consolidation	AZ	16.12	0	20.6	23	C	13 not verified
1 12000	36 months	12.69	10	OWN	75000	debt_consolidation	CA	10.78	0	67.1	34	0	22 verified
	36 months	13.49	0	RENT	30000	debt_consolidation	VA	10.08	0	91.7	9	1	7 verified
_	36 months	9.91	3	RENT	15000	credit_card	IL	12.56	0	43.1	11	C	8 verified
57 (53000000000000000000000000000000000000	36 months	10.65	3	RENT	100000	other	CA	7.06	0	55.5	29	1	20 verified
3000000000	36 months	16.29	0	RENT	28000	debt_consolidation	MO	20.31	0	81.5	23	0	4 not verified
	36 months	15.27	4	RENT	42000	home_improvement	CA	18.6	0	70.2	28	C	13 not verified
7 3600	36 months	6.03	10	MORTGAGE	110000	major_purchase	CT	10.52	0	16	42	0	18 not verified
100000000000000000000000000000000000000	36 months	11.71	1	MORTGAGE	84000	medical	UT	18.44	2	37.73	14	0	8 verified
The state of the s	36 months	6.03	6	RENT	77385.19	debt_consolidation	CA	9.86	0	23.1	28	C	10 not verified
85.62.222.200.000	36 months	12.42	10	RENT	105000	debt_consolidation	FL	13.22	0	90.3	38	1	28 verified
	36 months	11.71	10	OWN	50000	credit_card	TX	11.18	0	82.4	21	C	26 verified
UNIVERSAL SAME	36 months	11.71	5	RENT	50000	debt_consolidation	CA	16.01	0	91.8	17	C	8 not verified
	36 months	11.71	1	RENT	76000	major_purchase	CA	2.4	0	29.7	7	1	10 not verified
15000	36 months	9.91	2	MORTGAGE	92000	credit_card	IL	29.44	0	93.9	31	C	9 verified
	Data_Loa	n \oplus							(1)				



Yêu cầu 2.2.2: Đọc dữ liệu từ file Data_Loan.CSV vào 2 biến DataFrame tương ứng.

- df_number: Chỉ chứa các cột dữ liệu số
- **df_object**: Chỉ chứa các cột dữ liệu Object
- 1 #Hiển thị 5 dòng dữ liệu đầu tiên của biến df_number
- 2 df_number.head()

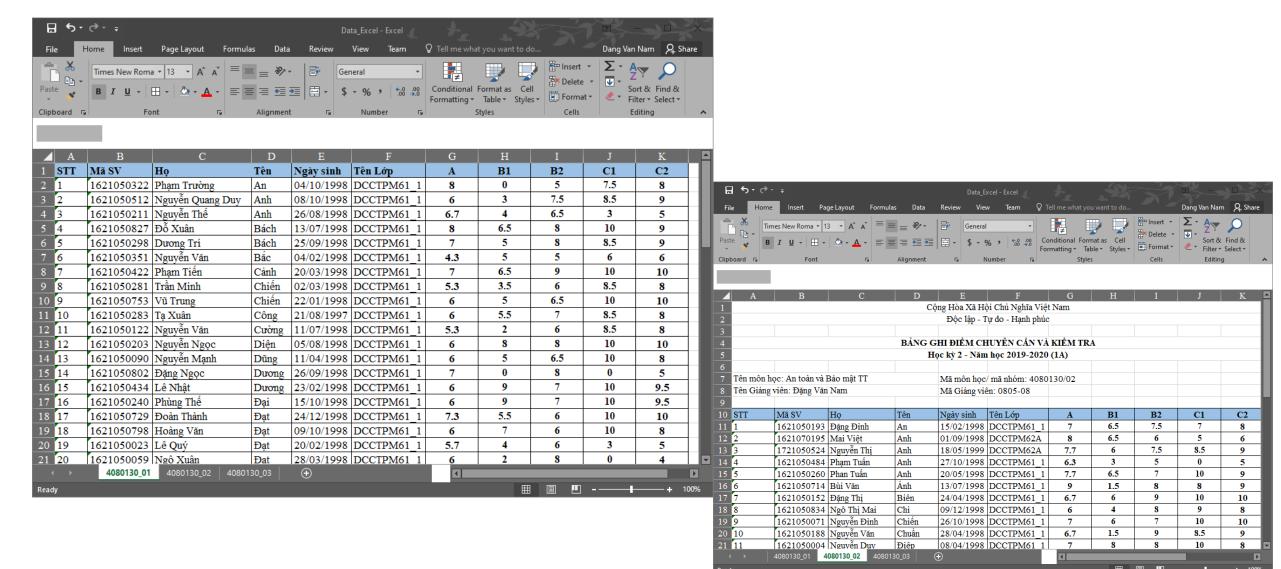
	loan_amnt	int_rate	emp_length	annual_inc	dti	delinq_2yrs	revol_util	total_acc	bad_loar	n longest_c	redit_length	_											
0	5000	10.65	10.0	24000.0	27.65	0.0	83.7	9.0	(0	26.0												
1	2500	15.27	0.0	30000.0	1.00	0.0	9.4	4.0		1	12.0												
2	2400	15.96	10.0	12252.0	8.72	0.0	98.5	10.0		-	_	ê	ệu đầu tiên	ệu đầu tiên của biến	ệu đầu tiên của biến df	ệu đầu tiên của biến df_ok	ệu đầu tiên của biến df_obje	ệu đầu tiên của biến df_objec	ệu đầu tiên của biến df_object	ệu đầu tiên của biến df_object	ệu đầu tiên của biến df_object	ệu đầu tiên của biến df_object	ệu đầu tiên của biến df_object
3	10000	13.49	10.0	49200.0	20.00	0.0	21.0	37.0	2 0	df_object.h	nead()												
4	5000	7.90	3.0	36000.0	11.20	0.0	28.3	12.0		term hon	ne_ownership		purpose	purpose addr_state	purpose addr_state v	purpose addr_state verifi	purpose addr_state verificat	purpose addr_state verification	purpose addr_state verification	purpose addr_state verification_	purpose addr_state verification_s	purpose addr_state verification_str	purpose addr_state verification_statu

	term	home_ownership	purpose	addr_state	verification_status
0	36 months	RENT	credit_card	AZ	verified
1	60 months	RENT	car	GA	verified
2	36 months	RENT	small_business	IL	not verified
3	36 months	RENT	other	CA	verified
4	36 months	RENT	wedding	AZ	verified

2.3 Đọc dữ liệu từ file Excel

2.3 Doc file Excel

• File dữ liệu Excel demo gồm 3 sheet:



memory usage: 5.8+ KB

- Sử dụng phương thức **pd.read_excel()** để đọc dữ liệu từ file excel.
 - Lưu ý 2 tham số sheetname=" xác định sheet muốn đọc dữ liệu (Mặc định là sheet đầu tiên)

```
import pandas as pd
    path excel = 'Data Excersice\Data Excel.xlsx'
 3 #Đọc dữ liệu từ file excel
    data ex = pd.read excel(path excel)
    data ex.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 66 entries, 0 to 65
Data columns (total 11 columns):
     Column
                 Non-Null Count Dtype
     STT
                 66 non-null
                                  int64
                                                     #Hiển thị 5 dòng dữ liệu đầu tiên
     Mã SV
                 66 non-null
                                  int64
                                                     data ex.head()
                 66 non-null
                                  object
     Но
     Tên
                 66 non-null
                                  object
     Ngày sinh
                                  object
                 66 non-null
                                                    STT
                                                            Mã SV
                                                                                  Tên Ngày sinh
                                                                                                   Tên Lớp A B1 B2
                                                                                                                      C1 C2
     Tên Lớp
                 66 non-null
                                  object
                                                      1 1621050322
                                                                      Pham Trường
                                                                                   An 04/10/1998 DCCTPM61 1 8.0 0.0 5.0 7.5 8.0
                 66 non-null
                                  float64
                                                      2 1621050512 Nguyễn Quang Duy
                                                                                  Anh 08/10/1998 DCCTPM61 1 6.0 3.0 7.5
                 66 non-null
                                  float64
     В1
                 66 non-null
                                  float64
     В2
                                                      3 1621050211
                                                                       Nguyễn Thế
                                                                                  Anh 26/08/1998 DCCTPM61 1 6.7 4.0 6.5
     C1
                 66 non-null
                                  float64
                                                                         Đỗ Xuân Bách 13/07/1998 DCCTPM61_1 8.0 6.5 8.0 10.0 9.0
                                                      4 1621050827
     C2
                 66 non-null
                                  float64
                                                                         Durong Trí Bách 25/09/1998 DCCTPM61 1 7.0 5.0 8.0 8.5 9.0
                                                      5 1621050298
dtypes: float64(5), int64(2), object(4)
```

• Một vài tham số quan trọng trong phương thức **pd.read_excel()** để đọc dữ liệu từ file excel.

Argument	Description
io	A string containing the pathname of the given Excel file.
sheet_name	The Excel sheet name, or sheet number, of the data you want to import. The sheet number can be an integer where 0 is the first sheet, 1 is the second, etc. If a list of sheet names/numbers are given, then the output will be a dictionary of DataFrames. The default is to read all the sheets and output a dictionary of DataFrames.
header	Row number to use for the list of column labels. The default is 0, indicating that the first row is assumed to contain the column labels. If the data does not have a row of column labels, None should be used.
names	A separate Python list input of column names. This option is None by default. This option is the equivalent of assigning a list of column names to the columns attribute of the output DataFrame.
index_col	Specifies which column should be used for row indices. The default option is None, meaning that all columns are included in the data, and a range of numbers is used as the row indices.
usecols	An integer, list of integers, or string that specifies the columns to be imported into the DataFrame. The default is to import all columns. If a string is given, then Pandas uses the standard Excel format to select columns (e.g. "A:C,F,G" will import columns A, B, C, F, and G).
skiprows	The number of rows to skip at the top of the Excel sheet. Default is 0. This option is useful for skipping rows in Excel that contain explanatory information about the data below it.

• Sử dụng phương thức **pd.read_excel()** với một số tham số cơ bản.

```
Int64Index: 66 entries, 1621050322 to 1621050013
Data columns (total 5 columns):
    Column Non-Null Count Dtype
            66 non-null
                           float64
            66 non-null
                           float64
    В1
        66 non-null
                        float64
    B2
    C1
            66 non-null float64
            66 non-null
                           float64
dtypes: float64(5)
memory usage: 3.1 KB
```

```
1 #Hiển thị 5 dòng dữ liệu đầu tiên
2 data_ex1.head()
```

	Α	B1	B2	C1	C2
Mã SV					
1621050322	8.0	0.0	5.0	7.5	8.0
1621050512	6.0	3.0	7.5	8.5	9.0
1621050211	6.7	4.0	6.5	3.0	5.0
1621050827	8.0	6.5	8.0	10.0	9.0
1621050298	7.0	5.0	8.0	8.5	9.0

- Sử dụng phương thức **pd.read_excel()** với một số tham số cơ bản.
 - Đọc dữ liệu sheet 2 ['4080130_02'], từ dòng 9.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 39 entries, 0 to 38
Data columns (total 11 columns):
    Column
               Non-Null Count Dtype
               39 non-null
    STT
                               int64
    Mã SV
               39 non-null
                               int64
                               object
               39 non-null
    Ηo
    Tên
               39 non-null
                               object
    Ngày sinh 39 non-null
                               object
               39 non-null
                               object
    Tên Lớp
               39 non-null
                               float64
                              float64
               39 non-null
               39 non-null
                             float64
               39 non-null
                               float64
                               float64
               39 non-null
dtypes: float64(5), int64(2), object(4)
memory usage: 3.5+ KB
```

1 #Hiển thị 5 dòng dữ liệu đầu tiên 2 data_ex3.head()

	STT	Mã SV	Họ	Tên	Ngày sinh	Tên Lớp	Α	B1	B2	C1	C2
0	1	1621050193	Đặng Đình	An	15/02/1998	DCCTPM61_1	7.0	6.5	7.5	7.0	8.0
1	2	1621070195	Mai Việt	Anh	01/09/1998	DCCTPM62A	8.0	6.5	6.0	5.0	6.0
2	3	1721050524	Nguyễn Thị	Anh	18/05/1999	DCCTPM62A	7.7	6.0	7.5	8.5	9.0
3	4	1621050484	Phạm Tuấn	Anh	27/10/1998	DCCTPM61_1	6.3	3.0	5.0	0.0	5.0
4	5	1621050260	Phan Tuấn	Anh	20/05/1998	DCCTPM61_1	7.7	6.5	7.0	10.0	9.0

- Sử dụng phương thức **pd.read_excel()** với một số tham số cơ bản.
 - Đọc dữ liệu sheet 3 ['4080130_03'], không có dòng header

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 39 entries, 0 to 38
Data columns (total 11 columns):
     Column Non-Null Count Dtype
             39 non-null
                             int64
             39 non-null
                             int64
             39 non-null
                             object
             39 non-null
                             object
             39 non-null
                             object
             39 non-null
                             object
             39 non-null
                             float64
                             float64
             39 non-null
                             float64
             39 non-null
             39 non-null
                             float64
             39 non-null
                             float64
dtypes: float64(5), int64(2), object(4)
memory usage: 3.5+ KB
```

```
#Hiển thi 5 dòng dữ liệu đầu tiên
    data ex4.head()
  0
             1
                                                                        9 10
                 Đào Tuấn
0 1 1621050041
                             Anh 22/10/1998 DCCTPM61_1 6.7 9.0 5.5
                                                                      8.5 8.0
                                 26/09/1998 DCCTPM61_1 6.7 7.0 9.0
     1621050262
                Vũ Thi Lan
                                                                      8.5 6.0
     1621050083
                 Trinh Như
                                 06/04/1998 DCCTPM61 1 7.3 8.5 9.5
     1621050113
                 Trần Văn
                                 19/06/1998 DCCTPM61 1 5.7 5.0 6.0
                          Cương
                Nguyễn Sỹ
     1621050384
                            Düng 02/10/1998 DCCTPM61 1 7.0 0.0 7.5
```

- Sử dụng phương thức **pd.read_excel()** với một số tham số cơ bản.
 - Đọc dữ liệu sheet 3 ['4080130_03'],
 - Không có dòng header

memory usage: 1.8 KB

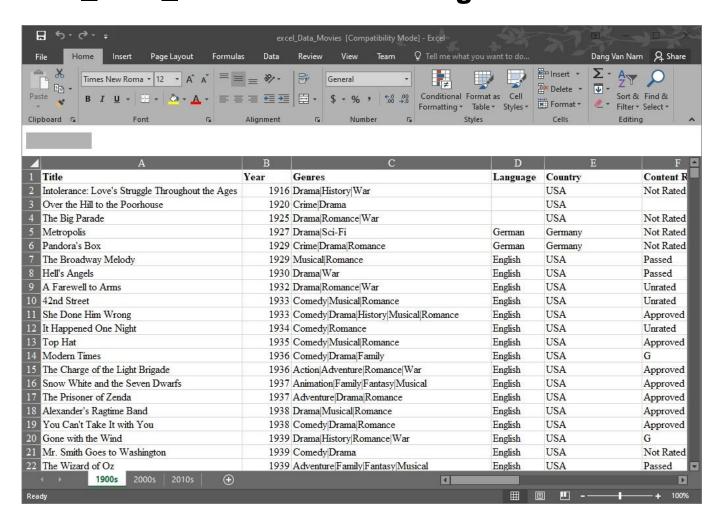
- Chỉ lấy dữ liệu cột 1,6,7,8,910
- Đặt tên cho các cột lần lượt là ['Mã SV', 'A', 'B1','B2','C1','C2']
- Thiết lập cột đầu tiên làm Index

```
1 #Hiển thị 5 dòng dữ liệu đầu tiên
2 data_ex41.head()
```

	Α	В1	B2	C1	C2
Mã SV					
1621050041	6.7	9.0	5.5	8.5	8.0
1621050262	6.7	7.0	9.0	8.5	6.0
1621050083	7.3	8.5	9.5	10.0	9.0
1621050113	5.7	5.0	6.0	10.0	5.0
1621050384	7.0	0.0	7.5	8.5	9.0



Yêu cầu: Sinh viên đọc dữ liệu dạng excel lưu trong file excel_Data_Movies.xls theo từng sheet



2.4 Đọc dữ liệu từ file Json

2.4 Đọc file Json

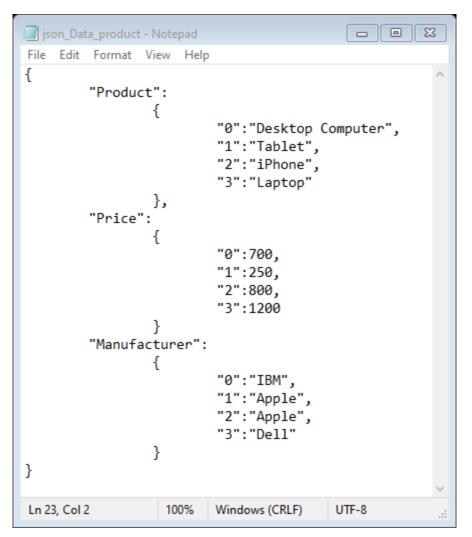
- JSON là chữ viết tắt của Javascript Object Notation, đây là một dạng dữ liệu tuân theo một quy luật nhất định mà hầu hết các ngôn ngữ lập trình đều có thể đọc được.
- JSON lưu trữ các dữ liệu theo cặp khóa (key) và giá trị (value). So với XML thì JSON có định dạng đơn giản, dễ sử dụng và nhẹ hơn.

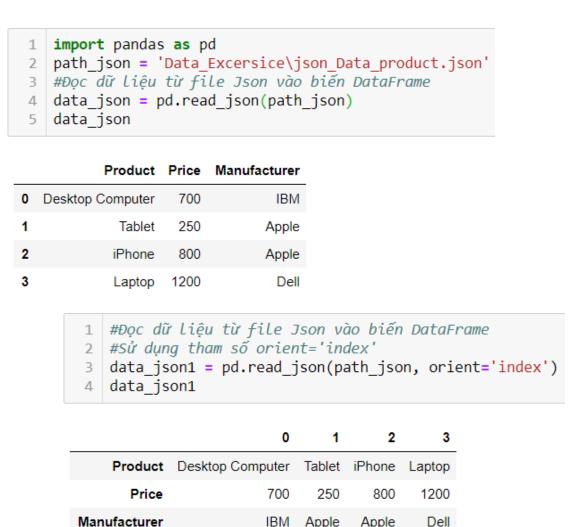


```
"orders": [
        "orderno": "748745375",
        "date": "June 30, 2088 1:54:23 AM",
        "trackingno": "TN0039291",
        "custid": "11045",
        "customer": [
                "custid": "11045",
                "fname": "Sue",
                "lname": "Hatfield",
                "address": "1409 Silver Street",
                "city": "Ashland",
                "state": "NE",
                "zip": "68003"
```

2.4 Đọc file Json

Sử dụng phương thức **pd.read_json()** để đọc dữ liệu từ file json vào dataframe.





2.4 Đọc file Json

Có thể sử dụng thư viện json để làm việc với file có định dạng json.

```
#Hoặc Sử dụng thư viện json làm việc với dữ liệu json
   import json
    with open(path json, 'r') as myfile:
        data=myfile.read()
 6 # Đọc dữ liệu vào biến obj
 7 obj = json.loads(data)
 8 type(obj)
dict
 1 #Lấy giá trị của key = 'Product'
 2 obj['Product']
{'0': 'Desktop Computer', '1': 'Tablet', '2': 'iPhone', '3': 'Laptop'}
 1 #Lấy giá trị của key = 'Price'
 2 obj['Price']
{'0': 700, '1': 250, '2': 800, '3': 1200}
```



Yêu cầu: Sinh viên đọc dữ liệu dạng json lưu trong file json_Data_flights.json vào trong DataFrame

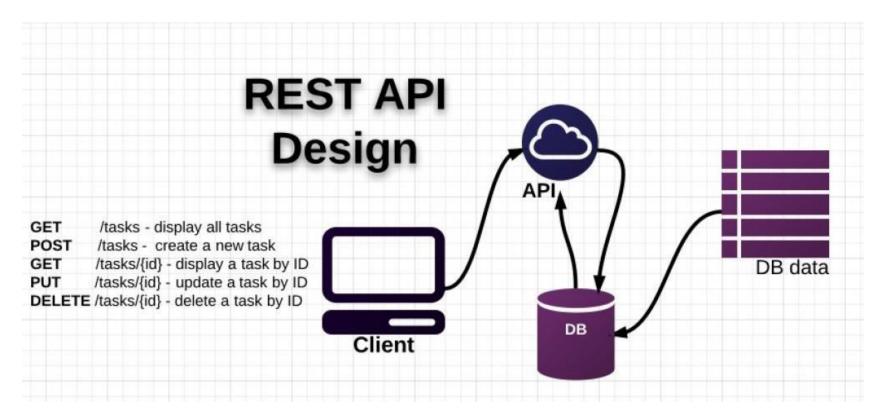
1 data_json.head()

	_id	dofW	carrier	origin	dest	crsdephour	crsdeptime	depdelay	crsarrtime	arrdelay	crselapsedtime	dist
0	ATL_BOS_2017- 01- 01_11_DL_1204	7	DL	ATL	BOS	11	1141	0	1409	0	148	946
1	ATL_BOS_2017- 01- 01_13_WN_3677	7	WN	ATL	BOS	13	1335	0	1600	0	145	946
2	ATL_BOS_2017- 01- 01_14_DL_1208	7	DL	ATL	BOS	14	1416	0	1644	0	148	946
3	ATL_BOS_2017- 01- 01_16_DL_980	7	DL	ATL	BOS	16	1616	15	1849	0	153	946
4	ATL_BOS_2017- 01- 01_18_WN_461	7	WN	ATL	BOS	18	1845	0	2110	0	145	946

3. Lấy dữ liệu từ API

3.1 Giới thiệu REST API

- REST API là một ứng dụng chuyển đổi cấu trúc dữ liệu có các phương thức để kếtnối với các thư viện và ứng dụng khác.
- Các trang web ngày nay thường sử dụng REST API để cho phép kết nối dữ liệu



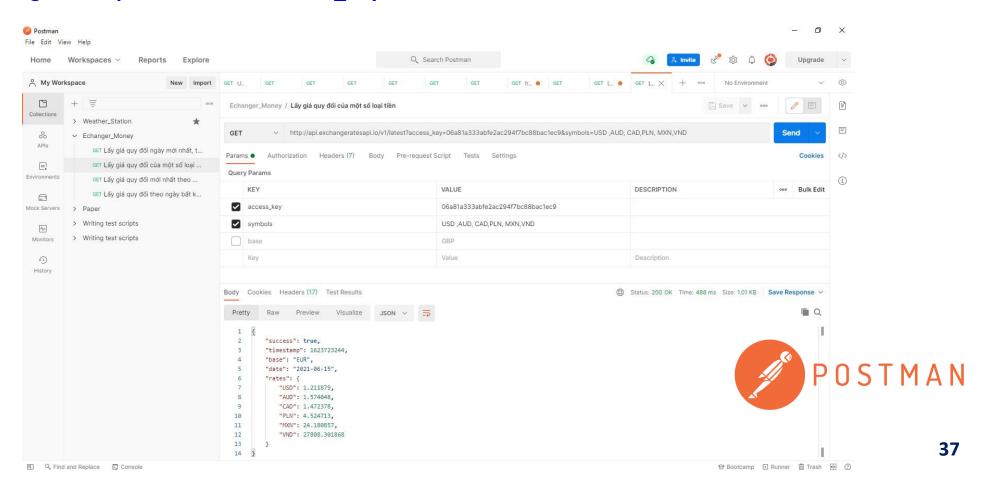


3.1 Giới thiệu REST API

- Vd API lấy dữ liệu tỷ giá ngoại tệ:
 - Lấy tỷ giá ngày cuối cùng, theo EUR:

http://api.exchangeratesapi.io/v1/latest?access key=06a81a333abfe2ac294f7bc88bac1ec9





3.2 Lấy dữ liệu từ API

- sử dụng thư viện requests, phương thức get() để lấy dữ liệu theo url
- Lấy tỷ giá mới nhất (chuyển đổi theo EUR (Default))

```
1 #Lấy dữ liệu tỷ lệ đổi tiền
    import pandas as pd
    import requests as rq
    url api = 'http://api.exchangeratesapi.io/v1/latest?access key=06a81a333abfe2ac294f7bc88bac1ec9'
    #lấy dữ liệu theo url api
    result money = rq.get(url api)
                                                                                                     1 #Chuyển đổi dữ liêu Json sang kiểu DataFrame
                                                                                                       df = pd.DataFrame(data json)
    #check status của request
                                                                                                       print(df.head(10))
    print(result money.status code)
                                                                                                        success timestamp base
                                                                                                                                     date
                                                                                                                                               rates
                                                                                                          True 1623723244 EUR 2021-06-15
                                                                                                   AED
                                                                                                                                            4.451469
200
                                                                                                   AFN
                                                                                                           True 1623723244
                                                                                                                           EUR
                                                                                                                               2021-06-15
                                                                                                                                            95.687275
                                                                                                          True 1623723244 EUR
                                                                                                                               2021-06-15 122.985167
                                                                                                   ALL
 1 #Chuyển đổi dữ liệu sang kiểu string
                                                                                                          True 1623723244
                                                                                                                           EUR
                                                                                                                               2021-06-15 630.187794
   data_text = result_money.text
                                                                                                   ANG
                                                                                                          True 1623723244 EUR 2021-06-15
                                                                                                                                            2.175393
 3 #chuyển đổi dữ liệu sang kiểu json
                                                                                                           True 1623723244
                                                                                                                           EUR 2021-06-15 778.026383
                                                                                                   AOA
    data json = result money.json()
                                                                                                           True 1623723244
                                                                                                                           EUR 2021-06-15 115.431336
    print('Text: ', data_text)
                                                                                                           True 1623723244
                                                                                                                           EUR 2021-06-15
                                                                                                                                            1.574048
                                                                                                   AUD
   print('-----
                                                                                                           True 1623723244
                                                                                                                           EUR 2021-06-15
                                                                                                                                            2.181988
   print('Json:', data json)
                                                                                                   AZN
                                                                                                           True 1623723244 EUR 2021-06-15
                                                                                                                                            2.059938
```

Text: {"success":true,"timestamp":1623723244,"base":"EUR","date":"2021-06-15","rates":{"AED":4.45146 9,"AFN":95.687275,"ALL":122.985167,"AMD":630.187794,"ANG":2.175393,"AOA":778.026383,"ARS":115.43133 6,"AUD":1.574048,"AWG":2.181988,"AZN":2.059938,"BAM":1.955651,"BBD":2.447085,"BDT":102.770573,"BGN": 1.955849,"BHD":0.456871,"BIF":2397.491166,"BMD":1.211879,"BND":1.607786,"BOB":8.380966,"BRL":6.13222

3.2 Lấy dữ liệu từ API

- Truyền tham số cho url: params={'key':'values'}
- Chỉ lấy Lấy tỷ giá một số loại tiền nhất định:

Mã NT	Tên Ngoại Tệ
USD	ĐÔ LA MỸ
AUD	ĐÔ LA ÚC
[◆] CAD	ĐÔ CANADA
CHF	FRANCE THỤY SĨ
EUR	EURO
∔ GBP	BÅNG ANH
JPY	YÊN NHẬT
SGD	ĐÔ SINGAPORE
= THB	BẠT THÁI LAN
MYR	RINGGIT MÃ LAY
∷ DKK	KRONE ĐAN MẠCH
MKD	ĐÔ HONGKONG
INR	RUPI ẤN ĐỘ
KRW	WON HÀN QUỐC

```
#Truyền tham số cho url
   #Lấy tỷ giá của một số ngoại tệ:
   #1.GBP (Bång Anh)
 4 #2.JPY (Yên Nhật)
   #3.MYR (Ringgit - Malaysia)
   #4. THB (Bat Thái Lan),
   #5.KRW (Won - Hàn Quốc)
    #6. VND (Viêt Nam Đồng)
    result_money1 = rq.get(url api,params={'symbols':'USD, JPY, THB, VND, MYR, GBP, KRW'})
    data json1 = result money1.json()
12 df1 = pd.DataFrame(data json1)
    print(df1.head(10))
              timestamp base
     success
                                    date
                                                 rates
       True 1623723244
GBP
                         EUR
                              2021-06-15
                                              0.859289
       True 1623723244
JPY
                         EUR
                              2021-06-15
                                            133.460542
       True 1623723244
KRW
                         EUR 2021-06-15
                                           1355.171091
       True 1623723244
                         EUR 2021-06-15
                                              4.985063
MYR
       True 1623723244
THB
                         EUR 2021-06-15
                                             37.749351
                                              1.211879
USD
       True 1623723244
                         EUR 2021-06-15
VND
       True 1623723244
                         EUR 2021-06-15 27808.301868
```

Thực hành 2.5

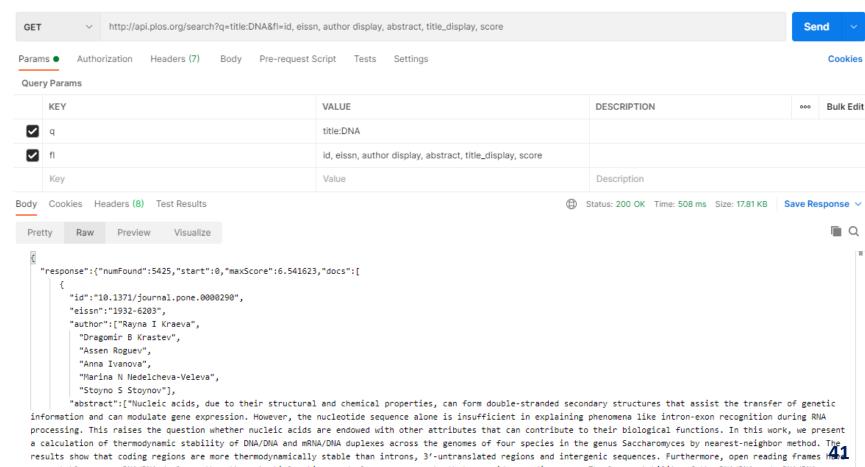
Thực hành 2.5



Yêu cầu: Sinh viên lấy dữ liệu các bài báo trong tiêu đề có từ Viruscung cấp tù đia

http://api.plos.org/search?q=title:VIRUS

- Tao DataFrame bao gồm các thông tin: id, eissn, author display, abstract, title_display, score.
- liêu từ Dataframe file Paper.csv

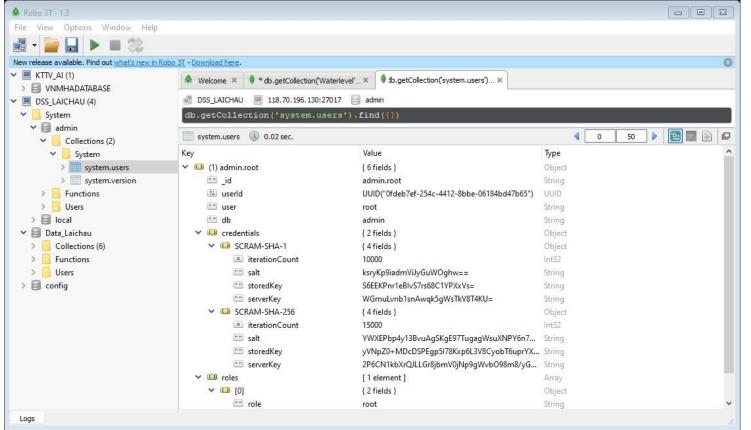


more stable sense mRNA/DNA duplexes than the potential antisense duplexes, a property that can aid gene discovery. The lower stability of the DNA/DNA and mRNA/DNA

4. Lấy dữ liệu từ MongoDB

4.1 CSDL NoSQL

Cơ sở dữ liệu NoSQL sử dụng nhiều mô hình dữ liệu để truy cập và quản lý dữ liệu. Cácloại cơ sở dữ liệu này được tối ưu hóa dành riêng cho các ứng dụng yêu cầu mô hình dữliệu linh hoạt có lượng dữ liệu lớn và độ trễ thấp, có thể đạt được bằng cách giảm bớt mộtsố hạn chế về tính nhất quán của dữ liệu của các cơ sở dữ liệu khác.



Is MongoDB NoSQL



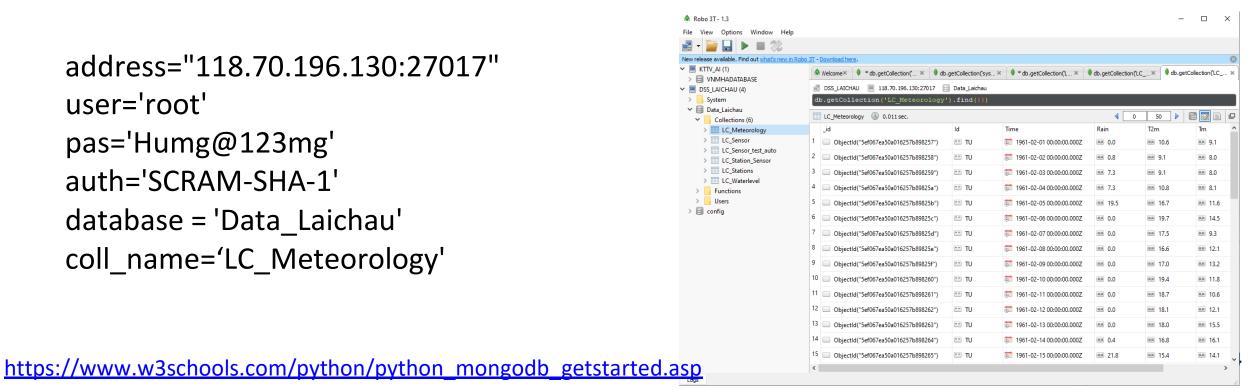
4.2 Truy cập dữ liệu Mongo

Để kết nối tới CSDL MongoDB sử dụng thư viện pymongo.

```
#Kiểm tra thư viện pymongo
    import pymongo
    print('Version: ',pymongo.version)
Version: 3.7.2
```

Thực hiện kết nối và lấy dữ liệu từ collection: LC_Meteorology với các tham số:

```
address="118.70.196.130:27017"
user='root'
pas='Humg@123mg'
auth='SCRAM-SHA-1'
database = 'Data Laichau'
coll name='LC_Meteorology'
```



4.2 Truy cập dữ liệu Mongo

• Kết nối tới CSDL MongoDB.

```
#Thực hiện truy vấn lấy dữ liệu
#Thiết lập câu truy vấn thực hiện:
#1. Lấy dữ liệu từ trạm có mã TU
#2. Sắp xếp dữ liệu theo thứ tự thời gian giảm dần
stationid = 'TU'
data_mg = col.find({"Id":stationid}).sort([('Time',-1)])
print(type(data_mg))
```

<class 'pymongo.cursor.Cursor'>

4.2 Truy cập dữ liệu Mongo

Chuyển đổi kết quả từ câu truy vấn về dữ liệu DataFrame.

```
#lấy dữ liệu chuyển vào các biến kiểu list

ttime,dt_rain,dt_t2m,dt_tm,dt_tx = [],[],[],[],[]

for i in data_mg:
    dt_time.append(str(i['Time']))
    dt_rain.append(i['Rain'])
    dt_t2m.append(i['T2m'])
    dt_tm.append(i['Tm'])
    dt_tx.append(i['Tx'])
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21143 entries, 0 to 21142
Data columns (total 5 columns):
# Column Non-Null Count Dtype

0 Time 21143 non-null object
1 Rain 21112 non-null float64
2 T2m 21142 non-null float64
3 Tm 21142 non-null float64
4 Tx 21142 non-null float64
dtypes: float64(4), object(1)
memory usage: 826.0+ KB
```

```
1 #hiển thị 5 dòng dữ liệu đầu tiên
2 df.head()
```

	Time	Rain	T2m	Tm	Tx
0	2018-12-31 00:00:00	2.6	9.7	8.3	11.5
1	2018-12-30 00:00:00	13.9	10.0	9.4	11.6
2	2018-12-29 00:00:00	1.6	14.1	11.5	18.4
3	2018-12-28 00:00:00	0.0	17.2	12.6	24.8
4	2018-12-27 00:00:00	0.0	17.0	12.0	25.8



Yêu cầu 4.1: Sinh viên tạo DataFrame đọc dữ liệu trong Collection LC_Meteorolory trạm có ld = 'LC' sắp xếp theo thứ tự thời gian giảm dần, chỉ lấy 1000 dòng dữ liệu đầu tiên và chỉ lấy dữ liệu các cột Time, Rain, T2m.

```
db.getCollection('LC_Meteorology')
.find({'Id':'LC'},{'Time':1,'Rain':1,'T2m':1})
.sort({'Time':-1})
.limit(1000)
```

```
1 #Hiển thị 10 hàng dữ liệu đầu tiên
2 df_lc.head(10)
```

	Time	Rain	T2m
0	2018-12-31 00:00:00	1.0	14.4
1	2018-12-30 00:00:00	11.9	15.1
2	2018-12-29 00:00:00	3.0	20.1
3	2018-12-28 00:00:00	0.0	19.3
4	2018-12-27 00:00:00	0.0	18.5
5	2018-12-26 00:00:00	0.0	17.8
6	2018-12-25 00:00:00	0.0	18.8
7	2018-12-24 00:00:00	0.0	18.9
8	2018-12-23 00:00:00	0.0	18.4
9	2018-12-22 00:00:00	0.0	19.2



Yêu cầu 4.2: Sinh viên tạo DataFrame df_station đọc toàn bộ dữ liệu trong Collection LC_Stations.

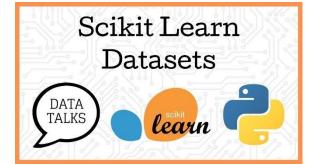
Idstation	Namestation	Timestart	Timeend
"." LC	"" Lai Chau	₹ 1961-01-01 00:00:00.000Z	≅ 2018-12-31 00:00:00.000Z
"" TU	"" Than Uyen	≅ 1961-02-01 00:00:00.000Z	≅ 2018-12-31 00:00:00.000Z
"" MT	"" Muong Te	≅ 1961-01-01 00:00:00.000Z	≅ 2018-12-31 00:00:00.000Z
"" SH	"" Sin Ho	≅ 1961-01-01 00:00:00.000Z	≅ 2018-12-31 00:00:00.000Z
"" TD	Tam Duong	≅ 1975-01-01 00:00:00.000Z	≅ 2018-12-31 00:00:00.000Z

5. Tải dữ liệu mẫu với Sklearn

<u>Dữ liệu mẫu trong Sklearn</u>

- Trong thư viện sklearn cung cấp một số bộ dữ liệu mẫu, kích thước nhỏ, rất hữu ích cho việc demo các thuật toán học máy.
- Sử dụng dữ liệu mẫu: import sklearn.datasets

load_boston(*[, return_X_y])	Load and return the boston house-prices dataset (regression).
<pre>load_iris(* [, return_X_y, as_frame])</pre>	Load and return the iris dataset (classification).
<pre>load_diabetes(* [, return_X_y, as_frame])</pre>	Load and return the diabetes dataset (regression).
<pre>load_digits(* [, n_class, return_X_y, as_frame])</pre>	Load and return the digits dataset (classification).
load_linnerud(* [, return_X_y, as_frame])	Load and return the physical excercise linnerud dataset.
load_wine(* [, return_X_y, as_frame])	Load and return the wine dataset (classification).
load_breast_cancer(* [, return_X_y, as_frame])	Load and return the breast cancer wisconsin dataset (classification).



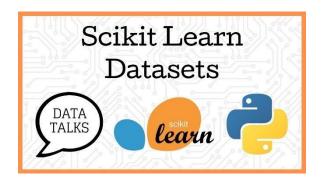
Read more: https://scikit-learn.org/stable/datasets/index.html#

1) Boston house-prices dataset

Tập dữ liệu về giá nhà tại thành phố Boston.

Data Set Characteristics:

Number of Instances:	506		
Number of Attributes:	13 numeric/categorical predictive. Median Value (attribute 14) is usually the target		
Attribute Information (in order):	 CRIM per capita crime rate by town ZN proportion of residential land zoned for lots over 25,000 sets. INDUS proportion of non-retail business acres per town CHAS Charles River dummy variable (= 1 if tract bounds river; NOX nitric oxides concentration (parts per 10 million) RM average number of rooms per dwelling AGE proportion of owner-occupied units built prior to 1940 DIS weighted distances to five Boston employment centres RAD index of accessibility to radial highways TAX full-value property-tax rate per \$10,000 PTRATIO pupil-teacher ratio by town B 1000(Bk - 0.63)^2 where Bk is the proportion of blacks by to LSTAT % lower status of the population MEDV Median value of owner-occupied homes in \$1000's 	0 other	#Đọc tậ
Missing Attribute Values:	None	2 3 4 5	import s X, y = 0 print(ty print(')
Creator:	Harrison, D. and Rubinfeld, D.L.	7	print('



Samples total	506
Dimensionality	13
Features	real, positive
Targets	real 5 50.

```
#Doc tập dữ liệu Boston House Prices Dataset
import sklearn.datasets as datask
X, y = datask.load_boston(return_X_y=True)

print(type(X))
print('Kích thước dữ liệu đầu vào (features):', X.shape)
print('Kích thước dữ liệu đầu ra (target):', X.shape)
```

<class 'numpy.ndarray'>
Kich thước dữ liệu đầu vào (features): (506, 13)
Kich thước dữ liệu đầu ra (target): (506, 13)

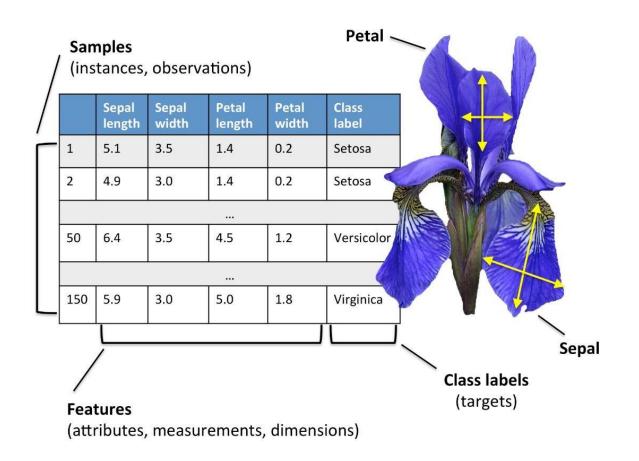
2) The iris dataset

• Tập dữ liệu về thông số chiều rộng, chiều dài của lá hóa và cánh hóa của 3 loại hoa Lan



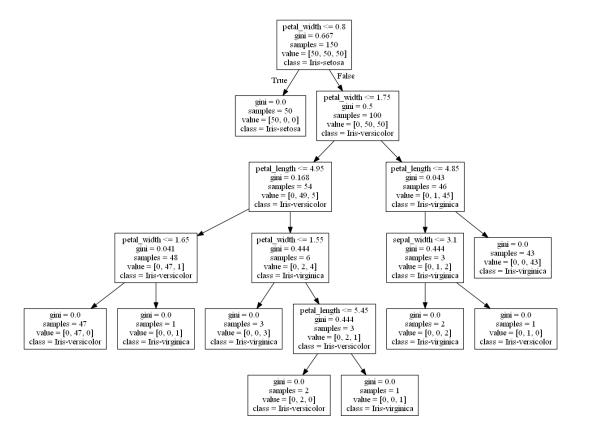
IRIS DATASET

Classes	3
Samples per class	50
Samples total	150
Dimensionality	4
Features	real, positive



2) The iris dataset

- Tham số:
 - return_X_y=True: Dữ liệu được đọc vào các biến array.
 - as_frame = True: Dữ liệu được đọc vào các biến DataFrame



```
##Doc tập dữ liệu Iris Dataset
import sklearn.datasets as datask
X_iris, y_iris = datask.load_iris(return_X_y=True)

print(type(X_iris))
print('Kích thước dữ liệu đầu vào (features):', X_iris.shape)
print('Kích thước dữ liệu đầu ra (target) :', y_iris.shape)
print('Bộ dữ liệu 1) ', X[1,:], '--',y[1])
print('Bộ dữ liệu 55) ', X[55,:], '--',y[55])
print('Bộ dữ liệu 111)', X[111,:], '--',y[111])
```

```
<class 'numpy.ndarray'>
Kích thước dữ liệu đầu vào (features): (150, 4)
Kích thước dữ liệu đầu ra (target) : (150,)
Bộ dữ liệu 1) [4.9 3. 1.4 0.2] -- 0
Bộ dữ liệu 55) [5.7 2.8 4.5 1.3] -- 1
Bộ dữ liệu 111) [6.4 2.7 5.3 1.9] -- 2
```

Minh họa xây dựng mô hình học máy sử dụng thuật toán Decision Tree cho bài toán phân lớp hoa lan sử dụng bộ dữ liệu mẫu để huấn luyện



Yêu cầu 5.1: Sinh viên tìm hiểu tập dữ liệu mẫu nhận dạng các loại rượu (Wine recognition dataset).

- Đọc dữ liệu từ dataset mẫu theo 2 dạng vào biến ndarray - dataframe.
- Hiển thị một số mẫu rượu trong tập dữ liệu

Data Set Characteristics:

Number of 178 (50 in each of three classes)	
Instances:	
Number of 13 numeric, predictive attributes and Attributes:	d the class
Attribute Information: • Alcohol • Malic acid • Ash • Alcalinity of ash • Magnesium • Total phenols • Flavanoids • Nonflavanoid phenols • Proanthocyanins • Color intensity • Hue • OD280/OD315 of diluted wines • Proline	

class:

- class_0
- o class_1
- o class_2

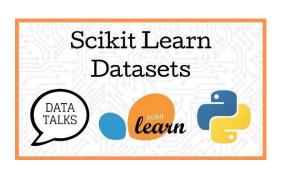


Yêu cầu 5.2: Sinh viên tìm hiểu tập dữ liệu mẫu còn lại trong bộ Datasets của thư viện Sklearn

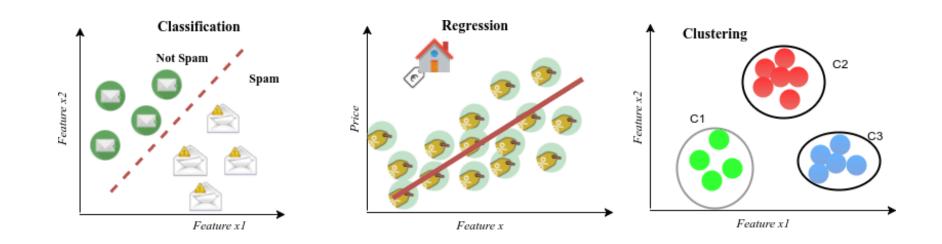
fetch_olivetti_faces(* [, data_home,])	Load the Olivetti faces data-set from AT&T (classification).
fetch_20newsgroups(* [, data_home, subset,])	Load the filenames and data from the 20 newsgroups dataset (classification).
<pre>fetch_20newsgroups_vectorized(* [, subset,])</pre>	Load the 20 newsgroups dataset and vectorize it into token counts (classification).
fetch_lfw_people(* [, data_home, funneled,])	Load the Labeled Faces in the Wild (LFW) people dataset (classification).
fetch_lfw_pairs(* [, subset, data_home,])	Load the Labeled Faces in the Wild (LFW) pairs dataset (classification).
fetch_covtype(*[, data_home,])	Load the covertype dataset (classification).
fetch_rcv1(* [, data_home, subset,])	Load the RCV1 multilabel dataset (classification).
fetch_kddcup99(* [, subset, data_home,])	Load the kddcup99 dataset (classification).
fetch_california_housing(* [, data_home,])	Load the California housing dataset (regression).

6. Tạo Simulated Dataset

• Sklearn cho phép tạo ra các bộ dữ liệu mô phỏng theo các tham số thiết lập để chạy cho các thuật toán:



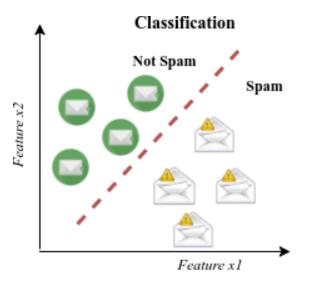
- Sklearn cho phép tạo các bộ dữ liệu cho 3 bài toán cơ bản trong học máy:
 - Phân lớp: sklearn.datasets.make_classification
 - Höi quy: sklearn.datasets.make_regression
 - Phân cụm: sklearn.datasets.make_blobs



- Tạo dữ liệu mẫu cho bài toán phân lớp:
 - sklearn.datasets.make_classification

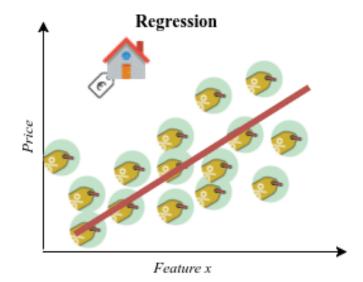
```
#load module datasets trong thư viện sklearn
import sklearn.datasets as datask

#Tạo dữ liệu mẫu cho bài toán PHÂN LỚP
#Với 200 mẫu, đầu vào 5 thuộc tính, phân thành 2 lớp
X_features,y_target = datask.make_classification(n_samples=200, n_features=5, n_classes=2)
print(X_features[:10,:])
print(y_target[:10])
```



- Tạo dữ liệu mẫu cho bài toán hồi quy:
 - sklearn.datasets.make_regression

```
[[-0.36956243 0.4972691 -2.14444405 0.2373327 -1.88141087]
  0.40890054 -0.03869551 1.12141771 -1.61577235 -1.31228341
  0.0607502 -0.43750898 0.92145007 0.09542509 -0.78191168
 -1.02188594 1.76795995 0.47761018 -0.47537288 -1.42655542
  0.6210827 -0.79726979 -0.82609743 0.28267571 -0.13597733
 -1.57915629 0.63019567 0.45194604 -0.4148469 -0.10534471
 -0.18695502 -2.277298
                         0.35387043 -0.06962454 -0.11644415
  0.44136444 -2.43483776 2.18697965 1.0388246 -0.31011677
  1.15528789 -0.70584051 -0.23794194 1.19268607 0.9561217
  0.7147896 2.13782807 -1.75592564 -0.785534
                                                0.56438286]]
 -42.26048415 -93.80258805 -49.39514348
                                          19.83790815 -37.36573177
  16.42912647 -118.52143931 -105.57802562
                                          32.55783001 107.22348747]
```



Read more: https://scikit-learn.org/stable/modules/generated/sklearn.datasets.make-regression.html

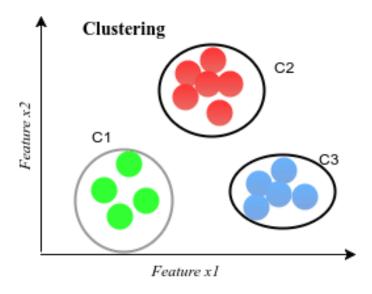
- Tạo dữ liệu mẫu cho bài toán PHÂN CỤM:
 - sklearn.datasets.make_blobs

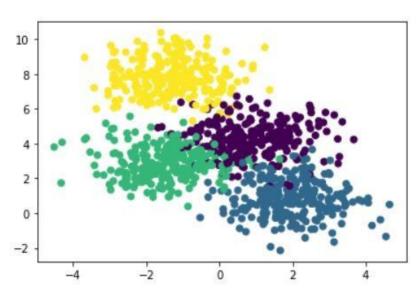
```
#load module datasets trong thư viện sklearn
import sklearn.datasets as datask

#Tạo dữ liệu mẫu cho bài toán PHÂN CỤM
X_features, y_target = datask.make_blobs(n_samples=1000, centers=4, n_features=2, random_state=0)

print(X_features[:10,:])
print(y_target[:10])
```

```
[[-2.50804312 7.86408052]
[ 0.4666179 3.86571303]
[-2.94062621 8.56480636]
[-2.89994656 1.85014025]
[ 2.34438803 1.31053448]
[ 0.87305123 4.71438583]
[ 0.84032038 5.44067869]
[ 1.66287852 -2.14847939]
[ 0.85810603 3.62360912]
[ -0.72183574 4.6910678 ]]
[ 3 0 3 2 1 0 0 1 0 0]
```







Yêu cầu 6.1: Sinh viên tạo một tập dữ liệu mẫu cho bài toán phân lớp theo yêu cầu sau:

- Dataset mẫu bao gồm: 1000 samples, 8 features, 2 classes.
- Xây dựng một mô hình ML theo theo thuật toán Decision Tree với Dataset tạo được.



Yêu cầu 6.2: Sinh viên tạo một tập dữ liệu mẫu cho bài toán phân cụm theo yêu cầu sau:

- 1. Dataset mẫu bao gồm: 10000 samples, 2 features, 5 centers, cluster_std = 0.5.
- 2. Trực quan hóa tập Dataset ở trên lên biểu đồ scatter.
- 3. Thay đổi tham số cluster_std trong khoảng từ [0-1.0], nhận xét kết quả thu được

Thank you!