

**ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH  
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN  
KHOA: CÔNG NGHỆ THÔNG TIN**



**BÁO CÁO  
LAB01 – Preprocessing**

**KHAI THÁC DỮ LIỆU VÀ ỨNG DỤNG – CQ2017**

**Sinh viên thực hiện:**

Nguyễn Văn Thìn – 1712787

Lê Bá Quyền – 1712713

Nguyễn Lê Trường Thành - 1712775

Q5, ngày 6 tháng 10 năm 2019

## I. THÔNG TIN NHÓM

Member	MSSV	Email	Tự đánh giá cá nhân
Nguyễn Văn Thìn	1712787	vanthin7111999@gmail.com	90%
Lê Bá Quyền	1712713	quyenleba2291@gmail.com	90%
Nguyễn Lê Trường Thành	1712775	provip218@gmail.com	90%

(Điểm thành viên chia đều)

Quá trình làm việc nhóm- Phân chia task cho từng member	
<b>Phản trả lời câu hỏi</b>	<b>Thành:</b> deadline 26/9 đã completed
<b>Phản code</b>	<b>Quyền:</b> Code BT1 <b>Thìn:</b> Code BT2 21/9 -> 27/9 >>> Tìm hiểu cú pháp python 28/9 -> 30/9 >>> Code 1 số yêu cầu 1/10 -> 5/10 >>> Hoàn thiện các phần còn lại 6/10 chỉnh sửa phần còn thiếu – summit 23h trên moodle

### Những phần đã làm được:

Trả lời câu hỏi tự luận

Câu hỏi từ 1->9

Nội dung thực hiện cài đặt(25 điểm)

#### 1. Tiền xử lý dữ liệu trên tập dữ liệu tổng quát với một số chức năng đơn giản (15 điểm)

- Chuẩn hóa min-max trên danh sách thuộc tính chỉ định.
- Chuẩn hóa Z-scores trên danh sách thuộc tính chỉ định.
- Rời rạc hóa dữ liệu bằng phương pháp chia giỏ theo độ rộng trên danh sách thuộc tính chỉ định.
- Rời rạc hóa dữ liệu bằng phương pháp chia giỏ theo độ sâu trên danh sách thuộc tính chỉ định.
- Xóa các mẫu dữ liệu thiếu giá trị trên danh sách thuộc tính chỉ định.
- Điền giá trị bị thiếu trên danh sách thuộc tính chỉ định, giá trị được điền là giá trị trung bình (mean) của thuộc tính nếu đó là thuộc tính số hoặc điền giá trị có tần số xuất hiện cao nhất (mode) nếu là thuộc tính rời rạc.

#### 2. Tiền xử lý dữ liệu trên tập dữ liệu cụ thể cho trước (10 điểm)

1. Xóa các mẫu rỗng.
2. Xóa các mẫu bị trùng lặp
3. Chuyển diện tích về km<sup>2</sup>
4. Sử dụng chương trình đã cài đặt ở phần B-1. để xóa các mẫu bị thiếu diện tích

### **Những phần chưa làm được:**

Câu 10 phần trả lời câu hỏi

---

## **II. TRẢ LỜI CÂU HỎI TỰ LUÂN**

### **YÊU CẦU 1 – TẠO TẬP TIN ARFF TỪ TẬP DỮ LIỆU COURSE RATINGS**

**1. Tập dữ liệu được đọc vào Weka thành công hay không? Nếu có, trả lời câu hỏi tiếp theo. Nếu không, cho biết lỗi gặp phải và bạn đã sửa lỗi đó như thế nào?**

- Tập dữ liệu được đọc vào weka thành công.

**2. Sau khi đọc dữ liệu thành công, quan sát thông tin thể hiện trên giao diện Explorer và trả lời những câu hỏi sau đây**

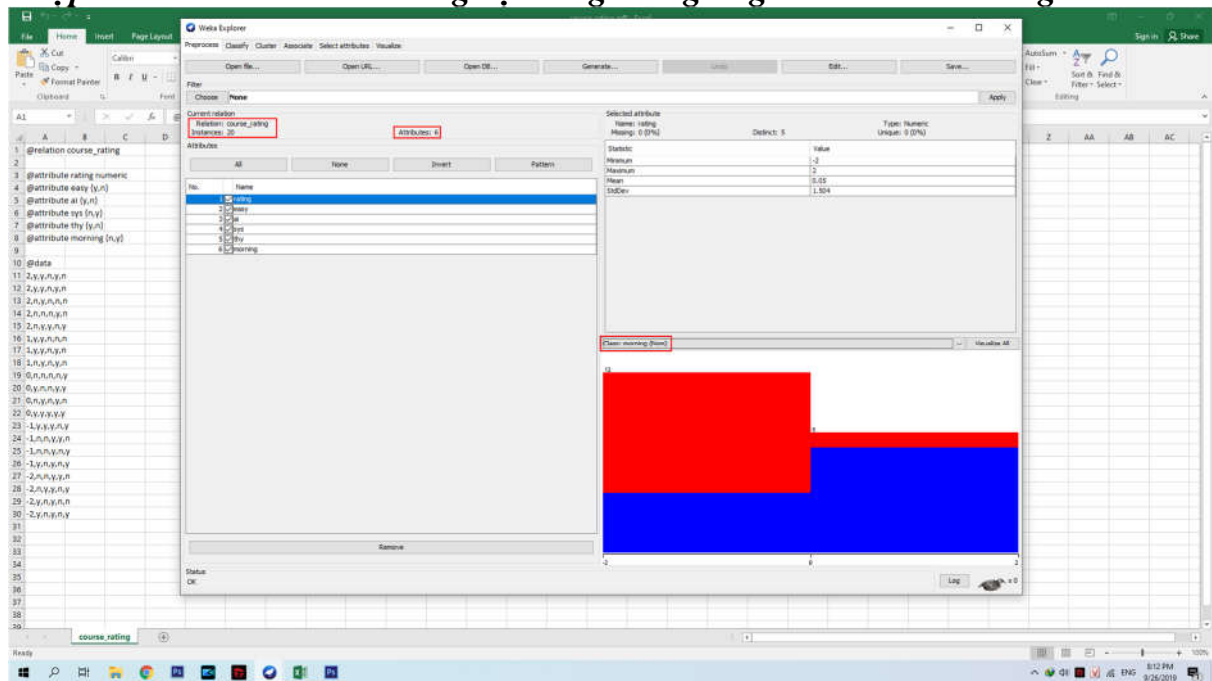
- **Tên của mối quan hệ (relation) trong dữ liệu là gì?**  
course\_rating

- **Tập dữ liệu có bao nhiêu mẫu (instances)?**  
instances: 20

- **Tập dữ liệu có bao nhiêu thuộc tính (attributes)?**  
attributes: 6

- **Thuộc tính nào trong tập dữ liệu là thuộc tính lớp (class)?**  
rating

**Chụp màn hình và dán những nội dung tương ứng để làm minh chứng.**



**3. Bạn có nhận thấy điều gì đáng chú ý khi quan sát các thông tin thống kê và đồ thị trình diễn?**

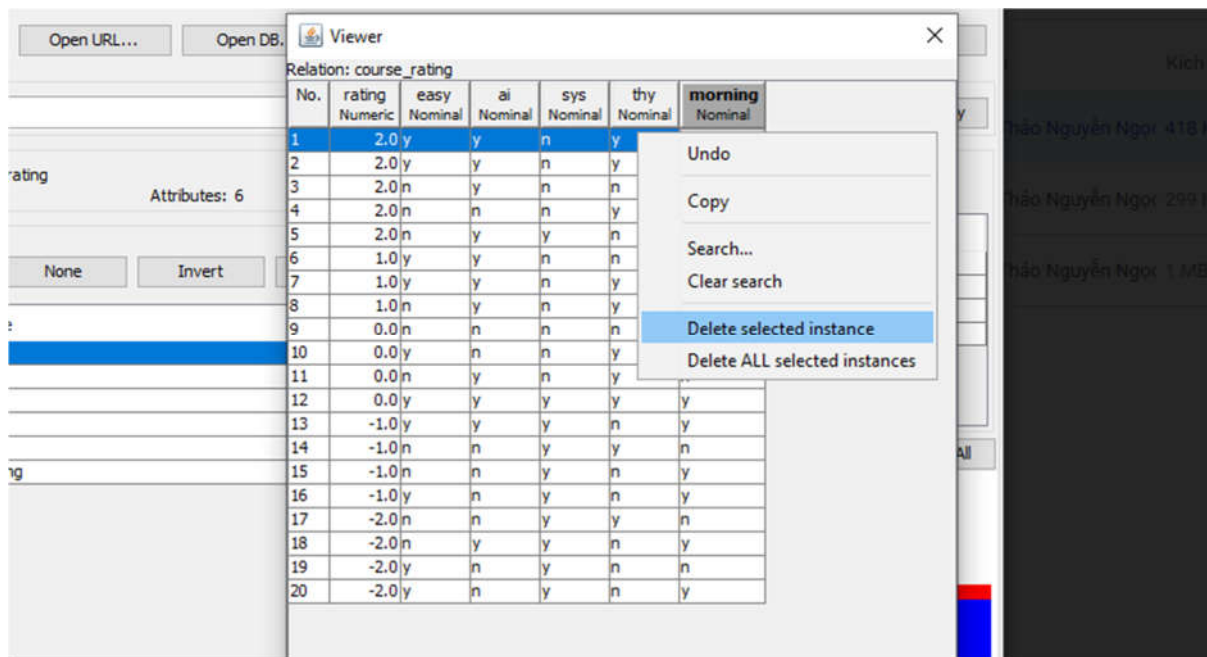
Đồ thị trình diễn ngoài chia ra theo giá trị thì còn được chia thành các màu (đỏ, xanh) trong từng giá trị cụ thể

**4. Điều gì xảy ra nếu thuộc tính lớp là thuộc tính đầu tiên bên trái? So sánh nội dung các đồ thị trình diễn do chức năng Visualize All (nằm góc dưới bên phải trong tab Preprocess) cung cấp, trong trường hợp thuộc tính lớp là cột đầu tiên và trong trường hợp thuộc tính lớp là cột cuối cùng. Lưu ý rằng tại bước này ta chỉ quan sát thông tin trình diễn dữ liệu chứ không làm bất kỳ thao tác gì ảnh hưởng đến nội dung dữ liệu.**

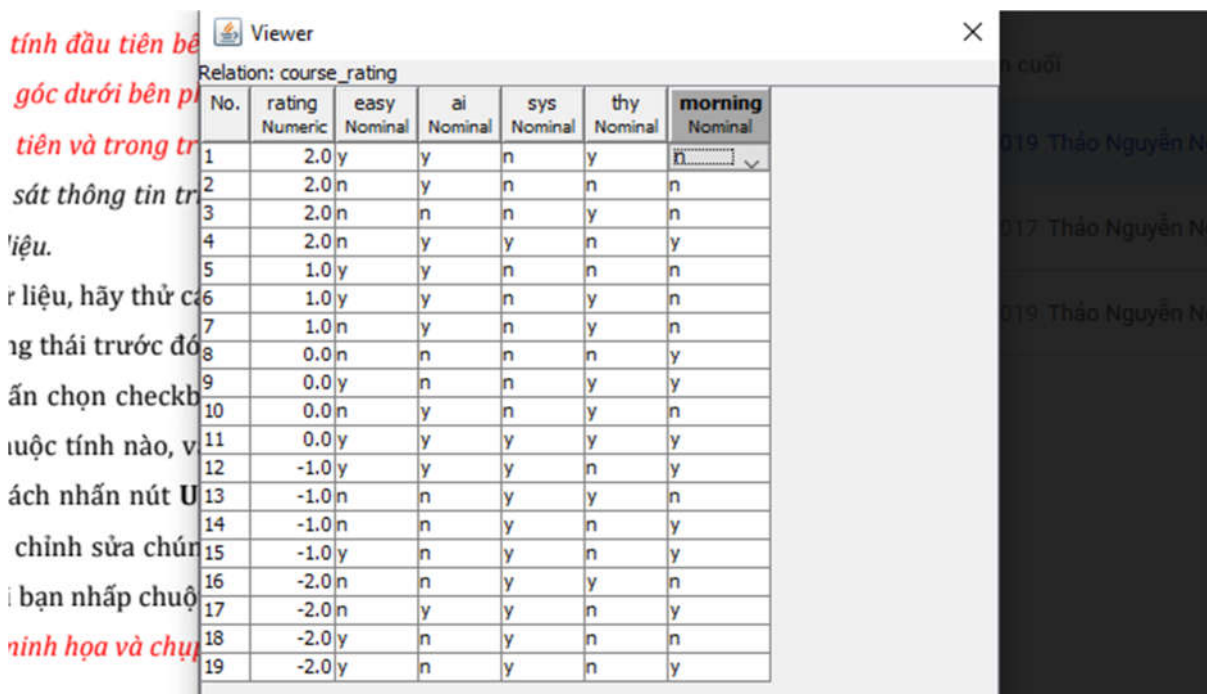
Nếu để thuộc tính lớp là cột đầu tiên thì thì **Visualize All** sẽ ra được đồ thị thì theo thuộc tính lớp morning. Còn nếu để ở cuối thì sẽ ko ra được đồ thị

**5. Với mỗi thao tác nêu trên, tùy ý thực hiện minh họa và chụp màn hình làm minh chứng**

Xóa:



Sau khi xóa:



Undo:

Viewer

Relation: course\_rating

No.	rating Numeric	easy Nominal	ai Nominal	sys Nominal	thy Nominal	morning Nominal
1	2.0	y	y	n	y	n
2	2.0	y	y	n	y	n
3	2.0	n	y	n	n	n
4	2.0	n	n	n	y	n
5	2.0	n	y	y	n	y
6	1.0	y	y	n	n	n
7	1.0	y	y	n	y	n
8	1.0	n	y	n	y	n
9	0.0	n	n	n	n	y
10	0.0	y	n	n	y	y
11	0.0	n	y	n	y	n
12	0.0	y	y	y	y	y
13	-1.0	y	y	y	n	y
14	-1.0	n	n	y	y	n
15	-1.0	n	n	y	n	y
16	-1.0	y	n	y	n	y
17	-2.0	n	n	y	y	n
18	-2.0	n	y	y	n	y
19	-2.0	y	n	y	n	n
20	-2.0	y	n	y	n	y

Undo OK Cancel

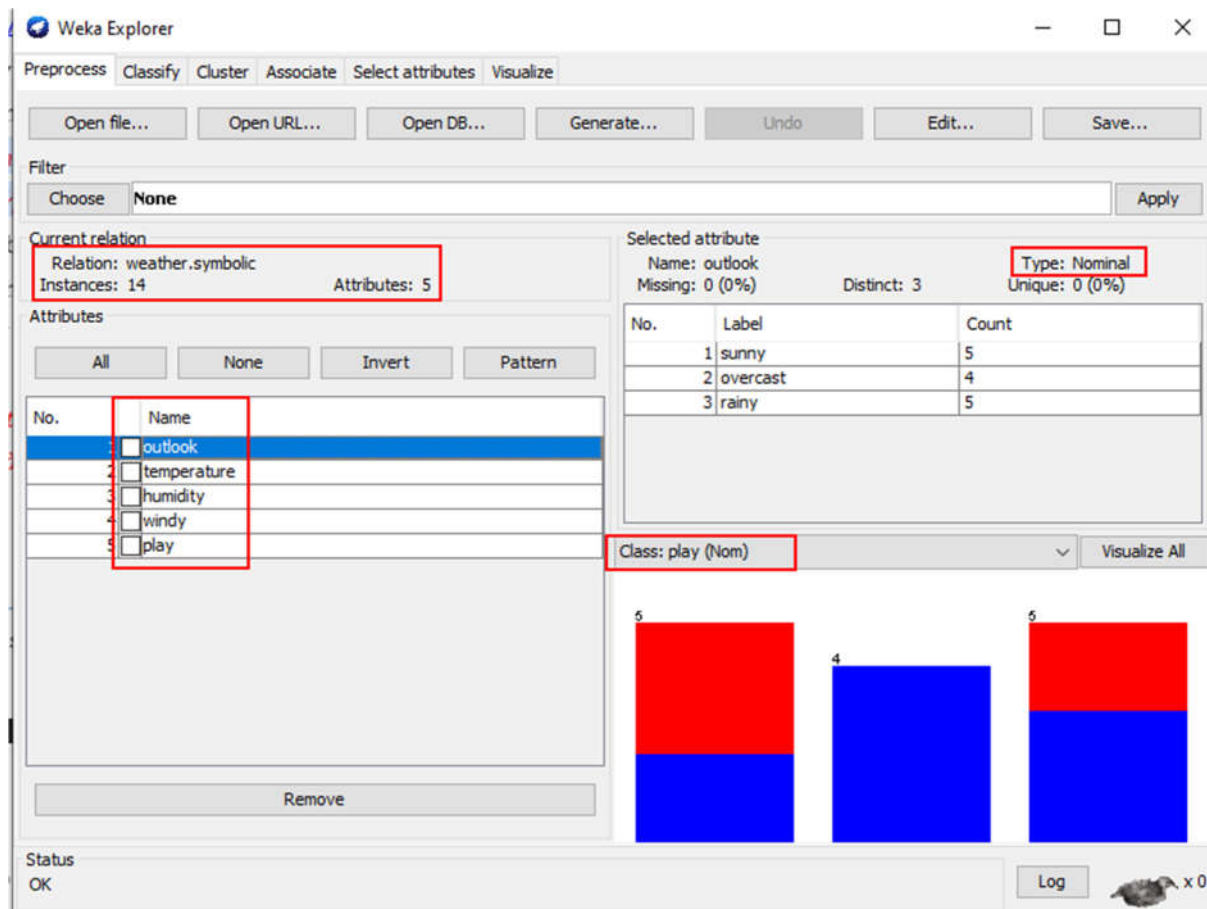
Chỉnh sửa:

Relation: course_rating						
No.	rating Numeric	easy Nominal	ai Nominal	sys Nominal	thy Nominal	morning Nominal
1	2.0	y	y	n	y	n
2	2.0	y	y	n	y	n
3	2.0	n	y	n	n	n
4	2.0	n	n	n	y	n
5	2.0	n	y	y	n	y
6	1.0	y	y	n	n	n
7	1.0	y	y	n	y	n
8	1.0	n	y	n	y	n
9	0.0	n	n	n	n	y
10	0.0	y	n	n	y	y
11	0.0	n	y	n	y	n
12	0.0	y	y	y	y	y
13	-1.0	y	y	y	n	y
14	-1.0	n	n	y	y	n
15	-1.0	n	n	y	n	y
16	-1.0	y	n	y	n	y
17	-2.0	n	n	y	y	n
18	-2.0	n	y	y	n	y
19	-2.0	y	n	y	n	n
20	-2.0	y	n	y	n	y

## Yêu cầu 2 – Khảo sát tập dữ liệu Weather (thuộc tính rời rạc)

**6. Tập dữ liệu có bao nhiêu mẫu? Bao nhiêu thuộc tính? Tên của các thuộc tính này là gì? Các thuộc tính này có loại gì? Thuộc tính nào là lớp?**

Tập dữ liệu có 14 mẫu, 5 thuộc tính: outlook, temperature, humidity, windy và play. Các thuộc tính này thuộc loại nominal. Thuộc tính lớp là play



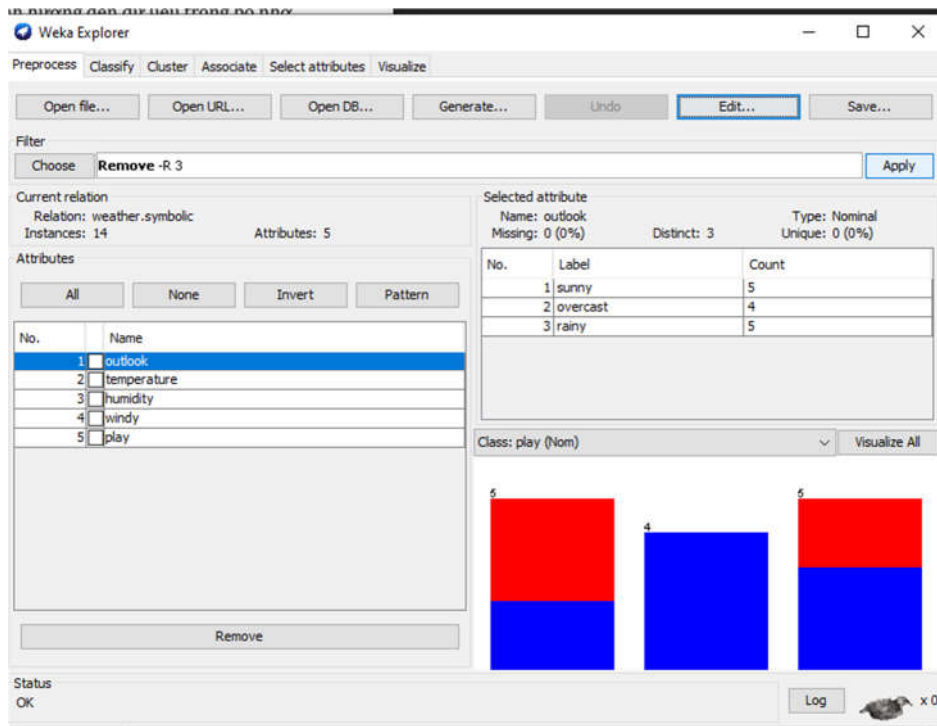
## 7. Chức năng của cột đầu tiên ở cửa sổ Viewer là gì?

Để đánh dấu số thứ tự cho các giá trị thuộc tính ở dưới

Lớp của mẫu thứ 8 trong tập dữ liệu là gì? No

Xóa thuộc tính 3





Sử dụng bộ lọc `weka.unsupervised.instance.RemoveWithValues` để loại bỏ mọi mẫu có giá trị thuộc tính `humidity` là `high`.

Trước khi áp dụng bộ lọc

Viewer

Relation: weather.symbolic

No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no

Sau khi áp dụng bộ lọc

Viewer

Relation: weather.symbolic-weka.filters.unsupervised.instance.RemoveWithValues

No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	rainy	cool	normal	FALSE	yes
2	rainy	cool	normal	TRUE	no
3	overcast	cool	normal	TRUE	yes
4	sunny	cool	normal	FALSE	yes
5	rainy	mild	normal	FALSE	yes
6	sunny	mild	normal	TRUE	yes
7	overcast	hot	normal	FALSE	yes

## Tham số áp dụng

weka.gui.GenericObjectEditor

weka.filters.unsupervised.instance.RemoveWithValues

About

Filters instances according to the value of an attribute. [More](#) [Capabilities](#)

attributeIndex: 3

dontFilterAfterFirstBatch: False

invertSelection: False

matchMissingValues: False

modifyHeader: False

nominalIndices: 1

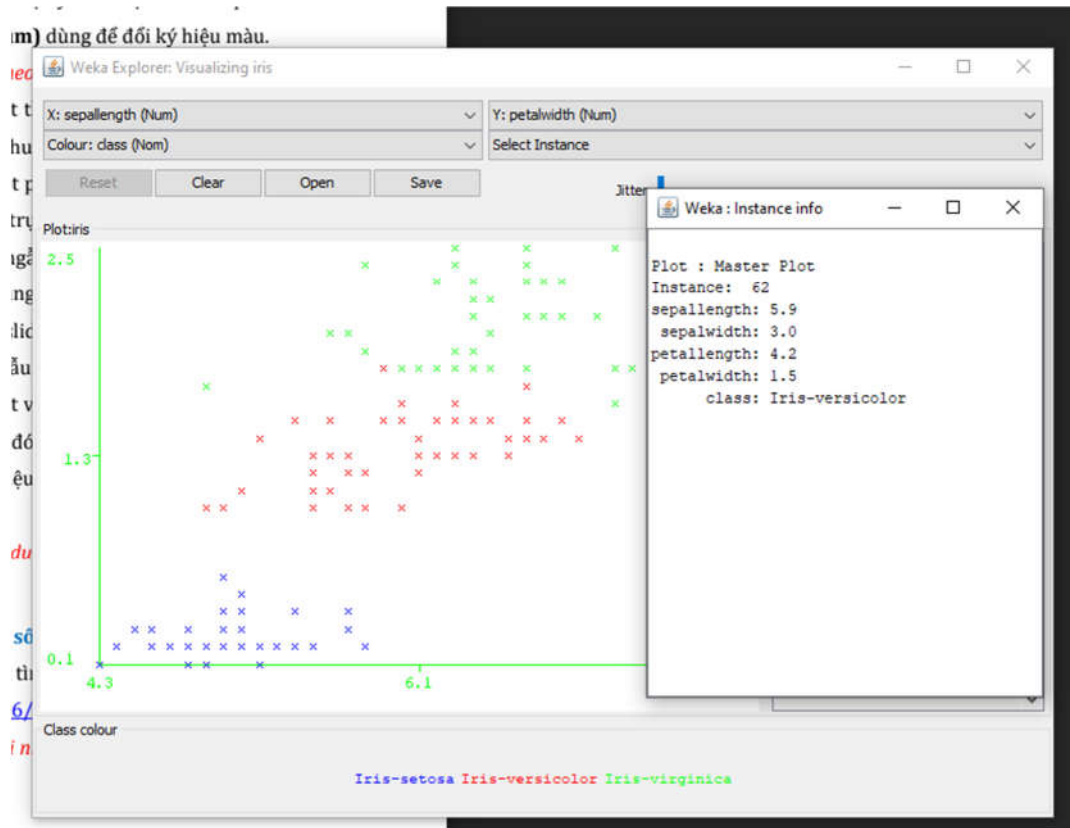
splitPoint: 0.0

[Open...](#) [Save...](#) [OK](#) [Cancel](#)

## 9. Mô tả tập dữ liệu. Nội dung của phần ghi chú (comment) nói về điều gì?

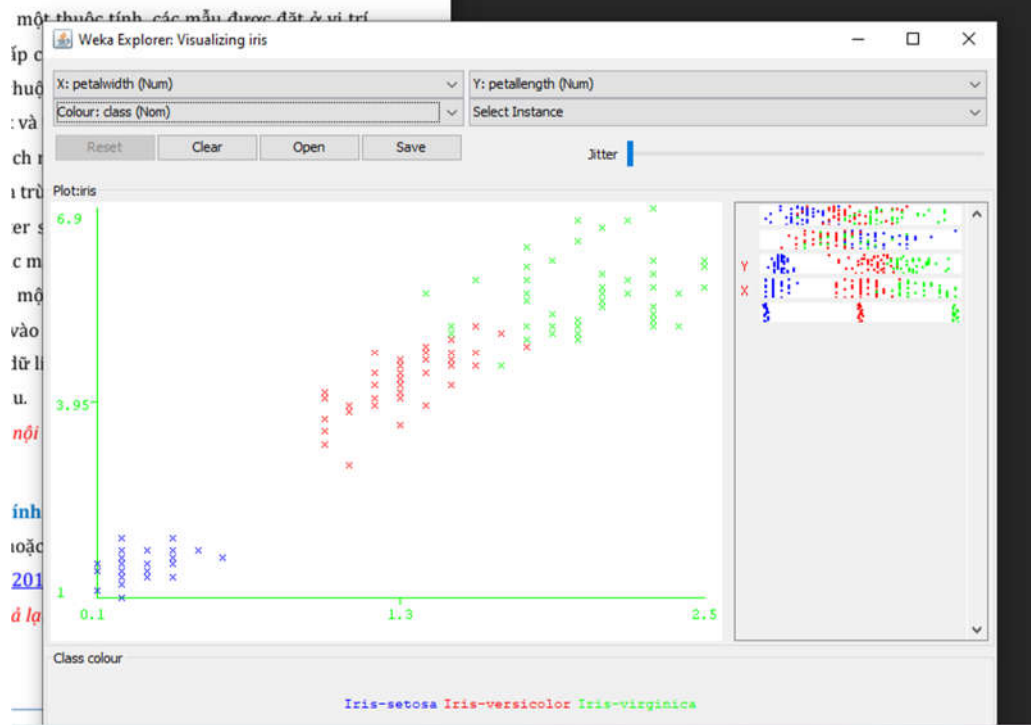
- Tập dữ liệu có bao nhiêu mẫu? 150 ( 50 thuộc tính mỗi 3 lớp)
- Bao nhiêu thuộc tính? 5 thuộc tính
- Miền giá trị của thuộc tính petallength là gì? 1.0 đến 6.9
- Tập dữ liệu có bao nhiêu thuộc tính số và bao nhiêu thuộc tính rời rạc? 4 thuộc tính số và 1 thuộc tính rời rạc

- Tên của thuộc tính lớp là gì? Iris Setosa, Iris Versicolour, Iris Virginica
- Đánh giá phân bố của các lớp, tức là cân bằng hay lệch về một lớp? Cân bằng (33% cho mỗi lớp)

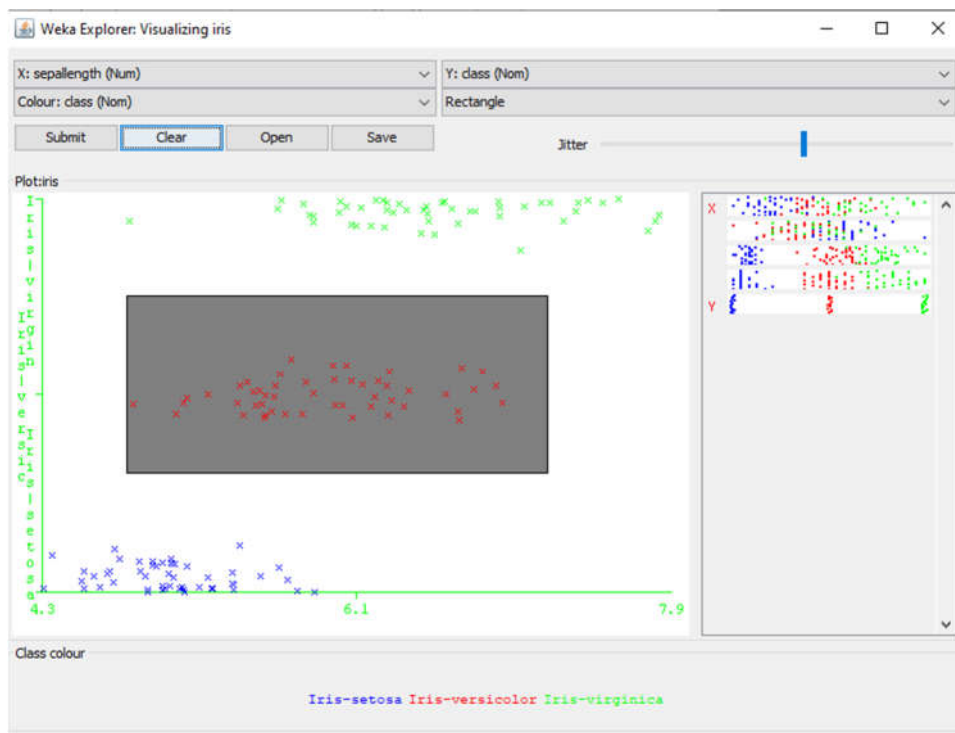


Đổi trục x thành petalwidth và trục y thành petallength.

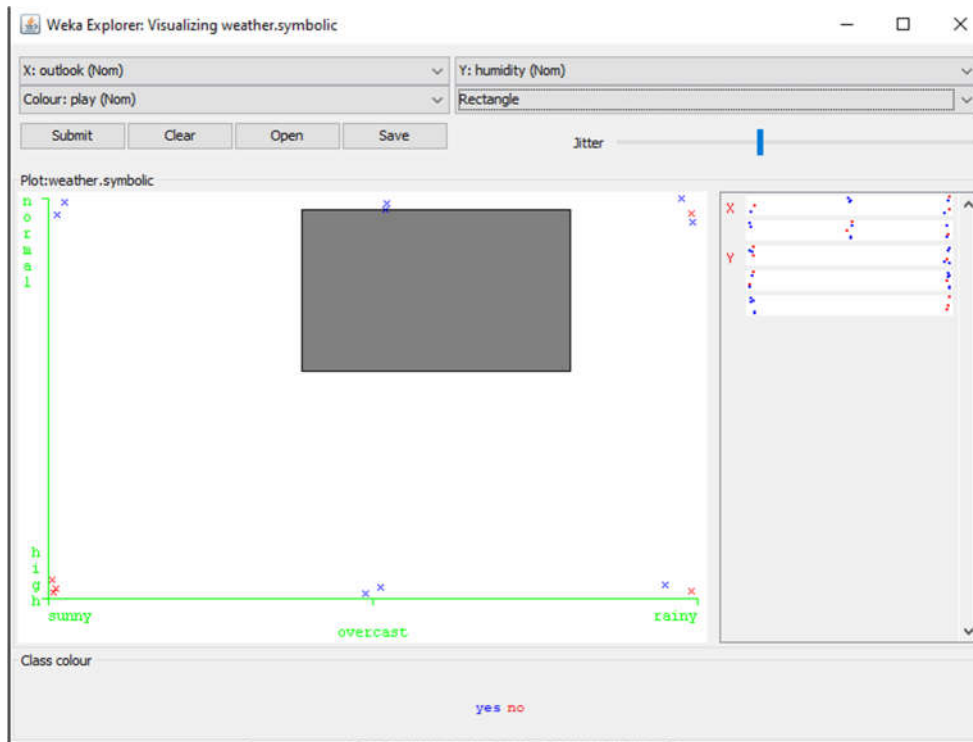
ig theo các bước trên để làm minh chứng.



## Tùy ý minh họa



## 10. Thực hiện lại những câu hỏi trong Yêu cầu 3.



## B – Nội dung thực hiện cài đặt (25 điểm)

### HƯỚNG DẪN CHẠY B1

```
File Edit Selection View Go Debug Terminal Help
1712713.py README.md
Data-Mining > README.md > # Bảng phân công task cho từng member trong mỗi tuần
Select Command Prompt
Microsoft Windows [Version 10.0.18362.256]
(c) 2019 Microsoft Corporation. All rights reserved.
C:\Users\PCIE>
E:\>cd E:\STUDENT\Năm 3\Khảo Thắc dữ liệu và ứng dụng\LT\BT\Git lap01\Data-Mining
E:\STUDENT\Năm 3\Khảo Thắc dữ liệu và ứng dụng\LT\BT\Git lap01\Data-Mining>python 1712713.py --input original.csv --output processed.csv --task min-max --nameList [passenger_numbers] --min 0 --max 1
E:\STUDENT\Năm 3\Khảo Thắc dữ liệu và ứng dụng\LT\BT\Git lap01\Data-Mining>python 1712713.py --input original.csv --output processed.csv --task z-score --nameList [passenger_numbers] --width 4 --depth 4
E:\STUDENT\Năm 3\Khảo Thắc dữ liệu và ứng dụng\LT\BT\Git lap01\Data-Mining>python 1712713.py --input original.csv --output processed.csv --task normMoth --nameList [passenger_numbers] --width 4 --depth 4
E:\STUDENT\Năm 3\Khảo Thắc dữ liệu và ứng dụng\LT\BT\Git lap01\Data-Mining>python 1712713.py --input original.csv --output processed.csv --task normDepth --nameList [passenger_numbers] --width 4 --depth 4
E:\STUDENT\Năm 3\Khảo Thắc dữ liệu và ứng dụng\LT\BT\Git lap01\Data-Mining>python 1712713.py --input original.csv --output processed.csv --task remove --nameList [passenger_numbers,Date]
E:\STUDENT\Năm 3\Khảo Thắc dữ liệu và ứng dụng\LT\BT\Git lap01\Data-Mining>python 1712713.py --input original.csv --output processed.csv --task insert --nameList [passenger_numbers,Date]
E:\STUDENT\Năm 3\Khảo Thắc dữ liệu và ứng dụng\LT\BT\Git lap01\Data-Mining>
python 1712713.py --input original.csv --output processed.csv --task min-max --nameList [passenger_numbers] --min 0 --max 1
python 1712713.py --input original.csv --output processed.csv --task z-score --nameList [passenger_numbers] --width 4 --depth 4
python 1712713.py --input original.csv --output processed.csv --task normMoth --nameList [passenger_numbers] --width 4 --depth 4
python 1712713.py --input original.csv --output processed.csv --task normDepth --nameList [passenger_numbers] --width 4 --depth 4
python 1712713.py --input original.csv --output processed.csv --task remove --nameList [passenger_numbers,Date]
python 1712713.py --input original.csv --output processed.csv --task insert --nameList [passenger_numbers,Date]
```

## KẾT QUẢ CHẠY

A1				A1				A1			
passenger_numbers				passenger_numbers				passenger_numbers			
passenger_numbers	Date			passenger_numbers	Date			passenger_numbers			
112	1/1/1949			112	1/1/1949			[104.0,115.0]			
118	2/1/1949			118	2/1/1949			[118.0,121.0]			
132	3/1/1949			132	3/1/1949			[125.0,132.0]			
129	4/1/1949			129	4/1/1949			[125.0,132.0]			
121	5/1/1949			148	7/1/1949			[118.0,121.0]			
157.2093023	6/1/1949			148	8/1/1949						
148	7/1/1949			119	10/1/1949			[145.0,148.0]			
148	8/1/1949			104	11/1/1949			[145.0,148.0]			
136	1/1/1949			118	12/1/1949			[135.0,141.0]			
119	10/1/1949			115	1/1/1950			[118.0,121.0]			
104	11/1/1949			126	2/1/1950			[104.0,115.0]			
118	12/1/1949			141	3/1/1950			[118.0,121.0]			
115	1/1/1950			125	5/1/1950			[104.0,115.0]			
126	2/1/1950			149	6/1/1950			[125.0,132.0]			
141	3/1/1950			170	7/1/1950			[135.0,141.0]			
135	1/1/1949			170	8/1/1950			[135.0,141.0]			
125	5/1/1950			158	9/1/1950			[125.0,132.0]			
149	6/1/1950			114	11/1/1950			[149.0,162.0]			
170	7/1/1950			140	12/1/1950			[163.0,171.0]			
170	8/1/1950			145	1/1/1951			[163.0,171.0]			
158	9/1/1950			150	2/1/1951			[149.0,162.0]			
157.2093023	10/1/1950			163	4/1/1951						
114	11/1/1950			172	5/1/1951			[104.0,115.0]			
140	12/1/1950			178	6/1/1951			[135.0,141.0]			
145	1/1/1951			199	8/1/1951			[145.0,148.0]			
150	2/1/1951			184	9/1/1951			[149.0,162.0]			
178	1/1/1949							[172.0,178.0]			

A1				A1				A1			
passenger_numbers				passenger_numbers				passenger_numbers			
passenger_numbers				passenger_numbers				passenger_numbers			
[104.0,138.5]				-1.383305132				1.057971014			
[104.0,138.5]				-1.19971834				1.101449275			
[104.0,138.5]				-0.771349158				1.202898551			
[104.0,138.5]				-0.863142554				1.18115942			
[104.0,138.5]				-1.107924944				1.123188406			
[138.5,173.0]				-0.281784379				1.31884058			
[138.5,173.0]				-0.281784379				1.31884058			
[104.0,138.5]				-0.648957963				1.231884058			
[104.0,138.5]				-1.169120541				1.108695652			
[104.0,138.5]				-1.628087522				1			
[104.0,138.5]				-1.19971834				1.101449275			
[104.0,138.5]				-1.291511736				1.079710145			
[138.5,173.0]				-0.95493595				1.15942029			
[104.0,138.5]				-0.49596897				1.268115942			
[104.0,138.5]				-0.679555762				1.224637681			
[138.5,173.0]				-0.985533749				1.152173913			
[138.5,173.0]				-0.25118658				1.326086957			
[138.5,173.0]				0.391367193				1.47826087			
[138.5,173.0]				0.391367193				1.47826087			
				0.024193608				1.391304348			
[104.0,138.5]											
[138.5,173.0]				-1.322109535				1.072463768			
[138.5,173.0]				-0.526566768				1.260869565			
[138.5,173.0]				-0.373577775				1.297101449			
[173.0,207.5]				-0.220588781				1.333333333			
				0.636149582				1.536231884			

c) b) a)

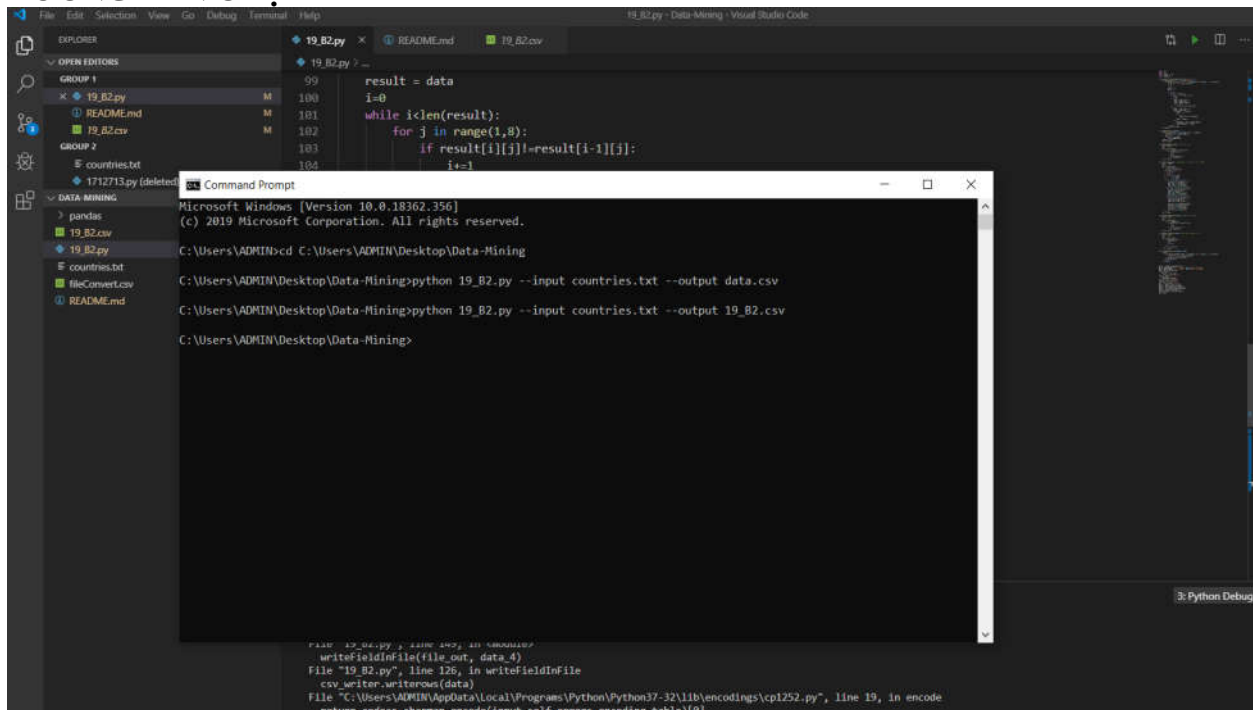
## CÁC LỆNH COMMAND LINE

a) `python 19_B1.py --input original.csv --output processed.csv --task min-max --nameList [passenger_numbers] --min 0 --max 1`

b) `python 19_B1.py --input original.csv --output processed.csv --task z-score --nameList [passenger_numbers]`

- c) python 19\_B1.py --input original.csv --output processed.csv --task normWidth --nameList [passenger\_numbers] --width 4
- d) python 19\_B1.py --input original.csv --output processed.csv --task normDepth --nameList [passenger\_numbers] --depth 4
- e) python 19\_B1.py --input original.csv --output processed.csv --task remove --nameList [passenger\_numbers,Date]
- f) python 19\_B1.py --input original.csv --output processed.csv --task insert --nameList [passenger\_numbers,Date]

## HƯỚNG DẪN CHẠY B2





## KẾT QUẢ CHẠY

	A	B	C	D	E	F	G	H	I	J
1	country	name	longName	foundinD	population	capital	largestCity	area		
2	14	Abkhazia	Republic o	#####	242862	Sukhumi		8660.968660968661km		
3	15	Abkhazia	Republic o	#####	242862	Sukhumi		8660.968660968661km		
4	16	Abkhazia	Republic o	#####	242862	Sukhumi		8660.968660968661km		
5	17	Abyei	Abyei Area	1/9/2005		Abyei (town)		10546.490546490546km		
6	18	Abyei	Abyei Area	#####		Abyei (town)		10546.490546490546km		
7	31	Adã©lie L	Adã©lie Land					432000km		
8	36	Adjara	Autonomous Republic		393700	Batumi		2900.8029008029007km		
9	38	Aerica	Aerican En	5/8/1987		Montreal		9000000km		
10	41	Afghanista	Islamic Re	#####	32564342	Kabul	Kabul	652232.5822325823km		
11	42	Afghanista	font-size:8	#####	32564342	Kabul	Kabul	652232.5822325823km		
12	44			6/3/1991	1.05E+09		Nigeria	29865860km		
13	45			6/3/1991	1.05E+09		Lagos	29865860km		
14	63	Akrotiri an	Sovereign Base Areas		7700	Episkopi Cantonment		253.8202538202538km		
15	64	Akrotiri an of			7700	Episkopi Cantonment		253.8202538202538km		
16	65	Akrotiri an	Akrotiri and Dhekelia		7700	Episkopi Cantonment		253.8202538202538km		
17	68	Ä...land	Ä...land Isl	5/7/1920	28666	Mariehamn		1579.90157990158km		
18	73		Republic o	4/7/1939	2893005	Tirana	Tirana	28749.02874902875km		
19	74		Republic o	#####	2893005	Tirana	Tirana	28749.02874902875km		
20	82	Alderney	Alderney		2013	Saint Anne, Alderney		7.77000777000777km		
21	84		People's D	7/3/1962	40400000	Algiers		2381753.431753432km		
22	85		People's D	7/5/1962	40400000	Algiers		2381753.431753432km		
23	106	American	Territory c	1889-06-1	54343	Pago Pago		198.989898989899km		
24	107	American	Territory c	1899-12-0	54343	Pago Pago		198.989898989899km		
25	116		Principality of Andorra		85470	Andorra la Vella		467.62496762496767km		
26	124		Republic o	#####		Luanda		1246708.1067081068km		
27	125	Anguilla	Anguilla		13600	The Valley, Anguilla		91km		
28	131	Anjouan	Autonomous Island o		277500	Mutsamudu		122.17042217042217km		
29	134		Antigua an	#####	91295	St. John's, Antigua and		440.3004403004403km		
30	135		Antigua an	#####	91295	St. John's, Antigua and		440.3004403004403km		

## CÁC LỆNH COMMAND LINE

```
python 19_B2.py --input countries.txt --output 19_B2.csv
```