

1. Unsupervised یک روش یادگیری ماشین است که به مدل اجازه می‌دهد تا روی الگوها و اطلاعاتی که قبلاً پیدا نشده بودند کار کند، این کار با داده‌های بدون label است.

2. یادگیری supervised روش یادگیری ماشین است که با استفاده از داده‌های آموزشی دارای برچسب استنباط می‌کند که داده تست یا داده‌ی جدید شامل کدام مجموعه‌های آموزش است.

3. ترکیب دو تایی بالایی است یعنی مقداری داده‌ی label خورده و مقدار بیشتری داده‌ی label نخورده. برای classify ام استفاده می‌شود.

4. Outlier به داده‌هایی می‌گویند که خارج از انتظار ما از پراکندگی داده‌ها (بیشتر داده‌ها)، هست.

5. تعداد متغیرهای ورودی یا ویژگی‌های یک مجموعه داده به عنوان بعد dimension تعریف می‌شود.

6. Training Dataset: نمونه داده‌ها تا مدل با آن fit شود و مدل آموزش داده شود، مدل داده‌ها را می‌بیند و یاد می‌گیرد.

Validation Dataset: نمونه داده‌ای از آموزش نگه داشته می‌شود به طور جداگانه برای استفاده در ارزیابی از یک مدل که با مجموعه داده‌های آموزش fit شده است، همزمان با تنظیم hyperparameterهای مدل.

Test Dataset: نمونه‌ای از داده‌های خارجی برای ارائه ارزیابی از مدل نهایی fit شده مجموعه داده‌های آموزشی استفاده شده است.

7. Data warehousing فرآیند ساخت و داده با ادغام داده های چندین منبع تا توانایی گزارش تحلیلی ، query و تصمیم گیرید داشته باشد.

8. داده های دنیای واقعی اغلب مقدار زیادی داده ی از دست رفته دارند. علت از دست رفتن مقادیر می تواند خرابی داده ها یا ضبط نکردن داده ها باشد. مدیریت داده های از دست رفته در طی پیش پردازش مجموعه داده بسیار مهم است.

9. Independent Variable: متغیر مستقل به عنوان متغیری که در یک تست تغییر یا کنترل می شود، تعریف می شود. این نشان دهنده علت یا دلیل یک نتیجه است. متغیرهای مستقل متغیرهایی هستند که تست کننده برای آزمایش متغیر وابسته خود تغییر می دهد.

(2) 1. Missing Values Ratio: بعید است ستون های داده با مقادیر زیادی از دست رفته اطلاعات مفیدی را حمل کنند. بنابراین ستون های داده با تعداد مقادیر از دست رفته بیشتر از یک آستانه داده شده را می توان حذف کرد. مثال:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
train=pd.read_csv("Train_UWu5bXk.csv")
# checking the percentage of missing values in each variable
train.isnull().sum()/len(train)*100
```

Item_Identifier	0.000000
Item_Weight	17.165317
Item_Fat_Content	0.000000
Item_Visibility	0.000000
Item_Type	0.000000
Item_MRP	0.000000
Outlet_Identifier	0.000000
Outlet_Establishment_Year	0.000000
Outlet_Size	28.276428
Outlet_Location_Type	0.000000
Outlet_Type	0.000000
Item_Outlet_Sales	0.000000
dtype:	float64

saving missing values in a variable

```
a = train.isnull().sum()/len(train)*100
```

saving column names in a variable

```
variables = train.columns
```

```
variable = [ ]
```

```
for i in range(0,12):
```

```
    if a[i]<=20: #setting the threshold as 20%
```

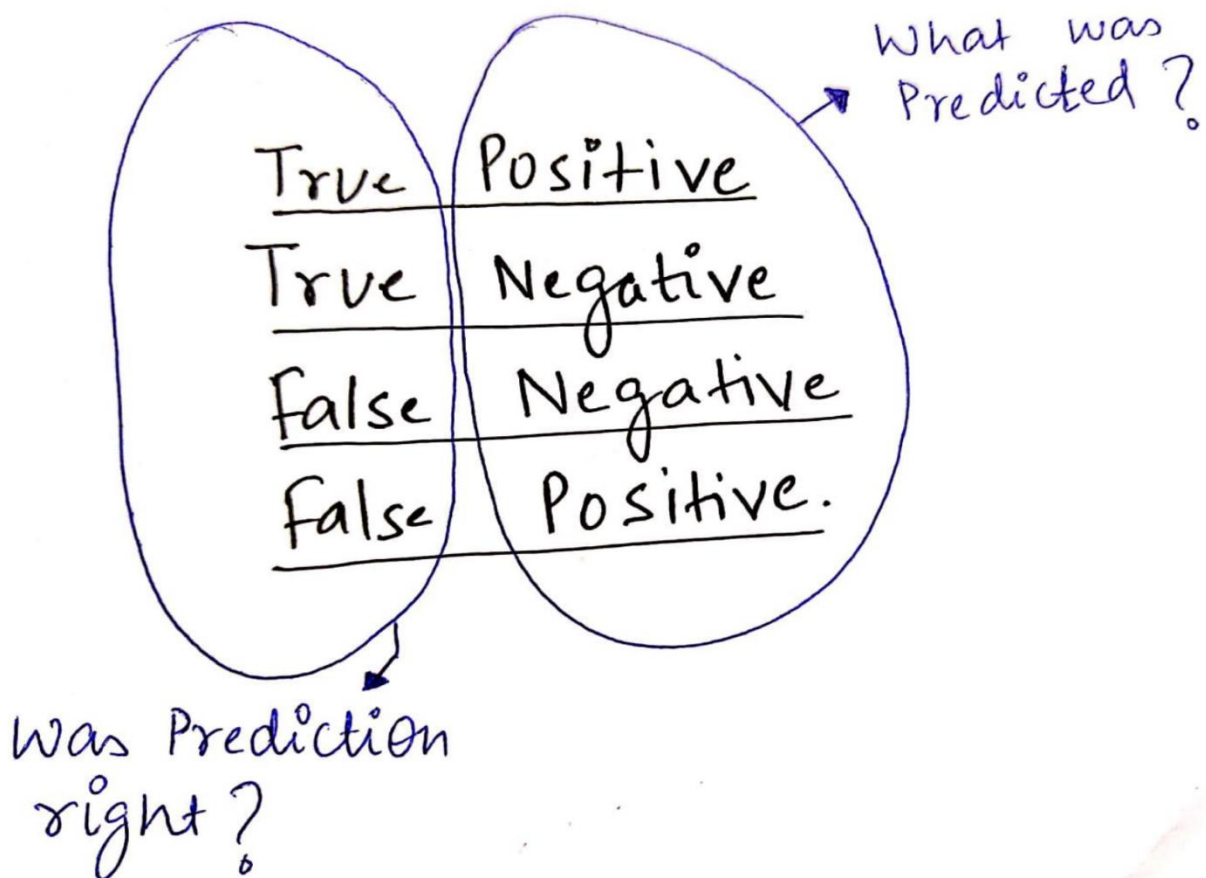
```
        variable.append(variables[i])
```

2. Low Variance Filter: به طور مشابه روش قبلی ، ستون های داده با کمی تغییر در داده ها ، اطلاعات با ارزش کمتری را حمل می کنند. بنابراین تمام ستون های داده با واریانس پایین تر از یک آستانه داده شده حذف می شوند. Variance به دامنه بستگی دارد. بنابراین قبل از استفاده از این روش normalization لازم است.

Feature Selection: فرایندی است که در آن به طور خودکار یا دستی آن ویژگی‌هایی را انتخاب می‌شود که بیشترین تأثیر را در متغیر پیش‌بینی یا خروجی مورد نظر دارند. داشتن ویژگی‌های بی‌ربط در داده‌های می‌تواند از دقت مدل‌ها بکاهد و مدل را بر اساس ویژگی‌های بی‌ربط بیاموزد.

Feature extraction: استخراج ویژگی فرایندی در کاهش ابعاد است که طی آن یک مجموعه اولیه از داده‌های خام به گروه‌های قابل کنترل بیشتری برای پردازش کاهش می‌یابد. ویژگی این مجموعه داده‌های بزرگ تعداد زیادی متغیر است که برای پردازش به منابع محاسباتی زیادی نیاز دارند.

3.



دقت: نسبت پیش بینی های مثبت صحیح به کل موارد مثبت پیش بینی شده.

$$P = \frac{TP}{TP + FP}$$

فراخوانی: نسبت پیش بینی های مثبت صحیح به کل نمونه های مثبت.

$$R = \frac{TP}{TP + FN}$$

F-score: دقت و فراخوانی را می توان در یک معیار واحد به نام f-score ترکیب کرد. اگر تأکیدی زیادی بر روی اهمیت بیشتر دقت یا فراخوانی نباشد، ترکیب مناسبی ست. فرمول زیر ترکیبی از هر دو معیار است.

$$F1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

4. همبستگی صفر نشان می دهد که دو متغیر رابطه ی خطی ندارند و معنای عدم وجود رابطه نیست. همبستگی صفر اغلب با استفاده از مخفف $r = 0$ نشان داده می شود.

درباره ی مستقل بودن آن ها نمیتوان به قطعیت گفت چون مستقل بودن یعنی هیچگونه رابطه ای با هم ندارند.

5.

پاکسازی داده ها: پر کردن مقادیر از دست رفته ، مرتب کردن داده های noisy ، شناسایی یا حذف کردن داده های noisy و outlier.

ادغام داده ها: ادغام چندین پایگاه داده یا file تا به کاربر ها یک view از تمام این داده ها نشان بدهد.

تبدیل داده ها: در data warehouse ، ما از تبدیل داده برای تبدیل داده ها از قالب داده منبع به داده های مقصد استفاده می کنیم مانند:

Data mapping: map کردن داده های منبع به داده های دیگر با تابع.

Smoothing: حذف داده های noisy

Normalization: scale کردن داده ها به یک range خاص.

Aggregation: ادغام داده ها.

Generalization: داده های سطح پایین مثل سن با مفاهیم سطح بالاتر مثل جوانی یا تجربه جایگزین می شوند.

6.

f1

Confidence = 60% , Support_Count $\rightarrow 6 \times 33\% = 2$

Sup_Count	
نان	2
الوجه	2
نهر	2
الوجه	3
نهر	2

k=1

Sup_Count > 2 \rightarrow

Sup_Count	
نان، الوجه	2
نان، نهر	1
نان، الوجه	2
نان، نهر	2
الوجه، نهر	1
الوجه، الوجه	0
الوجه، نهر	0
نهر، الوجه	0
نهر، نهر	1
نهر، الوجه	3

نان، الوجه	2
نان، نهر	2
نان، الوجه	2
نهر، الوجه	3

k=3

Subsets \rightarrow

not frequent

نان، الوجه	2
------------	---

k=4

P4PCG

نسبة اضمحلال

$$[\text{نهر، نان}] \rightarrow [\text{نهر}] = \frac{2}{2} = 100\% \quad (1)$$

$$[\text{نهر، نان}] \rightarrow [\text{نان}] = \frac{2}{2} = 100\% \quad (1)$$

$$[\text{نهر، نان}] \rightarrow [\text{الوجه}] = \frac{2}{3} = 66\% \quad (2)$$

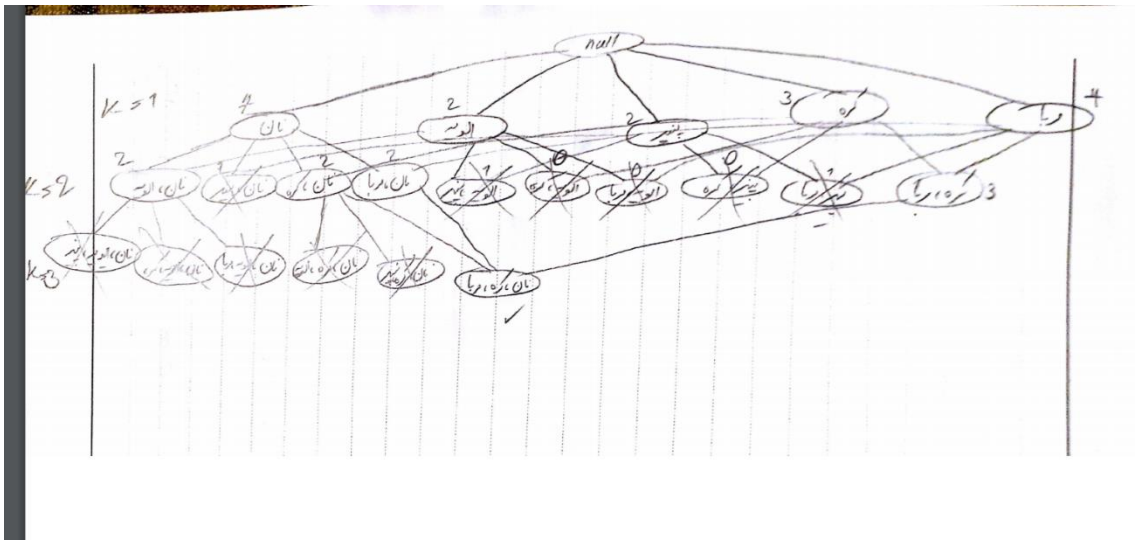
$$[\text{نهر}] \rightarrow [\text{نان، نهر}] = \frac{2}{2} = 100\% \quad (1)$$

$$[\text{نهر}] \rightarrow [\text{نان، الوجه}] = \frac{2}{2} = 100\% \quad (1)$$

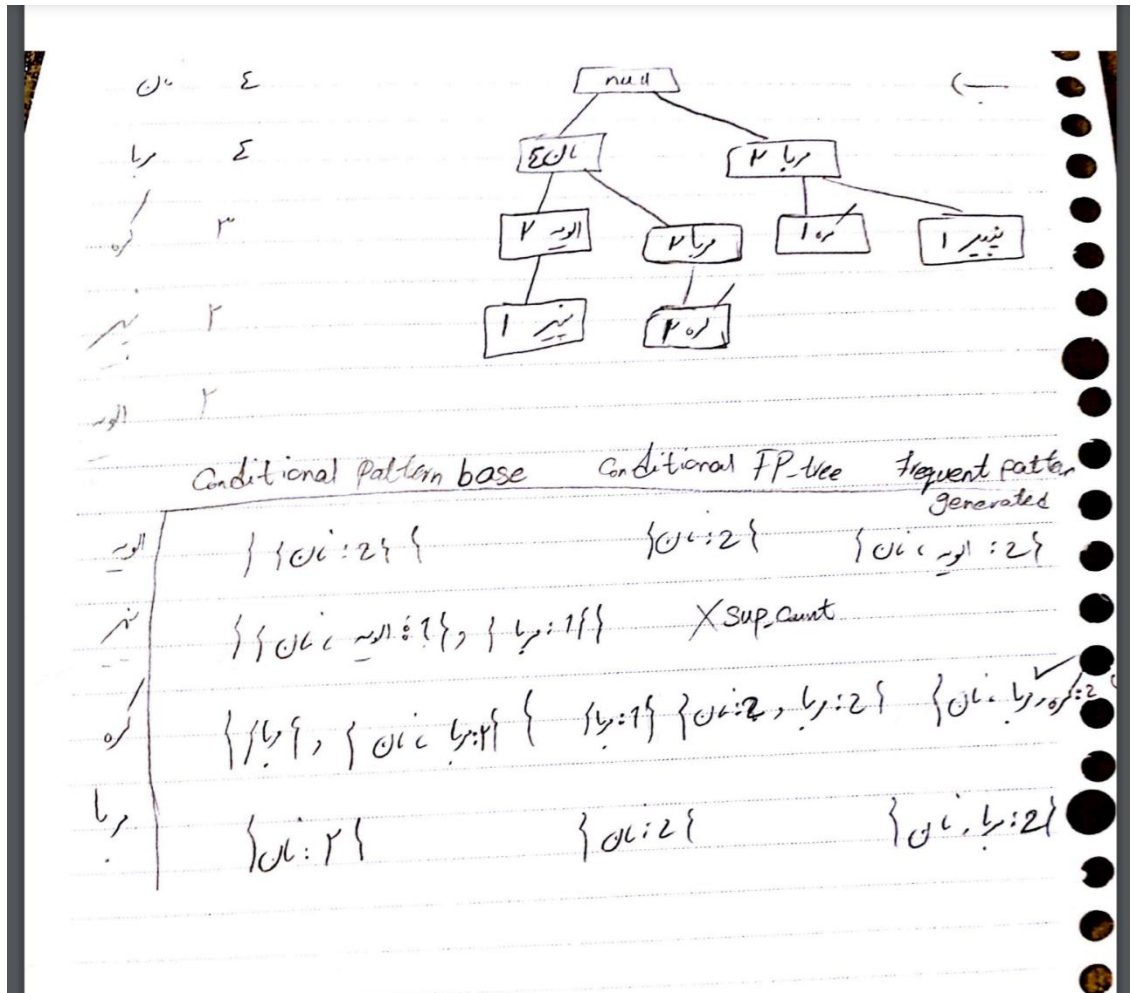
$$[\text{نان}] \rightarrow [\text{نهر، نان}] = \frac{2}{3} = 66\% \quad (2)$$

هذه النسبة

7. الف



(ب)



	Conditional Pattern base	Conditional FP-tree	Frequent pattern generated
الوجه	{الوجه: 2}	{الوجه: 2}	{الوجه: 2}
مربا	{مربا: 1}	X Sup. Count	
مربا	{مربا: 2}	{مربا: 2}	{مربا: 2}
مربا	{مربا: 2}	{مربا: 2}	{مربا: 2}

بخش عملی:

1.

The screenshot shows the Spyder Python IDE interface. The editor on the left contains a script named `temp.py` with the following code:

```
1 # -*- coding: utf-8 -*-
2 """
3 Spyder Editor
4 """
5 This is a temporary script file.
6 """
7 import pandas as pd
8 import seaborn as sns
9
10 data_set = pd.read_csv("pLayers.csv")
11 print(data_set.head())
12 print(data_set.tail())
```

The console on the right shows the output of the script. It displays the first 5 rows of the `pLayers.csv` file, which contains 90 columns. The output is as follows:

```
[3 rows x 90 columns]
In [6]: runfile('C:/Users/Amirhosein Khoshbin/Desktop/New folder/temp.py', wdir='C:/Users/Amirhosein Khoshbin/Desktop/New folder')
ID      Name  ... RBRating  GKRating
0  158073  L. Messi  ...      65         23
1  288801  Cristiano Ronaldo  ...      64         23
2  280389  J. Oblak  ...      35         92
3  192985  K. De Bruyne  ...      78         24
4  190871  Neymar Jr  ...      65         23

[5 rows x 90 columns]
ID      Name  FullName  ... CBRating  RBRating  GKRating
19015  257371  M. Nzongong  Mike Nzongong  ...      48         42         18
19016  259160  L. Bell  Lewis Bell  ...      35         41         13
19017  259157  Y. Arai  Yassin Arai  ...      39         44         17
19018  253763  R. Dinanga  Ricardo Dinanga  ...      30         34         16
19019  241493  S. Cartwright  Samuel Cartwright  ...      51         47         15

[5 rows x 90 columns]
In [7]:
```

2.

The screenshot shows the Spyder Python IDE interface. The editor on the left contains a script named `temp.py` with the following code:

```
1 # -*- coding: utf-8 -*-
2 """
3 Spyder Editor
4 """
5 This is a temporary script file.
6 """
7 import pandas as pd
8 import seaborn as sns
9
10 data_set = pd.read_csv("pLayers.csv")
11 #print(data_set.head())
12 #print(data_set.tail())
13 pd.set_option('display.max_columns', 5)
14 print(data_set[data_set.isna().any(axis=1)])
```

The console on the right shows the output of the script. It displays the first 5 rows of the `pLayers.csv` file, which contains 90 columns. The output is as follows:

```
raise ValueError("Must provide an even number of non-keyword arguments")
ValueError: Must provide an even number of non-keyword arguments

In [32]: runfile('C:/Users/Amirhosein Khoshbin/Desktop/New folder/temp.py', wdir='C:/Users/Amirhosein Khoshbin/Desktop/New folder')
ID      Name  ... RBRating  GKRating
4  190871  Neymar Jr  ...      65         23
6  208722  S. Mané  ...      69         23
11  212831  Alisson  ...      33         91
14  200145  Casenro  ...      84         24
16  105153  K. Benzema  ...      62         21
...      ...      ...      ...
19015  257371  M. Nzongong  ...      42         18
19016  259160  L. Bell  ...      41         13
19017  259157  Y. Arai  ...      44         17
19018  253763  R. Dinanga  ...      34         16
19019  241493  S. Cartwright  ...      47         15

[18118 rows x 90 columns]
In [33]:
```

3.

```

1 # -*- coding: utf-8 -*-
2 """
3 Spyder Editor
4 This is a temporary script file.
5 """
6
7 import pandas as pd
8 import seaborn as sns
9
10 data_set = pd.read_csv("players.csv")
11 #print(data_set.head())
12 #print(data_set.tail())
13 #pd.set_option('display.max_columns', 5)
14 #print(data_set[data_set.isna().any(axis=1)])
15 print("mean", data_set["weight"].mean())
16 print("max", data_set["weight"].max())
17 print("min", data_set["weight"].min())
  
```

Console I/O

```

75.05241850683491

In [34]: runfile('C:/Users/Amirhosein Khoshbin/Desktop/New folder/temp.py', wdir='C:/Users/Amirhosein Khoshbin/Desktop/New folder')
Traceback (most recent call last):

  File "C:/Users/Amirhosein Khoshbin/Desktop/New folder/temp.py", line 15, in <module>
    print("mean" + data_set["weight"].mean())
TypeError: can only concatenate str (not "float") to str

In [35]: runfile('C:/Users/Amirhosein Khoshbin/Desktop/New folder/temp.py', wdir='C:/Users/Amirhosein Khoshbin/Desktop/New folder')
mean 75.05241850683491

In [36]: runfile('C:/Users/Amirhosein Khoshbin/Desktop/New folder/temp.py', wdir='C:/Users/Amirhosein Khoshbin/Desktop/New folder')
mean 75.05241850683491
max 110
min 50

In [37]:
  
```

4. انگلیس: 1706، و تعداد زیادی از کشور ها که پایین مشخص است 1.

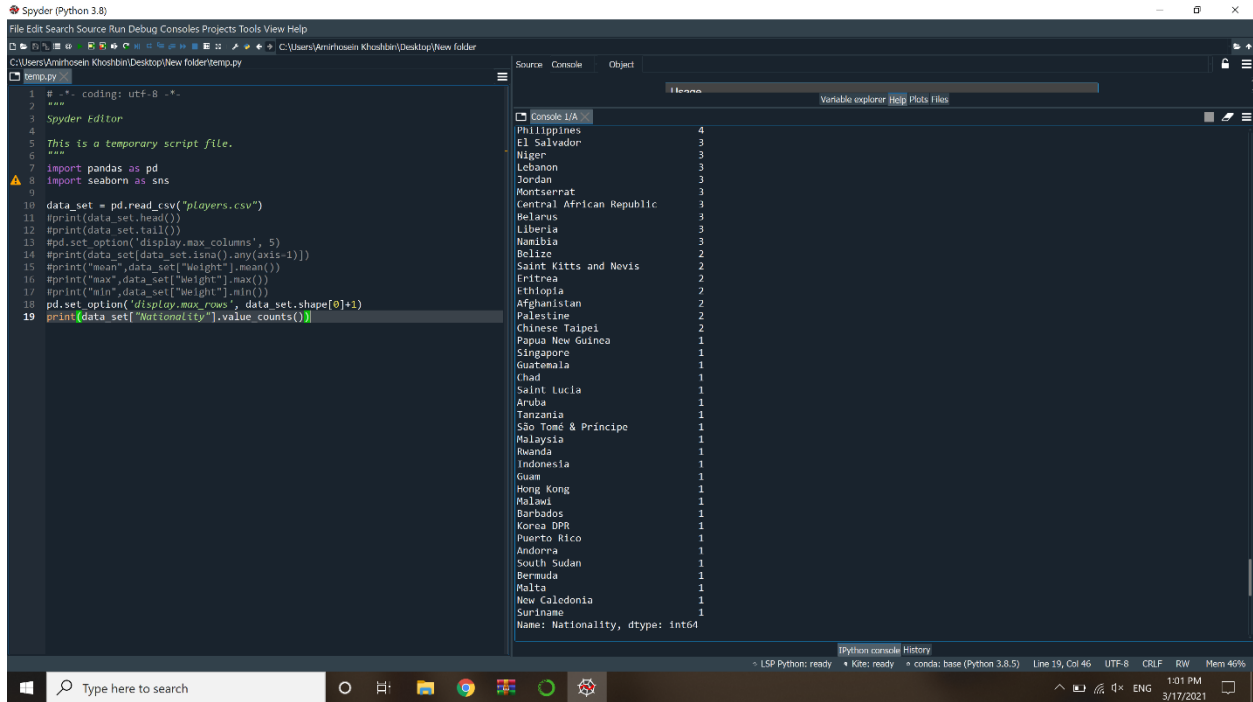
```

1 # -*- coding: utf-8 -*-
2 """
3 Spyder Editor
4 This is a temporary script file.
5 """
6
7 import pandas as pd
8 import seaborn as sns
9
10 data_set = pd.read_csv("players.csv")
11 #print(data_set.head())
12 #print(data_set.tail())
13 #pd.set_option('display.max_columns', 5)
14 #print(data_set[data_set.isna().any(axis=1)])
15 print("mean", data_set["weight"].mean())
16 print("max", data_set["weight"].max())
17 print("min", data_set["weight"].min())
18 pd.set_option('display.max_rows', data_set.shape[0]+1)
19 print(data_set["Nationality"].value_counts())
  
```

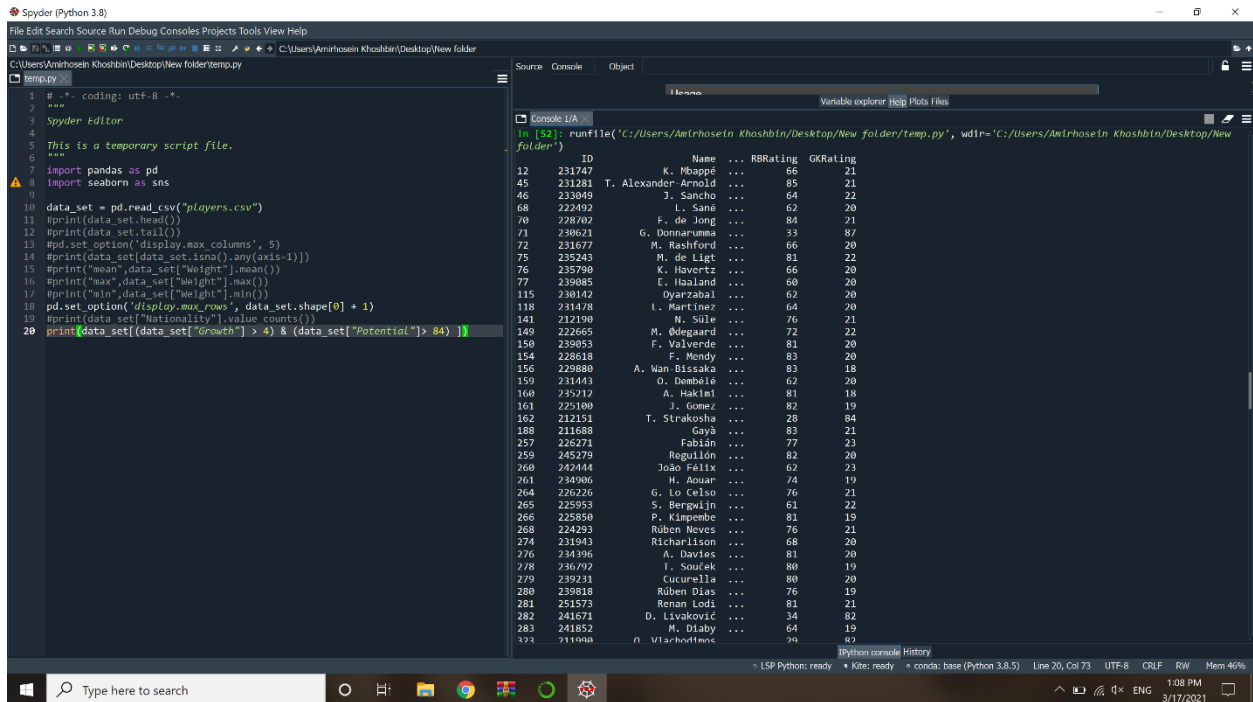
Console I/O

```

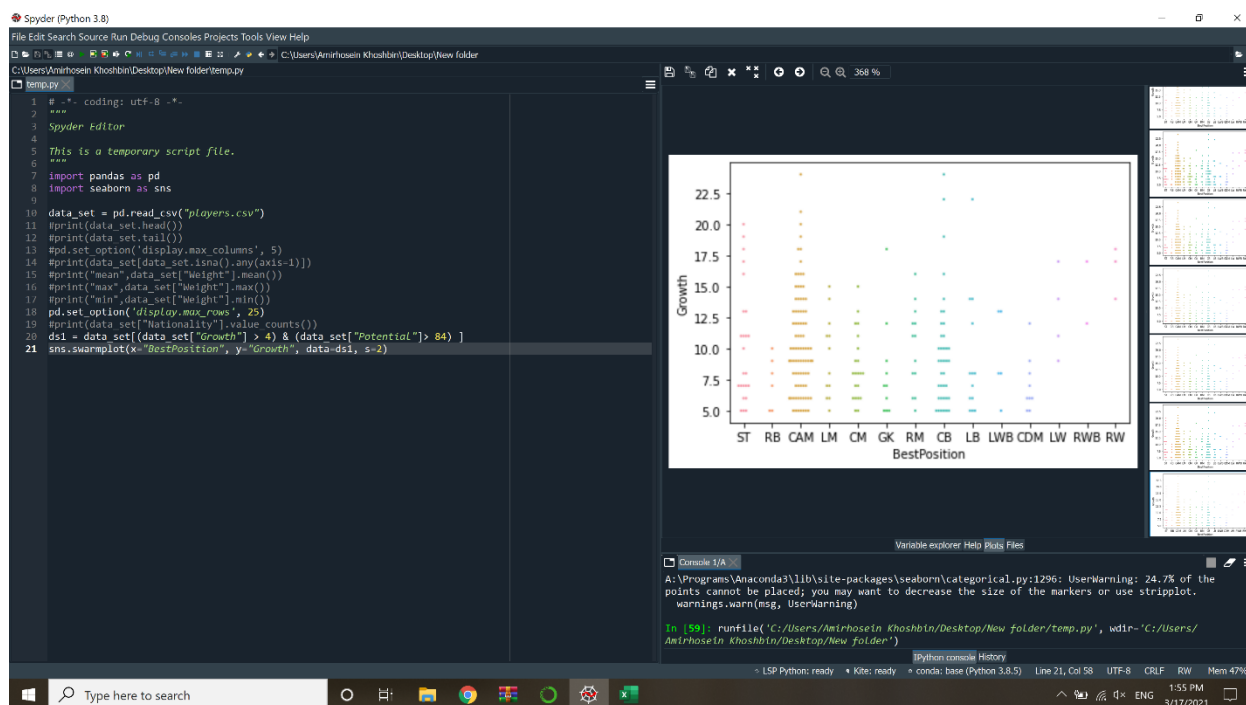
In [45]: runfile('C:/Users/Amirhosein Khoshbin/Desktop/New folder/temp.py', wdir='C:/Users/Amirhosein Khoshbin/Desktop/New folder')
Nationality
England      1706
Germany      1190
Spain        1084
France       1008
Argentina     945
Brazil        896
Japan         488
Netherlands  465
United States 390
Sweden        384
Italy         376
Mexico        366
Portugal      362
Norway        360
Poland        355
Uruguay       346
Korea Republic 342
Colombia      335
Republic of Ireland 334
China PR      330
Austria       329
Saudi Arabia  314
Turkey        314
Romania       310
Denmark       305
Belgium       299
Scotland      282
Ecuador       247
Australia     246
Paraguay      238
Switzerland   208
Venezuela     200
Chile         182
Peru          163
Bolivia       141
Senegal       138
  
```



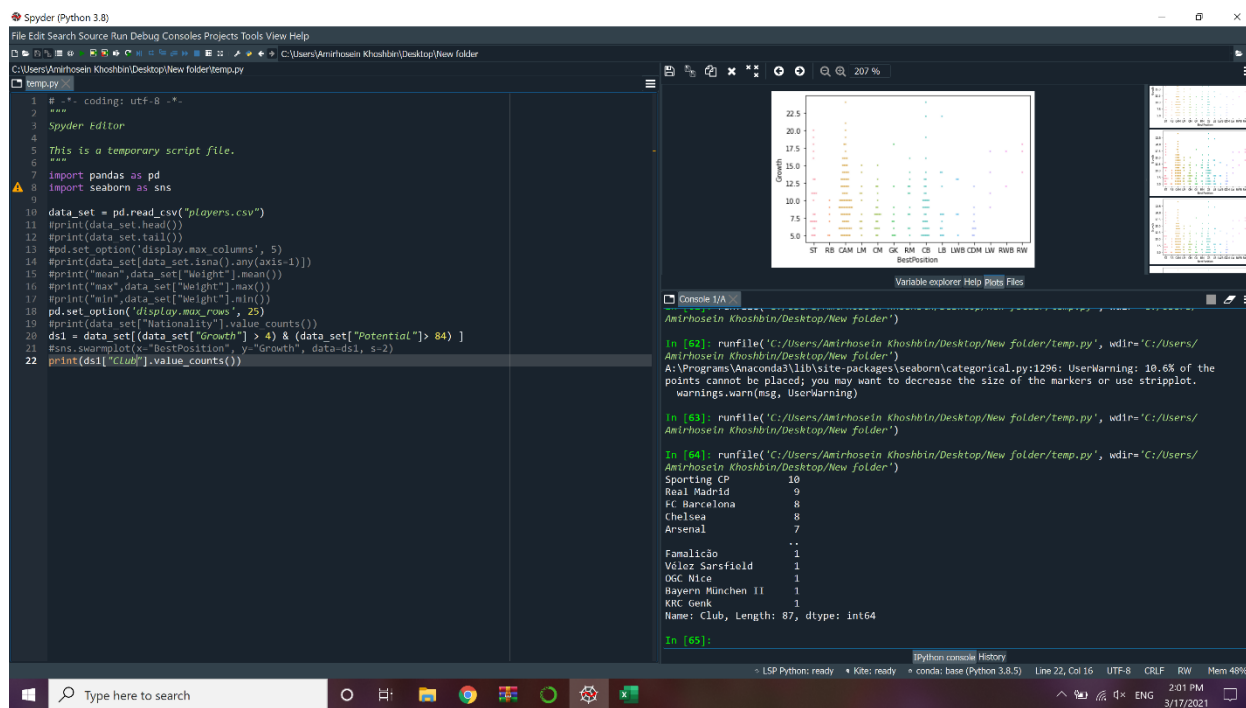
5. مشخصه به ترتیب، امپایه، ارنولد، سانچو و...



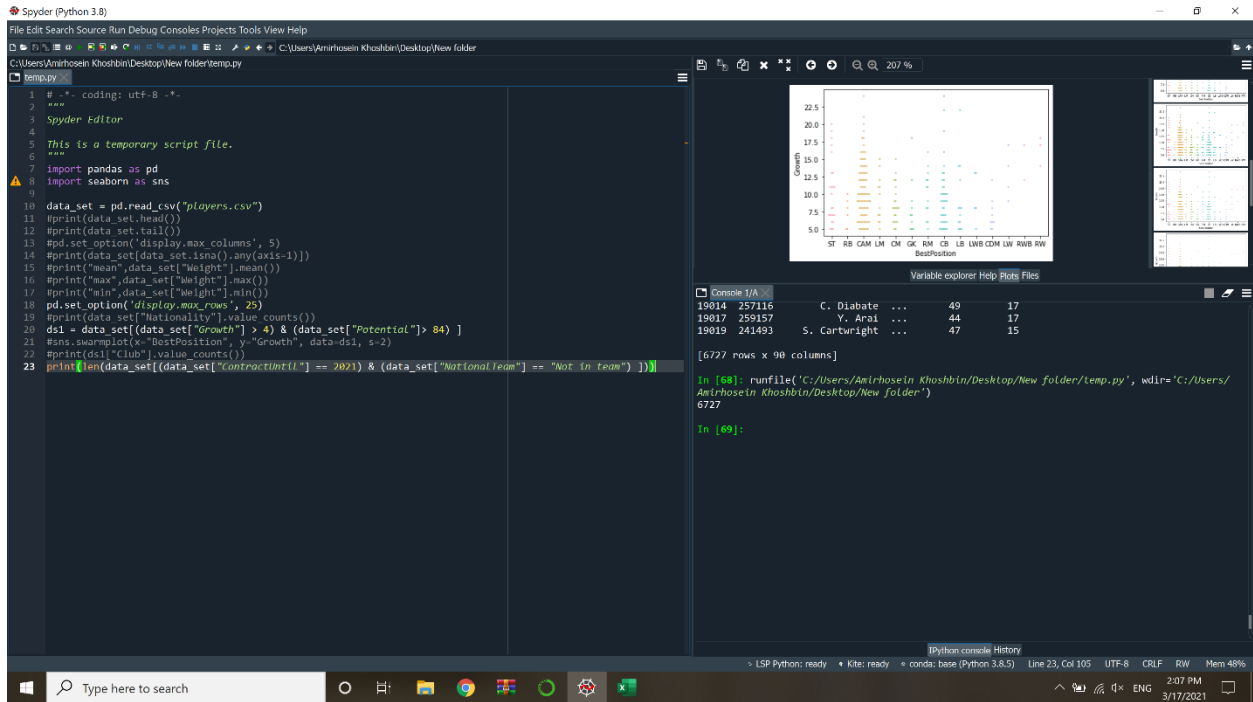
6.



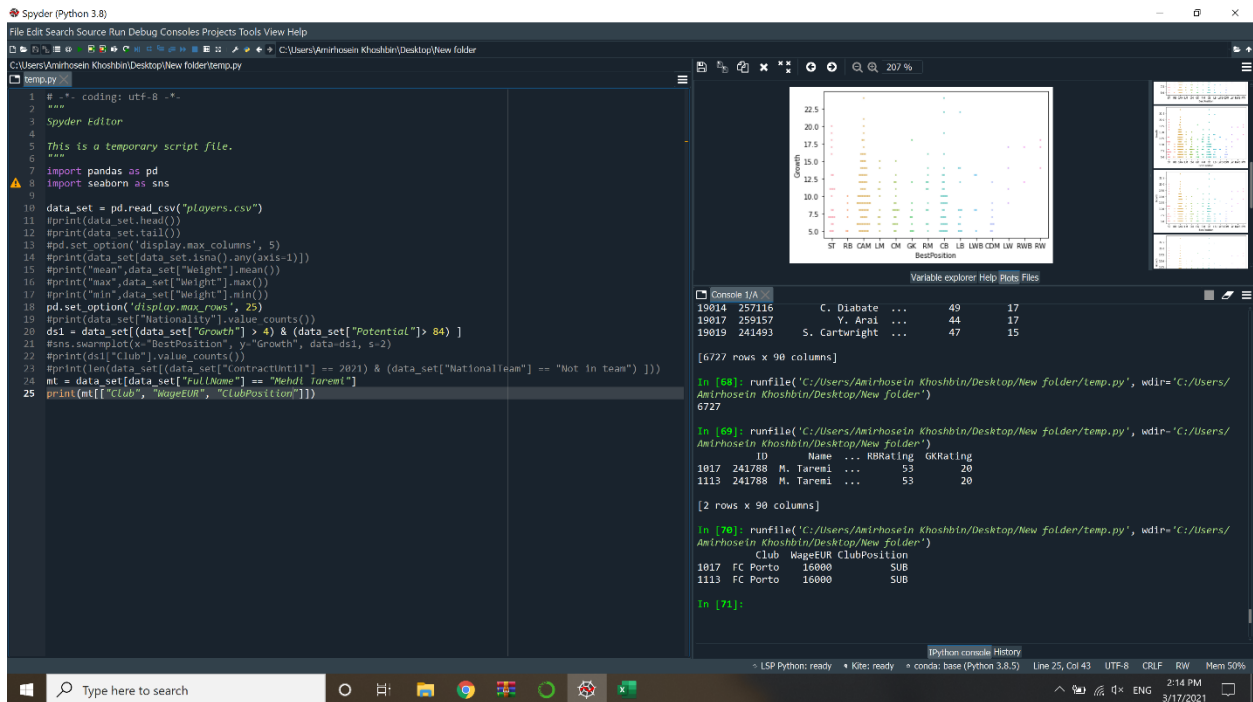
7. اسپورتینگ لیسبون با 10 بازیکن



9. 6727 تعداد این بازیکنان است.



10. موقعیت SUB، درآمد 16000 باشگاه پورتو



2. upperBoundMinSupport: الگوریتم با این معیار شروع میکند 100%
= 1.0 و با delta در هر دور کم میکند support رو. در واقع حد بالا برای حداقل support.

Delta: معیاری که در هر دور با آن حداقل support کم میشود.

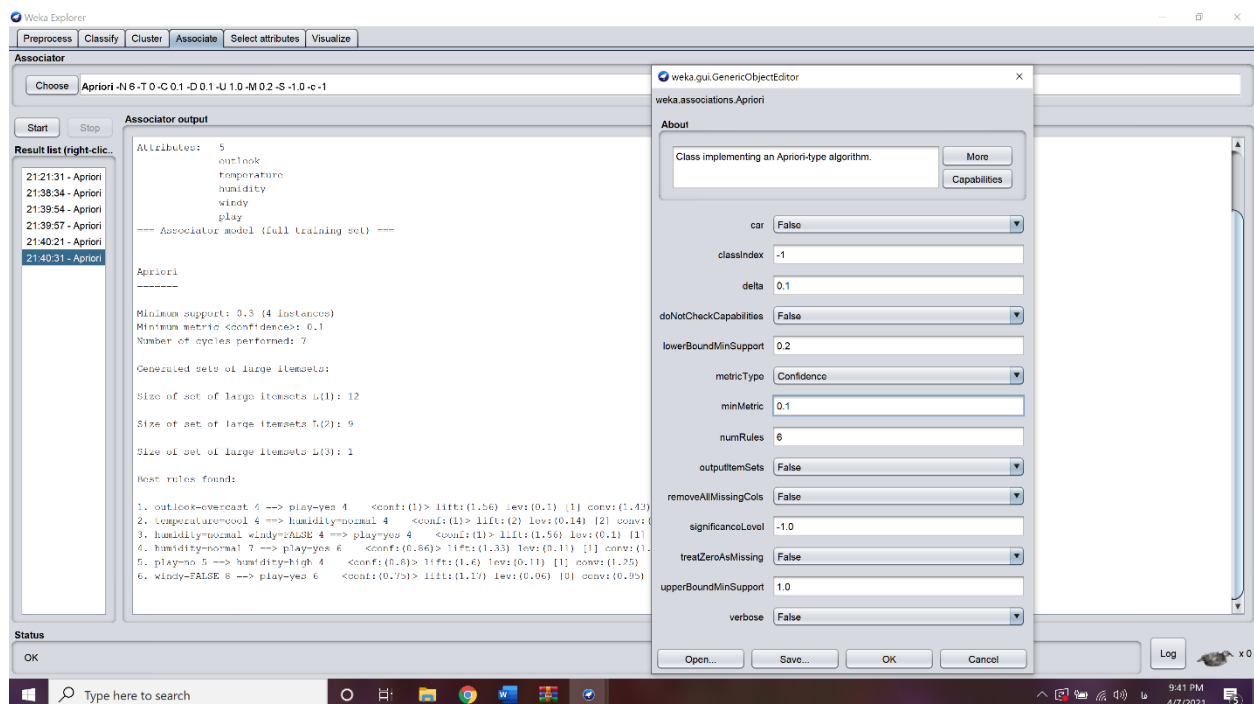
lowerBoundMinSupport: الگوریتم زمانی که support به این مقدار برسد، متوقف میشود.

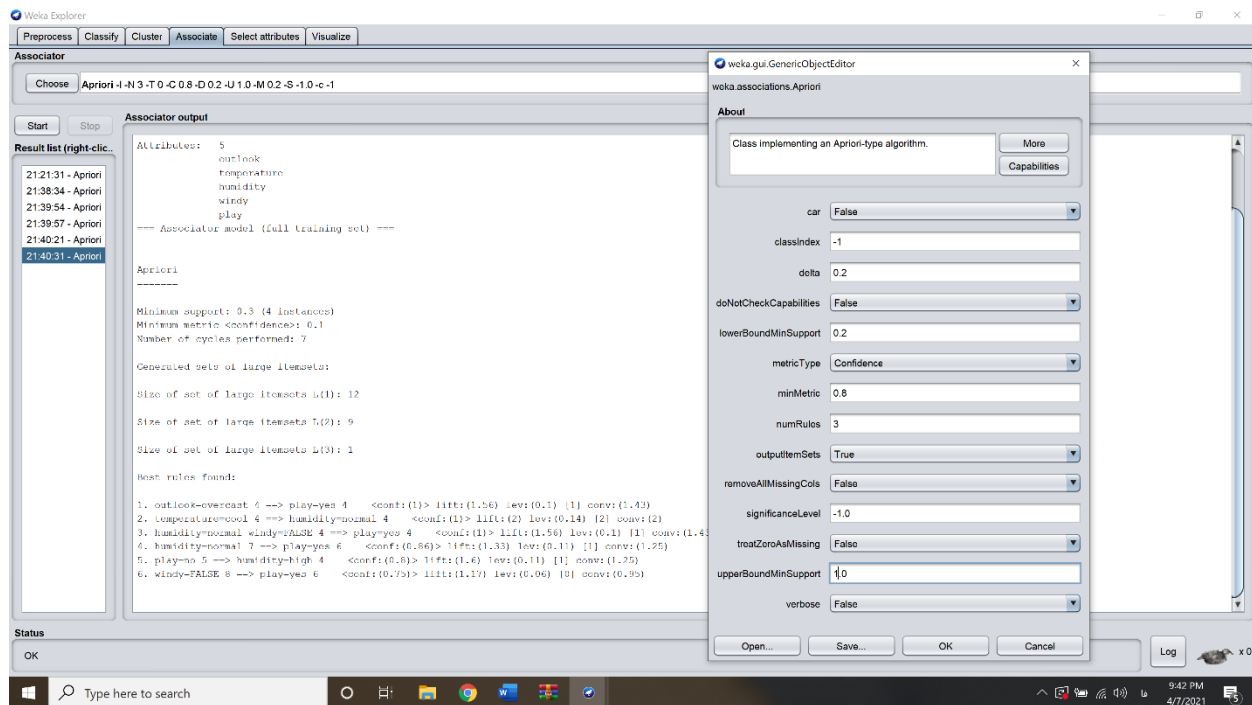
numRules: الگوریتم زمانی که به این تعداد بتواند قانون انجمنی تولید کند، متوقف میشود.

metricType: قوانین را اساس این معیار مرتب میشوند.

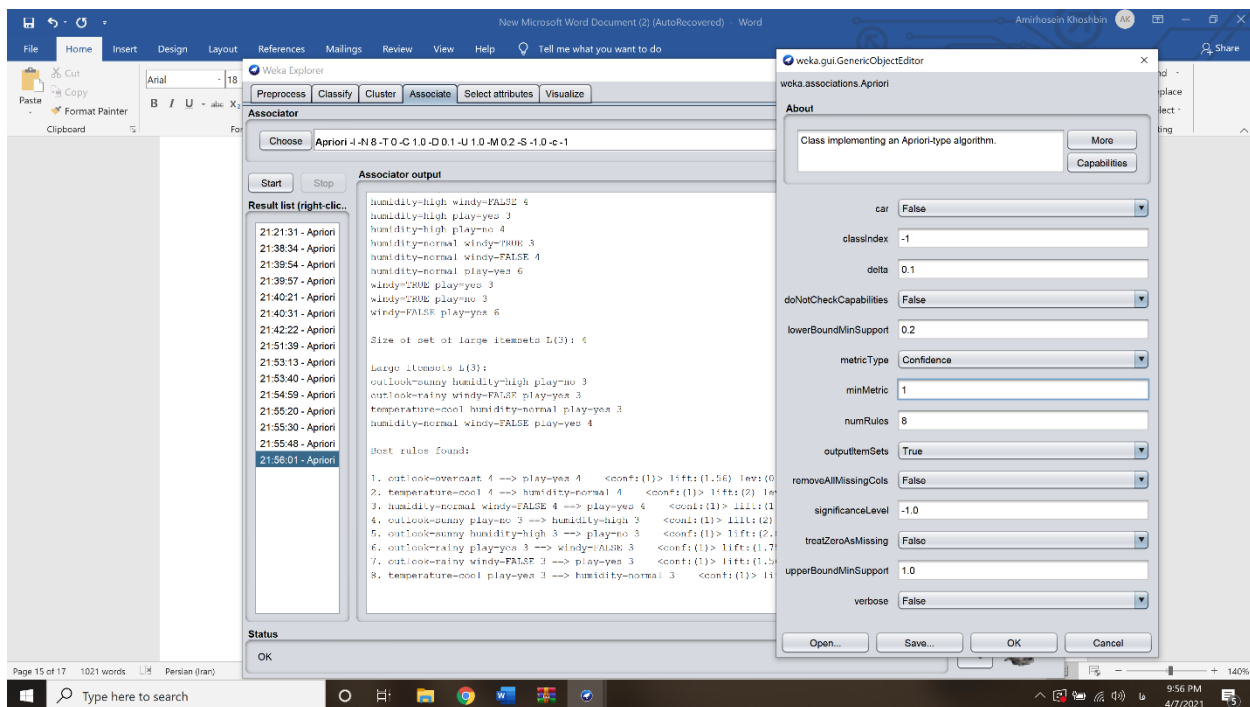
minMetric: قوانینی که از این عدد با معیار metricType بیشتر باشند به عنوان خروجی نمایش داده میشود.

outputItemSets: اگر میخواهیم تمام itemset های پرتکرار نمایش داده شوند، این باید true باشد.





الان بر اساس confidence فیلتر میکنیم و metric هم 1 میذاریم و میبینیم که confidence های بزرگتر از 1 و برای confidence فقط 1 را نشان میدهد چون بزرگ تر از 1 معنی نمیدهد.



وقتی numrules را به 3 تغییر دادیم، فقط سه قانون تولید شد، و بر اساس confidence مرتب شده و outputItemSet true شده و تمام set ها تولید شده نشان داده شده.

