

AI VIET NAM – COURSE 2024

Probability Exercise

(Naive Bayes Classifier)

Ngày 20 tháng 7 năm 2024

Giả sử X có các đặc trưng thuộc tính độc lập với nhau x_1, x_2, \dots, x_n , để phân loại X vào lớp một trong các lớp $C = c_1, c_2, \dots, c_m$, dựa vào công thức Bayes ta có:

$$P(c|X) = \frac{P(X|c).P(c)}{P(X)}$$

Dựa vào ước lượng tối đa xác suất hậu nghiệm (MAP - Maximum A Posterior) ta được:

$$P(c|X) \propto P(X|c).P(c)$$

$$P(c|X) \propto P(x_1|c).P(x_2|c)...P(x_n|c).P(c)$$

1. BINARY CLASSIFICATION - PLAY TENNIS

Cho tập dữ liệu huấn luyện mô hình phân loại nhị phân Naive Bayes gồm các thuộc tính "Outlook", "Temperature", "Humidity", "Wind":

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Overcast	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes

Bảng 1: Play Tennis - Tập dữ liệu huấn luyện

Cho sự kiện thử nghiệm:

$X = (\text{Outlook}=\text{Sunny}, \text{Temperature}=\text{Cool}, \text{Humidity}=\text{High}, \text{Wind}=\text{Strong})$

Câu hỏi 1: Xác suất xảy ra sự kiện "Play Tennis"="Yes" và sự kiện "Play Tennis"="No" lần lượt là:

- a) $P(\text{"Play Tennis"} = \text{"Yes"}) = 6/10, P(\text{"Play Tennis"} = \text{"No"}) = 4/10$
- b) $P(\text{"Play Tennis"} = \text{"Yes"}) = 4/10, P(\text{"Play Tennis"} = \text{"No"}) = 6/10$
- c) $P(\text{"Play Tennis"} = \text{"Yes"}) = 6/10, P(\text{"Play Tennis"} = \text{"No"}) = 6/10$
- d) $P(\text{"Play Tennis"} = \text{"Yes"}) = 4/10, P(\text{"Play Tennis"} = \text{"No"}) = 4/10$

Câu hỏi 2: Xác suất xảy ra sự kiện "Play Tennis"="Yes" khi sự kiện X xảy ra là:

- a) $P(\text{"Play Tennis"} = \text{"Yes"}|X) \propto 0.0014$
- b) $P(\text{"Play Tennis"} = \text{"Yes"}|X) \propto 0.0028$
- c) $P(\text{"Play Tennis"} = \text{"Yes"}|X) \propto 0.0188$
- d) $P(\text{"Play Tennis"} = \text{"Yes"}|X) \propto 0.0098$

Câu hỏi 3: Xác suất xảy ra sự kiện "Play Tennis"="No" khi sự kiện X xảy ra là:

- a) $P(\text{"Play Tennis"} = \text{"No"} | X) \propto 0.0014$
- b) $P(\text{"Play Tennis"} = \text{"No"} | X) \propto 0.0028$
- c) $P(\text{"Play Tennis"} = \text{"No"} | X) \propto 0.0188$
- d) $P(\text{"Play Tennis"} = \text{"No"} | X) \propto 0.0098$

Câu hỏi 4: Khi xảy ra sự kiện X, nhãn của "Play Tennis" sẽ là:

- a) "Play Tennis" = "Yes"
- b) "Play Tennis" = "No"

2. MULTI-LABEL CLASSIFICATION - TRAFFIC DATA Cho tập dữ liệu huấn luyện mô hình phân loại Naive Bayes gồm các thuộc tính "Day", "Season", "Fog", "Rain".

Day	Season	Fog	Rain	Class
Weekday	Spring	None	None	On Time
Weekday	Winter	None	Slight	On Time
Weekday	Winter	None	None	On Time
Holiday	Winter	High	Slight	Late
Saturday	Summer	Normal	None	On Time
Weekday	Autumn	Normal	None	Very Late
Holiday	Summer	High	Slight	On Time
Sunday	Summer	Normal	None	On Time
Weekday	Winter	High	Heavy	Very Late
Weekday	Summer	None	Slight	On Time
Saturday	Spring	High	Heavy	Cancelled
Weekday	Summer	High	Slight	On Time
Weekday	Winter	Normal	None	Late
Weekday	Summer	High	None	On Time
Weekday	Winter	Normal	Heavy	Vary Late
Saturday	Autumn	High	Slight	On Time
Weekday	Autumn	None	Heavy	On Time
Holiday	Spring	Normal	Slight	On Time
Weekday	Spring	Normal	None	On Time
Weekday	Spring	Normal	Heavy	On Time

Bảng 2: Traffic Data - Tập dữ liệu huấn luyện

Cho sự kiện thử nghiệm:

$X = (\text{Day}=\text{Weekday}, \text{Season}=\text{Winter}, \text{Fog}=\text{High}, \text{Rain}=\text{Heavy})$

Câu hỏi 5: Xác suất xảy ra sự kiện "Class"="On Time", sự kiện "Class"="Late", sự kiện "Class"="Very Late" và sự kiện "Class"="Cancelled" lần lượt là:

- (A) $P(\text{"Class"} = \text{"On Time"}) = 14/20$, $P(\text{"Class"} = \text{"Late"}) = 2/20$,
 $P(\text{"Class"} = \text{"Very Late"}) = 3/20$, $P(\text{"Class"} = \text{"Cancelled"}) = 1/20$
- (B) $P(\text{"Class"} = \text{"On Time"}) = 2/20$, $P(\text{"Class"} = \text{"Late"}) = 3/20$,
 $P(\text{"Class"} = \text{"Very Late"}) = 1/20$, $P(\text{"Class"} = \text{"Cancelled"}) = 14/20$
- (C) $P(\text{"Class"} = \text{"On Time"}) = 3/20$, $P(\text{"Class"} = \text{"Late"}) = 1/20$,
 $P(\text{"Class"} = \text{"Very Late"}) = 2/20$, $P(\text{"Class"} = \text{"Cancelled"}) = 14/20$
- (D) $P(\text{"Class"} = \text{"On Time"}) = 1/20$, $P(\text{"Class"} = \text{"Late"}) = 1/20$,
 $P(\text{"Class"} = \text{"Very Late"}) = 14/20$, $P(\text{"Class"} = \text{"Cancelled"}) = 3/20$

Câu hỏi 6: Xác suất xảy ra sự kiện "Class"="On Time" khi sự kiện X xảy ra là:

- (A) $P(\text{"Class"} = \text{"On Time"} \mid X) \propto 0.0222$
- (B) $P(\text{"Class"} = \text{"On Time"} \mid X) \propto 0.0013$

(C) $P(\text{"Class"} = \text{"On Time"} \mid X) \propto 0.0026$

(D) $P(\text{"Class"} = \text{"On Time"} \mid X) \propto 0.0000$

Câu hỏi 7: Xác suất xảy ra sự kiện "Class"="Late" khi sự kiện X xảy ra là:

(A) $P(\text{"Class"} = \text{"Late"} \mid X) \propto 0.0222$

(B) $P(\text{"Class"} = \text{"Late"} \mid X) \propto 0.0013$

(C) $P(\text{"Class"} = \text{"Late"} \mid X) \propto 0.0026$

(D) $P(\text{"Class"} = \text{"Late"} \mid X) \propto 0.0000$

Câu hỏi 8: Xác suất xảy ra sự kiện "Class"= "Very Late" khi sự kiện X xảy ra là:

(A) $P(\text{"Class"} = \text{"Very Late"} \mid X) \propto 0.0222$

(B) $P(\text{"Class"} = \text{"Very Late"} \mid X) \propto 0.0013$

(C) $P(\text{"Class"} = \text{"Very Late"} \mid X) \propto 0.0026$

(D) $P(\text{"Class"} = \text{"Very Late"} \mid X) \propto 0.0000$

Câu hỏi 9: Xác suất xảy ra sự kiện "Class"= Cancelled" khi sự kiện X xảy ra là:

(A) $P(\text{"Class"} = \text{"Cancelled"} \mid X) \propto 0.0222$

(B) $P(\text{"Class"} = \text{"Cancelled"} \mid X) \propto 0.0013$

(C) $P(\text{"Class"} = \text{"Cancelled"} \mid X) \propto 0.0026$

(D) $P(\text{"Class"} = \text{"Cancelled"} \mid X) \propto 0.0000$

Câu hỏi 10: Dự đoán "Class" của sự kiện X là:

(A) "On Time"

(B) "Late"

(C) "Very Late"

(D) "Cancelled"

3. IRIS CLASSIFICATION

Cho một tập dữ liệu huấn luyện phân loại hoa Iris dựa vào chiều dài cánh hoa như bảng dữ liệu bên dưới. Các bạn hãy trả lời các câu hỏi sau khi dùng Gaussian Naive Bayes cho data Iris này.

Length	1.4	1.0	1.3	1.9	2.0	1.8	3.0	3.8	4.1	3.9	4.2	3.4
Class	0	0	0	0	0	0	1	1	1	1	1	1

Bảng 3: Phân loại cánh hoa Iris dựa vào chiều dài cánh hoa - Tập dữ liệu huấn luyện

Câu hỏi 11: Giá trị mean và variance của biến đầu vào (Length) cho "Class"="0" lần lượt là:

- a) mean = 1.566 và variance = 0.128
- b) mean = 3.733 và variance = 0.172
- c) mean = 1.566 và variance = 0.172

Câu hỏi 12: Giá trị mean và variance của biến đầu vào (Length) cho "Class"="1" lần lượt là:

- a) mean = 1.566 và variance = 0.128
- b) mean = 3.733 và variance = 0.172
- c) mean = 1.566 và variance = 0.172

Câu hỏi 13: Cho dữ liệu kiểm thử $X = (\text{Length}=3.4)$. Xác suất dữ liệu kiểm thử X thuộc vào "Class"="0" và "Class"="1" lần lượt là:

- a) $P(\text{"Class"} = "0" | X) = 1.09 * 10^{-6}$ và $P(\text{"Class"} = "1" | X) = 0.3486$
- b) $P(\text{"Class"} = "0" | X) = 1.09 * 10^{-4}$ và $P(\text{"Class"} = "1" | X) = 0.3486$
- c) $P(\text{"Class"} = "0" | X) = 1.09 * 10^{-2}$ và $P(\text{"Class"} = "1" | X) = 0.3486$

4. PLAY TENNIS CLASSIFIER IMPLEMENTATION

Cho trước dữ liệu thời tiết của 10 ngày (D1-D10, như bảng 1). Hãy phát triển chương trình sử dụng mô hình phân loại Naive Bayes để dự đoán xem ngày thứ 11 (D11), AD có thể chơi tennis hay không?

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D11	Sunny	Cool	High	Strong	???

Bảng 4: Play Tennis - Dữ liệu testing

(a) "Play Tennis" = "Yes"

(b) "Play Tennis" = "No"

Để hoàn thành bài tập này bạn cần hoàn thành các function sau đây bằng cách sử dụng thư viện numpy:

4.1 Hoàn thiện function `create_train_dataset()` để tổ chức dữ liệu bảng 1 vào array 2 chiều như bên dưới.

```

1 #####
2 # Create data
3 #####
4 import numpy as np
5
6 def create_train_data():
7
8     #your code here *****
9
10    return np.array(data)
11
12 train_data = create_train_data()
13 print(train_data)
14
15 *****Sample Result when we print out train_data *****
16 [['Sunny' 'Hot' 'High' 'Weak' 'no']
17  ['Sunny' 'Hot' 'High' 'Strong' 'no']
18  ['Overcast' 'Hot' 'High' 'Weak' 'yes']
19  ['Rain' 'Mild' 'High' 'Weak' 'yes']
20  ['Rain' 'Cool' 'Normal' 'Weak' 'yes']
21  ['Rain' 'Cool' 'Normal' 'Strong' 'no']
22  ['Overcast' 'Cool' 'Normal' 'Strong' 'yes']
23  ['Overcast' 'Mild' 'High' 'Weak' 'no']
24  ['Sunny' 'Cool' 'Normal' 'Weak' 'yes']
25  ['Rain' 'Mild' 'Normal' 'Weak' 'yes']]

```

4.2 Hoàn thiện function `compute_prior_probability` tính $P(\text{"Play Tennis"} = \text{"Yes"})$ and tính $P(\text{"Play Tennis"} = \text{"No"})$ như bên dưới:

```

1 def compute_prior_probability(train_data):
2     y_unique = ['no', 'yes']
3     prior_probability = np.zeros(len(y_unique))
4     # your code here *****
5     return prior_probability
6
7 prior_probability = compute_prior_probability(train_data)
8 print("P(play tennis = No)", prior_probability[0])
9 print("P(play tennis = Yes)", prior_probability[1])

```

Câu hỏi 14: Kết quả nào sau đây là output từ chương trình trên:

- a) $P(\text{"Play Tennis"} = \text{"Yes"}) = 0.6$, $P(\text{"Play Tennis"} = \text{"No"}) = 0.4$
- b) $P(\text{"Play Tennis"} = \text{"Yes"}) = 0.3$, $P(\text{"Play Tennis"} = \text{"No"}) = 0.7$
- c) $P(\text{"Play Tennis"} = \text{"Yes"}) = 0.4$, $P(\text{"Play Tennis"} = \text{"No"}) = 0.8$
- d) $P(\text{"Play Tennis"} = \text{"Yes"}) = 0.4$, $P(\text{"Play Tennis"} = \text{"No"}) = 0.3$

4.3 Hoàn thiện function **compute_conditional_probability** để tính likelihood (The probability of "A" being True. Given "B" True, $P(A|B)$) như bên dưới:

```

1 def compute_conditional_probability(train_data):
2     y_unique = ['no', 'yes']
3     conditional_probability = []
4     list_x_name = []
5     for i in range(0, train_data.shape[1]-1):
6         x_unique = np.unique(data[:,i])
7         list_x_name.append(x_unique)
8
9         # your code here *****
10
11 conditional_probability.append(x_conditional_probability)
12 return conditional_probability, list_x_name
13

```

Câu hỏi 15: Hãy cho biết kết quả của đoạn chương trình sau đây:

```

1 train_data = create_train_data()
2 _, list_x_name = compute_conditional_probability(train_data)
3 print("x1 = ", list_x_name[0])
4 print("x2 = ", list_x_name[1])
5 print("x3 = ", list_x_name[2])
6 print("x4 = ", list_x_name[3])

```

- a) $x1 = [\text{'Cool'} \text{'Hot'} \text{'Mild'}]$
 $x2 = [\text{'Overcast'} \text{'Rain'} \text{'Sunny'}]$
 $x3 = [\text{'High'} \text{'Normal'}]$
 $x4 = [\text{'Strong'} \text{'Weak'}]$
- b) $x1 = [\text{'Overcast'} \text{'Rain'} \text{'Sunny'}]$
 $x2 = [\text{'Cool'} \text{'Hot'} \text{'Mild'}]$
 $x3 = [\text{'High'} \text{'Normal'}]$
 $x4 = [\text{'Strong'} \text{'Weak'}]$
- c) $x1 = [\text{'Strong'} \text{'Weak'}]$
 $x2 = [\text{'Cool'} \text{'Hot'} \text{'Mild'}]$
 $x3 = [\text{'High'} \text{'Normal'}]$
 $x4 = [\text{'Overcast'} \text{'Rain'} \text{'Sunny'}]$

```
d) x1 = ['Overcast' 'Rain' 'Sunny']
    x2 = ['Cool' 'Hot' 'Mild']
    x3 = ['Strong' 'Weak']
    x4 = ['High' 'Normal']
```

4.4 Hoàn thiện function `get_index_from_value` để tính trả về index tương ứng với feature name:

```
1 #This function is used to return the index of the feature name
2 def get_index_from_value(feature_name, list_features):
3     return np.where(list_eatures == feature_ame)[0][0]
```

Câu hỏi 16: Hãy cho biết kết quả của đoạn chương trình sau đây:

```
1 train_data = create_train_data()
2 _, list_x_name = compute_conditional_probability(train_data)
3 outlook = list_x_name[0]
4
5 i1 = get_index_from_value("Overcast", outlook)
6 i2 = get_index_from_value("Rain", outlook)
7 i3 = get_index_from_value("Sunny", outlook)
8
9 print(i1, i2, i3)
```

- a) 1 2 0
- b) 0 1 1
- c) 0 1 2
- d) 0 2 3

Câu hỏi 17: Hãy cho biết kết quả của đoạn chương trình sau đây:

```
1 train_data = create_train_data()
2 conditional_probability, list_x_name = compute_conditional_probability(train_data)
3 # Compute P("Outlook"="Sunny"|Play Tennis="Yes")
4 x1=get_index_from_value("Sunny",list_x_name[0])
5 print("P('Outlook'='Sunny'|Play Tennis'='Yes') = ", np.round(conditional_probability
    [0][1, x1],2))
```

- a) $P(\text{'Outlook'='Sunny'}|\text{Play Tennis'='Yes'}) = 0.27$
- b) $P(\text{'Outlook'='Sunny'}|\text{Play Tennis'='Yes'}) = 0.47$
- c) $P(\text{'Outlook'='Sunny'}|\text{Play Tennis'='Yes'}) = 0.37$
- d) $P(\text{'Outlook'='Sunny'}|\text{Play Tennis'='Yes'}) = 0.17$

Câu hỏi 18: Hãy cho biết kết quả của đoạn chương trình sau đây:

```
1 train_data = create_train_data()
2 conditional_probability, list_x_name = compute_conditional_probability(train_data)
3 # Compute P("Outlook"="Sunny"|Play Tennis="No")
4 x1=get_index_from_value("Sunny",list_x_name[0])
5 print("P('Outlook'='Sunny'|Play Tennis'='No') = ", np.round(conditional_probability
    [0][1, x1],2))
```


- a) $P(\text{'Outlook'='Sunny'}|\text{Play Tennis'='Yes'}) = 0.5$
 b) $P(\text{'Outlook'='Sunny'}|\text{Play Tennis'='Yes'}) = 0.4$
 c) $P(\text{'Outlook'='Sunny'}|\text{Play Tennis'='Yes'}) = 0.3$
 d) $P(\text{'Outlook'='Sunny'}|\text{Play Tennis'='Yes'}) = 0.2$

4.5 Hoàn thiện function **train_naive_bayes** như bên dưới:

```

1 #####
2 # Train Naive Bayes Model
3 #####
4 def train_naive_bayes(train_data):
5     # Step 1: Calculate Prior Probability
6     y_unique = ['no', 'yes']
7     prior_probability = compute_prior_probablity(train_data)
8
9     # Step 2: Calculate Conditional Probability
10    conditional_probability, list_x_name = compute_conditional_probability(train_data
11    )
12
13    return prior_probability, conditional_probability, list_x_name

```

4.6 Hoàn thiện function **prediction_play_tennis** để hỗ trợ AD có nên đi chơi tennis vào ngày D11 không:

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D11	Sunny	Cool	High	Strong	???

Bảng 5: Play Tennis - Dữ liệu testing

```

1 #####
2 # Prediction
3 #####
4 def prediction_play_tennis(X, list_x_name, prior_probability, conditional_probability)
5     :
6     x1=get_index_from_value(X[0],list_x_name[0])
7     x2=get_index_from_value(X[1],list_x_name[1])
8     x3=get_index_from_value(X[2],list_x_name[2])
9     x4=get_index_from_value(X[3],list_x_name[3])
10
11    p0 = 0
12    p1 = 0
13
14    # your code here *****
15
16    if p0>p1:
17        y_pred=0
18    else:
19        y_pred=1
20
21    return y_pred

```

Câu hỏi 19: Hãy cho biết kết quả của đoạn chương trình sau đây:

```
1 X = ['Sunny', 'Cool', 'High', 'Strong']
2 data = create_train_data()
3 prior_probability, conditional_probability, list_x_name = train_naive_bayes(data)
4 pred = prediction_play_tennis(X, list_x_name, prior_probability,
    conditional_probability)
5
6 if(pred):
7     print("Ad should go!")
8 else:
9     print("Ad should not go!")
```

a) Ad should not go!

b) Ad should go!

5. (OPTIONAL) IRIS CLASSIFIER IMPLEMENTATION

Cho trước dữ liệu chứa thông tin về hoa Iris gồm có sepal length, sepal width và petal length, và Species (bảng 6). Hãy phát triển chương trình sử dụng mô hình phân loại Gaussian Naive Bayes để dự đoán chủng loại của hoa Iris. Dữ liệu hoa iris được lưu trữ trong file iris_data.txt có thể được tải về [tại đây](#).

No.	Sepal length	Sepal width	Petal length	Petal width	Species
1	5.1	3.5	1.4	0.2	Iris-setosa
2	4.9	3.0	1.4	0.2	Iris-setosa
3	6.4	3.1	5.5	1.8	Iris-virginica
4	6.0	3.0	4.8	1.8	Iris-virginica
5	6.0	2.2	4.0	1.0	Iris-versicolora
...

Bảng 6: Iris flower - Tập dữ liệu huấn luyện

Dựa vào hướng dẫn dưới đây để thực thi mã nguồn cho bài toán phân loại.

```

1 # Example 1
2 # X=[sepal length, sepal width, petal length, petal width]
3 X = [6.3 , 3.3, 6.0, 2.5]
4 train_data = create_train_data_iris()
5 y_unique = np.unique(train_data[:,4])
6 prior_probability, conditional_probability = train_gaussian_naive_bayes(train_data)
7 pred = y_unique[prediction_iris(X, prior_probability, conditional_probability)]
8 assert pred == "Iris-virginica"
9
10 #Example 2 #####
11 # X=[sepal length, sepal width, petal length, petal width]
12 X = [5.0,2.0,3.5,1.0]
13 train_data = create_train_data_iris()
14 y_unique = np.unique(train_data[:,4])
15 prior_probability, conditional_probability = train_gaussian_naive_bayes(train_data)
16 pred = y_unique[prediction_iris(X, prior_probability, conditional_probability)]
17 assert pred == "Iris-versicolor"
18
19 #Example 3 #####
20 X = [4.9,3.1,1.5,0.1]
21 # X=[sepal length, sepal width, petal length, petal width]
22 train_data = create_train_data_iris()
23 y_unique = np.unique(train_data[:,4])
24 prior_probability, conditional_probability = train_gaussian_naive_bayes(train_data)
25 pred = y_unique[prediction_iris(X, prior_probability, conditional_probability)]
26 assert pred == "Iris-setosa"

```