

KHAI THÁC DỮ LIỆU & ỨNG DỤNG (*DATA MINING*)

GV: NGUYỄN HOÀNG TÚ ANH

1

BÀI 4 – PHẦN 2 PHÂN LỚP DỮ LIỆU

2

NỘI DUNG



1. ***Giới thiệu***
2. Phương pháp Naïve Bayes
3. Phương pháp dựa trên thể hiện
4. Đánh giá mô hình

3

GIỚI THIỆU



1. Phân lớp:

- ✚ Cho tập các mẫu đã phân lớp trước, xây dựng mô hình cho từng lớp
- ✚ ***Mục đích: Gán các mẫu mới vào các lớp với độ chính xác cao nhất có thể.***
- ✚ Cho CSDL $D=\{t_1, t_2, \dots, t_n\}$ và tập các lớp $C=\{C_1, \dots, C_m\}$, ***phân lớp*** là bài toán xác định ánh xạ $f: D \rightarrow C$ sao cho mỗi t_i được gán vào một lớp.

4

GIỚI THIỆU



Dữ liệu

Lượng giá, hồi qui, học, huấn luyện

Mô hình

Phân loại, ra quyết định

Hành động

5

Bài tập cá nhân



Customer	Age	Income (K)	No. cards	Response
Lâm	35	35	3	Yes
Hưng	22	50	2	No
Mai	28	60	1	Yes
Lan	45	100	2	No
Thủy	20	30	3	Yes
Tuấn	34	55	2	No
Minh	63	200	1	No
Vân	55	140	2	No
Thiện	59	170	1	No
Ngọc	25	40	4	Yes
<u>Châu</u>	30	45	1	???

Thời gian: 5'

Yêu cầu:

Trình bày ý tưởng xác định lớp cho mẫu cuối cùng (Châu) khi cho biết các mẫu còn lại.

6

NỘI DUNG



1. Giới thiệu
2. Phương pháp Naïve Bayes
3. Phương pháp dựa trên thể hiện
4. Đánh giá mô hình

7

GIỚI THIỆU



1. Phân lớp theo mô hình xác suất:

- ✦ Dự đoán xác suất hay dự đoán xác suất là thành viên của lớp
- ✦ *Nền tảng: dựa trên định lý Bayes*
 - ✦ Cho X, Y là các biến bất kỳ (rời rạc, số, cấu trúc, ...)
 - ✦ Dự đoán Y từ X
- ✦ Lượng giá các tham số của $P(X | Y)$, $P(Y)$ trực tiếp từ tập DL huấn luyện
- ✦ Sử dụng định lý Bayes để tính $P(Y | X=x)$ ⁸

GIỚI THIỆU



2. Định lý Bayes

$$P(y | x) = \frac{P(x | y) \cdot P(y)}{P(x)}$$

Cụ thể:

$$P(Y = y_i | X = x_j) = \frac{P(X = x_j | Y = y_i) P(Y = y_i)}{P(X = x_j)}$$

Biến bất kỳ

Giá trị thứ i

9

GIỚI THIỆU



2. Định lý Bayes

$$P(Y = y_i | X = x_j) = \frac{P(X = x_j | Y = y_i) P(Y = y_i)}{P(X = x_j)}$$

Tương đương:

$$(\forall i, j) P(Y = y_i | X = x_j) = \frac{P(X = x_j | Y = y_i) P(Y = y_i)}{\sum_k P(X = x_j | Y = y_k) P(Y = y_k)}$$

10

GIỚI THIỆU



3. Phân loại Bayes

Tập DL huấn luyện

X						Y
Sky	Temp	Humid	Wind	Water	Forecst	EnjoySpt
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes

XD mô hình: Lượng giá $P(X | Y)$, $P(Y)$

Phân lớp: Dùng định lý Bayes để tính

$$P(Y | X^{\text{new}})$$

11

GIỚI THIỆU



4. Độc lập điều kiện (Conditional independence)

Định nghĩa: X độc lập điều kiện với Y khi cho Z nếu phân bố xác suất trên X độc lập với các giá trị của Y khi cho các giá trị của Z.

$$(\forall i, j, k) P(X = x_i | Y = y_j, Z = z_k) = P(X = x_i | Z = z_k)$$

Ta thường viết:

$$P(X | Y, Z) = P(X | Z)$$

Ví dụ:

$$P(\text{Sấm sét} | \text{Mưa}, \text{Chớp}) = P(\text{Sấm sét} | \text{Chớp})$$

12

Thuật toán Naïve Bayes



Giả sử:

- D: tập huấn luyện gồm các mẫu biểu diễn dưới dạng $X = \langle x_1, \dots, x_n \rangle$
- $C_{i,D}$: tập các mẫu của D thuộc lớp C_i với $i = \{1, \dots, m\}$
- Các thuộc tính x_1, \dots, x_n độc lập điều kiện đôi một với nhau khi cho lớp C

Khi đó: ta cần xác định xác suất $P(C_i|X)$ lớn nhất

13

Thuật toán Naïve Bayes



Theo định lý Bayes:

$$P(C_i|X) = \frac{P(X|C_i)P(C_i)}{P(X)}$$

Theo tính chất độc lập điều kiện:

$$P(X|C_i) = \prod_{k=1}^n P(x_k | C_i) = P(x_1 | C_i) \times P(x_2 | C_i) \times \dots \times P(x_n | C_i)$$

Luật phân lớp cho $X^{\text{new}} = \{x_1, \dots, x_n\}$ là:

$$\arg \max_{C_k} P(C_i) \prod_{k=1}^n P(x_k | C_i)$$

14

Thuật toán Naïve Bayes



B1: Huấn luyện Naïve Bayes (trên tập DL huấn luyện)

Lượng giá $P(C_i)$

Lượng giá $P(X_k|C_i)$

B2: X^{new} được gán vào lớp cho giá trị công thức lớn nhất:

$$\arg \max_{C_k} P(C_i) \prod_{k=1}^n P(x_k | C_i)$$

15

Trường hợp X – giá trị rời rạc



Giả sử:

• $X = \langle X_1, \dots, X_n \rangle$

• x_i nhận các giá trị rời rạc

Khi đó: Lượng giá $P(C_i)$ và lượng giá $P(X_k|C_i)$ theo công thức

$$P(C_i) \approx \frac{|C_{i,D}|}{|D|} \quad P(x_k | C_i) \approx \frac{\#C_{i,D}\{x_k\}}{|C_{i,D}|}$$

16

Trường hợp X – giá trị rời rạc



- Để tránh trường hợp giá trị $P(X_k|C_i) = 0$ do không có mẫu nào trong DL huấn luyện thỏa mãn tử số, ta làm trơn bằng cách thêm một số mẫu ảo.

Khi đó:

- Làm trơn theo Laplace:

$$P(C_i) \approx \frac{|C_{i,D}|+1}{|D|+m} \quad P(x_k | C_i) \approx \frac{\#C_{i,D}\{x_k\}+1}{|C_{i,D}|+r}$$

với m – số lớp và r là số giá trị rời rạc của thuộc tính.

17

VÍ DỤ 1:



Cho tập dữ liệu huấn luyện:

Outlook	Temperature	Humidity	Windy	Play?
sunny	hot	high	weak	No
sunny	hot	high	strong	No
overcast	hot	high	weak	Yes
rain	mild	High	weak	Yes
rain	cool	Normal	weak	Yes
rain	cool	normal	strong	No
overcast	cool	normal	strong	Yes
sunny	mild	high	weak	No
sunny	cool	normal	weak	Yes
rain	mild	normal	weak	Yes
sunny	mild	normal	strong	Yes
overcast	mild	high	strong	Yes
overcast	hot	normal	weak	Yes
rain	mild	high	strong	No

18

VÍ DỤ 1:



B1: Ước lượng $P(C_i)$ với $C_1 = \text{"yes"}$, $C_2 = \text{"no"}$ và $P(x_k|C_i)$

Ta thu được $P(C_i)$:

$$P(C_1) = 9/14 = 0.643$$

$$P(C_2) = 5/14 = 0.357$$

Với thuộc tính Outlook, ta có các giá trị: sunny, overcast, rain. Trong đó $P(\text{sunny}|C_i)$ là:

Outlook	
$P(\text{sunny} \text{yes}) = 2/9$	$P(\text{sunny} \text{no}) = 3/5$

19

Bài tập theo nhóm



• Thời gian: 5'

Ước lượng $P(x_k|C_i)$ với $C_1 = \text{"yes"}$, $C_2 = \text{"no"}$

• $P(\text{Outlook}|C_i)$

• Nhóm

• $P(\text{Temperature}|C_i)$

• Nhóm

• $P(\text{Humidity}|C_i)$

• Nhóm

• $P(\text{windy}|C_i)$

• Nhóm

20

VÍ DỤ 1:



B2 : Phân lớp

$X^{new} = \langle \text{Outlook}=\text{sunny}, \text{Temp} = \text{cool}, \text{Humidity} = \text{high}, \text{Windy} = \text{strong} \rangle$

Ta cần tính:

$$P(C_1) * P(X|C_1) = P(C_1) * P(\text{sunny}|y) * P(\text{cool}|y) * P(\text{high}|y) * P(\text{strong}|y) = 0.005$$

$$P(C_2) * P(X|C_2) = P(C_2) * P(\text{sunny}|n) * P(\text{cool}|n) * P(\text{high}|n) * P(\text{strong}|n) = 0.021$$

→ X^{new} thuộc lớp C_2 ("no")

21

Bài tập cá nhân



Thời gian: 5'

Hãy xác định lớp cho mẫu mới sau:

$X^{new} = \langle \text{Outlook} = \text{overcast}, \text{Temp} = \text{cool}, \text{Humidity} = \text{high}, \text{Windy} = \text{strong} \rangle$

22

VÍ DỤ 1: Làm tròn Laplace



B1: Ước lượng $P(C_i)$ với $C_1 = \text{"yes"}$, $C_2 = \text{"no"}$ và $P(x_k|C_i)$ theo công thức làm tròn Laplace

$$P(C_1) = (9+1)/(14+2) = 10/16$$

$$P(C_2) = (5+1)/(14+2) = 6/16$$

Outlook	
$P(\text{sunny} y) = 3/12$	$P(\text{sunny} n) = 4/8$
$P(\text{overcast} y) = 5/12$	$P(\text{overcast} n) = 1/8$
$P(\text{rain} y) = 4/12$	$P(\text{rain} n) = 3/8$
Temperature	
$P(\text{hot} y) = 3/12$	$P(\text{hot} n) = 3/8$
$P(\text{mild} y) = 5/12$	$P(\text{mild} n) = 3/8$
$P(\text{cool} y) = 4/12$	$P(\text{cool} n) = 2/8$
Humidity	
$P(\text{high} y) = 4/11$	$P(\text{high} n) = 5/7$
$P(\text{normal} y) = 7/11$	$P(\text{normal} n) = 2/7$
Windy	
$P(\text{strong} y) = 4/11$	$P(\text{strong} n) = 4/7$
$P(\text{weak} y) = 7/11$	$P(\text{weak} n) = 3/7$

VÍ DỤ 1:



B2: Phân loại

$X^{new} = \langle \text{Outlook} = \text{overcast}, \text{Temp} = \text{cool}, \text{Humidity} = \text{high}, \text{Windy} = \text{strong} \rangle$

Ta tính theo công thức làm tròn Laplace :

$$P(C_1) * P(X|C_1) = P(C_1) * P(\text{overcast}|y) * P(\text{cool}|y) * P(\text{high}|y) * P(\text{strong}|y) = .011$$

$$P(C_2) * P(X|C_2) = P(C_2) * P(\text{overcast}|n) * P(\text{cool}|n) * P(\text{high}|n) * P(\text{strong}|n) = .005$$

$$P(C_2) * P(X|C_2) = P(C_2) * P(\text{overcast}|n) * P(\text{cool}|n) * P(\text{high}|n) * P(\text{strong}|n) = .005$$

$$P(\text{strong}|n) = .005$$

$\Rightarrow X^{new}$ thuộc lớp C_1 ("yes")

Trường hợp X – giá trị liên tục



- Nếu thuộc tính nhận giá trị liên tục thì xác suất $P(X_k|C_i)$ thường được tính dựa theo phân bố Gauss với giá trị trung bình μ và độ lệch σ :

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Và $P(X_k|C_i)$ là:

$$P(\mathbf{X} | C_i) = g(x_k, \mu_{C_i}, \sigma_{C_i})$$

25

Thuật toán Naïve Bayes



- Ưu điểm:

- Dễ dàng cài đặt
- Thời gian thi hành tương tự như cây quyết định
- Đạt kết quả tốt trong phần lớn các trường hợp

- Nhược điểm:

- Giả thiết về tính độc lập điều kiện của các thuộc tính làm giảm độ chính xác

26

NỘI DUNG



1. Giới thiệu
2. Phương pháp Naïve Bayes
3. Phương pháp dựa trên thể hiện
4. Đánh giá mô hình

27

GIỚI THIỆU



- Phương pháp phân lớp dựa trên thể hiện (Instance-based):
 - Lưu trữ các mẫu/đối tượng huấn luyện và chỉ xử lý khi có yêu cầu phân lớp mẫu/đối tượng mới
 - Đưa mẫu/đối tượng vào lớp mà gần với chúng nhất
- Các phương pháp:
 - Thuật toán k- láng giềng gần nhất (k-NN)
 - Hồi qui với trọng số cục bộ (Locally weighted regression)
 - Suy luận dựa trên trường hợp (Case-based reasoning)

28

K- LẮNG GIỀNG GẦN NHẤT

Hãy cho tôi biết bạn của bạn là ai, tôi sẽ nói bạn là người như thế nào.

- Một mẫu mới được gán vào lớp có nhiều mẫu giống với nó nhất trong số k mẫu gần nhất

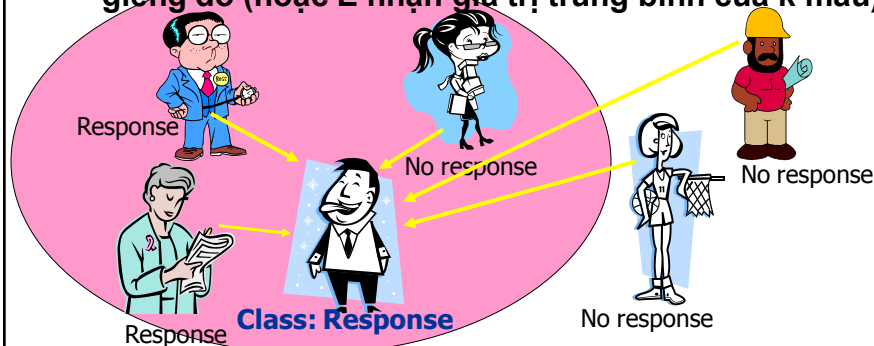


29

K- LẮNG GIỀNG GẦN NHẤT

- Thuật toán xác định lớp cho mẫu mới E:

- Tính khoảng cách giữa E và tất cả các mẫu trong tập huấn luyện
- Chọn k mẫu gần nhất với E trong tập huấn luyện
- Gán E vào lớp có nhiều mẫu nhất trong số k mẫu láng giềng đó (hoặc E nhận giá trị trung bình của k mẫu)



30

K- LĂNG GIỀNG GẦN NHẤT

- Tính khoảng cách giữa 2 mẫu/ đối tượng
 - Mỗi mẫu - tập thuộc tính số
 - Khoảng cách *Euclidean* giữa $X=(x_1, \dots, x_n)$ và $Y=(y_1, \dots, y_n)$ là:

$$D(X, Y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

- Khi thực hiện so sánh, có thể bỏ qua căn bậc 2

31

K- LĂNG GIỀNG GẦN NHẤT

- Ví dụ tính khoảng cách giữa John và Rachel



John:
Age=35
Income=95K
No. of credit cards=3



Rachel:
Age=41
Income=215K
No. of credit cards=2

$$D(\text{John}, \text{Rachel}) = \sqrt{(35-41)^2 + (95K-215K)^2 + (3-2)^2}$$

- Các thuộc tính **có giá trị lớn** sẽ ảnh hưởng nhiều đến khoảng cách giữa các đối tượng (VD: thuộc tính income)
 - Các thuộc tính có **miền giá trị khác nhau**
- > **Cần chuẩn hóa giá trị thuộc tính**

32

K- LĂNG GIỀNG GẦN NHẤT

- Cần phải chuẩn hoá dữ liệu: ánh xạ các giá trị vào đoạn $[0,1]$ theo công thức

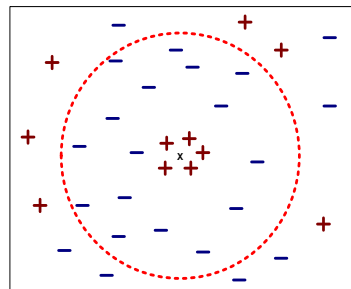
$$a_i = \frac{v_i - \min v_i}{\max v_i - \min v_i}$$

với: v_i là giá trị thực tế của thuộc tính i
 a_i là giá trị của thuộc tính đã chuẩn hóa

33

K- LĂNG GIỀNG GẦN NHẤT

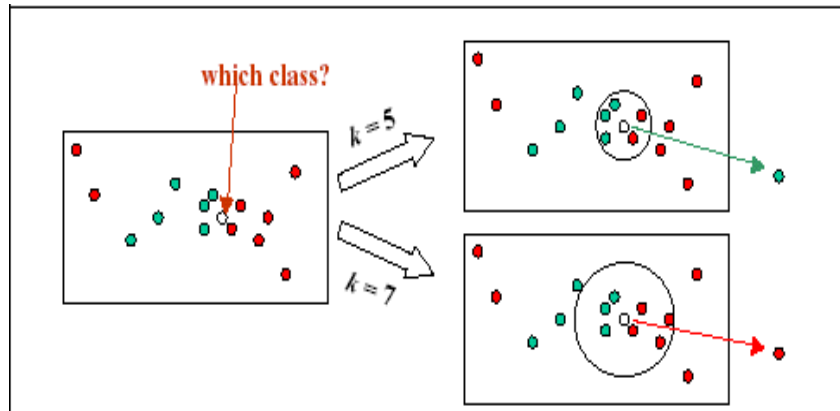
- Ưu điểm:
 - Dễ sử dụng và cài đặt
 - Xử lý tốt với dữ liệu nhiễu
- Khuyết điểm:
 - Cần lưu tất cả các mẫu
 - Cần nhiều thời gian để xác định lớp cho một mẫu mới (cần tính và so sánh khoảng cách đến tất cả các mẫu)
 - Phụ thuộc vào giá trị k do người dùng lựa chọn
 - Nếu k quá nhỏ, nhạy cảm với nhiễu
 - Nếu k quá lớn, vùng lân cận có thể chứa các điểm của lớp khác
 - Thuộc tính phi số ?



34

K- LĂNG GIỀNG GẦN NHẤT

- PHỤ THUỘC VÀO GIÁ TRỊ K DO NGƯỜI DÙNG LỰA CHỌN



NỘI DUNG

1. Giới thiệu
2. Phương pháp Naïve Bayes
3. Phương pháp dựa trên thể hiện
4. Đánh giá mô hình

Đánh giá mô hình



- Ngoài thuật toán học, sự thực thi của mô hình có thể phụ thuộc vào các yếu tố khác:
 - Sự phân bố của các lớp
 - Chi phí phân loại sai
 - Kích thước của tập huấn luyện và tập thử nghiệm
- Độ đo thực thi
- Phương pháp đánh giá

37

Đánh giá mô hình



- Đánh giá thực thi
 - Tập trung vào khả năng dự đoán của mô hình hơn là tốc độ phân loại hay xây dựng mô hình, khả năng co dẫn,...

		PREDICTED CLASS	
ACTUAL CLASS		Class=Yes	Class=No
	Class=Yes	a	b
	Class=No	c	d

a: TP (true positive) b: FN (false negative)
c: FP (false positive) d: TN (true negative)

38

Đánh giá thực thi

	PREDICTED CLASS	
	Class=Yes	Class=No
	Class=Yes	Class=No
ACTUAL CLASS	a (TP)	b (FN)
	c (FP)	d (TN)

- Độ chính xác của mô hình M, $\text{acc}(M)$

$$\text{Acc}(M) = \frac{a + d}{a + b + c + d} = \frac{TP + TN}{TP + TN + FP + FN}$$

39

Đánh giá thực thi

- Độ lỗi của mô hình M, $\text{error_rate}(M) = 1 - \text{acc}(M)$

- Một số độ đo khác:

$$\text{Precision (p)} = \frac{a}{a + c}$$

$$\text{Recall (r)} = \frac{a}{a + b}$$

$$\text{F-measure (F)} = \frac{2rp}{r + p} = \frac{2a}{2a + b + c}$$

40

Đánh giá thực thi



- Ví dụ

classes	buy_computer = yes	buy_computer = no	total	recognition(%)
buy_computer = yes	6954	46	7000	99.34
buy_computer = no	412	2588	3000	86.27
total	7366	2634	10000	95.42

- $\text{acc}(M) = (6954+2588)/10000=95.42\%$
- $\text{error_rate}(M) = 1-95.42\%=4.58\%$
- $\text{Precision}(M\text{-Yes}) = 6954/7366 = 94.41\%$
- $\text{Recall}(M\text{-Yes}) = 6954/7000 = 99.34\%$
- $\text{F-measure}(M\text{-Yes})= 96.81\%$

41

Đánh giá mô hình



- Phương pháp đánh giá

- Phương pháp Holdout:

- Phân chia ngẫu nhiên tập DL thành 2 tập độc lập:
 - Tập huấn luyện: 2/3 và tập thử nghiệm: 1/3
 - *Thích hợp cho tập DL nhỏ.*
 - Các mẫu có thể không đại diện cho toàn bộ DL: thiếu lớp trong tập thử nghiệm
 - Cải tiến:
 - Dùng phương pháp lấy mẫu sao cho mỗi lớp được phân bố đều trong cả 2 tập DL huấn luyện và thử nghiệm
 - Lấy mẫu ngẫu nhiên: thực hiện holdout k lần và độ chính xác $\text{acc}(M)$ = trung bình cộng k giá trị chính xác

Phương pháp đánh giá



• Phương pháp Cross-validation (k-fold)

- Phân chia DL thành k tập con có cùng kích thước
- *Tại mỗi vòng lặp sử dụng một tập con là tập thử nghiệm và các tập con còn lại là tập huấn luyện*
- Giá trị k thường là = 10
- *Leave-one-out: k=số mẫu trong DL (dành cho tập DL nhỏ)*
- Stratified cross-validation: dùng phương pháp lấy mẫu để phân bố các lớp trong từng tập con như trên toàn bộ DL.

43

TÓM TẮT



- Phân lớp là hình thức phân tích DL để rút ra các mô hình mô tả các lớp DL quan trọng
- Nhiều thuật toán hiệu quả được phát triển.
- Không thuật toán nào vượt trội nhất cho mọi tập DL
- Các vấn đề như độ chính xác, thời gian huấn luyện, tính linh hoạt, khả năng co giãn,... cần quan tâm và nghiên cứu sâu hơn .

44

Bài tập nhóm



Customer	Age	Income (K)	No. cards	Response
Lâm	35	35	3	Yes
Hưng	22	50	2	No
Mai	28	40	1	Yes
Lan	45	100	2	No
Thủy	20	30	3	Yes
Tuấn	34	55	2	No
Minh	63	200	1	No
Vân	55	140	2	No
Thiện	59	170	1	No
Ngọc	25	40	4	Yes
<u>Minh</u>	<u>39</u>	<u>41</u>	<u>2</u>	<u>???</u>

Thời gian: 15'
Sử dụng thuật toán k-NN với
1) $k = 3$ và
2) $k = 5$
Để xác định lớp cho “Minh”.

45

Qui định trình bày bài nộp



Bài tập nhóm

- Ngày nộp:
- Tên nhóm: liệt kê các thành viên tham gia
 - Họ và tên:
 - Mã số SV :
- Nội dung:

46

CÁC CÔNG VIỆC CẦN LÀM



1. Thảo luận và tự thực hiện các bài tập của chương 4 –Phần 1 và Phần 2 (không nộp)
2. Chuẩn bị bài Gồm nhóm dữ liệu
 - Xem nội dung các bài tập.
 - **Cách thực hiện :**
 - **Đọc slide, xem các ví dụ**
 - **Tham khảo trên Internet và tài liệu tham khảo**

47

BÀI TẬP PHẦN 2



1. Cho tập huấn luyện như trong ví dụ 1 của bài 5-P1 (“mua”, “không mua máy tính”). Áp dụng thuật toán Naïve Bayes cho ví dụ 1 và xác định lớp cho mẫu mới : $X = (<=30, \text{medium}, \text{yes}, \text{fair})$
So sánh với kết quả phân lớp sử dụng cây quyết định.
2. Cho tập huấn luyện như trong ví dụ 3 của bài 5-P1. Áp dụng phương pháp Naïve Bayes để tính các xác suất $P(C_i)$ và $P(x_k|C_i)$ với $C_1 = \text{“yes”}$, $C_2 = \text{“no”}$. Chuẩn hóa các xác suất bằng phương pháp làm trơn Laplace.

48

BÀI TẬP PHẦN 2



Tập DL huấn luyện ví dụ 1 – bài 5-P1

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

49

BÀI TẬP PHẦN 2



Tập DL huấn luyện ví dụ 3 – bài 5-P1

No	Size	Color	Shape	Decision
1	Vừa	Xanh dương	Hộp	Yes
2	Nhỏ	đỏ	Nón	No
3	Nhỏ	đỏ	Cầu	Yes
4	Lớn	đỏ	Nón	No
5	Lớn	Xanh lá cây	Trụ	Yes
6	Lớn	đỏ	Trụ	No
7	Lớn	Xanh lá cây	Cầu	Yes

50

BÀI TẬP PHẦN 2



3. Cho tập huấn luyện sau :

- Sử dụng thuật toán k-NN để xác định lớp cho “Tuyền” với $k = 3$, hoặc 5, hoặc 7. So sánh kết quả thu được.
 - Chuẩn hóa DL và xác định lớp cho “Dũng”. So sánh kết quả với câu a).
 - Tìm phương pháp biến đổi tập DL bên về dạng có thể áp dụng phương pháp cây quyết định, ILA, Naïve Bayes. Áp dụng một trong 3 phương pháp đó lên DL đã biến đổi để xác định lớp cho “Dũng”. So sánh kết quả với câu a).
4. So sánh ưu điểm, khuyết điểm của các phương pháp phân lớp dựa trên cây quyết định, dựa trên luật, xác suất và dựa trên thể hiện .

51

BÀI TẬP PHẦN 2



Customer	Age	Income (K)	No. cards	Response
Lâm	35	35	3	Yes
Hưng	22	50	2	No
Mai	28	40	1	Yes
Lan	45	100	2	No
Thủy	20	30	3	Yes
Tuấn	34	55	2	No
Minh	63	200	1	No
Vân	55	140	2	No
Thiện	59	170	1	No
Ngọc	25	40	4	Yes
<u>Tuyền</u>	<u>25</u>	<u>30</u>	<u>1</u>	<u>???</u>

52

TÀI LIỆU THAM KHẢO



1. T. M. Mitchell, **Machine Learning**. McGraw Hill, 1997
2. J.Han, M.Kamber, Chương 7 – Data mining : Concepts and Techniques
<http://www.cs.sfu.ca/~han/dmbook>
<http://www-faculty.cs.uiuc.edu/~hanj/bk2/slidesindex.html> : 2nd
3. P.-N. Tan, M. Steinbach, V. Kumar, Chương 4 - Introduction to Data Mining
<http://www-users.cs.umn.edu/~kumar/dmbook/ch4.pdf>

53

Q & A



54