

## Phương pháp học Bayes Bayesian classification

1

## Nội dung

- Giới thiệu về Bayesian classification
- Kiến thức về xác suất thống kê
- Giải thuật học của naive Bayes
- Kết luận và hướng phát triển

2

## Bayesian classification

Phương pháp học Bayes – bayesian classification

- Phân loại này được đặt theo tên của **Thomas Bayes** (1702-1761), người đề xuất các định lý Bayes
- Giải thuật học có giám sát (supervised learning) - xây dựng mô hình phân loại dựa trên dữ liệu tập học đã có nhãn (lớp)
- Mạng Bayes (Bayesian network), **Bayes ngây thơ (naive Bayes)**
- Giải quyết các vấn đề về phân loại, gom nhóm, etc.

3

## Bayesian classification

Phương pháp học Bayes ứng dụng thành công

- **Phân loại thư rác**  
Cho một email, dự đoán xem đó là thư rác hay không
- **Chẩn đoán y tế**  
Cho một danh sách các triệu chứng, dự đoán xem bệnh nhân có bệnh X hay không
- **Thời tiết**  
Dựa vào nhiệt độ, độ ẩm, vv ... dự đoán nếu nó sẽ mưa vào ngày mai

4

## Bayesian classification

- Phương pháp Bayesian là hệ thống **ham học**
- Dựa vào **các đặc trưng** đưa ra kết luận **nhãn** của đối tượng mới đến
- Khi đưa ra một tập huấn luyện, hệ thống **ngay lập tức** phân tích dữ liệu và **xây dựng một mô hình**. Khi cần phân loại một đối tượng mới đến, hệ thống sử dụng mô hình đã xây dựng để xác định đối tượng mới.
- Phương pháp Bayesian (ham học) có xu hướng phân loại các trường hợp nhanh hơn KNN (lười học)

## Kỹ thuật DM thành công (2011)

Which methods/algorithms did you use for data analysis in 2011? [311 voters]

Decision Trees/Rules (186)	59.8 %
Regression (180)	57.9 %
Clustering (163)	52.4 %
Statistics (descriptive) (149)	47.9 %
Visualization (119)	38.3 %
Time series/Sequence analysis (92)	29.6 %
Support Vector (SVM) (89)	28.6 %
Association rules (89)	28.6 %
Ensemble methods (88)	28.3 %
Text Mining (86)	27.7 %
Neural Nets (84)	27.0 %
Boosting (73)	23.5 %
Bayesian (68)	21.9 %
Bagging (63)	20.3 %
Factor Analysis (58)	18.7 %
Anomaly/Deviation detection (51)	16.4 %
Social Network Analysis (44)	14.2 %
Survival Analysis (29)	9.32 %
Genetic algorithms (29)	9.32 %
Uplift modeling (15)	4.82 %

### Top 10 DM algorithms (2015)



6

## Nội dung

- Giới thiệu về Bayesian classification
- **Kiến thức về xác suất thống kê**
- Giải thuật học của naive Bayes
- Kết luận và hướng phát triển

7

## Xác suất thống kê



Một vài ví dụ

- Khi tung 1 đồng xu, khả năng nhận mặt ngửa là bao nhiêu?
- Khi tung một hộp xúc xắc, khả năng xuất hiện mặt “6 nút” là bao nhiêu?

$P(h)$  : ký hiệu xác suất của giả thuyết  $h$

## Xác suất thống kê



Xác suất xuất hiện mặt ngửa:

$$P(\text{ngửa}) = 0.5$$

Xác suất xuất hiện mặt có 6 nút:

$$P(6) = 1/6$$

## Xác suất thống kê

name	laptop	phone
Kate	PC	Android
Tom	PC	Android
Harry	PC	Android
Annika	Mac	iPhone
Naomi	Mac	Android
Joe	Mac	iPhone
Chakotay	Mac	iPhone
Neelix	Mac	Android
Kes	PC	iPhone
B'Elanna	Mac	iPhone

➤ Xác suất mà một người được lựa chọn ngẫu nhiên sử dụng iPhone là bao nhiêu?

➤ Xác suất mà một người được lựa chọn ngẫu nhiên sử dụng iPhone khi người này có sử dụng một máy tính xách tay Mac là bao nhiêu?

## Xác suất thống kê

name	laptop	phone
Kate	PC	Android
Tom	PC	Android
Harry	PC	Android
Annika	Mac	iPhone
Naomi	Mac	Android
Joe	Mac	iPhone
Chakotay	Mac	iPhone
Neelix	Mac	Android
Kes	PC	iPhone
B'Elanna	Mac	iPhone

Xác suất mà một người được lựa chọn ngẫu nhiên sử dụng iPhone là bao nhiêu?

Xác suất mà một người được lựa chọn ngẫu nhiên sử dụng iPhone khi người này có sử dụng một máy tính xách tay Mac là bao nhiêu?

Xác suất của A với điều kiện B xảy ra được định nghĩa như sau :

$$P(A/B) = \frac{P(AB)}{P(B)}$$

## Xác suất thống kê

Xác suất của A với điều kiện B xảy ra được định nghĩa như sau :

$$P(A/B) = \frac{P(AB)}{P(B)}$$

name	laptop	phone
Kate	PC	Android
Tom	PC	Android
Harry	PC	Android
Annika	Mac	iPhone
Naomi	Mac	Android
Joe	Mac	iPhone
Chakotay	Mac	iPhone
Neelix	Mac	Android
Kes	PC	iPhone
B'Elanna	Mac	iPhone

Xác suất mà một người được lựa chọn ngẫu nhiên sử dụng iPhone?

$$P(\text{iPhone}) = 5/10 = 0.5$$

Xác suất mà một người được lựa chọn ngẫu nhiên sử dụng iPhone khi người này sử dụng một máy tính xách tay Mac?

$$P(\text{iPhone} | \text{mac}) = \frac{P(\text{mac} \cap \text{iPhone})}{P(\text{mac})}$$

$$P(\text{mac} \cap \text{iPhone}) = \frac{4}{10} = 0.4 \quad P(\text{mac}) = \frac{6}{10} = 0.6$$

$$P(\text{iPhone} | \text{mac}) = \frac{0.4}{0.6} = 0.667$$

## Định lý Bayes

**Định lý Bayes** bắt nguồn từ xác suất có điều kiện.

Định lý Bayes được đặt theo tên **Rev. Thomas Bayes** (/beɪz /; 1702-1761), người đầu tiên đã cho thấy làm thế nào để sử dụng thông tin mới để cập nhật những thông tin trước đó.

Xác suất của  $A$  với điều kiện  $B$  xảy ra được định nghĩa như sau :

$$P(A|B) = \frac{P(AB)}{P(B)}$$

$$\begin{aligned} P(A|B) &= P(AB)/P(B) \\ \Rightarrow P(AB) &= P(A|B) * P(B) \end{aligned}$$

$$\begin{aligned} P(B|A) &= P(AB)/P(A) \\ \Rightarrow P(AB) &= P(B|A) * P(A) \\ P(A|B) &= (P(B|A) * P(A)) / P(B) \end{aligned}$$

## Định lý Bayes

**Định lý Bayes** cho phép tính xác suất xảy ra của một sự kiện ngẫu nhiên  $A$  khi biết sự kiện liên quan  $B$  đã xảy ra. Xác suất này được ký hiệu là  $P(A|B)$ , và đọc là "xác suất của  $A$  nếu có  $B$ ".

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \frac{\text{likelihood} * \text{prior}}{\text{normalizing\_constant}}$$

## Định lý Bayes

Theo định lý Bayes, xác suất xảy ra  $A$  khi biết  $B$  sẽ phụ thuộc vào 3 yếu tố:

- Xác suất xảy ra  $A$  của riêng nó, không quan tâm đến bất kỳ thông tin nào của  $B$ . Ký hiệu là  $P(A)$ . Đại lượng này còn gọi là tiên nghiệm (**prior**)
- Xác suất xảy ra  $B$  của riêng nó, không quan tâm đến  $A$ . Ký hiệu là  $P(B)$ . Đại lượng này còn gọi là hằng số chuẩn hóa (**normalising constant**)
- Xác suất xảy ra  $B$  khi biết  $A$  xảy ra. Ký hiệu là  $P(B|A)$  và đọc là "xác suất của  $B$  nếu có  $A$ ". Đại lượng này gọi là khả năng (likelihood) xảy ra  $B$  khi biết  $A$  đã xảy ra.

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \frac{\text{likelihood} * \text{prior}}{\text{normalizing\_constant}}$$

## Định lý Bayes

**Định lý Bayes**  $P(H|E) = \frac{P(E|H).P(H)}{P(E)}$

**Evidence**  $E = [E_1, E_2, \dots, E_n]$  thuộc tính của dữ liệu cần dự báo

**Event**  $H$ : giá trị lớp/ nhãn của dữ liệu  $E$  cần sự báo

$H$	The probability of a hypothesis
$E$	Conditional on a new piece of evidence
$P(H E)$	The probability of a hypothesis conditional on a new evidence
$P(E H)$	The probability of the evidence given the hypothesis
$P(H)$	The prior probability of the hypothesis
$P(E)$	The prior probability of the evidence

## Giải thuật naive Bayes

### □ Ngây thơ

- các thuộc tính (biến) có độ quan trọng như nhau
- các thuộc tính (biến) độc lập thống kê

### □ Nhận xét

- Giả thiết các thuộc tính độc lập không bao giờ đúng
- nhưng trong thực tế, naive Bayes cho kết quả khá tốt

17

## Nội dung

- Giới thiệu về Bayesian classification
- Kiến thức về xác suất thống kê
- **Giải thuật học của naive Bayes**
- Kết luận và hướng phát triển

18

## Luật Bayes

Định lý xác suất Bayes

$$P(H|E) = \frac{P(E|H).P(H)}{P(E)}$$

**Evidence E** = [E1,E2,...,En] có n giá trị thuộc tính của dữ liệu cần dự báo

**Event H**: giá trị lớp/ nhãn của dữ liệu E cần dự báo

19

## Luật Bayes

Định lý xác suất Bayes

$$P[H | E] = \frac{P[E | H]P[H]}{P[E]}$$

**Do giả thiết: “ các thuộc tính độc lập nhau”**

$$\Rightarrow P(H|E) = \frac{P(E_1|H).P(E_2|H)....P(E_n|H).P(H)}{P(E)}$$

**Evidence E** = [E1,E2,...,En] có n thuộc tính của dữ liệu cần dự báo

**Event H**: giá trị lớp/ nhãn của dữ liệu E cần dự báo

20

## Bayes thơ ngây

**Bước 1:** học/ huấn luyện mô hình (learning Phase)

xây dựng mô hình sẵn dùng (tính sẵn xác suất xuất hiện của tất cả các trường hợp)

**Bước 2:** dự báo/ dự đoán

Khi có đối tượng/sự kiện mới xuất hiện cần phân loại : xác định nhãn của đối tượng mới đến thông qua giá trị xác suất lớn nhất tính được

## Ví dụ:

Outlook	Temp	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

22

**Ví dụ:** Dữ liệu weather, dựa trên các thuộc tính (Outlook, Temp, Humidity, Windy), quyết định (play/no)

Outlook	Temp	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

23

Dữ liệu weather, dựa trên các thuộc tính (Outlook, Temp, Humidity, Windy), quyết định (play/no)

**Bước 1**

$$P(H|E) = \frac{P(E_1|H).P(E_2|H)....P(E_n|H).P(H)}{P(E)}$$

Outlook	Temp	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

Outlook			Temperature		Humidity		Windy		Play	
	Yes	No					Yes	No	Yes	No
Sunny	2	3	Hot	2	2	High	3	4	6	2
Overcast	4	0	Mild	4	2	Normal	6	1	3	3
Rainy	3	2	Cool	3	1					
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	6/9	2/5
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	3/9	3/5
Rainy	3/9	2/5	Cool	3/9	1/5					

24

## Ví dụ

### Bước 2

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

← Evidence E

– Phần tử mới đến,

$\mathbf{x}' = (\text{Outlook}=\text{Sunny}, \text{Temperature}=\text{Cool}, \text{Humidity}=\text{High}, \text{Wind}=\text{True})$

Cần xác định: *xác suất của lớp “yes” và xác suất của lớp “no”*

$$P(H|E) = \frac{P(E_1|H).P(E_2|H)....P(E_n|H).P(H)}{P(E)}$$

25

## Ví dụ

$$P(H|E) = \frac{P(E_1|H).P(E_2|H)....P(E_n|H).P(H)}{P(E)}$$

### Bước 2

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

← Evidence E

– Phần tử mới đến,

$\mathbf{x}' = (\text{Outlook}=\text{Sunny}, \text{Temperature}=\text{Cool}, \text{Humidity}=\text{High}, \text{Wind}=\text{True})$

$$\begin{aligned} \text{Pr}[\text{yes} | E] &= \text{Pr}[\text{Outlook} = \text{Sunny} | \text{yes}] \\ &\quad \times \text{Pr}[\text{Temperature} = \text{Cool} | \text{yes}] \\ &\quad \times \text{Pr}[\text{Humidity} = \text{High} | \text{yes}] \\ &\quad \times \text{Pr}[\text{Windy} = \text{True} | \text{yes}] \\ &\quad \times \frac{\text{Pr}[\text{yes}]}{\text{Pr}[E]} \end{aligned}$$

*xác suất của lớp “yes”*

26

## Ví dụ

### Bước 2

$$\text{Pr}[\text{yes} | E] = \text{Pr}[\text{Outlook} = \text{Sunny} | \text{yes}]$$

*xác suất của lớp “yes”*

$$\times \text{Pr}[\text{Temperature} = \text{Cool} | \text{yes}]$$

$$\times \text{Pr}[\text{Humidity} = \text{High} | \text{yes}]$$

$$\times \text{Pr}[\text{Windy} = \text{True} | \text{yes}]$$

$$\times \frac{\text{Pr}[\text{yes}]}{\text{Pr}[E]}$$

$$= \frac{\frac{2}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{9}{14}}{\text{Pr}[E]}$$

$$P(\text{Outlook}=\text{Sunny} | \text{Play}=\text{Yes}) = 2/9$$

$$P(\text{Temperature}=\text{Cool} | \text{Play}=\text{Yes}) = 3/9$$

$$P(\text{Humidity}=\text{High} | \text{Play}=\text{Yes}) = 3/9$$

$$P(\text{Wind}=\text{True} | \text{Play}=\text{Yes}) = 3/9$$

$$P(\text{Play}=\text{Yes}) = 9/14$$

27

Dữ liệu weather, dựa trên các thuộc tính (Outlook, Temp, Humidity, Windy), quyết định (play/no)

Outlook			Temperature			Humidity			Windy			Play	
Yes		No	Yes		No	Yes		No	Yes		No	Yes	No
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/14	5/14
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5		
Rainy	3/9	2/5	Cool	3/9	1/5								

□ quyết định (play=yes/no)?

$$\begin{aligned} P[\text{Yes} | E] &= (2/9 \times 3/9 \times 3/9 \times 3/9 \times 9/14) / P[E] \\ &= 0.0053 / P[E] \end{aligned}$$

$$P[\text{No} | E] = 0.0206 / P[E]$$

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

=> yes/no?

28

Dữ liệu weather, dựa trên các thuộc tính (Outlook, Temp, Humidity, Windy), quyết định (play/no)

Outlook			Temperature			Humidity			Windy			Play	
Yes	No		Yes	No		Yes	No		Yes	No		Yes	No
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/14	5/14
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5		
Rainy	3/9	2/5	Cool	3/9	1/5								

□ quyết định (play=yes/no)?

$$\text{Likelihood(yes)} = 2/9 \times 3/9 \times 3/9 \times 3/9 \times 9/14 = 0.0053$$

$$\text{Likelihood(no)} = 3/5 \times 1/5 \times 4/5 \times 3/5 \times 5/14 = 0.0206$$

$$\text{Likelihood(yes)} = 0.0053 / (0.0053 + 0.0206) = 0.205$$

$$\text{Likelihood(no)} = 0.0206 / (0.0053 + 0.0206) = 0.795$$

=> yes/no?

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

29

Xác suất = 0

- giá trị của thuộc tính không xuất hiện trong tất cả các lớp sử dụng *Laplace estimator*
- xác suất không bao giờ có giá trị 0
- Cộng thêm cho tử một giá trị là  $p_i \mu$  và mẫu số giá trị  $\mu$  để tính xác suất.  $\mu$  hằng số dương và  $p_i$  là hệ số dương sao cho tổng các  $p_i = 1$  ( $i=1..n$ )

30

Laplace estimator – Ước lượng Laplace

□ VD: thuộc tính *outlook* cho lớp “no” =>  $p_1=p_2=p_3 = 1/3$ ;  $\mu=1$

$$\frac{3 + \mu / 3}{5 + \mu} \quad \frac{0 + \mu / 3}{5 + \mu} \quad \frac{2 + \mu / 3}{5 + \mu}$$

**Sunny**                      **Overcast**                      **Rainy**

Outlook			Temperature			Humidity			Windy			Play	
Yes No			Yes No			Yes No			Yes No			Yes No	
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/14	5/14
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5		
Rainy	3/9	2/5	Cool	3/9	1/5								

31

Laplace estimator – Ước lượng Laplace

□ ví dụ : thuộc tính *outlook* cho lớp “no”

$$\frac{3+1/3}{5+1} \quad \frac{0+1/3}{5+1} \quad \frac{2+1/3}{5+1}$$

**Sunny**                      **Overcast**                      **Rainy**

Outlook		
	Yes	No
Sunny	2	3
Overcast	4	0
Rainy	3	2
Sunny	2/9	3/5
Overcast	4/9	0/5
Rainy	3/9	2/5

$$p_1 = p_2 = p_3 = 1/3; \mu=1$$

$$\begin{aligned} \text{Sunny} &= 10/18 \\ \text{Overcast} &= 1/18 \\ \text{Rainy} &= 7/18 \end{aligned}$$

32



## Laplace estimator – Ước lượng Laplace

- trọng số có thể không bằng nhau, nhưng tổng phải là 1
- thuộc tính *outlook* cho lớp “Yes”

$$\frac{2 + \mu p_1}{9 + \mu} \quad \frac{4 + \mu p_2}{9 + \mu} \quad \frac{3 + \mu p_3}{9 + \mu}$$

**Sunny                  Overcast                  Rainy**

Đề xuất giá trị  $p_1, p_2, p_3$  và  $\mu$

33

## Laplace estimator – Ước lượng Laplace

Ước lượng Laplace cho trường hợp sau ( $\mu, p_i = ?$ )

	A	B	C
T1	1/7	2/10	5/13
T2	2/7	1/10	3/13
T3	1/7	2/10	0/13
T4	3/7	5/10	5/13

34

## Giá trị thuộc tính nhiều

- học : bỏ qua dữ liệu nhiều
- phân lớp : bỏ qua các thuộc tính nhiều
- ví dụ :

Outlook	Temp.	Humidity	Windy	Play
?	Cool	High	True	?

$$\begin{aligned} \text{Likelihood}(\text{yes}) &= 3/9 \times 3/9 \times 3/9 \times 9/14 = 0.0238 \\ \text{Likelihood}(\text{no}) &= 1/5 \times 4/5 \times 3/5 \times 5/14 = 0.0343 \\ \text{Likelihood}(\text{yes}) &= 0.0238 / (0.0238 + 0.0343) = 0.41 \\ \text{Likelihood}(\text{no}) &= 0.0343 / (0.0238 + 0.0343) = 0.59 \end{aligned}$$

35

## Bài tập- cho tập dữ liệu như bảng

R/D	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

Class: C1:buys\_computer= 'yes'                  C2:buys\_computer= 'no'

Dự đoán nhãn của phần tử X1 =(age=youth, Income=medium, Student=yes, Credit\_rating=Fair)

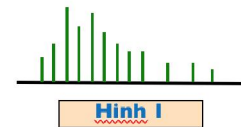
Dự đoán nhãn của phần tử X2 =(age=middle\_aged, Student=yes, Credit\_rating=Fair)

RID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

## Giá trị rời rạc và liên tục

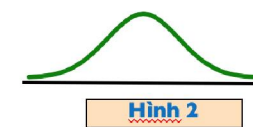
### ❖ Rời rạc

- Màu sắc
- Giới tính
- Tôn giáo



### ❖ Liên tục

- Chiều cao
- Cân nặng
- Thời gian hoàn thành công việc



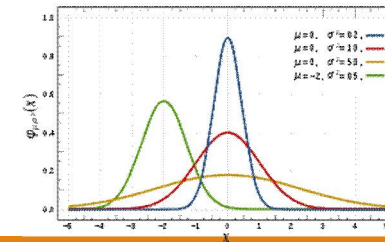
Xác định dữ liệu trong bảng kế tiếp, giá trị của các thuộc tính là giá trị rời rạc hay liên tục?

Outlook	Temp	Humidity	Windy	Play
Sunny	85	85	False	No
Sunny	80	90	True	No
Overcast	83	86	False	Yes
Rainy	70	96	False	Yes
Rainy	68	80	False	Yes
Rainy	65	70	True	No
Overcast	64	65	True	Yes
Sunny	72	95	False	No
Sunny	69	70	False	Yes
Rainy	75	80	False	Yes
Sunny	75	70	True	Yes
Overcast	72	90	True	Yes
Overcast	81	75	False	Yes
Rainy	71	91	True	No

## Dữ liệu liên tục

**Phân phối chuẩn**, còn gọi là **phân phối Gauss**, là một **phân phối xác suất** cực kì quan trọng trong nhiều lĩnh vực. Nó là họ phân phối có dạng tổng quát giống nhau, chỉ khác **tham số vị trí** (**giá trị trung bình**  $\mu$ ) và **tỉ lệ** (**phương sai**  $\sigma^2$ ).

**Phân phối chuẩn tắc** (*standard normal distribution*) là phân phối chuẩn với giá trị trung bình bằng 0 và phương sai bằng 1 (đường cong màu đỏ trong hình). Phân phối chuẩn còn được gọi là **đường cong chuông** (*bell curve*) vì đồ thị của **mật độ xác suất** có dạng **chuông**.



## Play tennis dataset

Outlook	Temp	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

## Dữ liệu liên tục

- Giả sử các thuộc tính có phân phối *Gaussian*
- hàm mật độ xác suất  $f(x)$  được tính như sau

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

➤ Mean  $\mu$

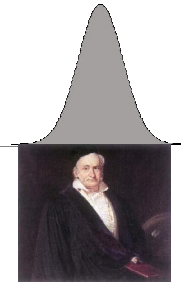
$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

➤ Phương sai (Variance)  $\sigma^2$

$$\sigma^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2$$

➤ Độ lệch chuẩn - standard deviation: căn bậc 2 của phương sai

$$\sigma = \sqrt{\sigma^2}$$



Karl Gauss, 1777-1855  
great German mathematician

<https://www.mathsisfun.com/data/standard->

42

Outlook	Temp.	Humidity	Windy	Play
Sunny	66	90	true	?

Outlook	Temp	Humidity	Windy	Play
Sunny	85	85	False	No
Sunny	80	90	True	No
Overcast	83	86	False	Yes
Rainy	70	96	False	Yes
Rainy	68	80	False	Yes
Rainy	65	70	True	No
Overcast	64	65	True	Yes
Sunny	72	95	False	No
Sunny	69	70	False	Yes
Rainy	75	80	False	Yes
Sunny	75	70	True	Yes
Overcast	72	90	True	Yes
Overcast	81	75	False	Yes
Rainy	71	91	True	No

## Bước 1: huấn luyện mô hình

Outlook	Temp	Humidity	Windy	Play
Sunny	85	85	False	No
Sunny	80	90	True	No
Overcast	83	86	False	Yes
Rainy	70	96	False	Yes
Rainy	68	80	False	Yes
Rainy	65	70	True	No
Overcast	64	65	True	Yes
Sunny	72	95	False	No
Sunny	69	70	False	Yes
Rainy	75	80	False	Yes
Sunny	75	70	True	Yes
Overcast	72	90	True	Yes
Overcast	81	75	False	Yes
Rainy	71	91	True	No



The numeric weather data with summary statistics												
outlook			temperature			humidity			windy			play
sunny	2/9	3/5	mean	73	74.6	mean	79.1	86.2	false	6/9	2/5	9/14 5/14
overcast	4/9	0/5	std dev	6.2	7.9	std dev	10.2	9.7	true	3/9	3/5	
rainy	3/9	2/5	$\sigma^2$	38..44								

Outlook	Temp.	Humidity	Windy	Play
Sunny	66	90	true	?

$$P(H|E) = \frac{P(E_1|H) \cdot P(E_2|H) \dots P(E_n|H) \cdot P(H)}{P(E)}$$

$$P[\text{Yes} | E] = (P(\text{Outlook}=\text{Sunny} | \text{Play}=\text{Yes}) \times P(\text{Temp.}=66 | \text{Play}=\text{Yes}) \times P(\text{Hum.}=90 | \text{Play}=\text{Yes}) \times P(\text{Wind}=\text{True} | \text{Play}=\text{Yes}) \times P(\text{Play}=\text{Yes})) / P[E]$$

$$P(\text{Outl}=\text{Sunny} | \text{Play}=\text{Yes}) = 2/9$$

$$P(\text{Temp.}=66 | \text{Play}=\text{Yes}) = ??$$

$$P(\text{Hum.}=90 | \text{Play}=\text{Yes}) = ??$$

$$P(\text{Wind}=\text{True} | \text{Play}=\text{Yes}) = 3/9$$

$$P(\text{Play}=\text{Yes}) = 9/14$$

The numeric weather data with summary statistics												
outlook			temperature			humidity			windy			play
sunny	2/9	3/5	mean	73	74.6	mean	79.1	86.2	false	6/9	2/5	9/14 5/14
overcast	4/9	0/5	std dev	6.2	7.9	std dev	10.2	9.7	true	3/9	3/5	
rainy	3/9	2/5	$\sigma^2$	38..44								

Outlook	Temp.	Humidity	Windy	Play
Sunny	66	90	true	?

$$P(\text{Outl}=\text{Sunny} | \text{Play}=\text{Yes}) = 2/9$$

$$P(\text{Temp.}=66 | \text{Play}=\text{Yes}) = ??$$

$$P(\text{Hum.}=90 | \text{Play}=\text{Yes}) = ??$$

$$P(\text{Wind}=\text{True} | \text{Play}=\text{Yes}) = 3/9$$

$$P(\text{Play}=\text{Yes}) = 9/14$$

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

The numeric weather data with summary statistics												
outlook			temperature			humidity			windy			play
sunny	2/9	3/5	mean	73	74.6	mean	79.1	86.2	false	6/9	2/5	9/14 5/14
overcast	4/9	0/5	std dev	6.2	7.9	std dev	10.2	9.7	true	3/9	3/5	
rainy	3/9	2/5	$\sigma^2$	38..44								

Outlook	Temp.	Humidity	Windy	Play
Sunny	66	90	true	?

$$P(\text{Outl}=\text{Sunny} | \text{Play}=\text{Yes}) = 2/9$$

$$P(\text{Temp.}=66 | \text{Play}=\text{Yes}) = 0.034$$

$$P(\text{Hum.}=90 | \text{Play}=\text{Yes}) = ??$$

$$P(\text{Wind}=\text{True} | \text{Play}=\text{Yes}) = 3/9$$

$$P(\text{Play}=\text{Yes}) = 9/14$$

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

$$f(\text{temperature} = 66 | \text{yes}) = \frac{1}{\sqrt{2\pi}6.2} e^{-\frac{(66-73)^2}{2 \cdot 6.2^2}} = 0.0340$$

The numeric weather data with summary statistics												
outlook			temperature			humidity			windy			play
sunny	2/9	3/5	mean	73	74.6	mean	79.1	86.2	false	6/9	2/5	9/14 5/14
overcast	4/9	0/5	std dev	6.2	7.9	std dev	10.2	9.7	true	3/9	3/5	
rainy	3/9	2/5	$\sigma^2$	38..44								

Outlook	Temp.	Humidity	Windy	Play
Sunny	66	90	true	?

$$P(\text{Outl}=\text{Sunny} | \text{Play}=\text{Yes}) = 2/9$$

$$P(\text{Temp.}=66 | \text{Play}=\text{Yes}) = 0.034$$

$$P(\text{Hum.}=90 | \text{Play}=\text{Yes}) = ??$$

$$P(\text{Wind}=\text{True} | \text{Play}=\text{Yes}) = 3/9$$

$$P(\text{Play}=\text{Yes}) = 9/14$$

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

$$f(\text{temp}=66/\text{Yes}) = ?$$

$$f(\text{temp}=66/\text{No}) = ?$$

$$f(\text{humidity}=90/\text{Yes}) = ?$$

$$f(\text{humidity}=90/\text{No}) = ?$$

The numeric weather data with summary statistics												
outlook			temperature			humidity			windy			play
sunny	2/9	3/5	mean	73	74.6	mean	79.1	86.2	false	6/9	2/5	9/14 5/14
overcast	4/9	0/5	std dev	6.2	7.9	std dev	10.2	9.7	true	3/9	3/5	
rainy	3/9	2/5	$\sigma^2$	38..44								

$$P(\text{Outl}=\text{Sunny} | \text{Play}=\text{Yes}) = 2/9$$

$$P(\text{Temp.}=66 | \text{Play}=\text{Yes}) = 0.034$$

$$P(\text{Hum.}=90 | \text{Play}=\text{Yes}) = 0.0221$$

$$P(\text{Wind}=\text{True} | \text{Play}=\text{Yes}) = 3/9$$

$$P(\text{Play}=\text{Yes}) = 9/14$$

$$f(\text{humidity}=90/\text{Yes}) =$$

Outlook	Temp.	Humidity	Windy	Play
Sunny	66	90	true	?

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

The numeric weather data with summary statistics												
outlook			temperature			humidity			windy			play
sunny	2/9	3/5	mean	73	74.6	mean	79.1	86.2	false	6/9	2/5	9/14 5/14
overcast	4/9	0/5	std dev	6.2	7.9	std dev	10.2	9.7	true	3/9	3/5	
rainy	3/9	2/5	$\sigma^2$	38..44								

$$P(\text{Outl}=\text{Sunny} | \text{Play}=\text{Yes}) = 2/9$$

$$P(\text{Temp.}=66 | \text{Play}=\text{Yes}) = 0.034$$

$$P(\text{Hum.}=90 | \text{Play}=\text{Yes}) = 0.0221$$

$$P(\text{Wind}=\text{True} | \text{Play}=\text{Yes}) = 3/9$$

$$P(\text{Play}=\text{Yes}) = 9/14$$

Outlook	Temp.	Humidity	Windy	Play
Sunny	66	90	true	?

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

$$f(\text{temp}=66/\text{Yes}) = 0.034$$

$$f(\text{humidity}=90/\text{Yes}) = 0.0221$$

$$f(\text{temp}=66/\text{No}) = 0.0291$$

$$f(\text{humidity}=90/\text{No}) = 0.0380$$

Nhãn????

Outlook	Temp.	Humidity	Windy	Play
Sunny	66	90	true	?

$$P(H|E) = \frac{P(E_1|H).P(E_2|H)....P(E_n|H).P(H)}{P(E)}$$

$$f(\text{temp}=66/\text{Yes}) = 0.034$$

$$f(\text{humidity}=90/\text{Yes}) = 0.0221$$

$$f(\text{temp}=66/\text{No}) = 0.0291$$

$$f(\text{humidity}=90/\text{No}) = 0.0380$$

The numeric weather data with summary statistics												
outlook			temperature			humidity			windy			play
sunny	2/9	3/5	mean	73	74.6	mean	79.1	86.2	false	6/9	2/5	9/14 5/14
overcast	4/9	0/5	std dev	6.2	7.9	std dev	10.2	9.7	true	3/9	3/5	
rainy	3/9	2/5	$\sigma^2$	38..44								

55

Dữ liệu liên tục

#### □ Bước 2- dự đoán

Outlook	Temp.	Humidity	Windy	Play
Sunny	66	90	true	?

$$\text{Likelihood}(\text{yes}) = 2/9 \times 0.0340 \times 0.0221 \times 3/9 \times 9/14 = 0.000036$$

$$\text{Likelihood}(\text{no}) = 3/5 \times 0.0291 \times 0.0380 \times 3/5 \times 5/14 = 0.000136$$

$$\text{Likelihood}(\text{yes}) = 0.000036 / (0.000036 + 0.000136) = 20.9\%$$

$$\text{Likelihood}(\text{no}) = 0.000136 / (0.000036 + 0.000136) = 79.1\%$$

56

## Multinomial Naive Bayes

- Mô hình này chủ yếu được sử dụng trong phân loại văn bản mà feature vectors được tính bằng [Bags of Words](#).
- Mỗi văn bản được biểu diễn bởi một vector có độ dài d chính là số từ trong từ điển.
- Giá trị của thành phần thứ i trong mỗi vector chính là số lần từ thứ i xuất hiện trong văn bản đó

$$p(x_i|c) = \frac{N_{ci}}{N_c}$$

- $N_{ci}$  là tổng số lần từ thứ i xuất hiện trong các văn bản của class c, nó được tính là tổng của tất cả các thành phần thứ i của các feature vectors ứng với class c.
- $N_c$  là tổng số từ (kể cả lặp) xuất hiện trong class c. Nói cách khác, nó bằng tổng độ dài của toàn bộ các văn bản thuộc vào class c.

## Bernoulli Naive Bayes

Mô hình này được áp dụng cho các loại dữ liệu mà mỗi thành phần là một giá trị binary - bằng **0** hoặc **1**. Ví dụ: cũng với loại văn bản nhưng thay vì đếm tổng số lần xuất hiện của 1 từ trong văn bản, ta chỉ cần quan tâm từ đó có xuất hiện hay không

Khi đó,  $p(x_i|c)$  được tính bằng:

$$p(x_i|c) = p(i|c)^{x_i} (1 - p(i|c))^{1-x_i}$$

$p(i|c)$  có thể được hiểu là xác suất từ thứ "i" xuất hiện trong các văn bản của lớp "c"

## Nội dung

- Giới thiệu về Bayesian classification
- Giải thuật học của naive Bayes
- Kết luận và hướng phát triển

## Kết luận

- naïve Bayes
  - cho kết quả tốt trong thực tế mặc dù chịu những giả thiết về tính độc lập thống kê của các thuộc tính
  - phân lớp không yêu cầu phải ước lượng một cách chính xác xác suất
  - dễ cài đặt, học nhanh, kết quả dễ hiểu
  - sử dụng trong phân loại text, spam, etc
  - tuy nhiên khi dữ liệu có nhiều thuộc tính dư thừa thì naïve Bayes không còn hiệu quả
  - dữ liệu liên tục có thể không tuân theo phân phối chuẩn (=> kernel density estimators)

## Hướng phát triển

- naïve Bayes
  - chọn thuộc tính con từ các thuộc tính ban đầu
  - chỉ sử dụng các thuộc tính con để học phân lớp
  - mạng Bayes : mối liên quan giữa các thuộc tính

61



Cám ơn !