

# Introduction

lecture 13

Note  $P_\theta(X, Y) = P(X, Y; \theta)$

$$(x_i, y_i) \sim P(X, Y)$$

$$P(X, Y; \theta) \leftarrow \text{learn} - \forall x \forall y P(X=x \wedge Y=y)$$

$$\text{Classifier } h(x_{\text{test}}) = \underset{y}{\operatorname{argmax}} P(Y=y | X=x_{\text{test}}; \theta)$$

Need to learn

$$P(Y|X)$$

assume it can access

Bayes

① optimal classifier  $h(x_t) = \underset{y}{\operatorname{argmax}} P(Y=y | X=x_t)$ ; probs ที่  $X=x_t$  แล้ว  $Y$  เป็นของใคร (เช่น salmon, mackael)  
 probs ของ salmon เมื่อ  $X=x_{\text{test}}$   
 probs ของ mackael เมื่อ  $X=x_{\text{test}}$

$$\textcircled{3} \begin{cases} P(Y|X) \\ \rightarrow P(Y=y | X=\vec{x}) \end{cases}$$

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \quad X = \begin{bmatrix} [X]_1 \\ [X]_2 \\ \vdots \\ [X]_d \end{bmatrix}$$

r.v. ที่อยู่ใน  $d=1$

$$P(Y=y | X=\vec{x}) = P(Y=y | [X]_1=x_1, [X]_2=x_2, \dots, [X]_d=x_d) \quad \text{way 1}$$

$$\text{Bayes rule: } P(Y=y | X=x) = \frac{P(X=x | Y=y) \cdot P(Y=y)}{P(X=x)} \quad \text{way 2} \rightarrow \text{Naive Bayes use.}$$

estimate  $P(X=x | Y=y)$

$$: P([X]_1=x_1, [X]_2=x_2, \dots, [X]_d=x_d | Y=y)$$

## Naive Bayes Assumption

lecture 14

: Assumes all feature values are independent given the label.

$$P(X=x | Y=y) = \prod_{i=1}^d P([X]_i=x_i | Y=y)$$

- ได้จากการใช้ Bayes classifier

## Bayes Classifier

$$h(x) = \underset{y}{\operatorname{argmax}}_\theta P(Y=y | X=x)$$

Goal: Estimate  $P(Y|X)$  estimate  $\forall x \forall y P(Y=y | X=x)$

$$\text{Chain Rule: } P(Y=y | X=x) = \frac{P(Y=y \wedge X=x)}{P(X=x)}$$

Job: ① Estimate  $\forall x \forall y P(Y=y \wedge X=x) \approx P(X, Y)$   
 ② Estimate  $P(X=x) \forall x$

Scenario use  $P_\theta(Y=y \wedge X=x) = \text{Bin}(n, \theta) ??$

Note: ตัว  $(x_i, y_i)$  ใน probability distribution  $X, Y$  ( $P(X, Y)$ ) ถ้าสมมติว่า  $X_{\text{test}}$  ที่ออกมาจาก  $P(X, Y)$  ซึ่งเราสามารถเข้าถึงได้ จึงต้องสร้าง modeling distribution ว่าเป็น  $P(X, Y; \theta)$

Scenario use  $P_\theta(Y=y \wedge X=x) = \text{Bin}(n, \theta)$  Way: 1

Coin tossing:  $n$  times

$n_H$ : จำนวนครั้งที่ได้ออก head

$$n_H \sim \text{Bin}(n, \theta)$$

$$\text{MLE} = \theta = \frac{n_H}{n}$$

Dice  $n_3$ : number of times that we obtain face 3

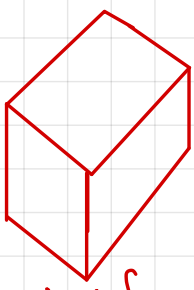
$X$  คือผลลัพธ์  $X \in \{1, 2, \dots, i\}$

$Y \in \{1, 2, \dots, j\}$

$I(E)$ : indicator random variable

Assign ค่าผลลัพธ์เป็น 1, 0

$$I(E) = \begin{cases} 1 & \text{if } E \text{ occur} \\ 0 & \text{otherwise} \end{cases}$$



→ การสุ่มค่าที่  
เกิดขึ้นคือ  $y, x$  เป็น  $Y, X$

$$\text{MLE } \theta = \frac{n_{y,x}}{n} = \frac{\sum_{i=1}^n I(Y_i=y \wedge X_i=x)}{n} \quad \textcircled{I}$$

↓  
correspondence  
 $P_\theta(Y=y \wedge X=x)$

$$\underline{P_\gamma(X=x) = \text{Bin}(n, \gamma)}$$

$$\text{MLE } \Rightarrow \gamma = \frac{\sum_{i=1}^n I(X_i=x)}{n}$$

II

Assume  $X, Y$  follow Binominal distribution

$$\text{MLE: } P_\theta(Y=y | X=x) = \frac{P_\theta(Y=y \wedge X=x)}{P_\gamma(X=x)}$$

$$= \frac{\sum_{i=1}^n I(Y_i=y \wedge X_i=x)}{n}$$

$$\frac{\sum_{i=1}^n I(X_i=x)}{n}$$

$$= \frac{\sum_{i=1}^n I(Y_i=y \wedge X_i=x)}{\sum_{i=1}^n I(X_i=x)}$$

# Problem with $P_\theta(Y=y|X=x)$

$$P_\theta(Y=y|X=x) = \frac{\sum_{i=1}^n I(Y_i=y \wedge X_i=x)}{\sum_{i=1}^n I(X_i=x)}$$

-  $X$  là 1 vector dimension  $\Rightarrow$  Vector ( $X$  là 1 vector)

$$P_\theta(Y=y|\vec{X}=\vec{x}) \Rightarrow P(Y=y|[X]_1=x^1 \wedge [X]_2=x^2 \wedge \dots \wedge [X]_d=x^d)$$

$$\vec{X} = \begin{bmatrix} [X]_1 \\ [X]_2 \\ \vdots \\ [X]_d \end{bmatrix} \quad \vec{x} = \begin{bmatrix} x^1 \\ x^2 \\ \vdots \\ x^d \end{bmatrix}$$

r.v.  $X$  = fixed  $X$  on each  $d$

$$\frac{\sum_{i=1}^n I(Y_i=y \wedge [X_i]_1=x^1 \wedge [X_i]_2=x^2 \wedge \dots \wedge [X_i]_d=x^d)}{\sum_{i=1}^n I([X_i]_1=x^1 \wedge \dots \wedge [X_i]_d=x^d)}$$

Note: When  $d \gg 0$  and  $n \rightarrow +\infty$

$$\Rightarrow P_\theta(Y=y \wedge X=x) = \frac{1}{n} = 0$$

$$\Rightarrow P_\theta(X=x) = \frac{1}{n} = 0$$

So:  $P_\theta(Y=y|X=x) = \frac{0}{0}$  undifind

## Apply Bayes' rule to Bayes Classifier ( $P_\theta(Y=y|X=x)$ ) way: 2

$$\text{Bayes' rule: } P(Y=y|X=x) = \frac{P(X=x|Y=y)P(Y=y)}{P(X=x)}$$

Bayes Classifier

$$= h(x) = \underset{y}{\operatorname{argmax}} P(Y=y|X=x)$$

$$= \underset{y}{\operatorname{argmax}} \frac{P(X=x|Y=y) \cdot P(Y=y)}{\cancel{P(X=x)}} \quad \text{binomial}$$

$$= \underset{y}{\operatorname{argmax}} P(X=x|Y=y) \cdot P(Y=y)$$

Estimate  $P(Y=y) \rightarrow$  binomial binary, multiclass classified

$$P_\theta(Y=y) = \frac{\sum_{i=1}^n I(Y_i=y)}{n}$$

- Estimate  $P(\vec{X}=\vec{x}|Y=y)$

$$= P([X]_1=x^1 \wedge [X]_2=x^2 \wedge \dots \wedge [X]_d=x^d | Y=y)$$

Can't estimate direct

So use Naive Bayes assumption.

Assume  $X, Y$  follow Binomial distributions

$$\Rightarrow P_\theta(Y=y) = \frac{\sum_{i=1}^n I(Y_i=y)}{n}$$

$$\Rightarrow P_\theta(X=x|Y=y) = ?$$

## Naive Bayes assumption:

- All feature values are Independent

$$P(\vec{X}=\vec{x} | Y=y) = \prod_{i=1}^d P([X]_i=x^i | Y=y) \quad ; \text{ Prob of } Y \text{ given } y \text{ is } P(Y=y) \text{ and } P([X]_i=x^i | Y=y) \text{ is } P([X]_i=x^i | Y=y)$$

Naive Bayes classifier

$$h(X) = \underset{y}{\operatorname{argmax}} \prod_{\alpha=1}^d P([X]_\alpha=x^\alpha | Y=y) P(Y=y)$$

## How to estimate $P([X]_\alpha | Y)$ 3 cases

- There are 3 notable cases:

Case 1: Categorical features: Categorical Naive Bayes Classifier

$$[X]_\alpha \in \{c_1, c_2, \dots, c_k\}$$

eg: {male, female}  
: {single, widowed, married}

We model  $P([X]_\alpha=j | Y=y) = [\theta_{jy}]_\alpha$  parameter  $\theta_{jy}$  is the probability of feature  $\alpha$  having value  $j$  given the label is  $y$

The probability of feature  $\alpha$  having value  $j$  given the label is  $y$

MLE estimate  $\Rightarrow [\theta_{jy}]_\alpha = \frac{\text{\# of sample with label } y \text{ that has feature } \alpha \text{ with value } j}{\text{\# of samples with label } y}$

$$\Rightarrow \frac{\sum_{i=1}^n I(Y_i=y) \cdot I(X_i^\alpha=j)}{\sum_{i=1}^n I(Y_i=y)} : \text{count of } Y_i=y \text{ and } X_i^\alpha=j$$

## Quiz 4

The following table is a result from observing the behavior of a person whether he went out or stayed home given the two weather conditions (sunny or rainy) and the two options regarding his car status (car-broken or car-working)

- Categorical feature {
- $y_i \in \{go-out, stayhome\}$
  - $x_i^1 \in \{sunny, rainy\}$
  - $x_i^2 \in \{car-broken, car-working\}$

Estimate  $P(\vec{x}_i | y)$

- $\rightarrow P(\vec{x}_i = sunny | y = go-out)$
- $\rightarrow P(\vec{x}_i = rainy | y = go-out)$
- $\rightarrow P(\vec{x}_i = sunny | y = stay)$
- $\rightarrow P(\vec{x}_i = rainy | y = stay)$

$$P(\vec{x}_i = sunny | y = go-out) = [\theta_{sunny, go-out}]_i = \frac{4}{5} = 0.8$$

$i$	$x_i^1$	$x_i^2$	$y_i$
1	sunny ✓	car-broken	go-out
2	rainy	car-working	go-out
3	sunny ✓	car-broken	go-out
4	sunny ✓	car-broken	go-out
5	sunny ✓	car-broken	go-out
6	sunny	car-working	stay home
7	rainy	car-working	stay home
8	rainy	car-broken	stay home
9	sunny	car-working	stay home
10	rainy	car-working	stay home

Assume that we are using Binomial distribution as the modeling distribution. You are to demonstrate solutions to the following questions.

1. Estimate  $P(y=go-out)$ .
2. Estimate  $P(y=stay home)$ .
3. What is the estimate of  $P(y)$ ?
4. What is the estimate of  $P(x)$ ?
5. Estimate  $P(x = (rainy, car-working) \text{ and } y=go-out)$ .
6. Estimate  $P(y=go-out | x = (rainy, car-working))$  directly.
7. Estimate  $P(x = (rainy, car-working) | y=go-out)$  using Naive Bayes assumption.
8. By using Naive Bayes assumption, what would be the return of  $h(x = (sunny, car-broken))$ ?

## Case 2: Multinomial feature: Multinomial Naive Bayes Classifier

eg. text data: "An ant is animal."  $\xrightarrow{\text{Back of word}}$   $\vec{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_j \\ \vdots \\ x_d \end{bmatrix}$   $\leftarrow$  # of word  $i$  appearing in the text  
: its count of word in sentence.  
coordinate

- we estimate  $P([X]_{\alpha} = j | Y = y)$  by using multinomial distribution

eg. Spam filter

-  $y \in \{\text{spam}, \text{ham}\}$

-  $\vec{X}$  represent text data (B.O.W)

$$\vec{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_j \\ \vdots \\ x_d \end{bmatrix} \begin{matrix} - w_1 \\ \\ - w_j \\ \\ - w_d \end{matrix}$$

Estimate  $P([X]_{\alpha} = j | Y = y)$ ;  $\alpha = 5, j = 10, y = \text{spam}$  Given  $W_5 = \text{Princess (เจ้าหญิง)}$

$P([X]_5 = 10 | Y = \text{spam})$ ; given spam ข้อความน่าจะเป็น coordinate 5 = 10

Modeling

$$P([X]_{\alpha} = j | Y = y) = \left[ \binom{m}{j} \cdot P(w_{\alpha} | Y = y) \right]^j$$

eg.  $P([X]_2 = 3 | Y = \text{spam}) = \left[ \binom{5}{3} \times P(w_2 | Y = \text{spam}) \right]^3$

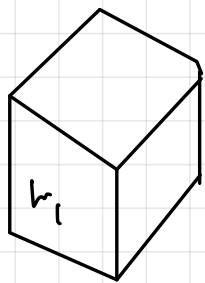
$\downarrow$   
word no 2

-  $P(w_2 | Y = y)$  is the prob. of selecting word  $w_{\alpha}$  given the label is  $y$ .

-  $m$  is the number of words in total ( $m = \sum_{\alpha=1}^d [X]_{\alpha}$ )



# Multinomial distribution



d faces  
(y = spam)

$$\tilde{x} = \begin{pmatrix} 1 & x^1 \\ 4 & x^2 \\ \vdots & \vdots \\ x^d & x^d \end{pmatrix} \sim W_1$$

Assume throw and get R1:

$$W_2, W_2, W_1, W_2, W_2, W_3, W_2$$

m times = 7

R2:

$$W_1, W_2, W_2, W_3, W_2, W_2, W_2$$

n times = 7

R1, R2 not a same words but  
with power of Back of words  
it give same feature Vector.  
(Order of word is useless)

- We model  $P(W_\alpha | Y = \text{spam}) = [\theta_{\text{spam}}]_\alpha$

$$P(W_\alpha | Y = \text{ham}) = [\theta_{\text{ham}}]_\alpha$$

- Estimate

$$\text{MLE: } [\theta_{\text{spam}}]_\alpha = \frac{\sum_{i=1}^n I(Y_i = \text{SPAM}) \cdot x_i^\alpha}{\sum_{i=1}^n I(Y_i = \text{SPAM}) \cdot \left( \sum_{d=1}^d x_i^d \right)}$$

จำนวนครั้งที่  $w_\alpha$  ปรากฏใน n sample  
ที่ใช่ spam

$$[\theta_y]_\alpha \forall y \forall \alpha$$

eg.  $[\theta_{\text{spam}}]_{\text{money}} = P(\text{money} | Y = \text{spam})$

จำนวนของ words ปรากฏ  
ใน n sample ที่ใช่ spam

$$P([X]_{\text{money}} = 2 | Y = \text{spam}) = \binom{m}{2} \cdot [P(\text{money} | Y = \text{spam})]^2$$

∴ เราหาว่าเงิน money 2 ครั้ง

1. ค่าที่อยู่ที่ตรงไหนบ้าง เราใช้  $\binom{m}{2}$  คำนวณ

## Summary

Bayes Rule  $\rightarrow h(x) = \arg \max_y P(X=x | Y=y) \cdot P(Y)$

Naive  $\rightarrow$   
Bayes  
Classifier

$$= \arg \max_y \prod_{\alpha=1}^d P([X]_\alpha = x^\alpha | Y=y) \cdot P(Y)$$

# Spam filter (text Classification)

"An ant is an animal" ; 5 words

Bag of word  $\rightarrow \vec{x} = \begin{bmatrix} 0 \\ 2 \\ \vdots \end{bmatrix}$   $\begin{matrix} w_1 = a \\ w_2 = an \\ \vdots \\ w_d = \end{matrix}$

eg:  $P([X]_2 = 2 \mid y = \text{spam}) = \binom{5}{2} [P(w_2 \mid y = \text{spam})]$

Estimate

$$P(w_2 \mid y = \text{spam}) = [\theta_{\text{spam}}]_2$$

MLE  $\rightarrow [\theta_{\text{spam}}]_2 = \frac{\sum_{i=1}^n I(y_i = \text{spam}) \cdot x_i^2}{\sum_{i=1}^n I(y_i = \text{spam}) \left( \sum_{b=1}^d x_i^b \right)}$

In general

$$P([X]_\alpha = j \mid y = \text{spam}) = \binom{m}{j} (P(w_\alpha \mid y = \text{spam}))^j$$

$$[\theta_y]_\alpha = \frac{\sum_{i=1}^n I(y_i = y) \cdot x_i^\alpha}{\sum_{i=1}^n I(y_i = y) \left( \sum_{b=1}^d x_i^b \right)}$$

$\rightarrow$  จน. ครั้งใดที่ word ปรากฏอยู่บ่อย

$$\sum_{i=1}^n I(y_i = y) \left( \sum_{b=1}^d x_i^b \right)$$

$\rightarrow$  จน. ของ word ทั้งหมดที่อยู่ใน y



$$h(X) = \underset{Y}{\operatorname{argmax}} \prod_{\alpha=1}^d P([X]_{\alpha} = X^{\alpha} | Y=Y) \cdot P(Y)$$

Note:  $m = \sum_{\alpha=1}^d X^{\alpha}$

$$= P([X]_1 = x^1 | Y=Y) \times \dots \times P([X]_d = x^d | Y=Y)$$

$$= \binom{m}{x^1} \times ([\Theta Y]_1)^{x^1} \cdot \left( \binom{m-x^1}{x^2} \times ([\Theta Y]_2)^{x^2} \right)$$

$$\times \dots \times \left( \binom{m-x^1-\dots-x^{d-1}}{x^d} \times ([\Theta Y]_d)^{x^d} \right)$$

$$\Rightarrow \left( \frac{m!}{(m-x^1)! x^1!} \times \frac{(m-x^1)!}{(m-x^1-x^2)! x^2!} \times \dots \times \frac{(m-x^1-x^2-\dots-x^{d-1})!}{(m-x^1-\dots-x^d)! x^d!} \right)$$

$$\left( \prod_{\alpha=1}^d ([\Theta Y]_{\alpha})^{x^{\alpha}} \right)$$

$$m = \sum_{i=1}^d x^i$$

↓

0!

$$\prod_{\alpha=1}^d P([X]_{\alpha} = x^{\alpha} | y = \gamma) = \frac{m!}{x^1! x^2! \dots x^d!} \cdot \prod_{\alpha=1}^d ([\theta_{\gamma}]_{\alpha})^{x^{\alpha}}$$

## Binary Classification

$$\frac{P(y = \text{spam}) \times \prod_{\alpha=1}^d P([X]_{\alpha} = x^{\alpha} | y = \text{spam})}{P(y = \text{ham}) \times \prod_{\alpha=1}^d P([X]_{\alpha} = x^{\alpha} | y = \text{ham})} = \frac{\frac{m!}{x^1! x^2! \dots x^d!} \prod_{\alpha=1}^d ([\theta_{\text{spam}}]_{\alpha})^{x^{\alpha}} \cdot P(y = \text{spam})}{\frac{m!}{x^1! x^2! \dots x^d!} \prod_{\alpha=1}^d ([\theta_{\text{ham}}]_{\alpha})^{x^{\alpha}} \cdot P(y = \text{ham})}$$

take log<sub>e</sub> into both sides

$$\frac{P(y = \text{spam}) \times \prod_{\alpha=1}^d P([X]_{\alpha} = x^{\alpha} | y = \text{spam})}{P(y = \text{ham}) \times \prod_{\alpha=1}^d P([X]_{\alpha} = x^{\alpha} | y = \text{ham})}$$

$$\frac{\left( \sum_{\alpha=1}^d x^{\alpha} \log_e([\theta_{\text{spam}}]_{\alpha}) \right) + \log_e(P(y = \text{spam}))}{\left( \sum_{\alpha=1}^d x^{\alpha} \log_e([\theta_{\text{ham}}]_{\alpha}) \right) + \log_e(P(y = \text{ham}))}$$

$$\underline{P([X]_{\alpha} | y)}$$

3 notable class

① Categorical feature  $[X]_{\alpha}$  เป็นประเภท (เช่น word) ไม่ใช่จำนวนจริง

$[X]_{\alpha} \in \{0, 1, \dots, k-1\}$  ของการปรากฏ (Categorical Naive Bayes Classifier)

② multinomial feature  $[X]_{\alpha}$  แทนค่าของของการปรากฏของคำ = word

$[X]_{\alpha} \in \{0, 1, \dots, m\}$   $m = \sum_{\alpha=1}^d x^{\alpha}$  เป็นค่าจริง (จำนวนของการปรากฏ)

③ Continuous feature

$[X]_{\alpha} \in \mathbb{R}$

( $[X]_{\alpha}$  เป็นค่าต่อเนื่อง)

→ multinomial Naive Bayes Classifier

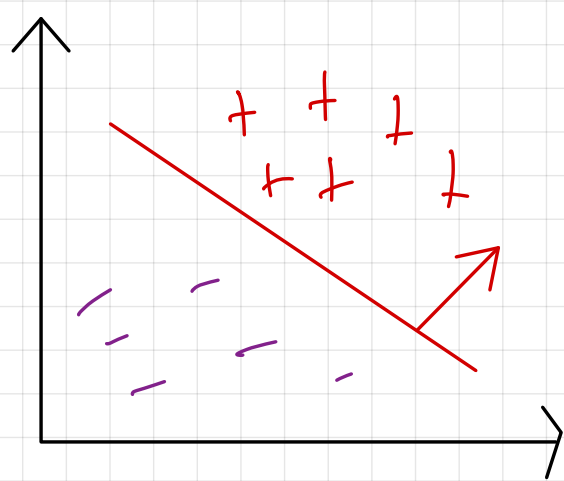
→ Gaussian Naive Bayes Classifier

# Summary of Naive Bayes

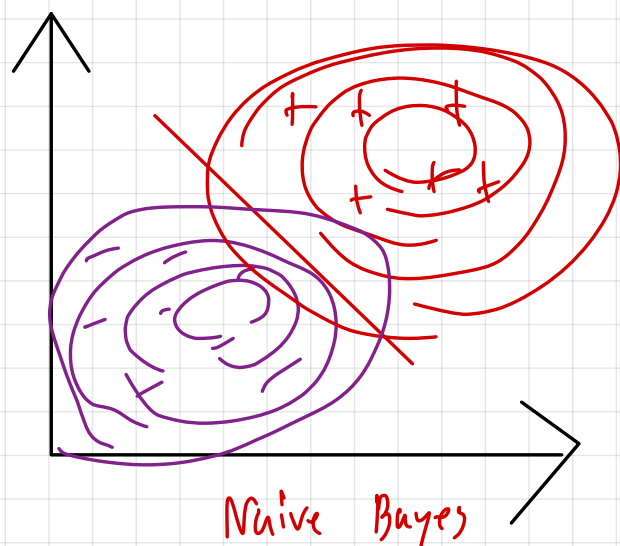
- Naive Bayes = Bayes classifier + Bayes rule + Naive Bayes assumption
- The assumption says "all feature values are Independent."

-  $\vec{w}$  is a label  $y = \text{spam}$   $P(\vec{w}_\alpha | y = \text{spam})$

- We may have data that violates the assumption.
- If our data follows multinomial distributions (features) and our task is binary classification, the Naive Bayes gives the linear decision boundary



Perceptron



Naive Bayes

- hyperplane separate our data

- find  $\vec{w}$  that linearly separate the data

-  $\vec{w}$  is distribution vs probs. in sample

- separate our distribution

find  $\vec{w}$  that separate the trained distributions

linearly separable

Discriminative learning: Try to model  $P(y|x)$  (eg. k-NN, Perceptron)

Generative learning: try to model  $P(x|y)$  and  $P(y)$  to estimate  $P(y|x)$

Both Base on Bayes Rule:  $P(y|x) = \frac{P(x|y) \times P(y)}{P(x)}$

# Prove multinomial Naive Bayes is a linear Classifier

Proof - Assume  $y \in \{-1, +1\}$

$$\begin{aligned} T &\rightarrow T = T \\ T &\rightarrow F = F \end{aligned}$$

$-h(x) = \pm 1$  iff  $P(y=+1|x) \geq P(y=-1|x)$  T.

iff  $P(x|y=+1) \times P(y=+1) \geq P(x|y=-1) \times P(y=-1)$

iff  $\prod_{\alpha=1}^d P(x_{\alpha}|y=+1) \times P(y=+1) \geq \prod_{\alpha=1}^d P(x_{\alpha}|y=-1) \times P(y=-1)$

iff  $\sum_{\alpha=1}^d \log_e P(x_{\alpha}|y=+1) + \log_e P(y=+1) \geq \sum_{\alpha=1}^d \log_e P(x_{\alpha}|y=-1) + \log_e P(y=-1)$

iff  $\sum_{\alpha=1}^d (\log_e P(x_{\alpha}|y=+1) - \log_e P(x_{\alpha}|y=-1)) + (\log_e P(y=+1) - \log_e P(y=-1)) \geq 0$

iff  $\sum_{\alpha=1}^d x_{\alpha} (\log_e [\theta_{+1}]_{\alpha} - \log_e [\theta_{-1}]_{\alpha}) + (\log_e P(y=+1) - \log_e P(y=-1)) \geq 0$

$$\vec{w}^T \vec{x}$$

$\therefore$  ในข้อนี้  $\vec{w}^T \vec{x} + b$  เป็น linear hyperplane ของ perceptron

$b$

$\vec{w} = \begin{bmatrix} \log_e [\theta_{+1}]_1 - \log_e [\theta_{-1}]_1 \\ \vdots \\ \log_e [\theta_{+1}]_d - \log_e [\theta_{-1}]_d \end{bmatrix}; \vec{x} = x_{\alpha}$  iff  $\vec{w}^T \vec{x} + b \geq 0$

### Case 3 : Continuous feature: : Gaussian Naive Bayes Classifier

-  $[X]_{\alpha} \in \mathbb{R}$

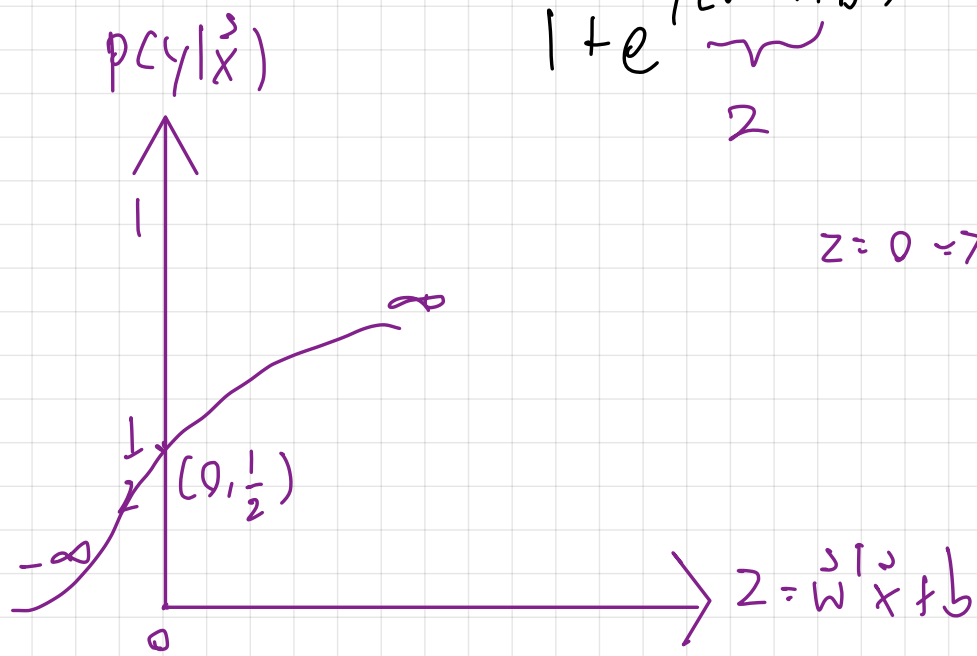
- Model  $P([X]_{\alpha} = j | Y = y)$

$$[X]_{\alpha} \sim N(\mu_y, [\sigma_y]_{\alpha})$$

- MLE  $\rightarrow [\mu_y] = \frac{\sum_{i=1}^n I(Y_i = y) \cdot X_i^{\alpha}}{n}$  same with  $\frac{\sum X}{n}$

- For Gaussian Naive Bayes, we will arrive the following expression by taking the same derivation

$$P(Y | \vec{x}) = \frac{1}{1 + e^{\frac{-y(\vec{w}^T \vec{x} + b)}{2}}}, \text{ for } y \in \{-1, +1\}$$



$$z = 0 \Rightarrow e^{-yz} = e^0 = 1$$

Recall before Prove

$$P(Y|x) \propto P(x|y) \times P(y)$$

Discriminative Learning: Try to model  $P(Y|x)$  (eg. k-NN, Perceptron)

Generative Learning: try to model  $P(x|y)$  and  $P(y)$  to estimate  $P(Y|x)$

Both Base on Bayes Rule:  $P(Y|x) = \frac{P(x|y) \times P(y)}{P(x)}$

- Discriminative Learning: Try to model  $P(Y|x)$  directly

- Generative Learning: Try to model  $P(x|y)$  and  $P(y)$

eg. Perceptron is a discriminative algorithm

$$P(Y|x) = \begin{cases} 1 & \text{if } \vec{w}^T \vec{x} \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

probability



$\therefore$   $\vec{w}^T \vec{x} \geq 0$  for  $y = +1$  and  $\vec{w}^T \vec{x} < 0$  for  $y = -1$

eg. Naive Bayes is a generative algorithm

try to model distribution  $\begin{cases} P(y) \\ P(x|y) = \prod_{\alpha=1}^d P(x_{\alpha}|y) \end{cases}$



linear classifier: A classifier  $h(x)$  is called 'linear' if  $h(x) = \pm 1$   
if and only if  $\exists w, b$  such that

$$w^T x + b \geq 0 \quad ; \text{ assume } y \in \{1, -1\}$$

$$\therefore \text{ với } b \gg \|w\| \|x\| \Rightarrow w^T x \geq 0$$

eg. Perceptron, Multinomial Naive Bayes  
are linear classifiers

Note:  $h(x)$  is linear iff  
 $h(x) = \pm 1$  iff  $\exists w, b$   
s.t.  $w^T x + b \geq 0$

- By taking the similar derivation, we can derive the following  
expression for Gaussian Naive Bayes

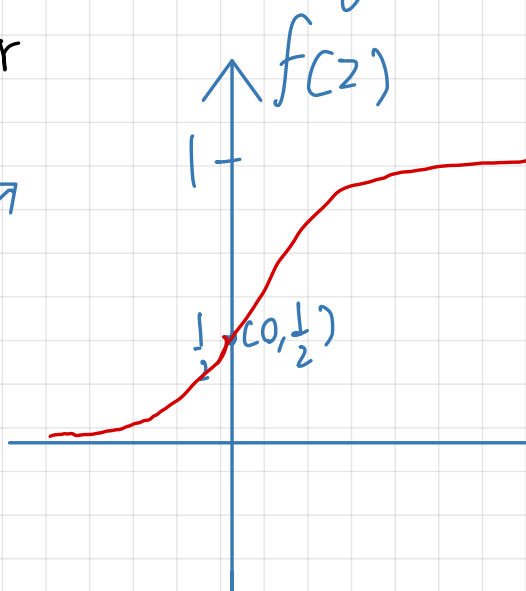
$$P(y|x) = \frac{1}{1 + e^{-y(\vec{w}^T \vec{x})}} \quad ; y \in \{-1, +1\}$$

sigmoid function ♥

Define  $z = w^T x$ ;  $z$  is scalar

$$f(z, y=+1) = \frac{1}{1 + e^{-z}}$$

$$f(z, y=-1) = \frac{1}{1 + e^z}$$



$$\text{if } (f(\infty) = \frac{1}{1 + e^{-\infty}} = 1)$$

$$\text{if } (f(-\infty) = 0)$$

$0 \leq f(z) \leq 1$   
same with  
probability

Recall #1 missclassification occurs when  $y(w^T x) < 0$   $P(y|x) < \frac{1}{2}$   
correct classification  $y(w^T x) = yz \geq 0$   $P(y|x) \geq \frac{1}{2}$

Recall #2 In correct classification (s.t.  $yz \geq 0$ )

$w^T x$  measures the distance from  $x$  to the hyperplane. and

$x$  is very far from the hyperplane, then  $w^T x$  will be large quantity (s.t.  $y=+1$ )

observations:

- If  $x$  lies on the right side of the hyperplane and very far from the hyperplane, then  $P(y|x) = 1$
- If  $x$  lies on the wrong side of the hyperplane and  $x$  is very far from the hyperplane, the  $P(y|x) =$