# Introduction    lecture 13

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

$(x_i, y_i) \sim P(X,Y)$   ②   $P(x,y;\theta) \leftarrow$ learn $- \forall_x \forall_y \ P(X=x \wedge Y=y)$

Classifier $h(x_{test}) = \arg\max\limits_y P(Y=y \mid X=x_{test}; \theta)$

↑ Bayes

assume it can access

$\leftarrow$ Need to learn $\boxed{P(Y \mid X)}$   **

① optimal classifier $h(x_t) = \arg\max\limits_y P(Y=y \mid X=x_t)$ ; pobs ที่ $X=x_t$ แล้ว $Y$ เป็นของอะไร (เช่น salmon, mackael)

pobs ของ salmon เมื่อ $X=x_{test}$
pobs ของ mackel เมื่อ $X=x_{test}$

Note: ถ้ารู้ $(x_i, y_i)$ ใน probability distribution $X, Y$ $(P(X,Y))$ ถ้าจุด๑ ทำนาย $X_{test}$ ต้องหาจาก $P(X,Y)$ ซึ่งไม่ สามารถเข้าถึงได้ จึง ต้องสร้าง modeling distribution ง่าย $P(X,Y;\theta)$

③ $P(Y \mid X)$
$\hookrightarrow P(Y=y \mid X=\vec{x})$

r.v ที่อยู่ใน d=1

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

$P(Y=y \mid X=\vec{x}) = P(Y=y \mid [X]_1 = x_1, [X]_2 = x_2, \ldots, [X]_d = x_d)$   <u>way 1</u>

<u>Bayes rule</u> : $P(Y=y \mid X=x) = \dfrac{P(X=x \mid Y=y) \cdot P(Y=y)}{P(X=x)}$   <u>way 2</u> $\rightarrow$ Navie Bayes use.

estimate $P(X=x \mid Y=y)$
$\qquad : P([X]_1 = x_1, [X]_2 = x_2, \ldots, [X]_d = x_d \mid Y=y)$

# Navie Bayes Assumption    lecture 14

: Assumes all feature values are independent give the label.

$P(X=x \mid Y=y) : \prod\limits_{i=1}^{d} P([X]_i = x_i \mid Y=y)$

- เกิดจากการเรียก Bayes Classifier

# Bayes Classifier

$h(x) = \arg\max\limits_y P_\theta(Y=y \mid X=x)$

<u>Goal</u>: Estimate $P(Y \mid X)$ estimate $\forall x \forall y \ P(Y=y \mid X=x)$

<u>Chain Rule</u>: $P(Y=y \mid X=x) = \dfrac{P(Y=y \wedge X=x)}{P(X=x)}$

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

②   Estimate $P(X=x) \ \forall x$    $P(X)$

Sinario use   $P_\theta(Y=y \wedge X=x) = Bin(n, \theta)$ ??

Coin tossing : n times

$n_H$ : จำนวนครั้งที่โยนแล้วได้ head.

$n_H \sim Bin(n,\theta)$

$MLE = \theta = \dfrac{n_H}{n}$
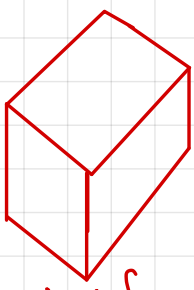
Dice $n_3$ : number of times that we obtian face 3

$X$ คือเป็นไปได้ $X \in \{1,2,\ldots,i\}$

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

i×j faces

→ โยน ลูกเต๋า
หน้าออกเป็น y,x ใช้ Y, X

$I(E)$ : indicator random variable

Assign ตัวแปรสุ่มให้มีค่า 1, 0

$I(E) = \begin{cases} 1 & \text{if E occor} \\ 0 & \text{other wise} \end{cases}$

$MLE \quad \theta = \dfrac{n_{y,x}}{n} = \dfrac{\sum\limits_{i=1}^{n} I(Y_i=y \wedge X_i=x)}{n}$ Ⓘ

ข้อกำหนด { ↓ correspionse

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

$P_\gamma(X=x) = Bin(n,\gamma)$

$MLE = \gamma = \sum\limits_{i=1}^{n} \dfrac{I(X_i=X)}{n}$
ⒾⒾ

Assume X,Y follow Binominal distribution

$MLE: P_\theta(Y=y \mid X=x) = \dfrac{P_\theta(Y=y \wedge X=x)}{P_\gamma(X=x)}$

$= \dfrac{\dfrac{\sum\limits_{i=1}^{n} I(Y_i=y \wedge X_i=X)}{\cancel{n}}}{\dfrac{\sum\limits_{i=1}^{n} I(X_i=X)}{\cancel{n}}}$

$= \dfrac{\sum\limits_{i=1}^{n} I(Y_i=y \wedge X_i=X)}{\sum\limits_{i=1}^{n} I(X_i=x)}$

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

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

- $X$ เด็กที่หลายdimension $\Rightarrow$ Vector ($X$ ข้อมูลกลายเป็น Vector)

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

$$\Downarrow$$

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

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

r.v. X = fixed X on each d

r.v.

---

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

---

# Apply Bayes' rule to Bayes Classifer $(P_\theta(Y=y \mid X=x))$ Way:2

Bayes' rule: $P(Y=y \mid X=x) = \dfrac{P(X=x \mid Y=y) \, P(Y=y)}{P(X=x)}$

Bayes Classifier

$$= h(X) = \underset{y}{\arg\max} \; P(Y=y \mid X=x)$$

$$= \underset{y}{\arg\max} \; \frac{P(X=x \mid Y=y) \cdot P(Y=y)}{\cancel{P(X=x)}} \text{ ไม่เกี่ยวกับ } y$$

$$= \underset{y}{\arg\max} \; P(X=x \mid Y=y) \cdot P(Y=y)$$

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

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

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

$$= P([\vec{X}]_1 = \vec{x}^1 \wedge [\vec{X}]_2 = \vec{x}_2 \wedge \ldots \wedge [\vec{X}]_d = \vec{x}_d \mid Y=y)$$ (can't estimate direct)

So use Naive Bayes assumption.

<u>Assume $X, Y$ follow Binomial distributions</u>

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

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

# Naive Bayes assumption:

- All feature values are Independent

$$P(\vec{X}=\vec{x} \mid Y=y) = \prod_{i=1}^{d} P([X]_i = x^i \mid Y=y) \quad ; \text{Prob} \text{ นี้ } Y \text{ given } y \text{ ค่า จ:ได้ } [X]_i = x^i \text{ จ:อิสระ:ต่อกัน}$$

<div style="border:1px solid red; padding:8px;">

Naive Bayes classifier

$$h(X) = \underset{y}{\text{argmax}} \prod_{\alpha=1}^{d} P([\vec{X}]_\alpha = \vec{x}^\alpha \mid Y=y) \, P(Y=y)$$

</div>

# <u>How to estimate $P([\vec{X}]_\alpha \mid Y)$</u> 3 cases

- There are 3 notable cases:

<u>Case1</u> : Categorical features : Categorical Naive Bayes Classifier

$$[\vec{x}]_\alpha \in \{c_1, c_2, \ldots, c_k\}$$

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

We model $P([\vec{X}]_\alpha = j \mid Y=y) = [\theta_{jy}]_\alpha$ parameter จำแพจะของ $j$ y ที่ codinate $\alpha$

$\uparrow$

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

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

$\Rightarrow \dfrac{\sum_{i=1}^{n} I(Y_i = y) \cdot I(x_i^\alpha = j)}{\sum_{i=1}^{n} I(Y_i = y)}$ : count เมื่อ เจอ $Y_i = y$ และ $X_i = j$ ที่ $\alpha$

# 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)

*Calagorical feature* $\{$

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

$\Rightarrow P([\vec{x}]_1 = sunny \mid y = go\text{-}out)$
$\rightarrow P([\vec{x}]_1 = rainy \mid y = go\text{-}out)$

Estimate $P([\vec{x}]_1 \mid Y)$

$\rightarrow P([\vec{x}]_1 = sunny \mid y = stay)$
$\rightarrow P([\vec{x}]_1 = rainy \mid y = stay)$

$P([\vec{x}]_1 = sunny \mid y = go\text{-}out) = [\theta_{sunny, goout}]_1 = \dfrac{4}{5} = 0.8$

| $i$ | $x_i^1$ | $x_i^2$ | $y_i$ |
|-----|---------|---------|-------|
| 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) 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))?

<u>Case 2</u>: Multinomial feature: Multinomial Naive Bayes Classifier

eg. text data: "An ant is animal." 

Back of word

$\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_j \\ \vdots \\ x_d \end{bmatrix}$ ← # of word i appearing in the text

: its count of word in Sentence.

$\int$ cordinate

- We estimate $P([\vec{x}]_\alpha = j \mid Y = y)$ by using multinomial distribution

eg. Spam filter

- $y \in \{$ spam, ham $\}$
       $+1$    $-1$

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

$\vec{x} = \begin{bmatrix} x^1 \\ \vdots \\ i \\ x^d \end{bmatrix} \begin{matrix} - w_1 \\ - w_2 \\ \\ - w_d \end{matrix}$

Estimate $P([\vec{x}]_\alpha = j \mid Y = y)$ ; $\alpha : 5$, $j = 10$, $y = $ spam  Given $W_5 = $ Princess (เกิด 10ครั้ง)

$P([\vec{x}]_5 = 10 \mid Y = $ spAM$)$; given spam ที่ตัวแปรสุ่ม cordinate 5 = 10

<u>Modeling</u>

$$\boxed{P([\vec{x}]_\alpha = j \mid Y = y) = \left[ \binom{m}{j} \cdot P(W_\alpha \mid Y = y) \right]^j}$$

eg. $P([\vec{x}]_2 = 3 \mid Y = $ spam$) = \left[ \binom{5}{3} \times P(W_2 \mid Y = $ spam$) \right]^3$
           $\downarrow$
        word ที่ 2

- $P(W_2 \mid Y = y)$ is the prob. of selecting word $W_\alpha$ given the label is y.

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

# Multinomial distribution



d faces
(y = spam)

$$\bar{x} = \begin{pmatrix} 1 & x^1 \\ 4 & x^2 \\ \vdots & x^d \end{pmatrix} \begin{array}{l} - W_1 \\ - W_2 \\ - W_d \end{array}$$

### Assume throw and get R1:

$$\underbrace{W_2, W_2, W_1, W_2, W_2, W_3, W_2}_{m \ times = 7}$$

R2:
$$\underbrace{W_1, W_2, W_2, W_3, W_2, W_2, W_2}_{n \ times = 7}$$

$R_1, R_2$ not a same words but with power of Back of words it give same feature Vector. (order of word is useless)

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

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

- Estimate

$\underline{[\theta_y]_\alpha \ \forall_y \forall_\alpha}$

- MIE : $[\theta \ spam]_\alpha = \dfrac{\sum\limits_{i=1}^{n} I(Y_i = SPAM) \cdot x_i^\alpha}{\sum\limits_{i=1}^{n} I(Y_i = SPAM) \cdot (\sum\limits_{\alpha=1}^{d} x_i^d)}$

$= \dfrac{\text{จำนวนครั้งที่ } w_\alpha \ \text{ปรากฎในทุก n Sample ที่เป็น spam}}{\text{จำนวนของ words รวมใน ทุก n Sample ที่เป็น spam}}$

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

$P([\overset{3}{\underset{money}{\bar{x}}}] = 2 \ | \ Y = \text{spam}) = \binom{m}{2} \cdot [P(\text{money} | y = \text{spam})]^2$

∴ โทร ต้องเน้น money 2 ครั้ง
เราไม่รู้อยู่ ตรงไหนบ้าง เลยใช้ $\binom{m}{2}$ ด้วย

---

# Summary

Bayes Rule → $h(x) = \underset{y}{\arg\max} \ P(X=x | Y=y) \cdot P(Y)$

Naive →
Bayes
Classifier

$\qquad = \underset{y}{\arg\max} \ \prod\limits_{\alpha=1}^{d} P([X]_\alpha = x^\alpha | Y=y) \cdot P(Y)$

Spam filter (text classification)

"An ant is an animal" ; 5 words

Bag of word $\searrow \tilde{x} = \begin{bmatrix} 0 \\ 2 \\ \vdots \end{bmatrix} \quad \begin{matrix} W_1 = a \\ W_2 = an \\ \vdots \\ W_d = \end{matrix}$

eg: $P([X]_2 = 2 \mid y = spam) = \binom{5}{2}\left[ P(W_2 \mid y = spam)\right]$

## Estimate

$$P(W_2 \mid y = spam) = [\Theta_{spam}]_2$$

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

## Ingeneral  $\cdots$

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

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

$\rightarrow$ จน. ครั้งแต่ละ word ปรากฏอยู่ใน y

$\rightarrow$ จน. ของ word ทั้งหมดที่อยู่ใน y

$$h(X) = \arg\max_Y \boxed{\prod_{\alpha=1}^{d} P([X]_\alpha = X^\alpha \mid y = y)} \cdot P(y)$$

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

$$= P([X]_1 = x^1 \mid y = y) \times \dots \times P([X]_\alpha = x^\alpha \mid y = y)$$

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

$$\times \dots \times \left( \binom{m - x_1 - \dots - x^{d-1}}{x^d} \left( [\theta_y]_d \right)^{x^d} \right)$$

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

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

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

$$0!$$

$$\prod_{\alpha=1}^{d} P([X]_\alpha = x^\alpha \mid y=y) = \frac{m!}{x^1! x^2! \dots x^\alpha!} \prod_{\alpha=1}^{d} \left([\theta_y]_\alpha\right)^{x^\alpha}$$

<u>Binary Classification.</u>

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

<u>take $log_e$ into both sides</u>

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

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

<u>$P([X]_\alpha \mid y)$</u>

3 notable class

① Categorical features ($[X]_\alpha$ เป็นประเภท(แทน word))ไม่ใช่ จำนวนครั้ง
  $[X]_\alpha \in \{0,1,\dots,k-1\}$  ของการปรากฏ) Catagorical Naive Bayes
  $\llcorner\rightarrow$ Classifier

② multinomial feature ($[X]_\alpha$ แทนความถี่ ของการปรากฏ ของ แต่ละ word)
  $[X]_\alpha \in \{0,1,\dots,m\}$  $m = \sum_{\alpha=1}^{d} x^\alpha$ เป็นความถี่ (จน.ครั้ง ของ การปรากฏ)
  $\llcorner\rightarrow$ multinomial Naive Bayes
  Classifier

③ Continous feature
  $[X]_\alpha \in \mathbb{R}$  ($[X]_\alpha$ เป็นค่าต่อเนื่อง)
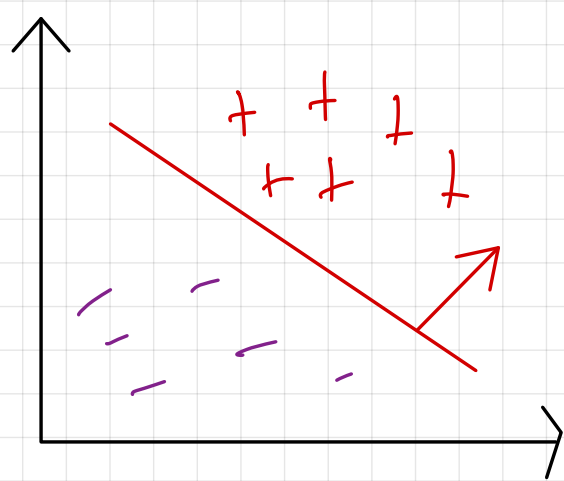  $\llcorner\rightarrow$ Guassian Naive Bayes Classifier

# Summary of Naive Bayes

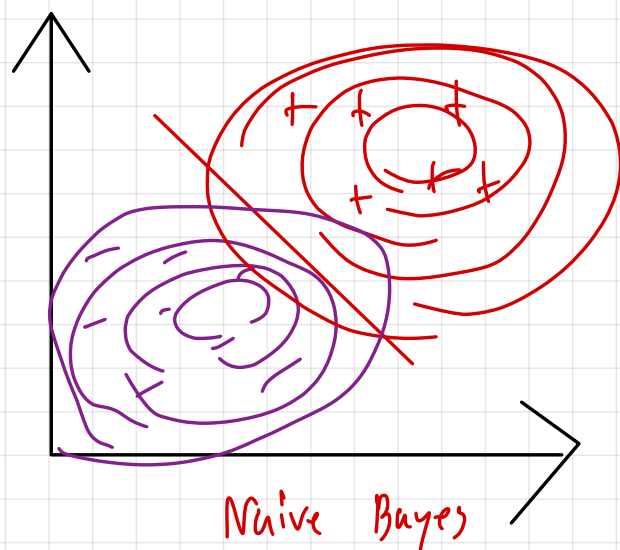- Naive Bayes = Bayes classifier + Bayes rule + Naive Bayes assumption

- The assumption says "all feature values are Indepedent."

  - พิจารณา $W_\alpha$ เมื่อให้ label $y = spam$  $P(W_\alpha | y = spam)$

- 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



Percepton

- hyperplane separate ตาม data

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

Naive Bayes
- สร้าง distribution vอง probs. ใน sample

- separate ตาม distribution

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

ก้าหนด $y$ โดยใช้ $x$

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 : $P(y|x) = \dfrac{P(x|y) \times P(y)}{P(x)}$

# Prove multinomial Naive Bayes is a linear Classifier

Proof - Assume $y \in \{-1, +1\}$    $T \to T = T$
$T \to F = F$

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

ด้านมากว่า

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

iff $\prod\limits_{\alpha=1}^{d} P(X_\alpha|y=+1) \times P(y=+1) \geq \prod\limits_{\alpha=1}^{d} P(X_\alpha|y=-1) \times P(y=-1)$

iff $\sum\limits_{\alpha=1}^{d} \log_e P(X_\alpha|y=+1) + \log_e P(y=+1) \geq \sum\limits_{\alpha=1}^{d} \log_e P(X_\alpha|y=-1)$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad + \log_e P(y=-1)$

$([\theta_{+1}]_\alpha)^{x_\alpha}$

iff $\sum\limits_{\alpha=1}^{d} (\log_e \underline{P(X_\alpha|y=+1)} - \log_e P(X_\alpha|y=-1))$

$\qquad + (\log_e P(y=+1) - \log_e P(y=-1)) \qquad \geq 0$

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

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

$b$

$\therefore$ เหมือนรูปการ $\vec{w}^T \vec{x} + b$ ในการหา hyper plane ของ perceptron

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

iff $\vec{w}^T \vec{x} + b \geq 0$

<u>Case 3</u> : Continuous feature: : Guassian Naive Bayes Classifier

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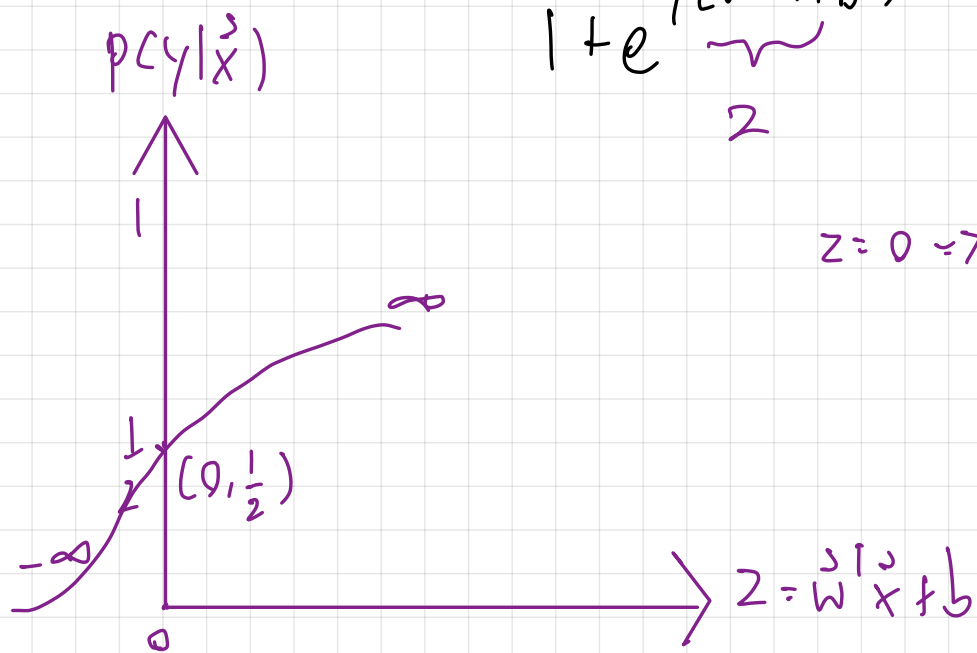
- Model $P([X]_\alpha = j | Y = y)$

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

- MLE $\rightarrow [\mu_y] = \dfrac{\sum\limits_{i=1}^{n} \mathbb{1}(Y_i = y) \cdot x_i^\alpha}{n}$   Same with $\dfrac{\sum x}{n}$

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

$$P(y|\vec{x}) = \dfrac{1}{1 + e^{\underbrace{-y(\vec{w}^T\vec{x} + b)}_{z}}} \quad , \quad \text{for } y \in \{-1, +1\}$$

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



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

$\left(0, \frac{1}{2}\right)$

$z = \vec{w}^T\vec{x} + b$

<u>Recall before Prove</u>

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

<span style="color:red">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 : $P(y|x) = \dfrac{P(x|y) \times P(y)}{P(x)}$</span>

-Discriminative learning : Try to model $P(y|x)$ directly
- Generative learning : Try to model $P(x|y)$ and $P(y)$

eg. Perceptron is a discrimative algorithm

$P(y|x) = \begin{cases} \text{① if } w^Tx \geq 0 \\ \quad \hookrightarrow \text{pobability} \\ 0 \quad \text{otherwise} \end{cases}$

เมื่อ $y = +1$

∴ มั่นใจได้ว่า โอกาส คือ 1 ถ้าไม่ใช่ก็เป็น 0

eg. Naive Bayes is a generative algorithm

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

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

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

$\therefore$ ถ้า $b$ ไม่ใช่ $w \Rightarrow w^T x \geq 0$

eg. Perceptron, Multinomial Naïve Bayes
are linear classifiers

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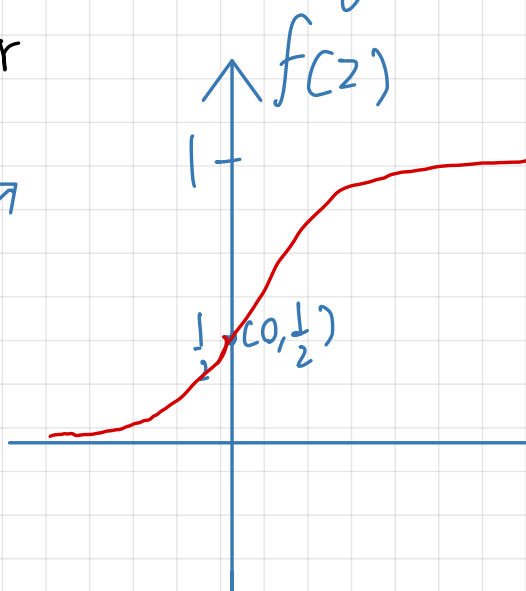expression for Guassian Naive Bayes

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

sigmond function ♡

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

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

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

$f(z)$

$1$

$(0, \tfrac{1}{2})$

if $(f(\infty) = \dfrac{1}{1 + z^{-\infty}} = 1)$

if $(f(-\infty) = 0\ )$

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

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

<u>Recall #2</u> In correct classification (กรณี $yz \geq 0$)
$\quad w^T x$ measures the distance from $x$ to the hyperplane. and
$\quad\quad x$ is very far from the hyperplane, then $w^T x$ will be large quantity
ค่ามาก (กรณี $y=+1$)

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 $P(y|x) =$