

Classification = Probabilistic Generative Model

- 找一个 function，它的 Input 是一个 object X ，
Output 是这个 object 属於哪个 class

for example:

1 金融学：要不要貸款給某個人

Input: 收入, 工作, 年紀, 有無欠債

Output: 要不要借錢

2. 醫療的病診斷

Input: 某一個人的症狀, 年紀, 性別, 救醫歷史

Output: 生的是哪一種病

3. 手寫文字辨識

Input: 手寫一個字

Output: 辨認出是什麼

4. 人臉辨識

Input: 人臉

Output: Who is the person

· 一隻 pokemon 是由 7个数字 所組成的一個 vector

輸入: 一个 function

Train data 400 ↓

Testing data 400 ↑

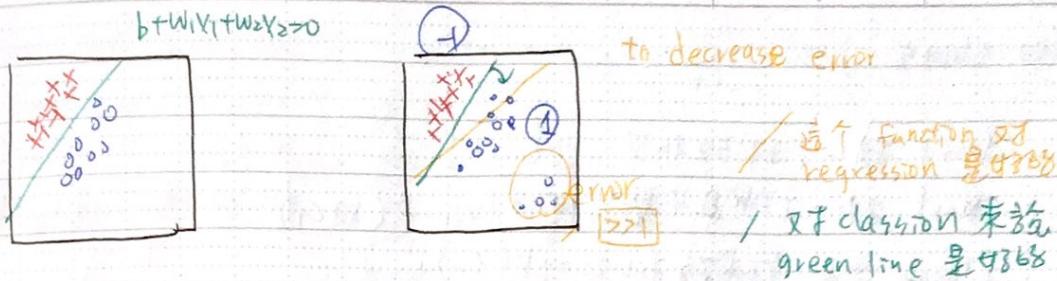
⇒ Classification as Regression

[binary 單二 example.)

Training $x \leftarrow$ class 1 ; Let class 1 = 1
class 2 ? class 2 = -1

Testing = Close to 1 → class 1

-1 → class 2.



* Regression 定義 function 去土壤的方式 对 classification 來說不適用

* Multiple class

if now have class 1, class 2, class 3.

如果把這個當作 Regression 的問題來處理
沒有辦法得到好的結果

理想的方法

→ classification 時 output 是離散的

Model $X \rightarrow f(X)$ [class 1 if $f(X) > 0$
class 2 o.w.]

Loss function

$$\Delta(f) = \sum_n \delta(f(x^n) \neq \hat{y}^n) \quad // \text{Index function 的概念}$$

$$\text{e.g. } f(x^n) \neq \hat{y}^n \Rightarrow 1$$

$$f(x^n) = \hat{y}^n \Rightarrow 0$$

Find the best function

example = Perceptron, SVM

Two class.

Class1 抽一个球的机率 = $P(C_1)$

Class1 抽一个球且它是blue的机率 $P(X|C_1)$

Class2 抽一个球的机率 = $P(C_2)$

Class2 抽一个且它是blue的机率 $P(X|C_2)$

\Rightarrow 通过 training data
求出



Generative Model.

某个 X 出现的机率

$$P(X) = P(C_1) P(X|C_1) + P(C_2) P(X|C_2)$$

\Rightarrow Can自己產生 X

Prior

water & normal type with ID $< 400 \rightarrow$ training, $> 400 \rightarrow$ test

Train 中 水 = 79

normal = 61

$$P(C_1) = \frac{79}{79+61}$$

$$P(C_2) = \frac{61}{79+61}$$

* 每个pokemon 都是用一堆 feature 来描述

Gaussian Distribution. $f_{X|\mu} = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right]$

Input = Vector X

Output = 机率密度 \propto 机率

function 的形状由 μ & Σ 决定

\downarrow
mean

covariance matrix

同样的 Σ , 不同的 μ

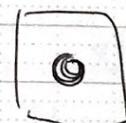
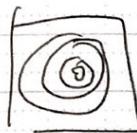
\Rightarrow 机率密度 (最高点的地方) 是 一样的



同样的 μ , 不同的 Σ

\Rightarrow 机率密度 (分布的最高点) 是 一样的

but 分散的程度 not same

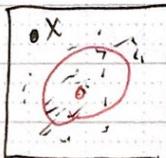


Probability from class.

- Assume point sampled from Gaussian Distribution.

- Find the Probability for new point X

$$\text{now know } \mu = \begin{bmatrix} 15 \\ 13 \end{bmatrix}, \Sigma = \begin{bmatrix} 814 & 321 \\ 321 & 929 \end{bmatrix}$$



If x 越靠近中心点 (μ) = sample 的机率越大

大

how to find μ and Σ

\Rightarrow Maximum Likelihood

#Maximum likelihood

* The Gaussian 在任何的一組 μ and Σ 都可以生成所需的某
 \Rightarrow 只是有些地方機率很高，有些地方機率很低
 But 沒有一個地方的機率是 0

\Rightarrow 其 likelihood 不一樣

\Rightarrow likelihood of a Gaussian with μ and Σ

= the probability of the Gaussian samples. x^1, \dots, x^{19}

$$L(\mu, \Sigma) = f_{\mu, \Sigma}(x^1) f_{\mu, \Sigma}(x^2) \dots f_{\mu, \Sigma}(x^{19})$$

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{\frac{n}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu))$$

$$\mu^*, \Sigma^* = \arg \max_{\mu, \Sigma} L(\mu, \Sigma)$$

$$\mu^* = \frac{1}{19} \sum_{n=1}^{19} x^n \quad (\text{or 把 } \mu^* \text{ 設為 } 0 \text{ 的解})$$

$$\Sigma^* = \frac{1}{19} \sum_{n=1}^{19} (x^n - \mu^*) (x^n - \mu^*)^T \quad (\text{或 } f \text{ 又 } \Sigma^* \text{ 微分，再 } k \text{ 跑})$$

#Now can do classification.

$$P(c_1|x) = \frac{P(x|c_1) P(c_1)}{P(x|c_1) P(c_1) + P(x|c_2) P(c_2)}$$

$$w) P(x|c_1) = f_{\mu^1, \Sigma^1}(x)$$

$$\mu^1 = \begin{bmatrix} 75 \\ 71.3 \end{bmatrix}, \Sigma^1 = \begin{bmatrix} 814 & 321 \\ 321 & 929 \end{bmatrix}$$

$$P(x|c_2) = f_{\mu^2, \Sigma^2}(x)$$

$$\mu^2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix}, \Sigma^2 = \begin{bmatrix} 841 & 422 \\ 422 & 685 \end{bmatrix}$$

$$P(c_1) = 0.56, P(c_2) = 0.44$$

$$\text{if } P(c_1|x) > 0.5 \Rightarrow x \in \text{class } (c_1)$$

A 答案 = If we want to do Pokemon 屬性 分類。

* 假設一個 機率 model < required probability

每個 class 自己的 distribution

每個 class 自己的 機率，Use Gaussian

比較常見的作法：不同的 class 可以 share 同一個 covariance matrix
 \Rightarrow less 傳數 \Rightarrow 不太會 overfitting

If now covariance matrix & feature size 很大時
 \Rightarrow 它的增長很快

If now have different Σ
 \Rightarrow model 傳數太多。
 \Rightarrow 可能 overfitting

Modifying Model (by Bishop ch 4.2.2)

- Maximum Likelihood

$$\text{water} = (\mu^1, \Sigma)$$

$$\text{normal} = (\mu^2, \Sigma)$$

$$L(\mu^1, \mu^2, \Sigma) = f_{\mu^1, \Sigma}(x^1) \times f_{\mu^2, \Sigma}(x^2) \times f_{\mu^3, \Sigma}(x^3) \times \dots \times f_{\mu^4, \Sigma}(x^4)$$

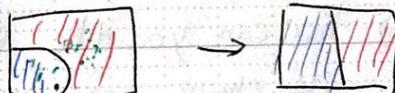
- If μ^1, μ^2 same, Σ^1, Σ^2 not same

$$\Rightarrow \Sigma = \frac{1}{40} \Sigma^1 + \frac{6}{40} \Sigma^2$$

If now 訊同 covariance

\Rightarrow The boundary is linear.

(原本是曲線) Gaussian



All = total, hp, att, dev, sp, de, speed

54% accuracy \rightarrow 73% accuracy

Three Step

回顧機率的 model, we know that machine learning is 3 step

now II 機率 model 也是 3 step

1. Function set (Model)

$X \rightarrow$

$$f(x) = \frac{P(c_1|x) P(c_1)}{P(c_1|x) P(c_1) + P(c_2|x) P(c_2)}$$

$$f(x) = \begin{cases} \text{class 1} & \text{if } P(c_1|x) > 0.5 \\ \text{class 2} & \text{o.w.} \end{cases}$$

choose different Σ, μ

probability distribution

required probability

\Rightarrow 3 model

各自不同

2. Goodness of a function

\Rightarrow evaluate function set by function 好壞

\Rightarrow 得到 probability distribution, 可以最大化這個 data by likelihood
 $\text{i.e. } \mu^*, \Sigma^* = \arg \max_{\mu, \Sigma} L(\mu, \Sigma)$

3. Find the best function

Probability distribution

\Rightarrow You can use the distribution you like

If choose easy probability model = 參數少, bias 大, variance 大
 complex = 大, 小 大

\Rightarrow Use data set 来決定 which is better

$$P(X|C_1) = P(x_1|c_1) P(x_2|c_1) \cdots P(x_k|c_1) \cdots$$

by feature 計算

$[x_1, \dots, x_k]$

and x_1, \dots, x_k is independent

(也從機率 model 產生出來的機率是 Independent)

1 維的 Gaussian

\Rightarrow covariance matrix \Rightarrow diagonal

\Rightarrow 不是對角線的地方 = 0

\Rightarrow 可以 decrease 共享數量

\Rightarrow can get easy mode

\Rightarrow 只有對角線的地方有值

\Rightarrow still need covariance

for binary feature (是 or 不是)

⇒ Use Bernoulli distribution

If now assume all dimension are independent

⇒ 不 model feature & feature PB covariance 的關係

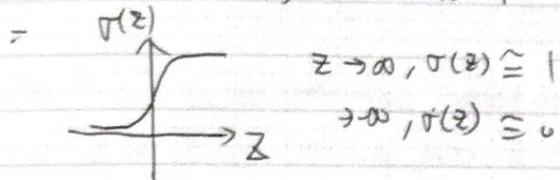
⇒ Using Naive Bayes classifier

(3 項不強 depend on 假設是不是精準的)

Posterior Probability

$$\begin{aligned} P(C_1|X) &= \frac{P(X|C_1)P(C_1)}{P(X|C_1)P(C_1) + P(X|C_2)P(C_2)} \\ &= \frac{1}{1 + \frac{P(X|C_2)P(C_2)}{P(X|C_1)P(C_1)}} \\ &= \frac{1}{1 + \exp(-z)} \quad \text{w) } Z = \ln \frac{P(X|C_1)P(C_1)}{P(X|C_2)P(C_2)} \\ &= \sigma(z) \end{aligned}$$

sigmoid function



$$\begin{aligned} Z &= \ln \frac{P(X|C_1)P(C_1)}{P(X|C_2)P(C_2)} \\ &= \ln \frac{P(X|C_1)}{P(X|C_2)} + \ln \frac{P(C_1)}{P(C_2)} \xrightarrow{\frac{N_1}{N_1+N_2}} \frac{\frac{N_1}{N_1+N_2}}{\frac{N_2}{N_1+N_2}} = \frac{N_1}{N_2} \\ &= \ln \left(\frac{|\Sigma^1|^{\frac{1}{2}}}{|\Sigma^2|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} [(\mathbf{x}-\mathbf{m}^1)^T (\Sigma^1)^{-1} (\mathbf{x}-\mathbf{m}^1) - (\mathbf{x}-\mathbf{m}^2)^T (\Sigma^2)^{-1} (\mathbf{x}-\mathbf{m}^2)] \right\} \right) \\ &\quad + \ln \frac{N_1}{N_2} \\ &= \ln \frac{|\Sigma^1|^{\frac{1}{2}}}{|\Sigma^2|^{\frac{1}{2}}} - \frac{1}{2} \mathbf{x}^T (\Sigma^1)^{-1} \mathbf{x} + (\mathbf{m}^1)^T (\Sigma^1)^{-1} \mathbf{x} - \frac{1}{2} (\mathbf{m}^1)^T (\Sigma^1)^{-1} \mathbf{m}^1 \\ &\quad + \frac{1}{2} \mathbf{x}^T (\Sigma^2)^{-1} \mathbf{x} - (\mathbf{m}^2)^T (\Sigma^2)^{-1} \mathbf{x} + \frac{1}{2} (\mathbf{m}^2)^T (\Sigma^2)^{-1} \mathbf{m}^2 + \ln \frac{N_1}{N_2} \\ \Sigma &= \Sigma_1 \cup \Sigma_2 \\ &= (\mathbf{m}^1 - \mathbf{m}^2)^T \sum_{i=1}^{N_1+N_2} \Sigma_i^{-1} \mathbf{x} - \frac{1}{2} (\mathbf{m}^1)^T (\Sigma^1)^{-1} + \frac{1}{2} (\mathbf{m}^2)^T (\Sigma^2)^{-1} \mathbf{m}^2 + \ln \frac{N_1}{N_2} \end{aligned}$$

$$P(C_1|X) = \sigma(z)$$

$$= \sigma(wX + b)$$

\Rightarrow We know that $\Sigma_1 = \Sigma_2 = \Sigma$ (share Σ)
the boundary is linear

\Rightarrow In generative model, we estimate $N_1, N_2, M^1, M^2, \Sigma$
Then we have w and b

\Rightarrow 直接找 w 和 b 找出来

\Rightarrow Logistic Regression

Logistic Regression

\rightarrow w & b 決定的 model

	Logistic Regression.	Linear Regression
Step 1	$f_{w,b}(x) = \sigma(\sum w_i x_i + b)$ output = between 0 & 1	$f_{w,b}(x) = \sum w_i x_i + b$ output = any value
Step 2.	Training data: (x^n, \hat{y}^n) $\hat{y}^n = 1$ for class 1 $= 0$ for class 2 $L(f) = \sum C(f(x^n), \hat{y}^n)$	Training data: (x^n, \hat{y}^n) $\hat{y}^n =$ a real number $L(f) = \frac{1}{2} \sum (f(x^n) - \hat{y}^n)^2$
Step 3:	$W_{\bar{x}} = W_{\bar{x}} - \eta \sum (\hat{y}^n - f_{w,b}(x^n)) \bar{x}$	same

\Rightarrow Cross entropy $C(f(x^n), \hat{y}^n) = -[\hat{y}^n \ln(f(x^n)) + (1-\hat{y}^n) \ln(1-f(x^n))]$

\downarrow is a Bernoulli distribution
 \Rightarrow want 這個 distribution 越接近越好.

Step 2:
 Training data = $\begin{bmatrix} x^1 & x^2 & x^3 & \dots & x^N \\ c_1 & c_1 & c_2 & \dots & c_1 \end{bmatrix} \rightarrow \begin{bmatrix} x^1 & x^2 & x^3 \\ \hat{y}^1 = 1 & \hat{y}^2 = 1 & \hat{y}^3 = 0 \\ \uparrow \text{Posterior probability} \end{bmatrix}$

Assume the data is generated based on $f_{w,b}(x) = P_{w,b}(C_i|x)$

(i.e. 這組 training data 是由 $P_{w,b}(C_i|x)$ 產生的)

now given w and b , what is the probability?

$$\Rightarrow L(w, b) = f_{w,b}(x^1) f_{w,b}(x^2) (1 - f_{w,b}(x^3)) \dots f_{w,b}(x^N).$$

most likely w^*, b^* (最大化這個機率 w^*, b^*)

$$w^*, b^* = \arg \max_{w, b} L(w, b) = \arg \min_{w, b} -\ln L(w, b)$$

$$-\ln L(w, b)$$

$$= -\ln f_{w,b}(x^1) - \ln f_{w,b}(x^2) - \ln(1 - f_{w,b}(x^3)) - \dots$$

$$= -[\hat{y}^1 \ln f_{w,b}(x^1) + (1 - \hat{y}^1) \ln(1 - f_{w,b}(x^1))]$$

$$- [\hat{y}^2 \ln f_{w,b}(x^2) + (1 - \hat{y}^2) \ln(1 - f_{w,b}(x^2))]$$

$$- [\hat{y}^3 \ln f_{w,b}(x^3) + (1 - \hat{y}^3) \ln(1 - f_{w,b}(x^3))]$$

$$= \sum_n -[\hat{y}^n \ln f_{w,b}(x^n) + (1 - \hat{y}^n) \ln(1 - f_{w,b}(x^n))]$$

\downarrow Bernoulli distribution to cross entropy

z distribution

Distribution P

$$\begin{aligned} P(x=1) &= \hat{y}_n \\ P(x=0) &= 1 - \hat{y}_n \end{aligned}$$

Distribution q:

$$q(x=1) = f(x^n)$$

$$q(x=0) = 1 - f(x^n)$$

$$H(P, q) = - \sum_x p(x) \ln(q(x))$$

* Cross entropy: 這兩個 distribution 有多接近

if 两个 distribution same 其 cross entropy = 0

$$\Rightarrow \hat{y}_n * \ln(f(x^n)) + (1 - \hat{y}_n) * \ln(1 - f(x^n))$$

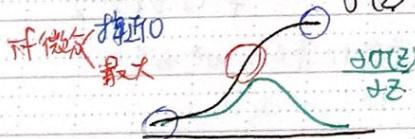
is the cross entropy

Step 3: Find the best function

$$-\frac{\partial \ln L(w, b)}{\partial w_i} = \sum_n \left[\hat{y}_n \frac{\partial \ln f_{w,b}(x^n)}{\partial w_i} + (1 - \hat{y}_n) \frac{\partial \ln(1 - f_{w,b}(x^n))}{\partial w_i} \right]$$

$$\begin{aligned} \frac{\partial \ln f_{w,b}(x)}{\partial w_i} &= \frac{\partial \ln f_{w,b}(x)}{\partial z} \frac{\partial z}{\partial w_i} & w_i \cdot f_{w,b}(x) = \sigma(z) = \frac{1}{1 + \exp(-z)} \\ &= (1 - \sigma(z)) X_i & z = w \cdot x + b \\ &\frac{\partial z}{\partial w_i} = X_i & = \sum_n w_i x_i + b \end{aligned}$$

$$\begin{aligned} \frac{\partial \ln(1 - f_{w,b}(x^n))}{\partial z} &= \frac{\partial \ln \sigma(z)}{\partial z} = \frac{1}{\sigma(z)} \cdot \frac{\partial \sigma(z)}{\partial z} = \frac{1}{\sigma(z)} \cdot \sigma(z)(1 - \sigma(z)) \\ &= 1 - \sigma(z) \end{aligned}$$



$$\begin{aligned} \frac{\partial \ln(1 - f_{w,b}(x^n))}{\partial w_i} &= \frac{\partial \ln(1 - f_{w,b}(x^n))}{\partial z} \frac{\partial z}{\partial w_i} \\ &= \frac{\partial \ln(1 - \sigma(z))}{\partial z} X_i \\ &= -\frac{1}{1 - \sigma(z)} \frac{\partial \sigma(z)}{\partial z} X_i \\ &= -\frac{1}{1 - \sigma(z)} \cdot \sigma(z)(1 - \sigma(z)) \cdot X_i \geq \sigma(z) \cdot X_i \end{aligned}$$

$$\begin{aligned}
 \Rightarrow \frac{-\ln(w, b)}{\sum w_i} &= \sum_n -[\hat{y}^n(1-f_{w, b}(x^n))x_i^n + (1-\hat{y}^n)f_{w, b}(x^n)x_i^n] \\
 &= \sum_n -[\hat{y}^n - \hat{y}^n f_{w, b}(x^n) - f_{w, b}(x^n) + \hat{y}^n f_{w, b}(x^n)]x_i^n \\
 &= \sum_n -[\hat{y}^n - f_{w, b}(x^n)]x_i^n
 \end{aligned}$$

If now use gradient descent to update

$$w_i = w_i - \eta \sum_n -(\hat{y}^n - f_{w, b}(x^n))x_i^n$$

w 的 update 取決於 3 件事

① η = learning rate = 自調的

x_i^n = 來自 data

$\hat{y}^n - f_{w, b}(x^n)$ ：現在的 output 和理想的目標差距有多大
目標 現在 model 的 output

Logistic Regression + Square Error.

$$\text{Step 1: } f_{w, b}(x) = \sigma(\sum w_i x_i + b)$$

$$\text{Step 2: Training data } (x^n, \hat{y}^n) \quad \hat{y}^n = \begin{cases} 1 & \text{class 1} \\ 0 & \text{class 2} \end{cases}$$

$$L(f) = \frac{1}{2} \sum_n (f_{w, b}(x^n) - \hat{y}^n)^2$$

$$\text{Step 3: } \frac{\partial (f_{w, b}(x) - \hat{y})^2}{\partial w_i} = 2(f_{w, b}(x) - \hat{y}) \cdot \frac{\partial f_{w, b}(x)}{\partial w_i} \cdot x_i$$

If now $\hat{y}^n = 0$ If $f_{w, b}(x^n) = 1$ far from the target $\rightarrow \frac{\partial L}{\partial w_i} \neq 0$

$\hat{y}^n = 0$ close to 1 $\rightarrow \frac{\partial L}{\partial w_i} = 0$

$$\begin{array}{lll}
 \hat{y}^n = 1 & f(x^n) = 1 & \text{close} \\
 & \approx 0 & \approx 0 \\
 & & \text{far} & \approx 0
 \end{array}$$

[If now is cross entropy, far from 目標，微分值越大 $\frac{\partial L}{\partial w_i}$ update 越快
square error \hat{y}^n 很小 $\frac{\partial L}{\partial w_i}$ 慢]

now we don't know gradient $\frac{\partial L}{\partial w_i}$ 的時候 (做什麼) is close to / far from 目標

Date

Logistic regression 和 discriminative
use Gaussian 来推导 posterior probability 和 Generative

Discriminative vs Generative

實際上這個 model., function set is same

If we f.e. covariance matrix 設成是 share

$$f(W \cdot x + b) = P(C|X)$$

If we use Logistic regression

→ 可以直接找出 w & b (use gradient descent)

→ 沒有对 probability distribution 有任何假設/描述

If now is generative model

→ 會先算 $\mu_1, \mu_2, \Sigma^{-1}$ (Σ 是 inverse)

$$WT = (\mu_1 - \mu_2)^T \Sigma^{-1}$$

$$b = -\frac{1}{2}(\mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2}(\mu_2^T \Sigma^{-1} \mu_2) + \ln \frac{N_1}{N_2})$$

→ 因為 probability distribution 假設它是 gaussian, bernoulli, no Naive Bayes

→ 因為 we use 不同的假設

所以根據同一組 training data 找出的參數會不一樣

(即找到的 w & b is different)

⇒ which is better (w & b)

⇒ discriminative is better

⇒ why?

Example:

每一筆 Data 有兩個 feature, 共 13 筆 data

Training data

$\begin{bmatrix} 1 \\ 0 \end{bmatrix} \times 1$	$\begin{bmatrix} 1 \\ 0 \end{bmatrix} \times 4$	$\begin{bmatrix} 0 \\ 1 \end{bmatrix} \times 4$	$\begin{bmatrix} 0 \\ 0 \end{bmatrix} \times 4$
Class1	Class2	Class2	Class2

Testing data

$\begin{bmatrix} 1 \\ ? \end{bmatrix}$ Class1?
Class2?

⇒ How about Naive Bayes?

$P(X|C_i) = P(X_1|C_i)P(X_2|C_i)$

↳ P of X 從某一個 class 產生出來的機率

$$P(C_1) = \frac{1}{3} \quad P(C_2) = \frac{2}{3}$$

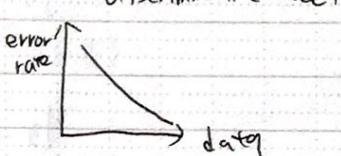
$$\begin{aligned} P(X_1=1|C_1) &= 1 & P(X_2=1|C_1) &= 1 \\ P(X_1=1|C_2) &= \frac{1}{3} & P(X_2=1|C_2) &= \frac{1}{3} \end{aligned}$$

$$\Rightarrow P(C_1|X) = \frac{P(X|C_1)P(C_1)}{P(X|C_1)P(C_1) + P(X|C_2)P(C_2)}$$

$$= \frac{1 \times 1 \times \frac{1}{3}}{\frac{1}{3} + \frac{1}{3} \times \frac{1}{3} \times \frac{2}{3}}$$

$$< 0.5$$

∈ Class 2.



⇒ If now training data 許少，use generative model

* generative model 有自己的假設

有時會無 data，遵從內心的假設

⇒ if label 有問題，is noise

⇒ Use generative model，可以把有問題的忽略掉

* 做 discriminative 時假設上是假設了 posterior probability

然後 find posterior 內的參數

but if now do generative；會把 formulation 扩大更精確

< prior (文字) ⇒ 取決於 proper estimator
class-dependent bay probability (voice & 文字) ⇒ 來自不同的來源

for example = 語音辨識 現在都是用 neural network

⇒ is a discriminative method

but 整個 system 是一個 generative 的 system

全部 use DNN 可能沒有那麼精確，still need to 算 prior 的 probability

↳ 某句話被說出來的機率 = if we want estimator, 並不需要有語音的 data
only need 許多文字 in 網路, then 以算出機率

Multi-class classification

- $C_1 : w^1, b_1 \quad z_1 = w^1 x + b_1$
- $C_2 : w^2, b_2 \quad z_2 = w^2 x + b_2$
- $C_3 : w^3, b_3 \quad z_3 = w^3 x + b_3$

可以是 $-\infty \sim +\infty$ 的任何值

对最大值做強化

softmax function (类似作 normalization)

→ 把 z_1, z_2, z_3 取 exponential

the 把 them summation $\times e \sum_{j=1}^3 e^{z_j}$

$$\Rightarrow \text{then } y_1 = \frac{e^{z_1}}{\sum_{j=1}^3 e^{z_j}}$$

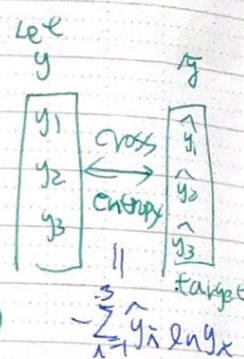
$$y_1 = P(C_1 | X)$$

$$y_2 = \frac{e^{z_2}}{\sum_{j=1}^3 e^{z_j}}$$

$$y_2 = P(C_2 | X)$$

$$y_3 = \frac{e^{z_3}}{\sum_{j=1}^3 e^{z_j}}$$

$$y_3 = P(C_3 | X)$$



$$\Rightarrow 0 < y_i < 1$$

$$\textcircled{3} \quad \sum_{j=1}^n y_j = 1$$

(Pf:

假設 3 個 class 都是 Gaussian distribution

and 共用同一 covariance matrix

P 209-210

we will get softmax function

google = maximum entropy

→ 和 Logistic Regression Same

If $x \in \text{class1}$

$$\hat{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$\in \text{class2}$

$$\hat{y} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$\in \text{class3}$

$$\hat{y} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

\Rightarrow cross entropy 這個 function (minimize cross entropy)
其實 maximize likelihood

多个 class is same

\rightarrow 把 max likelihood 的那个 function 列出来
then 整理 will get minimize cross entropy

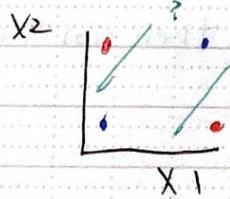
We want to use logistic regression 分类.

\Rightarrow We can't

the output we hope $\begin{cases} \text{class1} & y > 0.5 \\ \text{class2} & y < 0.5 \end{cases}$

\Rightarrow 因为 logistic regression 两个 class 之间的 boundary 是一条直线

\Rightarrow We can't 把 red put - 直, blue put - 曲



\Rightarrow if We still want to use logistic regression

\rightarrow We can use Feature Transformation

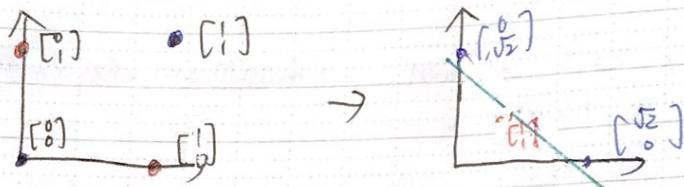
做一些转换去找一个比较好的 feature space

\rightarrow 其可以言之 logistic regression 原理

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{bmatrix} x'_1 \\ x'_2 \end{bmatrix}$$

Let x'_1 be the distance to $[0]$

x'_2 be the distance to $[1]$



$$[0] \rightarrow [x'_1]$$

$$[1] \rightarrow [x'_2]$$

but not always easy to find a good feature transformation
 \Rightarrow we hope this transformation 是由机器自己產生的

\Rightarrow 把很多个 Logistic Regression cascade 起來
 (接)
 feature transform

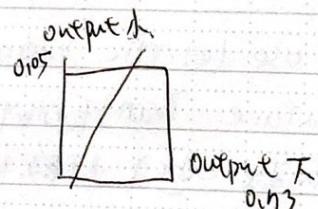
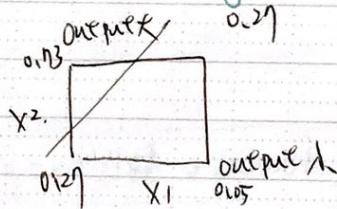
$$x_1 \times w + x_2 \times w \rightarrow z_1 \text{ by sigmoid function} \rightarrow \text{output is } x'_1$$

$\beta - 1$: logistic regression model

$$x_1 \times w + x_2 \times w \rightarrow z_2 \rightarrow \text{sigmoid function} \rightarrow \text{output is } x'_2$$

\hookrightarrow 接一个 logistic regression by model \rightarrow then can classification
 (input x'_1, x'_2)

$$\rightarrow y.$$



* 把每个 logistic
 regression by 一个
 Neuron

把这些 logistic regression
 串起来 就成这个 network
 就叫作 neural Network
 (神经元之网络)

