

Classification: Logistic Regression

Hung-yi Lee

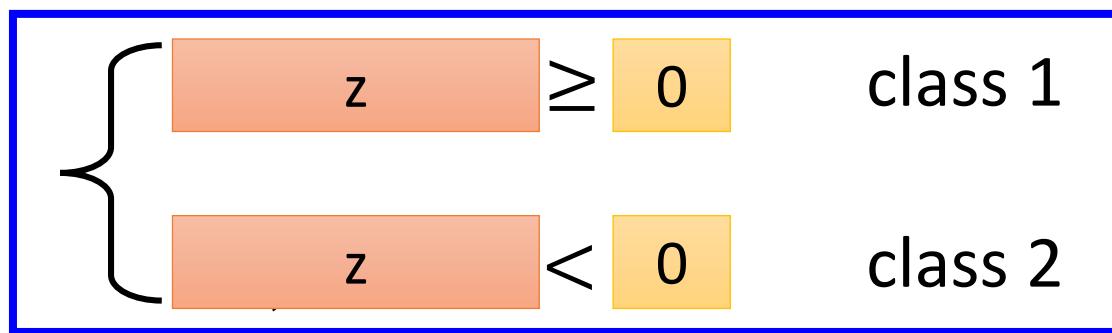
李宏毅

有關分組

- 作業以個人為單位繳交
- 期末專題才需要分組
- 找不到組員也沒有關係，期末專題公告後找不到組員的同學助教會幫忙湊對

Step 1: Function Set

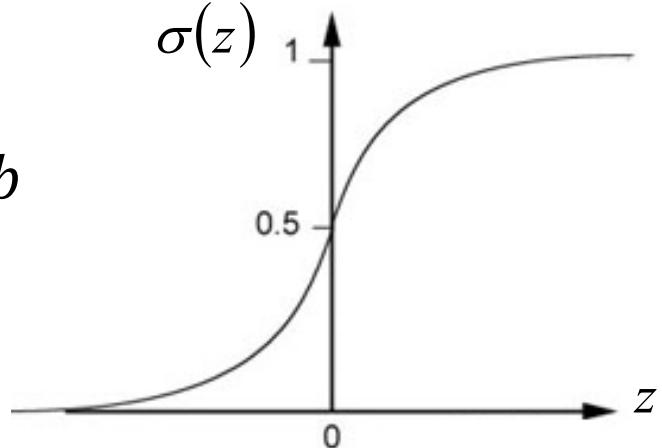
Function set: Including all different w and b



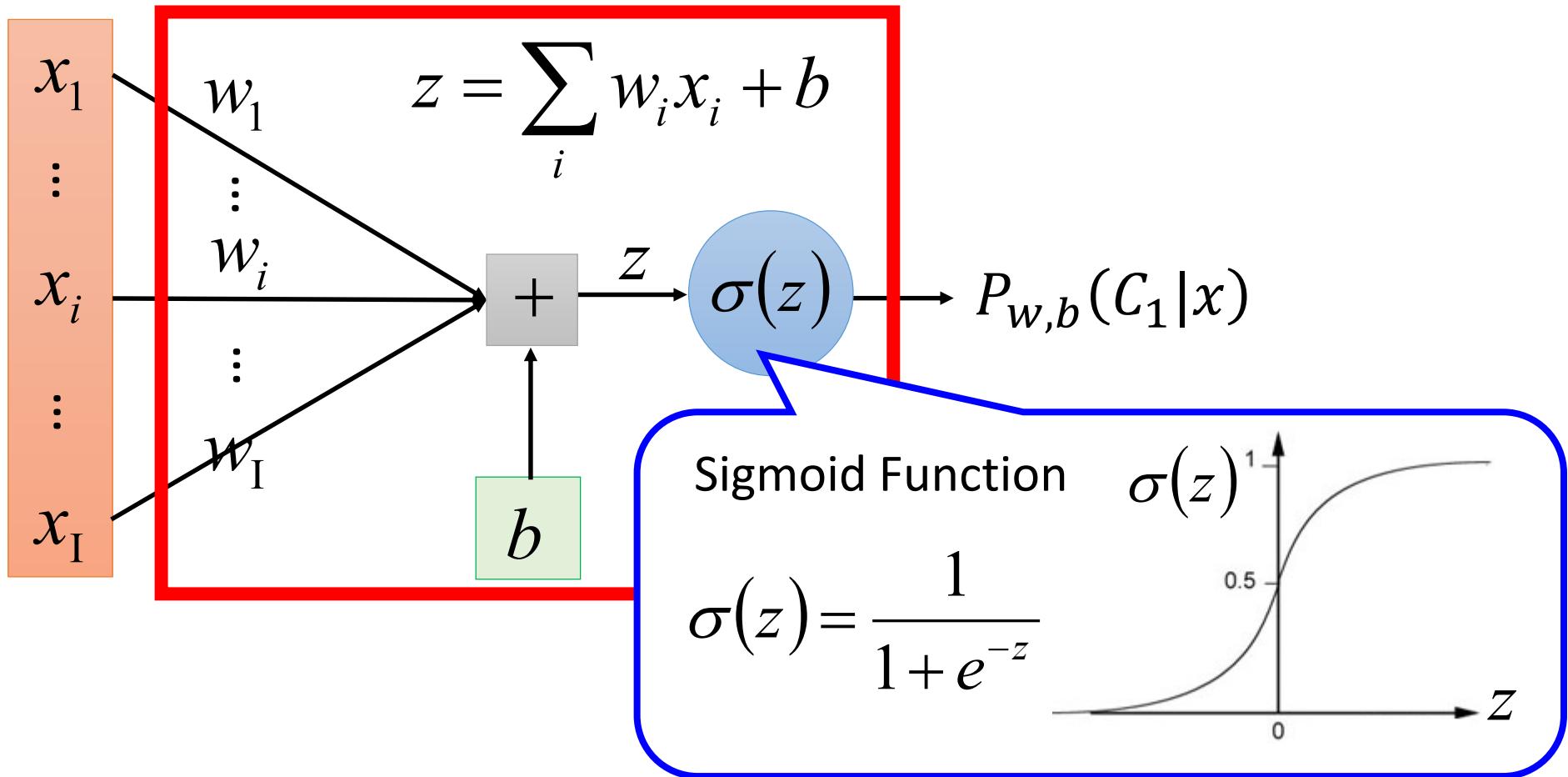
$$P_{w,b}(C_1|x) = \sigma(z)$$

$$z = w \cdot x + b = \sum_i w_i x_i + b$$

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$



Step 1: Function Set



Step 2: Goodness of a Function

Training
Data

	x^1	x^2	x^3	x^N
	C_1	C_1	C_2	C_1

Assume the data is generated based on $f_{w,b}(x) = P_{w,b}(C_1|x)$

Given a set of w and b, what is its probability of generating the data?

因為只有兩個class，所以用一扣掉即為class2之機率

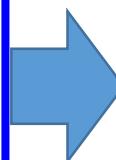
$$L(w, b) = f_{w,b}(x^1)f_{w,b}(x^2)\left(1 - f_{w,b}(x^3)\right)\cdots f_{w,b}(x^N)$$

The most likely w^* and b^* is the one with the largest $L(w, b)$.

maximum likelihood

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

$$\begin{matrix} x^1 & x^2 & x^3 \\ C_1 & C_1 & C_2 \end{matrix} \quad \dots \dots$$



$$\begin{matrix} x^1 & x^2 & x^3 \\ \hat{y}^1 = 1 & \hat{y}^2 = 1 & \hat{y}^3 = 0 \end{matrix} \quad \dots \dots$$

Notation 改變

\hat{y}^n : 1 for class 1, 0 for class 2

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

取log order不變，負號將max改成min

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

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

$$\begin{aligned} -\ln L(w, b) &= -\ln f_{w,b}(x^1) \rightarrow -[1 \ln f(x^1) + 0 \ln(1 - f(x^1))] \\ &\quad \hat{y} \quad 1 - \hat{y} \\ -\ln f_{w,b}(x^2) &\rightarrow -[1 \ln f(x^2) + 0 \ln(1 - f(x^2))] \\ -\ln(1 - f_{w,b}(x^3)) &\rightarrow -[0 \ln f(x^3) + 1 \ln(1 - f(x^3))] \\ &\quad \vdots \end{aligned}$$

Step 2: Goodness of a Function

$$L(w, b) = f_{w,b}(x^1)f_{w,b}(x^2)\left(1 - f_{w,b}(x^3)\right)\cdots f_{w,b}(x^N)$$

$$-lnL(w, b) = lnf_{w,b}(x^1) + lnf_{w,b}(x^2) + ln\left(1 - f_{w,b}(x^3)\right)\cdots$$

loss function \hat{y}^n : 1 for class 1, 0 for class 2

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

Cross entropy between two Bernoulli distribution

Distribution p:

$$p(x = 1) = \hat{y}^n$$

$$p(x = 0) = 1 - \hat{y}^n$$

Distribution q:

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

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

cross
entropy

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

cross entropy 表示
這兩筆資訊有多接近

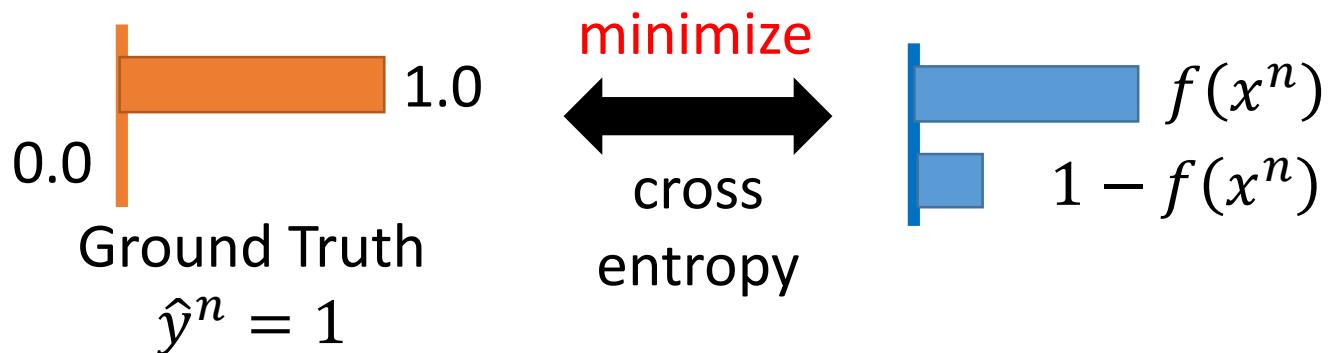
Step 2: Goodness of a Function

找max $L(w, b) = f_{w,b}(x^1)f_{w,b}(x^2)\left(1 - f_{w,b}(x^3)\right)\cdots f_{w,b}(x^N)$

找min $-lnL(w, b) = lnf_{w,b}(x^1) + lnf_{w,b}(x^2) + ln\left(1 - f_{w,b}(x^3)\right)\cdots$

loss function \hat{y}^n : 1 for class 1, 0 for class 2

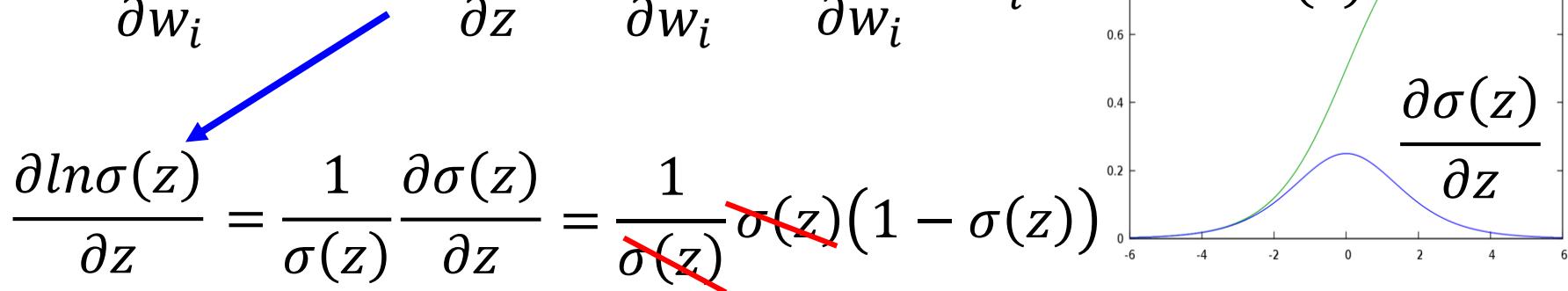
$$= \sum_n - \underbrace{\left[\hat{y}^n \ln f_{w,b}(x^n) + (1 - \hat{y}^n) \ln \left(1 - f_{w,b}(x^n)\right) \right]}_{\text{Cross entropy between two Bernoulli distribution}}$$



Step 3: Find the best function

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

$$\frac{\partial \ln f_{w,b}(x)}{\partial w_i} = \frac{\partial \ln f_{w,b}(x)}{\partial z} \frac{\partial z}{\partial w_i} \quad \frac{\partial z}{\partial w_i} = x_i$$



$$\begin{aligned} f_{w,b}(x) &= \sigma(z) \\ &= 1/(1 + \exp(-z)) \end{aligned}$$

$$z = w \cdot x + b = \sum_i w_i x_i + b$$

Step 3: Find the best function

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

$$\frac{\partial \ln(1 - f_{w,b}(x))}{\partial w_i} = \frac{\partial \ln(1 - f_{w,b}(x))}{\partial z} \frac{\partial z}{\partial w_i} \quad \frac{\partial z}{\partial w_i} = x_i$$

$$\frac{\partial \ln(1 - \sigma(z))}{\partial z} = -\frac{1}{1 - \sigma(z)} \frac{\partial \sigma(z)}{\partial z} = -\frac{1}{1 - \sigma(z)} \sigma(z)(1 - \sigma(z))$$

$$f_{w,b}(x) = \sigma(z)$$

$$= 1 / 1 + \exp(-z)$$

$$z = w \cdot x + b = \sum_i w_i x_i + b$$

Step 3: Find the best function

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

$$= \sum_n - \left[\hat{y}^n \underbrace{(1 - f_{w,b}(x^n))}_{\textcolor{blue}{-}} x_i^n - (1 - \hat{y}^n) \underbrace{f_{w,b}(x^n)}_{\textcolor{blue}{-}} x_i^n \right]$$

$$= \sum_n - \left[\hat{y}^n - \cancel{\hat{y}^n f_{w,b}(x^n)} - f_{w,b}(x^n) + \cancel{\hat{y}^n f_{w,b}(x^n)} \right] \underbrace{x_i^n}_{\textcolor{blue}{-}}$$

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

Larger difference, larger update

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

w的gradient

Logistic Regression + Square Error

Step 1: $f_{w,b}(x) = \sigma\left(\sum_i w_i x_i + b\right)$

Step 2: Training data: (x^n, \hat{y}^n) , \hat{y}^n : 1 for class 1, 0 for class 2

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

Step 3:

$$\frac{\partial (f_{w,b}(x) - \hat{y})^2}{\partial w_i} = 2(f_{w,b}(x) - \hat{y}) \frac{\partial f_{w,b}(x)}{\partial z} \frac{\partial z}{\partial w_i}$$
$$= 2(f_{w,b}(x) - \hat{y}) f_{w,b}(x) (1 - f_{w,b}(x)) x_i$$

$$\hat{y}^n = 1 \quad \text{If } f_{w,b}(x^n) = 1 \text{ (close to target)} \rightarrow \partial L / \partial w_i = 0$$

error!!

$$\text{If } f_{w,b}(x^n) = 0 \text{ (far from target)} \rightarrow \partial L / \partial w_i = 0$$

Logistic Regression + Square Error

Step 1: $f_{w,b}(x) = \sigma\left(\sum_i w_i x_i + b\right)$

Step 2: Training data: (x^n, \hat{y}^n) , \hat{y}^n : 1 for class 1, 0 for class 2

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

Step 3:

$$\frac{\partial (f_{w,b}(x) - \hat{y})^2}{\partial w_i} = 2(f_{w,b}(x) - \hat{y}) \frac{\partial f_{w,b}(x)}{\partial z} \frac{\partial z}{\partial w_i}$$
$$= 2(f_{w,b}(x) - \hat{y}) f_{w,b}(x) (1 - f_{w,b}(x)) x_i$$

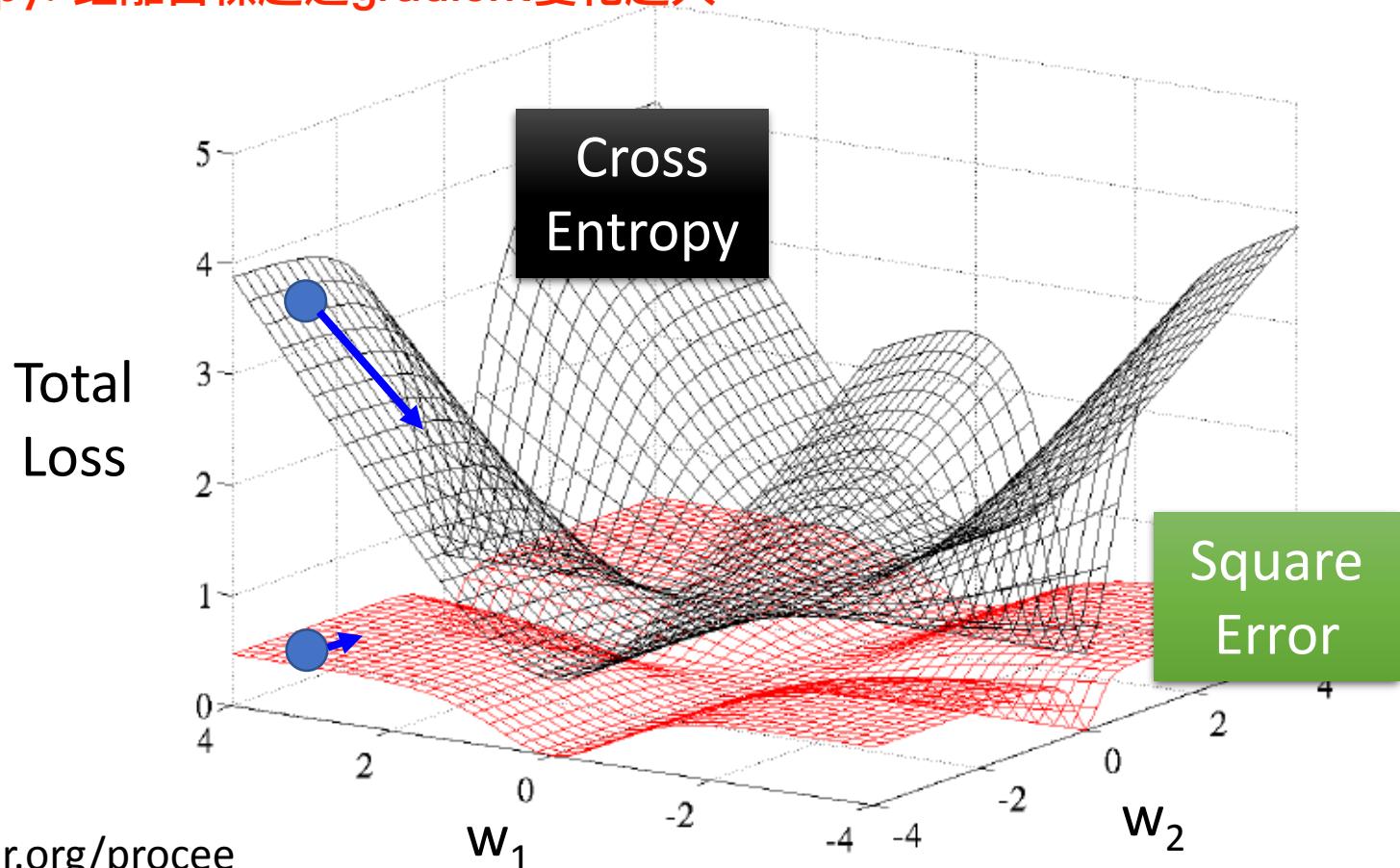
error!!

$$\hat{y}^n = 0 \quad \boxed{\text{If } f_{w,b}(x^n) = 1 \text{ (far from target)} \rightarrow \partial L / \partial w_i = 0}$$

$$\text{If } f_{w,b}(x^n) = 0 \text{ (close to target)} \rightarrow \partial L / \partial w_i = 0$$

Cross Entropy v.s. Square Error

cross entropy: 距離目標越遠gradient變化越大



<http://jmlr.org/proceedings/papers/v9/glorot10a/glorot10a.pdf>

square error: 距離目標即使很遠但gradient變化很小

Logistic Regression

Step 1: $f_{w,b}(x) = \sigma\left(\sum_i w_i x_i + b\right)$

Output: between 0 and 1

Linear Regression

$$f_{w,b}(x) = \sum_i w_i x_i + b$$

Output: any value

Step 2:

Step 3:

Logistic Regression

Step 1: $f_{w,b}(x) = \sigma\left(\sum_i w_i x_i + b\right)$

Output: between 0 and 1

Linear Regression

$$f_{w,b}(x) = \sum_i w_i x_i + b$$

Output: any value

Training data: (x^n, \hat{y}^n)

Step 2: \hat{y}^n : 1 for class 1, 0 for class 2

$$L(f) = \sum_n l(f(x^n), \hat{y}^n)$$

Training data: (x^n, \hat{y}^n)

\hat{y}^n : a real number

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

Cross entropy:

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

Logistic Regression

Step 1: $f_{w,b}(x) = \sigma\left(\sum_i w_i x_i + b\right)$

Output: between 0 and 1

Training data: (x^n, \hat{y}^n)

Step 2: \hat{y}^n : 1 for class 1, 0 for class 2

$$L(f) = \sum_n l(f(x^n), \hat{y}^n)$$

Linear Regression

$$f_{w,b}(x) = \sum_i w_i x_i + b$$

Output: any value

Training data: (x^n, \hat{y}^n)

\hat{y}^n : a real number

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

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

Step 3:

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

logistic regression Gaussian

Discriminative v.s. Generative

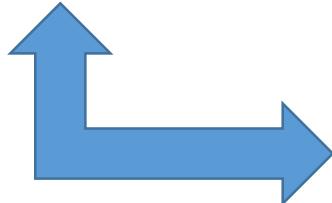
performance 較好

$$P(C_1|x) = \sigma(w \cdot x + b)$$



directly find w and b

利用gradient descend



Will we obtain the same set of w and b ? **No, 因為做了不同的假設**



Find $\mu^1, \mu^2, \Sigma^{-1}$

$$w^T = (\mu^1 - \mu^2)^T \Sigma^{-1}$$

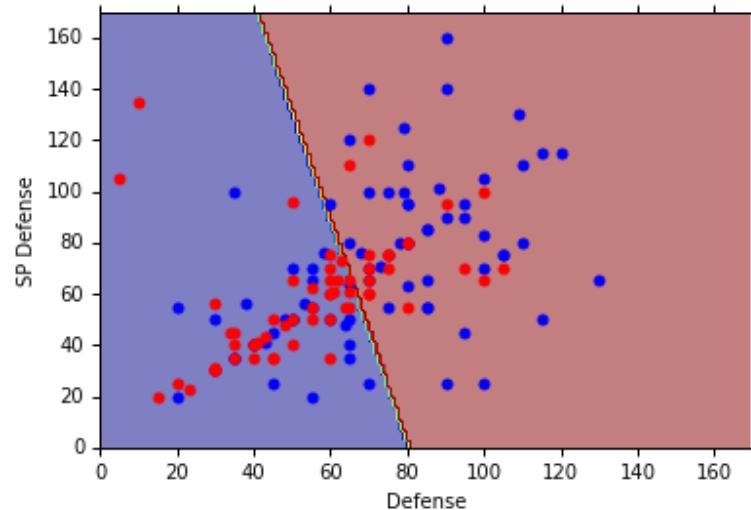
$$b = -\frac{1}{2}(\mu^1)^T (\Sigma^1)^{-1} \mu^1$$

$$+ \frac{1}{2}(\mu^2)^T (\Sigma^2)^{-1} \mu^2 + \ln \frac{N_1}{N_2}$$

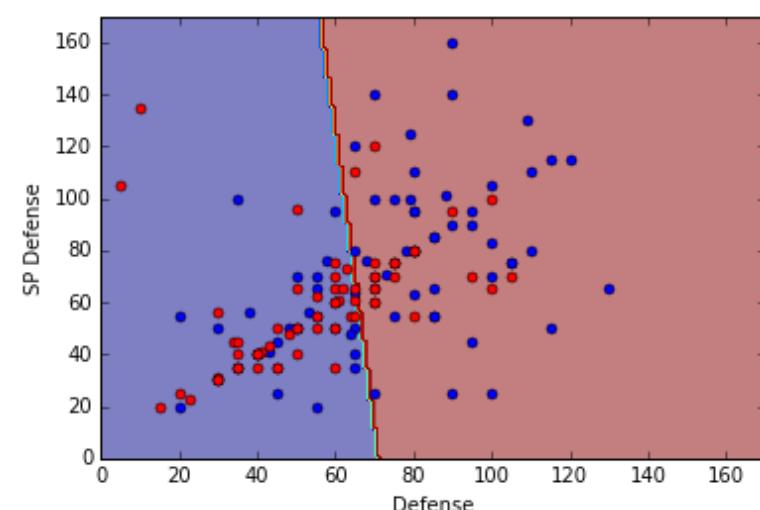
The same model (function set), but different function may be selected by the same training data.

Gaussian logistic regression
Generative v.s. Discriminative

Generative



Discriminative



All: hp, att, sp att, de, sp de, speed

73% accuracy

79% accuracy

Generative做了部分假設，以下例而言，他將兩個dimension視為
indépendant，假設資料sample夠多的話(1,1)也有可能是class2

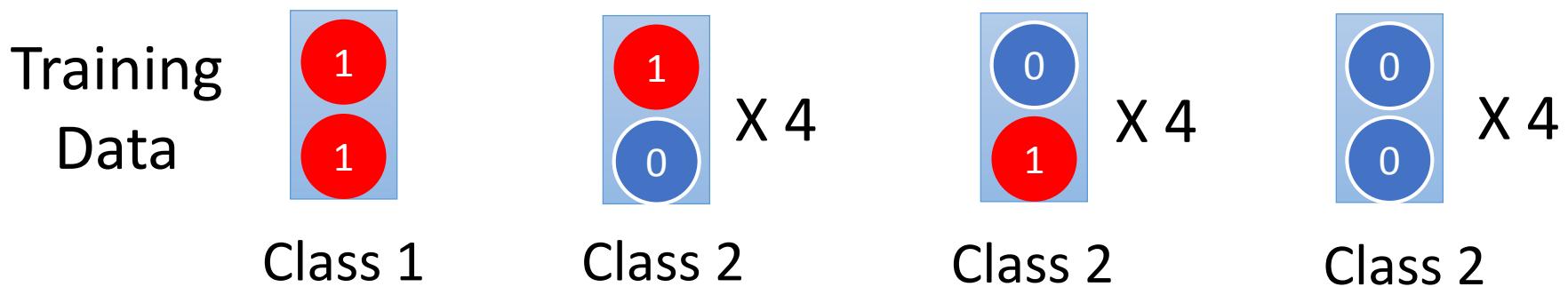
Generative v.s. Discriminative

- Example

Training Data	 X 4	 X 4	 X 4	 X 4
	Class 1	Class 2	Class 2	Class 2
Testing Data		Class 1? Class 2?	假設所有feature為independant How about Naïve Bayes?	
			$P(x C_i) = P(x_1 C_i)P(x_2 C_i)$	

Generative v.s. Discriminative

- Example



$$P(C_1) = \frac{1}{13}$$

$$P(x_1 = 1|C_1) = 1$$

$$P(x_2 = 1|C_1) = 1$$

$$P(C_2) = \frac{12}{13}$$

$$P(x_1 = 1|C_2) = \frac{1}{3}$$

$$P(x_2 = 1|C_2) = \frac{1}{3}$$

Training
Data



Class 1



X 4



X 4



X 4

Testing
Data



對naive bayse來說不考慮不同
dimension之間的corelation

$$P(C_1|x) < 0.5$$

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

$$1 \times 1$$

$$\frac{1}{13}$$

$$\frac{1}{3} \times \frac{1}{3}$$

$$\frac{1}{13}$$

$$\frac{12}{13}$$

$$P(C_1) = \frac{1}{13}$$

$$P(x_1 = 1|C_1) = 1$$

$$P(x_2 = 1|C_1) = 1$$

$$P(C_2) = \frac{12}{13}$$

$$P(x_1 = 1|C_2) = \frac{1}{3}$$

$$P(x_2 = 1|C_2) = \frac{1}{3}$$

假設data來自某個機率模型

Generative v.s. Discriminative

naive bayse

- Usually people believe discriminative model is better
- Benefit of generative model
 - With the assumption of probability distribution
 - less training data is needed
 - more robust to the noise 因為自己已經做了部分假設 (independant)
比較不容易受到data的影響
 - Priors and class-dependent probabilities can be estimated from different sources.

Multi-class Classification

(3 classes as example)

$$C_1: w^1, b_1$$

$$z_1 = w^1 \cdot x + b_1$$

$$C_2: w^2, b_2$$

$$z_2 = w^2 \cdot x + b_2$$

$$C_3: w^3, b_3$$

$$z_3 = w^3 \cdot x + b_3$$

Probability:

■ $1 > y_i > 0$

■ $\sum_i y_i = 1$

$$y_i = P(C_i | x)$$

強化資料的差距，並將其做normalization

Softmax

$$z_1 = 3$$

$$z_2 = 1$$

$$z_3 = -3$$

$$e^{z_1}$$

$$e^{z_2}$$

$$e^{z_3}$$

normalization

$$0.88$$

$$0.12$$

$$\approx 0$$

$$y_1 = e^{z_1} / \sum_{j=1}^3 e^{z_j}$$

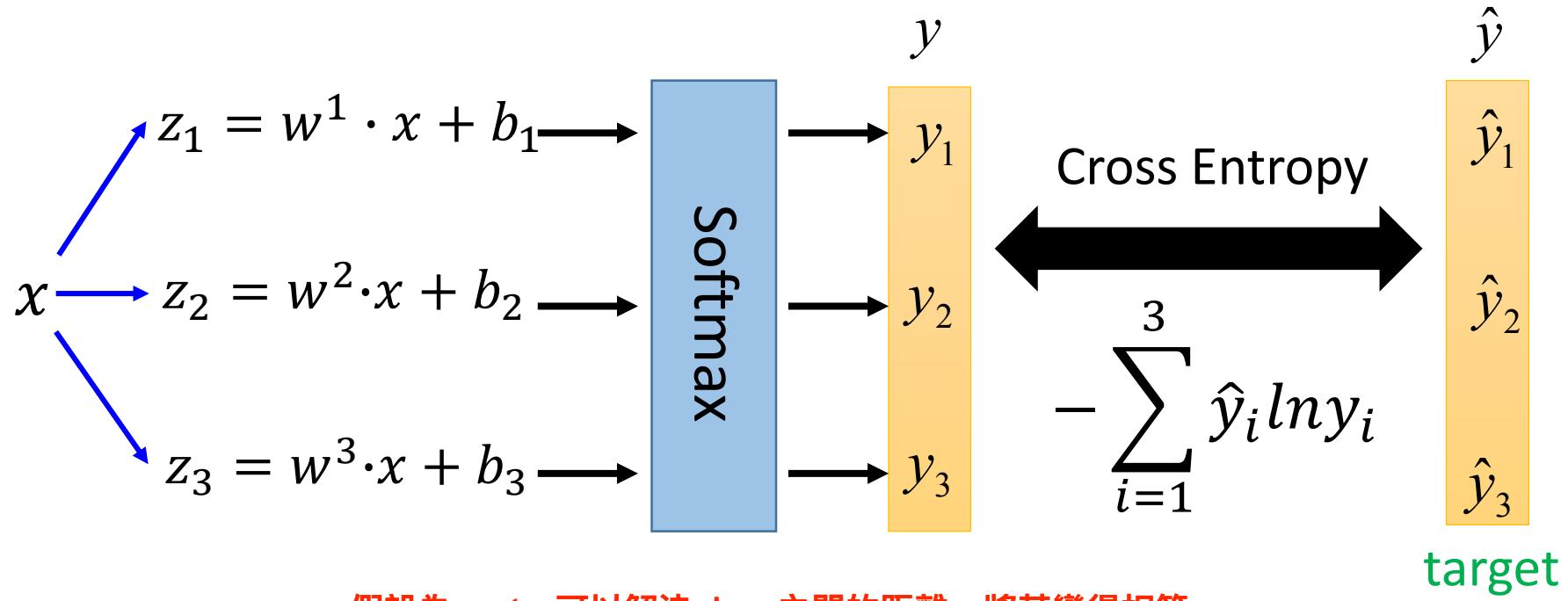
$$y_2 = e^{z_2} / \sum_{j=1}^3 e^{z_j}$$

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

$$+ \sum_{j=1}^3 e^{z_j}$$

Multi-class Classification

(3 classes as example)

If $x \in$ class 1

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

$$-\ln y_1$$

If $x \in$ class 2

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

$$-\ln y_2$$

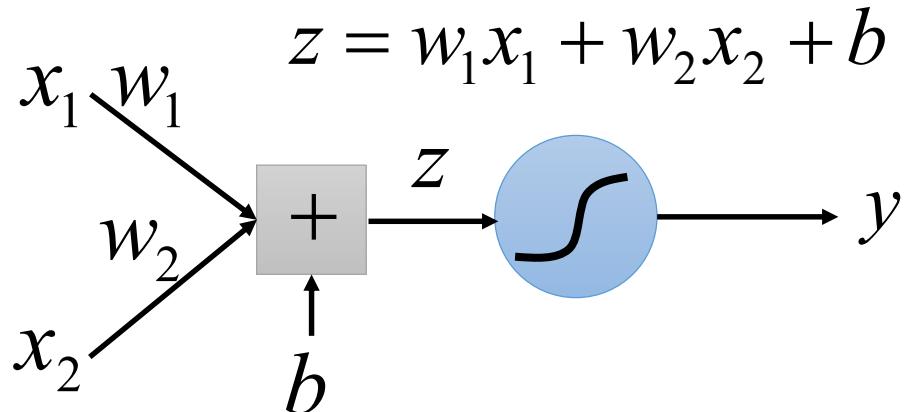
If $x \in$ class 3

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

$$-\ln y_3$$

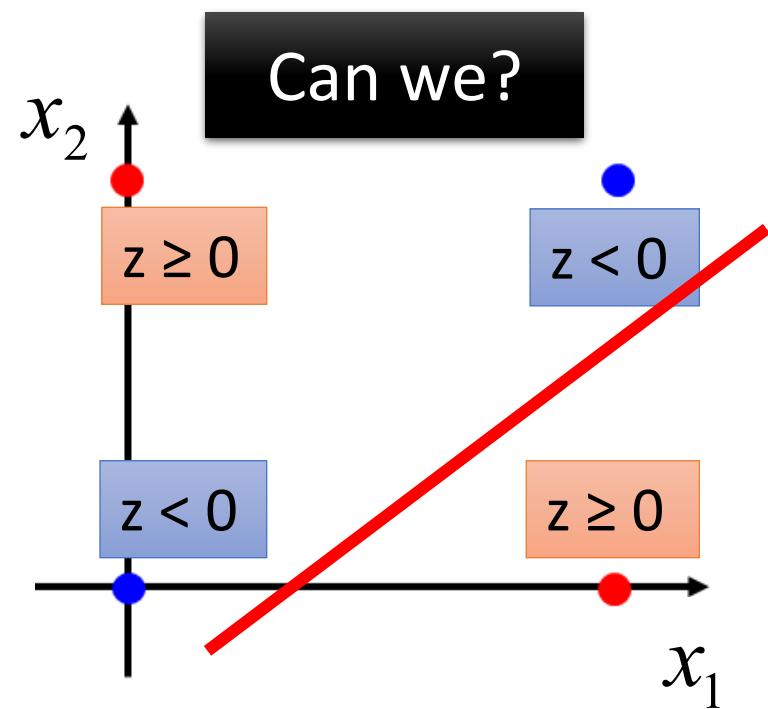
假設為vector可以解決class之間的距離，將其變得相等

Limitation of Logistic Regression



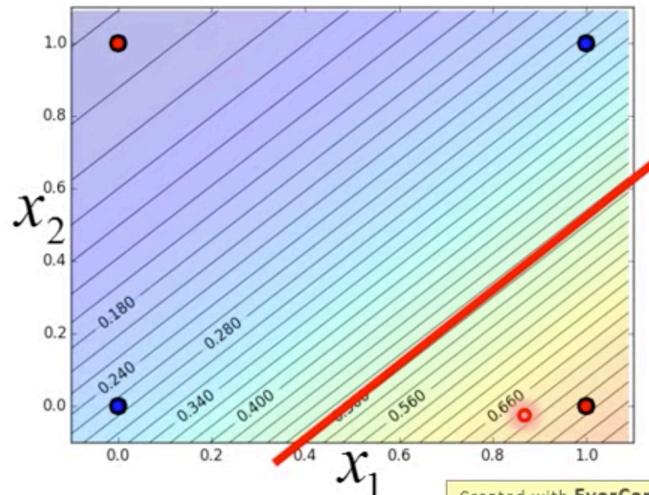
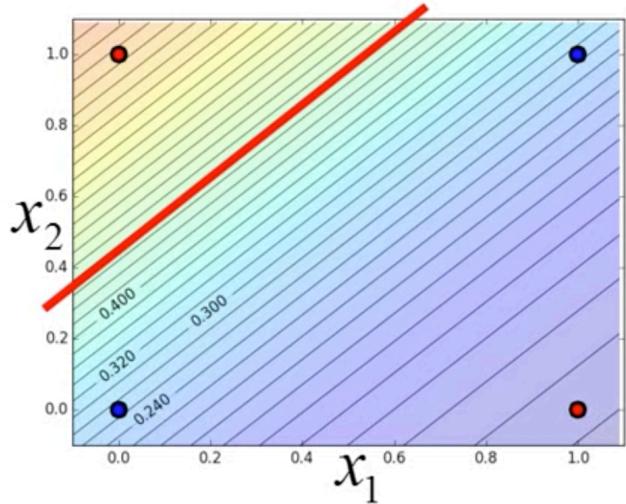
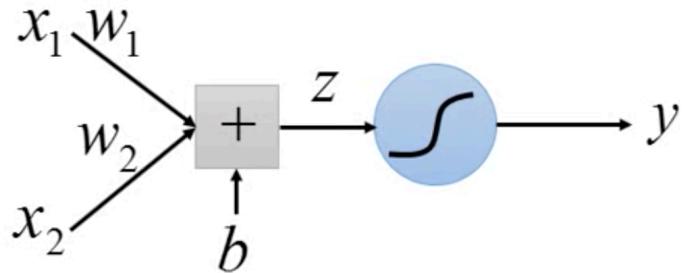
Input Feature		Label
x_1	x_2	
0	0	Class 2
0	1	Class 1
1	0	Class 1
1	1	Class 2

$$\begin{cases} \text{Class 1} & y \geq 0.5 \quad (z \geq 0) \\ \text{Class 2} & y < 0.5 \quad (z < 0) \end{cases}$$



Limitation of Logistic Regression

- No, we can't

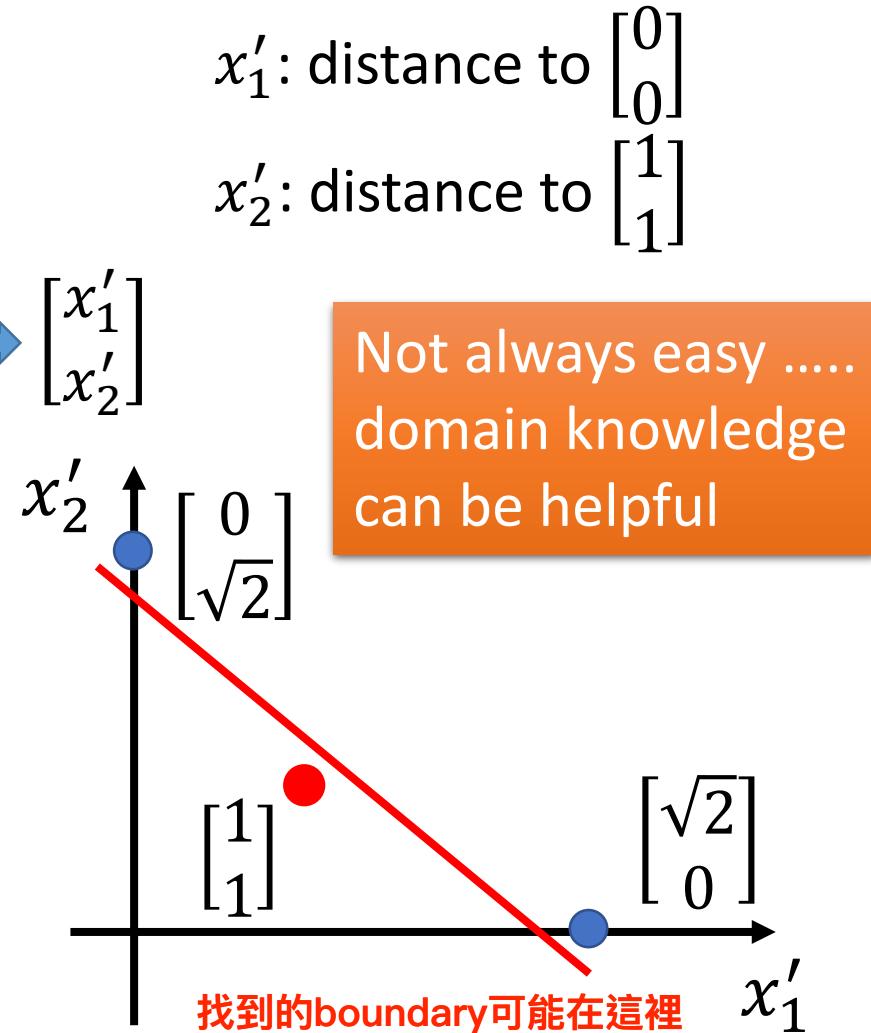
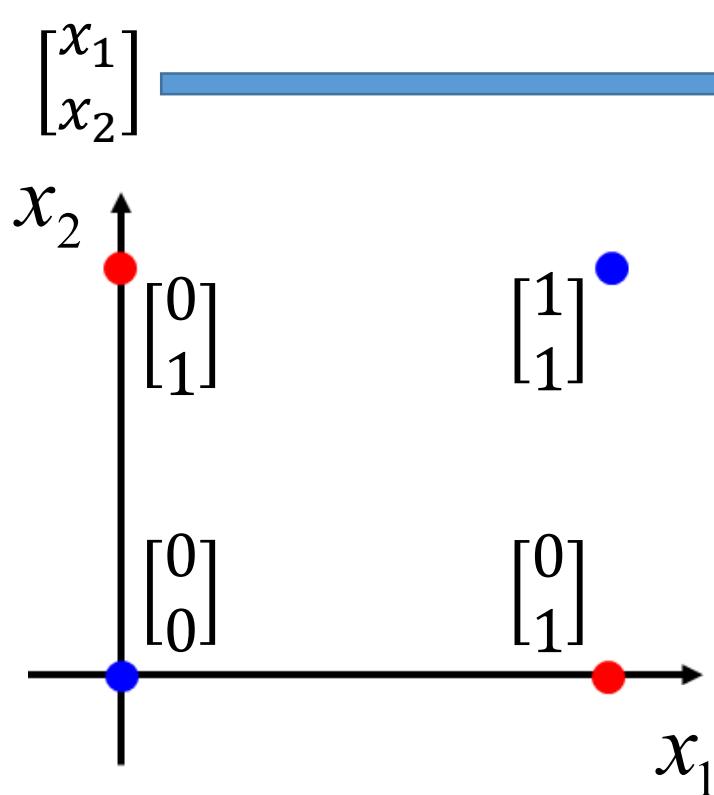


logistic regression分類必須畫一條boundary(直線)

Created with EverCam.
<http://www.camdemey.com>

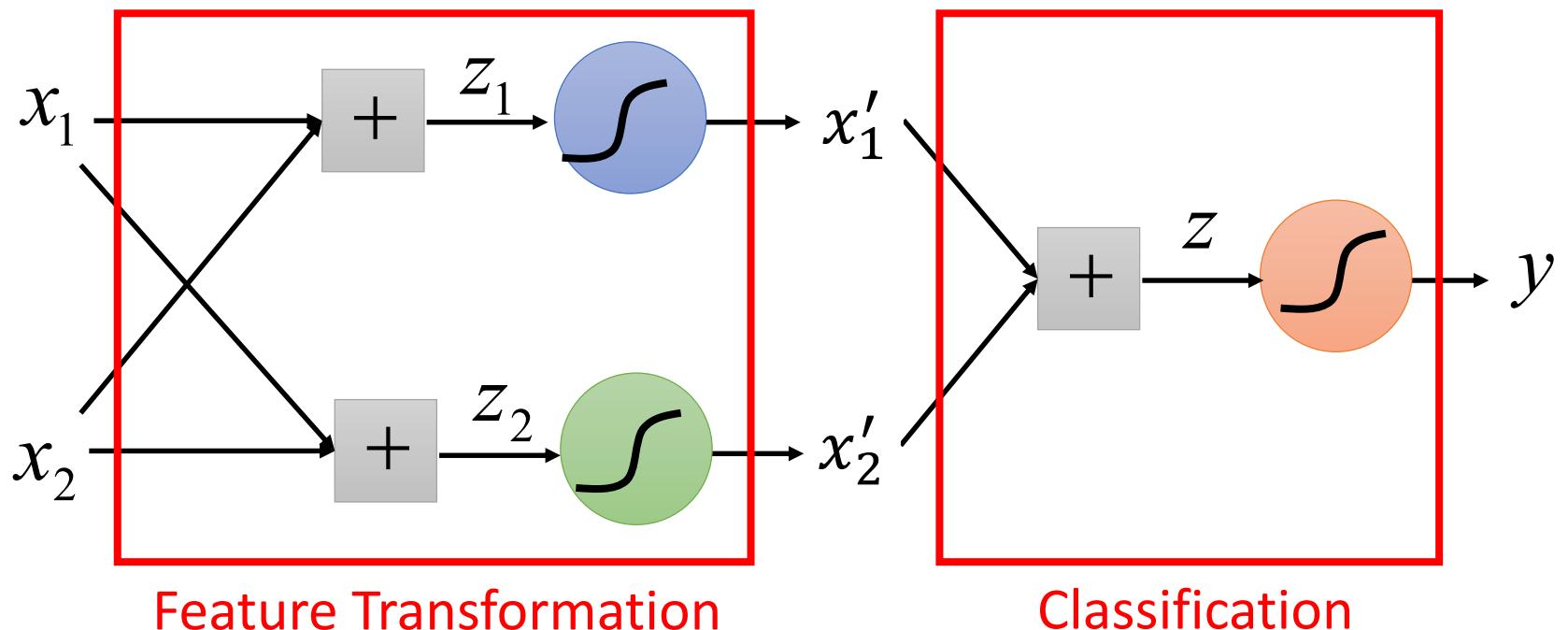
Limitation of Logistic Regression

- Feature transformation

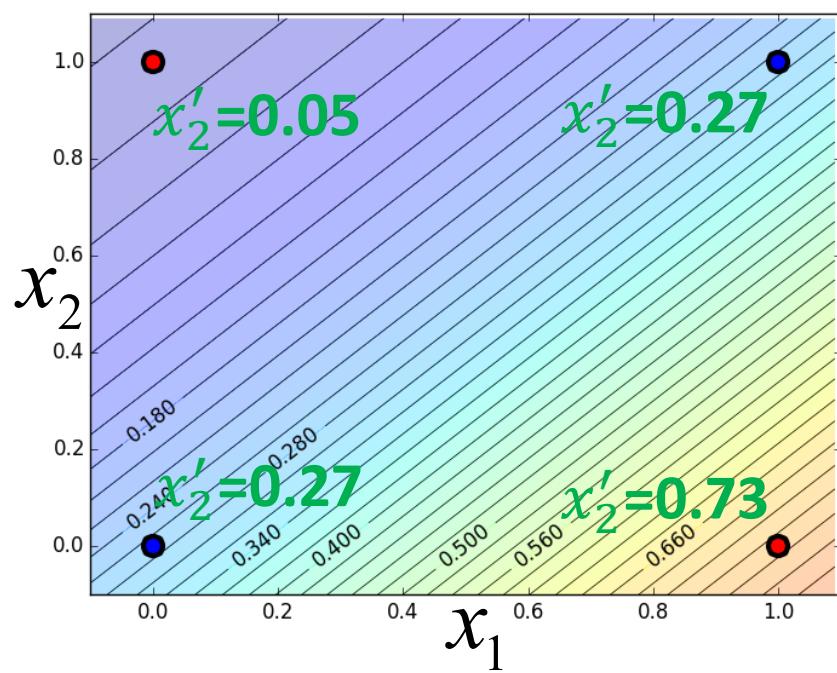
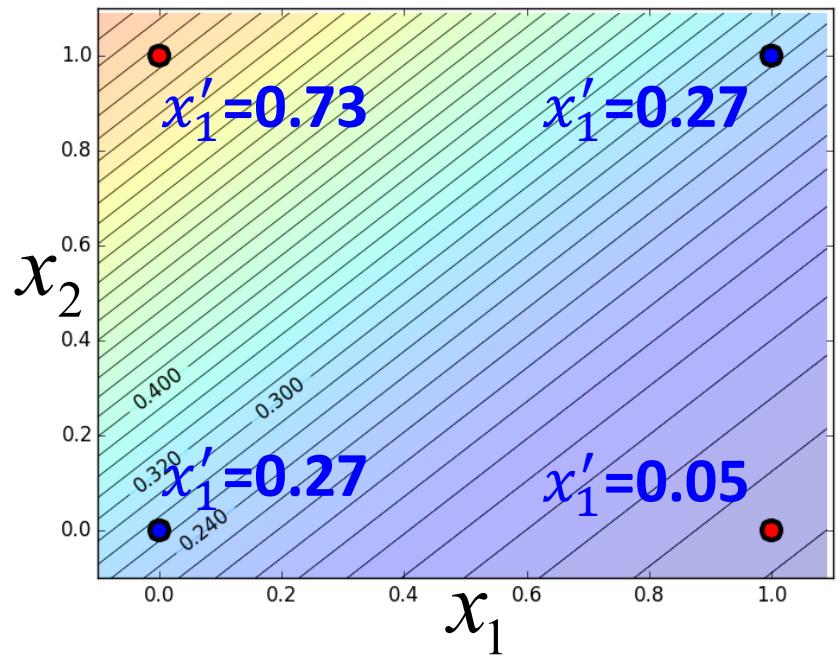
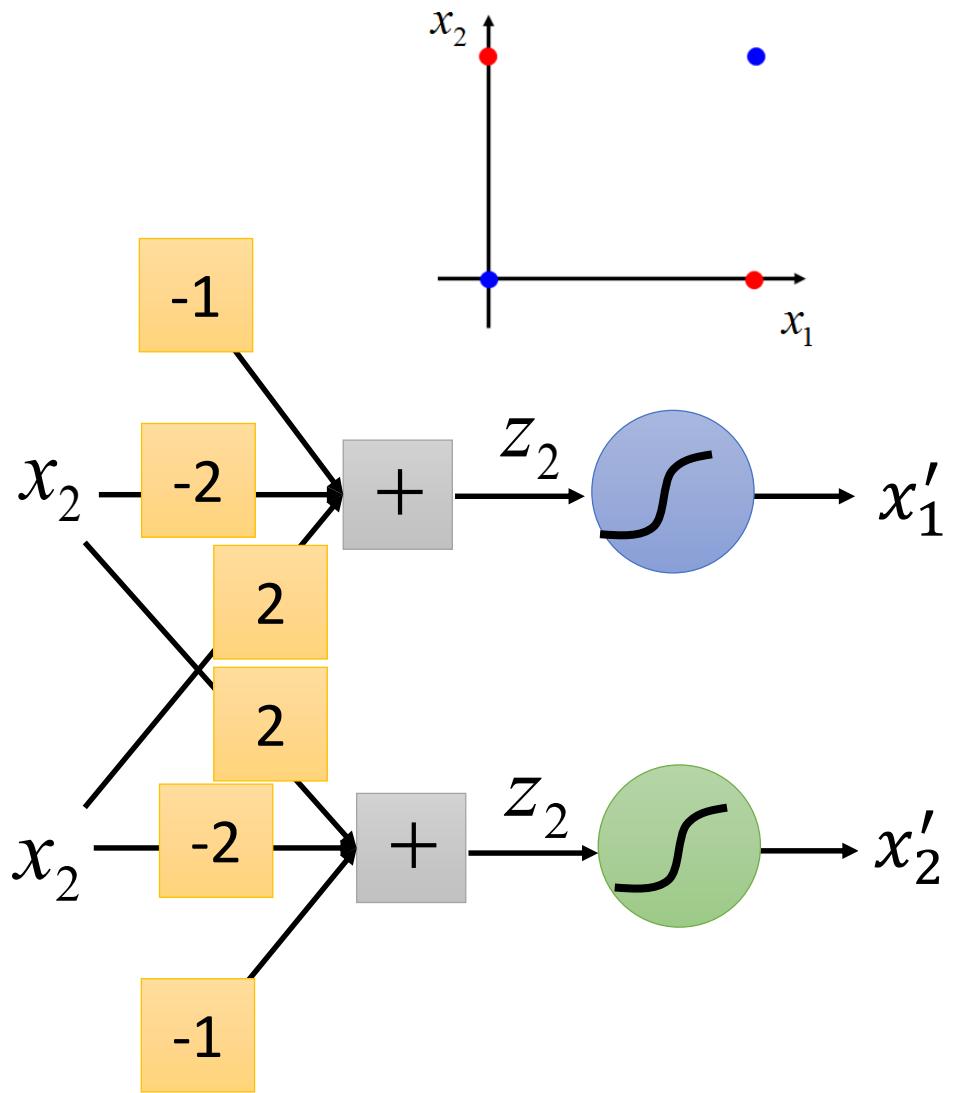


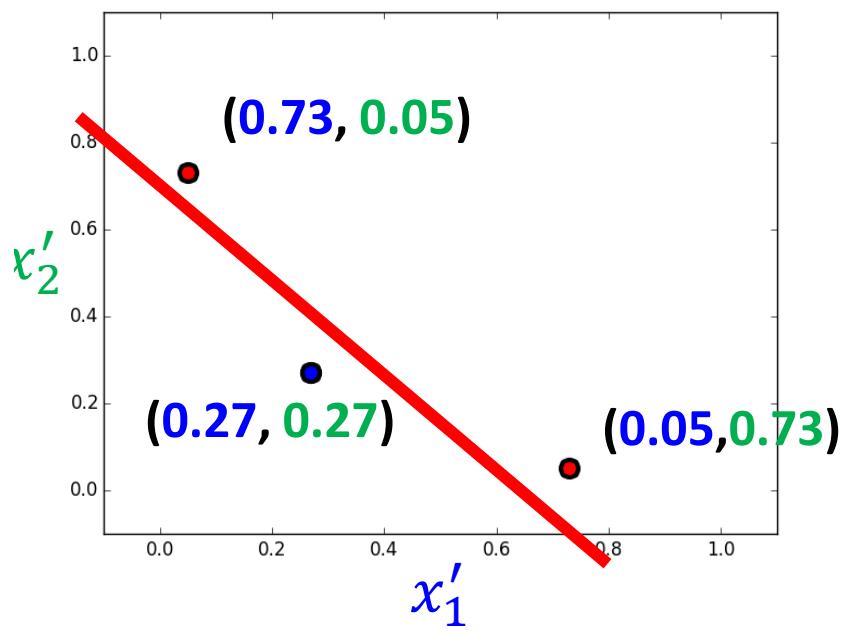
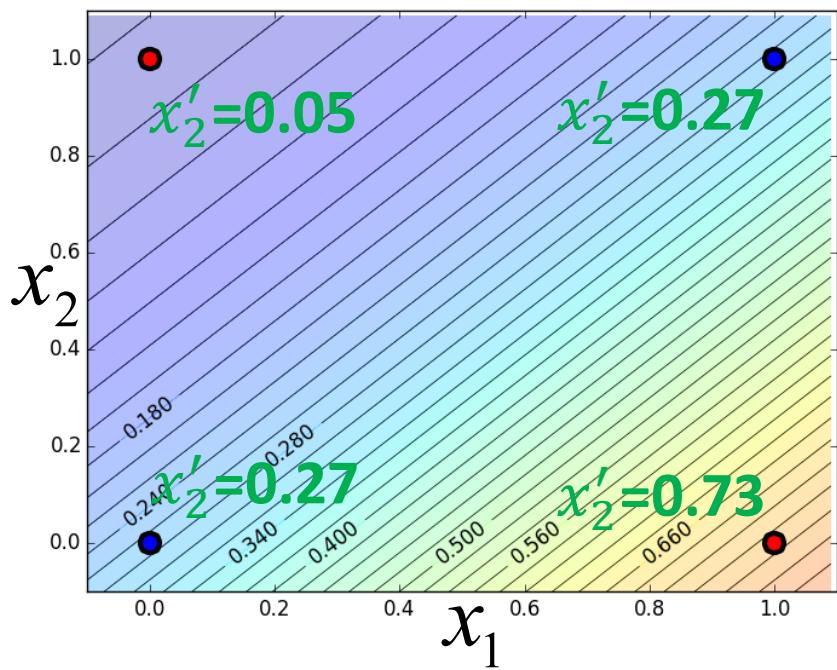
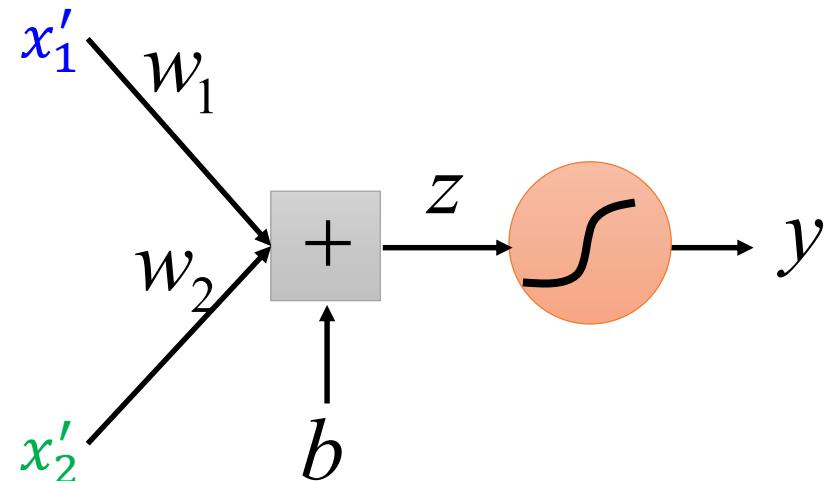
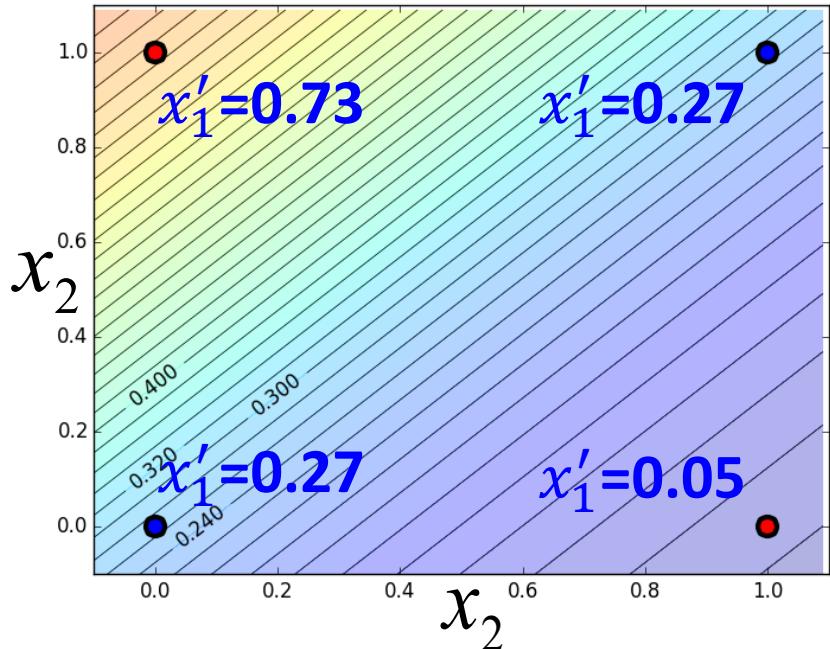
Limitation of Logistic Regression

- Cascading logistic regression models



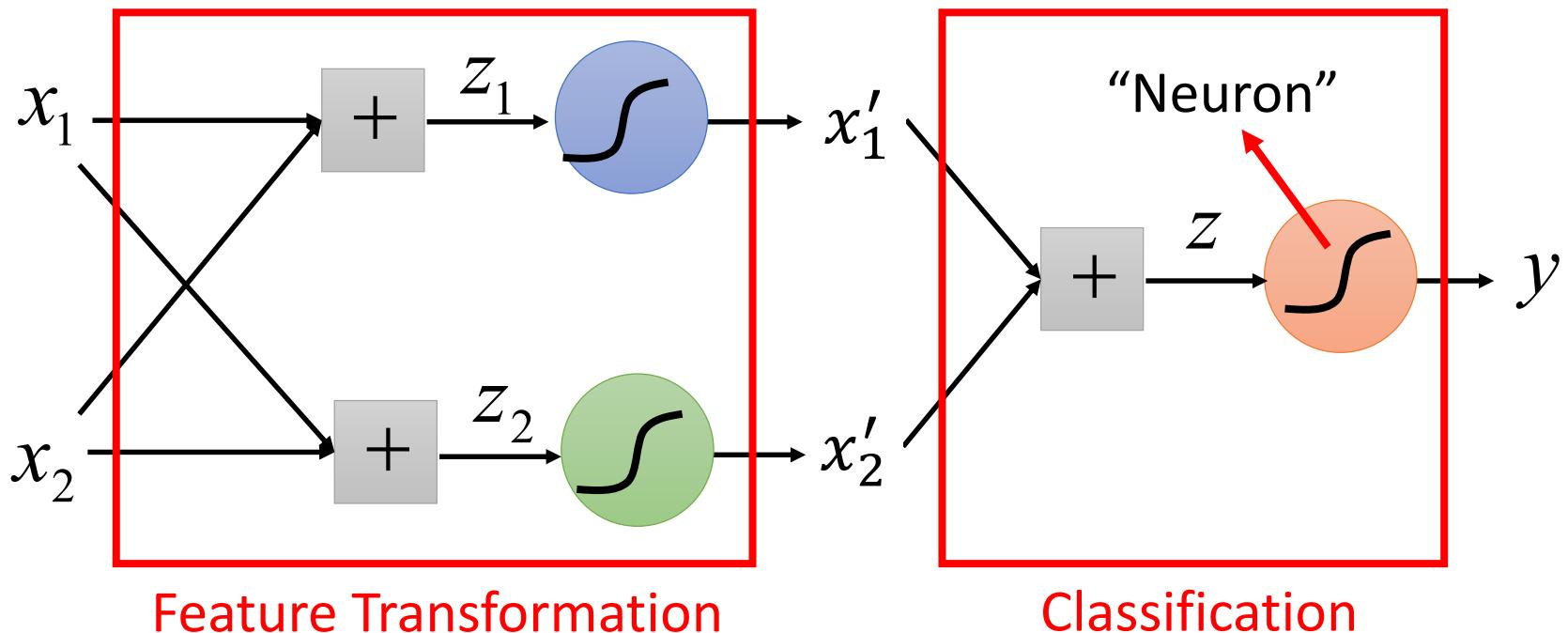
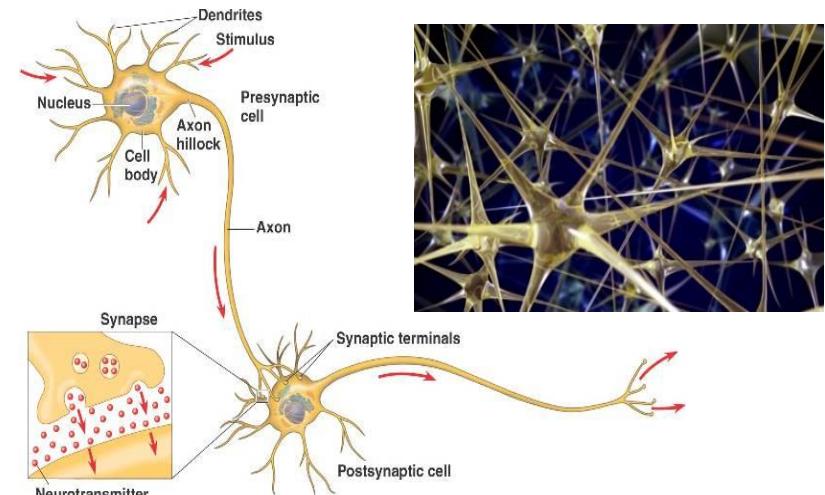
(ignore bias in this figure)





Deep Learning!

All the parameters of the logistic regressions are jointly learned.



Reference

- Bishop: Chapter 4.3

Acknowledgement

- 感謝 林恩妤 發現投影片上的錯誤

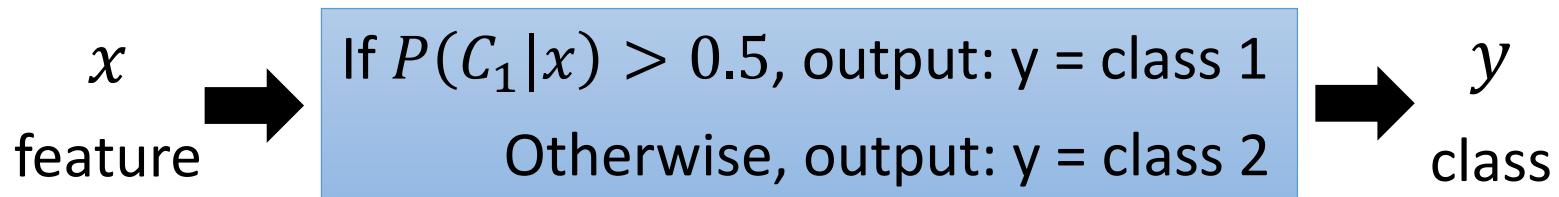
Appendix

Three Steps

x^1	x^2	x^3	x^n
\hat{y}^1	\hat{y}^2	\hat{y}^3	\hat{y}^n

$$\hat{y}^n = \text{class 1, class 2}$$

- Step 1. Function Set (Model)



$$P(C_1|x) = \sigma(w \cdot x + b)$$

w and b are related to $N_1, N_2, \mu^1, \mu^2, \Sigma$

- Step 2. Goodness of a function

$$L(f) = \sum_n \delta(f(x^n) \neq \hat{y}^n) \rightarrow L(f) = \sum_n l(f(x^n) \neq \hat{y}^n)$$

- Step 3. Find the best function: gradient descent

Step 2: Loss function

$$f_{w,b}(x) = \begin{cases} +1 & z \geq 0 \\ -1 & z < 0 \end{cases}$$

Ideal loss:

$$L(f) = \sum_n \delta(f(x^n) \neq \hat{y}^n)$$

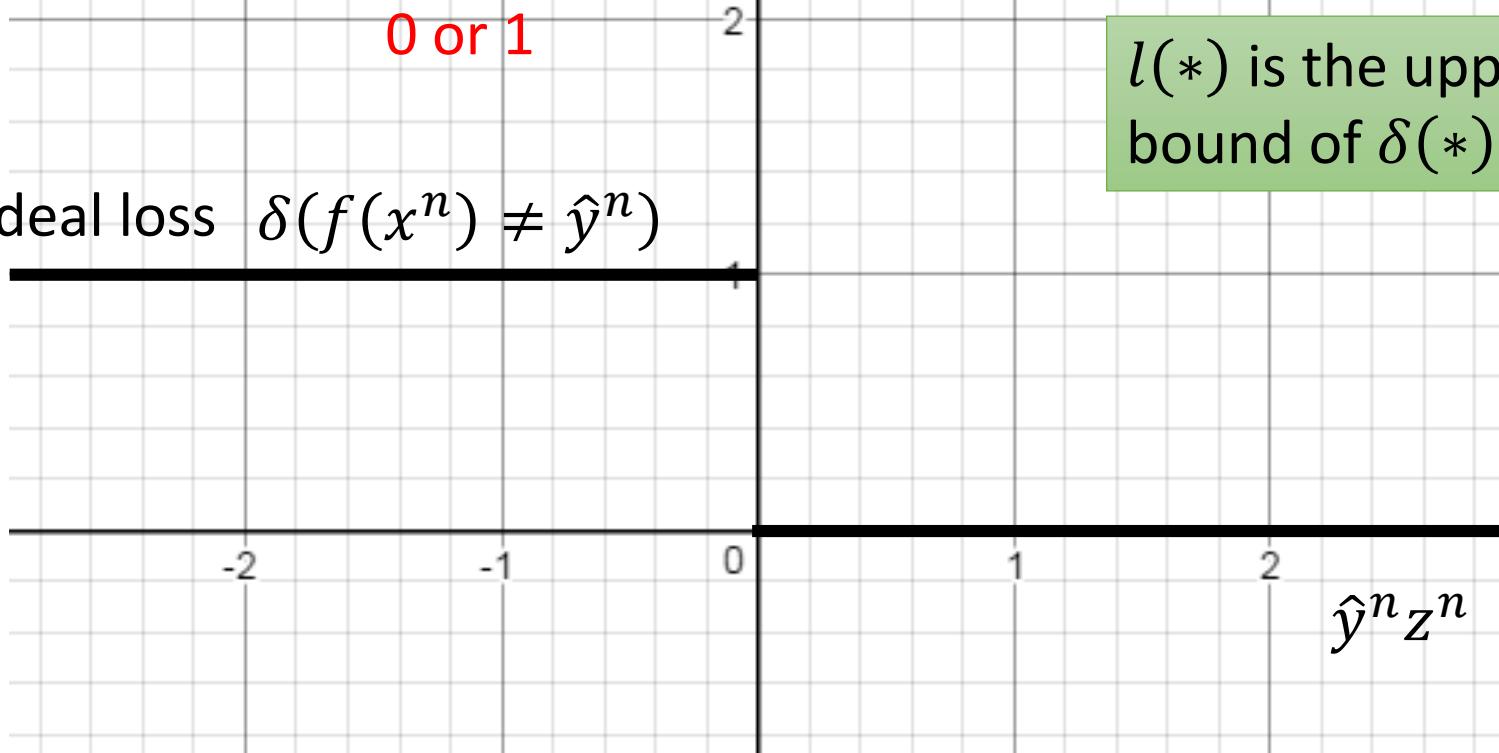
0 or 1

Ideal loss $\delta(f(x^n) \neq \hat{y}^n)$

Approximation:

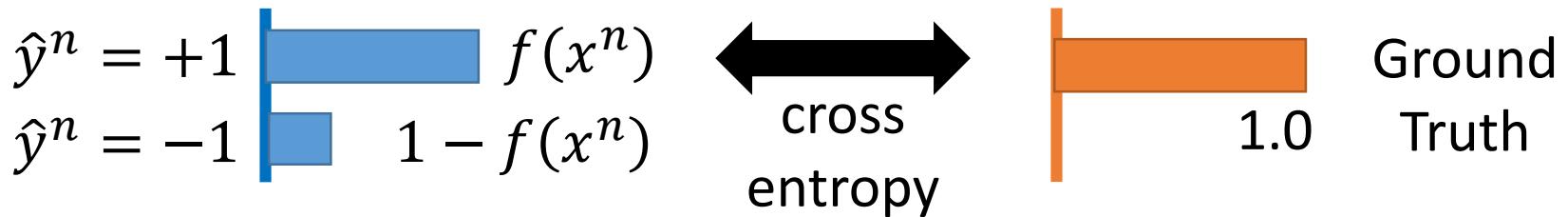
$$L(f) = \sum_n l(f(x^n), \hat{y}^n)$$

$l(\cdot)$ is the upper bound of $\delta(\cdot)$



Step 2: Loss function

$l(f(x^n), \hat{y}^n)$: cross entropy



If $\hat{y}^n = +1$:

$$\begin{aligned} l(f(x^n), \hat{y}^n) &= -\ln f(x^n) = -\ln \sigma(z^n) = -\ln \frac{1}{1 + \exp(-z^n)} \\ &= \ln(1 + \exp(-z^n)) = \underline{\ln(1 + \exp(-\hat{y}^n z^n))} \end{aligned}$$

If $\hat{y}^n = -1$:

$$\begin{aligned} l(f(x^n), \hat{y}^n) &= -\ln(1 - f(x^n)) \\ &= -\ln(1 - \sigma(x^n)) = -\ln \frac{\exp(-z^n)}{1 + \exp(-z^n)} = -\ln \frac{1}{1 + \exp(z^n)} \\ &= \ln(1 + \exp(z^n)) = \underline{\ln(1 + \exp(-\hat{y}^n z^n))} \end{aligned}$$

Step 2: Loss function

$l(f(x^n), \hat{y}^n)$: cross entropy

$$l(f(x^n), \hat{y}^n) = \ln(1 + \exp(-\hat{y}^n z^n))$$

