

# Classification: ↴ Logistic Regression

Hung-yi Lee

李宏毅

We just found out that  $P(C_i | x) = \sigma(wx+b)$ , so it means that we can use regression to implement classification.

↵  
Add sigmoid

# Step 1: Function Set

Function set: Including all different  $w$  and  $b$

$$\left\{ \begin{array}{ll} P_{w,b}(C_1|x) \geq 0.5 & \text{class 1} \\ P_{w,b}(C_1|x) < 0.5 & \text{class 2} \end{array} \right.$$

# Step 1: Function Set

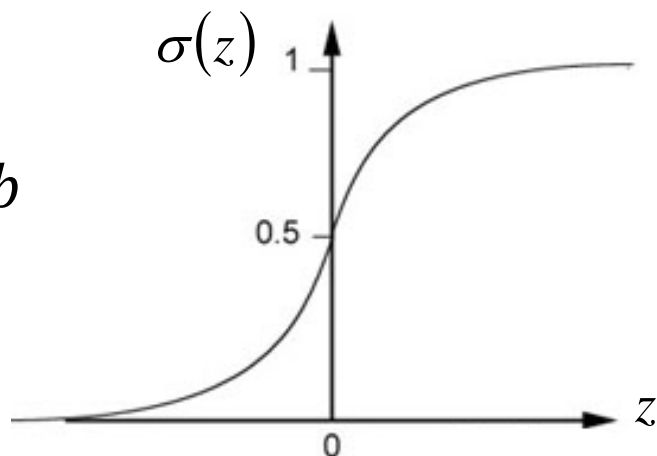
Function set: Including all different  $w$  and  $b$

$$\begin{array}{l} \sigma(z) \geq 0.5 \\ \sigma(z) < 0.5 \end{array} \left\{ \begin{array}{ll} z \geq 0 & \text{class 1} \\ z < 0 & \text{class 2} \end{array} \right.$$

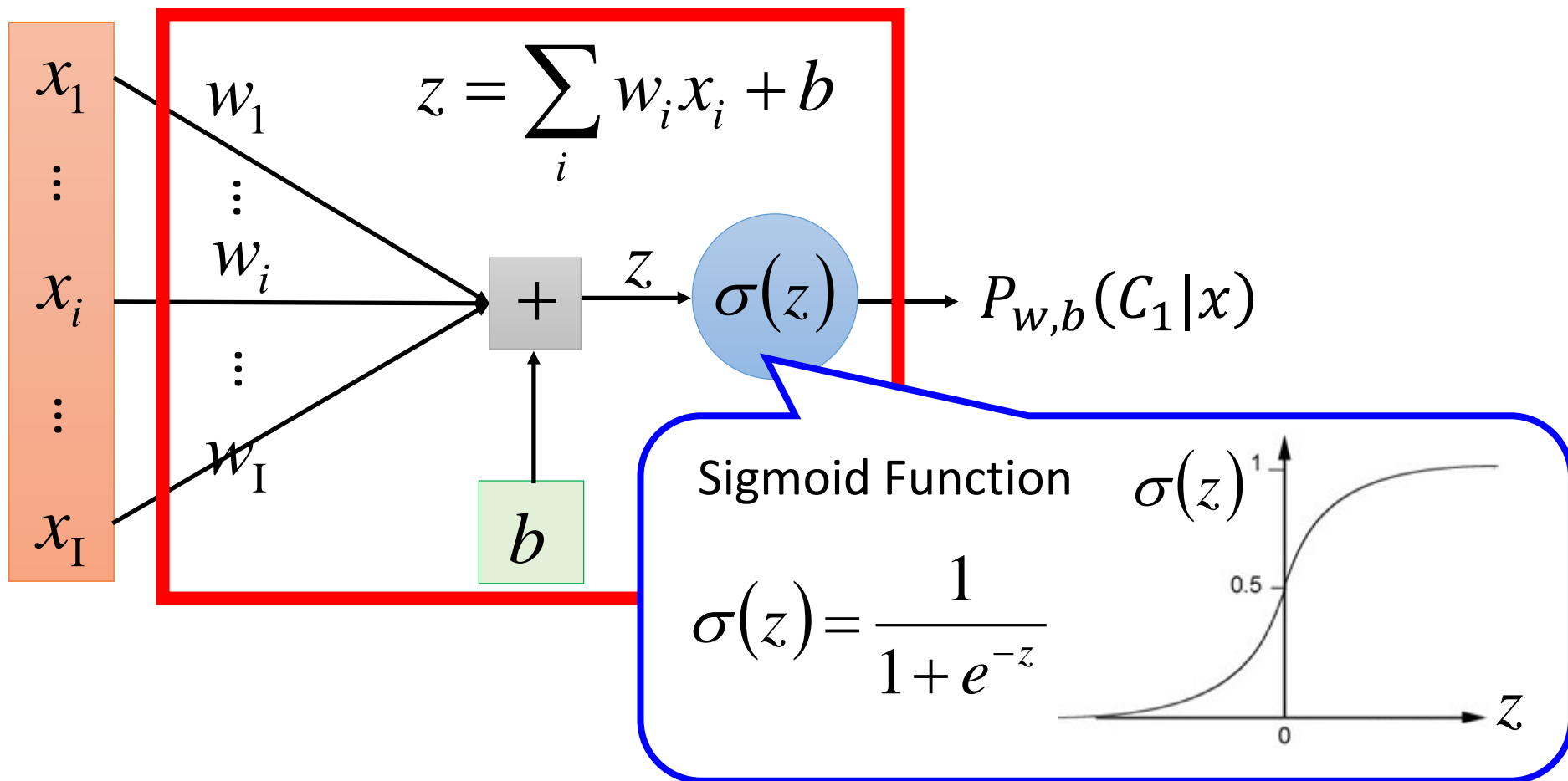
$$P_{w,b}(C_1|x) = \sigma(z)$$

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

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



# Step 1: Function Set



## **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

Logistic Regression = Linear Regression + sigmoid

Step 2:

Step 3:

# Step 2: Goodness of a Function

## Binary classification

Training  
Data

$x^1$	$x^2$	$x^3$	...	$x^N$
<u><math>C_1</math></u>	<u><math>C_1</math></u>	<u><math>C_2</math></u>	...	<u><math>C_1</math></u>

Notice that we  
are considering  
 $C_1$  instead of  $C_2$

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

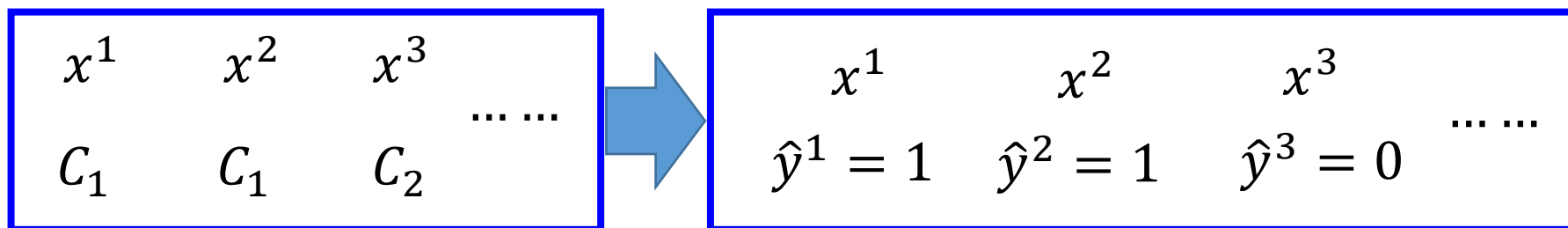
Given a set of  $w$  and  $b$ , what is its probability of generating the data?

likelihood

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

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

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



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

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

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

We want to change from "max" to "min" in order to do gradient "descent".

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

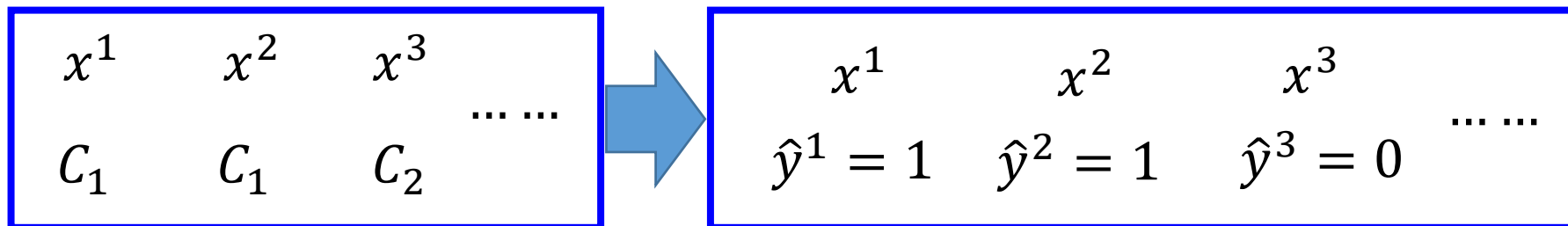
$$= -\ln f_{w,b}(x^1) \longrightarrow -[\hat{y}^1 \ln f(x^1) + (1 - \hat{y}^1) \ln(1 - f(x^1))]$$

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

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

⋮

We need to change to the unified format so that we can do the summation.



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

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

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

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



## Step 2: Goodness of a Function

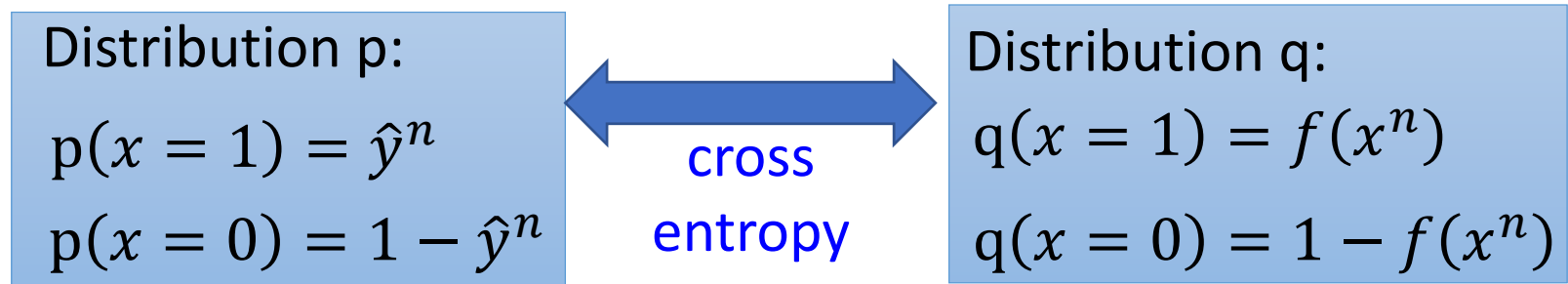
$$L(w, b) = f_{w,b}(x^1)f_{w,b}(x^2)(1 - f_{w,b}(x^3)) \cdots f_{w,b}(x^N)$$

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

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

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

Cross entropy between two Bernoulli distribution



If two distributions were identical, the cross entropy of these two distributions would be zero.

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

## Step 2: Goodness of a Function

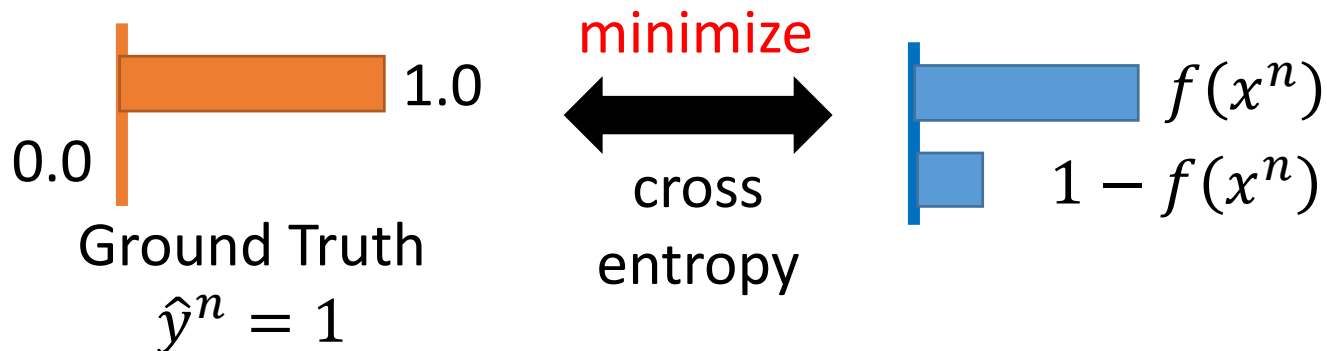
$$L(w, b) = f_{w,b}(x^1)f_{w,b}(x^2)(1 - f_{w,b}(x^3)) \cdots f_{w,b}(x^N)$$

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

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

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

Cross entropy between two Bernoulli distribution



## 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

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

## 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

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

Why don't we simply use square error as in linear regression?

Cross entropy:

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

$\therefore \sigma(z)$  + square error is not suitable for gradient descent.

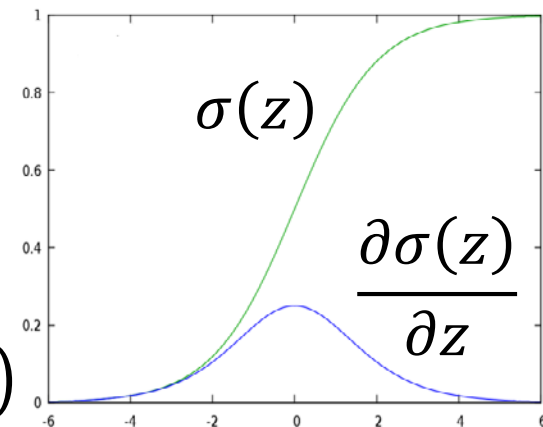
P.16

# Step 3: Find the best function

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

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

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



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

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

## Step 3: Find the best function

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

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

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

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

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

# Step 3: Find the best function

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

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

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

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

Larger difference, larger update

Conclusion:

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

## 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

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

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

Step 3:

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

## 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

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

# Logistic Regression + Square Error

Why can't we use square error with  $\sigma(z)$  when we are doing gradient descent?

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

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

Step 3:

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

$$= 2(f_{w,b}(x) - \hat{y}) \boxed{f_{w,b}(x) (1 - f_{w,b}(x))} \boxed{x_i}$$

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

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

Normally, when we are far from the target, the gradient should be very large.



## 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:

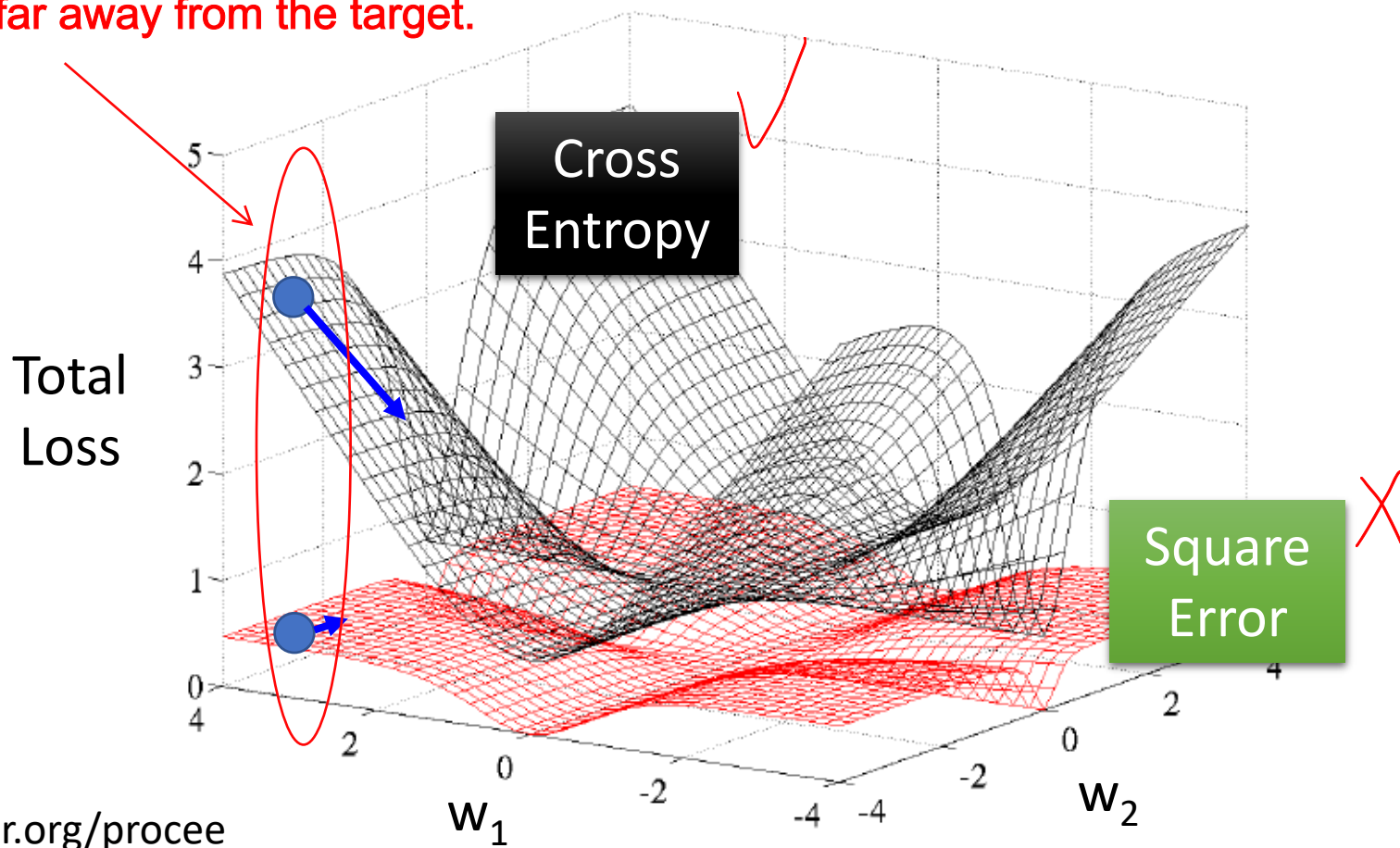
$$\frac{\partial (f_{w,b}(x) - \hat{y})^2}{\partial w_i} = 2(f_{w,b}(x) - \hat{y}) \frac{\partial f_{w,b}(x)}{\partial z} \frac{\partial z}{\partial w_i}$$
$$= 2(f_{w,b}(x) - \hat{y}) f_{w,b}(x) (1 - f_{w,b}(x)) x_i$$

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

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

# Cross Entropy v.s. Square Error

When we are far away from the target.



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

logistic regression

probabilistic model

# Discriminative v.s. Generative

without any assumptions

with some assumptions

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

Gradient descent



directly find  $w$  and  $b$



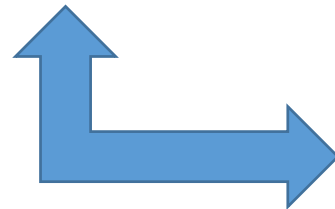
MLE (Gaussian distribution)

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

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

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

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



Will we obtain the same set of  $w$  and  $b$ ? **No**

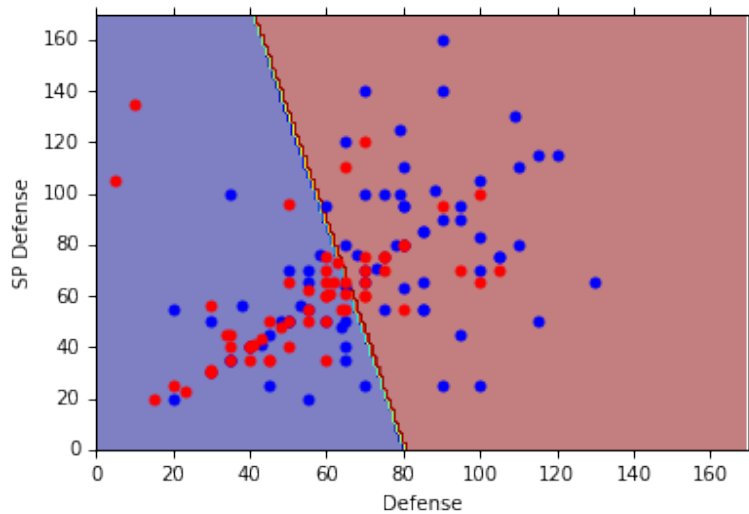
The same model (function set), but different function may be selected by the same training data. **different parameters ( $w, b$ )**

## Example 1:

# Generative v.s. Discriminative

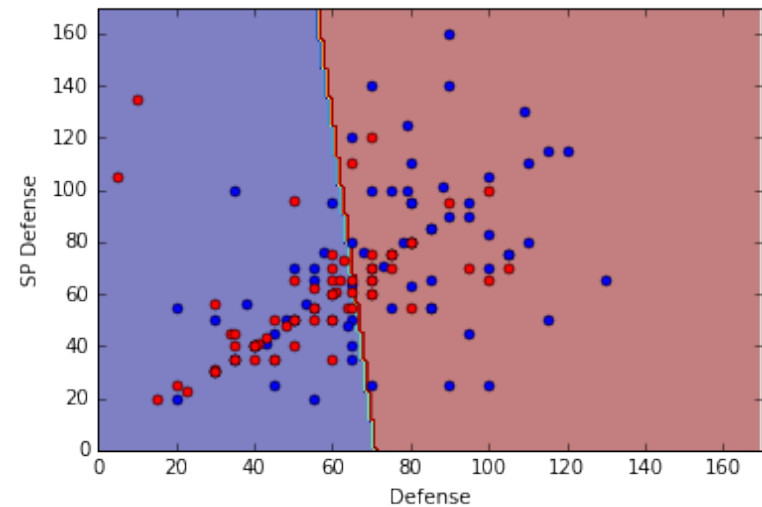
Normally, discriminative should have better performance.

## Generative



73% accuracy

## Discriminative



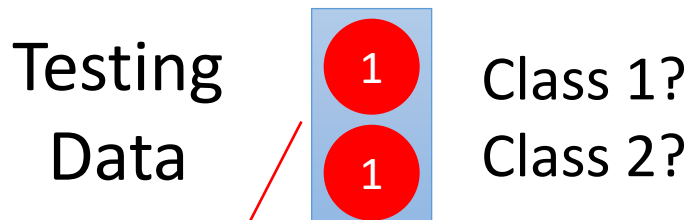
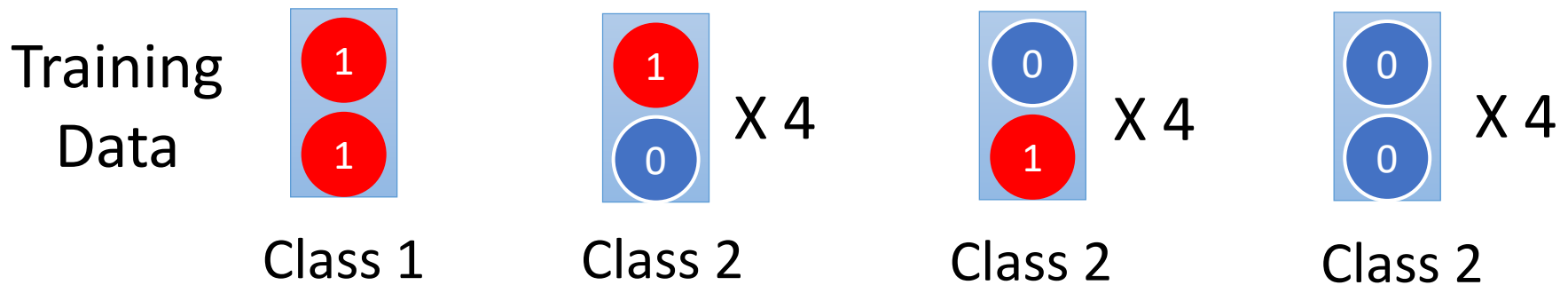
79% accuracy

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

Example 2: (This case exemplifies the disadvantage of the generative model.)

# Generative v.s. Discriminative

- Example



The assumption is that every attribute is independent.

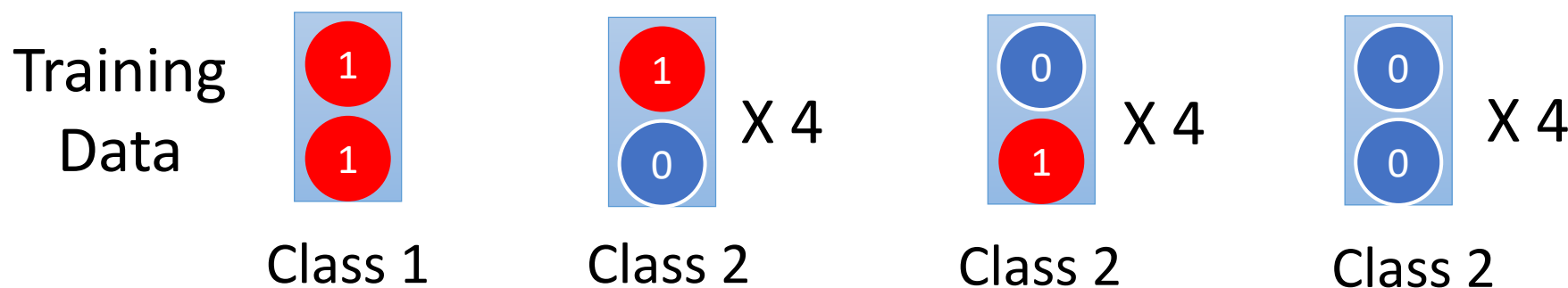
How about Naïve Bayes?

$$P(x|C_i) = P(x_1|C_i)P(x_2|C_i)$$

It's very clear that we should classify this testing data into Class 1.

# Generative v.s. Discriminative

- Example



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

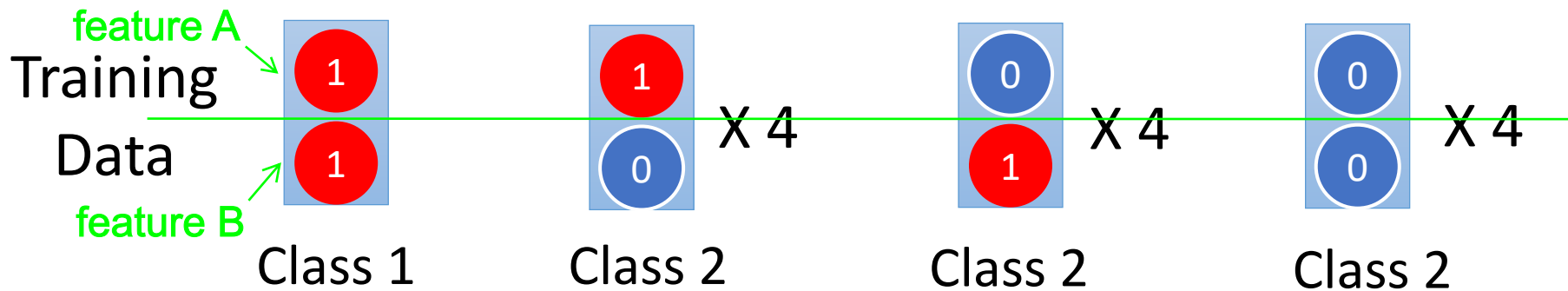
$$P(x_1 = 1|C_1) = 1$$

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

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

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

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



In Class 2, both feature A and feature B have appeared 1, and because feature A and feature B are independent, the Naïve Bayes model thinks that it's possible to appear 1 on both feature A and feature B. Also, the amount of Class 2's data is much bigger than Class 1, so the model classifies the testing data into Class 2.

Testing Data

$P(C_1|x) < 0.5$

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

Diagram illustrating the calculation of the posterior probability  $P(C_1|x)$  for the testing data (feature A=1, feature B=1). The numerator is  $1 \times 1 \times \frac{1}{13}$ . The denominator is  $1 \times 1 \times \frac{1}{13} + \frac{1}{3} \times \frac{1}{3} \times \frac{12}{13}$ .

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

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

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

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

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

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

mind-picturing

# Generative v.s. Discriminative

- 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
  - Priors and class-dependent probabilities can be estimated from different sources.

class independent

ex: In speech recognition, we can collect a lot of text files and compute its probability. It is independent with the speech files.



# Multi-class Classification (3 classes as example)

Final output of

$$C_1: w^1, b_1 \quad z_1 = w^1 \cdot x + b_1$$

$$C_2: w^2, b_2 \quad z_2 = w^2 \cdot x + b_2$$

$$C_3: w^3, b_3 \quad z_3 = w^3 \cdot x + b_3$$

Gaussian distribution with shared covariance matrix. (Skip the proof here.)

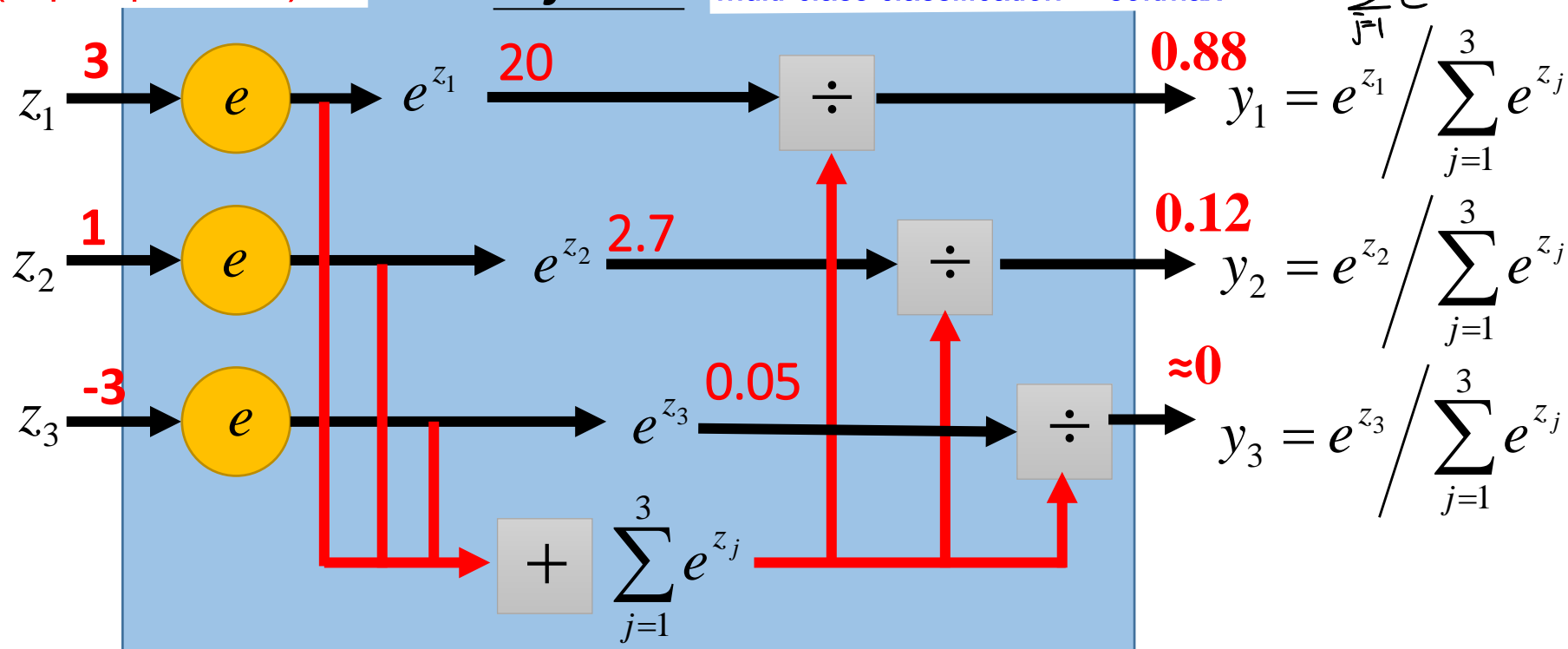
**Probability:**

- $1 > y_i > 0$
- $\sum_i y_i = 1$

If the number of class is two in multi-class classification, we will find that the softmax is equal to sigmoid in logistic regression. (and  $w_1, w_2$  can merge into  $w$ )

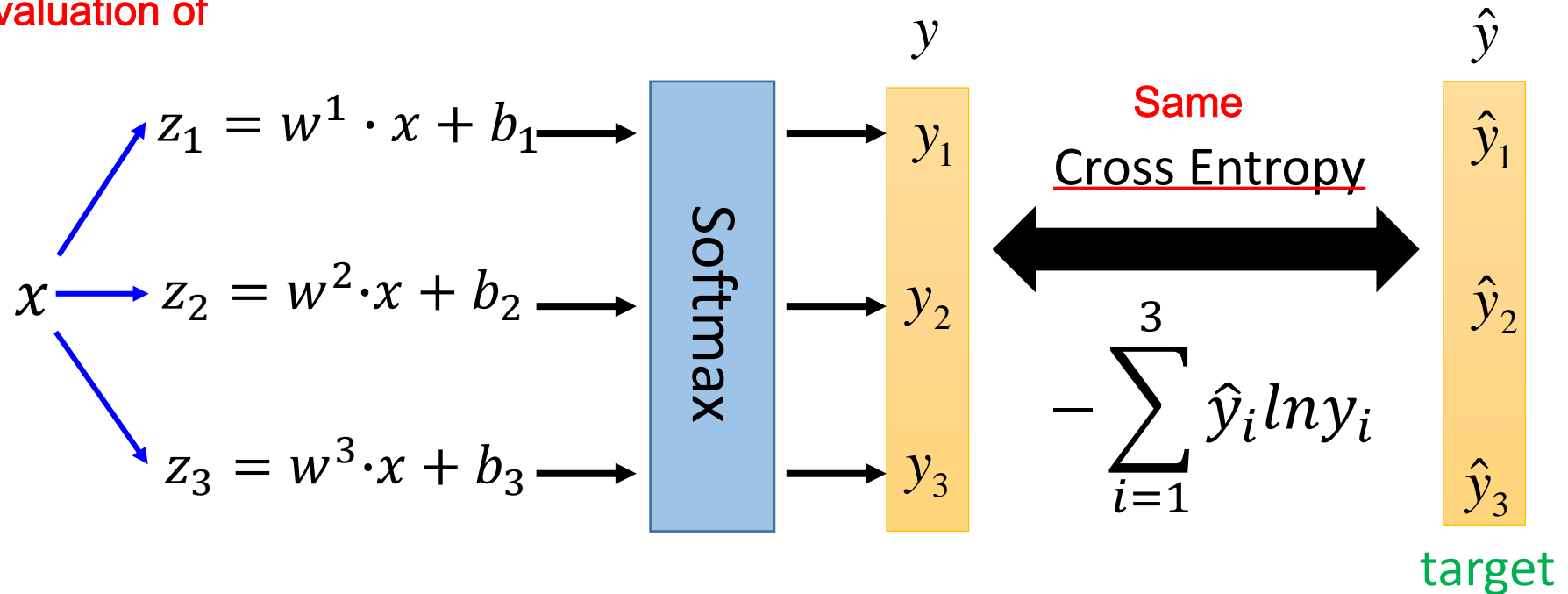
Normalize the final result:  
binary classification  $\Rightarrow$  sigmoid  
multi-class classification  $\Rightarrow$  softmax

**Softmax**



# Multi-class Classification (3 classes as example)

Evaluation of



If  $x \in \text{class 1}$

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

$$-\ln y_1$$

If  $x \in \text{class 2}$

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

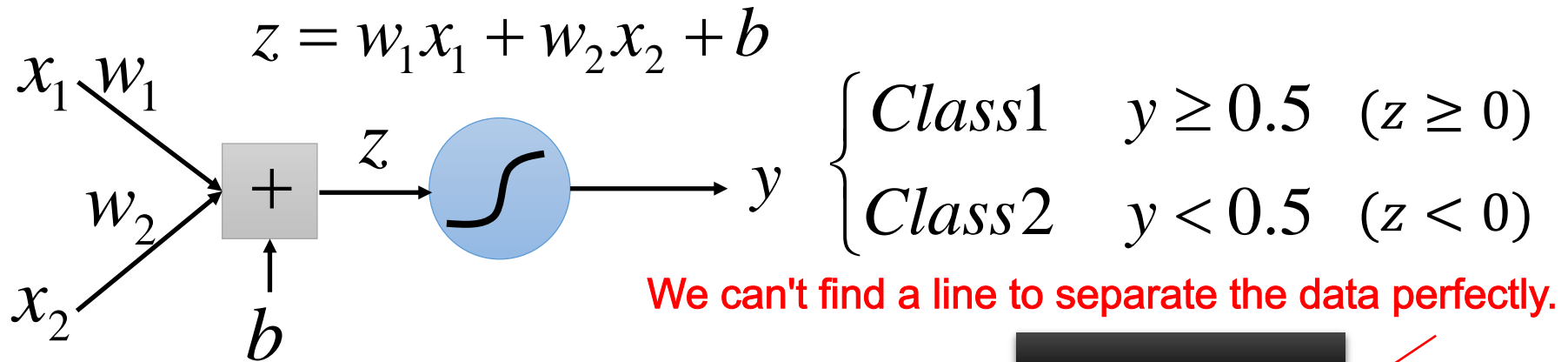
$$-\ln y_2$$

If  $x \in \text{class 3}$

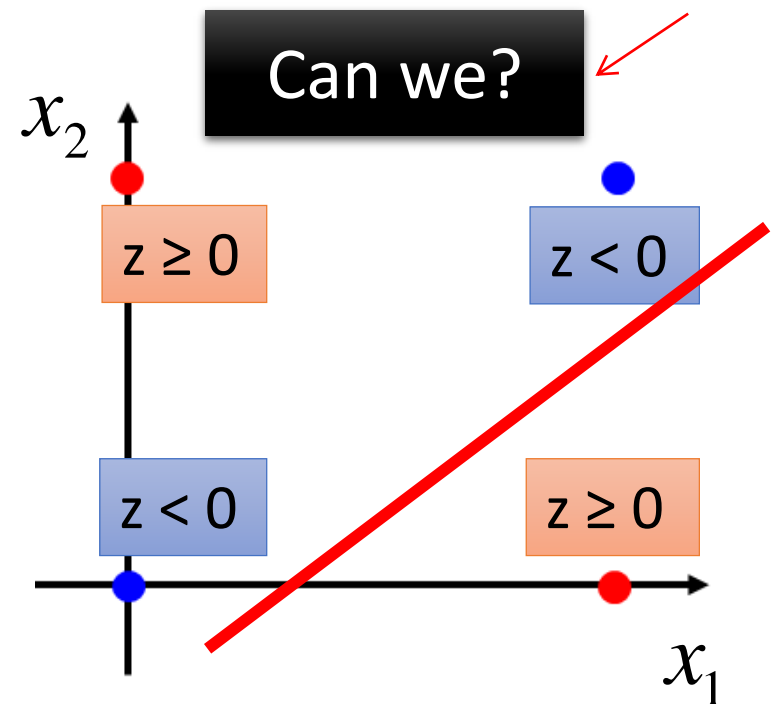
$$\hat{y} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$-\ln y_3$$

# Limitation of Logistic Regression



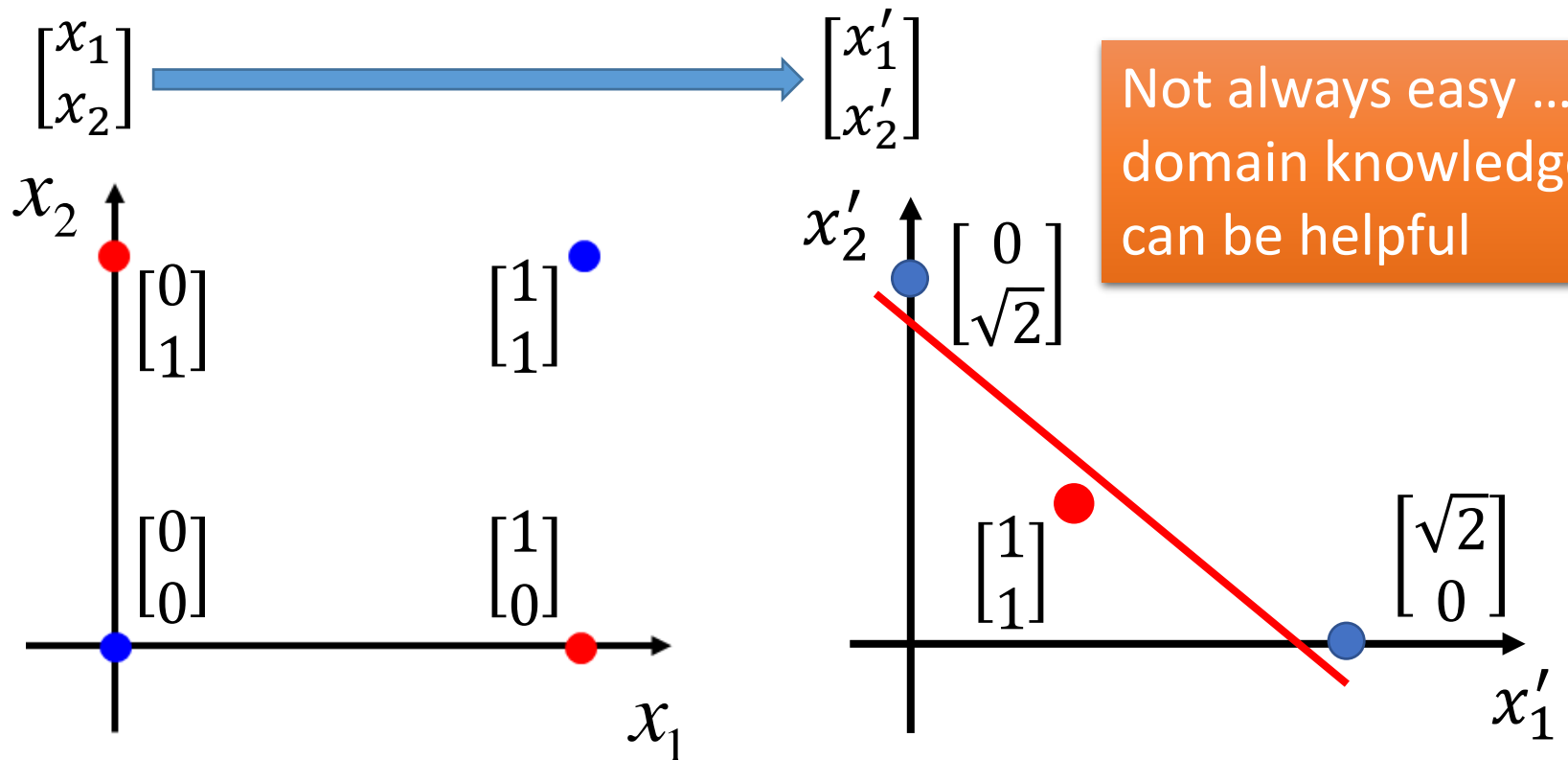
Input Feature		Label
$x_1$	$x_2$	
0	0	Class 2
0	1	Class 1
1	0	Class 1
1	1	Class 2



# Limitation of Logistic Regression

- **Feature transformation**

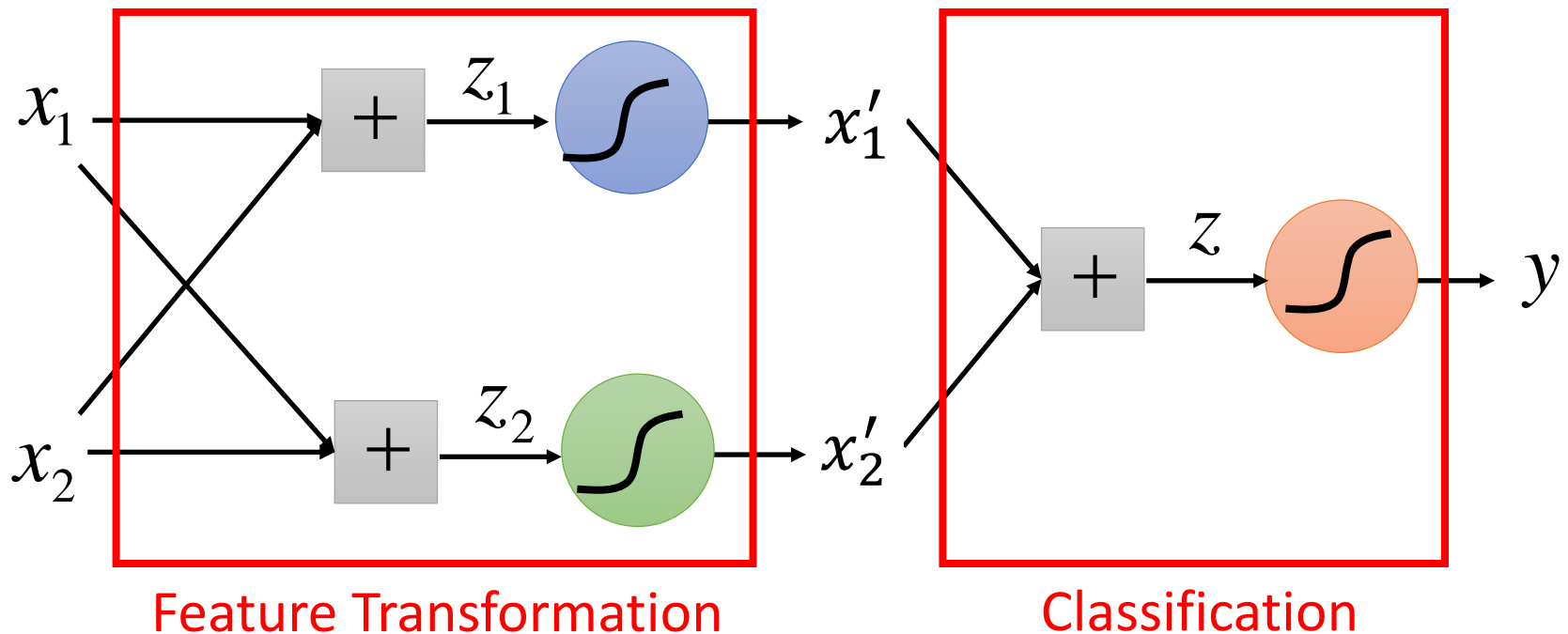
$x'_1$ : distance to  $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$   
 $x'_2$ : distance to  $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$



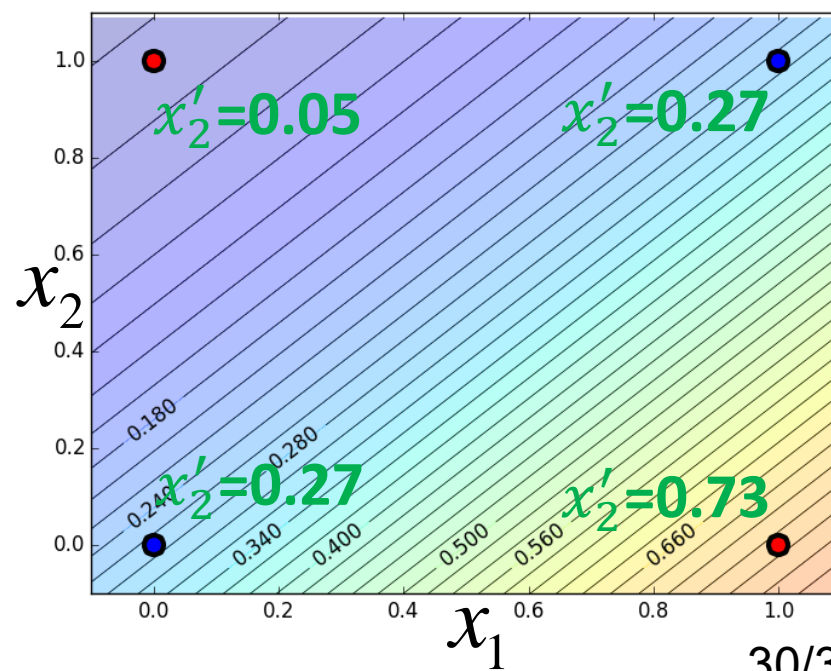
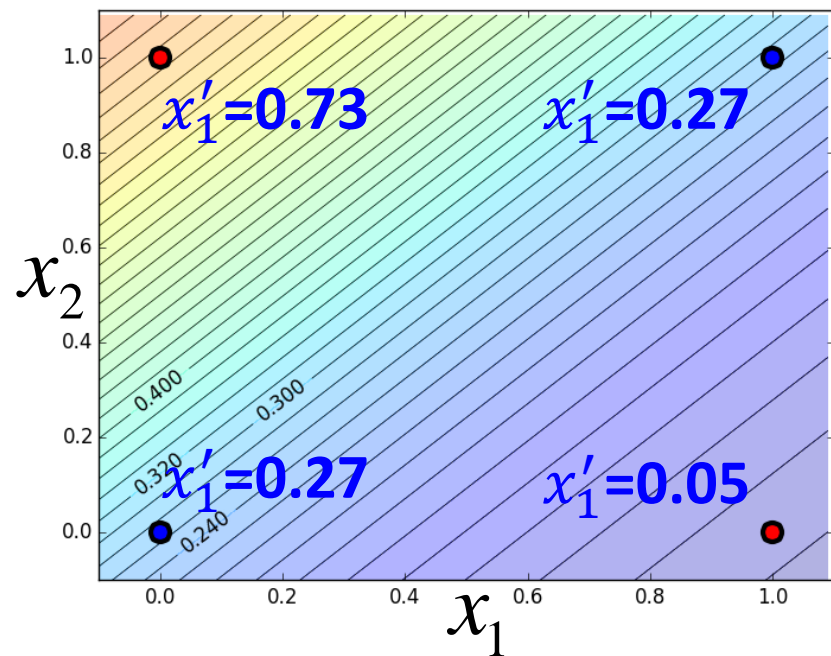
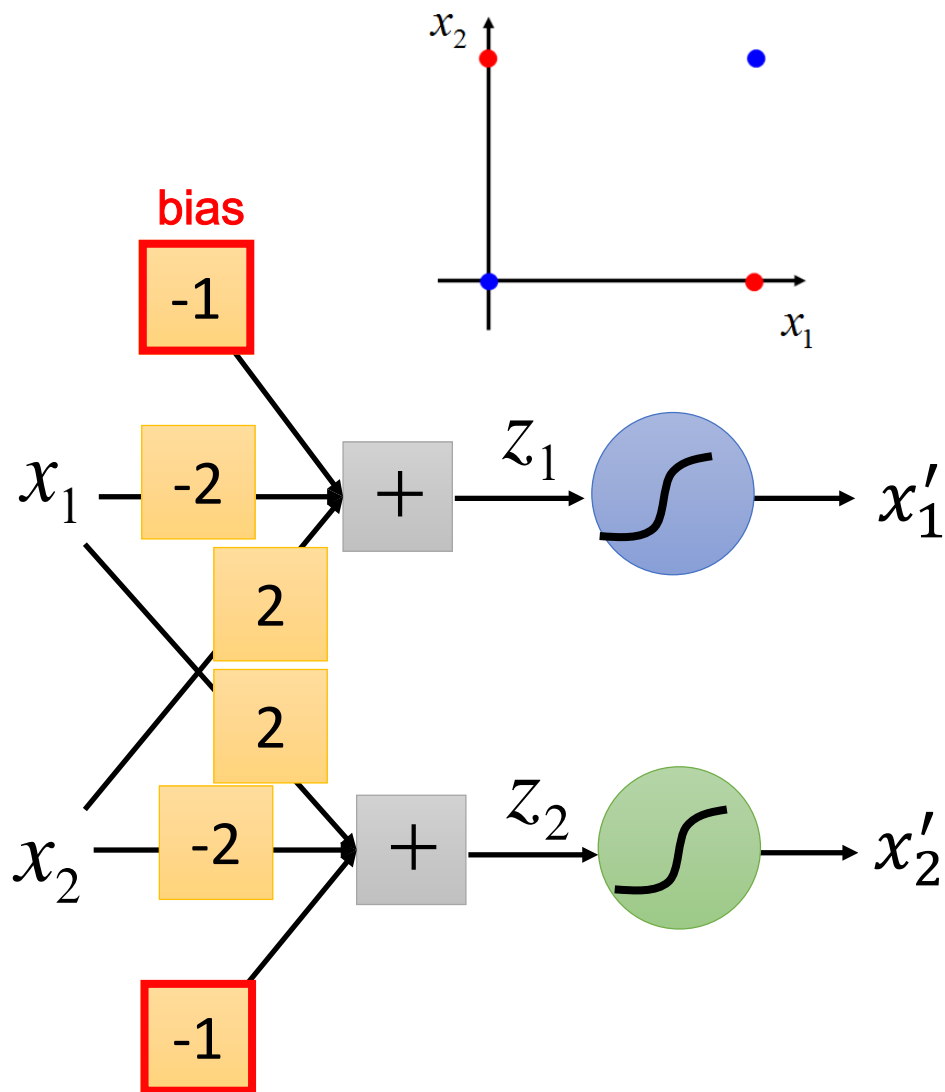
# Limitation of Logistic Regression

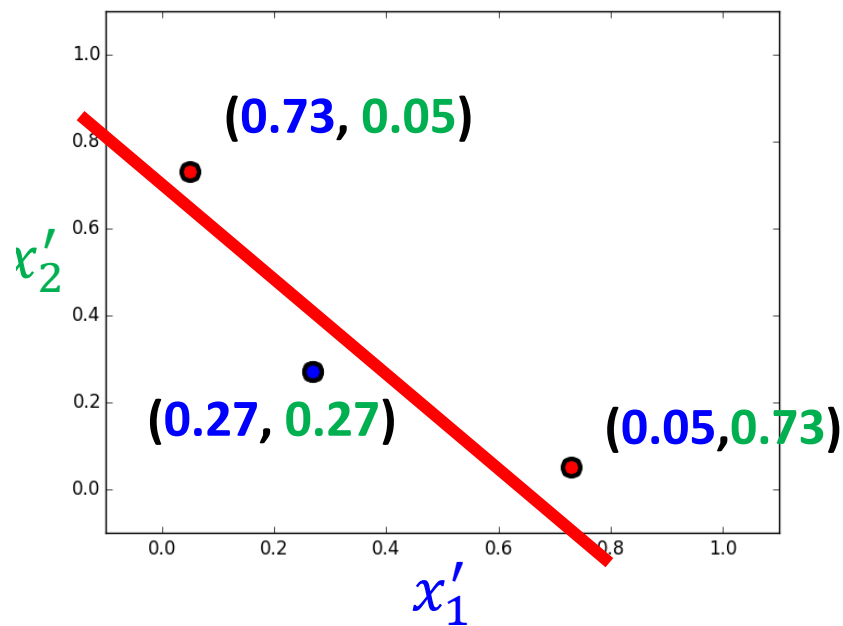
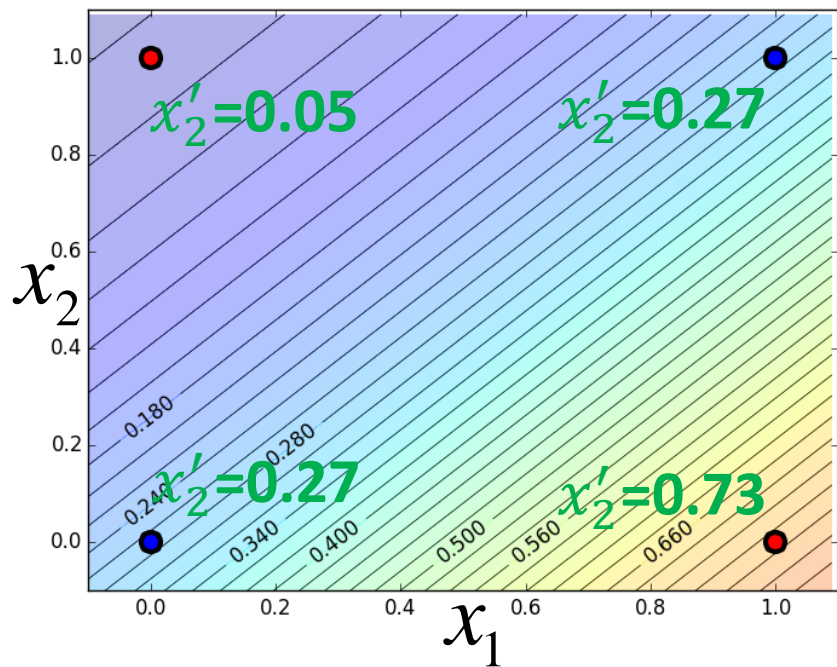
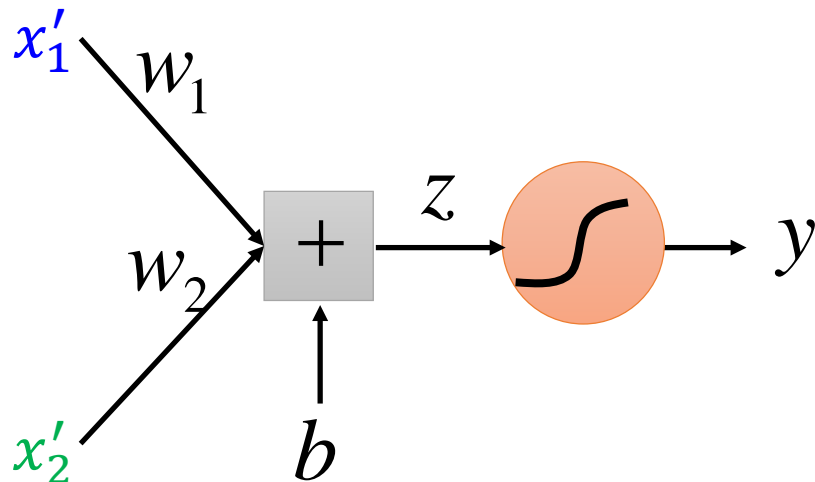
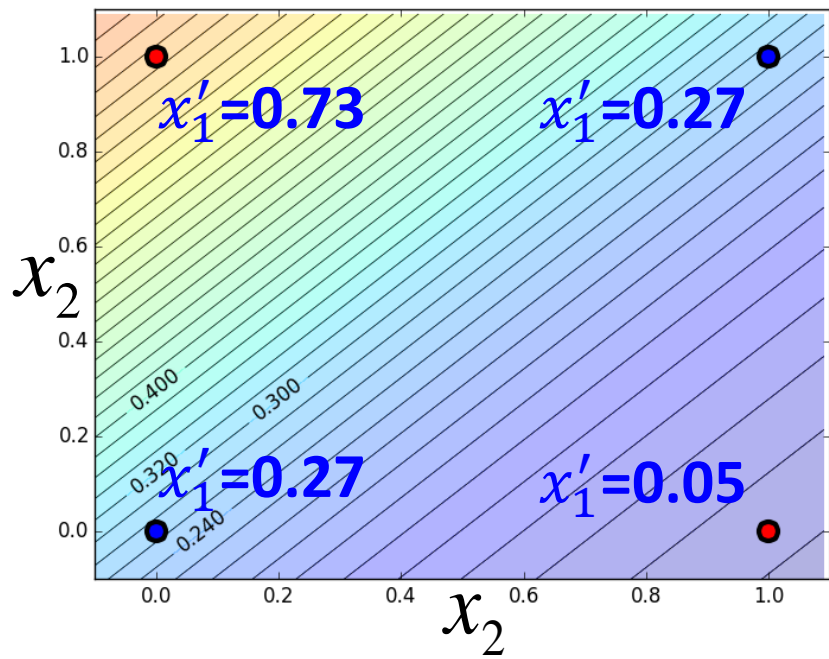
We use logistic regression to do feature transformation and classification respectively.

- Cascading logistic regression models



We ignore the bias here to simplify the illustration.

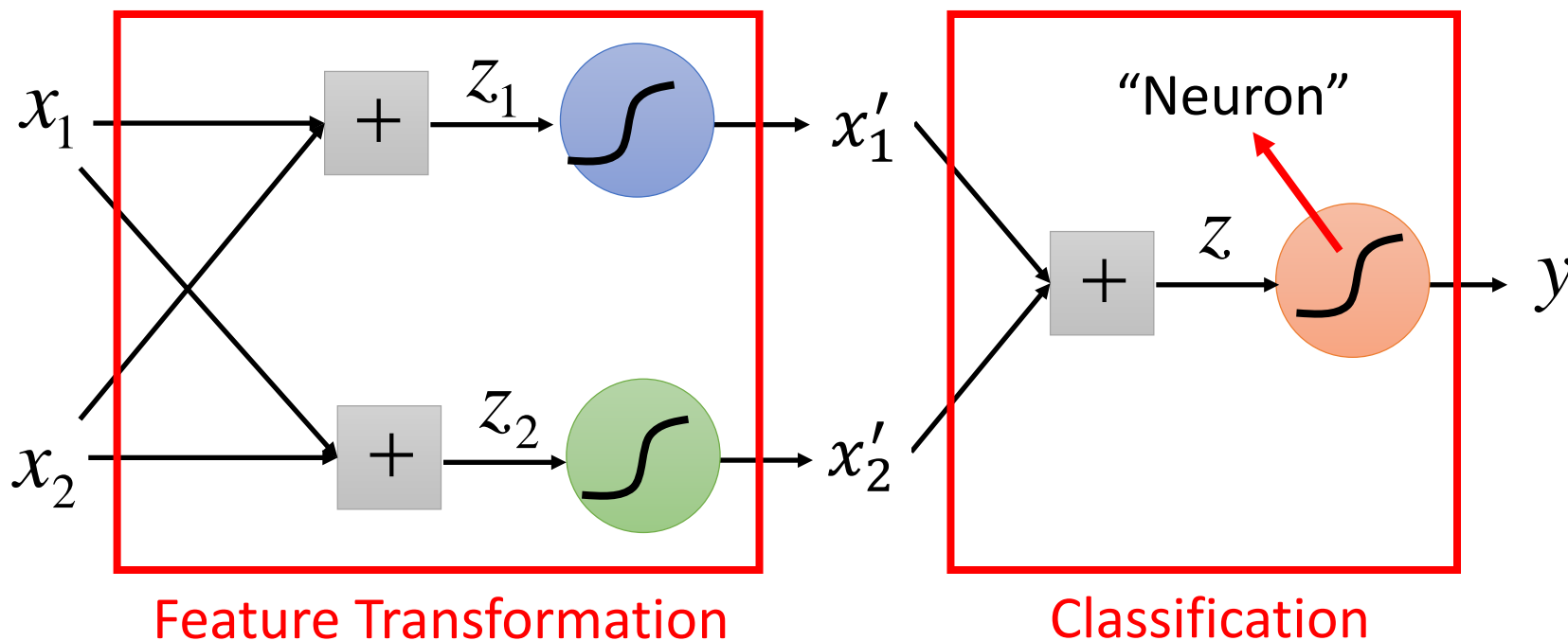
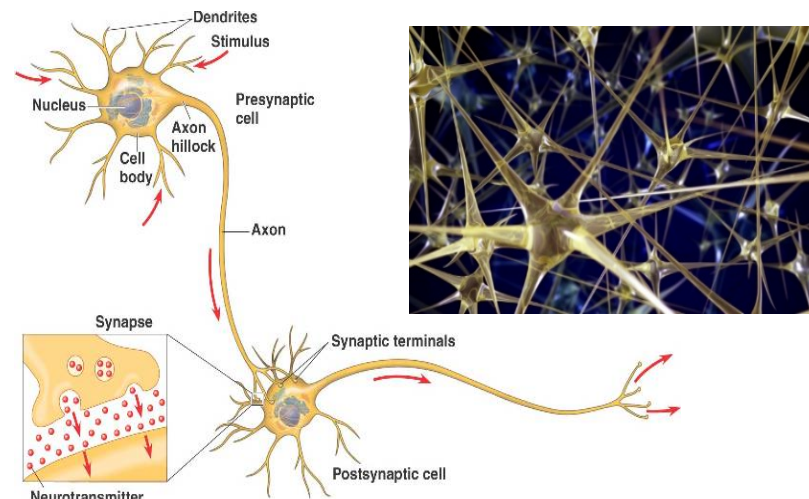




# Deep Learning!

All the parameters of the logistic regressions are jointly learned.

using gradient descent



Neural Network



# Reference

- Bishop: Chapter 4.3

# Appendix

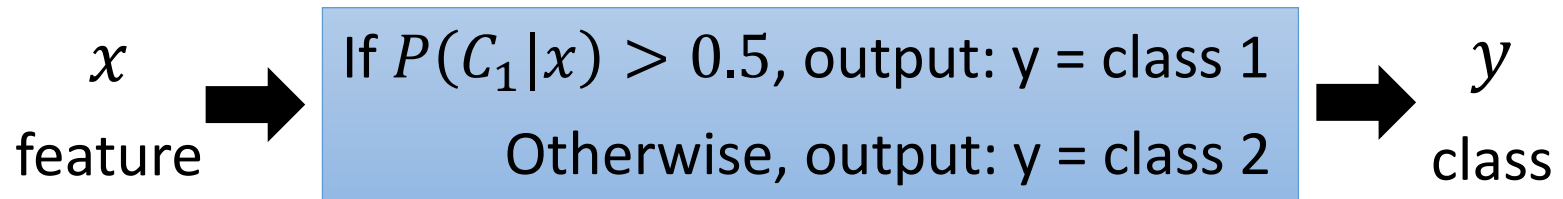
# Three Steps

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

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

of logistic regression

- Step 1. Function Set (Model)



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

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

ideal

- Step 2. Goodness of a function **Couldn't use gradient descent**

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

- Step 3. Find the best function: gradient descent

## Step 2: Loss function

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

Ideal loss:

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

0 or 1

Approximation:

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

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

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

Bad

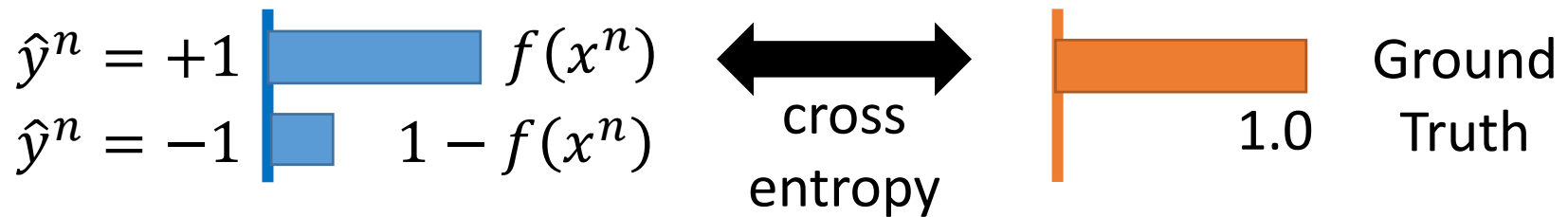
Good

$\hat{y}^n z^n$

They should have the same sign.  
(Larger value, smaller loss)

## Step 2: Loss function

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



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

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

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

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

## Step 2: Loss function

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

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

