

Logistic Regression

① Step 1 Function set

Function set: Including all different w and b

$$\begin{cases} z > 0 & \text{class 1} \\ z < 0 & \text{class 2} \end{cases}$$

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

$$z = w^T x + b = \sum w_i x_i + b$$

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

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 $f_{w,b}(x) = P_{w,b}(C_1|x)$

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

The most likely w^* and b^* is the one with the largest $L(w, b)$: $w^*, b^* = \arg \max_{w, b} L(w, b)$

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

$$\hat{y} = 1 \text{ class 1 ; } 0 \text{ class 2}$$

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

$$\begin{aligned} -\ln L(w, b) &= -\ln f_{w,b}(x^1) - \ln f_{w,b}(x^2) - \ln [1 - f_{w,b}(x^3)] \cdots \\ &= \sum [-\hat{y}^n \ln f_{w,b}(x^n) + (1 - \hat{y}^n) \ln [1 - f_{w,b}(x^n)]] \end{aligned}$$

Cross entropy between two Bernoulli distribution

Distribution p:		Distribution q:
$p(x = 1) = \hat{y}^n$	↔	$q(x = 1) = f(x^n)$
$p(x = 0) = 1 - \hat{y}^n$	cross entropy	$q(x = 0) = 1 - f(x^n)$
$H(p, q) = - \sum_x p(x) \ln(q(x))$		
minimize		

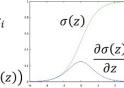
$$f_{w,b}(x) = \delta(z) = \frac{1}{1 + e^{-z}}$$

$$z = w^T x + b = \sum w_i x_i + b$$

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

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

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



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

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

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

Step 3 Find the best function

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

Larger difference, larger update
 $w_i \leftarrow w_i - \eta \sum_n -(\hat{y}^n - f_{w,b}(x^n)) x_i^n$

f_w(x^n) 越离目标越远，参数W更新的越大

thinking: 为什么用 square error 计算 Loss

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

更新 W

$\hat{y}^n = 1$ 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) $\rightarrow \partial L / \partial w_i = 0$

$\hat{y}^n = 0$ 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) $\rightarrow \partial L / \partial w_i = 0$

Logistic Regression and Linear Regression

<u>Logistic Regression</u>	<u>Linear Regression</u>
Step 1: $f_{w,b}(x) = \sigma \left(\sum_i w_i x_i + b \right)$	$f_{w,b}(x) = \sum_i w_i x_i + b$
Output: between 0 and 1	Output: any value
Training data: (x^n, \hat{y}^n)	Training data: (x^n, \hat{y}^n)
Step 2: $\hat{y}^n: 1$ for class 1, 0 for class 2	\hat{y}^n : a real number
$L(f) = \sum_n l(f(x^n), \hat{y}^n)$	$L(f) = \frac{1}{2} \sum_n (f(x^n) - \hat{y}^n)^2$
Step 3: Logistic regression: $w_i \leftarrow w_i - \eta \sum_n -(\hat{y}^n - f_{w,b}(x^n)) x_i^n$	Linear regression: $w_i \leftarrow w_i - \eta \sum_n -(\hat{y}^n - f_{w,b}(x^n)) x_i^n$

② Discriminative vs Generative

$$P(C|x) = \frac{1}{Z} (\omega x + b)$$

Generative

directly find w and b

Will we obtain the same set of w and b ?

No

Discriminative

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

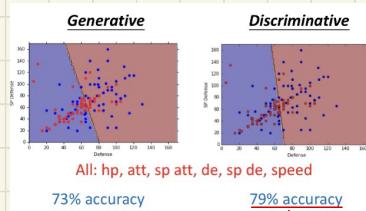
$$w^T = (\mu^1 - \mu^2)^T \Sigma^{-1}$$

$$b = -\frac{1}{2}(\mu^1)^T (\Sigma^{-1})^{-1} \mu^1 + \frac{1}{2}(\mu^2)^T (\Sigma^{-1})^{-1} \mu^2 + \ln \frac{n_1}{n_2}$$

Generative: 假設後, 實際上是怎樣的 eg: 特徵同時

Discriminative: 計算 w, b 沒有假設 eg: KNN, SVM, 決策樹

Generative vs Discriminative



better why?

example:

Training Data	Class 1	Class 2	Class 2	Class 2
	1 1	1 0 0 X 4	0 1 X 4	0 0 X 4
$P(C_1) = \frac{1}{13}$	$P(x_1 = 1 C_1) = 1$	$P(x_2 = 1 C_1) = 1$		
$P(C_2) = \frac{12}{13}$	$P(x_1 = 1 C_2) = \frac{1}{3}$	$P(x_2 = 1 C_2) = \frac{1}{3}$		

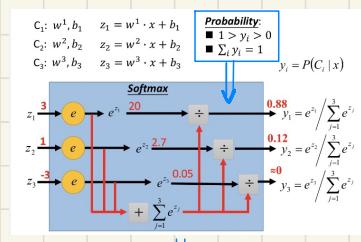
$$\begin{aligned} P(C_1|x) &= \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)} \\ &= \frac{\frac{1}{13} \times 1}{\frac{1}{13} \times 1 + \frac{1}{3} \times \frac{12}{13}} \\ &= \frac{1}{13} \end{aligned}$$

Benefit

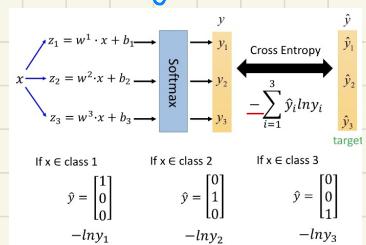
- Usually people believe discriminative model is better
- Benefit of generative model

- 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.

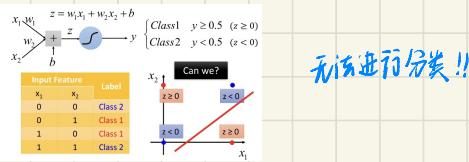
③ Multi-class classification



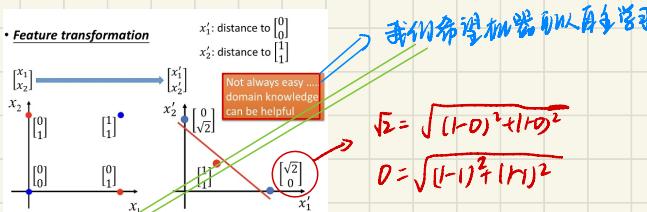
softmax function: 归一化指数函数，二分类函数 sigmoid 在多分类上的推广。目的是将多分类结果以概率形式展示。



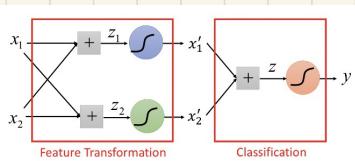
④ Limitation of logistic Regression



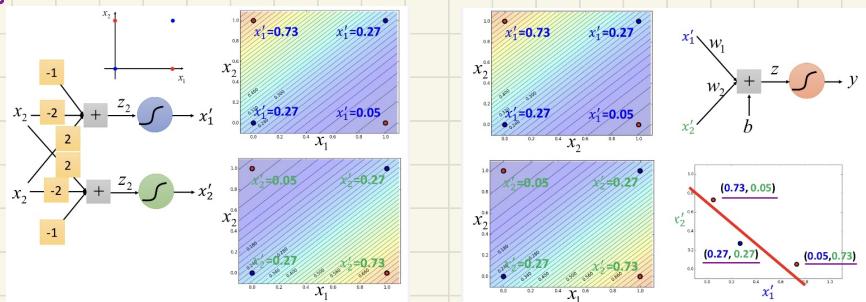
feature transformation



cascading logistic regression models

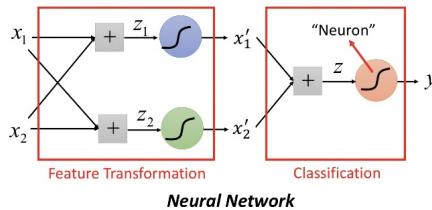
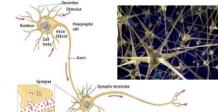


eg:



Deep Learning!

All the parameters of the logistic regressions are jointly learned.



31 Logistic Regression