

Linear Classification Model:

Apply discrete function on linear regression model s.t. output is discrete. find a linear bound in space that best split the sample

$$g(f(x;w)) = \begin{cases} 1 & f(x;w) > 0 \\ 0 & f(x;w) < 0 \end{cases}$$

$f$ : discriminant fn  
 $g$ : activation/decision fn

$$f(x;w) = w^T x + b \quad \underline{x}, w \in \mathbb{R}^{p \times 1}; \quad b \in \mathbb{R}$$

individual observation

learning criterion: loss function minimized

• Binary Classification  $y \in \{0, 1\}$

• one thought on loss function: 0-1 loss

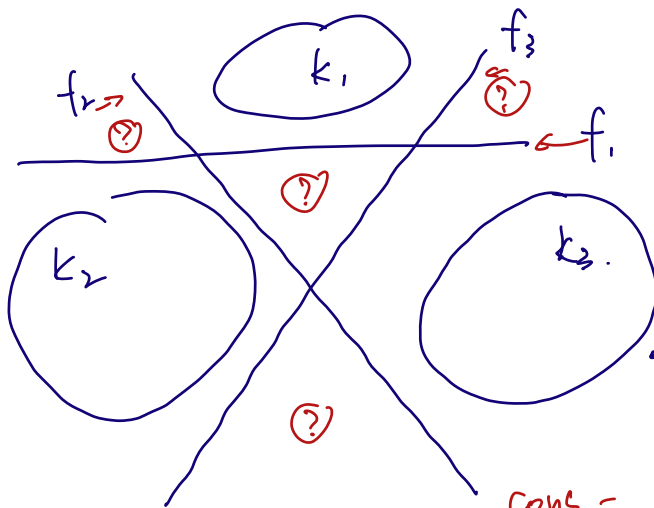
$$L_0 = \mathbb{I}(y \neq \hat{y}) = \mathbb{I}(y \neq g(f(x;w)))$$

prob is  $L_0$  is not diff. so cannot develop optimization prob on it; need smoother loss fn as criterion.

- Multi classification  $y \in \{1, \dots, C\}$

potential model:

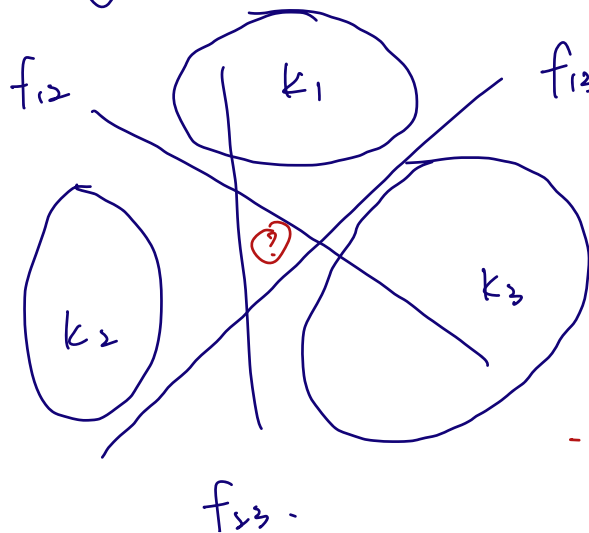
①. 一对其余.



- 3 discriminant  $f$ 's
- each  $f_k$  decide whether  $x^{(i)} \in C_k$  or not.

• then use majority voting to decide final output  
cons - (?) : indistinguishable area

②. 一对一



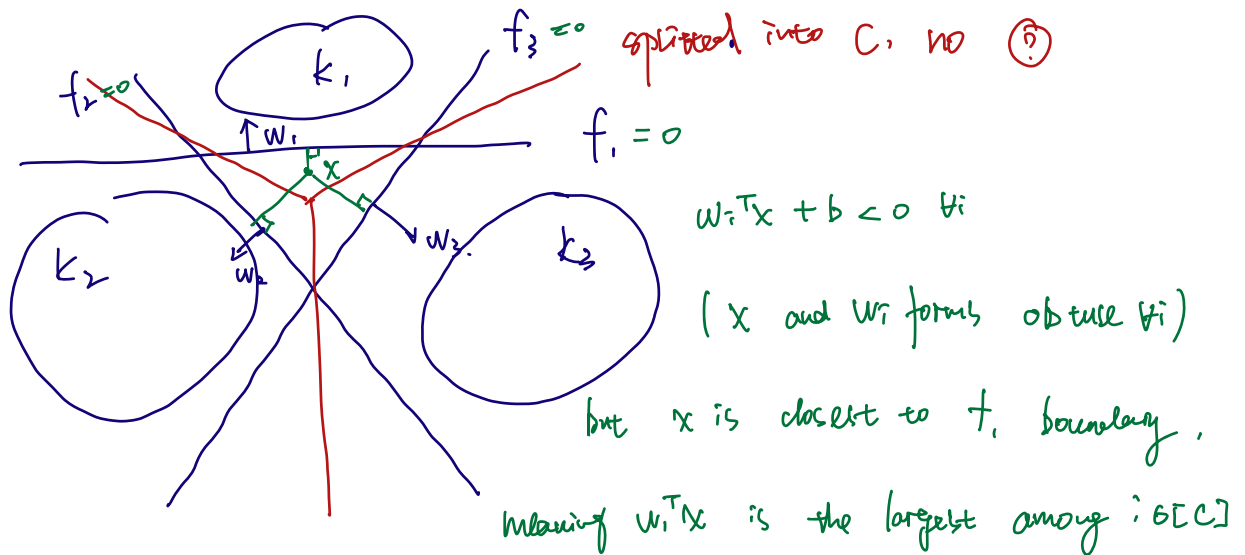
- 每两个 class 间建立一个分类器.
- total  $C^2$  classifiers.
- use majority voting

- cons: (?) (?)

② computationally costly.

② argmax.

为对“一对一策略”的一种改进.



• need  $C$   $f_i$ 's

• decision fn:  $\arg\max_c f_c(x; w) = 1$  for  $x$  in fig.

Cross-entropy: measures how similar the predicted dist

is to the real dist; equiv to MLE; the smaller the H(p, q) the closer

$$L(\theta) = \prod_{x \in X} P_\theta(X=x)^{n \cdot P_r(X=x)}$$

$$\begin{aligned} \frac{1}{n} \cdot \log(L(\theta)) &= \sum_{x \in X} P_r(X=x) \cdot \log(P_\theta(X=x)) \\ &= -H(P_r, P_\theta). \end{aligned}$$

Intuition: if  $P_r, P_\theta$  really close, then larger  $n$  will have greater power than smaller, s.t. the  $\log(L)$  is larger

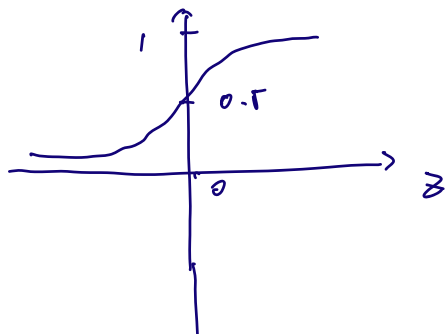
## Logistic Regression.

As said 0-1 loss not diff. cannot do gradient descent.  
So let  $g$  outputs the conditional prob. of the obs.  
being in class 1.

$$\text{i.e. } \underbrace{P_{\theta}(y=1|x)}_{\substack{\downarrow \\ \text{模型预测} = \hat{y}}} = \underbrace{g(f(x;w))}_{G[0,1]}$$

Sigmoid, logistic fn.  
 $\sigma(z)$

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



true probability for  $(x, y^*)$

$$P(y=1|x) = y^*$$

$$P(y=0|x) = 1 - y^*$$

Intuition: want to train  $w$  and  $b$

$$\text{s.t. } \hat{y} = P_{\theta}(y=1|x) = \sigma(w^T x + b) \sim y^* = P(y=1|x)$$

We cross-entropy as loss function to measure how similar two prob. distributions are -

$$H(p, q) = - [y^* \log(\hat{y}) + (1-y^*) \log(1-\hat{y})]$$

$$\Rightarrow R(w; b) = -\frac{1}{n} \sum_{i=1}^n [y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$$

$\text{argmin } R(w; b) + \lambda \|w\|^2 \stackrel{\text{equiv}}{\iff} \text{argmax log likelihood}$

$$\text{i.e. } \text{argmax}_{\sum_{i=1}^n} P(\hat{y} = y^{(i)} | x = x^{(i)})$$

intuition:  $y^* = 0$  ; want  $1-\hat{y}$  large  $\rightarrow \hat{y}$  small

$y^* = 1$  ; want  $\hat{y}$  large

Gradient:

$$\frac{\partial R}{\partial w} = -\frac{1}{n} \sum_{i=1}^n \left( y^{(i)} \frac{1}{\hat{y}^{(i)}} \frac{\partial \hat{y}^{(i)}}{\partial w} - (1-y^{(i)}) \frac{1}{1-\hat{y}^{(i)}} \frac{\partial \hat{y}^{(i)}}{\partial w} \right)$$

$$\frac{\partial \hat{y}^{(i)}}{\partial w} = \frac{\partial}{\partial w} \sigma(w^T x^{(i)} + b) = \frac{1}{(1+e^{-z})^2} \cdot e^{-z} \cdot x$$

$$= \frac{1}{1+e^{-z}} \cdot \frac{e^{-z}}{1+e^{-z}} \cdot x$$

$$= \hat{y}^{(i)} \cdot (1-\hat{y}^{(i)}) \cdot x^{(i)}$$

$$\frac{\partial R}{\partial w} = -\frac{1}{n} \sum_{i=1}^n \frac{y^{(i)}(1-\hat{y}^{(i)}) - \hat{y}^{(i)}(1-y^{(i)})}{\hat{y}^{(i)} \cdot (1-\hat{y}^{(i)})} \cdot \hat{y}^{(i)} \cdot (1-\hat{y}^{(i)}) \cdot x^{(i)}$$

$$= -\frac{1}{n} \sum_{i=1}^n \left( \underbrace{y^{(i)}}_{\in \mathbb{R}} - \underbrace{\hat{y}^{(i)}}_{P \times 1} \right) \cdot \underbrace{x^{(i)}}_{P \times 1}$$

迭代  $W_{t+1} \leftarrow W_t + \alpha \cdot \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)}) \cdot x^{(i)}$

Softmax 模型.

$$f_c(x; w_c) = w_c^T x + b_c \in \mathbb{R}, \quad c \in \{1, \dots, C\}.$$

Softmax 函数:

$$\text{softmax}(x_k) = \frac{\exp\{x_k\}}{\sum_{i=1}^C \exp\{x_i\}}.$$

all classes involved in

In this classification model:

$$\hat{y} = P_{\theta}(y=k|x) = \text{softmax}(w_k^T x) = \frac{\exp\{w_k^T x\}}{\sum_{i=1}^C \exp\{w_i^T x\}}$$

向量表示:  $\hat{y} = \frac{\exp(W^T x)}{\mathbf{1}_C^T \exp(W^T x)}$

C ↑ classifier coeffs stacked together

$$W = \begin{bmatrix} | & | & & | \\ w_1 & w_2 & \dots & w_C \\ | & | & & | \end{bmatrix}$$

$P \times C$

eg.  $\begin{bmatrix} 0.3 \\ 0.2 \\ 0.5 \end{bmatrix} \leftarrow P_{\theta}(y=1|x)$   
 $\leftarrow P_{\theta}(y=3|x)$

Loss function (for one 样本) is same as logistic:

$$- \sum_{y=1}^C P_{\theta}(y|x) \log(P_{\theta}(y|x))$$

$$P_{\theta}(y|x) = \begin{cases} 1 & \text{if } y = y^* \\ 0 & \text{otherwise.} \end{cases}$$

intuition: find  $W$  that maximize.  
 $P_{\theta}(y=y^*|x)$

Criterion, cost fn: cross-entropy

$$R(w; b) = -\frac{1}{n} \sum_{i=1}^n (y^{(n)})^T \log(\hat{y}^{(n)})$$

$y^{(n)}$  is a one hot vector

$$y^{(n)} = \begin{bmatrix} \mathbb{1}_{y^*=1} \\ \vdots \\ \mathbb{1}_{y^*=c} \end{bmatrix}$$

Gradient:

$$\frac{\partial R(w; b)}{\partial w} = -\frac{1}{n} \sum_{i=1}^n \underbrace{x^{(n)}}_{p \times 1} \underbrace{(y^{(n)} - \hat{y}^{(n)})^T}_{c \times 1}$$

$p \times c$

$w: p \times c$

