

# Linear Discrimination:

Multiple classes ( $K > 2$ )

$$\mathcal{X} = \{(x_i, y_i)\}_{i=1}^N$$

$x_i \in \mathbb{R}^D$   $y_i \in \{1, 2, \dots, K\}$

reference class  $\leftarrow$

$$\log \left[ \frac{p(x|y=c)}{p(x|y=K)} \right] = \underline{w_c^T \cdot x + w_{c0}}$$

$$\exp \left[ \log \left[ \frac{p(y=c|x)}{p(y=K|x)} \right] \right] = \log \left[ \frac{p(x|y=c) P(y=c)}{p(x|y=K) P(y=K)} \right]$$

is not a function of  $x$

$$= \log \left[ \frac{p(x|y=c)}{p(x|y=K)} \right] + \log \left[ \frac{P(y=c)}{P(y=K)} \right]$$

$\underline{w_c^T x + w_{c0}}$

$\Downarrow$

$$\exp \left[ w_c^T \cdot x + w_{c0} \right]$$

$$w_{c0} = \dot{w}_{c0} + \log \left[ \frac{P(y=c)}{P(y=K)} \right]$$

$$\frac{P(y=c|x)}{P(y=K|x)} = \exp[w_c^T x + w_{c0}]$$

$$P(y=1|x) + P(y=2|x) + \dots + P(y=K|x) = 1$$

$$P(y=1|x) + P(y=2|x) + \dots + P(y=K-1|x) = 1 - P(y=K|x)$$

$$\sum_{c=1}^{K-1} \frac{P(y=c|x)}{P(y=K|x)} = \frac{1 - P(y=K|x)}{P(y=K|x)} = \sum_{c=1}^{K-1} \exp[w_c^T x + w_{c0}]$$

$$P(y=K|x) = ? \Rightarrow \frac{1}{P(y=K|x)} = 1 + \sum_{c=1}^{K-1} \exp[w_c^T x + w_{c0}]$$

$$\underbrace{P(y=K|x)}_{\text{class } K} = \frac{1}{1 + \sum_{c=1}^{K-1} \exp[w_c^T x + w_{c0}]}$$

$$\frac{P(y=c|x)}{P(y=K|x)} = \exp[w_c^T x + w_{c0}]$$

$$\underbrace{P(y=c|x)}_{\text{classes } 1, 2, \dots, K-1} = \frac{\exp(w_c^T x + w_{c0})}{1 + \sum_{d=1}^{K-1} \exp[w_d^T x + w_{d0}]}$$

$$\theta = \left\{ \underbrace{\frac{D \times 1}{w_1}, \frac{1 \times 1}{w_{10}}}_{\text{Dx1, 1x1}}, \underbrace{\frac{D \times 1}{w_2}, \frac{1 \times 1}{w_{20}}}, \dots, \underbrace{\frac{D \times 1}{w_{(k-1)}}, \frac{1 \times 1}{w_{(k-1)0}}} \right\}$$

$$\text{total \# of parameters} = (K-1)(D+1)$$

$$\begin{array}{l} \text{all classes except} \\ \text{the reference class} \end{array} \left[ \begin{array}{l} P(y=1|x) = \frac{\exp[w_1^T x + w_{10}]}{1 + \exp[w_1^T x + w_{10}] + \dots + \exp[w_{(k-1)}^T x + w_{(k-1)0}]} \\ \vdots \\ P(y=k-1|x) = \frac{\exp[w_{(k-1)}^T x + w_{(k-1)0}]}{1 + \exp[w_1^T x + w_{10}] + \dots + \exp[w_{(k-1)}^T x + w_{(k-1)0}]} \end{array} \right]$$

$$\begin{array}{l} \text{reference} \\ \text{class} \end{array} \left[ P(y=k|x) = \frac{1}{1 + \exp[w_1^T x + w_{10}] + \dots + \exp[w_{(k-1)}^T x + w_{(k-1)0}]} \right]$$

$$P(y=c|x) = \frac{\exp[w_c^T \cdot x + w_{c0}]}{\sum_{d=1}^K \exp[w_d^T \cdot x + w_{d0}]}$$

$\left. \begin{array}{l} \geq 0 \\ \leq 1 \end{array} \right\} \sum_{c=1}^K P(y=c|x) = 1$

$x^*$  is a new data point

SOFT MAX

$$\begin{array}{lcl} w_1^T \cdot x^* + w_{10} & = & +2 \\ w_2^T \cdot x^* + w_{20} & = & -2 \\ w_3^T \cdot x^* + w_{30} & = & +1 \end{array} \left. \vphantom{\begin{array}{l} \\ \\ \end{array}} \right\} \text{pick the maximum one}$$

$$P(y=1|x) = \frac{\exp(2)}{\exp(2) + \exp(-2) + \exp(1)}$$

] maximum

$$P(y=2|x) = \frac{\exp(-2)}{\exp(2) + \exp(-2) + \exp(1)}$$

$$P(y=3|x) = \frac{\exp(1)}{\exp(2) + \exp(-2) + \exp(1)}$$

$$w_1^T x^* + w_{10} = 1000 - 3000 = -2000$$

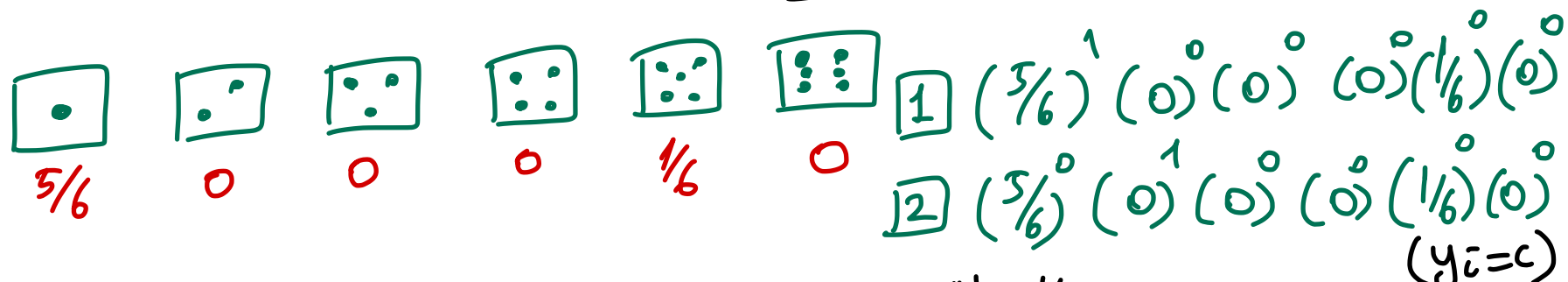
$$w_2^T x^* + w_{20} = 2000 - 3000 = -1000$$

$$w_3^T x^* + w_{30} = 3000 - 3000 = 0$$

$$\frac{\text{Inf}}{\text{Inf}} = \frac{\exp(1000) [\exp(-3000)]}{(\exp(1000) + \exp(2000) + \exp(3000)) \exp(-3000)} = P(y=1|x)$$

$$0 \approx \frac{\exp(-2000)}{\underbrace{\exp(-2000)}_0 + \underbrace{\exp(-1000)}_0 + \underbrace{\exp(0)}_1} = P(y=1|x)$$

$$y_i | x_i \sim \text{Multinomial}(y_i; \underline{1}, \sum_{c=1}^K P(y=c|x))$$



$$\text{likelihood}(\sum w_c, w_{c0} \mid \mathcal{X}) = \prod_{i=1}^N \prod_{c=1}^K P(y_i=c|x_i)$$

$$\underbrace{y_{i1} \log(\hat{y}_{i1}) + (1-y_{i1}) \log(1-\hat{y}_{i1})}_{y_{i1} \log(\hat{y}_{i1}) + y_{i2} \log(\hat{y}_{i2})} \Rightarrow Y = \begin{bmatrix} \textcircled{1} & 0 & 0 \\ 2 & 1 & 0 \\ 2 & 0 & 0 \\ 1 & 0 & 0 \\ 3 & 0 & 1 \end{bmatrix} = \prod_{i=1}^N \prod_{c=1}^K P(y_i=c|x_i)^{y_{ic}}$$

one-hot encoding

i<sup>th</sup> row, c<sup>th</sup> column of Y

$$\text{log-likelihood}(\sum w_c, w_{c0} \mid \mathcal{X}) = \sum_{i=1}^N \sum_{c=1}^K y_{ic} \log \left[ \underbrace{P(y_i=c|x_i)}_{\hat{y}_{ic}} \right]$$

$$\text{Error}(\sum w_c, w_{c0} \mid \mathcal{X}) = - \sum_{i=1}^N \sum_{c=1}^K y_{ic} \log(\hat{y}_{ic})$$

$$\text{Error}(\{w_c, w_{c0}\}_{c=1}^K | \mathcal{X}) = - \sum_{i=1}^N \sum_{c=1}^K y_{ic} \cdot \log(\hat{y}_{ic})$$

where  $\hat{y}_{ic} = \frac{\exp(w_c^T \cdot x_i + w_{c0})}{\sum_{d=1}^K \exp(w_d^T \cdot x_i + w_{d0})}$

Exercise #6

$$\frac{\partial \text{Error}}{\partial w_c} = ?$$

$$\frac{\partial \text{Error}}{\partial w_{c0}} = ?$$

$$\left. \begin{aligned} w_c^{(t+1)} &= w_c^{(t)} - \underbrace{\eta}_{\Delta w_c} \cdot \frac{\partial \text{Error}}{\partial w_c} \\ w_{c0}^{(t+1)} &= w_{c0}^{(t)} - \underbrace{\eta}_{\Delta w_{c0}} \cdot \frac{\partial \text{Error}}{\partial w_{c0}} \end{aligned} \right] \text{gradient descent updates}$$

$$\begin{aligned} \Delta w_d &= \eta \cdot \sum_{i=1}^N \sum_{c=1}^K \frac{y_{ic}}{\hat{y}_{ic}} \cdot \hat{y}_{ic} \cdot [\underbrace{\delta_{cd}}_{1(c=d)} - \hat{y}_{id}] \cdot x_i \\ &= \eta \sum_{i=1}^N (y_{id} - \hat{y}_{id}) \cdot x_i \\ \Delta w_{d0} &= \eta \sum_{i=1}^N (y_{id} - \hat{y}_{id}) \end{aligned}$$

binary  
class.

$$\begin{aligned} &\eta \sum_{i=1}^N (y_i - \hat{y}_i) \cdot x_i \\ &\eta \sum_{i=1}^N (y - \hat{y}_i) \end{aligned}$$

Hint:

$$2 \sum_{i=1}^N \sum_{c=1}^K y_{ic} (\delta_{cd} - \hat{y}_{icd}) \cdot x_i = 2 \sum_{i=1}^N (y_{id} - \hat{y}_{id}) \cdot x_i$$

$$2 \sum_{i=1}^N \sum_{c=1}^K y_{ic} (\delta_{cd} - \hat{y}_{icd}) \cdot x_i = 2 \sum_{i=1}^N (y_{id} - \hat{y}_{id}) \cdot x_i$$

$$2 \sum_{i=1}^N \sum_{c=1}^K y_{ic} (\delta_{cd} - \hat{y}_{icd}) \cdot x_i = 2 \sum_{i=1}^N (y_{id} - \hat{y}_{id}) \cdot x_i$$

ALGORITHM:  
STEP #1: initialize  $\{w_1, w_{10}, w_2, w_{20}, \dots, w_k, w_{k0}\}$  randomly  
 $\hookrightarrow$  uniform  $(-0.001, +0.001)$

ALGORITHM:  
STEP #1: initialize  $\{w_1, w_{10}, w_2, w_{20}, \dots, w_k, w_{k0}\}$  randomly  
 $\hookrightarrow$  uniform  $(-0.001, +0.001)$

STEP #2: calculate gradients

STEP #3: update  $\{w_1, w_{10}, w_2, w_{20}, \dots, w_k, w_{k0}\}$  using gradients

STEP #4: Go to Step #2 if there is enough change in the gradients.