

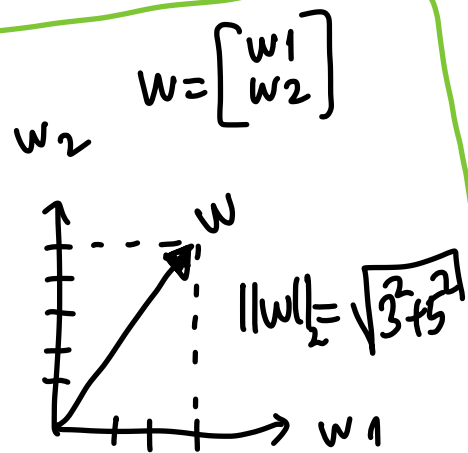
STEP #1: Initialize $\{w, w_0\}$ and decide η .
initialize them to very small values
for example $\text{Uniform}(-0.001, +0.001)$

STEP #2: Calculate Δw and Δw_0

STEP #3: Update w and w_0 using Δw and Δw_0

$$w^{(t+1)} = w^{(t)} + \Delta w^{(t)}$$
$$w_0^{(t+1)} = w_0^{(t)} + \Delta w_0^{(t)}$$

STEP #4: Go to STEP#2 if there is a change in the parameters [i.e., $\|\Delta w\|_2 \neq 0, |\Delta w_0| \neq 0$]



If $\|\Delta w\|_2 < \epsilon$ & $|\Delta w_0| < \epsilon$ where ϵ is a very small # such as 10^{-10} , we should stop the algorithm.

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

$$\exp \left[\log \left[\frac{\Pr(y=c|x)}{\Pr(y=K|x)} \right] \right] = \exp \left[\underbrace{w_c^T \cdot x + w_{c0}}_{\text{is a function of } x} \right] \quad \underbrace{\text{constant with respect to } x}$$

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

$$\text{i) } \Pr(y=c|x) \geq 0$$

$$\text{ii) } \sum_{c=1}^K \Pr(y=c|x) = 1$$

$$\Pr(y=1|x) + \Pr(y=2|x) + \dots + \Pr(y=K-1|x) + \boxed{\Pr(y=K|x)} = 1$$

$$\frac{\Pr(y=1|x) + \Pr(y=2|x) + \dots + \Pr(y=K-1|x)}{\Pr(y=K|x)} = \frac{1 - \Pr(y=K|x)}{\Pr(y=K|x)}$$

$$\sum_{c=1}^{K-1} \frac{\Pr(y=c|x)}{\Pr(y=K|x)} = \frac{1 - \Pr(y=K|x)}{\Pr(y=K|x)}$$

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

$$\left\{ \begin{array}{l} \text{reference} \\ \text{nonreference} \end{array} \right. \Pr(y=k|x) = \frac{1}{1 + \sum_{c=1}^{K-1} \exp[w_c^T x + w_{c0}]} \quad (1)$$

$$\Pr(y=c|x) = \frac{\exp[w_c^T x + w_{c0}]}{1 + \sum_{d=1}^{K-1} \exp[w_d^T x + w_{d0}]} \quad (2)$$

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

of parameters
= (K-1)(D+1)

$$\rightarrow \Pr(y=c|x) = \frac{\exp[w_c^T x + w_{c0}]}{\sum_{d=1}^K \exp[w_d^T x + w_{d0}]} \quad \left. \vphantom{\frac{\exp[w_c^T x + w_{c0}]}{\sum_{d=1}^K \exp[w_d^T x + w_{d0}]}} \right\} \text{softmax function}$$

x_{N+1} is a new data point

$$\left. \begin{array}{l} g_1(x_{N+1}) \\ g_2(x_{N+1}) \\ \vdots \\ g_K(x_{N+1}) \end{array} \right\} \text{pick the maximum one}$$

~~ideal~~
1

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

→ 0.7214
0.9999

$$\begin{aligned} g_1(x) &= +2 & 20 \\ g_2(x) &= -2 & -20 \\ g_3(x) &= +1 & 10 \end{aligned}$$

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

→ 0.0132
0.0000

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

→ 0.2654
0.0000

$\exp(20)$

$$\frac{\exp(20)}{\exp(20) + \exp(-20) + \exp(10)}$$

$$\frac{\exp(-20)}{\exp(20) + \exp(-20) + \exp(10)}$$

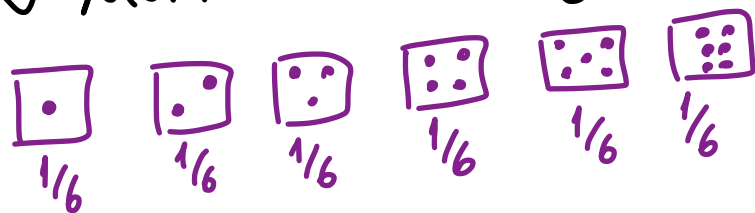
$$\frac{\exp(10)}{\exp(20) + \exp(-20) + \exp(10)}$$

$$\frac{\exp(2000)}{\exp(2000) + \exp(-2000) + \exp(1000)}$$

$$\frac{\exp(-2000)}{\exp(-2000)}$$

$$= \frac{\exp(0)}{\underbrace{\exp(0)}_1 + \underbrace{\exp(-4000)}_{\approx 0} + \underbrace{\exp(-1000)}_{\approx 0}} \approx 1$$

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



$\text{Multinomial}(y_i; 1, \{1/6, \dots, 1/6\})$

$$\Pr[\text{die}] = (1/6)^0 (1/6)^1 (1/6)^0 (1/6)^0 (1/6)^0 (1/6)^0$$

$$\Pr[\text{die}, \text{die}, \text{die}] = (1/6)^2 (1/6)^1 (1/6)^0 (1/6)^0 (1/6)^0 (1/6)^0$$

$$\text{likelihood} \left(\sum_{c=1}^K w_c, w_{c0} \middle| \mathcal{X} \right) = \prod_{i=1}^N \prod_{c=1}^K \Pr(y_i=c | x_i)$$

$y_1 = 1 \rightarrow \text{TRUE}$
 $y_1 = 2 \rightarrow \text{FALSE}$
 $y_1 = 3 \rightarrow \text{FALSE}$

$$y = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 1 \end{bmatrix}$$

one-hot encoding

$$Y = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

$$= \prod_{i=1}^N \prod_{c=1}^K \Pr(y_i=c | x_i)$$

$$\text{log-likelihood} = \sum_{i=1}^N \sum_{c=1}^K y_{ic} \log \left[\underbrace{P_r(y_i=c | x_i)}_{\hat{y}_{ic}} \right]$$

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

$$0 \cdot \log(0) + 1 \cdot \log(1) + 0 \cdot \log(0)$$

$$\frac{s_1}{s_1 + s_2 + s_3}$$

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

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

$$= - \sum_{i=1}^N \sum_{c=1}^K y_{ic} \cdot \log \left[\frac{\exp[w_c^T x_i + w_{c0}]}{\sum_{d=1}^K \exp[w_d^T x_i + w_{d0}]} \right]$$

$$\frac{\partial f(s_1, s_2, s_3)}{\partial s_1} = \frac{1 \cdot (s_1 + s_2 + s_3) - s_1 \cdot (1)}{(s_1 + s_2 + s_3)^2}$$

$$\frac{\partial f(s_1, s_2, s_3)}{\partial s_2} = \frac{0 \cdot (s_1 + s_2 + s_3) - s_1 \cdot 1}{(s_1 + s_2 + s_3)^2}$$

$$\frac{\partial f(s_1, s_2, s_3)}{\partial s_i} = \frac{1(i=1) \cdot (s_1 + s_2 + s_3) - s_1}{(s_1 + s_2 + s_3)^2}$$

$$w_c^{(t+1)} = w_c^{(t)} + \Delta w_c^{(t)} \rightarrow -\eta \cdot \frac{\partial \text{Error}}{\partial w_c}$$

$$w_{c0}^{(t+1)} = w_{c0}^{(t)} + \Delta w_{c0}^{(t)} \rightarrow -\eta \cdot \frac{\partial \text{Error}}{\partial w_{c0}}$$

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

$$= \eta \sum_{i=1}^N (y_{id} - \hat{y}_{id}) \cdot x_i$$

$$\Delta w_{d0} = \eta \sum_{i=1}^N (y_{id} - \hat{y}_{id})$$

ALGORITHM: STEP#1: Initialize $\{w_1, w_{10}, \dots, w_k, w_{k0}\}$ randomly using Uniform $(-0.001, +0.001)$

STEP#2: Calculate gradients

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

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