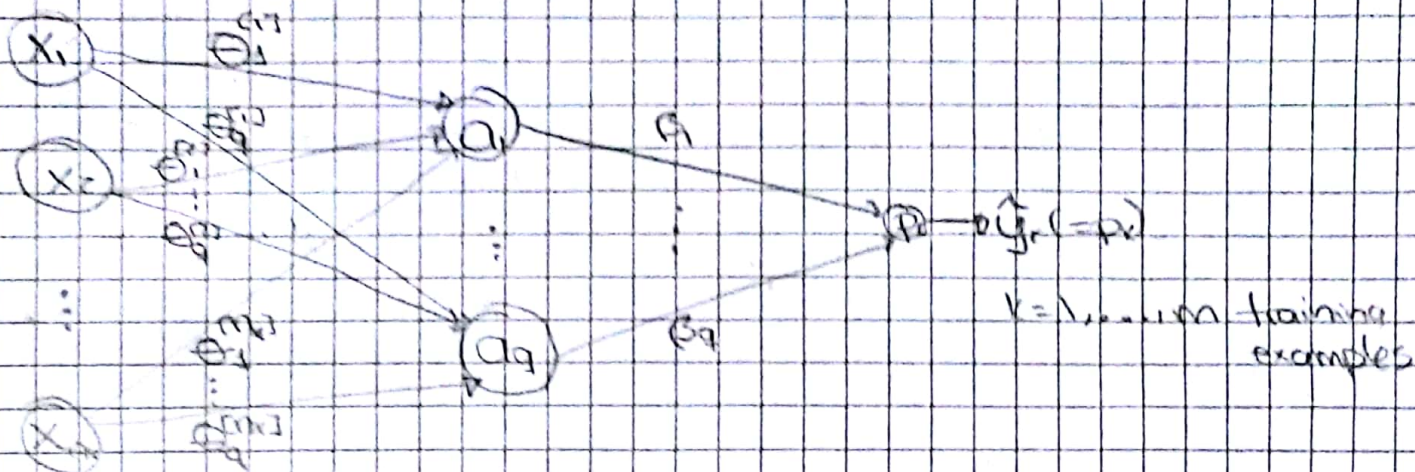


## Backprop

Simple neural network (1 hidden layer) for classification

Architecture



$$a_j = \sigma_1 \left( \sum_{i=1}^m \theta_j^{i,j} x_i \right) = \sigma_1(z_j); \sigma_1(z) = \tanh(z) \text{ or } \sigma_1(z) = \frac{1}{1+e^{-z}} \quad x_i \in \mathbb{R}^m$$

$$a_0 = x_0 = 1$$

$$p = \sigma_2 \left( \sum_{j=0}^g a_j \beta_j \right) = \sigma_2 \left( \sum_{j=0}^g \beta_j \sigma_1 \left( \sum_{i=1}^m \theta_j^{i,j} x_i \right) \right)$$

$$\text{with } \sigma_2(z) = \frac{1}{1+e^{-z}} = (1+e^{-z})^{-1} = \frac{e^z}{1+e^z}$$

$$p = \sigma_2(w) = \sigma_2 \left( \sum_{j=0}^g a_j \beta_j \right) \quad w = (y-1) \log(1-\sigma_2(w)) - y \log(\sigma_2(w))$$

$$\text{Distance loss: } l_k(\hat{y}_k, y_k) = -(y_k \log(\hat{y}_k) + (1-y_k) \log(1-\hat{y}_k)) = *$$

$$J(\theta, \alpha) = \frac{1}{m} \sum_{k=1}^m l_k(\hat{y}_k, y_k)$$

$$\text{We want } \frac{\partial J}{\partial \theta_j^{i,j}} \text{ and } \frac{\partial J}{\partial \beta_j} = l_k(\hat{y}_k, y_k) = (y-1) \log(1-p) - y \log(p)$$

$$\frac{\partial J}{\partial \theta_j^{i,j}} = \frac{\partial J}{\partial p} \cdot \frac{\partial p}{\partial w} = \frac{\partial J}{\partial p} \cdot \frac{\partial p}{\partial w} = \left[ \frac{y-1}{1-p} - \frac{y}{p} \right] \sigma_2'(w) \left[ \sum_{j=0}^g a_j \beta_j \right]$$

$$\text{but } \sigma_2'(w) = -(1+e^w)^{-2} (-e^w) = \frac{e^{-w}}{(1+e^{-w})^2} = \sigma_2(w) (1-\sigma_2(w)) = p_k(1-p_k)$$

$$\text{Then } \frac{\partial p_k}{\partial \beta_j} = \left[ \frac{1-y_k}{1-p_k} - \frac{y_k}{p_k} \right] p_k(1-p_k) \alpha_k = \frac{(p_k - y_k)(1-p_k)}{(1-p_k)(p_k)} \alpha_k = (p_k - y_k) \alpha_k$$

$$\text{And } \frac{\partial J}{\partial \beta_j} = \frac{1}{m} \sum_{k=1}^m (p_k - y_k) \alpha_k \quad \text{with } \alpha_k = \sigma_1 \left( \sum_{i=1}^m \theta_j^{i,j} x_i \right), x_i \in \mathbb{R}^m$$



$$\frac{\partial L_k}{\partial \Theta_k^{(l)}} = \frac{\partial L_k}{\partial P_k} \cdot \frac{\partial P_k}{\partial w_k} \cdot \frac{\partial w_k}{\partial z_k} \cdot \frac{\partial z_k}{\partial \Theta_k^{(l)}} \quad \text{with } \sigma_k(z_k) = \sigma_k(z_k)$$

$$= \left[ \frac{P_k - y_k}{(1 - P_k)P_k} \right] \left[ P_k(1 - P_k) \right] \left[ \beta_l \right] \left[ \sigma'_k(z_k) \right] X_{nk}$$

$$= (P_k - y_k) [\beta_l] [\sigma_k(z_k)] [(1 - \sigma_k(z_k))] X_{nk}$$

$\sigma_k(z_k) = \frac{1}{1 + e^{-z_k}}$   
 $\sigma'_k(z_k) = \sigma_k(z_k)(1 - \sigma_k(z_k))$

$$\therefore \frac{\partial L}{\partial \Theta_k^{(l)}} = \frac{1}{n} \sum_{k=1}^n \left[ (P_k - y_k) [\beta_l] [\sigma_k(z_k)] [(1 - \sigma_k(z_k))] X_{nk} \right]$$

$n = 1, 2, \dots, n$   
 $l = 1, 2, \dots, l$   
 $k = 1, 2, \dots, m$

$$Z_k = \sum_{i=0}^{r_k} \Theta_k^{(i)} X_{ik} \Rightarrow \frac{\partial Z_k}{\partial \Theta_k^{(l)}} = X_{lk}$$

$A \in \mathbb{R}^{100 \times 3}$

for  $l = 1, 2, 3$

for  $n = 1, 2, 3, 4$

$$\frac{\partial L_k}{\partial \Theta_k^{(l)}} = \underbrace{(P - y)}_{100 \times 1} \underbrace{\beta_l}_{1 \times 1} \underbrace{A_{l, l}}_{100 \times 1} \underbrace{(1 - A_{l, l})}_{100 \times 1} \underbrace{X_0}_{100 \times 1}$$

$$X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

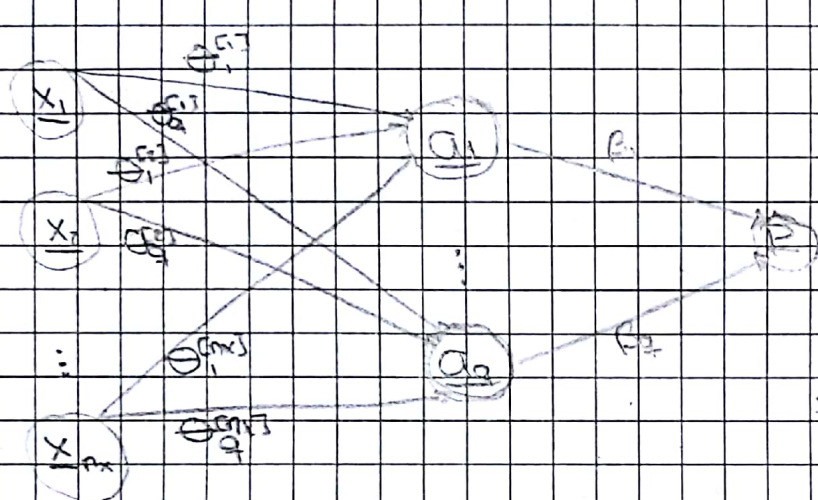
$m = \#$  training examples  
 $n = \#$  of variables

$x_i = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n$

$$\begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_n^{(1)} \\ 1 & x_1^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & \dots & x_n^{(m)} \end{bmatrix}$$

$x_i \in \mathbb{R}^n$   
 $x = \frac{1}{n}$

$$y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \in \{0, 1\}^m$$



$$\Theta = \begin{bmatrix} \theta_1^{(1)} & \theta_1^{(2)} & \dots & \theta_1^{(n)} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_q^{(1)} & \theta_q^{(2)} & \dots & \theta_q^{(n)} \end{bmatrix}$$

$$z_1^{(1)} = \sum_{i=1}^n \theta_1^{(1)} x_i^{(1)}$$

$$a_1^{(1)} = \sigma(z_1^{(1)})$$

$$\mathbb{R}^m \rightarrow \mathbb{Q}_0 = \sigma \left( \sum_{i=1}^n \theta_1^{(1)} x_i^{(1)} \right) = \sigma(z_1^{(1)})$$

$h=1, \dots, q$

with  $z_2 = \sum_{i=1}^n \theta_2^{(1)} x_i^{(1)}$   
 and  $\sigma(w) = \frac{1}{1 + e^{-w}}$

$$\mathbb{R}^m \rightarrow \hat{p} = \sigma \left( \sum_{i=1}^q \beta_i a_i^{(i)} \right) = \sigma(w)$$

$\beta_i \in \mathbb{R}$

with  $w = \sum_{i=1}^q \beta_i a_i^{(i)}$   
 and  $w = \sum_{i=1}^q \beta_i a_i^{(i)}$   
 $\hat{p} = \sigma(w)$

Loss function:  $\mathcal{L}(\Theta, \beta) = -(y \log(\hat{p}) + (1-y) \log(1-\hat{p}))$

$$\mathcal{L}(\Theta, \beta) = \text{mean}(\mathcal{L}(\Theta, \beta))$$

$$\mathcal{L}_k(\Theta, \beta) = -(y^{(k)} \log(\hat{p}_k) + (1-y^{(k)}) \log(1-\hat{p}_k))$$

$$\frac{\partial \mathcal{L}}{\partial \beta_k} = \frac{\partial \mathcal{L}}{\partial \hat{p}_k} \cdot \frac{\partial \hat{p}_k}{\partial w_k} \cdot \frac{\partial w_k}{\partial \beta_k} = \left[ \frac{1-y^{(k)}}{1-\hat{p}_k} - \frac{y^{(k)}}{\hat{p}_k} \right] \left[ \frac{\sigma(w_k)(1-\sigma(w_k))}{\hat{p}_k} \right] \left[ a_k^{(k)} \right]$$

$$= \left[ \frac{\hat{p}_k - y^{(k)}}{(1-\hat{p}_k)\hat{p}_k} \right] a_k^{(k)} = (\hat{p}_k - y^{(k)}) a_k^{(k)}$$

$$\frac{\partial \mathcal{L}}{\partial \beta} = (\hat{p} - y) a \in \mathbb{R}^m$$

$\mathbb{R}^m$



$$\frac{\partial l}{\partial \Theta_l^{(n)}} = \frac{\partial l}{\partial p} \cdot \frac{\partial p}{\partial w} \cdot \frac{\partial w}{\partial z_l^{(n)}} \cdot \frac{\partial a_l^{(k)}}{\partial z_l^{(k)}} \cdot \frac{\partial z_l^{(k)}}{\partial \Theta_l^{(n)}}$$

$$= [\hat{p} - y] [\beta_l] \left[ \frac{\sigma(z_l^{(k)})}{\sigma_l^{(k)}} (1 - \sigma(z_l^{(k)})) \right] x_n^{(k)}$$

$l = 1, \dots, q$   
 $n = 1, \dots, n$   
 $k = 1, \dots, m$

$$\frac{\partial l}{\partial \Theta_l^{(n)}} = \underbrace{(\hat{p} - y)}_{\mathbb{R}^m} \underbrace{\beta_l}_{\mathbb{R}} \underbrace{[a_l (1 - a_l)]}_{\mathbb{R}^m} \underbrace{x_n}_{\mathbb{R}^m}$$

$n \times q \times m$  elements

Let  $\text{temp} = (\hat{p} - y) \beta_l [a_l (1 - a_l)] \in \mathbb{R}^m \Rightarrow \frac{\partial l}{\partial \Theta_l^{(n)}} = \frac{\text{temp}}{\mathbb{R}^m} \cdot \frac{x_n}{\mathbb{R}^m} \in \mathbb{R}^m$

Then  $T_l = \begin{bmatrix} \text{temp} & \text{temp} & \dots & \text{temp} \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix} \in \mathbb{R}^{m \times n}$

for  $l$  in  $1, \dots, q$

$$\underbrace{T_l}_{m \times n} * \underbrace{\overline{X}}_{m \times n}$$

element-wise