20170209　　Neural　Network

— Previous Notes

- Gradient Descent

$$\theta_i \leftarrow \theta_i - \gimel \frac{\partial}{\partial \theta_i} J(\theta)$$

$$\theta \leftarrow \theta - \gimel \nabla J(\theta)$$
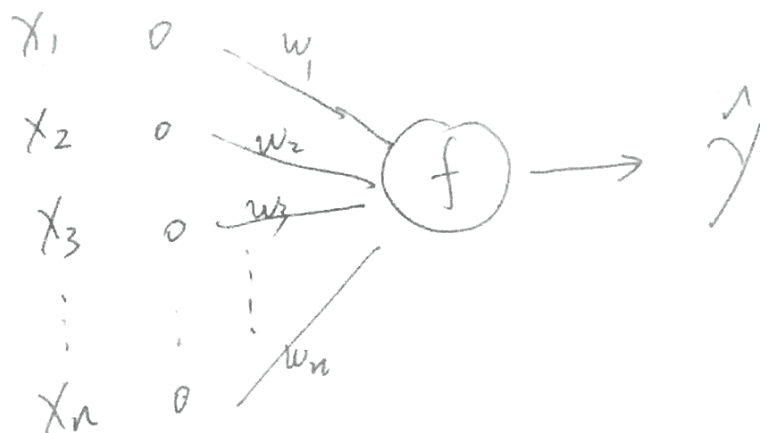
where $\nabla J(\theta) = \left( \frac{\partial}{\partial \theta_1} J(\theta), \frac{\partial}{\partial \theta_2} J(\theta), \cdots \frac{\partial}{\partial \theta_n} J(\theta) \right)$

- Mean Square Error

$$E = \frac{1}{N} \sum_{i=1}^{N} \left( y^{(i)} - \hat{y}^{(i)} \right)^2$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left( y^{(i)} - f(x^{(i)}) \right)^2$$

- Neuron



$$z = \sum_{i=1}^{n} x_i W_i$$

$$\hat{y} = f(z)$$

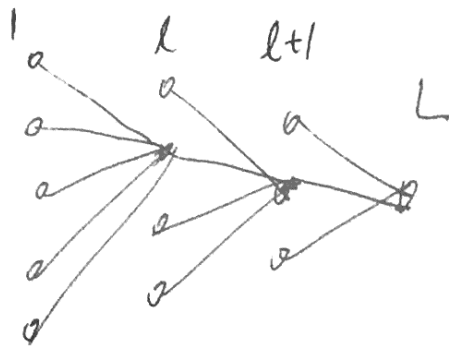$$\hat{y} = f\left( \sum_{i=1}^{n} x_i W_i \right)$$

when the activation function is sigmoid $\left( x \to \frac{1}{1+e^{-x}} \right)$

$$f'(x) = f(x)(1 - f(x))$$

because $f'(x) = \frac{d}{dx}(1+e^{-x})^{-1}$

$$= (-1)(1+e^{-x})^{-2}\frac{d}{dx}(1+e^{-x})$$

$$= (1+e^{-x})^{-2}(e^{-x})$$

$$= \frac{e^{-x}}{1+e^{-x}} \cdot \frac{1}{1+e^{-x}} = \frac{1}{1+e^{-x}} \cdot \left(1 - \frac{1}{1+e^{-x}}\right)$$

$$= f(x)(1-f(x))$$

- Neural Network



- $L$: number of layers

- $n_\ell$: number of nodes in layer $\ell$.

- $a_i^\ell$: output of $i$th node in layer $\ell$.

$$a^\ell = \begin{pmatrix} a_1^\ell \\ a_2^\ell \\ \vdots \\ a_{n_\ell}^\ell \end{pmatrix} \qquad \text{shape } n_\ell \times 1$$

- $z_i^\ell$: output of $i$th node in layer $\ell$

$$z^\ell = \begin{pmatrix} z_1^\ell \\ z_2^\ell \\ \vdots \\ z_{n_\ell}^\ell \end{pmatrix} \qquad \text{shape } n_\ell \times 1$$

- $W_{ji}^\ell$: ~~only~~ weights connecting $i$th node of layer $\ell$ and $j$th node of layer $\ell+1$)

$$W^\ell = \begin{pmatrix} w_{11}^\ell & w_{12}^\ell \cdots & w_{1n_\ell}^\ell \\ w_{21}^\ell & w_{22}^\ell \cdots & w_{2n_\ell}^\ell \\ \vdots & \vdots & \vdots \\ w_{n_{\ell+1}1}^\ell & w_{n_{\ell+1}2}^\ell \cdots & w_{n_{\ell+1}n_\ell}^\ell \end{pmatrix} \qquad \text{shape } n_{\ell+1} \times n_\ell$$

# Forward Propagation

Consider the $j$th node of layer $l+1$

$$z_j^{l+1} = \sum_{i=1}^{n_l} w_{ij}^l \, a_i^l$$

Therefore. $\quad z^{l+1} = w^l \cdot a^l$

$$a_j^{l+1} = f(z_j^{l+1})$$

we denote $f\left(\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}\right) = \begin{pmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_n) \end{pmatrix}$
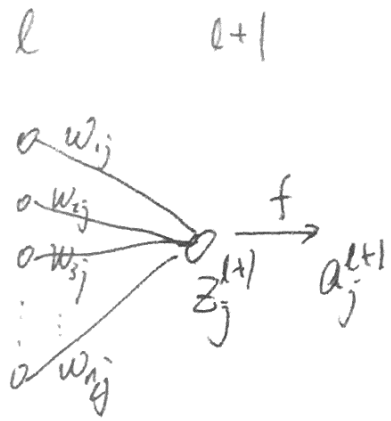
Therefore $\quad a^{l+1} = f(z^{l+1})$

we initialize $\quad a^1 = x$

then apply $\quad z^{l+1} = w^l \cdot a^l$

$$a^{l+1} = f(z^{l+1})$$

finally we have $\quad \hat{y} = a^L$

Backward Propagation.

$$E(w^1, w^2 \dots w^L; x) = \frac{1}{2} \| y - \hat{y} \|^2.$$

we compute $\nabla_{w^1} E, \nabla_{w^2} E, \dots \nabla_{w^L} E$

i.e. $\frac{\partial}{\partial w_{ji}^l} E$ for every $i, j, l$.

$$\frac{\partial}{\partial w_{ji}^l} E = \frac{\partial E}{\partial z_j^{l+1}} \frac{\partial z_j^{l+1}}{\partial w_{ji}} \qquad \text{(chain rule)}$$

Define $\delta_j^l = \frac{\partial E}{\partial z_j^l}$ $\qquad \delta^l = \begin{pmatrix} \frac{\partial E}{\partial z_1^l} \\ \frac{\partial E}{\partial z_2^l} \\ \vdots \\ \frac{\partial E}{\partial z_{d_l}^l} \end{pmatrix}$

while $\frac{\partial}{\partial w_{ji}} z_j^{l+1} = \frac{\partial}{\partial w_{ji}} \sum_{k=1}^{n} w_{jk}^l a_k^l = a_i^l$

$$\frac{\partial}{\partial w_{ji}^l} E = \delta_j^{l+1} \cdot a_i^l$$

Then we have in matrix form.

$$\nabla_{w^l} E = \delta^{l+1} \cdot (a^l)^T$$

We only need to compute $\delta^\ell$ now

if $L = \ell$, i.e. the output layer.

$$\delta_j^L = \frac{\partial}{\partial z_j^L} E = \frac{\partial}{\partial z_j^L} \frac{1}{2} \| \gamma - \hat{\gamma} \|^2$$

$$= \frac{\partial}{\partial z_j^L} \frac{1}{2} \sum_{k=1}^{n_L} \left( \gamma_k - f(z_k^L) \right)^2$$

$$= \frac{\partial}{\partial z_j^L} \frac{1}{2} \left( \gamma_j - f(z_j^L) \right)^2$$

$$= -\left( \gamma_j - f(z_j^L) \right) f'(z_j^L)$$

Therefore.

$$\delta^L = -\left( \gamma - f(z^L) \right) * f'(z_j^L)$$

$$= -\left( \gamma - a^L \right) * f'(z_j^L)$$

we denote $\begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix} * \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix} = \begin{pmatrix} a_1 * b_1 \\ a_2 b_2 \\ \vdots \\ a_n b_n \end{pmatrix}$

if $\ell < L$ i.e. the ~~input~~ hidden layer.

$$\delta_j^\ell = \frac{\partial}{\partial z_j^\ell} E$$

$$= \sum_{i=1}^{n_{\ell+1}} \frac{\partial E}{\partial z_i^{\ell+1}} \frac{\partial z_i^{\ell+1}}{\partial z_j^\ell} \quad \left( \text{multi variable} \atop \text{chain rule.} \right)$$
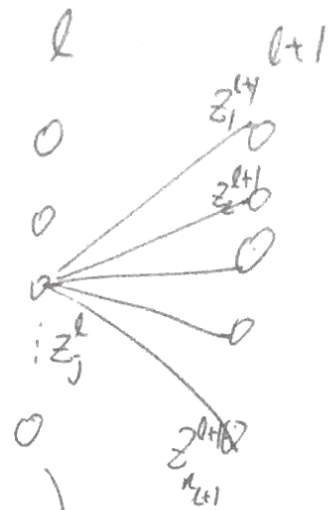
$$= \sum_{i=1}^{n_{\ell+1}} \delta_i^{\ell+1} \frac{\partial}{\partial z_j^\ell} \sum_{k=1}^{n_\ell} w_{ik}^\ell a_k^\ell$$

$$= \sum_{i=1}^{n_{\ell+1}} \delta_i^{\ell+1} \frac{\partial}{\partial z_j^\ell} w_{ij}^\ell a_j^\ell$$

$$= \sum_{i=1}^{n_{\ell+1}} \delta_i^{\ell+1} w_{ij}^\ell f'(z_j^\ell)$$

In matrix form

$$\delta^\ell = (w^\ell)^T \delta^{\ell+1} * f'(z^\ell)$$

cheat sheet (all column vectors)

Initialization:
$$a^\ell = \begin{pmatrix} a_1^\ell \\ a_2^\ell \\ \vdots \\ a_{n_\ell}^\ell \end{pmatrix} \quad z^\ell = \begin{pmatrix} z_1^\ell \\ z_2^\ell \\ \vdots \\ z_{n_\ell}^\ell \end{pmatrix} \quad \delta^\ell = \begin{pmatrix} \delta_1^\ell \\ \delta_2^\ell \\ \vdots \\ \delta_{n_\ell}^\ell \end{pmatrix}$$

shape
$$n_\ell \times 1$$

$$W^\ell = \begin{pmatrix} w_{11}^\ell & w_{12}^\ell & \cdots & w_{1n_\ell}^\ell \\ w_{21}^\ell & w_{22}^\ell & \cdots & w_{2n_\ell}^\ell \\ \vdots & \vdots & & \vdots \\ w_{n_{\ell+1}1}^\ell & w_{n_{\ell+1}2}^\ell & \cdots & w_{n_{\ell+1}n_\ell}^\ell \end{pmatrix}$$

shape
$$n_{\ell+1} \times n_\ell$$

**Forward**   $a^1 = X$

$$z^{\ell+1} = W^\ell a^\ell$$
$$a^{\ell+1} = f(z^{\ell+1})$$

$$\hat{y} = a^L$$

**Backward**   $\delta^L = -(y - a^L) * f'(z^L)$

$$\delta^\ell = (W^\ell)^T \delta^{\ell+1} * f'(z^\ell)$$

$$\nabla_{W^\ell} E(x) = \delta^{\ell+1} \otimes (a^\ell)^T$$

$$f'(z^\ell) = a^\ell(1 - a^\ell) \quad \text{if } f \text{ is sigmoid}$$
$$f'(z^\ell) = 1 \quad \text{if } f(x) = X$$

**Update**

$$W^\ell \leftarrow W^\ell - \eta \frac{1}{N} \sum_{i=1}^{N} \nabla_{W^\ell} E(x^{(i)})$$

cheat sheet  (all row vectors)

## Initialization

$$a^\ell = (a_1^\ell \; a_2^\ell \cdots a_{n_\ell}^\ell) \qquad z^\ell = (z_1^\ell \; z_2^\ell \cdots z_{n_\ell}^\ell)$$

$$\delta^\ell = (\delta_1^\ell, \; \delta_2^\ell \cdots \delta_{n_\ell}^\ell) \qquad \text{shape} \quad 1 \times n_\ell$$

$$w^\ell = \begin{pmatrix} w_{11}^\ell & w_{12}^\ell & & w_{1 n_{\ell+1}}^\ell \\ w_{21}^\ell & w_{22}^\ell & & w_{2 n_{\ell+1}}^\ell \\ \vdots & & & \\ w_{n_\ell 1}^\ell & w_{n_\ell 2}^\ell & \cdots & w_{n_\ell n_{\ell+1}}^\ell \end{pmatrix} \qquad \text{shape} \quad n_\ell \times n_{\ell+1}$$

## Forward

$$a^1 = x \qquad z^{\ell+1} = a^\ell w^\ell \qquad \hat{y} = a^L$$
$$a^{\ell+1} = f(z^{\ell+1})$$

## Backward

$$\delta^L = -(y - a^L) * f'(z^L)$$

$$\delta^\ell = \delta^{\ell+1} (w^\ell)^T * f'(z^\ell)$$

$$\nabla_{w^\ell} E(x) = (a^\ell)^T \delta^{\ell+1}$$

$$f'(z^\ell) = a^\ell (1 - a^\ell) \qquad \text{if } f \text{ is sigmoid}$$
$$f'(z^\ell) = 1 \qquad\qquad \text{if } f(x) = x$$

## Upadte

$$w^\ell \leftarrow w^\ell - \gamma \frac{1}{N} \sum_{i=1}^{N} \nabla_{w^\ell} E(x^{(i)})$$

Question:

1. if we add bias term $b^\ell$

$$z^{\ell+1} = w^\ell a^\ell + b^\ell$$

compute $\nabla_{b^\ell} E$.

2. if we add regulization term in the cost function.

$$E = \frac{1}{2N} \sum_{r=1}^{N} \left\| \hat{F}(x^{(r)}) - y \right\|^2 + \frac{\lambda}{2} \sum_{\ell=1}^{L} \sum_{i=1}^{n_\ell} \sum_{j=1}^{n_{\ell+1}} \left( w_{ji}^\ell \right)^2$$

compute $\nabla_{w^\ell} E$   $\nabla_{b^\ell} E$.

3. compute $f'(z^\ell)$ when $f$ is ReLU

$$f(x) = \max(0, x)$$