

Lecture 15 : Stochastic Gradient Descent + Neural Networks

Stochastic Gradient Descent

Recall gradient descent:

$$\underline{w}^* = \arg \min_{\underline{w}} f(\underline{w})$$

$$\underline{w}^{(k+1)} = \underline{w}^{(k)} - \gamma \nabla f(\underline{w}^{(k)})$$

Imagine $f(\underline{w}) = \frac{1}{n} \sum_{i=1}^n f_i(\underline{w})$

e.g. if $f(\underline{w}) = \frac{1}{n} \|y - X\underline{w}\|_2^2$
 $= \frac{1}{n} \sum_{i=1}^n (y_i - \langle x_i, \underline{w} \rangle)^2$
 $\Rightarrow f_i(\underline{w}) = (y_i - \langle x_i, \underline{w} \rangle)^2$

Then Gradient Descent =

$$\underline{w}^{(t+1)} = \underline{w}^{(t)} - \gamma \sum_{i=1}^n \nabla f_i(\underline{w}^{(t)})$$

Today:

@ iteration t , choose $i \in \{1, 2, \dots, n\}$

$$\underline{w}^{(t+1)} = \underline{w}^{(t)} - \gamma \nabla f_{i_t}(\underline{w}^{(t)})$$

- each iteration easier/faster to compute
- need more iterations

How to choose i_t ?

A. cyclical ("Incremental Gradient Descent")

$$i_t = t \bmod n$$

e.g. $n=3$: i_t 's = 1, 2, 3, 1, 2, 3, 1, 2, ...

B. random permutations

every n rounds, reshuffle

e.g. $n=3$: i_t 's = 1, 3, 2, 3, 1, 2, 2, 1, 3, ...

C. choose i_t uniformly at random

"stochastic gradient descent"

$$i_t \sim \text{unif}(1, \dots, n)$$

e.g. $n=3$: i_t 's = 1, 3, 3, 2, 3, 1, 2, 2, 2, ...

note: expected value $\mathbb{E}[\nabla f_{i_t}(\underline{w})] = \nabla f(\underline{w})$

$$\text{Ex: } f(\underline{w}) = \frac{1}{n} \sum_{i=1}^n (y_i - \langle \underline{x}_i, \underline{w} \rangle)^2 + \lambda \|\underline{w}\|_2^2$$

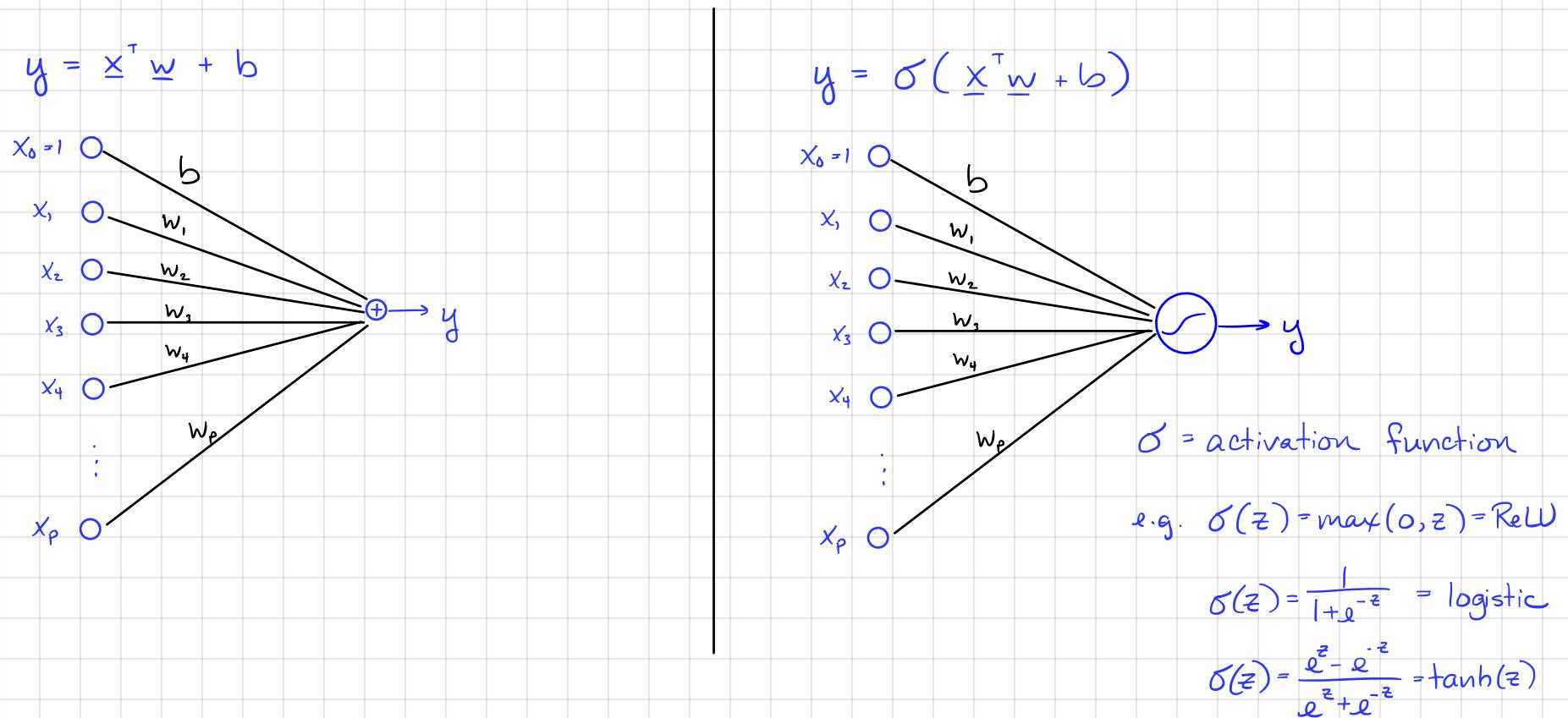
$$f_i(\underline{w}) = (y_i - \langle \underline{x}_i, \underline{w} \rangle)^2 + \lambda \|\underline{w}\|_2^2$$

$$\text{check: } \frac{1}{n} \sum_{i=1}^n f_i(\underline{w}) = f(\underline{w})$$

$$\nabla f_i(\underline{w}) = -2(y_i - \langle \underline{x}_i, \underline{w} \rangle) \underline{x}_i + 2\lambda \underline{w}$$

$$\text{SGD: } \underline{w}^{(t+1)} = \underline{w}^{(t)} + 2\tau (y_{i_t} - \langle \underline{x}_{i_t}, \underline{w}^{(t)} \rangle) \underline{x}_{i_t} - 2\tau \lambda \underline{w}^{(t)}$$

A Simple Neural Network



for $\hat{y} = \sigma(\underline{x}^\top \underline{w} + b)$ and $\sigma(z) = \frac{1}{1+e^{-z}}$, how do we learn weights?

$$\text{loss } f(\underline{w}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(\underbrace{\langle \underline{x}_i, \underline{w} \rangle + b}_{z_i}))^2$$

learn \underline{w} via Stochastic Gradient Descent!

@ iteration t

- choose i_t uniformly at random

- set $\underline{w}^{(t+1)} = \underline{w}^{(t)} - \tau \nabla f_{i_t}(\underline{w}^{(t)})$

$$b^{(t+1)} = b^{(t)} - \tau \nabla f_{i_t}(b^{(t)})$$

What is $\nabla f_{i_t}(\underline{w}^{(t)})$?

$$\frac{df_{i_t}}{dw_j} \Big|_{\underline{w}^{(t)}} = \frac{df_{i_t}}{d\hat{y}_i} \cdot \frac{d\hat{y}_i}{dz_i} \cdot \frac{dz_i}{dw_j} \Big|_{\underline{w}^{(t)}}$$

$$= 2(\hat{y}_i - y_i) \cdot \sigma'(z_i) \cdot x_{ij} \Big|_{\underline{w}^{(t)}}$$

$$= 2(\hat{y}_i - y_i) \cdot \sigma(z_i)(1 - \sigma(z_i)) x_{ij} \Big|_{\underline{w}^{(t)}}$$

$$= 2(\hat{y}_i - y_i) \hat{y}_i (1 - \hat{y}_i) x_{ij}$$

scalar, independent of j
call $S_i = S_i(\underline{w}^{(t)})$

$$\hat{y}_i = \sigma(\underline{x}_i^\top \underline{w} + b) = \sigma(z_i)$$

$$z_i = \underline{x}_i^\top \underline{w} + b$$

aside : $\frac{d\sigma}{dz} = \sigma'(z)$

$$= \sigma(z)(1 - \sigma(z))$$

What is $\nabla f_{i_t}(b^{(t)}) = \frac{df_{i_t}}{db} \Big|_{b^{(t)}}$?

$$\frac{df_{i_t}}{db} \Big|_{b^{(t)}} = \frac{df_{i_t}}{d\hat{y}_i} \cdot \frac{d\hat{y}_i}{dz_i} \cdot \frac{dz_i}{db} \Big|_{b^{(t)}}$$

$$= S_i \cdot 1$$

$$\Rightarrow \nabla f_{i_t}(\underline{w}^{(t)}) = S_i \underline{x}_i$$

$$\Rightarrow \text{SGD: } \underline{w}^{(t+1)} = \underline{w}^{(t)} - \tau S_{i_t} \underline{x}_{i_t}, b^{(t+1)} = b^{(t)} - \tau S_{i_t}$$

3-layer network (1 hidden layer)

w_{kj} = weight on j^{th} element of x_i on hidden node k

Hidden

$\Rightarrow h_{ik} = \text{output of } k^{\text{th}} \text{ hidden node if } x_i = \text{input}$

$$= \sigma(W x_i) = \sigma\left(\sum_{j=1}^p w_{kj} x_{ij}\right)$$

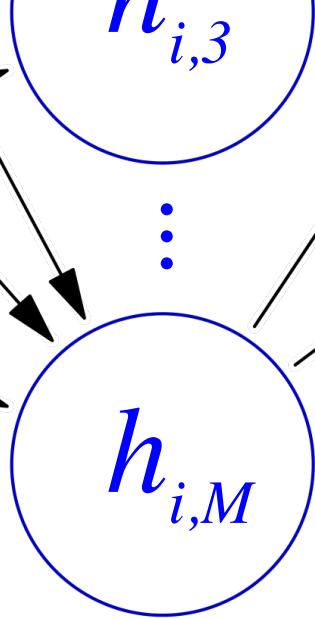
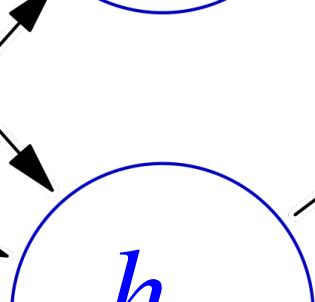
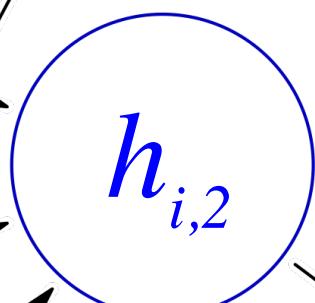
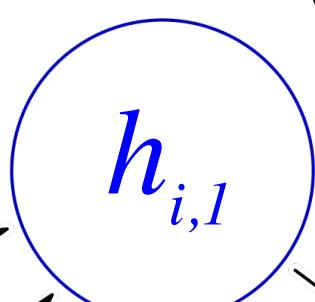
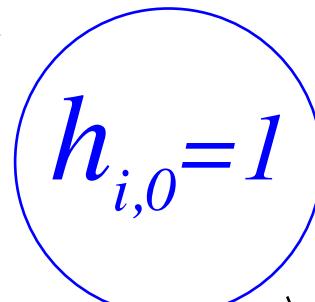
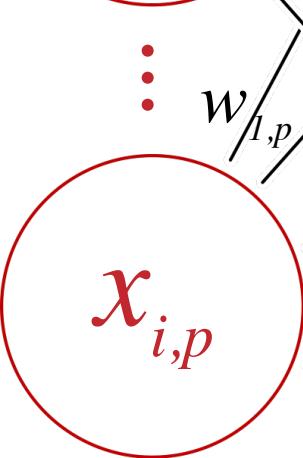
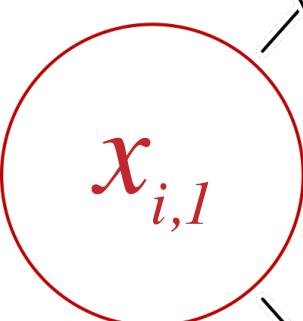
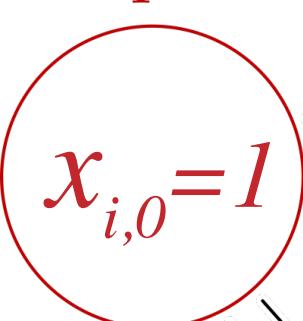
$$\Rightarrow \underline{h}_i = \begin{bmatrix} h_{i,0} \\ h_{i,1} \\ \vdots \\ h_{i,M} \end{bmatrix}$$

v_{kj} = weight on j^{th} hidden node output on predictor k

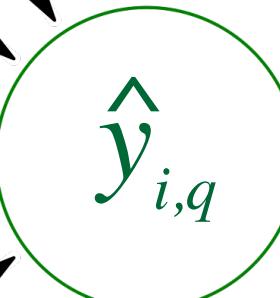
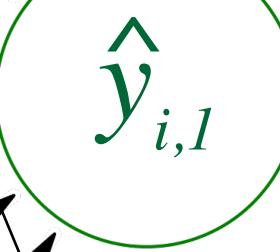
$\Rightarrow \hat{y}_{ik} = \text{output } k \text{ for input } i$

$$= \sigma(V \underline{h}_i) = \sigma\left(\sum_{m=0}^M v_{km} h_{im}\right)$$

Input



Output



$$w_{1,0}$$

$$w_{2,0}$$

$$w_{1,p}$$

$$w_{M,p}$$

$$v_{1,0}$$

$$v_{1,1}$$

$$v_{q,1}$$

$$v_{2,q}$$

$$v_{1,M}$$

$$v_{q,M}$$

SGD for 3-layer network:

$$\hat{y}_{ik} = \sigma(V h_i) = \sigma(V \sigma(W x_i)) = \sigma\left(\sum_{m=0}^M v_{km} \sigma\left[\sum_{j=0}^P w_{mj} x_{ij}\right]\right)$$

Backpropagation algorithm

- choose initial weights $W^{(l)}$ and $V^{(l)}$
- for $t=1, 2, \dots$
 - choose i_t
 - calculate \hat{y}_{itk} and h_{itm} using $W^{(t)}$ and $V^{(t)}$ (forward pass)
 - update weights for each layer, starting with deepest layer (closest to output) and working back to shallowest

To update $\mathbf{V}^{(t)}$:

$$\frac{df_i}{dV_{km}} = \frac{df_i}{d\hat{y}_i} \cdot \frac{d\hat{y}_i}{dV_{k,m}}$$

$$= 2(\hat{y}_{ik} - y_{ik}) \sigma'(\mathbf{V}^{(t)} \underline{h}_i) h_{im}$$

$$= \underbrace{2(\hat{y}_{ik} - y_{ik}) \hat{y}_{ik} (1 - \hat{y}_{ik})}_{= S_{i,k}} h_{im}$$

$$= S_{ik} h_{im}$$

$$\Rightarrow V_{k,m}^{(t+1)} = V_{k,m}^{(t)} - 2\tau S_{i_t k} h_{i_t m}$$

$$\Rightarrow \underline{V}_k^{(t+1)} = \underline{V}_k^{(t)} - 2\tau S_{i_t k} \underline{h}_{i_t} \quad \text{where } \underline{V}_k = \left(\begin{matrix} k^{\text{th}} \text{ row of } \mathbf{V} \end{matrix} \right)^T = \text{vector of weights determining } k^{\text{th}} \text{ predictor}$$

$$\Rightarrow \mathbf{V}^{(t+1)} = \mathbf{V}^{(t)} - 2\tau \underline{h}_{i_t} \underline{S}_i^T \quad \text{where } \underline{S}_i^T = [S_{i,0} \ S_{i,1} \ \dots \ S_{i,k}]$$

To update $W^{(t)}$:

$$\frac{df_i}{dW_{m,j}} = \sum_{k=1}^q \frac{df_i}{d\hat{y}_{ik}} \cdot \frac{d\hat{y}_{ik}}{dh_{im}} \cdot \frac{dh_{im}}{dW_{m,j}}$$

$$= 2 \sum_{k=1}^q (\hat{y}_{ik} - y_{ik}) \sigma' \left((\underline{v}_k^{(t)})^\top \underline{h}_i \right) v_{k,m} \cdot \sigma' \left((\underline{w}_m^{(t)})^\top \underline{x}_i \right) x_{ij}$$

$$= \sum_{k=1}^q \underbrace{2(\hat{y}_{ik} - y_{ik}) \hat{y}_{ik} (1 - \hat{y}_{ik})}_{=: S_{i,k}} v_{k,m} h_{im} (1 - h_{im}) x_{ij}$$

$$= \sum_{k=1}^q \underbrace{S_{i,k} v_{k,m} h_{im} (1 - h_{im})}_{=: \gamma_{im}} x_{ij}$$

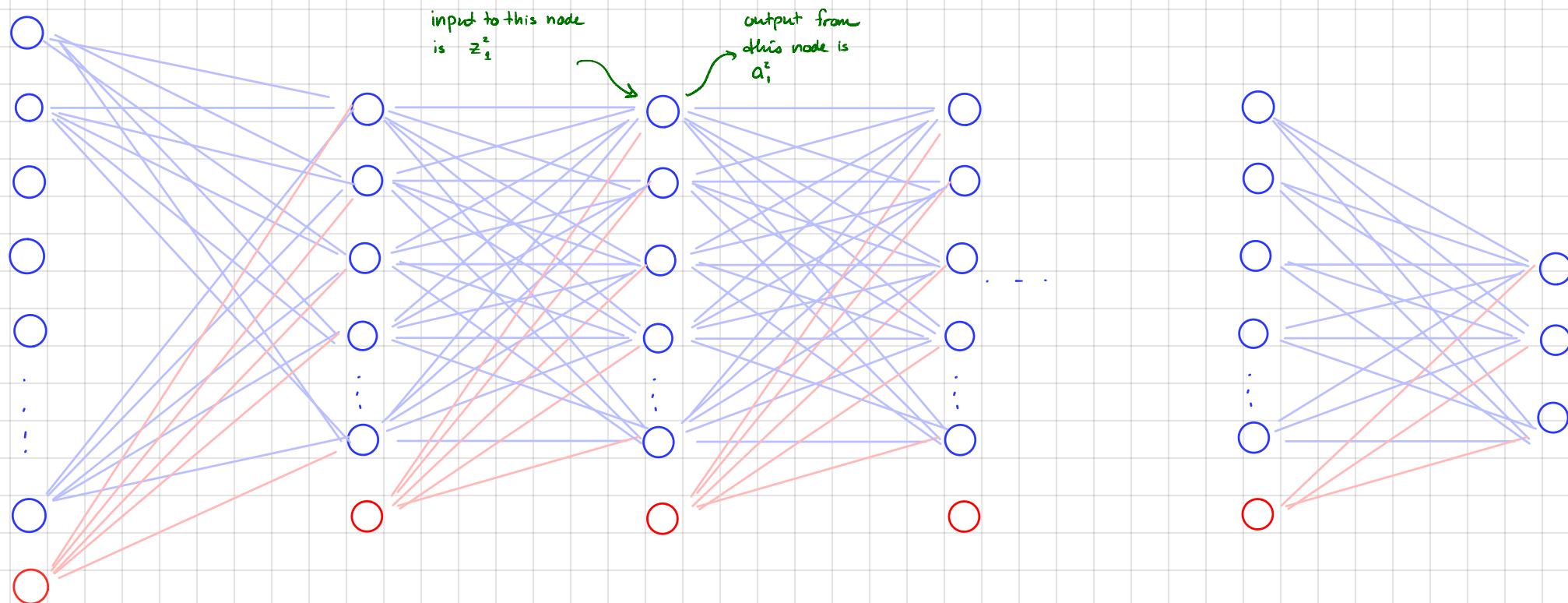
$$= \gamma_{im} x_{ij}$$

$$\Rightarrow W_{m,j}^{(t+1)} = W_{m,j}^{(t)} - \tau \gamma_{i_t m} x_{i_t j}$$

$$\Rightarrow \underline{w}_m^{(t+1)} = \underline{w}_m^{(t)} - \tau \gamma_{i_t m} \underline{x}_{i_t} \quad \text{where } \underline{w}_m = (\text{m}^{\text{th}} \text{ row of } W)^\top = \text{weights going into } m^{\text{th}} \text{ hidden node.}$$

$$\Rightarrow W^{(t+1)} = W^{(t)} - \tau \underline{x}_{i_t} \underline{\gamma}_{i_t}^\top \quad \text{where } \underline{\gamma}_{i_t}^\top = [\gamma_{i_t, 1}, \dots, \gamma_{i_t, M}]$$

More generally



X

W^1, b^1

$= a^0$

$$z^1 = W^1 a^0 + b^1$$

$$a^1 = \sigma(z^1)$$

W^2, b^2

$$z^2 = W^2 a^1 + b^2$$

$$a^2 = \sigma(z^2)$$

W^3, b^3

$$z^3 = W^3 a^2 + b^3$$

$$a^3 = \sigma(z^3)$$

\dots

W^L, b^L

$$z^L = W^L a^{L-1} + b^L$$

$$a^L = \sigma(z^L) = \hat{y}$$

$$W^l \sim M^l \times M^{l-1}, \quad S^l \sim M^l \times 1, \quad a^l \sim M^l \times 1, \quad b^l \sim M^l \times 1$$

z^l is input to nodes in layer l , a^l is output of those nodes

$W_{j,k}^l$ = weight in layer l , applied to layer $l-1$ output (a_k^{l-1}), feeding into next layer (z_j^l)

To update W^l , need to compute $\nabla f(W^l)$

$$\frac{df}{dW_{j,k}^l} = \frac{df}{dz_j^l} \cdot \frac{dz_j^l}{dW_{j,k}^l} = S_j^l \cdot a_k^{l-1} \Rightarrow \nabla f(W^l) = S^l (a^{l-1})^\top$$

$$S_j^l := \frac{df}{dz_j^l} = \begin{cases} \frac{df}{da_j^l} \cdot \frac{da_j^l}{dz_j^l} = \left[(W^{l+1})^\top S^{l+1} \right]_j \cdot \sigma'(z_j^l) & \text{for } l < L \\ \frac{df}{a_j^l} \sigma'(z_j^l) & \text{for } l = L \end{cases} \Rightarrow S^l = \begin{cases} ((W^{l+1})^\top S^{l+1}) \odot \sigma'(z^l) & l < L \\ \nabla f(a^L) \odot \sigma'(z^l) & l = L \end{cases}$$

$$\frac{df}{da_j^l} = \sum_k \frac{df}{dz_k^{l+1}} \frac{dz_k^{l+1}}{da_j^l} = \sum_k S_k^{l+1} W_{kj}^{l+1} = \left[(W^{l+1})^\top S^{l+1} \right]_j$$

$$\frac{da_j^l}{dz_j^l} = \sigma'(z_j^l)$$

$$\frac{dz_j^l}{dW_{j,k}^l} = a_k^{l-1} \text{ because } z_j^l = (W^l a^{l-1} + b^l)_j = \sum_k W_{jk}^l a_k^{l-1} + b_j^l$$

$$\frac{df}{db_j^l} = \frac{df}{dz_j^l} \cdot \frac{dz_j^l}{db_j^l} = S_j^l \cdot 1 \Rightarrow \nabla f(b^l) = S^l_j$$

Backpropagation Algorithm

forward pass:

$$a^0 = \sigma(\underline{x}_{i_t})$$

for $l = 1, 2, \dots, L$

$$\underline{z}^l = W^l \underline{a}^{l-1} + \underline{b}^l$$

$$\underline{\alpha}^l = \sigma(\underline{z}^l)$$

end

$$\underline{s}^L = \nabla f(a^L) \cdot \sigma'(\underline{z}^L)$$

backprop

for $l = L-1, L-2, \dots, 1$

$$\underline{s}^l = [(W^{l+1})^\top \underline{s}^{l+1}] \odot \sigma'(\underline{z}^l)$$

$$\nabla f(W^l) = \underline{s}^l (\underline{\alpha}^{l-1})^\top$$

$$\nabla f(b^l) = \underline{s}^l$$

$$W^l \leftarrow W^l - \tau \nabla f(W^l)$$

$$b^l \leftarrow b^l - \tau \nabla f(b^l)$$

end