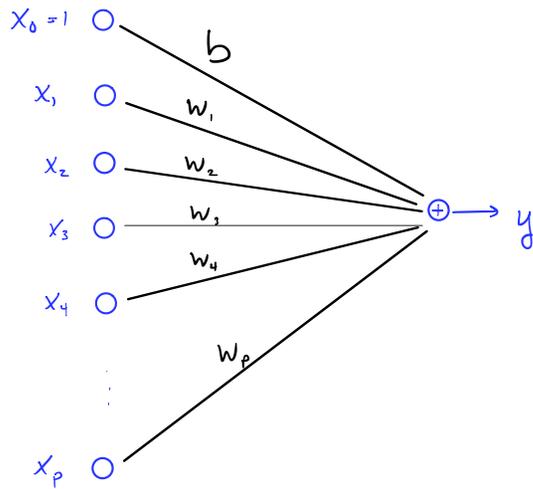


Lecture 14

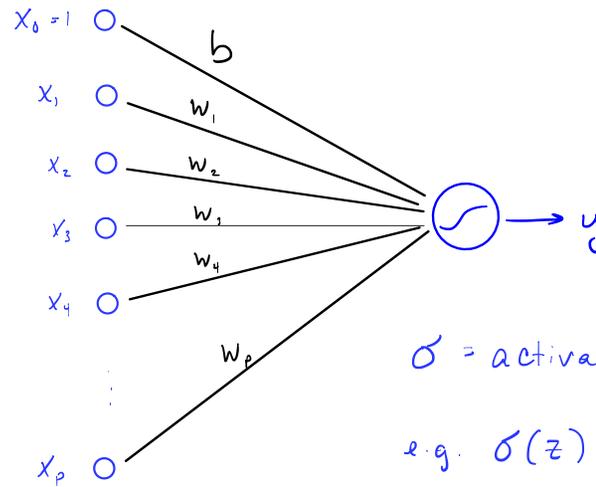
Backpropagation

A Simple Neural Network

$$y = \underline{x}^T \underline{w} + b$$



$$y = \sigma(\underline{x}^T \underline{w} + b)$$



σ = activation function

e.g. $\sigma(z) = \max(0, z) = \text{ReLU}$

$$\sigma(z) = \frac{1}{1+e^{-z}} = \text{logistic}$$

$$\sigma(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \text{tanh}(z)$$

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

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

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

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

$$\left. \frac{df_i}{dw_j} \right|_{\underline{w}^{(t)}} = \frac{df_i}{d\hat{y}_i} \cdot \frac{d\hat{y}_i}{dz_i} \cdot \left. \frac{dz_i}{dw_j} \right|_{\underline{w}^{(t)}}$$

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

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

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

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

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

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

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

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

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

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

$$\frac{d}{dz} (1 + e^{-z})^{-1} = + (1 + e^{-z})^{-2} (e^{-z})$$

$$= \frac{1}{1 + e^{-z}} \cdot \frac{e^{-z}}{1 + e^{-z}}$$

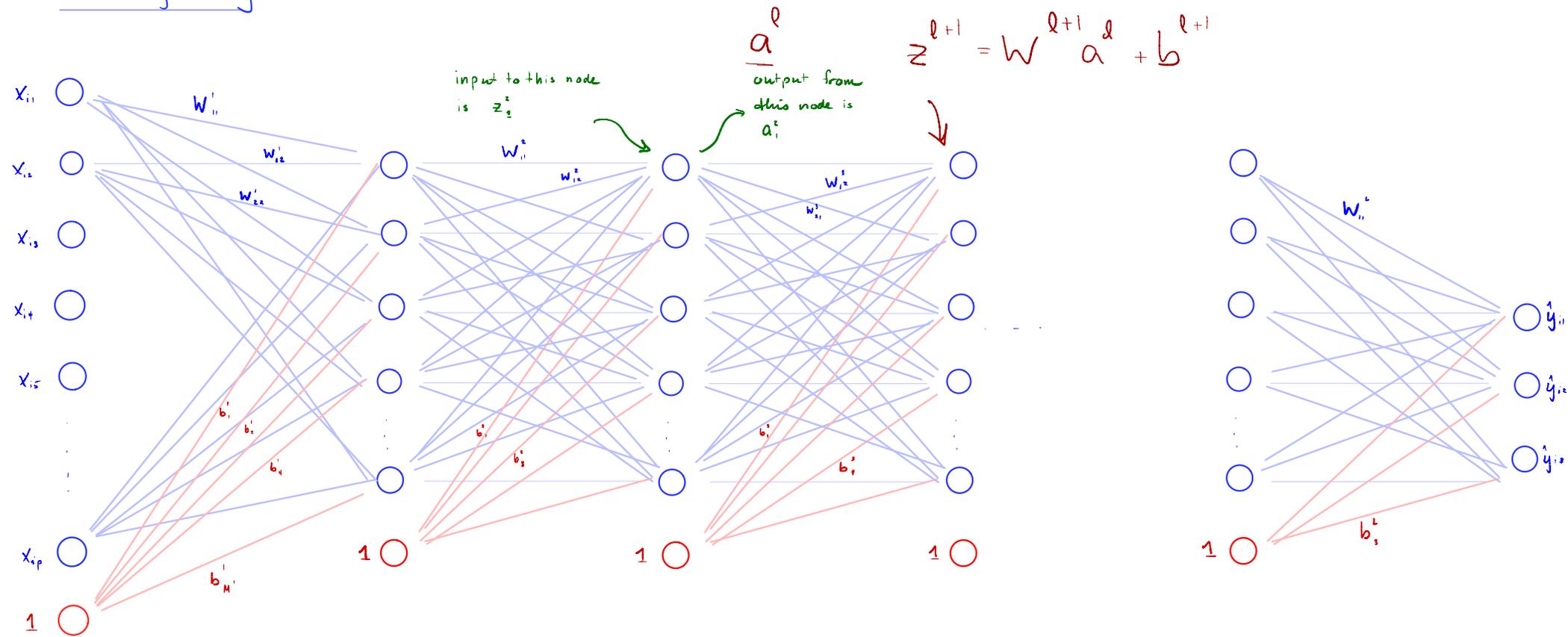
$$= \frac{1}{\sigma(z)} \cdot \frac{1 - \sigma(z)}{1 - \sigma(z)}$$

What is $\nabla_b f_{i_t}(b^{(t)}) = \left. \frac{df_{i_t}}{db} \right|_{b^{(t)}}$?

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

$$= \delta_i \cdot 1$$

More generally



a^l
output from this node is a^l

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

\underline{x}
 $= a^0$

W^1, b^1

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

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

W^2, b^2

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

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

W^3, b^3

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

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

W^L, b^L

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

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

Let $M^l = \#$ of blue nodes in layer l

$W^l \sim M^l \times M^{l-1}$,

$S^l \sim M^l \times 1$,

$\underline{a}^l, \underline{z}^l \sim M^l \times 1$,

$\underline{b}^l \sim M^l \times 1$

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 $\nabla f(W^l)$

$$\frac{df}{dW_{j,k}^l} = \frac{df}{dz_j^l} \cdot \frac{dz_j^l}{dW_{j,k}^l} = \delta_j^l \cdot a_k^{l-1} \implies \nabla f(W^l) = \underline{\delta}^l (\underline{a}^{l-1})^T$$

$$\delta_j^l := \frac{df}{dz_j^l} = \begin{cases} \frac{df}{da_j^l} \cdot \frac{da_j^l}{dz_j^l} = \left[(W^{l+1})^T \delta^{l+1} \right]_j \cdot \sigma'(z_j^l) & l < L \\ \frac{df}{da_j^l} \cdot \frac{da_j^l}{dz_j^l} & l = L \end{cases}$$

$$\underline{\delta}^l = \left\{ \begin{array}{l} ((W^{l+1})^T \underline{\delta}^{l+1}) \odot \sigma'(\underline{z}^l) \\ \nabla f(\underline{a}^l) \odot \sigma'(\underline{z}^l) \end{array} \right.$$

← element-wise multiplication

$$\frac{df}{da_j^l} = \sum_{k=1}^{M^{l+1}} \frac{df}{dz_k^{l+1}} \frac{dz_k^{l+1}}{da_j^l} = \sum_{k=1}^{M^{l+1}} \delta_k^{l+1} W_{kj}^{l+1} = \left[(W^{l+1})^T \delta^{l+1} \right]_j$$

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

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

$$\frac{df}{db_j^l} = \frac{df}{dz_j^l} \cdot \frac{dz_j^l}{db_j^l} = \delta_j^l \cdot 1 \implies \nabla f(b^l) = \delta_j^l$$

consider $f(x_i) = \sum_{k=1}^p (\hat{y}_{i,k} - y_{i,k})^2 = \|\hat{y}_i - y_i\|_2^2$

$$\begin{aligned} \frac{df}{da_j^l} &= \frac{d}{da_j^l} \left(\sum_{k=1}^p (a_k^l - y_{i,k})^2 \right) \\ &= \frac{d}{da_j^l} (a_j^l - y_{i,j})^2 \\ &= 2(a_j^l - y_{i,j}) \end{aligned}$$

$$a^l = \hat{y}_i \implies a_j^l = \hat{y}_{i,j}$$

Backpropagation Algorithm

for $t=1, 2, 3, \dots$

select $i_t \sim \text{unif}(1, 2, \dots, n)$

forward pass:

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

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

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

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

end

forward pass

backprop

$$\underline{\delta}^L = \nabla_{f_{i_t}}(a^L) \cdot \sigma'(\underline{z}^L)$$

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

$$\underline{\delta}^l = [(\underline{W}^{l+1})^T \underline{\delta}^{l+1}] \circ \sigma'(\underline{z}^l)$$

$$\nabla f(\underline{W}^{l+1}) = \underline{\delta}^l (\underline{a}^{l+1})^T$$

$$\nabla f(\underline{b}^{l+1}) = \underline{\delta}^l$$

$$\underline{W}^{l+1} = \underline{W}^{l+1} - \tau \nabla f(\underline{W}^{l+1})$$

$$\underline{b}^{l+1} = \underline{b}^{l+1} - \tau \nabla f(\underline{b}^{l+1})$$

end

backprop.

The following slides show backpropagation for a 2-layer network using slightly different notation. The concepts are no different than the above slides, but some students find it easier to study the 3-layer case first.

2-layer network (1 hidden layer)

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

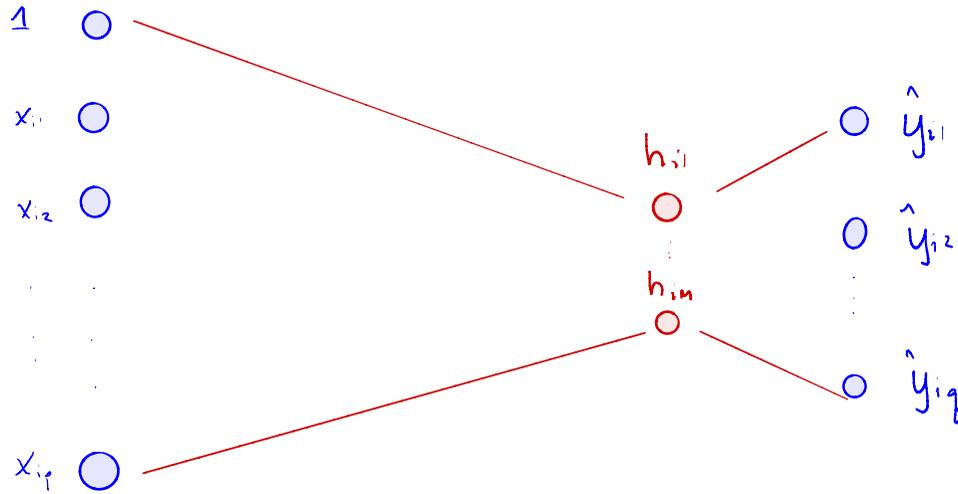
$$\begin{aligned}\Rightarrow h_{ik} &= \text{output of } k^{\text{th}} \text{ hidden node if } x_i = \text{input} \\ &= \sigma(w_{k \cdot} x_i) = \sigma\left(\sum_{j=1}^P w_{kj} x_{ij}\right)\end{aligned}$$

$$\Rightarrow \underline{h}_i = \begin{bmatrix} h_{i1} \\ h_{i2} \\ \vdots \\ h_{iM} \end{bmatrix}$$

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

$$\begin{aligned}\Rightarrow \hat{y}_{ik} &= \text{output } k \text{ for input } i \\ &= \sigma(v_{k \cdot} \underline{h}_i) = \sigma\left(\sum_{m=1}^M v_{km} h_{im}\right)\end{aligned}$$

Ex: 2-layer network



When M is small (i.e. small # hidden nodes @ some layer), then we can take the vector

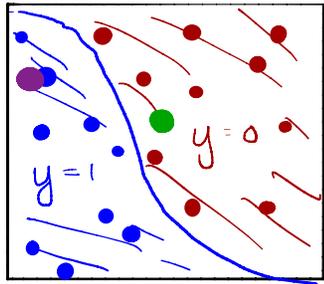
$$\begin{bmatrix} h_{i1} \\ \vdots \\ h_{im} \end{bmatrix}$$

and this can be considered as latent low-dimensional features.

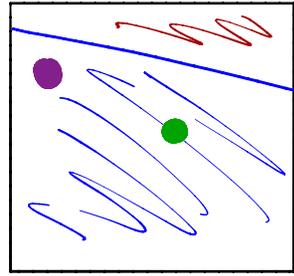
Output of hidden node k : $x_i \rightarrow h_k$

1st set of weights
define these
intermediate
linear classifiers

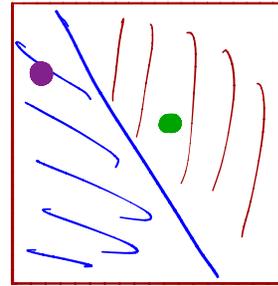
x_{i2} ↑



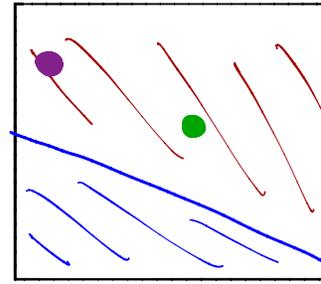
→ x_{i1}



$k=1$

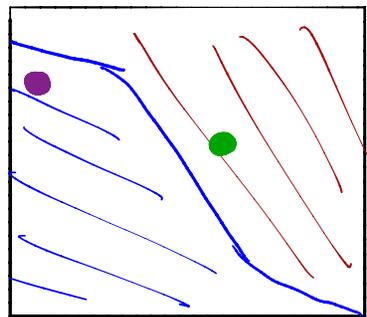


$k=2$



$k=3$

final output = weighted sum of h_k 's + thresholding



2nd set of weights
let us combine the linear
classifiers

SGD for 2-layer network:

$$\hat{y}_{ik} = \sigma(V h_i) = \sigma\left(V \sigma\left(W \underline{x}_i\right)\right) = \sigma\left(\sum_{m=0}^M V_{km} \sigma\left[\sum_{j=0}^P W_{mj} x_{ij}\right]\right)$$

Backpropagation algorithm

· choose initial weights $W^{(i)}$ and $V^{(i)}$

· for $t=1, 2, \dots$

· choose i_t

· calculate $\hat{y}_{i_t k}$ and $h_{i_t m}$ using $W^{(t)}$ and $V^{(t)}$ (forward pass)

$$\hat{y}_{i_t k} = \sigma(V^{(t)} h_{i_t}) = \sigma\left(V^{(t)} \sigma\left(W^{(t)} \underline{x}_{i_t}\right)\right) = \sigma\left(\sum_{m=0}^M V_{km}^{(t)} \sigma\left[\sum_{j=0}^P W_{mj}^{(t)} x_{i_t j}\right]\right)$$

· update weights for each layer, starting with deepest layer (closest to output) and working back to shallowest

$$V^{(t+1)} = V^{(t)} - \tau \underline{h}_{i_t} \underline{\delta}_{i_t}^T \quad \text{where } \underline{\delta}_{i_t}^T = [\delta_{i_t,0}, \delta_{i_t,1}, \dots, \delta_{i_t,k}]$$

$$W^{(t+1)} = W^{(t)} - \tau \underline{x}_{i_t} \underline{\gamma}_{i_t}^T \quad \text{where } \underline{\gamma}_{i_t}^T = [\gamma_{i_t,0}, \dots, \gamma_{i_t,M}]$$

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

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

$$= 2(\hat{y}_{i,k} - y_{i,k}) \sigma'(v^{(t)} \underline{h}_i) h_{i,m}$$

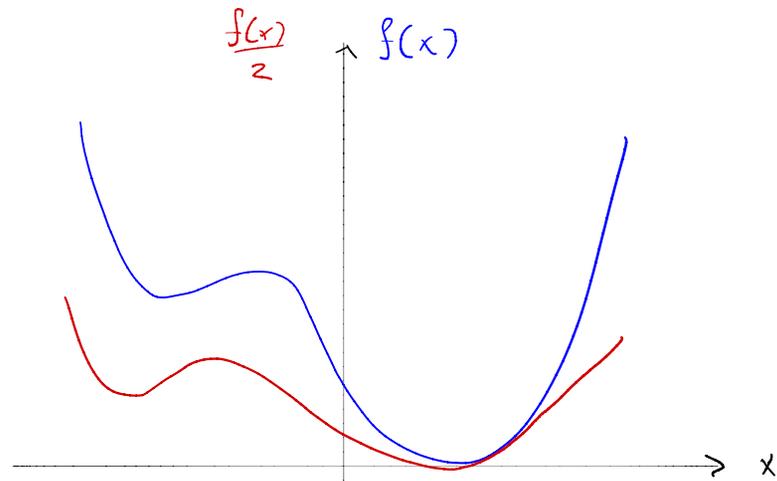
$$= \underbrace{2(\hat{y}_{i,k} - y_{i,k}) \hat{y}_{i,k} (1 - \hat{y}_{i,k})}_{= \delta_{i,k}} h_{i,m}$$

$$= \delta_{i,k} h_{i,m}$$

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

$$\Rightarrow \underline{V}_k^{(t+1)} = \underline{V}_k^{(t)} - \tau \delta_{i_t, k} \underline{h}_{i_t} \quad \text{where } \underline{V}_k = (k^{\text{th}} \text{ row of } V)^T = \text{vector of weights determining } k^{\text{th}} \text{ predictor}$$

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



if $f_i = (y_i - \hat{y}_i)^2$

then $\frac{df_i}{d\hat{y}_i} = -2(y_i - \hat{y}_i)$

if $f_i = \frac{1}{2} (y_i - \hat{y}_i)^2$
then $\frac{df_i}{d\hat{y}_i} = -(y_i - \hat{y}_i)$

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

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

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

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

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

$$= \gamma_{i,m} x_{ij}$$

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

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

$$\Rightarrow W^{(t+1)} = W^{(t)} - \tau \underline{x}_{i_t} \underline{\gamma}_{i_t}^T \quad \text{where } \underline{\gamma}_{i_t}^T = [\gamma_{i_t, 0}, \dots, \gamma_{i_t, M}]$$