

Feed Forward Neural Networks

Input: n dimensional vector (0^{th} layer)

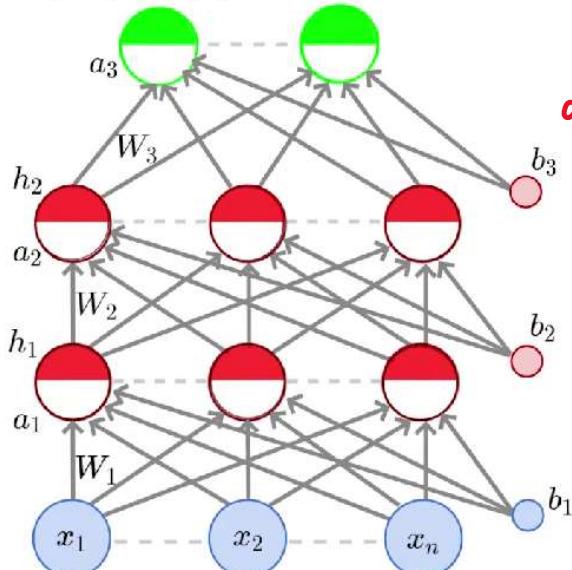
Hidden layers: $L-1$ hidden layers with 2 neurons each
(2 in our case)

Output layer: 1 output layer containing k neurons (for k classes)
(L^{th} layer)

Each layer in the hidden & output layer will be split into two parts: pre-activation(a_i) & activation(h_i)

Weights: $W_i \in \mathbb{R}^{n \times n}$ & $W_L \in \mathbb{R}^{k \times n}$
Bias: $b_i \in \mathbb{R}^n$ & $b_L \in \mathbb{R}^k$
b/w layers $i-1$ & i
($0 \leq i \leq L$)

$$h_L = \hat{y} = \hat{f}(x)$$



$$a_i(x) = b_i + W_i h_{i-1}(x)$$

$$a_i = \begin{bmatrix} a_{i1} \\ a_{i2} \\ a_{i3} \end{bmatrix} = \begin{bmatrix} b_{i1} \\ b_{i2} \\ b_{i3} \end{bmatrix} + \begin{bmatrix} w_{i11} & w_{i12} & w_{i13} \\ w_{i21} & w_{i22} & w_{i23} \\ w_{i31} & w_{i32} & w_{i33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$h_i(x) = g(a_i(x))$$

$$h_i = g\left(\begin{bmatrix} a_{i1} \\ a_{i2} \\ a_{i3} \end{bmatrix}\right) = \begin{bmatrix} g(a_{i1}) \\ g(a_{i2}) \\ g(a_{i3}) \end{bmatrix}$$

Taking g as a sigmoid function:

$$\begin{cases} \frac{1}{1+e^{-a_{i1}}} \\ \vdots \end{cases}$$

g is called activation function

e.g. logistic, tanh, linear etc. =

unction

eg - logistic, tanh, linear etc. =

$$\begin{bmatrix} \dots \\ \frac{1}{1+e^{-a_{12}}} \\ \frac{1}{1+e^{-a_{13}}} \end{bmatrix}$$

The activation in the output layer is given by,

$$f(x) = h_L(x) = O(a_L(x))$$

O is the output activation function
eg - softmax, linear, etc.

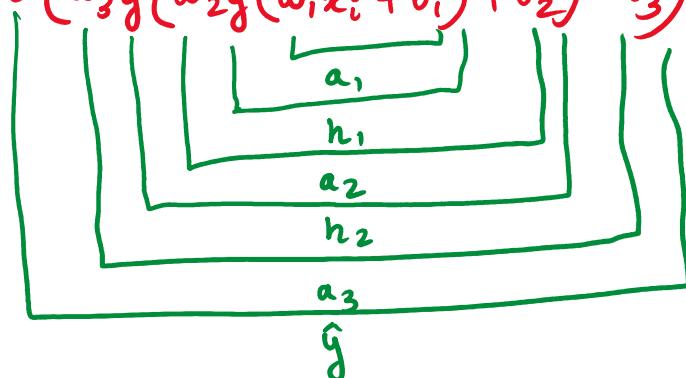
For simplicity, we are going to write $a_i(x)$ as a_i & $h_i(x)$ as h_i :

$$\begin{aligned} \therefore a_i &= b_i + W_i h_{i-1} \\ h_i &= g(a_i) \\ f(x) &= h_L = O(a_L) \end{aligned}$$

In a typical machine learning setup,

Data: $\{x_i, y_i\}_{i=1}^N$ (For this case of 3 inputs)

Model: $\hat{y}_i = f(x_i) = O(W_3 g(W_2 g(W_1 x_i + b_1) + b_2) + b_3)$



Parameters: $\theta = w_1, \dots, w_L, b_1, \dots, b_L$ ($L=3$ in this case)

Algorithm: Gradient Descent or Backpropagation

Algorithm: Gradient Descent or Backpropagation

Loss Function: say $\min \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^k (\hat{y}_{ij} - y_{ij})^2$

in general $\min L(\theta)$ Some loss function

Learning Parameters (intuition)

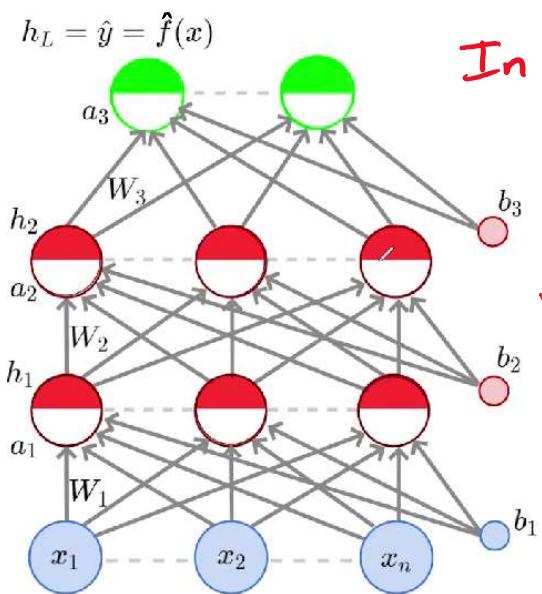
Writing the gradient descent algorithm more concisely:

```

 $t \leftarrow 0;$ 
max_iterations  $\leftarrow 1000;$ 
Initialize  $\Theta_t = [\omega_0, b_0];$ 
while  $t++ < \text{max\_iterations}$  do
     $\Theta_{t+1} \leftarrow \Theta_t - \eta \nabla \Theta_t$ 
end

```

where $\nabla \Theta_t = \begin{bmatrix} \frac{\partial L(\Theta)}{\partial \omega_t} & \frac{\partial L(\Theta)}{\partial b_t} \end{bmatrix}^T$



In this feed forward neural network we have,

$$\Theta = [\omega_1, \dots, \omega_L, b_1, \dots, b_L]$$

Now,

$$\nabla \Theta_t = \begin{bmatrix} \frac{\partial L(\Theta_t)}{\partial \omega_{1,t}} & \dots & \frac{\partial L(\Theta_t)}{\partial \omega_{L,t}}, \dots & \frac{\partial L(\Theta_t)}{\partial b_{1,t}}, \dots & \frac{\partial L(\Theta_t)}{\partial b_{L,t}} \end{bmatrix}^T$$

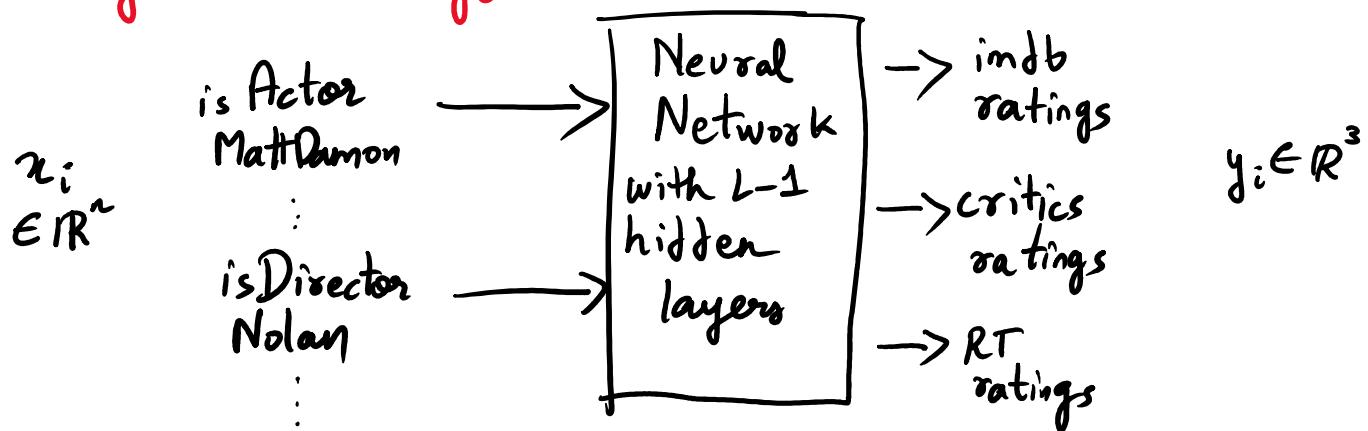
$\nabla \Theta$ is now composed of

$$\nabla \omega_1, \nabla \omega_2, \dots, \nabla \omega_{L-1} \in \mathbb{R}^{n \times n}, \nabla \omega_L \in \mathbb{R}^{k \times n},$$

$$\nabla b_1, \nabla b_2, \dots, \nabla b_{L-1} \in \mathbb{R}^n, \nabla b_L \in \mathbb{R}^k$$

Output Function: & Loss Functions

The choice of loss function depends upon the problem
Consider the movie example but this time we are predicting movie ratings.

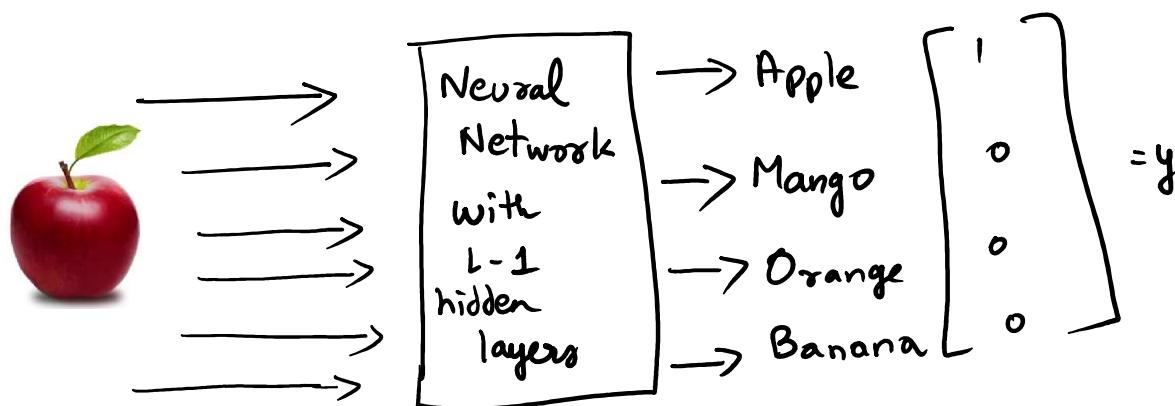


The loss function should capture how much \hat{y}_i deviates from y_i

If $y \in \mathbb{R}^n$ then the squared loss can capture this deviation

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^3 (\hat{y}_{ij} - y_{ij})^2$$

Let's take another example:



We want to classify the image into 1 of K classes

We want to classify the image into 1 of k classes

We could use Savare error loss to capture deviation

Notice that y is a probability distribution

$\therefore \hat{y}$ should be one as well

For this, let O be a Softmax function

$$\hat{y}_j = O(a_L)_j = \frac{e^{a_{L,j}}}{\sum_{i=1}^k e^{a_{L,i}}}$$

$O(a_L)_j$ is the j^{th} element of \hat{y} & $a_{L,j}$ is the j^{th} element of vector a_L

Cross-entropy : $L(\Theta) = \sum_{i=1}^k y_c \log \hat{y}_c$

$$y_c = \begin{cases} 1 & \text{if } c = l \text{ (true class label)} \\ 0 & \text{otherwise} \end{cases}$$

$$L(\Theta) = -\log \hat{y}_l$$

For classification, we use the following objective function

$$\underset{\Theta}{\text{minimize}} \ L(\Theta) = -\log \hat{y}_l$$

or

$$\underset{\Theta}{\text{maximize}} \ -L(\Theta) = \log \hat{y}_l$$

\hat{y}_l is a function of all parameters $\Theta = [w_1, \dots, w_L, b_1, \dots, b_L]$

$$\hat{y}_l = \left[O(w_3 g(w_2 g(w_1 x + b_1) + b_2) + b_3) \right]_l$$

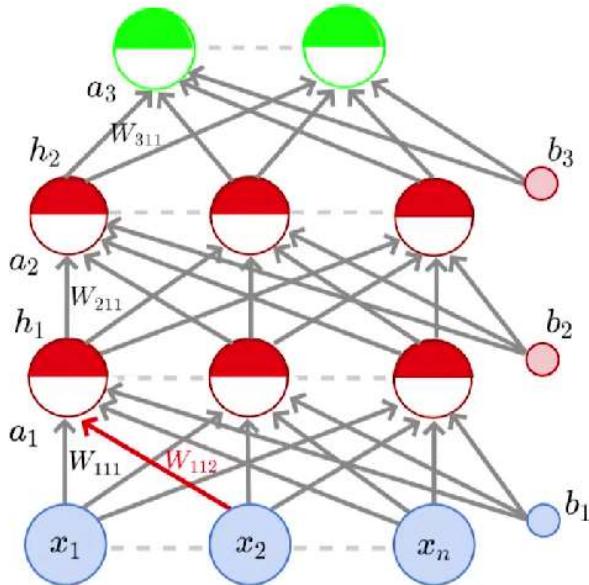
It is the probability that x belongs to l^{th} class
We want it to be as close to 1

It is the probability that x belongs to a class
We want it to be as close to 1

Outputs		
	Real Values	Probabilities
Output Function	Linear	Softmax
Loss Function	Squared Error	Cross Entropy

Backpropagation (Intuition)

In the feed forward neural network, let's focus on one weight w_{112} .



To learn this weight using SGD, we need a formula for $\frac{\partial L(\theta)}{\partial w_{112}}$

We want to calculate all partial derivatives of all weights in that layer in one go & other layers as well.

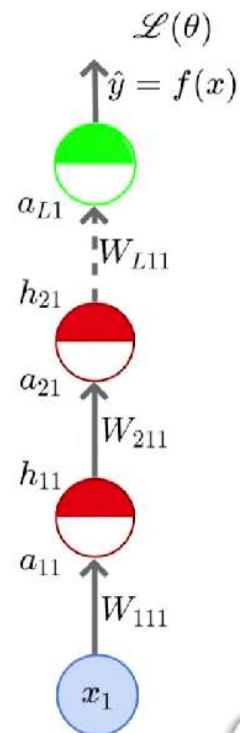
Let us consider the example of a deep but thin network
It is easy to find derivative by chain rule:

$$\frac{\partial L(\theta)}{\partial w_{111}} = \frac{\partial L(\theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{11}} \frac{\partial a_{11}}{\partial h_{21}} \frac{\partial h_{21}}{\partial a_{21}} \frac{\partial a_{21}}{\partial h_{11}} \frac{\partial h_{11}}{\partial a_{11}} \frac{\partial a_{11}}{\partial w_{111}}$$

$$\frac{\partial L(\theta)}{\partial w_{111}} = \frac{\partial L(\theta)}{\partial h_{11}} \frac{\partial h_{11}}{\partial w_{111}} \quad (\text{just compressing chain rule})$$

$$\frac{\partial L(\theta)}{\partial w_{211}} = \frac{\partial L(\theta)}{\partial h_{21}} \frac{\partial h_{21}}{\partial w_{211}}$$

$$\frac{\partial L(\theta)}{\partial h_{11}} \quad \frac{\partial L(\theta)}{\partial h_{21}}$$



$$\frac{\partial L(\theta)}{\partial w_{L''}} = \frac{\partial L(\theta)}{\partial h_L} \frac{\partial h_L}{\partial w_{L''}}$$

Quantities of interest (roadmap for the remaining part):

Gradient w.r.t. output units

Gradient w.r.t. hidden units

Gradient w.r.t. weights and biases

$$\underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial W_{111}}}_{\text{Talk to the weight directly}} = \underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_3}}_{\text{Talk to the output layer}} \underbrace{\frac{\partial a_3}{\partial h_2} \frac{\partial h_2}{\partial a_2}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_1}{\partial W_{111}}}_{\text{and now talk to the weights}}$$

Gradient wrt Output Units

Let us first consider the partial derivative wrt i^{th} output : $\nabla_{a_i} L(\theta)$

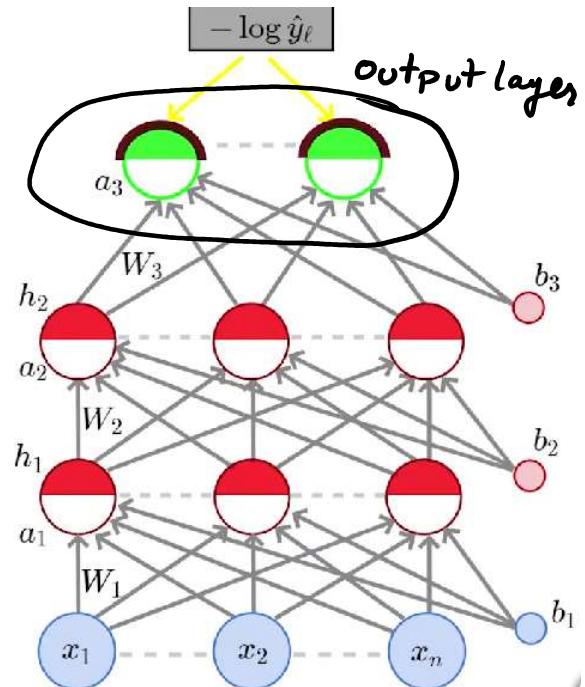
$$L(\theta) = -\log \hat{y}_l \quad (l = \text{true class label})$$

$$\frac{\partial}{\partial \hat{y}_i} (L(\theta)) = \frac{\partial}{\partial \hat{y}_i} (-\log \hat{y}_l)$$

$$= -\frac{1}{\hat{y}_l} \quad \text{if } i=l$$

$$= 0 \quad \text{otherwise}$$

$$\boxed{\frac{\partial}{\partial \hat{y}_i} L(\theta) = -\frac{1}{\hat{y}_l}}$$



We can now talk about the gradient wrt vector \hat{y}

$$\nabla_{\hat{y}} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial \hat{y}_1} \\ \vdots \\ \frac{\partial L(\theta)}{\partial \hat{y}_K} \end{bmatrix} = -\frac{1}{\hat{y}_l} \begin{bmatrix} 1_{l=1} \\ 1_{l=2} \\ \vdots \\ 1_{l=K} \end{bmatrix}$$

this vector will have 1 in only one place for that value of l & 0 everywhere else

These kinds of vectors are called one-hot vectors denoted by e_l

eg: if $l=3$ & $K=5$ then:

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \checkmark$$

denoted by e_l

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\nabla_{y_l} L(\theta) = -\frac{1}{g_l} e_l$$

e_l is a k^{th} dimension vector,
 l^{th} element is 1 & every other
 element is 0

Now,

$$\begin{aligned} \frac{\partial L(\theta)}{\partial a_{L_i}} &= \frac{\partial (-\log \hat{y}_l)}{\partial a_{L_i}} \\ &= \frac{\partial (-\log \hat{y}_l)}{\partial \hat{y}_l} \cdot \frac{\partial \hat{y}_l}{\partial a_{L_i}} \end{aligned}$$

Does \hat{y}_l depend on a_{L_i} ? Yes.

$$\hat{y}_l = \frac{e^{a_{L_i}}}{\sum_{i=1}^k e^{a_{L_i}}}$$

Now we can derive the full expression

$$\begin{aligned} \frac{\partial}{\partial a_{L_i}} -\log \hat{y}_l &= \frac{-1}{\hat{y}_l} \frac{\partial}{\partial a_{L_i}} \hat{y}_l \\ &= \frac{-1}{\hat{y}_l} \frac{\partial}{\partial a_{L_i}} \text{softmax}(a_L)_l \end{aligned}$$

$$\left| \begin{array}{l} a_L = [a_{L_1}, a_{L_2}, \dots, a_{L_k}] \\ \hat{y} = \text{softmax}(a_L) \\ \hat{y} = \left[\frac{e^{L_1}}{\sum e^{L_i}}, \frac{e^{L_2}}{\sum e^{L_i}}, \dots, \frac{e^{L_k}}{\sum e^{L_i}} \right] \end{array} \right.$$

$$= -\frac{1}{\hat{y}_l} \frac{\partial}{\partial a_{L_i}} \frac{e^{a_{L_i}}}{\sum e^{a_{L_i}}}$$

$$= -\frac{1}{\hat{y}_l} \left[\frac{\frac{\partial}{\partial a_{L_i}} e^{a_{L_i}}}{\sum e^{a_{L_i}}} - \frac{e^{a_{L_i}} \left(\frac{\partial}{\partial a_{L_i}} \sum_{i'} e^{a_{L_i'}} \right)}{\left(\sum e^{a_{L_i}} \right)^2} \right]$$

$$\begin{aligned}
&= -\frac{1}{\hat{y}_l} \left[\frac{\frac{\partial a_{l,i}}{\partial a_{l,i}}}{\sum_{i'} e^{a_{l,i'}}} - \frac{e^{-\frac{(a_{l,i}-\hat{y}_l)}{\sum_{i'} e^{a_{l,i'}}}}}{(\sum_{i'} e^{a_{l,i'}})^2} \right] \\
&= \frac{-1}{\hat{y}_l} \left[\frac{\mathbb{1}_{l=i} e^{a_{l,i}}}{\sum_{i'} e^{a_{l,i'}}} - \frac{e^{a_{l,i}}}{\sum_{i'} e^{a_{l,i'}}} \cdot \frac{e^{a_{l,i}}}{\sum_{i'} e^{a_{l,i'}}} \right] \\
&= -\frac{1}{\hat{y}_l} \left[\mathbb{1}_{l=i} \text{softmax}(a_l)_l - \text{softmax}(a_l)_l \text{softmax}(a_l)_i \right] \\
&= -\frac{1}{\hat{y}_l} \left[\mathbb{1}_{l=i} \hat{y}_l - \hat{y}_l \hat{y}_i \right] \\
&= -\frac{1}{\hat{y}_l} \hat{y}_l \left[\mathbb{1}_{l=i} - \hat{y}_i \right]
\end{aligned}$$

$$\nabla_{a_{l,i}} L(\theta) = \frac{\partial L(\theta)}{\partial a_{l,i}} = -[\mathbb{1}_{l=i} - \hat{y}_i]$$

Partial derivative of loss function w.r.t one of the elements of output vector

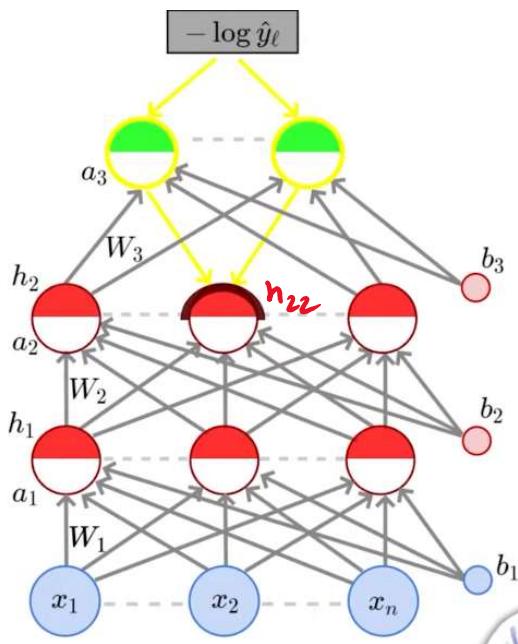
Now for the entire output vector a_l :

$$\nabla_{a_l} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial a_{l,1}} \\ \frac{\partial L(\theta)}{\partial a_{l,2}} \\ \vdots \\ \frac{\partial L(\theta)}{\partial a_{l,k}} \end{bmatrix} = \begin{bmatrix} -(\mathbb{1}_{l=1} - \hat{y}_1) \\ -(\mathbb{1}_{l=2} - \hat{y}_2) \\ \vdots \\ -(\mathbb{1}_{l=k} - \hat{y}_k) \end{bmatrix}$$



$$\nabla_{\alpha_L} L(\theta) = \frac{\partial L(\theta)}{\partial \alpha_L} = -(\hat{e}_L - \hat{y})$$

Gradient wrt Hidden Units



$$a_{i+1} = W_{i+1} h_i + b_{i+1}$$

The derivative of the loss $-\log \hat{y}_l$ depends on h_{22} . There are multiple paths (say k paths) from h_{22} to loss

Let $p(z)$ be loss function of $L(\theta)$

$$z = h_{ij}$$

$$\nabla_m(z) = a_{Lm}$$

If a function $p(z)$ can be written as function of intermediate results $\nabla_i(z)$ then we have:

$$\frac{\partial p(z)}{\partial z} = \sum_m \frac{\partial p(z)}{\partial \nabla_m} \frac{\partial \nabla_m}{\partial z}$$

We are interested in $\nabla_{h_2} L(\theta)$ but since h_2 is a vector, we'll focus on one element of h_2 i.e. h_{22} we'll write it $\boxed{h_{ij}}$ where i is the layer number & j is the neuron

$$\frac{\partial L(\theta)}{\partial h_{ij}} = \sum_{m=1}^k \frac{\partial L(\theta)}{\partial a_{i+1,m}} \cdot \frac{\partial a_{i+1,m}}{\partial h_{ij}}$$

$$= \sum_{m=1}^k \frac{\partial L(\theta)}{\partial a_{i+1,m}} w_{i+1,m,j}$$

We know $a_3 = W_3 h_2 + b_3$

$$\begin{bmatrix} a_{31} \\ a_{32} \end{bmatrix} = \begin{bmatrix} W_{311} & W_{312} & W_{313} \\ W_{321} & W_{322} & W_{323} \end{bmatrix} \begin{bmatrix} h_{21} \\ h_{22} \\ h_{23} \end{bmatrix} + \begin{bmatrix} b_{31} \\ b_{32} \end{bmatrix}$$

$$\begin{bmatrix} a_{31} \\ a_{32} \end{bmatrix} = \begin{bmatrix} w_{311} & w_{312} & w_{313} \\ w_{321} & w_{322} & w_{323} \end{bmatrix}_{(2 \times 3)} \begin{bmatrix} h_{21} \\ h_{22} \\ h_{23} \end{bmatrix}_{(3 \times 1)} + \begin{bmatrix} b_{31} \\ b_{32} \end{bmatrix}_{2 \text{ biases}}$$

only two neurons in current layer

3 neurons in previous layer

$$a_{31}^{i+1} = w_{311} h_{21} + w_{312} h_{22} + w_{313} h_{23} + b_{31}$$

$$\frac{\partial a_{31}}{\partial h_{22}} = 0 + w_{312} + 0 + 0 = w_{312}^j$$

$i \uparrow j$

We know, $\nabla_{a_{i+1}} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial a_{i+1,1}} \\ \frac{\partial L(\theta)}{\partial a_{i+1,2}} \\ \vdots \\ \frac{\partial L(\theta)}{\partial a_{i+1,k}} \end{bmatrix}$

$w_{i+1,:j} = \begin{bmatrix} w_{i+1,1,j} \\ w_{i+1,2,j} \\ \vdots \\ w_{i+1,k,j} \end{bmatrix}$

; j^{th} column of w_{i+1}

$$(w_{i+1,:j})^T \nabla_{a_{i+1}} L(\theta) = \sum_{m=1}^k \frac{\partial L(\theta)}{\partial a_{i+1,m}} w_{i+1,m,j}$$

$$\therefore \text{We have } \frac{\partial L(\theta)}{\partial h_{ij}} = (w_{i+1,:j})^T \nabla_{a_{i+1}} L(\theta)$$

Now we can compute gradient w.r.t any hidden layer h_i

$$\left[\frac{\partial L(\theta)}{\partial h_{ii}} \right] \quad \left[(w_{i+1,:i})^T \nabla_{a_{i+1}} L(\theta) \right]$$

$$\nabla_{h_i} L(\theta) = \begin{bmatrix} \frac{\partial h_{i1}}{\partial \theta} \\ \frac{\partial L(\theta)}{\partial h_{i2}} \\ \vdots \\ \frac{\partial L(\theta)}{\partial h_{in}} \end{bmatrix} = \begin{bmatrix} (\omega_{i+1, \cdot, 1})^T \nabla_{a_{i+1}} L(\theta) \\ (\omega_{i+1, \cdot, 2})^T \nabla_{a_{i+1}} L(\theta) \\ \vdots \\ (\omega_{i+1, \cdot, n})^T \nabla_{a_{i+1}} L(\theta) \end{bmatrix}$$

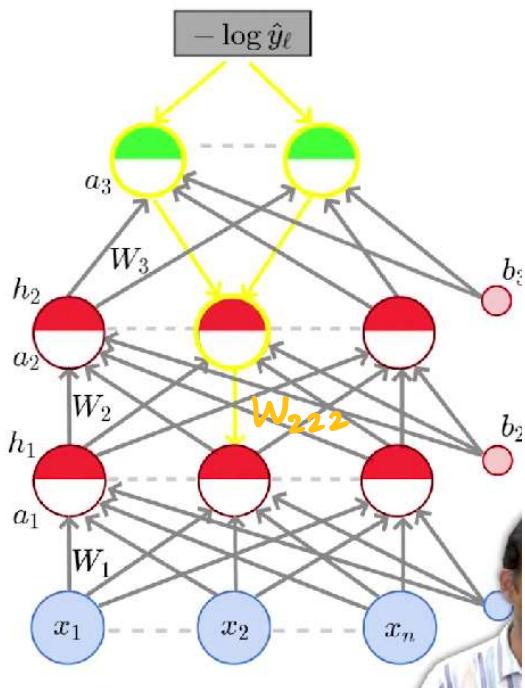
$\nabla_{h_i} L(\theta) = (\omega_{i+1})^T (\nabla_{a_{i+1}} L(\theta))$

Now we need to know how to calculate $\nabla_{a_{i+1}} L(\theta)$
for $i < L-1$

$$\begin{aligned} \frac{\partial L(\theta)}{\partial a_{ij}} &= \frac{\partial L(\theta)}{\partial h_{ij}} \cdot \frac{\partial h_{ij}}{\partial a_{ij}} \\ &= \frac{\partial L(\theta)}{\partial h_{ij}} g'(a_{ij}) \quad [\because h_{ij} = g(a_{ij})] \end{aligned}$$

$$\nabla_{a_i} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial h_{i1}} g'(a_{i1}) \\ \frac{\partial L(\theta)}{\partial h_{i2}} g'(a_{i2}) \\ \vdots \\ \frac{\partial L(\theta)}{\partial h_{in}} g'(a_{in}) \end{bmatrix}$$

Gradient wrt Parameters



Now we want to calculate derivative of loss function wrt one of the weight matrices.

$$\nabla_{W_k} L(\theta)$$

Let's start by computing derivative wrt one element.

$$\frac{\partial L(\theta)}{\partial W_{kij}}$$

$$\text{We know, } a_k = b_k + W_k h_{k-1}$$

$$\frac{\partial L(\theta)}{\partial W_{kij}} = \frac{\partial L(\theta)}{\partial a_{ki}} \frac{\partial a_{ki}}{\partial W_{kij}}$$

$$a_k = \begin{bmatrix} a_{k1} \\ a_{k2} \\ a_{k3} \end{bmatrix} = \begin{bmatrix} W_{k11} & W_{k12} & W_{k13} \\ W_{k21} & W_{k22} & W_{k23} \\ W_{k31} & W_{k32} & W_{k33} \end{bmatrix} \begin{bmatrix} h_{11} \\ h_{12} \\ h_{13} \end{bmatrix} \quad (\text{ignoring bias})$$

$$\therefore a_{k2} = W_{k21} \cdot h_{11} + W_{k22} \cdot h_{12} + W_{k23} \cdot h_{13}$$

$$\frac{\partial a_{k2}}{\partial W_{k22}} = h_{12}^{(k)}$$

$$\frac{\partial L(\theta)}{\partial W_{kij}} = \frac{\partial L(\theta)}{\partial a_{ki}} h_{k-1,j}$$

$$\partial w_{kij}$$

$$\partial a_{ki}$$

∴ for entire weight w_k the gradient will be:

$$\nabla_{w_k} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial w_{k11}} & \frac{\partial L(\theta)}{\partial w_{k12}} & \dots & \frac{\partial L(\theta)}{\partial w_{k1n}} \\ \vdots & \ddots & & \vdots \\ \dots & \dots & \dots & \frac{\partial L(\theta)}{w_{knn}} \end{bmatrix}$$

For example let's take $W \in \mathbb{R}^{3 \times 3}$ then,

$$\nabla_{w_k} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial w_{k11}} & \frac{\partial L(\theta)}{\partial w_{k12}} & \frac{\partial L(\theta)}{\partial w_{k13}} \\ \frac{\partial L(\theta)}{\partial w_{k21}} & \frac{\partial L(\theta)}{\partial w_{k22}} & \frac{\partial L(\theta)}{\partial w_{k23}} \\ \frac{\partial L(\theta)}{\partial w_{k31}} & \frac{\partial L(\theta)}{\partial w_{k32}} & \frac{\partial L(\theta)}{\partial w_{k33}} \end{bmatrix}$$

$$\frac{\partial L(\theta)}{\partial w_{kij}} = \frac{\partial L(\theta)}{\partial a_{ki}} \frac{\partial a_{ki}}{\partial w_{kij}}$$

$$\nabla_{w_k} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial a_{k1}} h_{k-1,1} & \frac{\partial L(\theta)}{\partial a_{k1}} h_{k-1,2} & \frac{\partial L(\theta)}{\partial a_{k1}} h_{k-1,3} \\ \frac{\partial L(\theta)}{\partial a_{k2}} h_{k-1,1} & \frac{\partial L(\theta)}{\partial a_{k2}} h_{k-1,2} & \frac{\partial L(\theta)}{\partial a_{k2}} h_{k-1,3} \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$$\left[\frac{\partial L(\theta)}{\partial a_{k3}} h_{k-1,1} \quad \frac{\partial L(\theta)}{\partial a_{k3}} h_{k-1,2} \quad \frac{\partial L(\theta)}{\partial a_{k3}} h_{k-1,3} \right]$$

$$\nabla_{w_k} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial a_{k1}} \\ \frac{\partial L(\theta)}{\partial a_{k2}} \\ \frac{\partial L(\theta)}{\partial a_{k3}} \end{bmatrix} \cdot \begin{bmatrix} h_{k-1,1} & h_{k-1,2} & h_{k-1,3} \end{bmatrix}$$

$$\boxed{\nabla_{w_k} L(\theta) = \nabla_{a_k} L(\theta) \cdot h_{k-1}^T}$$

Coming to the biases,

$$a_{ki} = b_{ki} + w_{kij} h_{k-1j}$$

$$\frac{\partial L(\theta)}{\partial b_{ki}} = \frac{\partial L(\theta)}{\partial a_{ki}} \cdot \frac{\partial a_{ki}}{\partial b_{ki}}$$

$$\frac{\partial a_{ki}}{\partial b_{ki}} = 1 \quad \therefore \frac{\partial L(\theta)}{\partial b_{ki}} = \frac{\partial L(\theta)}{\partial a_{ki}}$$

Now we can write the gradient wrt b_k

$$\nabla_{b_k} L(\theta) = \nabla_{a_k} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial a_{k1}} \\ \frac{\partial L(\theta)}{\partial a_{k2}} \end{bmatrix}$$

$$\begin{matrix} \frac{\partial L(\theta)}{\partial a_{k2}} \\ \vdots \\ \frac{\partial L(\theta)}{\partial a_{kn}} \end{matrix}$$

Backpropagation: Pseudocode

We have the gradients wrt :

Output layer: $\nabla_{a_L} L(\theta)$

Hidden layers: $\nabla_{h_k} L(\theta), \nabla_{a_k} L(\theta) [1 \leq k < L]$

Weights & Biases: $\nabla_{w_k} L(\theta), \nabla_{b_k} L(\theta) [1 \leq k < L]$

Algorithm

$t \leftarrow 0$

$\text{max_iterations} \leftarrow 1000$

Initialize $\Theta_0 = [w_1^0, w_2^0, \dots, w_L^0, b_1^0, b_2^0, \dots, b_L^0]$

while $t++ < \text{max_iterations}$ do

 calculate $h_1, \dots, h_{L-1}, a_1, \dots, a_L, \hat{y} = \text{forward_propagation}(\Theta_t)$

$\nabla \Theta_t = \text{back_propagation}(h_1, \dots, h_{L-1}, a_1, \dots, a_L, y, \hat{y})$

$\nabla \Theta_{t+1} \leftarrow \Theta_t - \eta \nabla \Theta_t$

end

Forward Propagation: Algorithm

For $k=1$ to $L-1$ do

$a_k = b_k + W_k h_{k-1}$

$h_k = g(a_k)$

end

end

$$a_L = b_L + W_L h_{L-1}$$
$$\hat{y} = \sigma(a_L)$$

Back Propagation

// Compute output gradient

$$\nabla_{a_L} L(\theta) = - (e(y) - \hat{y})$$

For k=L-1 to 1 do

// compute gradient wrt parameters

$$\nabla_{W_k} L(\theta) = \nabla_{a_k} L(\theta) h_{k-1}^T$$

$$\nabla_{b_k} L(\theta) = \nabla_{a_k} L(\theta)$$

// compute gradient wrt layer below

$$\nabla_{h_{k-1}} L(\theta) = W_k^T \nabla_{a_k} L(\theta)$$

// compute gradient wrt layer below (post-activation)

$$\nabla_{a_{k-1}} L(\theta) = \nabla_{h_{k-1}} L(\theta) \odot [\dots, g'(a_{k-1,j}), \dots]$$

end

This is how we can compute g'

Logistic Function

$$g(z) = \sigma(z)$$

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

Tanh Function

$$g(z) = \tanh(z)$$

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

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

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

$$\frac{e^z + e^{-z}}{e^z + e^{-z}}$$

$$g(z) = \frac{1}{1+e^{-z}} \left(\frac{e^z - e^{-z}}{e^z + e^{-z}} \right)$$

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

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

$$g(z) = 1 - (g(z))^2$$

$$g'(z) = g(z)(1-g(z))$$