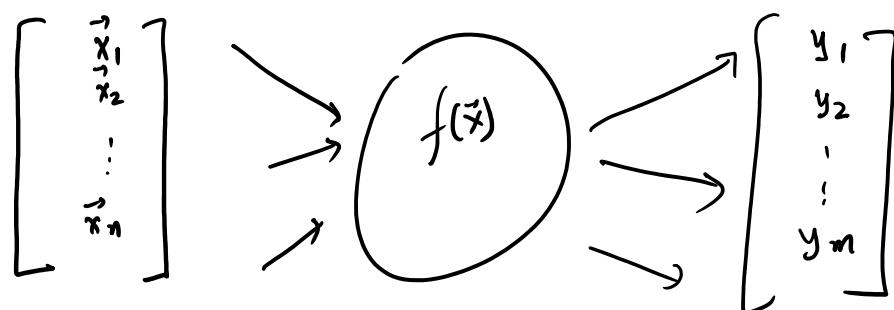
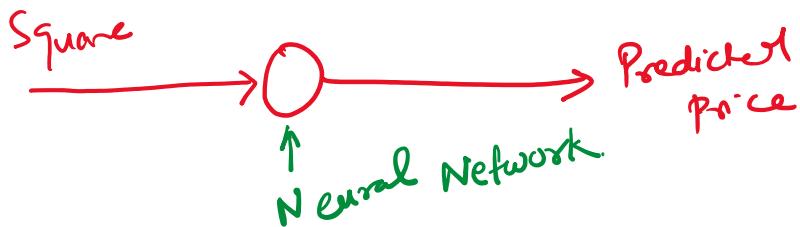
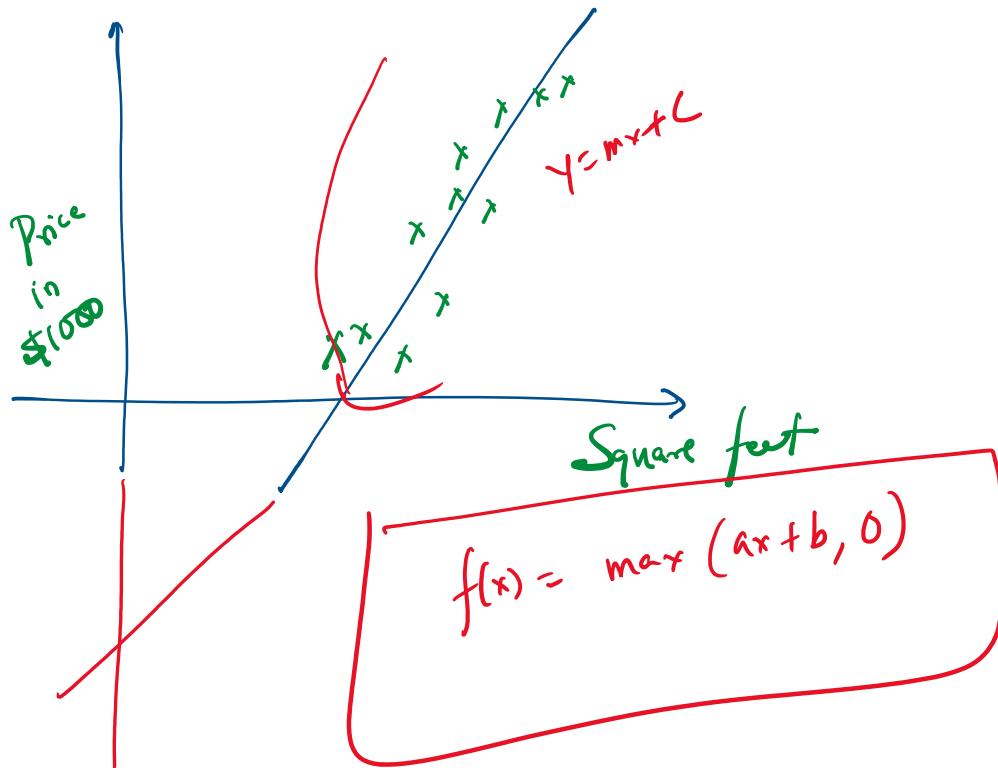
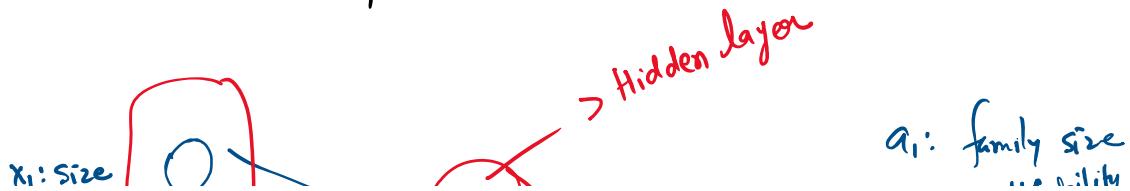


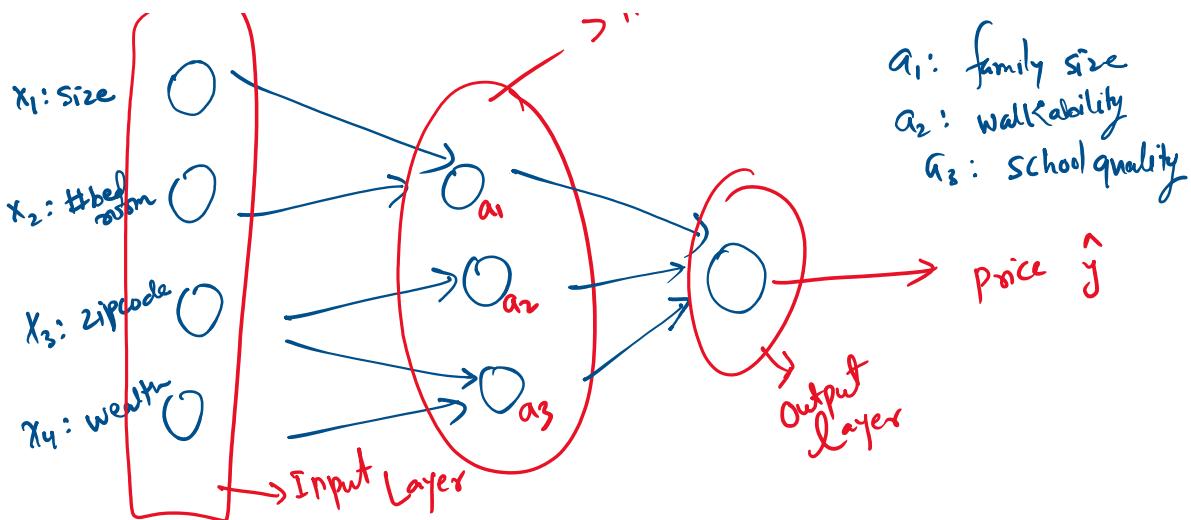
Neural Network

Monday, February 17, 2025 6:16 PM

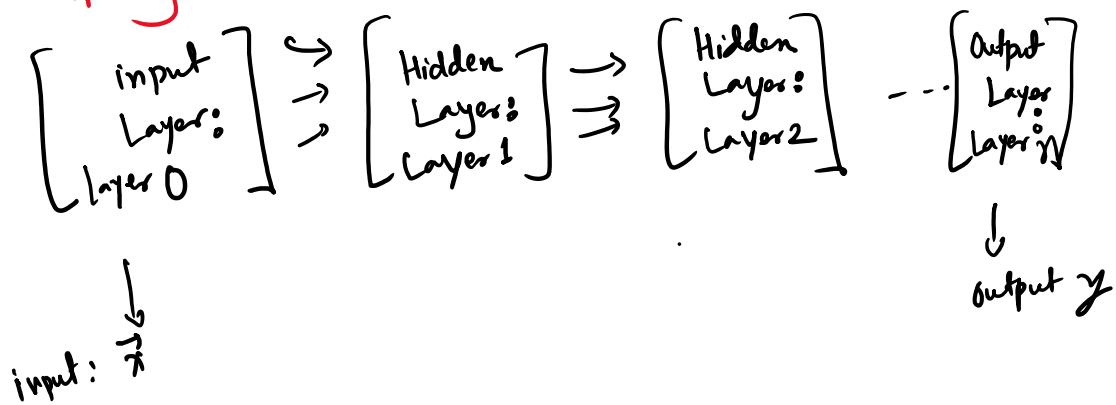


A neural network can process multiple inputs and provides multiple outputs.





A general architecture of a neural network:



When intermediate features are also learnt then we call it end-to-end learning.

Example to drive Motivation

Consider input features: x_1, x_2, \dots, x_3

Output : y

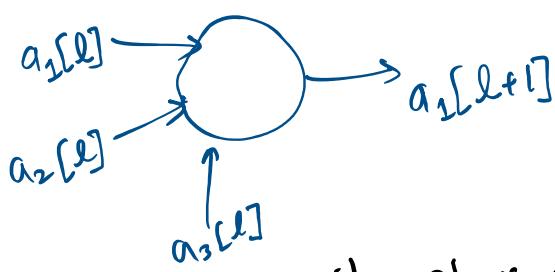
Hidden Layer contains 4 units
Output layer : 1 neurons.

A single hidden unit looks like this:



$l = \text{index of layer}$

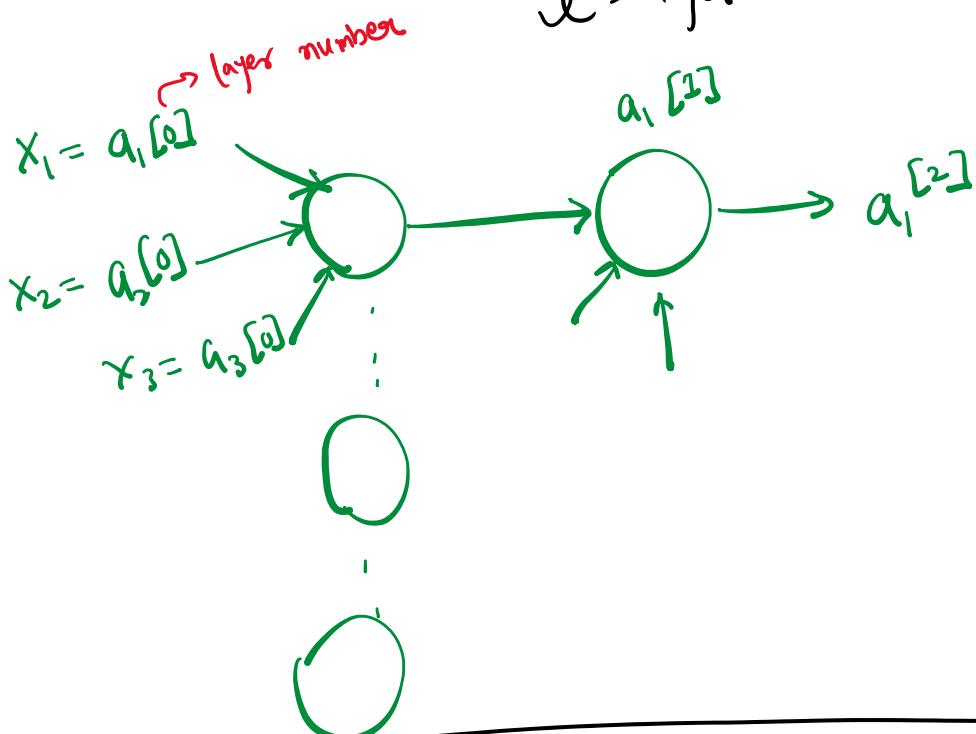
A single neuron ...



l = index of layer

a = activation of neuron.

This particular unit shows one hidden unit in l^{th} layer.



Any $a_j[l]$ is associated with the activation of the j^{th} unit of layer l .

Let's first consider the problem of classification

$$g(\vec{x}) = \frac{1}{1 + \exp(-\vec{w}^T \vec{x})}$$

$$\vec{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$y = \begin{cases} 0 & g(x) < 0.5 \\ 1 & g(x) \geq 0.5 \end{cases}$$

$g(\vec{x})$ can be broken into two different mutation

r17

$g(\vec{z})$ can go into two different computation for this $\vec{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$

$$\textcircled{1} \quad \vec{z} = \vec{w}^T \vec{x} + \vec{b}$$

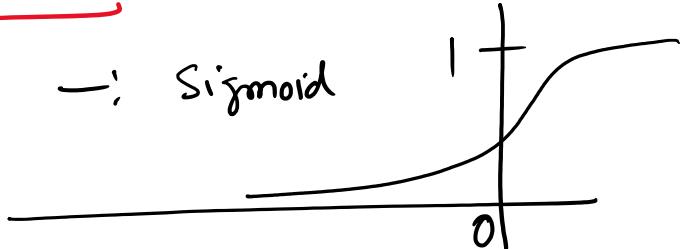
$$\textcircled{2} \quad \vec{a} = \sigma(\vec{z}) \quad \text{where } \sigma(\vec{z}) = \frac{1}{1 + \exp(-\vec{z})}$$

↳ activation function

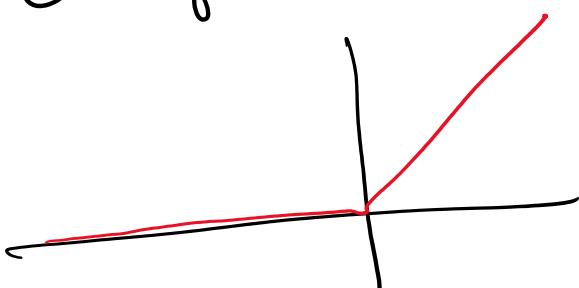
The sigmoid function shown is an example of an activation function.

Activation function

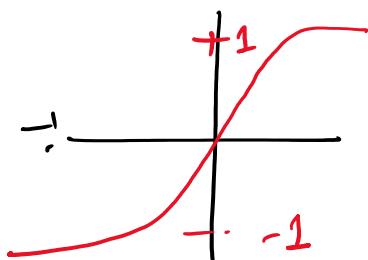
$$\textcircled{1} \quad g(z) = \frac{1}{1 + e^{-z}} \quad \text{=: Sigmoid}$$



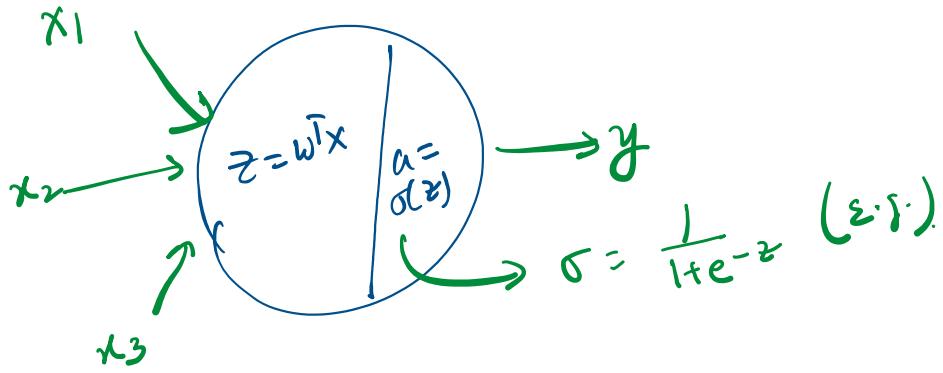
$$\textcircled{2} \quad g(z) = \max(z, 0) \rightarrow \text{Rectified Linear Unit (ReLU)}$$



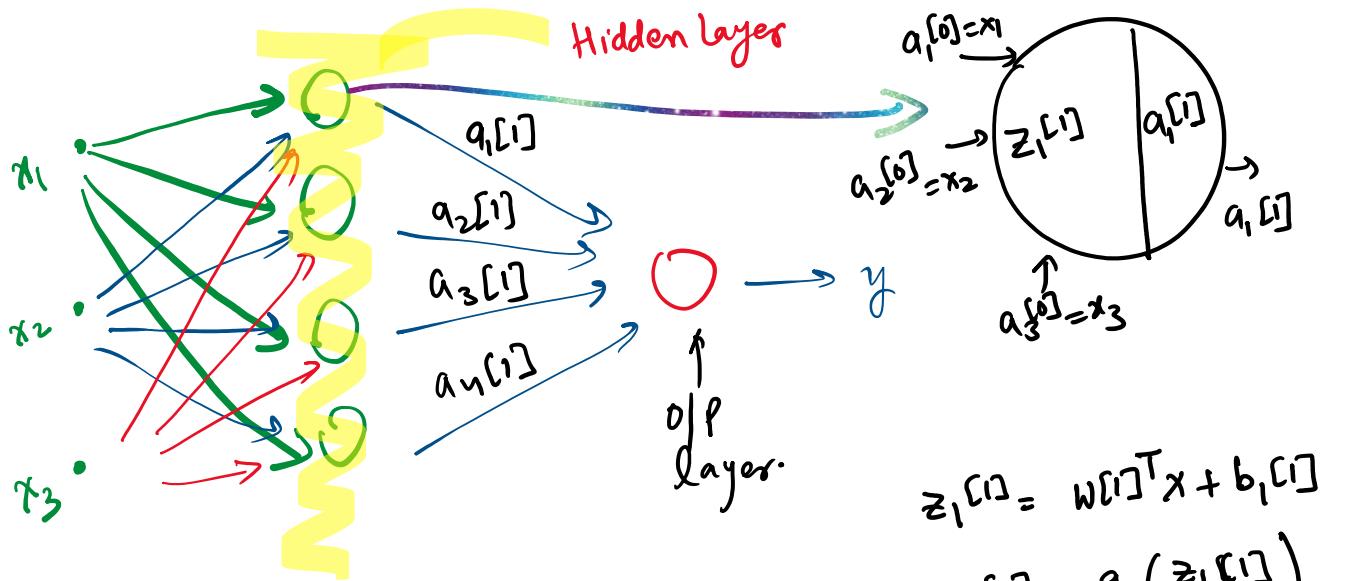
$$\textcircled{3} \quad g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



But in general $g(z)$ has to be nonlinear.



So, in the first hidden unit of the first layer,



Here g is unit function.

$$\begin{aligned} z_1[1] &= w^{[1]T}x + b_1[1] \\ a_1[1] &= g(z_1[1]) \\ &= W^{[1]T}x + b_1[1] \end{aligned}$$

g can be some function.

where $W^{[1]}$ is a matrix of parameter for layer $[1]$

so that $w_{1[1]}^T$ is the first row of this matrix.

$$W^{[1]} = \left[\begin{array}{c} w_{1[1]}^T \\ w_{2[1]}^T \\ \vdots \\ w_{4[1]}^T \end{array} \right]$$

and thus, $w_1^{[1]T} = [w_{11}^{[1]}, w_{21}^{[1]}, w_{31}^{[1]}]$

$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}$$

$$= w_{11}^{[1]} x_1 + w_{21}^{[1]} x_2 + w_{31}^{[1]} x_3 + b_1^{[1]}$$

so $w_1^{[1]} \in \mathbb{R}^3$
 $b_1^{[1]} \in \mathbb{R}^1$ then

For each of the hidden units in layer 1.

$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}$$

$$a_1^{[1]} = g(z_1^{[1]})$$

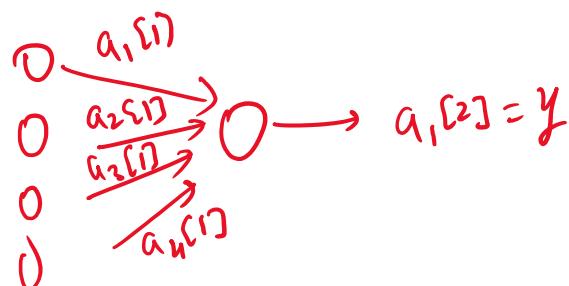
$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}$$

:

$$\vdots$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}$$

$$a_4^{[1]} = g(z_4^{[1]})$$



and then, at O/P

$$z_1^{[2]} = w_1^{[2]T} a^{[1]} + b_1^{[2]}$$

$$\text{and } a_1^{[2]} = g(z_1^{[2]})$$

where $a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix}$

↳ Concatenation of all
the first layer

↳ Concatenation of in
the first layer
activation.

and we are free to choose $g(z)$ to be different depending on the task.

Vectorize

$$z^{[1]} = W^{[1]T} \bar{x} + b^{[1]}$$

$$\begin{bmatrix} z_1^{[1]} \\ \vdots \\ z_n^{[1]} \end{bmatrix} = \left[\begin{array}{c} W_1^{[1]T} \\ W_2^{[1]T} \\ \vdots \end{array} \right] + \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_m \end{bmatrix}$$

$x \in \mathbb{R}^{3m}$

$$z^{[1]} \in \mathbb{R}^{n \times 1} \quad W^{[1]} \in \mathbb{R}^{n \times 3} \quad b^{[1]} \in \mathbb{R}^{n \times 1}$$

What about $g(z)$ as a vector?

Although $g(z)$ is a scalar operation,

for $g(z) = \frac{1}{1 + \exp(-z)}$, we can define

the vector as elementwise computation which can be executed concurrently.

Then the output layer comes out as

$$z^{[2]} = \underbrace{W^{[2]} a^{[1]}}_{1 \times 4} + \underbrace{b^{[2]}}_{1 \times 1}$$

$$a^{[2]} = \underbrace{g(z^{[2]})}_{1 \times 1}$$

$$\underbrace{\mathbf{w}}_{1 \times 1} - \underbrace{\mathbf{b}}_{1 \times 1}$$

When including more than one training sample
 (Consider 3 training samples)

$$z^{1} = w^{[1]} \vec{x}(1) + b^{[1]}$$

$$z^{[1](2)} = w^{[1]} \vec{x}(2) + b^{[1]}$$

$$z^{[1](3)} = w^{[1]} \vec{x}(3) + b^{[1]}$$

\vec{x} contains all features.

In matrix representation:

$$\mathbf{z}^{[1]} = \begin{bmatrix} | & | & | \\ z^{1} & z^{[1](2)} & z^{[1](3)} \\ | & | & | \end{bmatrix}$$

$$\vec{\mathbf{X}} = \begin{bmatrix} | & | & | \\ \vec{x}(1) & \vec{x}(2) & \vec{x}(3) \\ | & | & | \end{bmatrix}$$

Then

$$\begin{aligned} \mathbf{z}^{[1]} &= \begin{bmatrix} | & | & | \\ z^{1} & z^{[1](2)} & z^{[1](3)} \\ | & | & | \end{bmatrix} \\ &= w^{[1]} \begin{bmatrix} | & | & | \\ \vec{x}(1) & \vec{x}(2) & \vec{x}(3) \\ | & | & | \end{bmatrix} + b^{[1]} \end{aligned}$$

Training Samples are columns.

$$z^{[i]} = w^{[i]} X + b^{[i]}$$

To train, the process is :

- ① initialize $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$ to small random numbers.
- ② for all $i \in [1, n]$, compute the op probability from $a^{[2]}(i)$ and then

then

$$\max_{w,b} \sum_{i=1}^n \left(y^{(i)} \log a^{[2]}(i) + (1-y^{(i)}) \log (1-a^{[2]}(i)) \right)$$

↳ use gradient ascent
find w, b that maximizes the above cross entropy function means training a neural network.

Back propagation

Neural network has 2 related components

① Network architecture

⑤ # hidden layers.

⑥ # neurons

- .. - number of neurons

- (b) # neurons
- (c) # connections b/w neurons
- (d) activation functions

② Parameters / Weights

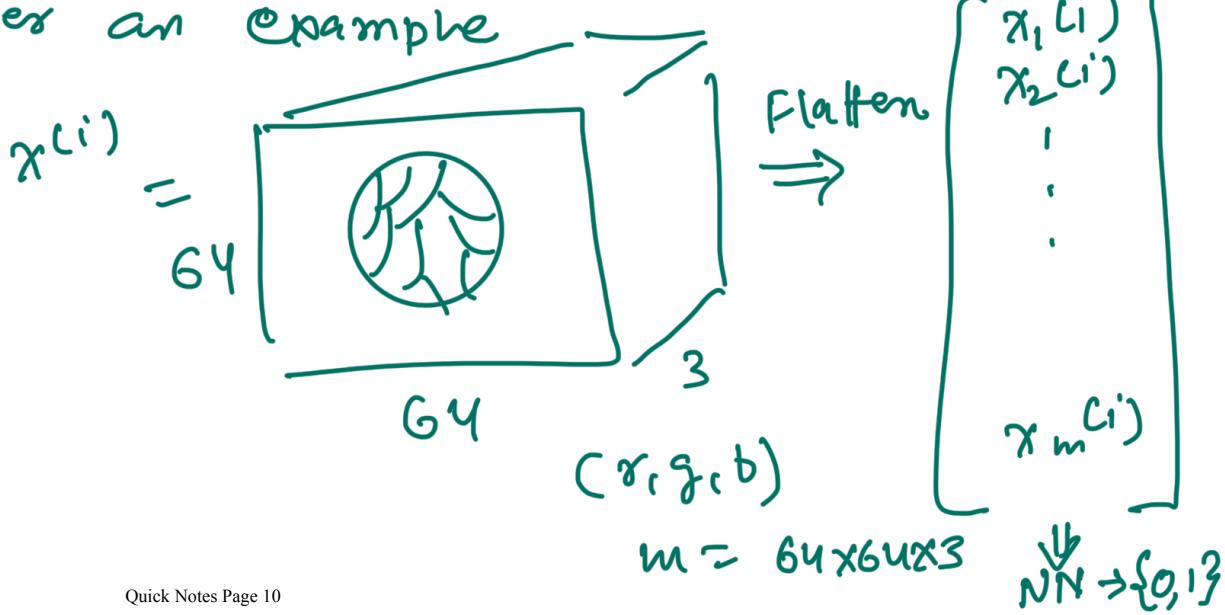
For individual neurons (aka, for individual optimizations), we've determined how to select parameters.

But activations of layer l change the training set for layer $l+1$. So how do we find w, b !

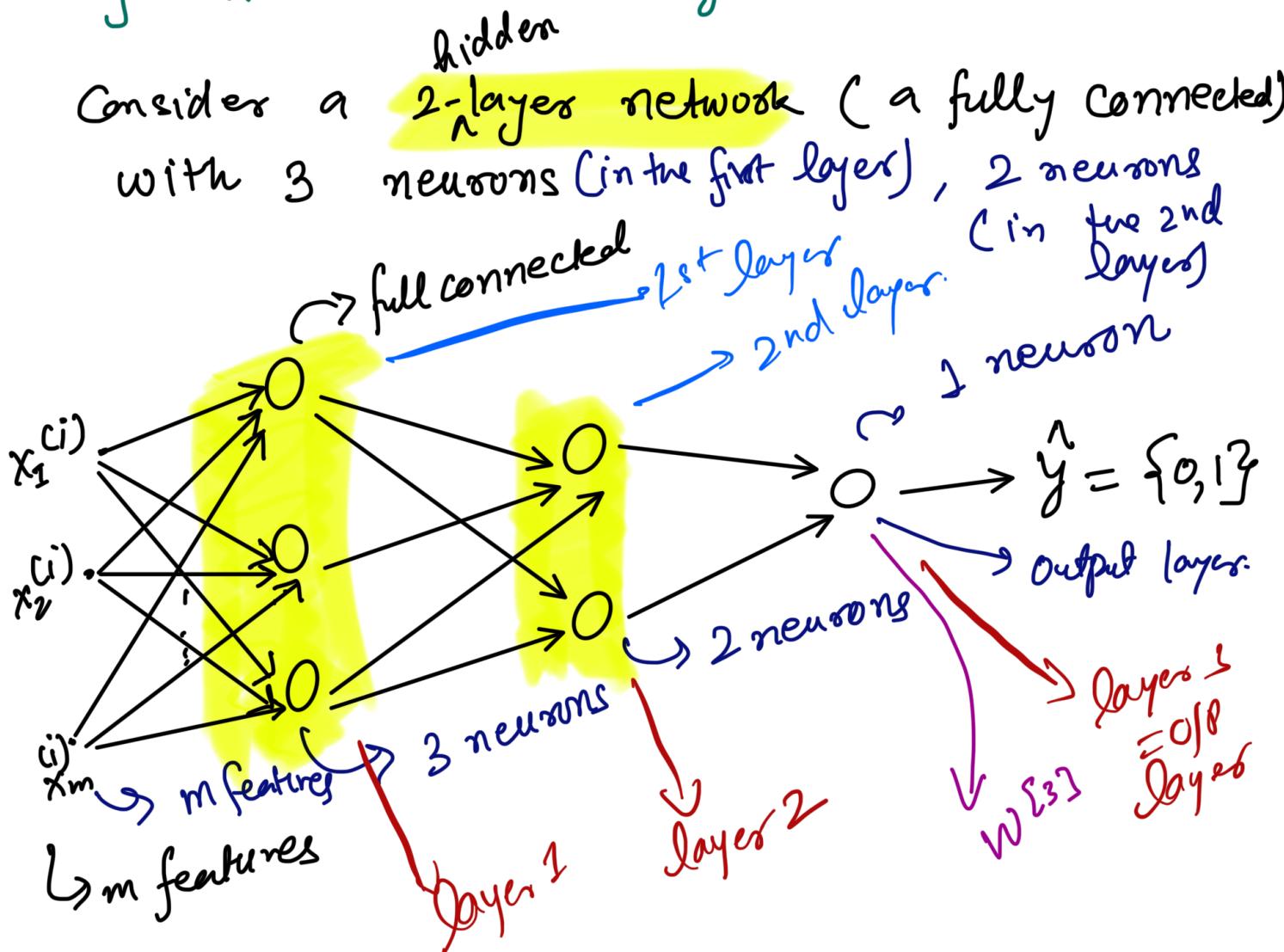
We need a new kind of solution, that can help in finding the concatenated solution:

One such solution is the "Back propagation algorithm" because in this weights are updated backward from the last layer to the first layer.

Consider an example



For this kind of data, it will require a lot of neurons (Hidden Layer).



Let's write the parameters by hand:

$$z^{[1]} = w^{[1]} \cdot x^{(i)} + b^{[1]}$$

$$a^{[1]} = g(z^{[1]})$$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = g(z^{[2]})$$

$$z^{[3]} = w^{[3]} a^{[2]} + b^{[3]}$$

$$\hat{y}^{(i)} = a^{[3]} = g(z^{[3]})$$

$$z^{[1]}, a^{[1]} \in \mathbb{R}^{3 \times 1}$$

$$z^{[2]}, a^{[2]} \in \mathbb{R}^{2 \times 1}$$

$$z^{[3]}, a^{[3]} \in \mathbb{R}^{1 \times 1}$$

What about the dimension of W

$$\underbrace{z^{[1]}}_{3 \times 1} = \underbrace{W^{[1]} x^{[1]}}_{(3 \times ?) \times 1} + \underbrace{b^{[1]}}_{3 \times 1}$$

$$\begin{matrix} \swarrow \\ 3 \end{matrix} \quad \begin{matrix} \searrow \\ m \end{matrix}$$

$$\hookrightarrow W^{[1]} \in \mathbb{R}^{3 \times m}$$

Similarly,

$$W^{[2]} \in \mathbb{R}^{2 \times 3} \quad b^{[2]} \in \mathbb{R}^{2 \times 1}$$
$$W^{[3]} \in \mathbb{R}^{1 \times 2} \quad b^{[3]} \in \mathbb{R}^{1 \times 1}$$

In Total:

Layer 1: $W^{[1]}, b^{[1]}$ $3 \cdot m + 3 \cdot 1 = 3m + 3$

Layer 2: $W^{[2]}, b^{[2]}$ $2 \cdot 3 + 2 \cdot 1 = 8$

Layer 3 (Output Layer):

$$W^{[3]}, b^{[3]} \quad 1 \cdot 2 + 1 \cdot 1 = 3$$

$$3m + 3 + 8 + 3$$

$$= 3m + 14 \text{ Parameters}$$

How do we initialize these parameters?

Can be set them to 0?

$$\begin{aligned}z^{[i]} &= w^{[i]} \cdot x^{(i)} + b^{[i]} \\&= [0] \cdot x^{(i)} + 0 = 0\end{aligned}$$

In practice, we do initialization near zero using a normal distribution, $\mathcal{N}(0, 0.1)$ with a small standard deviation.

Our loss is

$$L(\hat{y}, y) = -[(1-y) \log(1-\hat{y}) + y \log \hat{y}]$$

$\log(1-(\hat{y}+\epsilon))$ $\hat{y}+\epsilon$

$$\text{As, } \hat{y} = g(z^{[2]})$$

So, we see that it will never be exactly 0 or 1.

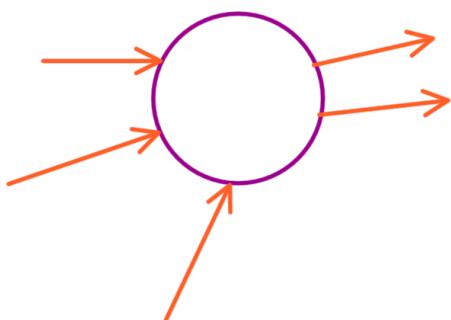
Hence, we see that we never calculate $\log(1)$, rather we calculate $\log(1-\epsilon)$,

$$0 < \epsilon \ll 1$$

With this idea of our loss function, and what we learnt in regression, what can we do about N.N.?

Let's look at the update function during gradient descent.

Consider that we are at layer l :



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

$$a^{[l]} = g(z^{[l]})$$

at layer l .

From our logistic regression

$$\theta_j = \theta_j - \alpha \frac{\partial L}{\partial \theta_j}$$

$$W^{[l]} = W^{[l]} - \alpha \frac{\partial L}{\partial W^{[l]}} \quad \approx$$

$$b^{[l]} = b^{[l]} - \alpha \frac{\partial L}{\partial b^{[l]}}$$

For linear regression

$$\Theta^T = [w_0, w_1, \dots, w_n]$$

$$x = [1, \underbrace{x}_b, x_2, \dots, x_n]$$

Then, if we can calculate $\frac{\partial L}{\partial w^{[l]}}$ & $\frac{\partial L}{\partial b^{[l]}}$

then we can use gradients-based method.

Pick layer $l=3$ (last layer)

$$L = -[(1-y) \log(1-\hat{y}) + y \log(\hat{y})]$$

$$\hat{y} = g(z^{[2]}) \quad g(\cdot) = \sigma(z^{[2]})$$

a sigmoid

$$\begin{aligned} \frac{\partial L}{\partial w^{[3]}} &= - \frac{\partial}{\partial w^{[2]}} \left[(1-y) \log(1-g(z^{[3]})) + y \log(g(z^{[3]})) \right] \\ &= -(1-y) \left[\frac{\partial}{\partial w^{[2]}} \log(1-g(z^{[3]})) \right] \\ &\quad - y \frac{\partial}{\partial w^{[3]}} \log(g(z^{[2]})) \end{aligned}$$

(1)

Some notes

$$\textcircled{1} \quad \sigma' = \sigma(1-\sigma) \quad \text{for } \sigma = \text{sigmoid}$$

$$\textcircled{2} \quad \frac{\partial}{\partial x} \log(f(x)) = \frac{1}{f(x)} \cdot f'(x)$$

$$\textcircled{3} \quad f(x) = h(g(x)) \Rightarrow f'(x) = h'(g(x)) \cdot g'(x)$$

with this in mind:

$$\begin{aligned} \frac{\partial}{\partial w^{[3]}} \log(1-g(z^{[3]})) &= \frac{1}{1-g(z^{[3]})} \cdot \\ &\quad \cdot \frac{\partial(1-g(z^{[3]}))}{\partial w^{[3]}} \frac{\partial z^{[2]}}{\partial w^{[2]}} \end{aligned}$$

$$= \frac{1}{1-g(z^{[2]})} (-1) \underbrace{\frac{\partial g(z^{[2]})}{\partial w^{[3]}}}_{g'(z^{[2]})} \cdot \frac{\partial z^{[2]}}{\partial w^{[3]}}$$

Rewriting ①

$$-(-y)(1A) - y(1B)$$

let's calculate a few items:

$$z^{[3]} = w^{[3]}a^{[2]} + b^{[3]}$$

$$\checkmark \quad \frac{\partial z^{[3]}}{\partial w^{[3]}} = \frac{\partial (w^{[3]}a^{[2]} + b^{[3]})}{\partial w^{[3]}}$$

$$= \frac{\partial (w^{[3]} \cdot a^{[2]})}{\partial w^{[3]}} = a^{[2]T}$$

$$\checkmark \quad \frac{\partial z^{[3]}}{\partial b^{[3]}} = \frac{\partial (w^{[3]} \cdot a^{[2]} + b^{[3]})}{\partial b^{[3]}} = 1$$

$$\text{and } g'(z^{[2]}) = g(z^{[2]}) (1-g(z^{[2]}))$$

①A becomes:

$$= (-1) \frac{1}{1-g(z^{[2]})} g(z^{[2]}) (1-g(z^{[2]})) - a^{[2]T}$$

$$= -g(z^{[2]}) a^{[2]T}$$

(B)

becomes.

$$\frac{1}{g(z^{[2]})} \cancel{g(z^{[2]})} (1 - g(z^{[2]})) a^{[2]T}$$

$$= (1 - g(z^{[2]})) a^{[2]T}$$

Hence eqn ①

$$\frac{\partial L}{\partial W^{[2]}} = -(1-y)(-1) g(z^{[2]}) a^{[2]T}$$

$$- y (1 - g(z^{[3]})) a^{[2]T}$$

$$= [g(z^{[3]}) - y \cancel{g(z^{[2]})} - y + y \cancel{g(z^{[2]})}] \cdot a^{[2]T}$$

$$= [g(z^{[3]}) - y] a^{[2]T} \quad \text{(Note: } a^{[3]} = g(z^{[3]})\text{)}$$

$$\textcircled{1} \Rightarrow \frac{\partial L}{\partial W^{[3]}} = (a^{[2]} - y) a^{[2]T}$$

Given this, can we calculate $\frac{\partial L}{\partial W^{[2]}}$?

Chain Rule:

$$\frac{\partial L}{\partial W^{[3]}} = \frac{\partial L}{?} \cdot \frac{?}{\partial W^{[2]}} \cdot ?$$

$$a^{[2]} = g(z^{[3]})$$

Some important relationship:

① $z^{[3]}$ is related to $a^{[2]}$, $\frac{\partial z^{[3]}}{\partial a^{[2]}}$

$$\underline{(z^{[3]} = w^{[2]} \cdot a^{[2]} + b^{[3]})}$$

② $a^{[2]}$ depends on $z^{[2]}$

$$a^{[2]} = g(z^{[2]}) \Rightarrow \frac{\partial a^{[2]}}{\partial z^{[2]}}$$

③ $z^{[2]}$ depends on $w^{[2]}$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \Rightarrow \frac{\partial z^{[2]}}{\partial w^{[2]}}$$

④ h depends on $a^{[3]}$ $\rightarrow \frac{\partial h}{\partial a^{[3]}}$

⑤ $a^{[3]}$ depends on $z^{[3]}$. $\rightarrow \frac{\partial a^{[3]}}{\partial z^{[3]}}$

So

$$\frac{\partial h}{\partial w^{[2]}} = \underbrace{\frac{\partial h}{\partial a^{[3]}}} \cdot \underbrace{\frac{\partial a^{[3]}}{\partial z^{[3]}}} \cdot \underbrace{\frac{\partial z^{[3]}}{\partial a^{[2]}}} \cdot \underbrace{\frac{\partial a^{[2]}}{\partial z^{[2]}}} \cdot \underbrace{\frac{\partial z^{[2]}}{\partial w^{[2]}}}$$

(IV) (V) (I) (II) (III)

Since

$$\frac{\partial h}{\partial w^{[2]}} = (a^{[2]} - y) a^{[2] T}$$

Recall (IA) $\frac{\partial z^{[3]}}{\partial w^{[3]}} = a^{[2] T}$

$$\frac{\partial L}{\partial w^{(2)}} = (a^{(2)} - y) \frac{\partial z^{(2)}}{\partial w^{(2)}}$$

$$\frac{\partial L}{\partial z^{(3)}} = \frac{\partial L}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(2)}}$$

$$\frac{\partial L}{\partial w^{(2)}} = (a^{(3)} - y) \cdot \frac{\partial z^{(3)}}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial w^{(1)}}$$

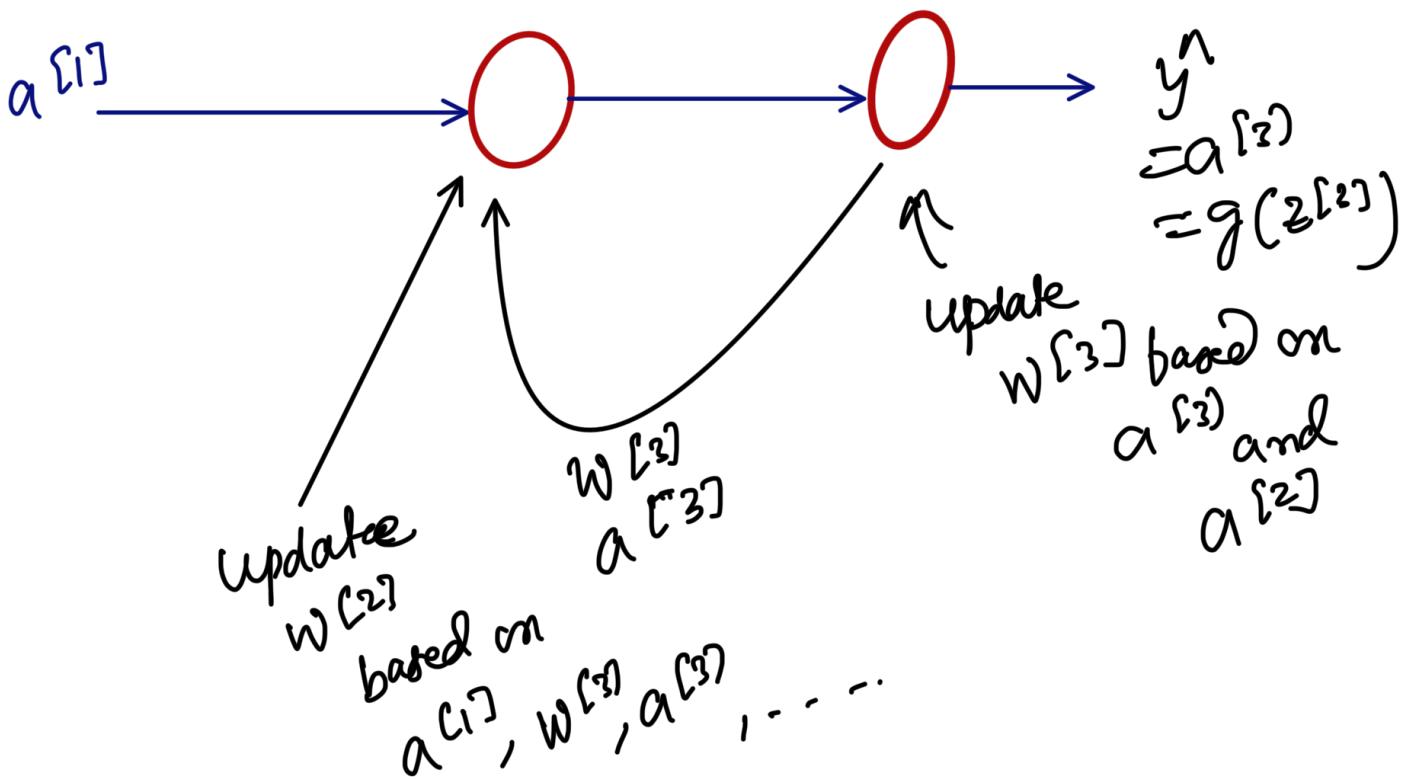
↗ $g'(z^{(2)})$
↙ $a^{(1)T}$
↙ $a^{(1)T}$ from ΔA

$$w^{(3)} \left(\frac{\partial z^{(2)}}{\partial a^{(2)}} = \frac{\partial}{\partial a^{(2)}} (w^{(3)} a^{(2)} + b^{(3)}) \right)$$

$$\Rightarrow \frac{\partial L}{\partial w^{(2)}} = (\underbrace{a^{(3)} - y}_{1 \times 1}) \underbrace{w^{(2)T}}_{1 \times 2} \underbrace{g'(z^{(2)})}_{2 \times 1} \underbrace{a^{(1)T}}_{1 \times 3}$$

$$\Rightarrow \frac{\partial L}{\partial w^{(2)}} = \underbrace{w^{(3)T}}_{2 \times 1} \cdot \underbrace{g'(z^{(2)})}_{2 \times 1} \cdot \underbrace{(a^{(3)} - y)}_{1 \times 1} \cdot \underbrace{a^{(1)T}}_{1 \times 3}$$

Or in essence, changes to the cost w.r.t
 $w^{(2)}$ are based on the downstream
weights $w^{(3)}$



But $w^{[2]}$ and $a^{[2]}$ propagate backward to do this optimization!

Full Gradient descent

For any layer

$$w^{[l]} = w^{[l]} - \alpha \frac{\partial J}{\partial w^{[l]}}$$

$$\text{where } J = \frac{1}{m} \sum_{i=1}^m L(a^{[i]})$$

$m = \# \text{features}$
 $i = \text{training samples.}$

In this kind of scenario,
 $L^{(i)}$ will be noisy

J will be computationally difficult.

Solution: Minibatch Gradient Descent

$$J_{MB} = \frac{1}{B} \sum_{i=1}^B L^{(i)}$$

where $B \ll m$ is the no. of samples in the minibatch

Momentum

For any layer l

$$\mathcal{V}_{dw}[l] = \beta \mathcal{V}_{dw}[l] + (1 - \beta) \frac{\partial J}{\partial w[l]}$$

↗ momentum parameter

$$w[l] = w[l] - \alpha \mathcal{V}_{dw}[l].$$

↓ learning rate