# units.

## Network of units.

Layered architecture.



$X$

input
layer

Output
layer

$Y$

Neural
Network.

$$Y = f(x)$$

## Universal Approximators.

Unit



$g(net)$

Thresholded unit

Perpec
Perceptron

Sigmoid

tanh =

RELU = ( —— )

$g \rightarrow \mathbb{1}( \ )$

$g \rightarrow \sigma( \ )$

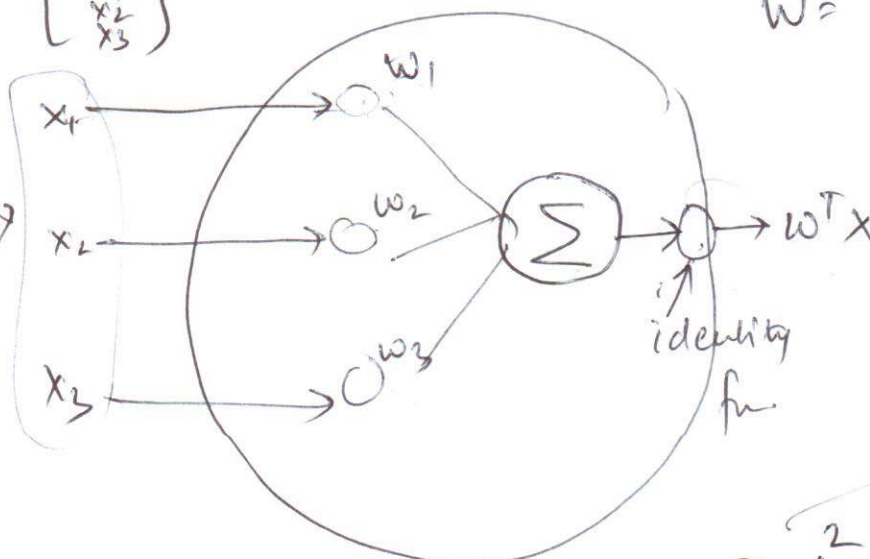$= \dfrac{1}{1 + exp(-)}$

$L$

Feed forward NN Mode

Training Mode
Back Propagation

$X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$



$W = \begin{bmatrix} W_1 \\ W_2 \\ W_3 \end{bmatrix}$

$\Sigma = W^T x$

$= W_1 x_1 + W_2 x_2 + W_3 x_3$

identity fn.

$L(W) = \dfrac{1}{2} (y - \underline{W^T x})^2 \longrightarrow$ We have only one data point in training

$= \dfrac{1}{2} (y - W_1 x_1 - W_2 x_2 - W_3 x_3)^2$

$\boxed{\dfrac{1}{2} \Sigma (y_i - W^T x_i)^2}$

$\dfrac{\partial L}{\partial W_1} = (y - W_1 x_1 - W_2 x_2 - W_3 x_3)(-x_1)$

$= -(y - W^T x) x_1$

$\dfrac{\partial L}{\partial W_2} = -(y - W^T x) x_2$

$\dfrac{\partial L}{\partial W_3} = -(y - W^T x) x_3$

$\nabla = \begin{bmatrix} \dfrac{\partial L}{\partial W_1} \\ \dfrac{\partial L}{\partial W_2} \\ \dfrac{\partial L}{\partial W_3} \end{bmatrix}$

$$W \leftarrow W^{old} - \eta \nabla$$

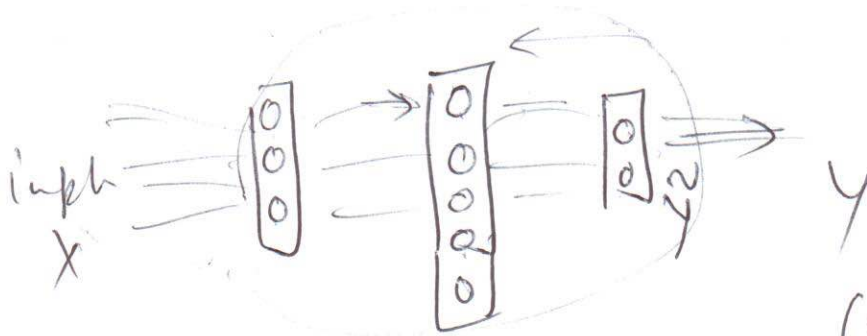$$W_1 = W_1^{old} - \eta \left(\frac{\partial L}{\partial w_1}\right)$$

$$= W_1^{old} + \eta \, (y - w^T x) \, x_1$$

→ Error
between true label $(y)$
and $(w^T x) \rightarrow$ predictia.

$$W_2 = W_2^{old} + \eta \, (y - w^T x) \, x_2$$

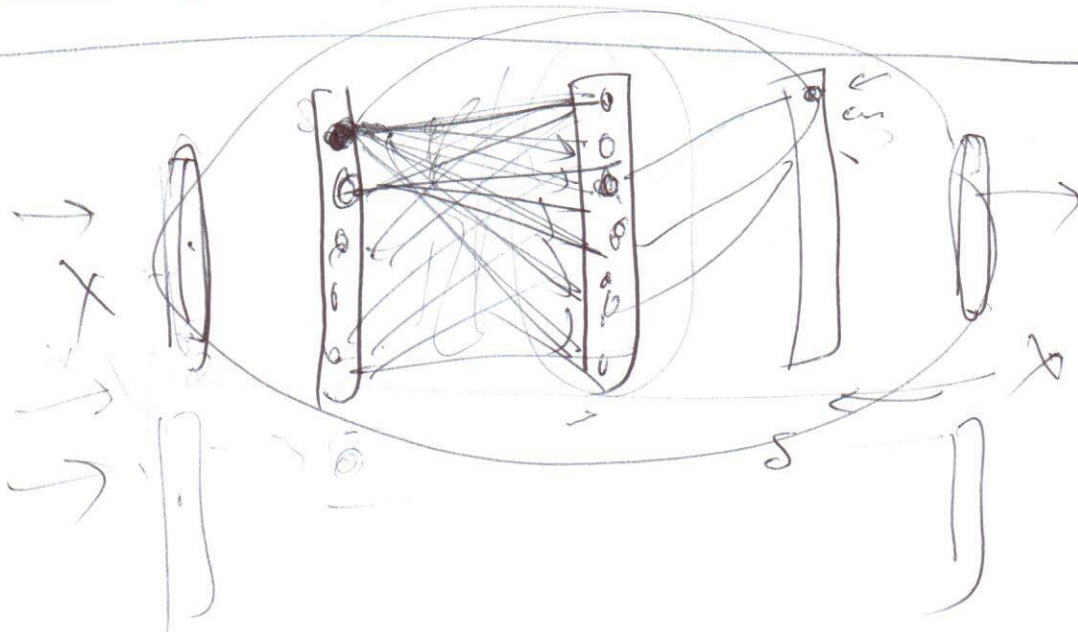$$W_3 = W_3^{old} + \eta \, (y - w^T x) \, x_3$$

→ learning
rate
(eta) 
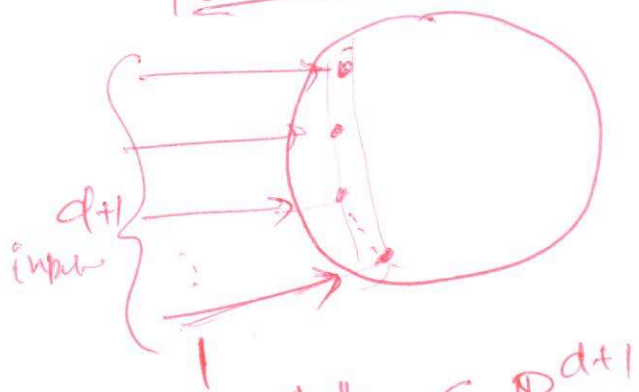Step-size



inpt
$X$

$Y$

$(y - \tilde{y})$

error

Back propogation

# Neural Networks.

1 hidden layer $\rightarrow$ m units.
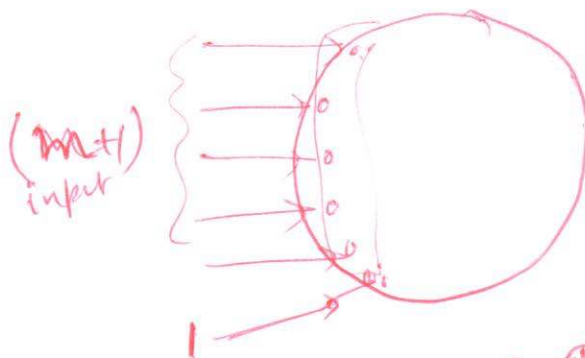
1 output layer $\rightarrow$ k units (K classes)

Output unit

### Hidden unit

$d+1$ input

1

$$W_j \in \mathbb{R}^{d+1}$$

weight vector for $j^{th}$ hidden unit.

$(m+1)$ input

1

$$W_\ell \in \mathbb{R}^{m+1}$$

weight vector for the $\ell^{th}$ output unit.

$\boxed{1 \leq j \leq m}$ $\rightarrow$ total # hidden units.

$\boxed{1 \leq \ell \leq k}$ $\rightarrow$ total # classes

Let **x** be the input. $x \in \mathbb{R}^d$ (we will add a bias term to this)

After adding bias term $\rightarrow$ $x \in \mathbb{R}^{d+1}$

$$z_1 = \sigma(w_1^T x)$$

$$z_2 = \sigma(w_2^T x)$$

$$z_m = \sigma(w_m^T x)$$

activation

output nodes

Consider the first output node.



$$O_1 \quad \sigma(W_1^T \mathbf{z})$$

$$\boxed{1 \mid z_1 \mid z_2 \cdots \mid z_l} \quad \to z_l$$

$$O_2 \quad \sigma(W_2^T z)$$
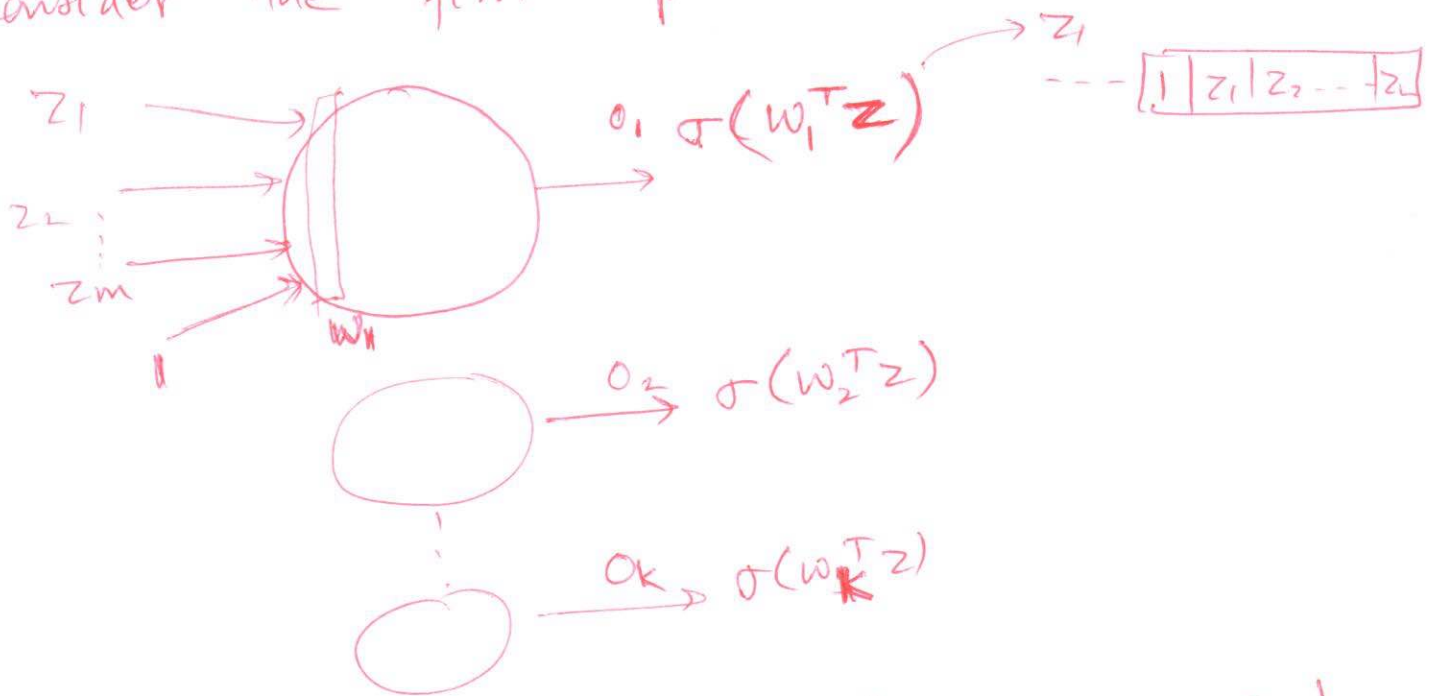
$$O_k \quad \sigma(W_k^T z)$$

Predicted Class $\longrightarrow \underset{K}{\mathrm{argmax}} \left( O_1, O_2, \cdots O_k \right)$

$$\tilde{O_1} = \frac{O_1}{\sum O_l}$$

---

Loss function for LinRegression

$$J(W) = \frac{1}{2} \sum \left( y_i - W^T x_i \right)^2$$

---

Grad Descent.

$$W^{(t)} \quad \leftarrow W^{(t-1)} - \eta \, \nabla J(W)$$

$$\boxed{W_i^{(t)} \leftarrow W_i^{(t-1)} - \eta \left( \frac{\partial J}{\partial W_i} \right) \longrightarrow}$$
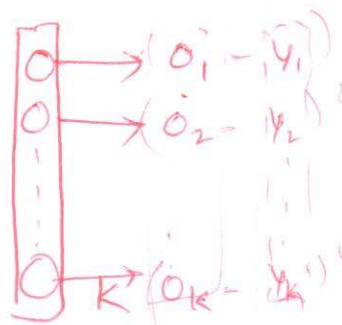
Hidden          Output



At output layer

$$\frac{\partial J}{\partial w_{kj}}$$

At ~~the~~ hidden layer

$$\frac{\partial J}{\partial w_{ji}}$$

??

---

At the Output Layer:



$(O_1 - y_1)^2$
$(O_2 - y_2)^2$
$K (O_k - y_k)^2$

$y \in \{0, \dots k\}$

1 of k Encoding (Dummy).

$\bar{y} \in R^k$

Vector with k entries

For example, if. $K = 10$.

Let $y = 4$

# Neural Networks

## Squared Loss

$$J = \sum_{i=1}^{N} J_i$$

$$\frac{\partial J}{\partial w_{ij}} = \sum_{i=1}^{N} \frac{\partial J_i}{\partial w_{ij}}$$

$$= \sum_{i=1}^{N} \sum_{\ell=1}^{K} (y_{i\ell} - O_{i\ell})^2$$

True label for $i^{th}$ training point. $\longrightarrow$ Output of the NN at the $\ell^{th}$ output unit for the $i^{th}$ training point.

For any $y_i$ ($\in \{1, \cdots, k\}$)

| 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|

K length vector

Such that $y_{i\ell} = 1$ if $\ell = y_i$
$y_{i\ell} = 0$ else

e.g. $y_1 = 4$

| 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|

1 of k encoding.

e.g. ~~k=10~~ k=4

Training point $\xrightarrow{} x_i$

$$
\begin{aligned}
O &\to O_1 &&= &0 \\
O &\to O_2 &&= &0 \\
O &\to O_3 &&= &1 \quad \leftarrow y_i = 3 \to \text{label}\\
O &\to O_4 &&= &0
\end{aligned}
$$

Output layer

$i$ — tracks training examples $\quad i \in 1, \ldots, N$

$x_i \rightarrow i^{th}$ training example.

$p \rightarrow$ track features $\quad p \in 1, \ldots, D$

$x_{ip} \rightarrow p^{th}$ feature of $x_i$

~~$l$ — tra~~

$j$ — track hidden units $\quad j \in 1, \ldots, M$

$net_j$ — Summation for $j^{th}$ hidden unit.

$$net_j = W_j^T \, x$$

$$z_j = \sigma(net_j) \rightarrow \text{output of the } j^{th} \text{ hidden unit.}$$

$W_j$ — weight vector for $j^{th}$ hidden unit.

$W_j \rightarrow$ length $D$ [assume $D$ includes the bias]

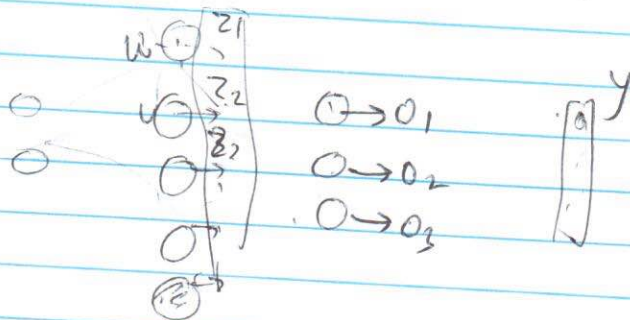$W_{jp} \rightarrow p^{th}$ entry of the weight vector of $j^{th}$ hidden unit.

$l$ — track output units $\quad l \in 1, \ldots, k$

~~$net$~~ $W_l$ — weight vector for $l^{th}$ output unit

$\hookrightarrow$ length $(M+1)$

$W_{lj} \rightarrow j^{th}$ entry of the weight vector of the $l^{th}$ output unit

$$net_l = W_l^T \, z \qquad z \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_M \end{bmatrix}$$

$$O_l = \sigma(net_l) \rightarrow \text{output at } l^{th} \text{ output unit.}$$

$r^{th}$ unit.



$$\frac{\partial J}{\partial w_{r_1}} = \frac{\partial J}{\partial net_r} \frac{\partial net_r}{\partial w_{r_1}}$$

$$net_r = u_{r_1} w_{r_1} + u_{r_2} w_{r_2} + \cdots$$

$$\frac{\partial net_r}{\partial w_{r_1}} = u_{r_1}$$

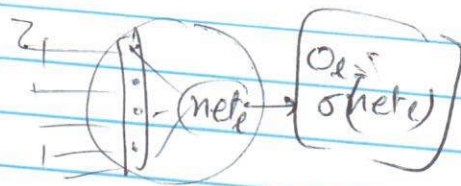$$\frac{\partial J}{\partial w_{rq}} = \frac{\partial J}{\partial net_r} u_{rq}$$

Output unit.

$$\frac{\partial J}{\partial w_{jp}} = \frac{\partial J}{\partial net_j} x_p$$

$$\frac{\partial J}{\partial w_{\ell j}} = \frac{\partial J}{\partial net_\ell} z_j$$

At the output node $(\ell)$

$$\frac{\partial J}{\partial net_\ell} = \frac{\partial J}{\partial O_\ell} \frac{\partial O_\ell}{\partial net_\ell}$$



$$O_\ell = \sigma(net_\ell)$$

What is $\dfrac{\partial O_\ell}{\partial net_\ell} = \dfrac{\partial}{\partial net_\ell} \sigma(net_\ell)$

$\sigma(a)$
$$\frac{d\sigma(a)}{da} = \sigma(a)(1 - \sigma(a))$$

$$\frac{\partial O_\ell}{\partial net_\ell} = O_\ell(1 - O_\ell)$$

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

$$\frac{d}{da}\sigma(a) = \frac{-\exp(-a)}{(1 + \exp(-a))^2}$$

$$= \left(\frac{1}{1 + \exp(-a)}\right)\left(1 - \frac{1}{1 + \exp(-a)}\right)$$

$$\frac{\partial J}{\partial net_\ell} = \frac{\partial J}{\partial O_\ell} \left( \frac{\partial O_\ell}{\partial net_\ell} \right) = O_\ell(1-O_\ell) \frac{\partial J}{\partial O_\ell}$$

What is $\frac{\partial J}{\partial O_\ell}$ ?

$$J = \frac{1}{2}(y_\ell - O_\ell)^2$$

$$\frac{\partial J}{\partial O_\ell} = -(y_\ell - O_\ell)$$

$$\boxed{\frac{\partial J}{\partial net_\ell} = -O_\ell(1-O_\ell)(y_\ell - O_\ell)}$$

$$\boxed{\frac{\partial J}{\partial w_{\ell j}} = -O_\ell(1-O_\ell)(y_\ell - O_\ell)z_j}$$

)