| AND | | | OR | | | XOR | | |
|---|---|---|---|---|---|---|---|---|
| x1 | x2 | y | x1 | x2 | y | x1 | x2 | y |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 |
| 1 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 1 |
| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 |

# McCulloch-Pitts model of a neuron



$$\pi_i = \frac{1}{1+e^{-x_i \cdot \Theta}} \qquad i=1:n$$

$$\underline{\Theta} = [\Theta_0, \Theta_1, \ldots]^T$$

a) $x_1$ AND $x_2$    b) $x_1$ OR $x_2$    c) $x_1$ XOR $x_2$

# MLP – 1 neuron

$$\hat{y}_i = P(y_i = 1 \mid x_i, \theta)$$

We are given the data $\{x_i, y_i\}_{i=1}^{n}$

eg.

| | $X_{i1}$ | $X_{i2}$ | $Y_i$ |
|---|---|---|---|
| $i=1$ | 0.2 | 6 | 0 |
| $i=2$ | 0.3 | 22 | 1 |
| $i=3$ | 0.6 | -0.6 | 1 |
| $i=4$ | -0.4 | 58 | 0 |
| $\vdots$ | | | |

separating hyper-plane

$x \hookleftarrow Y=1$

$Y=0$

# MLP – 1 neuron



$$\hat{y}_i = P(y_i = 1 \mid x_i, \theta)$$

$$0 < \hat{y}_i < 1$$

$$U = \theta_1 + \theta_2 x_{i1} + \theta_3 x_{i2}$$

$$\hat{y}_i = \frac{1}{1+e^{-u}} = \frac{1}{1+e^{-\theta_1 - \theta_2 x_{i1} - \theta_3 x_{i2}}} = P(y_i = 1 \mid \underline{x}_i, \underline{\theta})$$

$$P(y_i \mid \underline{x}_i, \underline{\theta}) = \hat{y}_i^{y_i}(1 - \hat{y}_i)^{1-y_i} = \begin{cases} \hat{y}_i & \text{when } y_i = 1 \\ 1 - \hat{y}_i & \text{otherwise} \end{cases}$$
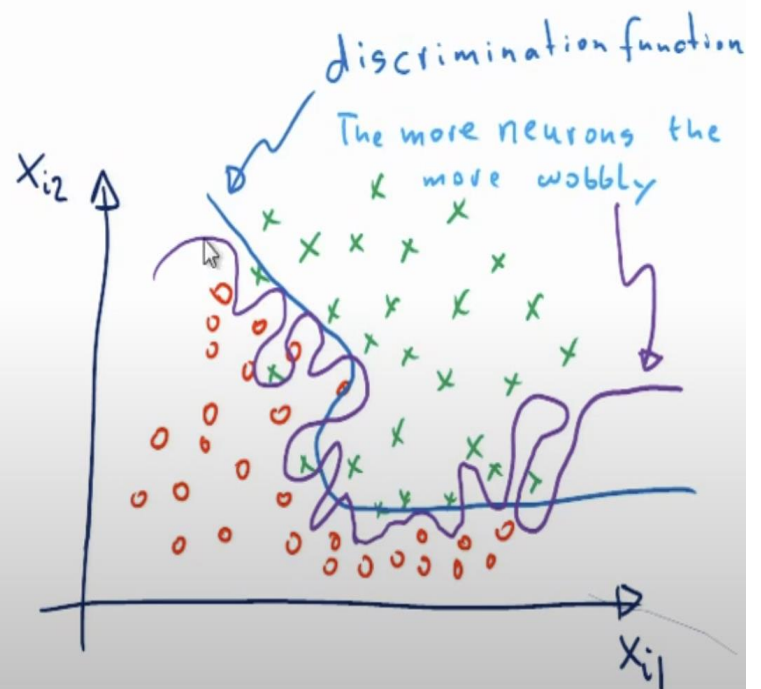
For $n$ independent observations (Bernoulli)

$$P(Y \mid \underline{\underline{x}}, \underline{\theta}) = \prod_{i=1}^{n} P(y_i \mid \underline{x}_i, \underline{\theta})$$

# MLP – 3 neurons, 2 layers

1 $\theta_1$

$\theta_2$

$x_{i1}$

$\theta_3$

$\theta_9$

$x_{i2}$

$\theta_8$

1 $\theta_4$

$\Sigma$  $u_{11}$  $o_{11}$

$\Sigma$  $u_{12}$  $o_{12}$

1 $\theta_5$

$\theta_6$

$\theta_7$

$\Sigma$  $u_{21}$

$\hat{y}_i = P(y_i = 1 | x_i, \theta)$

discrimination function

The more neurons the more wobbly

$x_{i2}$

$x_{i1}$

Data:

| | $x_{i1}$ | $x_{i2}$ | $y_i$ |
|---|---|---|---|
| i=1 | 6 | 8 | 1 |
| i=2 | 0.2 | -5 | 0 |
| | -100 | 3.1 | 1 |
| | 6 | 9 | 0 |
| | 5 | 8 | 0 |

MLP – 3 neurons, 2 layers

$$\hat{y}_i = P(y_i = 1 | x_i, \theta)$$
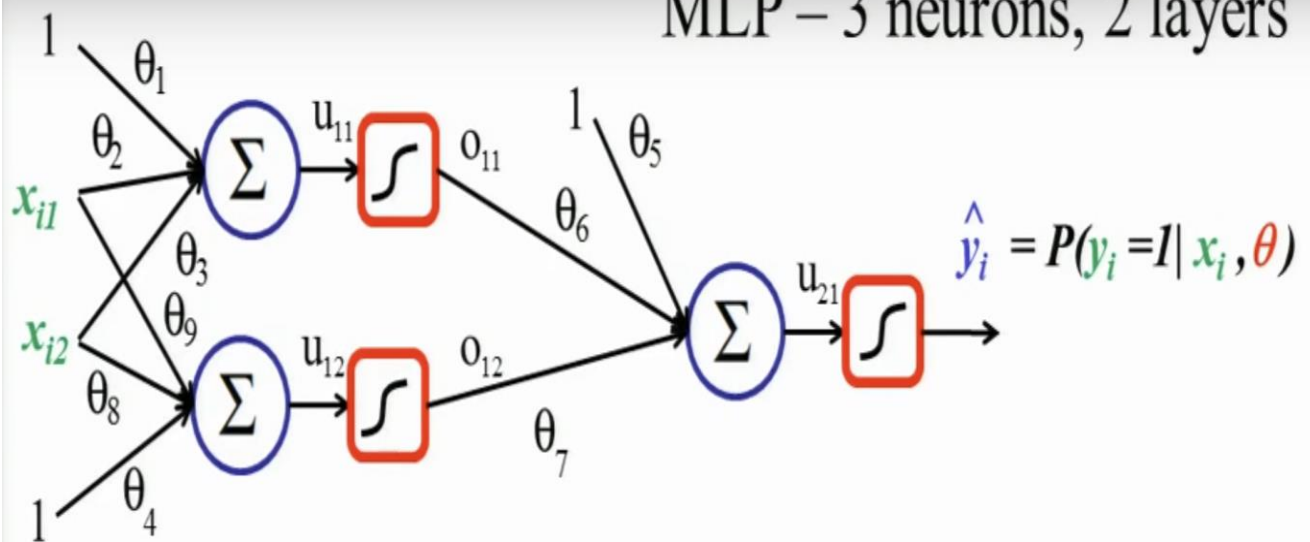
$$u_{11} = \theta_1 + \theta_2 x_{i1} + \theta_3 x_{i2}$$

$$O_{11} = \frac{1}{1 + e^{-u_{11}}}$$

$$\hat{y}_i = \frac{1}{1 + e^{-u_{21}}}$$

$$u_{12} = \theta_4 + \theta_9 x_{i1} + \theta_8 x_{i2}$$

$$O_{12} = \frac{1}{1 + e^{-u_{12}}}$$

$$u_{21} = \theta_5 + \theta_6 O_{11} + \theta_7 O_{12}$$

$$P(y_i | x_i, \theta) = \hat{y}_i^{\,y_i} (1 - \hat{y}_i)^{1 - y_i}$$

MLP – 3 neurons, 2 layers

$\hat{y}_i = P(y_i = 1 | x_i, \theta)$

For $n$ independent observations.

$$P(Y | \underline{x}, \underline{\theta}) = \prod_{i=1}^{n} \hat{y}_i^{Y_i} (1-\hat{y}_i)^{1-Y_i} = \prod_{i=1}^{n} P(y_i | x_i, \theta)$$
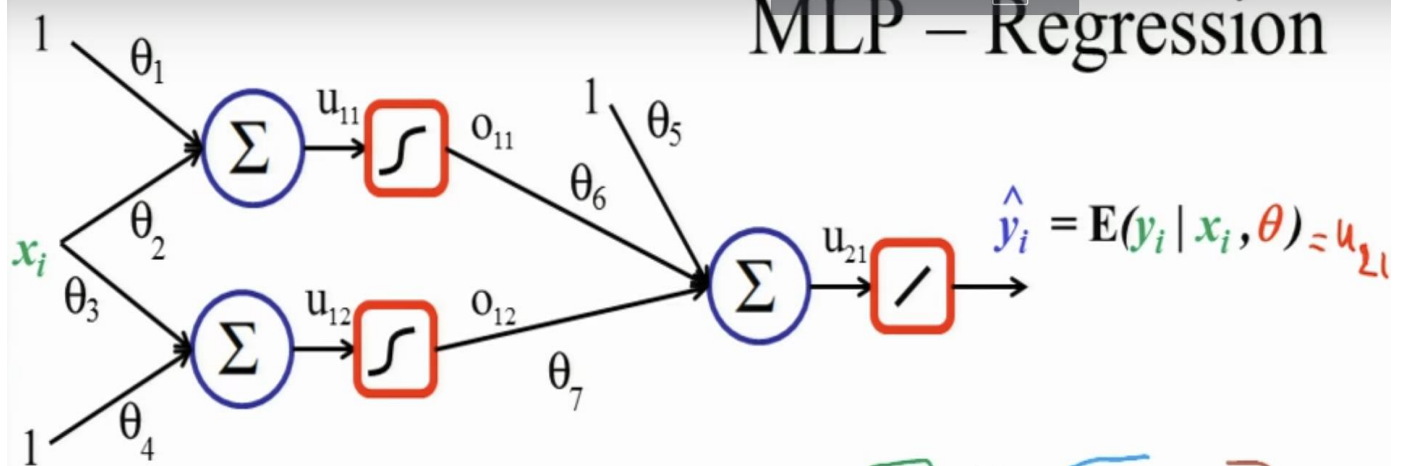
Cost:

$$C(\theta) = -\log P(Y | \underline{x}, \theta) = -\sum_{i=1}^{n} Y_i \log \hat{y}_i + (1-Y_i)\log(1-\hat{y}_i)$$
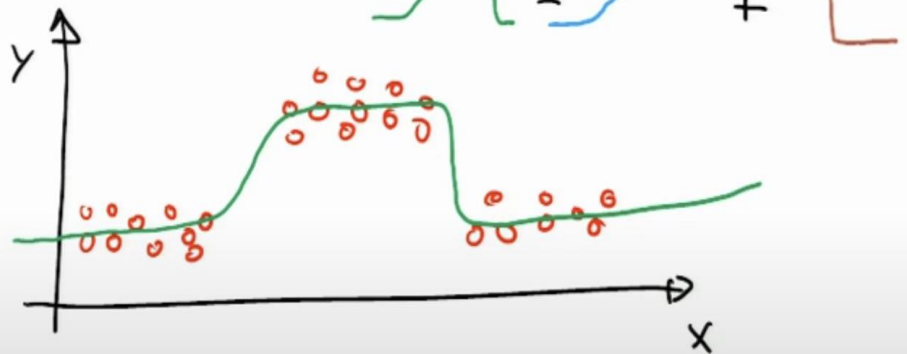
i.e minimize the cross-entropy error.

Cross-entropy measures uncertainty. By minimizing cross-entropy, we maximize the information gained about the data as the model learns

# MLP – Regression



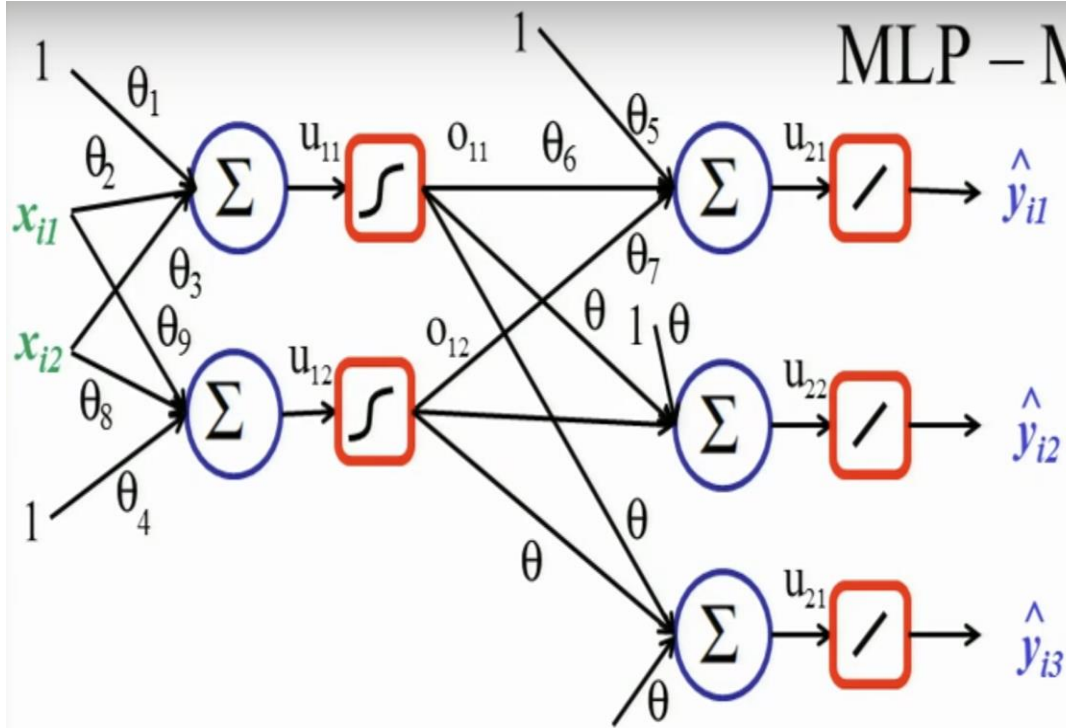$$\hat{y}_i = E(y_i \mid x_i, \theta) = u_{21}$$

Data:

| $x_i$ | $y_i$ |
|-------|-------|
| 0.2 | 0.6 |
| 0.9 | 0.4 |
| -0.6 | 5 |
| -0.3 | 6.2 |

$$\hat{y}_i = \theta_5 + \frac{\theta_6}{1+e^{-\theta_1 - \theta_2 x_i}} + \frac{\theta_7}{1+e^{-\theta_4 - \theta_3 x_i}}$$
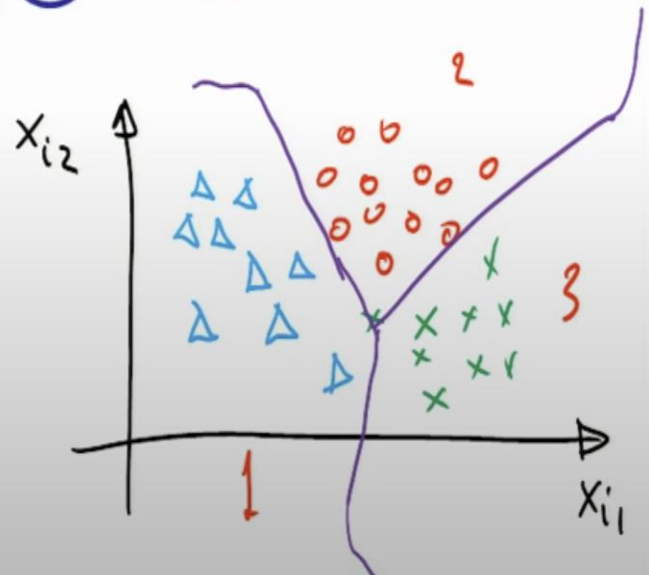
Theta five shifts the curve up and down. Theta six controls the height of the S-shaped curve; if theta six is large, the curve is tall, and vice versa. If theta six is negative, it flips the curve. Theta one shifts the curve left and right, while theta two controls its width, making it either wider or thinner.
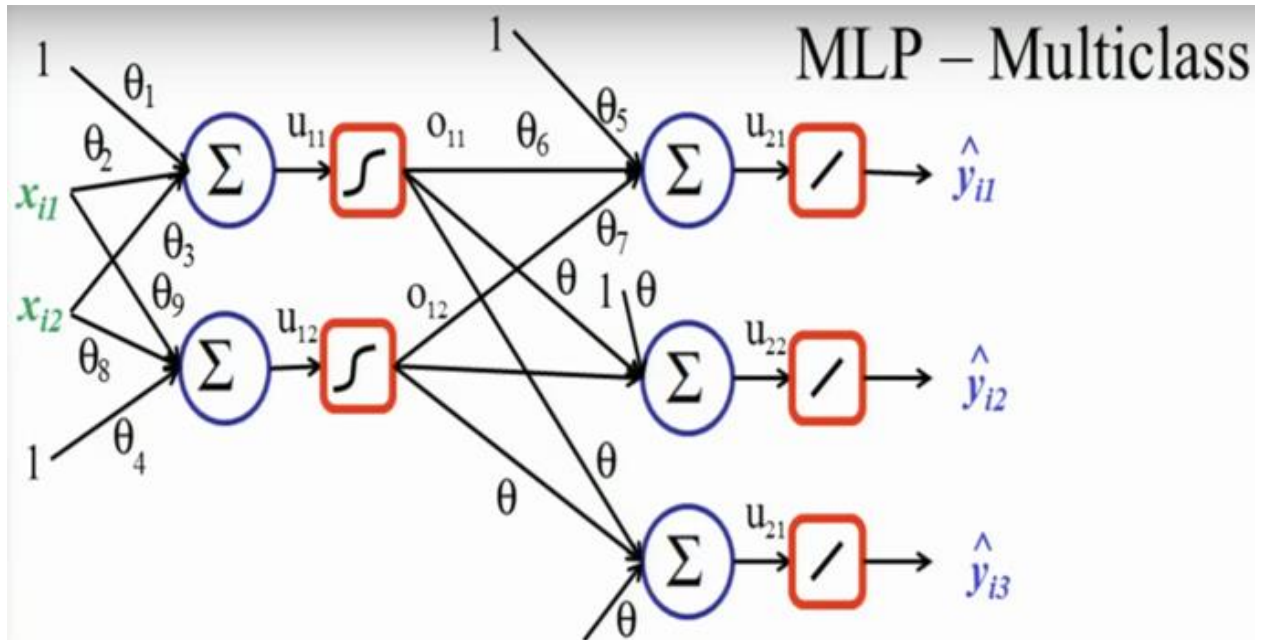
MLP – Multiclass

$1$  $\theta_1$

$\theta_2$

$x_{i1}$

$\theta_3$
$\theta_9$

$x_{i2}$

$\theta_8$

$1$  $\theta_4$

$\Sigma$  $u_{11}$  $\int$  $o_{11}$  $\theta_6$

$\Sigma$  $u_{12}$  $\int$  $o_{12}$

$1$

$\theta_5$

$\Sigma$  $u_{21}$  $/$  $\rightarrow$  $\hat{y}_{i1}$

$\theta_7$

$\theta$  $1$  $\theta$

$\Sigma$  $u_{22}$  $/$  $\rightarrow$  $\hat{y}_{i2}$

$\theta$

$\theta$

$\Sigma$  $u_{21}$  $/$  $\rightarrow$  $\hat{y}_{i3}$

$\theta$

Data :

| $x_{i1}$ | $x_{i2}$ | $y_{i1}$ | $y_{i2}$ | $y_{i3}$ | |
|---|---|---|---|---|---|
| 0.2 | 0.3 | 0 | 1 | 0 | Class 2 |
| -5 | -6 | 1 | 0 | 0 | Class 1 |
| -20 | 4 | 1 | 0 | 0 | " " |
| 42 | 6.8 | 0 | 0 | 1 | Class 3 |

$x_{i2}$

2

3

1

$x_{i1}$

# MLP – Multiclass

To get a probabilistic model, define:   SOFTMAX

$$P(y_i = (010)|x_i, \theta) = P(y_i = 2 | x_i, \theta) = \frac{e^{\hat{y}_2}}{e^{\hat{y}_1} + e^{\hat{y}_2} + e^{\hat{y}_3}}$$

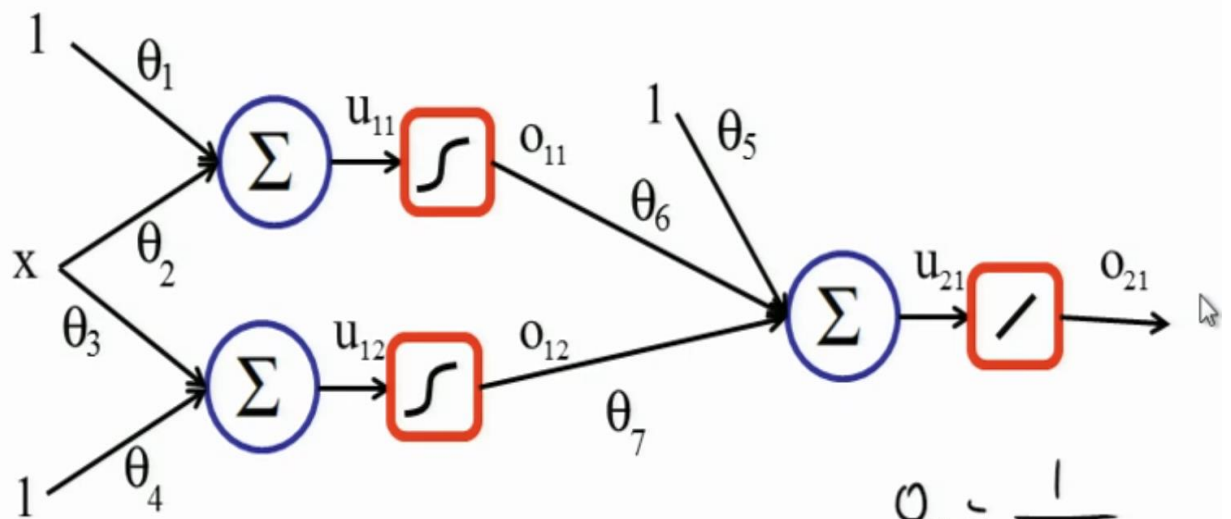$$\mathbb{I}_2(Y_i) = \begin{cases} 1 & Y_i = 2 \\ 0 & 0.\omega. \end{cases}$$

Then,

$$P(Y_i | x_i, \theta) = \overbrace{\left[ \frac{e^{\hat{Y}_{i1}}}{\underbrace{e^{\hat{Y}_{i1}} + e^{\hat{Y}_{i2}} + e^{\hat{Y}_{i3}}}_{sum}} \right]}^{\mathbb{I}_1(Y_i)} \overbrace{\left[ \frac{e^{\hat{Y}_{i2}}}{e^{\hat{Y}_{i1}} + e^{\hat{Y}_{i2}} + e^{\hat{Y}_{i3}}} \right]}^{\mathbb{I}_2(Y_i)} \overbrace{\left[ \frac{e^{\hat{Y}_{i3}}}{e^{\hat{Y}_{i1}} + e^{\hat{Y}_{i2}} + e^{\hat{Y}_{i3}}} \right]}^{\mathbb{I}_3(Y_i)}$$

$$= \begin{cases} e^{\hat{Y}_{i1}}/sum & Y_i = 1 \\\\ e^{\hat{Y}_{i2}}/sum & Y_i = 2 \\\\ e^{\hat{Y}_{i3}}/sum & Y_i = 3 \end{cases}$$

Cost:

$$C(\theta) = -\log P(Y | X, \theta) = -\sum_{i=1}^{n} \sum_{j=1}^{3} \mathbb{I}_j(Y_i) \log \frac{e^{\hat{Y}_{ij}}}{sum}$$
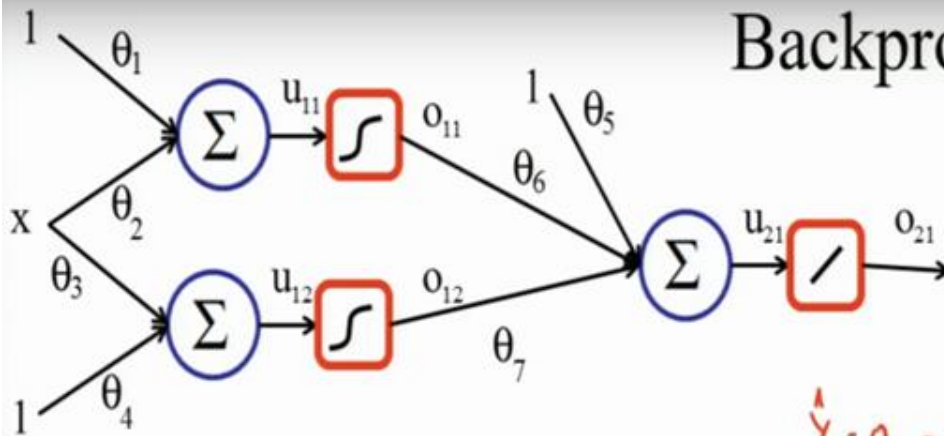
# Backpropagation



$$\hat{y} = O_{21} = u_{21} = \Theta_5 + \Theta_6 O_{11} + \Theta_7 O_{12}$$

$$u_{11} = \Theta_1 + \Theta_2 x$$

$$u_{12} = \Theta_4 + \Theta_3 x$$

$$O_{11} = \frac{1}{1 + e^{-u_{11}}}$$

$$O_{12} = \frac{1}{1 + e^{-u_{12}}}$$

# Backpropagation

$$\hat{Y} = O_{21} = \theta_5 + \theta_6 O_{11} + \theta_7 O_{12}$$
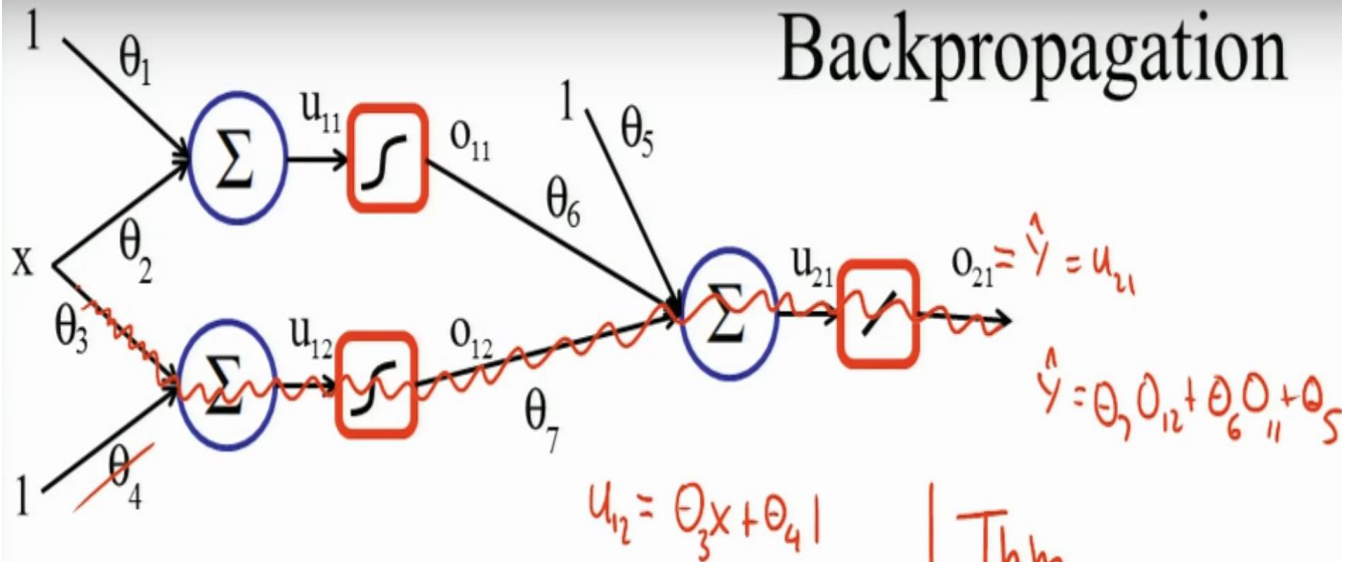
$$E(\theta) = \left( Y_i - \hat{Y}_i(x_i, \theta) \right)^2$$

$$\frac{\partial E(\theta)}{\partial \theta_j} = -2 \left( Y_i - \hat{Y}_i(x_i, \theta) \right) \frac{\partial \hat{Y}_i(x_i, \theta)}{\partial \theta_j}$$

$$\frac{\partial \hat{Y}_i}{\partial \theta_5} = 1 \qquad \frac{\partial \hat{Y}_i}{\partial \theta_6} = O_{11} \qquad \frac{\partial \hat{Y}_i}{\partial \theta_7} = O_{12}$$

# Backpropagation



$$\hat{y} = u_{21}$$

$$\hat{y} = \theta_7 O_{12} + \theta_6 O_{11} + \theta_5$$

$$u_{12} = \theta_3 x + \theta_4 1$$

$$\frac{\partial \hat{y}}{\partial \theta_3} = \frac{\partial \hat{y}}{\partial O_{12}} \frac{\partial O_{12}}{\partial u_{12}} \frac{\partial u_{12}}{\partial \theta_3}$$
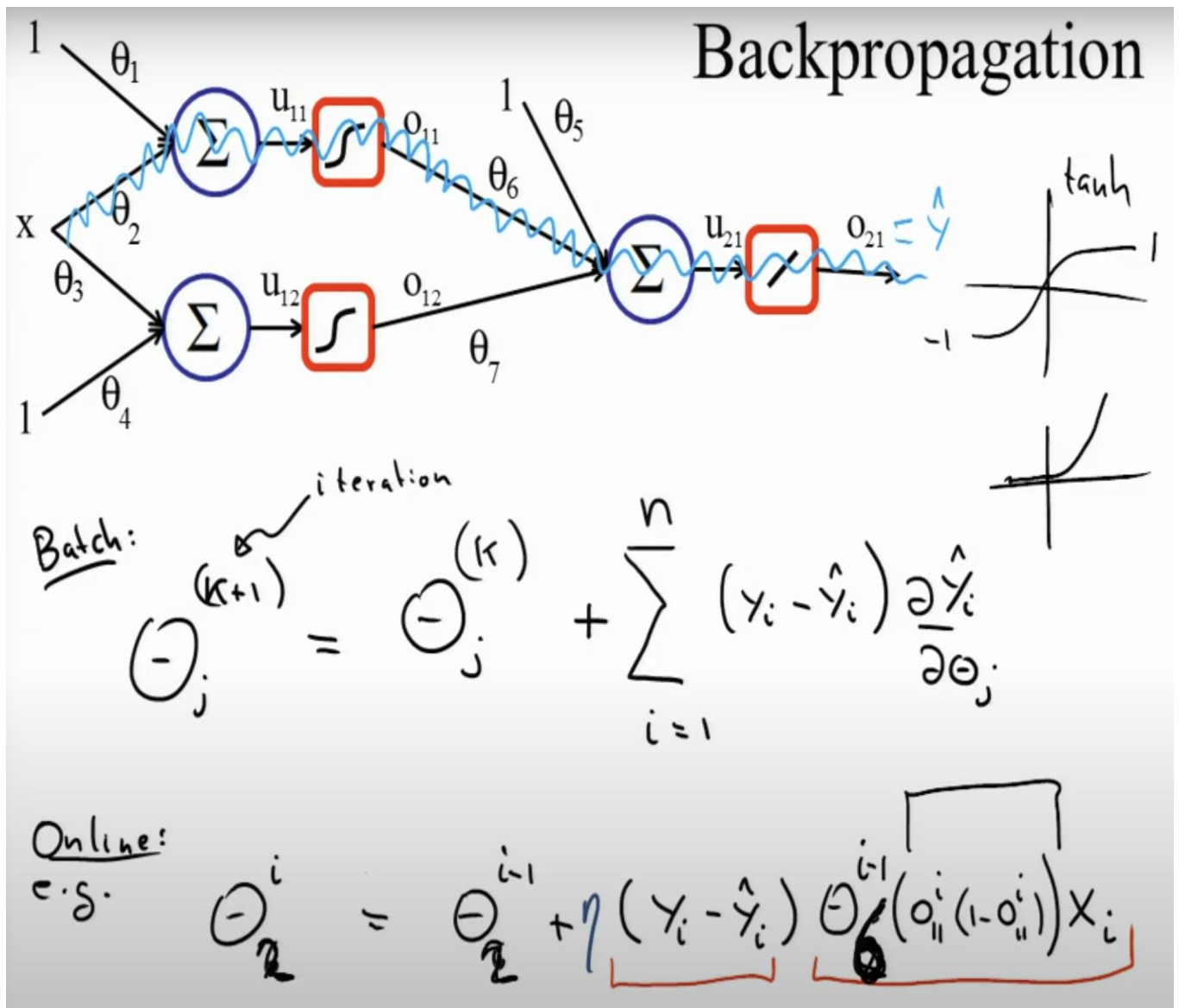
$$= \theta_7 O_{12} \left[ 1 - O_{12} \right] x$$

## Thm

$$\frac{\partial O_{12}}{\partial u_{12}} = O_{12} (1 - O_{12})$$

i.e.

$$\frac{\partial}{\partial x} \frac{1}{1 + e^{-x}} = \left( \frac{1}{1 + e^{-x}} \right) \left( 1 - \frac{1}{1 + e^{-x}} \right)$$

# Backpropagation

$\tanh$

Batch:

$$\Theta_j^{(k+1)} = \Theta_j^{(k)} + \sum_{i=1}^{n} (y_i - \hat{y}_i) \frac{\partial \hat{y}_i}{\partial \Theta_j}$$

iteration

Online:

e.g.

$$\Theta_2^i = \Theta_2^{i-1} + \eta (y_i - \hat{y}_i) \Theta_6^{i-1} (o_{11}^i (1 - o_{11}^i)) x_i$$

Hyperbolic tangent function or rectified unit to deal with vanishing gradients as the network gets deeper.