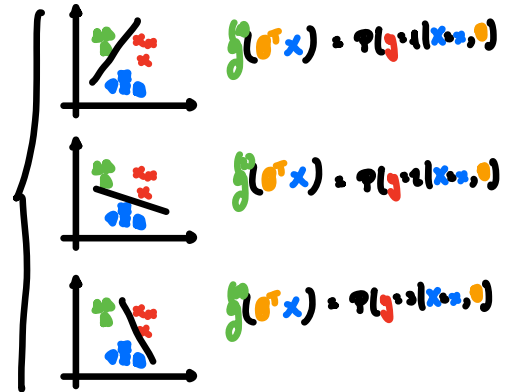
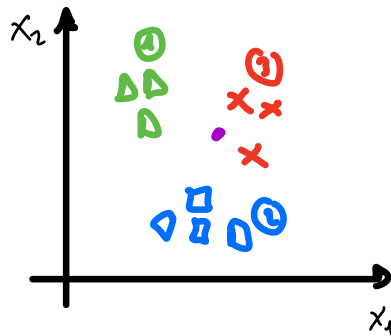


⑥ Supervised learning II

* Multi-classification

* One-vs-all:

$y \in \{1, 2, 3\}$



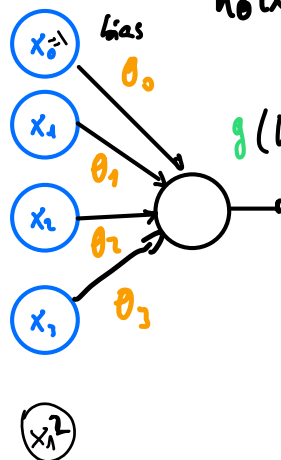
* How is the cost?

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[\sum_{k=1}^K (1 - y_k^i) (-\log(1 - g(h_\theta(x^i)))) + y_k^i (-\log(g(h_\theta(x^i)))) \right] + \lambda \|\theta\|^2$$

$$y_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}; y_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}; \dots$$

• Neural model: Logistic unit

Linear case



$$h_\theta(x) = X\theta = x_1\theta_1 + x_2\theta_2 + x_0\theta_0$$

$$x = [x_0, x_1, \dots, x_n]$$

$$\theta = [\theta_0, \theta_1, \dots, \theta_n]$$

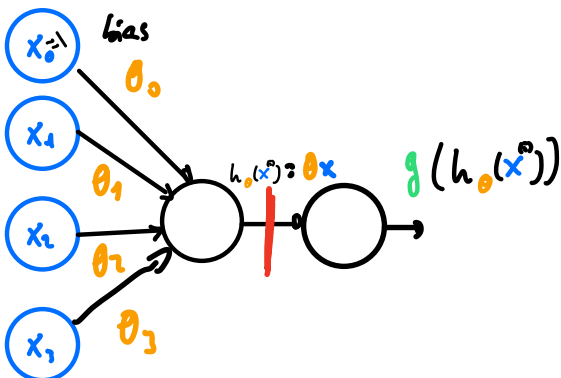
$g \equiv$ sigmoid activation function (logistic)

1

↓

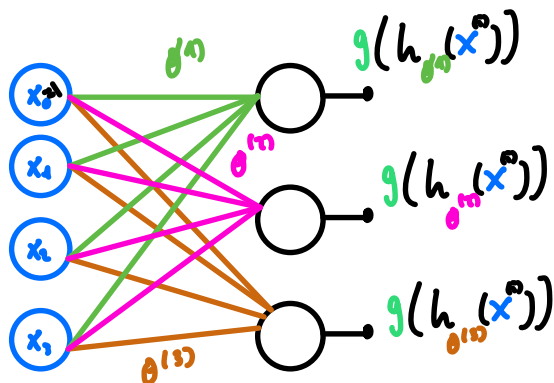
θ_0

or more
schematically



In one or all :

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}; \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$



$$\Rightarrow \begin{bmatrix} g(h_{\theta^1}(x)) \\ g(h_{\theta^2}(x)) \\ g(h_{\theta^3}(x)) \end{bmatrix} \sim \begin{bmatrix} y \end{bmatrix}$$

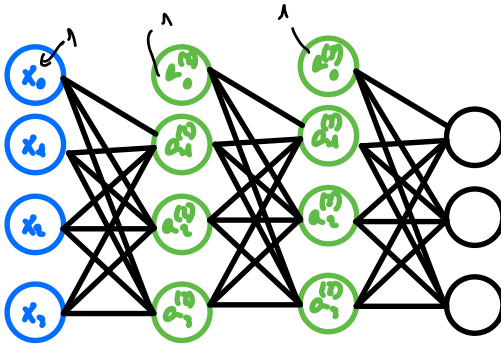
$$\theta = [\theta^1, \theta^2, \theta^3] \Rightarrow \begin{bmatrix} \boxed{} \\ \boxed{} \\ \boxed{} \end{bmatrix}$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[\sum_{k=1}^K (1 - y_k^i) (-\log(1 - g(h_{\theta_k}(x^i))) + y_k^i (-\log(g(h_{\theta_k}(x^i)))) \right] + \lambda \|\theta\|^2$$

$$= J_1(\theta_1) + J_2(\theta_2) + J_3(\theta_3) + \lambda \|\theta\|^2$$

$$\left\{ \begin{array}{l} \nabla J_{\theta_1} = \frac{1}{m} \sum_{i=1}^m [g(h_{\theta_1}(x^i)) - y_1^i] x^{(i)} \\ \nabla J_{\theta_2} = \frac{1}{m} \sum_{i=1}^m [g(h_{\theta_2}(x^i)) - y_2^i] x^{(i)} \\ \nabla J_{\theta_3} = \frac{1}{m} \sum_{i=1}^m [g(h_{\theta_3}(x^i)) - y_3^i] x^{(i)} \end{array} \right.$$

* Neural network (classification)



$$\{(x^1, y^1), (x^2, y^2), (x^3, y^3), (x^4, y^4)\}$$

L = total number of layers

s_L = number of units in layer L

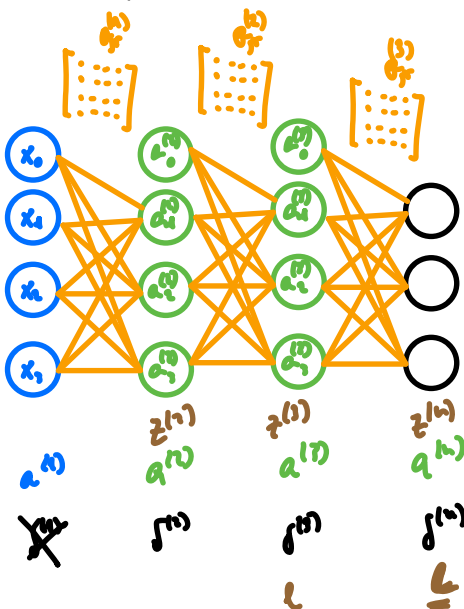
* Cost function

$$h_\theta(x)$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[\sum_{k=1}^K (1 - y_k^i) (-\log(1 - g(h_{\theta}(x^i) - y_k^i))) + y_k^i (-\log(g(h_{\theta}(x^i) - y_k^i))) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{r=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\theta_{jr}^{(l)})^2$$

but not 0

* Forward propagation



$$\left\{ \begin{array}{l} a^{(0)} = x \\ z^{(1)} = \theta^{(0)} a^{(0)} \\ a^{(1)} = g(z^{(1)}) \\ z^{(2)} = \theta^{(1)} a^{(1)} \\ a^{(2)} = g(z^{(2)}) \\ z^{(3)} = \theta^{(2)} a^{(2)} \\ a^{(3)} = g(z^{(3)}) \end{array} \right.$$

Remember:

$$\left\{ \begin{array}{l} g(z) = \frac{1}{1 + e^{-z}} \\ g'(z) = g(z)(1 - g(z)) \end{array} \right.$$

* Gradient.

Before: one layer architecture

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta_j^k} &= \frac{1}{m} \sum_{i=1}^m (1 - y_k^i) \frac{-1}{(1 - g(h_{\bullet}(x^i)))} \overbrace{-(g(h_{\bullet}(x^i))(1 - g(h_{\bullet}(x^i)))}^{g'(h_{\bullet}(x^i))} x^{(i)} \frac{\partial z^k}{\partial \theta_j^k} \\ &\quad + (y_k^i) \frac{-1}{(g(h_{\bullet}(x^i)))} \overbrace{-(g(h_{\bullet}(x^i))(1 - g(h_{\bullet}(x^i)))}^{g'(h_{\bullet}(x^i))} x^{(i)} \frac{\partial z^k}{\partial \theta_j^k} \\ &= -\frac{1}{m} \sum_{i=1}^m [g(h_{\bullet}(x^i)) - y_k^i] x^{(i)} \frac{\partial z^k}{\partial \theta_j^k} \end{aligned} \quad h = \theta x$$

Now: several layer architecture

* For the last layer: $(L-1)$

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta^{L-1}} &= \frac{1}{m} \sum_{i=1}^m (1 - y^i) \frac{-1}{(1 - g(z^L))} \overbrace{-(g(z^L)(1 - g(z^L)))}^{g'(z^L)} \frac{\partial z^L}{\partial \theta^{L-1}} \\ &\quad + (y^i) \frac{-1}{(g(z^L))} \overbrace{-(g(z^L)(1 - g(z^L)))}^{g'(z^L)} \frac{\partial z^L}{\partial \theta^{L-1}} \\ &= \frac{1}{m} \sum_{i=1}^m [g(z^L) - y^i] \frac{\partial z^L}{\partial \theta^{L-1}} = \frac{1}{m} \sum \delta^L \frac{\partial z^L}{\partial \theta^{L-1}} \end{aligned}$$

$$\frac{\partial J(\theta)}{\partial \theta^{L-1}} = \frac{1}{m} \sum_{i=1}^m \delta^L \frac{\partial z^L}{\partial \theta^{L-1}} = \frac{1}{m} \sum_{i=1}^m \delta^L a^{(L-1)}(i) \quad \text{with } z^{(L)} = \theta^{(L-1)} a^{(L-1)}$$

$$\boxed{\frac{\partial J(\theta)}{\partial \theta^{L-1}} = \frac{1}{m} \sum_{i=1}^m \delta^L a^{(L-1)}(i)}$$

$$\text{with } \boxed{\delta^L = g(z^L) - y^i}$$

* For one before the last layer: $(L-2)$

$$\frac{\partial J(\theta)}{\partial \theta^{(L-2)}} = \frac{1}{n} \sum_{i=1}^n \delta^{(L)} \frac{\partial z^{(L)}}{\partial \theta^{(L-2)}} = \frac{1}{n} \sum_{i=1}^n \delta^{(L)} \frac{\partial a^{(L)}}{\partial \theta^{(L-2)}}$$

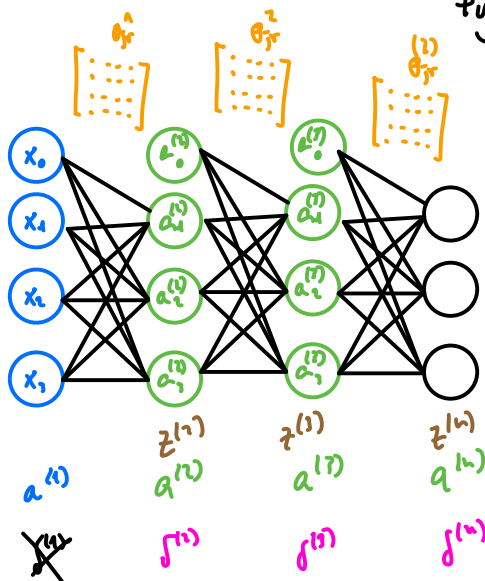
$$\begin{aligned} a^{(L-1)} &= g(z^{(L-1)}) \\ &\downarrow \\ &= \frac{1}{n} \sum_{i=1}^n \delta^{(L)} \frac{\partial a^{(L)}}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial \theta^{(L-2)}} = \frac{1}{n} \sum_{i=1}^n \delta^{(L)} \frac{\partial a^{(L)}}{\partial z^{(L-1)}} g'(z^{(L-1)}) a^{(L-1)} \end{aligned}$$

$$\boxed{\frac{\partial J(\theta)}{\partial \theta^{(L-2)}} = \frac{1}{n} \sum_{i=1}^n \delta^{(L-1)} a^{(L-1)}}$$

$$\text{with } \delta^{(L-1)} = \delta^{(L)} \frac{\partial a^{(L)}}{\partial z^{(L-1)}} g'(z^{(L-1)})$$

* For one before the last layer: $(L-3)$ works similarly

Backpropagation



Forward propagation

$$\begin{cases} a^{(1)} = x \\ z^{(2)} = \theta^{(1)} a^{(1)} \\ a^{(2)} = g(z^{(2)}) \\ z^{(3)} = \theta^{(2)} a^{(2)} \\ a^{(3)} = g(z^{(3)}) \\ z^{(4)} = \theta^{(3)} a^{(3)} \\ a^{(4)} = g(z^{(4)}) \end{cases}$$

Backward propagation

start on layer 4

$$\begin{aligned} \delta^{(4)} &= a^{(4)} - y \\ \delta^{(3)} &= \theta^{(3)} \delta^{(4)} * g'(z^{(3)}) \\ \delta^{(2)} &= \theta^{(2)} \delta^{(3)} * g'(z^{(2)}) \\ \delta^{(1)} &= \theta^{(1)} \delta^{(2)} * g'(z^{(1)}) \end{aligned}$$

No $\delta^{(0)}$

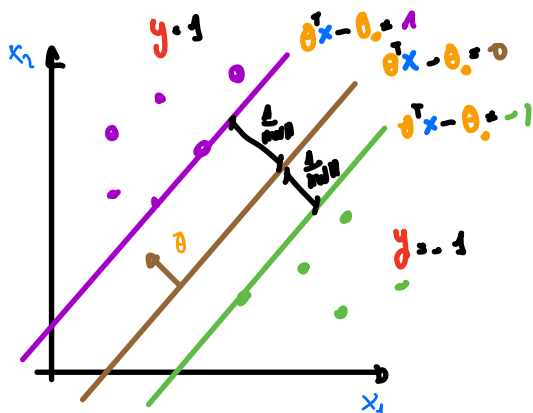
"We propagate the error"

$$\frac{\partial J(\theta)}{\partial \theta_j^{(l)}} = \frac{1}{n} \sum_{i=1}^n a_j^{(l)} \delta_i^{(l+1)}$$

* Remarks: $\left\{ \begin{array}{l} * \text{ Now people use automatic differentiation} \\ * \text{ Also use regularization} \end{array} \right.$

Support vector machines

SVH (Supp. vect. machine)



Linearly separable:

$$\begin{cases} \theta^T x - \theta_0 \geq 1 & \text{when } y_i = 1 \\ \theta^T x - \theta_0 \leq -1 & \text{when } y_i = -1 \end{cases}$$

$$y_i (\theta^T x_i - \theta_0) \geq 1$$

$$\max_{\theta, \theta_0} \frac{2}{\|\theta\|}$$

$$\text{s.t. } y_i (\theta^T x_i - \theta_0) \geq 1$$

Non-linearly separable

Primal

$$\begin{cases} \min_{\theta, \theta_0} \|\theta\| \\ \text{s.t. } y_i (\theta^T x_i - \theta_0) - 1 \geq 0 \quad \forall i \end{cases}$$

$$f_i(\theta) \leq 0 \rightarrow 1 - y_i (\theta^T x_i - \theta_0) \leq 0 \quad \forall i$$

Primal (relaxed)

$$\begin{cases} \min_{z_i, \theta, \theta_0} \frac{1}{n} \sum_{i=1}^n z_i + \mu \|\theta\|^2 \\ \text{s.t. } y_i (\omega^T x_i - b) \geq 1 - z_i \quad (\lambda_i) \end{cases}$$

Dual

$$\begin{cases} \max_c \quad \sum \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i y_i (x_i^T x_j) y_j \lambda_j \\ \text{s.t.} \quad \sum \lambda_i y_i = 0 \quad \text{with } 0 \leq \lambda_i \leq \frac{1}{2n\mu} \end{cases}$$