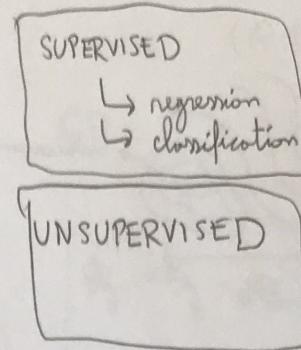


1

ANDREW NG - MACHINE LEARNING

WEEK 1

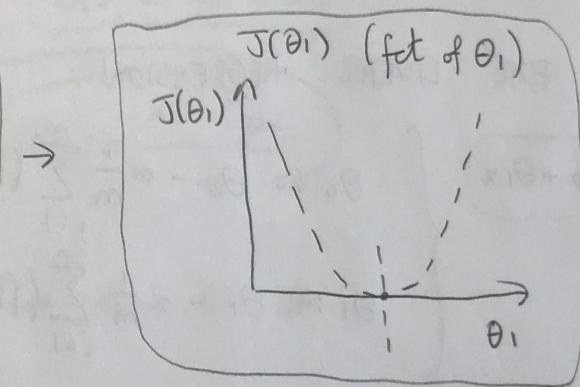
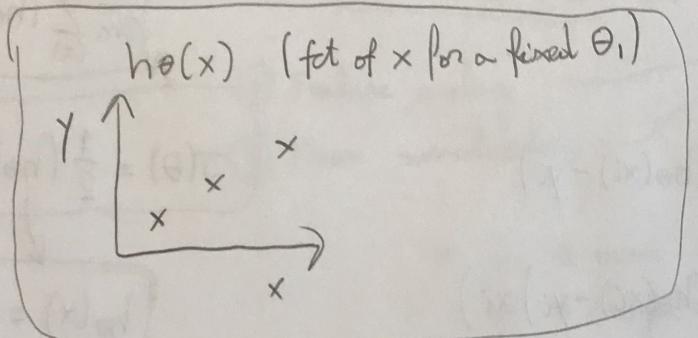
- ML :
- learn from experience E
 - task T
 - performance measure P



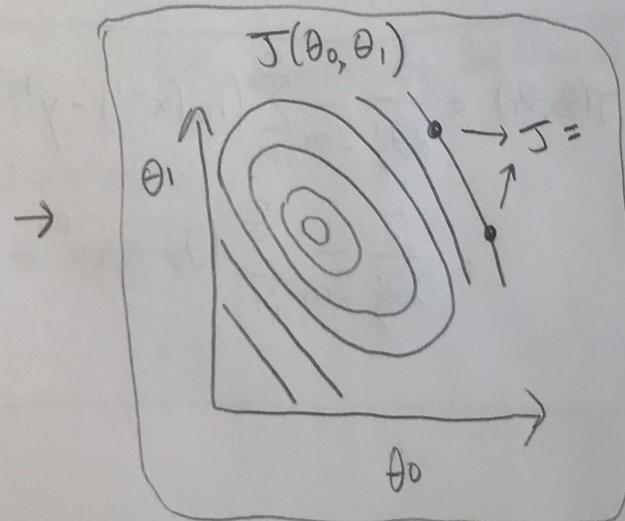
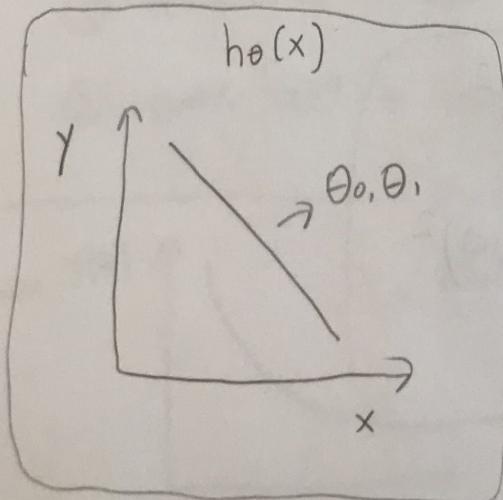
COST FUNCTION :

$$\overbrace{J(\theta_0, \theta_1)} = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2$$

to minimize

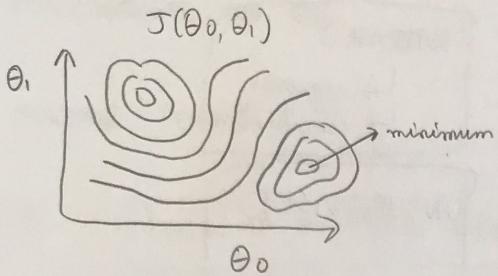


CONTOUR PLOTS



(2)

GRADIENT DESCENT



$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

* simultaneously

- α
- α too small \rightarrow GD can be slow
 - α too large \rightarrow GD may miss the min, fail to converge or even diverge

* No need to decrease α over time
because GD will automatically take smaller steps
as we approach a local min.

FOR LINEAR REGRESSION :

$$\boxed{h_\theta(x) = \theta_0 + \theta_1 x}$$

$$\left\{ \begin{array}{l} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x_i) - y_i) \\ \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m ((h_\theta(x_i) - y_i) x_i) \end{array} \right.$$

$$\frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\boxed{J(\theta) = \frac{1}{2} (h_\theta(x) - y)^2}$$

$$\boxed{h_\theta(x) = \sum_{i=0}^n \theta_i x_i}$$

* BATCH GD \rightarrow looks through all training set

NOTE:

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \\ &= \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2 \end{aligned}$$

MULTIVARIATE LIN. REG.

MULTIPLE FEATURES:

GD:

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$h_{\theta}(x) = [\theta_0 \ \theta_1 \dots \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} = \theta^T x$$

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

$$j := 0 \dots n$$

FEATURE SCALING:

→ Aim to have $-1 \leq x^{(i)} \leq 1$ OR $-0,5 \leq x^{(i)} \leq 0,5$

① Feature scaling ② Mean normalization

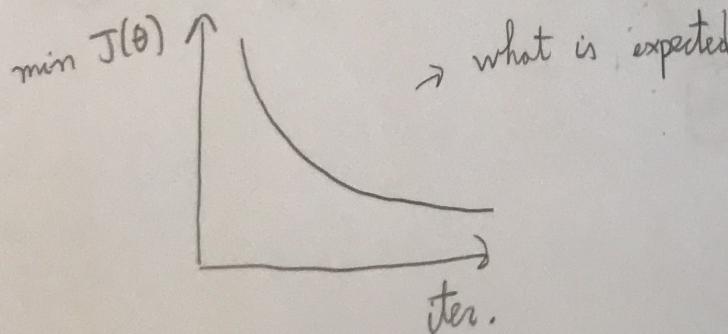
$$x_i := \frac{x_i - \mu_i}{s_i}$$

* $\mu_i \rightarrow \text{mean}$
 $s_i \rightarrow \text{range}(\max - \min) \text{ OR std}$

GD TIPS:

① DEBUGGING → Plot $J(\theta)$ vs nbr of iteration

② CONV. TEST → Converges if $J(\theta)$ decreases by less than E (e.g. 10^{-3})



4

POLYNOMIAL REGRESSION:

→ CHANGE features depending on problem

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$\begin{aligned} x_2 &= x_1^2 \\ x_3 &= x_1^3 \end{aligned}$$

NORMAL EQUATION:

→ Find directly optimal θ without iterations

$$\theta = (X^T X)^{-1} X^T y$$

- * No need to do feature scaling
- If n too large (many features) → better with GD (if $n > 10000$)

(5)

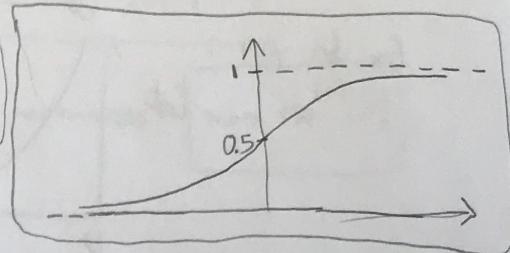
WEEK (3) - LOGISTIC REGRESSION

CLASSIFICATION

→ We want $0 \leq h_{\theta}(x) \leq 1$

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

↳ logistic fct
↳ sigmoid fct



→ $h_{\theta}(x)$ → Estimated probability that $y=1$ on input x

$$h_{\theta}(x) = P(y=1 | x; \theta)$$

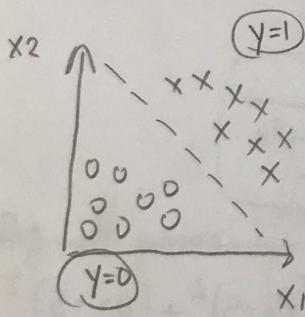
$$P(y=0 | x; \theta) + P(y=1 | x; \theta) = 1$$

$$h_{\theta}(x) = g(\theta^T x) \quad \& \quad g(z) = \frac{1}{1 + e^{-z}}$$

DECISION BOUNDARY

$$\begin{aligned} &\rightarrow y=1 \text{ if } h_{\theta}(x) \geq 0.5 \rightarrow g(z) \geq 0.5 \rightarrow z \geq 0 \rightarrow \theta^T x \geq 0 \\ &\rightarrow y=0 \text{ if } h_{\theta}(x) < 0.5 \rightarrow g(z) < 0.5 \rightarrow z < 0 \rightarrow \theta^T x < 0 \end{aligned}$$

Ex : LINEAR DECISION BOUNDARY



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

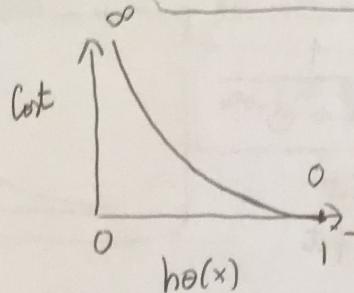
$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

Predict $y=1$ if $-3 + x_1 + x_2 \geq 0$
 $\theta^T x$
 $x_1 + x_2 \geq 3$

6

COST FUNCTION:

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

Ex for $y=1$:Cost = 0 if $y=1$ and $h_{\theta}(x)=1$ Cost $\rightarrow \infty$ if $y=1$ and $h_{\theta}(x) \rightarrow 0$

$$\Rightarrow J(\theta) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)})) \right]$$

fct to minimize to fit θ

new prediction

$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}} \rightarrow P(y=1 | x; \theta)$$

GD:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

$\hookrightarrow \theta_j := \theta_j - \alpha \sum_{j=1}^m \underbrace{(h_{\theta}(x^{(i)}) - y^{(i)})}_{\text{residual}} x_j^{(i)} \rightarrow \text{identical to lin. reg.}$

$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$$

IMPLEMENTATION

- Compute
- ① $J(\theta)$
 - ② $\frac{\partial}{\partial \theta_j} J(\theta)$

EX:

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

$$J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

$$\begin{cases} \frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5) \\ \frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5) \end{cases}$$

 $\rightarrow f_{\min}$

(7)

MULTICLASS CLASSIFICATION

→ ONE VS. ALL

$$h_{\theta}^{(i)}(x) = P(y=i \mid x; \theta) \quad (i=1,2,3,\dots,n)$$

↑ nbr of classes

at prediction → Pick class i that maximizes $\max_i h_{\theta}^{(i)}(x)$

OVERFITTING

- Underfitting = high bias
- Overfitting = high variance

too many features

Fits the training set very well but fails to generalize to new samples

SOLUTIONS

- ① Reduce number of features
- ② Regularization

REGULARIZATION

$$\text{Initial: } \theta_0 + \theta_1 x + \underbrace{\theta_2 x^2}_{+} + \underbrace{\theta_3 x^3}_{+} + \underbrace{\theta_4 x^4}_{+}$$

$$\hookrightarrow \text{REG: } \min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \underbrace{1000 \theta_3^2}_{\downarrow} + \underbrace{1000 \theta_4^2}_{\downarrow}$$

we need θ_3 and θ_4 very close to 0

→ Reg. all params. :

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2$$

$\lambda \rightarrow$ neg. parameter

$\lambda \uparrow \rightarrow$ underfitting

$\lambda \downarrow \rightarrow$ overfitting (no neg)

(8)

REG., LIN. REG.

New GD:

$$\begin{aligned}\theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_j &:= \theta_j - \alpha \left[\left(\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \right) + \frac{\lambda}{m} \theta_j \right] \\ j &\in [1, 2, \dots, n]\end{aligned}$$

New Normal Eqn:

$$\theta = (X^T X + \lambda L)^{-1} X^T Y$$

$$L = \begin{bmatrix} 0 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 0 \end{bmatrix}_{n+1}$$

$$\theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\hookrightarrow \left(1 - \frac{\alpha \lambda}{m}\right) \text{ always } < 1$$

\hookrightarrow reducing value of θ_j on every update

REG., LOG. REG.

New GD:

\rightarrow Same GD as LIN.REG.

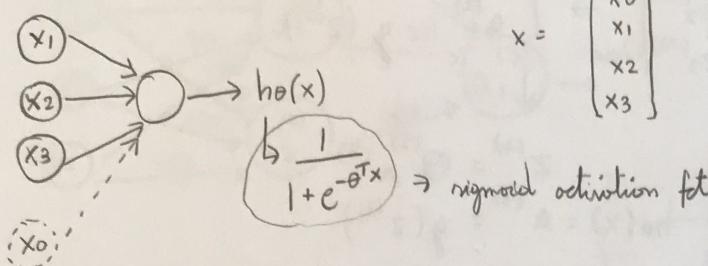
$$\rightarrow \text{But } h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$\hookrightarrow J(\theta) = \left[-\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

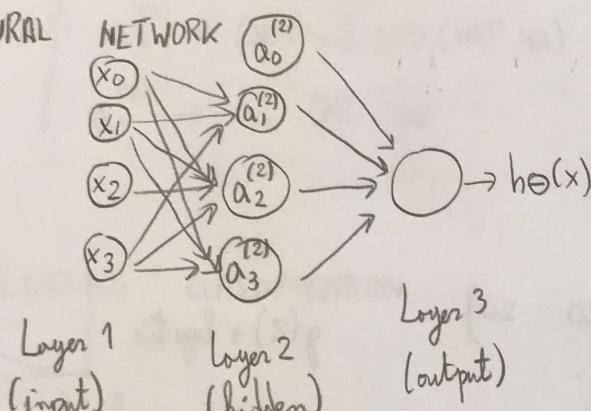
9

WEEK 4

NEURON MODEL: LOGISTIC UNIT



NEURAL



$\Theta^{(j)}$ \rightarrow matrix of weights mapping from layer j to $j+1$

$$\text{Ex: } a_2^{(2)} = g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3)$$

$$h_\theta(x) = a_1^{(3)} = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})$$

$$\Theta^{(1)} \in \mathbb{R}^{3 \times 4}$$

\rightarrow Network with s_j units in layer j and s_{j+1} units in layer $j+1 \rightarrow \Theta^{(j)}$ dim $s_{j+1} \times (s_j + 1)$

EX

$$\Theta^{(1)} \rightarrow 3 \times 4$$

$$\begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = \begin{bmatrix} \Theta_{10} & \Theta_{11} & \Theta_{12} & \Theta_{13} \\ \Theta_{20} & \Theta_{21} & \Theta_{22} & \Theta_{23} \\ \Theta_{30} & \Theta_{31} & \Theta_{32} & \Theta_{33} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$3 \times 1 \quad 3 \times 4 \quad 4 \times 1$

$$\Theta^{(2)} \rightarrow [a_1] = [\Theta_{10} \ \Theta_{11} \ \Theta_{12} \ \Theta_{13}] \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \end{bmatrix}$$

$1 \times 4 \quad 1 \times 4 \quad 4 \times 1$

(10)

FORWARD PROPAGATION

$$\left\{ \begin{array}{l} a_1^{(2)} = g(z_1^{(2)}) \\ a_2^{(2)} = g(z_2^{(2)}) \\ a_3^{(2)} = g(z_3^{(2)}) \end{array} \right. \rightarrow z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} \rightarrow \left\{ \begin{array}{l} z^{(2)} = \Theta^{(1)} x = \Theta^{(1)} a^{(1)} \\ a^{(2)} = g(z^{(2)}) \end{array} \right.$$

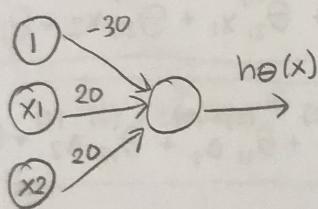
↓

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

→ NNs learn their own features

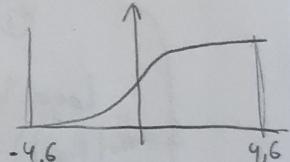
EXAMPLES:

LOGICAL AND



$$\Theta^{(1)} = [-30 \quad 20 \quad 20]$$

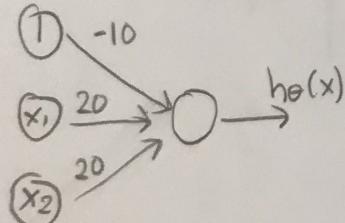
$g(z) \rightarrow$ logistic



$$h_{\theta}(x) = g(-30 + 20x_1 + 20x_2)$$

x_1	x_2	$h_{\theta}(x) \approx$
0	0	0
0	1	0
1	0	0
1	1	1

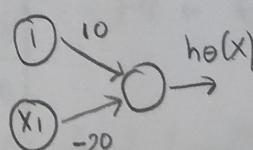
LOGICAL OR



x_1	x_2	$h_{\theta}(x) \approx$
0	0	0
0	1	1
1	0	1
1	1	1

$$h_{\theta}(x) = g(-10 + 20x_1 + 20x_2)$$

NOT

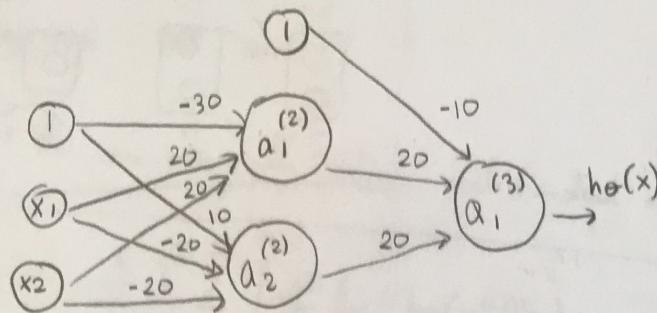


x_1	$h_{\theta}(x) \approx$
0	1

$$h_{\theta}(x) = g(10 - 20x_1)$$

(11)

XNOR

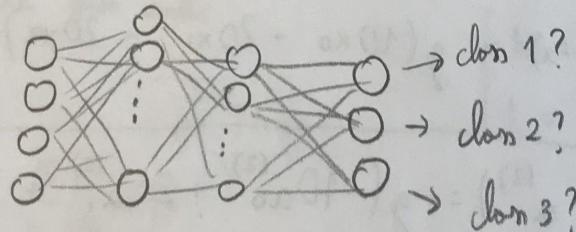


$$\left\{ \begin{array}{l} a_1^{(2)} \rightarrow x_1 \text{ AND } x_2 \\ a_2^{(2)} \rightarrow (\text{NOT } x_1) \text{ AND } (\text{NOT } x_2) \\ a_3^{(3)} \rightarrow x_1 \text{ OR } x_2 \end{array} \right.$$

x_1	x_2	$a_1^{(2)}$	$a_2^{(2)}$	$h_\theta(x) \approx$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

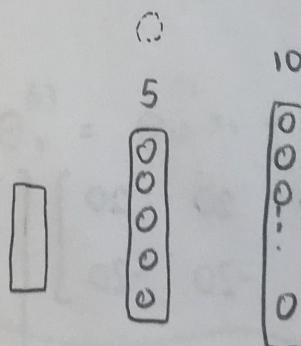
MULTICLASS CLASSIFICATION

3 classes



$$h_\theta(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \rightarrow \text{class 2}$$

Example : { 5 hidden units
10 classes (output) } \rightarrow

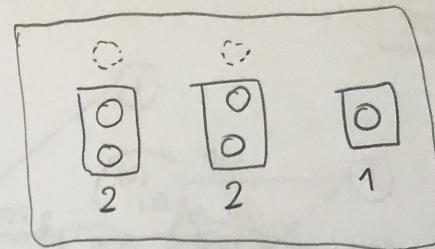


$$\Theta^{(2)} \rightarrow 10 \times (5+1) \rightarrow 60 \text{ elements (weights)}$$

(12)

XNOR IN DETAIL

$$\left\{ \begin{array}{l} \text{AND} \rightarrow \Theta^{(1)} = [-30 \quad 20 \quad 20] \\ \text{NOR} \rightarrow \Theta^{(1)} = [10 \quad -20 \quad -20] \\ \text{OR} \rightarrow \Theta^{(1)} = [-10 \quad 20 \quad 20] \end{array} \right.$$



$$\text{XNOR: } \left\{ \begin{array}{l} \Theta^{(1)} = \begin{bmatrix} -30 & 20 & 20 \\ 10 & -20 & -20 \end{bmatrix} \\ \Theta^{(2)} = [-10 \quad 20 \quad 20] \end{array} \right. \rightarrow \left\{ \begin{array}{l} a^{(2)} = g(\Theta^{(1)} \cdot x) \\ a^{(3)} = g(\Theta^{(2)} \cdot a^{(2)}) \\ h_\Theta(x) = a^{(3)} \end{array} \right.$$

$$\rightarrow a_1^{(2)} = g(\Theta_1^{(1)} \cdot x) = g(-30x_0 + 20x_1 + 20x_2)$$

$$\rightarrow a_2^{(2)} = g(\Theta_2^{(1)} \cdot x) = g(10x_0 - 20x_1 - 20x_2)$$

$$\rightarrow a_1^{(3)} = g(\Theta_1^{(2)} \cdot a^{(2)}) = g(-10a_0^{(2)} + 20a_1^{(2)} + 20a_2^{(2)})$$

Ex:

$$Z^{(2)} = \Theta^{(1)} x = \begin{bmatrix} -30 & 20 & 20 \\ 10 & -20 & -20 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

$$Z_1^{(2)} = -30 + 20x_1 + 20x_2 \rightarrow a_1^{(2)} = g(Z_1^{(2)})$$