

overfitting \rightarrow too many
 ↘ high bias \rightarrow linear
 ↗ high variance \rightarrow high order
 reduce features.

regularization: penalize the params.

(penalize term) regularization term

$$\left[\lambda \sum_{j=1}^n \theta_j^2 \right] \xrightarrow{\text{regulation parameters}}$$

low \rightarrow ineffective

high \rightarrow underfitting

$$\begin{aligned} \theta_j &= \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta - y) x_j + \frac{\lambda}{m} \theta_j \right] \\ &= \theta_j \left(1 - \alpha \frac{\lambda}{m} \right) - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta - y) x_j \end{aligned}$$

shrinkage

gradient
descent

normal equation:

$$\theta = (\mathbf{X}^T \mathbf{X} + \lambda I)^{-1} \mathbf{X}^T y$$

invertible

logistic classification

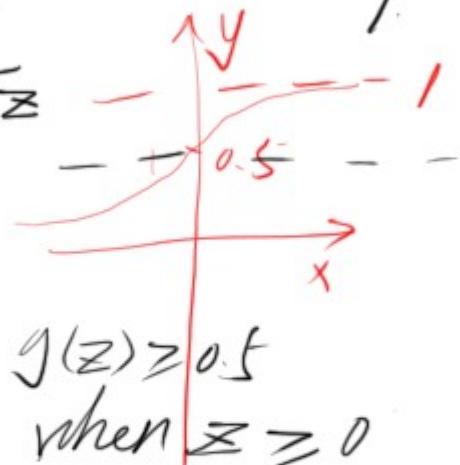
$y \in \{0, 1\}$ multo $y \in \{0, 1, 2, \dots, k\}$

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

estimate P of $y=1$

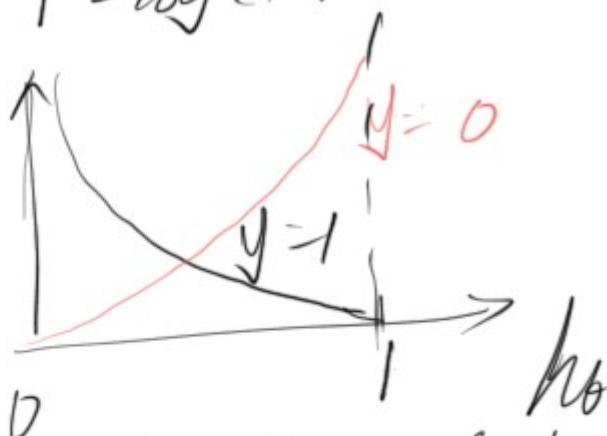
$$h_\theta(x) = P(y=1 | x \theta)$$

$$\text{For } \theta^T x \geq 0 \Rightarrow h_\theta \geq 0.5$$



cost function

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



simplified cost function

$$\text{cost}(h_\theta(x), y) = -y \log(h_\theta) - (1-y) \log(1-h_\theta)$$

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \text{cost} \right]$$

$$\text{Gradient Descent: } \theta_j = \theta_j - \alpha \sum_{i=1}^m (h_\theta - y) x_j$$

Neuron model:
logistic unit



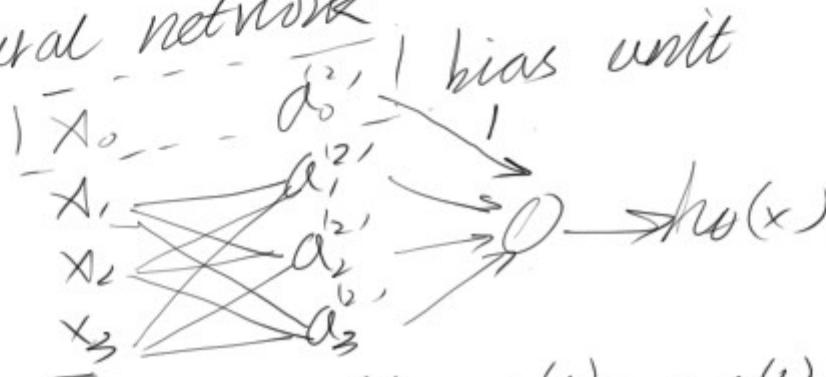
$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

↓

"weights"

activation function

neural network



$$a^{(2)} = g(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3)$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

$$z^{(2)} = \theta^{(1)} x$$

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

add $a_0^{(2)} = 1$ ← bias unit

$$z^{(3)} = \theta^{(2)} a^{(2)}$$

$$h_0(x) = a^{(3)} = g(z^{(3)})$$

$$x \rightarrow a \rightarrow h_0$$

logR logR

layer₁ → layer₂ → ... h₀(x)

$$x_1 \rightarrow a_1$$

$$x_2 \rightarrow a_2$$

$$x_3 \rightarrow a_3$$

final result

multiclassification

neural network with multiple outputs

$$h_{\theta} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad h_{\theta} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

Type 1

Type 2

Training set $(x^{(i)}, y^{(i)})$
 $y^{(i)}$ one of $\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$

cost function

Training set, L , total no of layers

s_L : No of unit in layer L

$h_{\theta}(x) \in R^k$

$$J_{\theta} = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^k y_k^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y_k^{(i)}) \log(1 - h_{\theta}(x^{(i)}))_k \right]$$

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_L} \sum_{j=1}^{s_{l+1}} (\theta_{ji}^{(l)})^2$$

layer no bias term

Neural network example
non-linear classification

$$\begin{array}{c} \uparrow x_2 \\ \phi \end{array}$$

\times

x_1 and x_2 are binary

$$y = x_1 \text{ XOR } x_2$$

$$x_1 \text{ XNOR } x_2$$

$$\begin{array}{c} \rightarrow x_1 \\ \oplus \end{array}$$

$$\begin{array}{c} \oplus \rightarrow -30 \\ \oplus \end{array}$$

$$x_1, x_2 \in \{0, 1\}$$

$$\begin{array}{c} \oplus \rightarrow +20 \\ \oplus \end{array} \rightarrow 0 \rightarrow h_0(x)$$

$$y = x_1 \text{ AND } x_2$$

$$\begin{array}{c} \oplus \rightarrow +20 \\ \oplus \end{array}$$

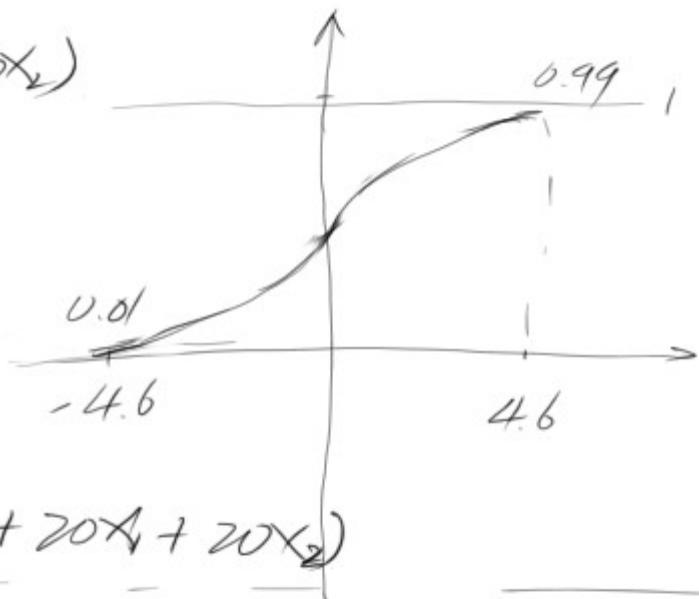
$$h_0(x) = g(-30 + 20x_1 + 20x_2)$$

input all x_1 and x_2

the result will be

the same as

$$x_1 \text{ AND } x_2$$



Negation

$$\begin{array}{c} \oplus \rightarrow 10 \\ \oplus \end{array} \rightarrow 0 \rightarrow h_0(x)$$

$$\begin{array}{c} \oplus \rightarrow -20 \\ \oplus \end{array}$$

$$h_0 = g(10 - 20x_1)$$

large negating factor

$$\begin{array}{c} \oplus \rightarrow -30 \\ \oplus \end{array}$$

$$\begin{array}{c} \oplus \rightarrow 10 \\ \oplus \end{array}$$

$$\begin{array}{c} \oplus \rightarrow 20 \\ \oplus \end{array}$$

$$\begin{array}{c} \oplus \rightarrow 20 \\ \oplus \end{array}$$

$$\begin{array}{c} \oplus \rightarrow -20 \\ \oplus \end{array}$$

$$\begin{array}{c} \oplus \rightarrow -10 \\ \oplus \end{array}$$

$$\begin{array}{c} \oplus \rightarrow 20 \\ \oplus \end{array} \rightarrow h_0(x)$$

NOR

AND

$$\begin{array}{c} \oplus \rightarrow 10 \\ \oplus \end{array}$$

OR

$$(Wx_1) \text{ AND } (Wx_2)$$

Backpropagation

$J(\theta)$, $\frac{\partial}{\partial \theta_{ij}} J(\theta)$ required.

$\delta_j^l = \text{"error" of node } j \text{ in layer } l$

$\delta_j^4 = \hat{a}_j^4 - y_j \rightarrow$ the value observed, actual
calculated by the hypothesis

$$\delta^3 = (\theta^3)^T \delta^4 * g'(\hat{z}^3) \leftarrow \hat{a}^3 * (1 - \hat{a}^3)$$

$$\frac{\partial}{\partial \theta_{ij}^l} J(\theta) = \hat{a}_j^{(l)} \delta_i^{(l+1)}$$

procedure: $\begin{cases} a' = \hat{a}' \\ \text{forward propagation, using hypo} \\ \text{compute } \delta^{(L-1)} \dots \delta^2 \end{cases}$

$$\begin{aligned} \Delta_{ij}^l &= \Delta_{ij}^l + \hat{a}_j^l \delta_i^{l+1} \\ &= \Delta^l + \delta^{l+1} (\hat{a}^l)^T \end{aligned}$$

$$D_{ij}^l = \frac{1}{m} \Delta_{ij}^l + \underbrace{\frac{\partial \theta_{ij}^l}{\partial \theta_{ij}^l}}_{1.} \text{ for } i \neq 1.$$

$$\frac{\partial}{\partial \theta_{ij}^l} J(\theta) = D_{ij}^l$$

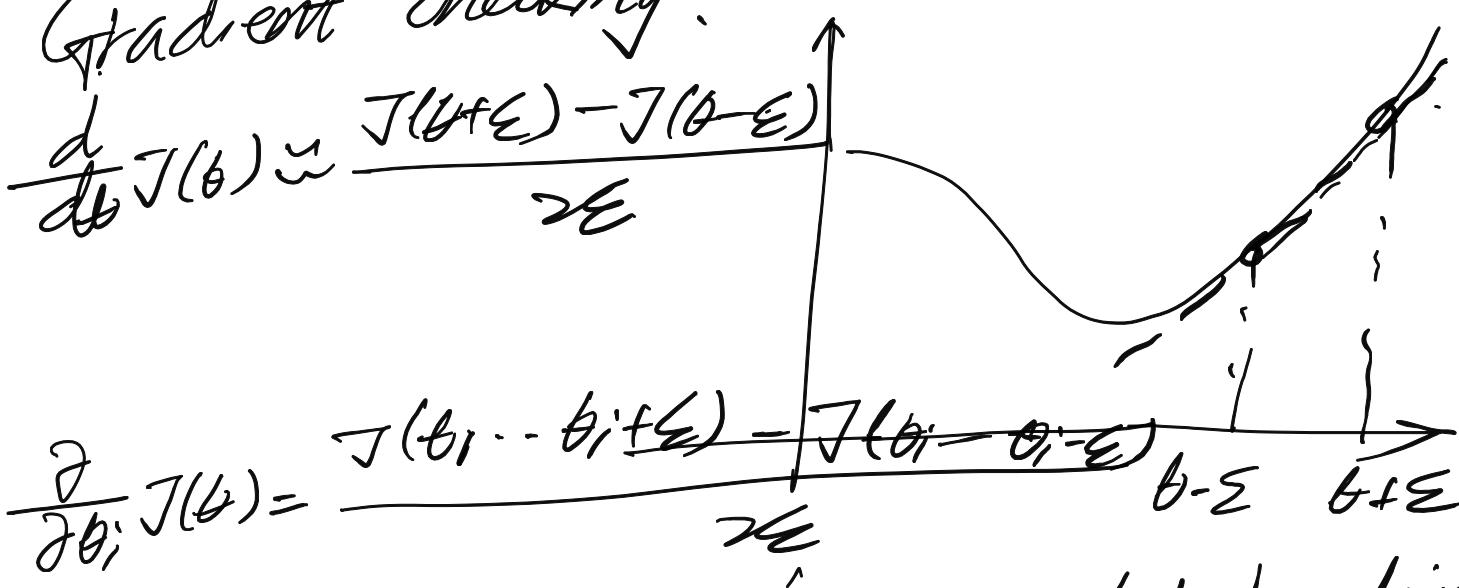
unrolling parameters

matrix \Rightarrow vector for optimization

$\text{ThetaVec} = [\theta_{01}(:) \dots \theta_{0n}(:) \dots \theta_{l1}(:) \dots \theta_{ln}(:)]$

$\text{thetaVec} \Rightarrow \text{Theta}$
 $= \text{reshape}(\text{thetaVec}(1:\text{total_#}),$
 $n_j, n_j)$

Gradient checking:



Random Initialization: Symmetry breaking
Small value

1. pick a network architecture
2. random initial weight
3. PP
4. cost function
5. BP $\frac{\partial}{\partial \theta} J(\theta)$ — loop $\Delta = \Delta + \alpha^T$
6. Gradient checking, numeric checking
7. minimize cost function

machine learning diagnostic:
evaluating hypothesis

split into "training set" and "test set"
"shuffle" the data

$J(\theta)$ training $J(\theta)$ test both low

$$\frac{1}{m_{\text{test}}} \sum (h_{\theta}(x) - y_{\text{test}})$$

For logistic regression some residual.
misclassification error, prediction accuracy.

Model selection:

