

SCHEDULE

	MON	TUE	WED	THU	FRI

Important dates

Time Zones (US)



International

New York	10 am
Los Angeles	9 am
London	8 am
Moscow	7 am
Tokyo	6 am
Beijing	5 am
Rome	4 am
Hong Kong	3 am
Bangkok	2 am
Sydney	1 am
Parrs	12 am
Caro	11 am
Mumbai	10 am
Shanghai	9 am
Dhaka	8 am
Calcutta	7 am
Dubai	6 am
Yokohama	5 am

*Times not take DST daylight savings into account

this is so thin, what the fuck
check check

Week 1 notes

$h_\theta(x)$ → prediction model

y → actual model

minimize $\frac{1}{2m} \sum_{i=1}^m [h_\theta(x^{(i)}) - y^{(i)}]^2$ most commonly used
squared error
function

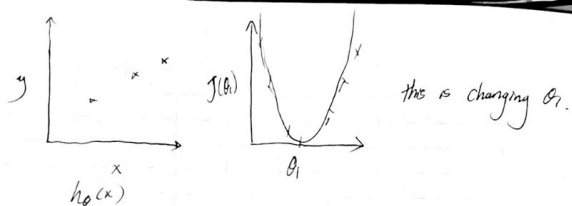
$x^{(i)}$ training set $\rightarrow h_\theta(x^{(i)}) = \theta_0 + \theta_1 \cdot x^{(i)}$

cost function = $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m [h_\theta(x^{(i)}) - y^{(i)}]^2$
minimize $J(\theta_0, \theta_1)$

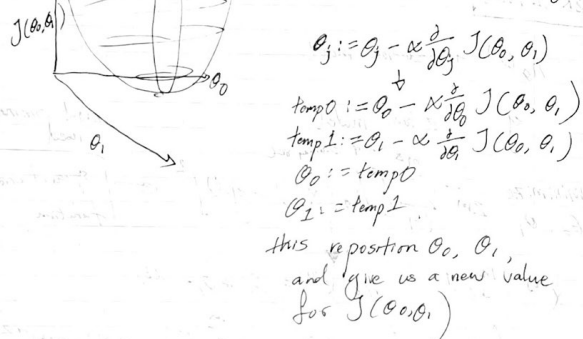
Hypothesis $h_\theta(x) = \theta_0 + \theta_1 \cdot x$ Parameters θ_0, θ_1

cost function $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$

Goal : minimize $J(\theta_0, \theta_1)$



if we change multiple θ , then we see high dimensional graphs



$A^T \rightarrow$ transpose

$\text{inv}(A) \rightarrow$ inverse

Week 2 Notes

multiple features n . $x^{(i)}$ = input (features) of i^{th} training example
 $x_j^{(i)}$ = value of feature j

size	# of rooms, x_2	floors, x_3	age of home, x_4	price, y
x_1				
2104	4	1	45	460
1416	3	2	40	232
1410				
852				

$x_1^{(4)} = 852 \rightarrow 4^{\text{th}}$ row of feature 1.

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

$$\text{example} = 80 + 0.1 \cdot x_1 + 0.01 \cdot x_2 + 3 \cdot x_3 - 2 \cdot x_4$$

$$x_0^{(i)} = x_0 = 1$$

$$\text{if } x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}, \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1} \text{ then } h_0(x) = \theta^T x$$

$(n+1) \times (1)$ $(n+1) \times (1)$

Hypothesis $h_0(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

Parameters $\theta_0, \theta_1, \theta_2, \dots$ aka θ

Cost function $J(\theta_0, \theta_1, \theta_2, \dots, \theta_n) = J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$

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

$$= \frac{1}{2m} \sum_{i=1}^m \left(\sum_{j=0}^n \theta_j x_j^{(i)} - y^{(i)} \right)^2$$

Gradient descent

Repeat ∞

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

Partial Derivative Demonstration

\sum

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

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} \left[\frac{1}{2m} \sum_{i=1}^m \left(\sum_{j=0}^n \theta_j x_j^{(i)} - y^{(i)} \right)^2 \right]$$

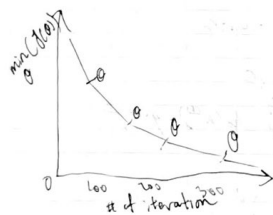
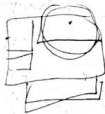
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} \left[\frac{1}{2m} \sum_{i=1}^m \left(\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n - y^{(i)} \right)^2 \right]$$

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m \left(\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n - y^{(i)} \right) \cdot x_j^{(i)} \right]$$

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

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



W

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

W bruhs this isnt trash.

$$\text{pinv}(X' * X) * X' * y$$

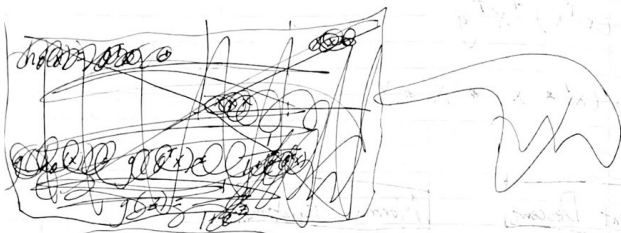
Gradient Descent

\sum , need α ,
works well n large.

Normal Equation

matrix, find θ .
no good if n large.

Week 3 Notes



$$h_{\theta}(x) = g(\theta^T x)$$

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

example:

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

predict $y=1$ if $-3 + x_1 + x_2 \geq 0$

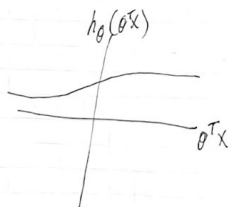
example:

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

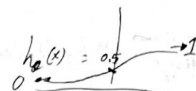
$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{bmatrix} \quad \text{predict } y=1 \text{ if } \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 \geq 0$$

$$h_{\theta}(x) \geq 0.5 \rightarrow y=1$$

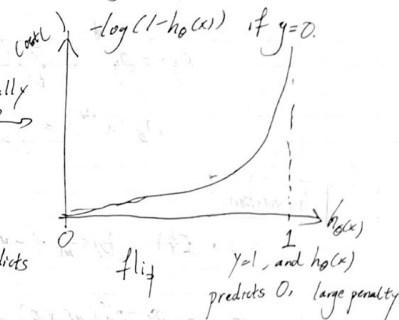
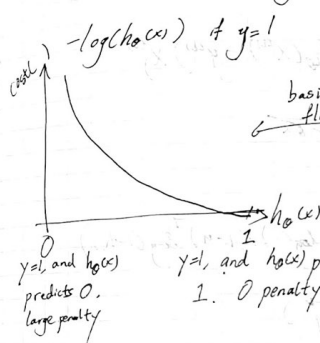
$$h_{\theta}(x) < 0.5 \rightarrow y=0$$



Logistic Regression Cost Function



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



$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

$$-y \cdot \log(h_{\theta}(x)) - (1-y) \cdot \log(1-h_{\theta}(x))$$

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

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

$$\min_{\theta} J(\theta) \quad \text{repeat } \xi$$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

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

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

Vectorization

$$h = g(\mathbb{X}\theta) \quad J(\theta) = \frac{1}{m} \cdot (-y^T \log(h) - (1-y)^T \log(1-h))$$

$$\theta := \theta - \frac{\alpha}{n} \mathbb{X}^T (g(\mathbb{X}\theta) - \vec{y})$$

Regularization

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

$$\begin{aligned} &\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \\ &\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \end{aligned} \quad \left. \vphantom{\begin{aligned} &\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \\ &\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \end{aligned}} \right\} \text{same format}$$

Repeat ξ

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)} \quad j=0$$

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad j=1, 2, 3, \dots, n$$

ξ

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

can be rewritten as

$$\mathbb{X} = \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(m)} \end{bmatrix}^T \rightarrow \text{first data.}$$

$m \leq n$
examples # features

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

if $\lambda > 0$

$$\theta = (X^T X + \lambda \begin{bmatrix} 0 & & \\ & 1 & \\ & & \ddots \\ & & & 1 \end{bmatrix})^{-1} X^T y$$

use this if no inverse pinv vs inv

Logistic regression without regularization

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

With regularization

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

derivative:

$$\theta_0 \rightarrow \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_j \rightarrow \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j$$

EX2 workspace

$$X = \begin{bmatrix} 1 & 0.1 & 0.02 \\ \vdots & \vdots & \vdots \\ 1 & 0.1 & 0.02 \end{bmatrix} \quad Y = \begin{bmatrix} 0 \\ \vdots \\ 1 \end{bmatrix}$$

EX2.m

100 x 2

initial theta zero (3, 1)

theta (1, 1)

$$h = \text{sigmoid}(X \cdot \text{theta});$$

$$h = \text{sigmoid} \left[\begin{matrix} \theta_0 + \theta_1 x_1 + \theta_2 x_2 \\ \vdots \end{matrix} \right] \rightarrow h = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

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

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

$h = \text{sigmoid}(X \cdot \text{theta})$

active code

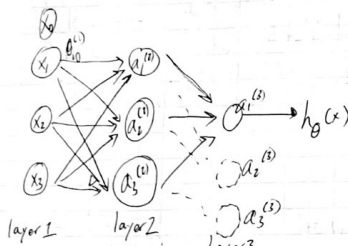
$$J(\theta) = \frac{1}{m} \sum [-y \cdot \log(h) - (1 - y) \cdot \log(1 - h)]$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \left(\frac{1}{m} \sum [(h - y) \cdot X] \right)_j$$

$$X \cdot Y = \begin{bmatrix} 2 \\ 6 \end{bmatrix} \quad X \cdot Y = 17$$

$$X = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad Y = \begin{bmatrix} 2 \\ 6 \\ 7 \end{bmatrix}$$

Week 4 Notes



$a_i^{(j)}$ = activation of unit i in layer j

$a_1^{(2)}$ = activation of unit 1 in layer 2.

$\theta^{(j)}$ = matrix of weights controlling function mapping from layer j to layer $j+1$.

$$a_1^{(2)} = g(\theta_{10}^{(2)} x_0 + \theta_{11}^{(2)} x_1 + \theta_{12}^{(2)} x_2 + \theta_{13}^{(2)} x_3) \quad z_1^{(2)}$$

$$a_2^{(2)} = g(\theta_{20}^{(2)} x_0 + \theta_{21}^{(2)} x_1 + \theta_{22}^{(2)} x_2 + \theta_{23}^{(2)} x_3) \quad z_2^{(2)}$$

$$a_3^{(2)} = g(\theta_{30}^{(2)} x_0 + \theta_{31}^{(2)} x_1 + \theta_{32}^{(2)} x_2 + \theta_{33}^{(2)} x_3) \quad z_3^{(2)}$$

$$h_\theta(x) = a_1^{(2)} = g(\theta_{10}^{(2)} a_0^{(1)} + \theta_{11}^{(2)} a_1^{(1)} + \theta_{12}^{(2)} a_2^{(1)} + \theta_{13}^{(2)} a_3^{(1)}) \quad z_1^{(2)}$$

$$a_1^{(2)} = g(z_1^{(2)}) \quad z_1^{(2)} = \theta_{10}^{(2)} x_0 + \theta_{11}^{(2)} x_1 + \theta_{12}^{(2)} x_2 + \theta_{13}^{(2)} x_3$$

$$a_2^{(2)} = g(z_2^{(2)})$$

$$a_3^{(2)} = g(z_3^{(2)}) \quad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} \quad a^{(2)} = g(z^{(2)})$$

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

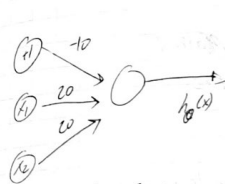
$$z_k^{(j)} = \theta_{k,0}^{(j)} x_0 + \theta_{k,1}^{(j)} x_1 + \dots + \theta_{k,n}^{(j)} x_n$$

layer $j=2$ node k , variable z .

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad z^{(j)} = \begin{bmatrix} z_1^{(j)} \\ z_2^{(j)} \\ \vdots \\ z_n^{(j)} \end{bmatrix}$$

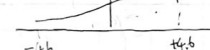
$$z_1^{(2)} = \theta_{10}^{(2)} a_0^{(1)} + \theta_{11}^{(2)} a_1^{(1)} + \theta_{12}^{(2)} a_2^{(1)} + \theta_{13}^{(2)} a_3^{(1)}$$

$$z^{(j)} = \theta^{(j-1)} a^{(j-1)}$$

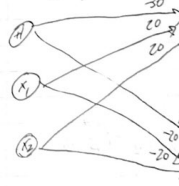


$$\theta^1 = [10, 20, 20]$$

$x_1, x_2 \in \{0, 1\}$	$g(-b + 10x_1 + 20x_2)$
$x_1 \backslash x_2$	
0 0	$g(-10) \rightarrow 0$
0 1	$g(0) \rightarrow 1$
1 0	$g(10) \rightarrow 1$
1 1	$g(30) \rightarrow 1$



$(x_1) XOR (x_2)$



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

Ex3 Workspace

over all $X, y, \text{num_labels}, \text{lambda}$
 ↓
 image data whether it is X or y
 5000 x 400 5000 x 1 $m = 5000$
 ↓ each image $n = 400$
 28x28 (gray scale pixel)
 all theta = zeros Matrix
 10 x 402
 $X = 5000 \times 401$

10 classes for 10 numbers
 $\theta_0^{(0)}$
 $\theta_1^{(1)}$
 $\theta_2^{(1)}$
 $\theta_3^{(1)}$
 $\theta_4^{(1)}$

$$a_1^{(1)} = g(\theta_{10}^{(0)} X_0 + \theta_{11}^{(0)} X_1 + \theta_{12}^{(0)} X_2 + \theta_{13}^{(0)} X_3)$$

classifier

lr cost function (theta, X, y, lambda)

10 classifier for each number
different

lr cost function (theta, X, y, lambda)

@ t t X y = C lambda

$X = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \end{bmatrix}$ $X(:,1) = \begin{bmatrix} 1 \\ 5 \end{bmatrix}$ $X(:,2) = \begin{bmatrix} 2 \\ 6 \end{bmatrix}$
 $X(:,3) = \begin{bmatrix} 3 \\ 7 \end{bmatrix}$ $X(:,4) = \begin{bmatrix} 4 \\ 8 \end{bmatrix}$
 $X(:,1) = 1, 2, 3, 4$
 $X(:,2) = 5, 6, 7, 8$

options = optimset('GradObj', 'on', 'MaxIter', 50);
 min iteration 50

for $c = 1 : \text{num_labels}$ = for $c = 1 : 10$ all theta = zeros(10, 401)
 all theta (c,:) = ... this means continue
 fmincg(c) 401 is for each pixel intensity

10 x 401
 return the θ parameters for each of the 10 classes

for $p = 1 : \text{size}(X, 1)$
 $X(\text{Peep}, :) =$
 1 x 401

all theta 10 x 401

std: 1, 16, 7

167

Predict One Vs All (all-theta, X)

$$X = 5000 \times 400$$

↓

$$X = 5000 \times 401 \rightarrow 1 \times 401$$

$$\text{all-theta} = 10 \times 401$$

$$X = 5000 \times 400$$

$$m = 5000$$

$$\text{num-labels} = 10$$

$$p = 5000 \times 1$$

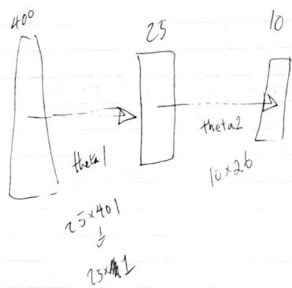
$$\text{theta1} = 25 \times 401$$

$$\text{theta2} = 10 \times 26$$

X' all-theta'

$$p(1,:) = X(1,:) * (\text{all-theta})'$$

gives you a 1x10 matrix of probabilities for the first image 20x20.



$$25 \times \begin{bmatrix} 401 \\ \vdots \end{bmatrix}$$

↓

$$25 \times \begin{bmatrix} 1 \\ \vdots \end{bmatrix} \rightarrow 26 \times \begin{bmatrix} \text{one}(\cdot) \\ \vdots \end{bmatrix}$$

$$\begin{bmatrix} 401 \times 1 \\ \vdots \end{bmatrix} \rightarrow X(\text{biggens}, :)'$$

$$\text{theta1} = 25 \times 401$$

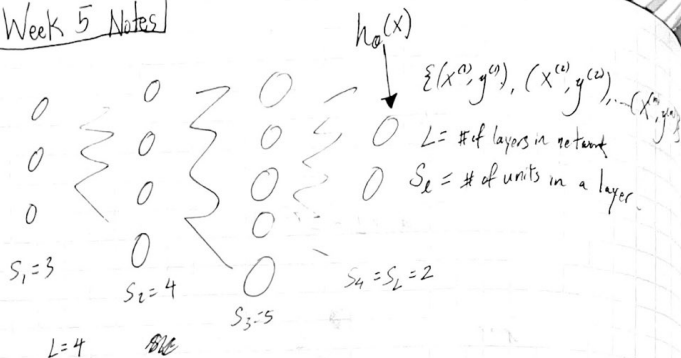
$$\text{sig} \left(25 \times \begin{bmatrix} 1 \\ \vdots \end{bmatrix} \right) \Rightarrow 25 \times \begin{bmatrix} 1 \\ \vdots \end{bmatrix}$$

$$\text{sigmoid} \left[\begin{matrix} 25 \times 401 \\ \vdots \end{matrix} , X(\text{biggens}, :) \right]$$

$$\text{theta2} = 10 \times \begin{bmatrix} 26 \\ \vdots \end{bmatrix}$$

$$(\text{theta2} * \text{biggens}) = 10 \times 1 \times \begin{bmatrix} 1 \\ \vdots \end{bmatrix}$$

Week 5 Notes



Binary Classification

$y = 0$ or 1

$$S_L = 1 \quad h_\theta(x) \in \mathbb{R}$$

Multi-Class Classification (k classes)

$$y \in \mathbb{R}^K \quad \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

$S_L = K \quad k \geq 3$

$K \times K \quad K \times K$

Logistic Regression

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

Neural Network

$$h_\theta(x) \in \mathbb{R}^K \quad (h_\theta(x))_i = i^{\text{th}} \text{ output}$$

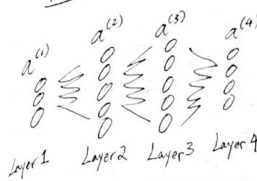
$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)}))_k + (1-y_k^{(i)}) \log(1-(h_\theta(x^{(i)}))_k) + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{ji}^{(l)})^2$$

Gradient Computation

min $J(\theta)$ compute $J(\theta)$

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

Forward Propagation



activation values

$$\begin{aligned} a^{(1)} &= x \\ z^{(2)} &= \theta^{(1)} a^{(1)} \\ a^{(2)} &= g(z^{(2)}) \quad (\text{add } a_0^{(1)}) \\ z^{(3)} &= \theta^{(2)} a^{(2)} \\ a^{(3)} &= g(z^{(3)}) \quad (\text{add } a_0^{(2)}) \\ z^{(4)} &= \theta^{(3)} a^{(3)} \\ a^{(4)} &= h_\theta(x) = g(z^{(4)}) \end{aligned}$$

Backward Propagation

$\delta_j^{(l)}$ = the error of $a_j^{(l)}$ node

$$\delta_j^{(l)} = a_j^{(l)} - y_j \quad \text{vectorized} \rightarrow \delta^{(l)} = a^{(l)} - y$$

$$\delta^{(l)} = (\theta^{(l)})^T \delta^{(l+1)} * \frac{a'(z^{(l)})}{\text{derivative}} \quad \delta^{(l)} = \begin{bmatrix} \delta_1^{(l)} \\ \delta_2^{(l)} \\ \delta_3^{(l)} \end{bmatrix} \rightarrow \theta_1^{(l)} \delta_1^{(l+1)} + \theta_2^{(l)} \delta_2^{(l+1)} + \theta_3^{(l)} \delta_3^{(l+1)}$$

$$\delta^{(l)} = (\theta^{(l)})^T \delta^{(l+1)} * g'(z^{(l)}) \quad a^{(l)} * (1-a^{(l)})$$

$$\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}} = a_j^{(l)} \theta_i^{(l+1)} \quad \text{ignore } \lambda; \text{ if } \lambda=0$$

$$\Delta_{ij}^{(l)} \leftarrow \text{current } \delta, \text{ use to capture } \frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}}$$

for $i=1$ to m
 get $a^{(1)} = x^{(i)}$
 forward propagation to compute $a^{(l)}$ $l=2,3,\dots,L$
 compute $\delta^{(L)} = a^{(L)} - y^{(i)}$ \rightarrow actual answer.
 then $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(1)}$ \rightarrow going backwards.

$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$$

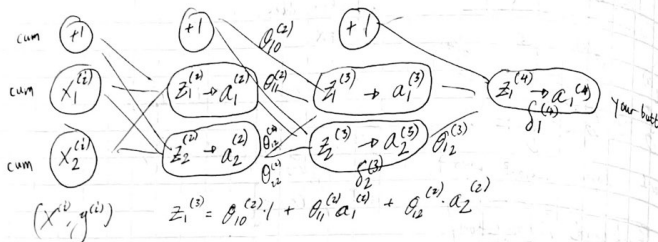
Two cases

$$\begin{aligned} D_{ij}^{(l)} &:= \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)} \quad \text{if } i \neq 0 \\ D_{ij}^{(l)} &:= \frac{1}{m} \Delta_{ij}^{(l)} \quad \text{if } i = 0 \end{aligned}$$

complex math proof shows that $\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}} = D_{ij}^{(l)}$

Forward Propagation

act. activation value



$\delta_j^{(l)}$ = "error" of cost for unit j layer l

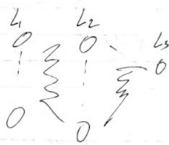
$$\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \cdot \text{cost}(i) \quad \text{where } \text{cost}(i) = y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log h_\theta(x^{(i)})$$

$$\delta_1^{(2)} = g^{(1)} - a_1^{(1)} \quad \delta_2^{(2)} = \theta_{12}^{(2)} \cdot \delta_1^{(1)} \quad \delta_2^{(2)} = \theta_{12}^{(2)} \cdot \delta_1^{(1)} + \theta_{22}^{(2)} \cdot \delta_2^{(1)}$$

whatever output you choose

$$\delta_2^{(2)} = \delta_2^{(2)} + \theta_{22}^{(2)} \cdot \delta_2^{(1)} \quad \delta_2^{(2)} = \delta_2^{(2)} (1 + \theta_{22}^{(2)})$$

(θ) parameter (D) gradients



$$S_1 = 10 \quad S_2 = 10 \quad S_3 = 1$$

$$\theta^{(1)} \in \mathbb{R}^{10 \times 11} \quad \theta^{(2)} \in \mathbb{R}^{10 \times 11} \quad \theta^{(3)} \in \mathbb{R}^{1 \times 11}$$

$$S_2 = 10 \text{ nodes } S_1 = 11 \text{ nodes } + 1$$

$$D^{(1)} \in \mathbb{R}^{10 \times 11} \quad D^{(2)} \in \mathbb{R}^{10 \times 11} \quad D^{(3)} \in \mathbb{R}^{1 \times 11}$$

$$\text{thetaVec} = [\text{Theta1}(:); \text{Theta2}(:); \text{Theta3}(:)];$$

$$DVec = [D1(:); D2(:); D3(:)];$$

$$\text{Theta1} = \text{reshape}(\text{thetaVec}(1:10), 10, 11);$$

$$\text{Theta2} = \text{reshape}(\text{thetaVec}(11:20), 10, 11);$$

$$\text{Theta3} = \text{reshape}(\text{thetaVec}(21:25), 1, 11);$$

→ reshape back to original

Gradient could have a bug. So use gradient checking



grad approx $\approx DVec$

Putting it together

$$\frac{d}{d\theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

$$\text{gradApprox} = \frac{J(\text{theta} + \text{EPSILON}) - J(\text{theta} - \text{EPSILON})}{2 \cdot \text{EPSILON}}$$

$$\frac{d}{d\theta_j} J(\theta) \approx \frac{J(\theta_1, \dots, \theta_j + \epsilon, \dots, \theta_n) - J(\theta_1, \dots, \theta_j - \epsilon, \dots, \theta_n)}{2\epsilon}$$

$$g = \begin{bmatrix} \frac{\partial}{\partial \theta_1} J(\theta) \\ \frac{\partial}{\partial \theta_2} J(\theta) \\ \frac{\partial}{\partial \theta_3} J(\theta) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial \theta_1} J(\theta) \\ \frac{\partial}{\partial \theta_2} J(\theta) \\ \frac{\partial}{\partial \theta_3} J(\theta) \end{bmatrix}$$