

~~WEEK 1~~
this is so thin what the fuck
check due if

Week 1 notes

$h_{\theta}(x) \rightarrow$ prediction model

$y \rightarrow$ actual model

$$\text{minimize}_{\theta_0, \theta_1} \frac{1}{2m} \sum_{i=1}^{m \# \text{ of training set}} [h_{\theta}(x^{(i)}) - y^{(i)}]^2 \rightarrow \begin{cases} \text{most commonly used} \\ \text{squared error function} \end{cases}$$

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

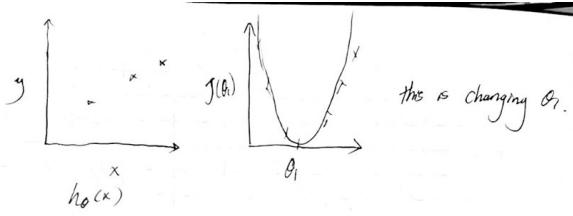
$$\text{cost function} = J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}]^2$$

$$\text{minimize}_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

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

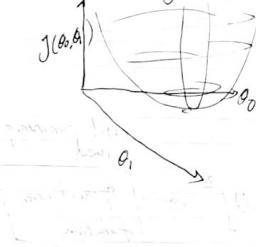
Cost function $J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Goal : $\text{minimize}_{\theta_0, \theta_1} J(\theta_0, \theta_1)$



this is changing θ_1 .

If we change multiple θ , then we get high dimensional graphs



$$\begin{aligned}\theta_j &:= \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \\ \text{temp}_0 &:= \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \text{temp}_1 &:= \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \\ \theta_0 &:= \text{temp}_0 \\ \theta_1 &:= \text{temp}_1\end{aligned}$$

This reposition θ_0, θ_1 , and give us a new value for $J(\theta_0, \theta_1)$

$A' \rightarrow$ 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	floors	age of home	price
x_1	x_2	x_3	x_4	y
2.104	5	1	45	460
1.116	3	2	40	232
1.107	2	2	30	150
852				

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

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

$$\text{example} = \theta_0 + 0.1 \cdot x_1 + 0.01 \cdot x_2 + 3 \cdot x_3 - 2 \cdot x_4$$

$$\theta_0^{(i)} = \theta_0 = 1$$

$$\text{if } x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \text{then } h_0(x) = \theta^T x$$

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

$$\boxed{\text{Parameters}} \quad \theta_0, \theta_1, \theta_2, \dots \text{aka } \theta$$

$$\begin{aligned}\boxed{\text{Cost function}}: J(\theta_0, \theta_1, \theta_2, \dots, \theta_n) &= J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2 \\ &= \frac{1}{2m} \sum_{i=1}^m \left(\left(\sum_{j=0}^n \theta_j x_j^{(i)} \right) - y^{(i)} \right)^2\end{aligned}$$

Gradient descent

Repeat θ_j

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

Partial Derivative Demonstration

3

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

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

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

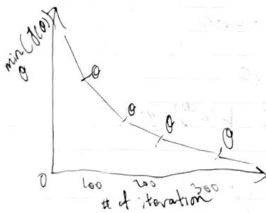
$$\theta_j := \theta_j - \alpha \cdot \left[\frac{1}{m} \sum_{i=1}^m ((\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n) - y^{(i)}) \cdot x_j^{(i)} \right]$$

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



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

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



MV

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

$$\text{pinv}(X^T * X) * X^T * y$$

✓ bruh, this shit trash.

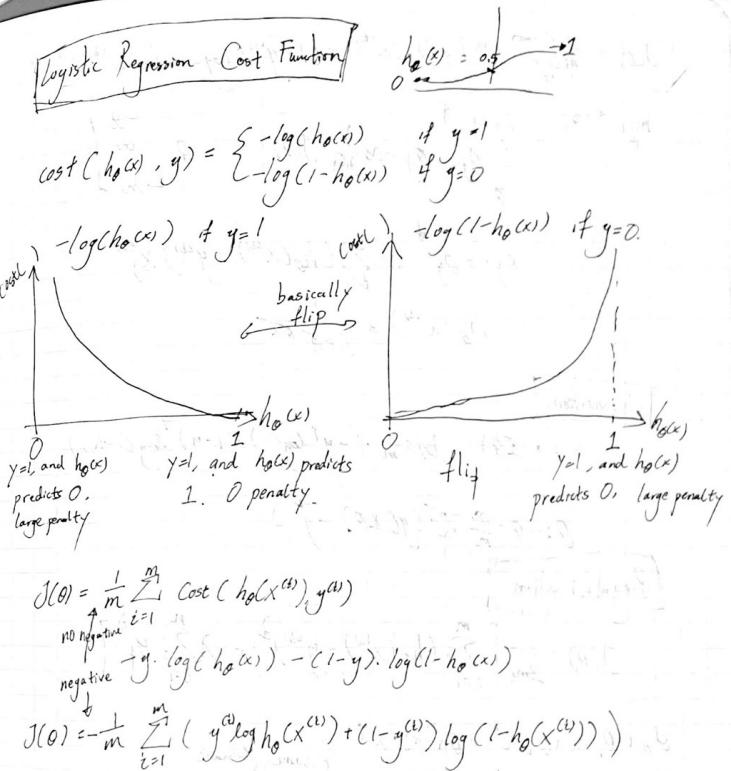
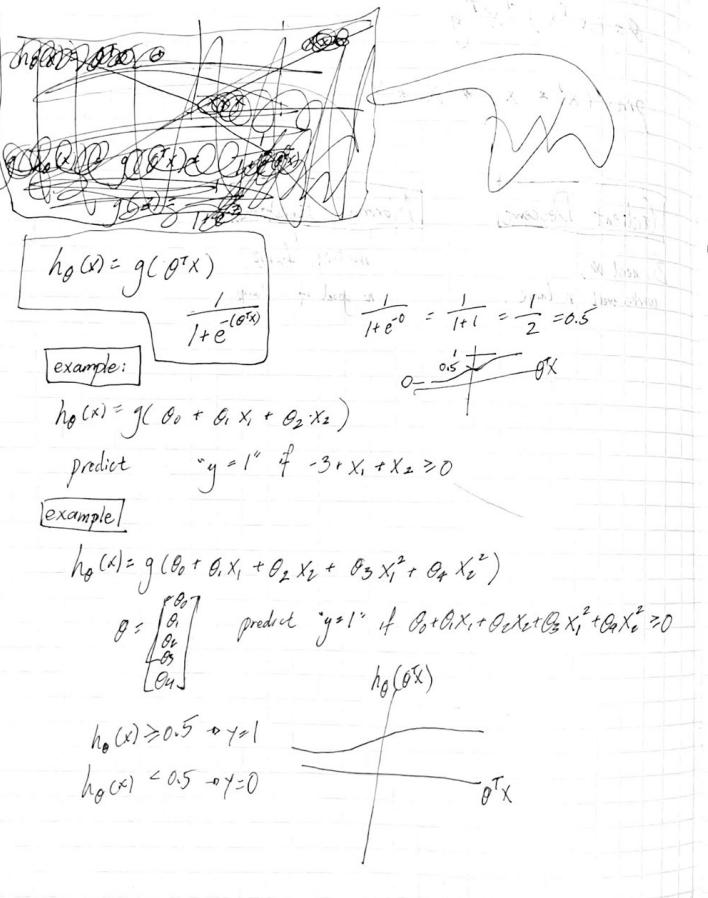
Gradient Descent

θ , need n ,
works well if n large.

Normal Equation

matrix, find θ .
no good if n large.

Week 3 Notes



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

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

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

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

Vectorization

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

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

Regularization

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

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \quad \left. \begin{array}{l} \text{same format} \\ \text{if no intercept} \end{array} \right\}$$

repeat ξ

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

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

↓

can be rewritten

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

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

$m \leq n$
examples # features

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

if $\lambda > 0$

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

use this
if no intercept

Logistic regression without regularization

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

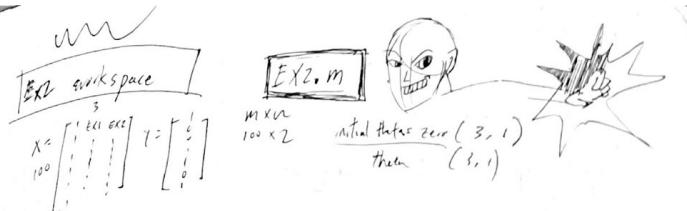
With regularization

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

derivative:

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

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



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

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

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

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

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

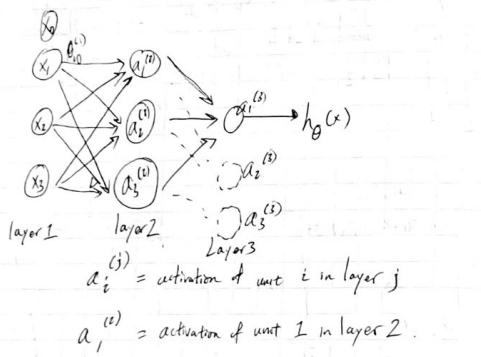
Octave code

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum [y \cdot \log(h) - (1-y) \cdot \log(1-h)] \\ \frac{\partial J(\theta)}{\partial \theta_j} &\rightarrow \left(\frac{1}{m} \right) \sum [(h - y) \cdot x_j] \end{aligned}$$

$$x \cdot y = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad x' \cdot y = 17$$

$$y = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad y' = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Week 4 Notes



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

$$a_1^{(2)} = g(\theta_{10}^{(1)} X_0 + \theta_{11}^{(1)} X_1 + \theta_{12}^{(1)} X_2 + \theta_{13}^{(1)} X_3) z_1^{(2)}$$

$$a_2^{(2)} = g(\theta_{20}^{(1)} X_0 + \theta_{21}^{(1)} X_1 + \theta_{22}^{(1)} X_2 + \theta_{23}^{(1)} X_3) z_2^{(2)}$$

$$a_3^{(2)} = g(\theta_{30}^{(1)} X_0 + \theta_{31}^{(1)} X_1 + \theta_{32}^{(1)} X_2 + \theta_{33}^{(1)} X_3) z_3^{(2)}$$

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

$$a_1^{(1)} = g(z_1^{(1)}) \quad z_1^{(1)} = \theta_{10}^{(0)} X_0 + \theta_{11}^{(0)} X_1 + \theta_{12}^{(0)} X_2 + \theta_{13}^{(0)} X_3$$

$$a_2^{(1)} = g(z_2^{(1)})$$

$$a_3^{(1)} = g(z_3^{(1)}) \quad z^{(1)} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \end{bmatrix} \quad a = g(z)$$

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

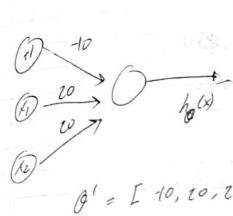
$$z_k^{(n)} = \theta_{k,0}^{(n)} X_0 + \theta_{k,1}^{(n)} X_1 + \dots + \theta_{k,n}^{(n)} X_n$$

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

$$X = \begin{bmatrix} X_0 \\ X_1 \\ \vdots \\ X_n \end{bmatrix} \quad Z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ \vdots \\ z_n^{(2)} \end{bmatrix}$$

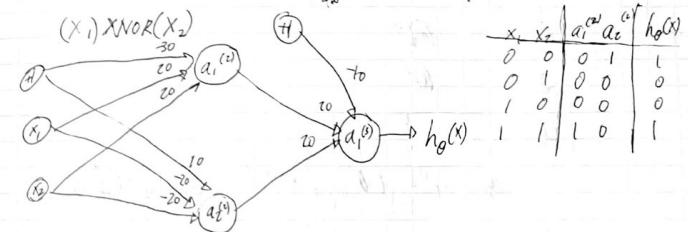
$$z_1^{(2)} = \theta_{10}^{(1)} a_0^{(1)} + \theta_{11}^{(1)} a_1^{(1)} + \theta_{12}^{(1)} a_2^{(1)} + \theta_{13}^{(1)} a_3^{(1)}$$

$$Z^{(2)} = \theta^{(j-1)} a^{(j-1)}$$



$x_1, x_2 \in \{0, 1\}$	$g(-b + \theta_0 x_1 + \theta_1 x_2)$
0 0	$g(-20) \rightarrow 0$
0 1	$g(10) \rightarrow 1$
1 0	$g(10) \rightarrow 1$
1 1	$g(30) \rightarrow 1$

$$(x_1) XNOR (x_2)$$



EX3 Workspace

overall I , x, y , num_labels, lambda
 images data whether it is 10 or 0.1
 5000×400 5000×1 $M=5000$
 each image 20×20 $n=400$
 grayscale pixel 10×400
 $X = 5000 \times 400$

①

10 classes for 10 numbers

②

$\theta^{(0)}$

③

$\theta^{(1)}$

④

$\theta^{(2)}$

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

classifier

lrCostFunction(theta, x, y, lambda)

10 classifier for each number
 different

lrCostFunction(theta, x, y, lambda)

@ t $t \in X$ $y=c$ lambda.

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

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

for c = 1: num_labels = for c = 1: 10. all_theta = zeros(10, 401)

all_theta(c, :) = ... this means continue.

fmincg C

401 by
 for each pixel
 intensity

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

for i = 1: size(X, 1)
 $X(PePe, :)$ all_theta 10 x 401

Predict One vs All (all-theta, X)

$$X = 5000 \times 400$$

↓

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

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

$$X = 5000 \times 400$$

$$m = 5000$$

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

$$P = 5000 \times 1$$

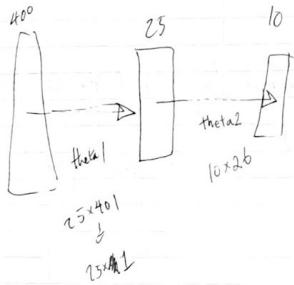
$$\theta_{01} = 25 \times 401$$

$$\theta_{02} = 10 \times 26$$

$$X' = \text{all-theta}'$$

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

gives you a 1×10 matrix of probabilities for the first image 20×20 .



$$25 \begin{bmatrix} \text{sig}(\cdot) \end{bmatrix}$$

$$\begin{bmatrix} 5000 \\ X \end{bmatrix} \xrightarrow{400 \times 1} 25 \begin{bmatrix} \text{sig}(\text{big penis}, :) \end{bmatrix}$$

$$\text{theta2} \rightarrow 25 \times 401$$

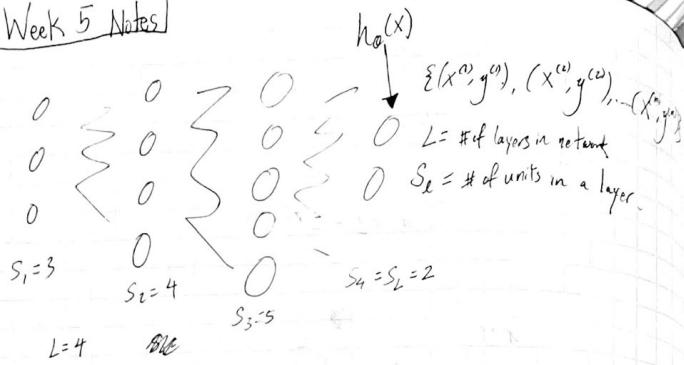
$$\text{sig}\left(25 \begin{bmatrix} \text{big penis} \end{bmatrix}\right) \Rightarrow 25 \begin{bmatrix} 1 \end{bmatrix}$$

$$\text{sigmoid}\left[25 \times 401 \times \begin{bmatrix} \text{big penis}, : \end{bmatrix}^T\right]$$

$$\begin{array}{c} \text{text} \\ \text{big penis} \end{array} \begin{bmatrix} 25 \\ 1 \end{bmatrix} \quad \begin{array}{c} \text{theta2} \\ \text{big penis} \end{array} \begin{bmatrix} 10 \\ 26 \end{bmatrix}$$

$$(\text{Theta2} * \text{big penis}) = 10 \times 26 \begin{bmatrix} 1 \end{bmatrix}$$

Week 5 Notes



Binary Classification

$$y = 0 \text{ or } 1$$

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

Multi-class classification (k classes)

$$y \in \mathbb{R}^k \quad \left[\begin{matrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{matrix} \right] \quad \left[\begin{matrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \\ 0 \end{matrix} \right] \quad \left[\begin{matrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \\ 0 \end{matrix} \right] \quad \dots \quad \left[\begin{matrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{matrix} \right]$$

$$S_L = k \quad k \geq 3$$

Logistic Regression

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

Neural Network

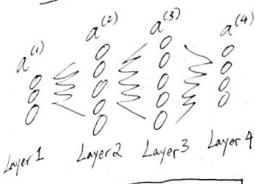
$$h_\theta(x) \in \mathbb{R}^k \quad (h_\theta(x))_i = i^{\text{th}} \text{ output}$$

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

Gradient computation

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

Given one training set (x, y)
Forward propagation



Backward Propagation

$\delta_j^{(l)}$ = the error of a j^{th} node

$$\delta_j^{(4)} = a_j^{(4)} - y_j \quad \text{rectified} \rightarrow \delta_j^{(4)} = a_j^{(4)} - y_j$$

$$\delta_j^{(3)} = (\theta^{(3)})^T \delta^{(4)} * \frac{\partial g'(z^{(3)})}{\partial z^{(3)}} \quad z^{(3)} = \begin{bmatrix} z_1^{(3)} \\ z_2^{(3)} \\ \vdots \\ z_n^{(3)} \end{bmatrix} \rightarrow \theta_{13}^{(2)} a_{13}^{(2)} + \theta_{23}^{(2)} a_{23}^{(2)}$$

$$\delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} * \frac{\partial g'(z^{(2)})}{\partial z^{(2)}} \quad a^{(2)} * (1-a^{(2)})$$

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

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

for $i = 1 \text{ to } m$

set $a^{(0)} = x^{(i)}$

forward propagation to compute $a^{(l)}$ $l = 2, 3, \dots, L$

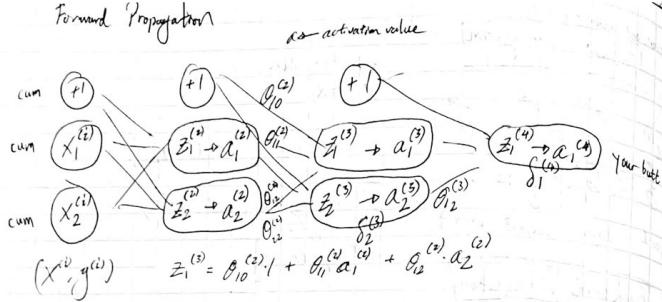
compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(1)}$ \times actual answer

then $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(1)}$ \leftarrow going backwards.

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

Two cases $\Delta_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)}$ if $i \neq 0$ $\left\{ \begin{array}{l} \text{complex} \\ \text{math proof shows that} \end{array} \right.$

$$\Delta_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \quad \text{if } j = 0 \quad \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = \Delta_{ij}^{(l)}$$



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

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

$$\delta_j^{(l)} = g^{(l)} - a_j^{(l)} \quad \delta_2 = \theta_{12}^{(3)} \cdot \delta_1^{(2)} \quad \delta_2 = \theta_{12}^{(2)} \delta_1^{(2)} + \theta_{22}^{(2)} \cdot \delta_2^{(2)}$$

whatever output you choose

$$\delta_2^{(2)} = \delta_2^{(3)} + \theta_{22}^{(2)} \cdot \delta_2^{(3)}$$

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

(θ) parameter (D) gradients

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

$$\theta^{(0)} \in \mathbb{R}^{10 \times 11} \quad \theta^{(1)} \in \mathbb{R}^{10 \times 11} \quad \theta^{(2)} \in \mathbb{R}^{10 \times 11}$$

$$S_2 + 10 \text{ nodes } S_3 + 1 \text{ node } + 1$$

$$D^{(0)} \in \mathbb{R}^{10 \times 11} \quad D^{(1)} \in \mathbb{R}^{10 \times 11} \quad D^{(2)} \in \mathbb{R}^{10 \times 11}$$

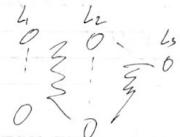
thetaVec = [Theta1(:, :); Theta2(:, :); Theta3(:, :)];

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

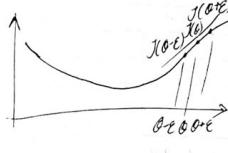
Theta1 = reshape(thetaVec(1:10, :), 10, 11);

Theta2 = reshape(thetaVec(11:20, :), 10, 11); → reshape back to original

Theta3 = reshape(thetaVec(21:23, :), 1, 11);



Gradient could have a bug. So use gradient checking



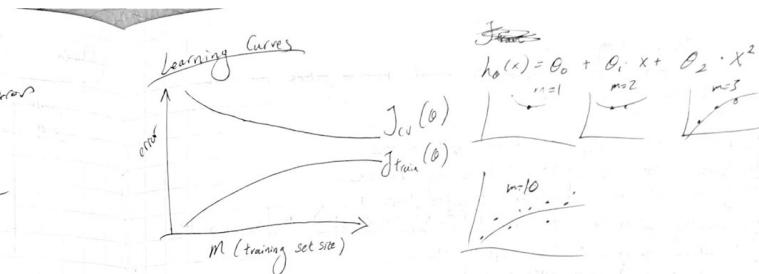
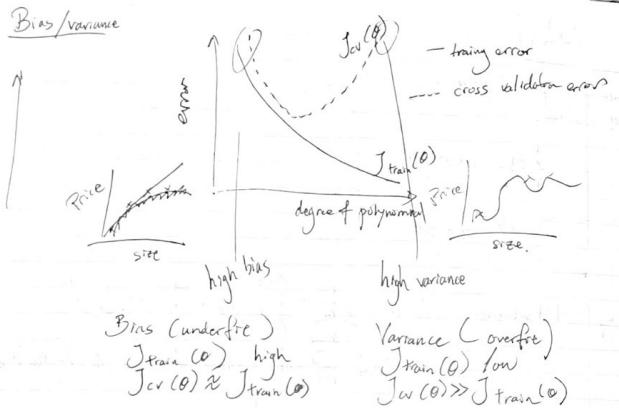
$$\frac{d}{d\theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

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

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

Putting it together

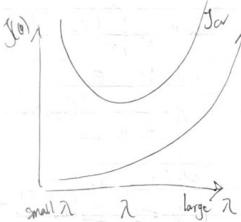
$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad g = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad g = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$



Linear Regression with Regularization

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

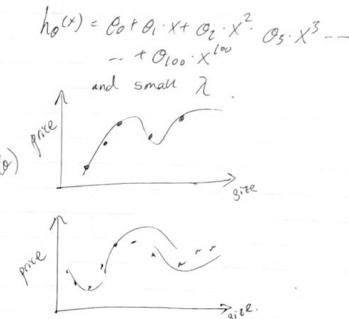
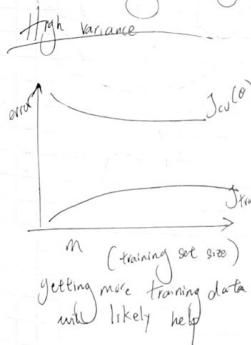
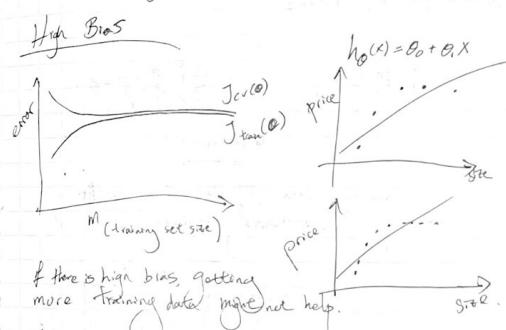
λ big λ just right λ small



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

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

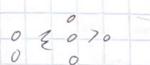
$$J_{cv}(\theta) = \frac{1}{M_{cv}} \sum_{i=1}^{M_{cv}} (h_{\theta}(x_i^{(1)}) - y_i^{(1)})^2$$



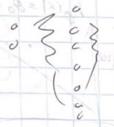
Debugging learning algorithm

Implement regularized linear regression to predict housing prices on new dataset. Not working, now what?

- Get more training examples \rightarrow fixes high variance
- Try smaller sets of features \rightarrow fixes high variance
- Try getting additional features \rightarrow fixes high bias
- Try adding polynomial features (X_1^2, X_2^2, X_1X_2 , etc) \rightarrow fixes high bias
- Try decreasing λ \rightarrow fixes high bias
- Try increasing λ \rightarrow fixes high variance



Small neural network
prone to underfit,
computationally cheap



large neural network
use regularization (λ) to address overfitting
computationally expensive



Ex 05 workspace

(12, 1) (12, 1) (2, 1) (1, 1)

Linear Reg Cost Function (X, Y, θ, λ)

Cost ($m, 1$) X_j

C_{ij}

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

theta = 2×1

$12 \times 2 \cdot 2 \times 1 = 12 \times 1$

$(\frac{1}{m} \sum ((X \cdot \theta) - y)^2) + (\frac{\lambda}{2 \cdot m}) (\text{sum})$

grad₁ $\begin{bmatrix} (1) \\ (2) \\ \vdots \\ (n) \end{bmatrix} \rightarrow \text{grad } (1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) X_j^{(i)}$

grad₂ $\begin{bmatrix} (1) \\ (2) \\ \vdots \\ (n) \end{bmatrix} \rightarrow \text{grad } (2) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) X_j^{(i)} + \frac{\lambda}{m} \theta_j$

$12 \times 2 \cdot 2 \times 1 = 12 \times 1 \quad X(t) \cdot \theta(t)$

$(X \cdot \theta - y) \cdot (X \cdot \theta - y)^T$

$= 12 \times 1$

2-1 learning curves
task: plot learning curve, training set and cross validation set
for different training set.

learning curve ($X, y, X_{val}, Y_{val}, \lambda$) \rightarrow error_train
 \downarrow
 $[ones, X] \quad y \quad [ones, X_{val}] \quad Y_{val} \quad 0$
 \downarrow why doesn't X_{val} contain x ?
 \downarrow
 X_{val}

$X: 3 \times 1$
 $X_{-val}: 3 \times 8$
 $Y: 3 \times 1$

$X \quad X_{test} \quad X_{val}$
 $y \quad y_{test} \quad y_{val}$
 $1 \times 2 \quad 2 \times 1 \quad 2 \times 1$
 $\uparrow \quad \uparrow \quad \uparrow$
 $Mod \quad train \quad cross \quad validation$

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

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

use train linearly

Training set used to train parameters
cross validation set used to fine tune and pick equation

Test to test final equation

X_{val}	2×2	$x(1:1, :)$	1×2	1×2
y_{val}	2×1	$y(1:1)$	1×1	1×1
x	12×2	$x(1:2, :)$	2×2	2×3
y	12×1	$y(1:2)$	2×1	2×1

$$\text{Precision} = \frac{\text{true positive}}{\text{true positive} + \text{false positive}}$$

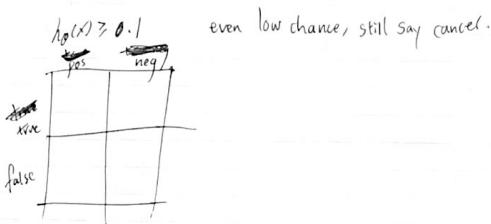
$$\text{Recall} = \frac{\text{true positive}}{\text{true positive} + \text{false negative}}$$

high precision and high recall = good algorithm.
means true positive is high.

evaluation for classification with skewed classes

how to choose threshold? predict 1 if $h_{\theta}(x) \geq \text{threshold}$

$$F_1 \text{ score: } 2 \cdot \frac{P \cdot R}{P+R}$$



$$(a + \frac{1}{m} b)$$

$$\frac{1}{m} a + \frac{1}{m} b$$

$$\sum_i y^{(i)}$$

$$x$$

$$\lim_{x \rightarrow 0} \lambda$$

Week 7 Notes

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

θ k x matrices

if $y=1$, we want $h_{\theta}(x) \approx 1$, $\theta^T x \gg 0$

if $y=0$, we want $h_{\theta}(x) \approx 0$, $\theta^T x \ll 0$

cost function: $-\left[y \log(h_{\theta}(x)) + (1-y) \log(1-h_{\theta}(x))\right]$

$$\text{cost}_1(\theta^T x^{(i)})$$

$$-y \log(h_{\theta}(x^{(i)}))$$

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

Support vector machine:

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

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

$$\min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

Find the θ that will give you the min

$$\text{hypothesis: } \begin{cases} 1 & \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$\theta^T x$ is your function

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

makes it between 0 or 1

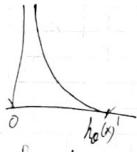
$h_{\theta}(x) > 0.5$ then predicts $y=1$
 $h_{\theta}(x) < 0.5$ then predicts $y=0$

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

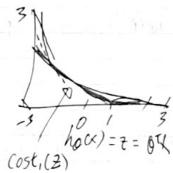
cost function for prediction.

$\text{arg} \min \ln.$

$$-\log(h_{\theta}(y)) \quad -\log\left(\frac{1}{1+e^{-\theta^T x}}\right)$$



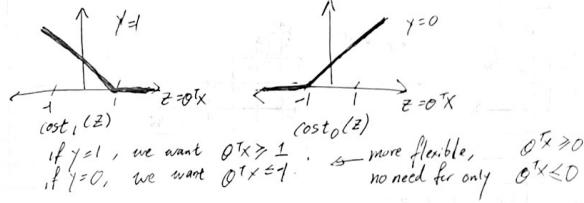
for $y=1$,
if $h_{\theta}(x)$ predicts 1,
receives 0 penalty



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

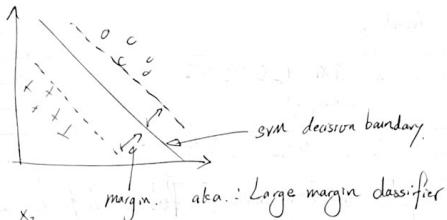
$$\min_C \sum_{i=1}^n [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^d \theta_j^2$$

find the θ for θ^T that will give the lowest sum cost function with regularization.

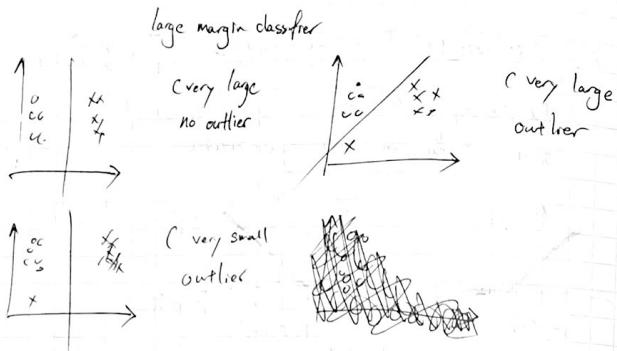


if $y=1$, we want $\theta^T x \geq 0$. more flexible, $\theta^T x \geq 0$
if $y=0$, we want $\theta^T x \leq 0$. no need for only $\theta^T x \leq 0$

set C to very large, then if we make prediction mistake, then function will be penalized greatly.

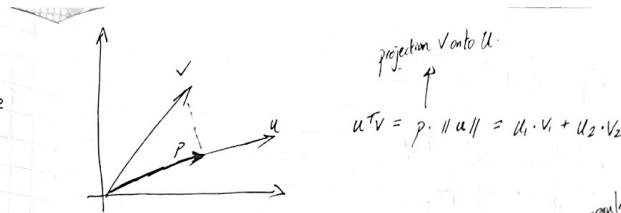


$$\begin{aligned} &\theta_3 + \theta_2 x_2 + \theta_1 x_1 \geq 0 \\ &\theta_3 = -3 \\ &\theta_2 = 0 \\ &\theta_1 = 1 \end{aligned}$$



$$u^T v = p \cdot \|u\|$$

projection of v onto u



Regularization cost function

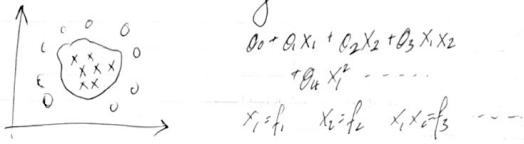
$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^n \theta_j^2 \right]$$

let say features: x_1, x_2, \dots, x_{100}
Parameters: $\theta_0, \theta_1, \dots, \theta_{100}$

regularization parameter
 λ
+ $\frac{\lambda}{2} \sum_{j=1}^n \theta_j^2$
Shrinks all parameters
 θ_1 to θ_{100}
skip θ_0

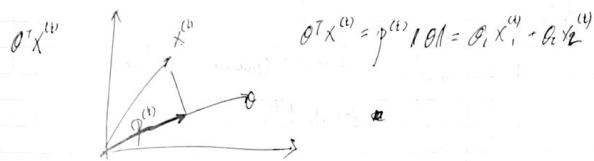
keeping $h_\theta(x)$ the hypothesis simple to prevent overfitting.

Non linear decision boundary



is there other/better choices for features f_1, f_2, f_3, \dots ?

$$\begin{aligned} & \min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 \quad \text{let say } n=2 \\ & \theta^T \theta = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} (\sqrt{\theta_1^2 + \theta_2^2})^2 = \frac{1}{2} \|\theta\|^2 \\ & \text{s.t. } \theta^T x^{(1)} \geq 1 \quad \text{if } y^{(1)} = 1 \\ & \quad \theta^T x^{(2)} \leq -1 \quad \text{if } y^{(2)} = 0 \end{aligned}$$



$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} \|\theta\|^2$$

s.t. $\theta^{(1)} \cdot \|\theta\| \geq 1$ if $\gamma^{(1)} = 1$
 $\theta^{(2)} \cdot \|\theta\| \leq -1$ if $\gamma^{(2)} = 1$

$$f_i = \text{similarity}(x, \ell^{(i)}) = \exp\left(\frac{-\|x - \ell^{(i)}\|^2}{2\sigma^2}\right) = \exp\left(\frac{-\sum_{j=1}^n (x_j - \ell_j^{(i)})^2}{2\sigma^2}\right)$$

if $x \approx \ell^{(i)}$
 $f_i \approx \exp\left(\frac{-\sigma^2}{2\sigma^2}\right) \approx 1$

if x is far from $\ell^{(i)}$

$$f_i \approx \exp\left(\frac{-(\text{large number})^2}{2\sigma^2}\right) \approx 0.$$

Non-linear Decision Boundary

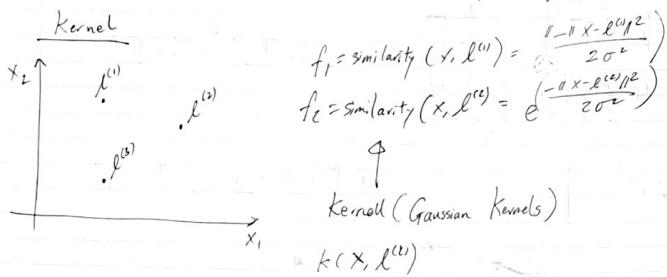


Predict $y=1$ if

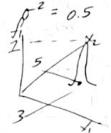
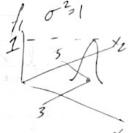
$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 + \dots \geq 0.$$

$$f_1 = x_1, f_2 = x_2, f_3 = x_1 x_2, f_4 = x_1^2, f_5 = x_2^2$$

better / different choice for f_1, f_2, f_3, \dots ?

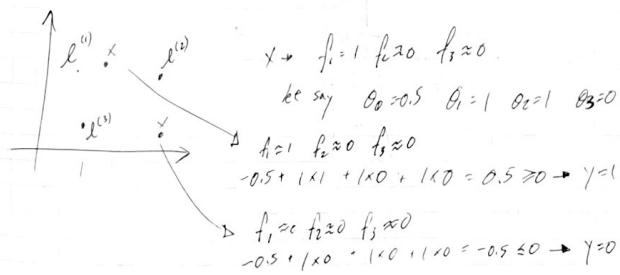


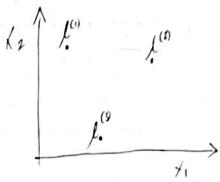
$$\ell^{(w)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix} \quad f_i = e^{\frac{-\|x - \ell^{(i)}\|^2}{2\sigma^2}}$$



f_i falls faster when moving away from $\ell^{(i)}$

f_i falls slower when moving away from $\ell^{(i)}$





Given X :

$$f_i = \text{similarity}(x, \ell^{(i)}) = e^{\left(\frac{-\|x - \ell^{(i)}\|^2}{2\sigma^2}\right)}$$

how to pick $\ell^{(i)}$?



Given $(x^{(1)}, y^{(1)})$, $(x^{(2)}, y^{(2)})$, ..., $(x^{(m)}, y^{(m)})$
choose $\ell^{(1)} = x^{(1)}$, $\ell^{(2)} = x^{(2)}$, ..., $\ell^{(m)} = x^{(m)}$

$$f = \begin{bmatrix} f_0 \\ f_1 \\ \vdots \\ f_m \end{bmatrix} \quad f_0 = 1 \quad X^{(i)} \rightarrow f_1^{(i)} = \text{sim}(x^{(i)}, \ell^{(1)}) \quad x^{(i)} \& \ell^{(1)} \text{ will be the same} \\ f_2^{(i)} = \text{sim}(x^{(i)}, \ell^{(2)}) \\ \vdots \\ f_m^{(i)} = \text{sim}(x^{(i)}, \ell^{(m)}) \quad f_i = e^{\left(\frac{-\|x^{(i)} - \ell^{(i)}\|^2}{2\sigma^2}\right)}$$

$$f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix}$$

Training:

$$\min_{\theta} C \sum_{t=1}^m f^{(t)} (\text{cost}_1(\theta^T f^{(t)}) + (1-y^{(t)}) \text{cost}_0(\theta^T f^{(t)}) + \frac{1}{2} \sum_{j=1}^m \theta_j^2)$$

Hypothesis: Given x , compute features $f \in \mathbb{R}^{m+1}$
Predict " $y=1$ " if $\theta^T f > 0$ $\theta \in \mathbb{R}^{m+1}$

$$\theta_0 + \theta_1 f_1 + \dots + \theta_m f_m$$

Linear

$C \in \mathbb{R}$) large C : low bias, high variance (small λ)
lower C : high bias, low variance (big λ)



If SVM overfits, decrease C , increase σ^2

Linear kernel vs Gaussian kernel

$$\text{function } f = \text{kernel}(x_i, x_o)$$

not all similarities are valid kernels, need to satisfy Mercer's theorem.

Logistic Regression vs SVMs

$n = \# \text{ of features}$ $m = \# \text{ of training sets}$

If n large relative to m , use logistic regression or svm without kernel. (linear)

If n small, m intermediate ($1-1000$), ($10-10000$), use SVM with gaussian kernel.

If n small, mistakes, ($1-1000$), ($10-10000$), ($m=50,000$), use logistic/SVM without kernel.

- Most cases, neural network better, but slower.

Week 8 Notes



K means algorithm → find lowest distance for cluster centroid

Input:

K (# of clusters)

Training set $\{X^{(1)}, X^{(2)}, \dots, X^{(m)}\}$ $X^{(i)} \in \mathbb{R}^n$ drop $x_0 = 1$

Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

Repeat

 for $i = 1$ to m

$c^{(i)}$ $i = 1$ to K min $\|x^{(i)} - \mu_k\|$ $c^{(3)} = 5$
 closest to $x^{(i)}$ k closest μ_5

 for $k = 1$ to K

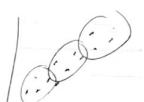
$\mu_k :=$ average of $x^{(i)}$ assigned to μ_k

 example $\mu_2 = \frac{1}{4} [x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)}]$

$k=3$



$k=3$



$c^{(i)}$ = index of cluster assigned to $x^{(i)}$
 μ_k = cluster centroid \leftarrow it is a location
 $\mu_{c^{(i)}}$ = cluster assigned to $x^{(i)}$
 $x^{(1)} \rightarrow 5 \quad c^{(1)} = 5 \quad \mu_{c^{(1)}} = \mu_5$

Optimization Objective:

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

$$\min_{c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

How to randomly initialize?

local optima issues

pick random $x^{(i)}$'s
 $k < m$



choosing k .

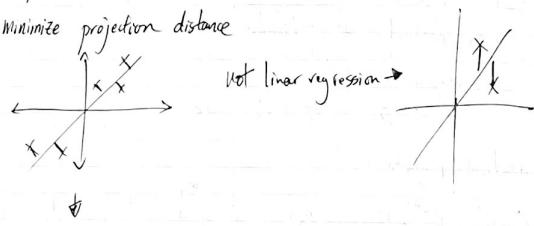


Dimension Compression

3D → 2D graph

PCA Principal Component Analysis

Minimize projection distance



not linear regression \rightarrow

$$x^{(i)} \in \mathbb{R}^2 \rightarrow x^{(i)} \in \mathbb{R}$$

Reduce data from n to k dimensions

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)})(x^{(i)})^T \quad [U, S, V] = \text{svd}(\Sigma)$$

U is also $n \times n$.

$$U = \begin{bmatrix} u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(n)} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$$Z = \begin{bmatrix} u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(k)} \end{bmatrix}^T \cdot X$$

U_{reduce}

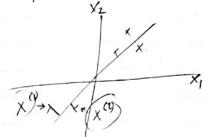
$$X = \begin{bmatrix} x^{(1)^T} \\ \vdots \\ x^{(n)^T} \end{bmatrix}$$

$$\text{SVD} = (\frac{1}{m}) \cdot X^T \cdot X$$

$$U_{\text{reduce}} = U(:, 1:k)$$

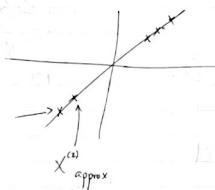
$$Z = U_{\text{reduce}}^T \cdot X$$

Compression Reconstruction



$$Z = U_{\text{reduce}}^T \cdot X$$

$$X_{\text{approx}} = U_{\text{reduce}} \cdot Z$$



K : # of principal components

average squared projection error: $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - X_{\text{approx}}^{(i)}\|^2$
total variation in the data: $\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - X_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01 \quad \Rightarrow \quad 99\% \text{ of variance is retained}$$



choose $k \rightarrow$ and repeat until
 Try PCA $k=1$ if fail, $k=2$, if fail, $k=3$...
 compute U and Σ , $Z^{(1)}, \dots, Z^{(m)}$, $X_{approx}^{(1)}, \dots, X_{approx}^{(m)}$
 check if:

$$\frac{1}{m} \sum_{i=1}^m \|X^{(i)} - X_{approx}^{(i)}\|^2 \leq 0.01?$$

$$\frac{1}{m} \sum_{i=1}^m \|X^{(i)}\|^2$$

not efficient

$$[U, \Sigma, V] = svd(\text{Sigma})$$

$$\begin{bmatrix} S_1 \\ S_2 \\ S_3 \\ \vdots \\ S_n \end{bmatrix} \quad \left| \quad 1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \leq 0.01 \right.$$

$n \times n$

Application of PCA

$$X^{(1)}, X^{(2)}, \dots, X^{(m)} \in \mathbb{R}^{10,000} \quad \downarrow \text{PCA} \quad \left[\begin{array}{c|c} & 100 \\ \hline & 100 \end{array} \right]$$

$$Z^{(1)}, Z^{(2)}, \dots, Z^{(m)} \in \mathbb{R}^{1000} \quad h_\theta(z) = \frac{1}{1 + e^{-\theta^T z}}$$

use PCA for compression and visualization

But don't use PCA to prevent overfitting!

when designing ML system, try implementing it first before applying PCA. Only use PCA when you have a really good reason.

[ex 7 workspace]

$\text{idx} = \text{findClosestCentroids}(X, initial_centroids);$

$(300, 1)$

$range(1:k)$

$\text{range}(1:m)$

$\text{initial_centroids}(j, :)$

$\text{centroids} = \text{computeCentroids}(X, idx, K)$

$3[1]$

$3[2]$

$3[3]$

$3[4]$

$$\mu_k = \frac{1}{|C_k|} \sum_{i \in C_k} X^{(i)}$$

$2^8 = 256$

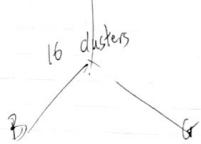
24 bit color representation = three 8 bits

$0-255$ red

$0-255$ green

$0-255$ blue

Thousands of combinations compress 16 combo of colors



original image (128) $128^2 = 16384$

$$\Sigma = \frac{1}{m} X^T X \quad \text{covariance matrix}$$

Sigma
n by n matrix

$$Z \quad U$$

$50 \times 1 \quad 2 \times 2$

$$X_{\text{recover}} = Z \cdot U$$

Week 9 Notes

Anomaly Detection

Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

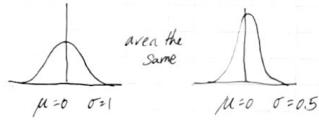
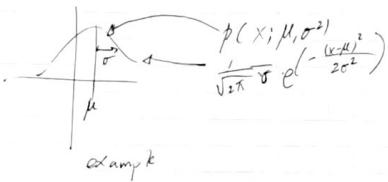
Model $p(x)$

$\begin{cases} p(x) < \epsilon & \text{flag for anomaly.} \\ p(x) \geq \epsilon & \text{flag OK.} \end{cases}$



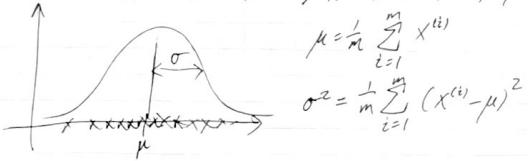
Gaussian (Normal) Distribution

$$X \sim N(\mu, \sigma^2)$$



Parameter Estimation

Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \subset \mathbb{R}^n \quad X^{(i)} \sim N(\mu, \sigma^2)$



Density Estimation

training set: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \subset \mathbb{R}^n$

x_i \leftarrow raw feature

$$\begin{aligned} p(x) &= p(x_1; \mu_1, \sigma_1^2) \cdot p(x_2; \mu_2, \sigma_2^2) \cdots p(x_n; \mu_n, \sigma_n^2) \\ &\downarrow \quad \downarrow \quad \downarrow \\ &x_1 \sim N(\mu_1, \sigma_1^2) \quad x_2 \sim N(\mu_2, \sigma_2^2) \\ &= \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) \end{aligned}$$

Anomaly Detection Algorithm

(1) choose features x_i different features

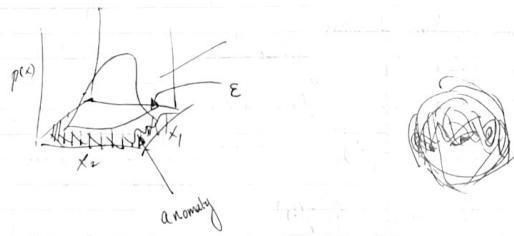
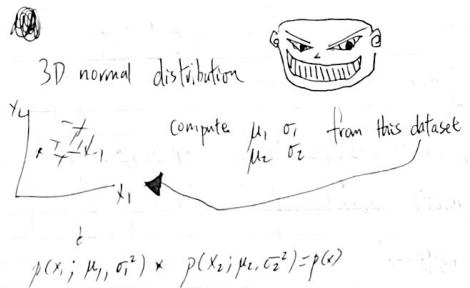
(2) Fit parameter $\{\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2\}$

$$\begin{aligned} \mu_j &= \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \longrightarrow \mu = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_n \end{bmatrix} = \frac{1}{m} \sum_{i=1}^m X^{(i)} \\ \sigma_j^2 &= \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2 \end{aligned}$$

(3) Given new example x , compute $p(x)$:

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}}$$

Anomaly if $p(x) < \epsilon$



Aircraft engine

Developing and evaluating anomaly detection systems.

10,000 good engines

20 flawed engines

training set: 6000 good engine.

CV set: 2000 good engine ($y=0$), 10 anomalous ($y=1$)

Test set: 2000 good engine ($y=0$), 10 anomalous ($y=1$)

repeated bad data practice	600	10
	4000	10
	4000	10

Algorithm Evaluation

fit model $p(x)$ on training set $\{x^{(0)}, \dots, x^{(m)}\}$
on cross validation / test example x ,

predict:

$$g \begin{cases} 1 & \text{if } p(x) \leq E \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq E \text{ (normal)} \end{cases}$$

classification no good, because dataset is skewed because $y=0$ is more common.

Anomaly Detection vs Supervised learning

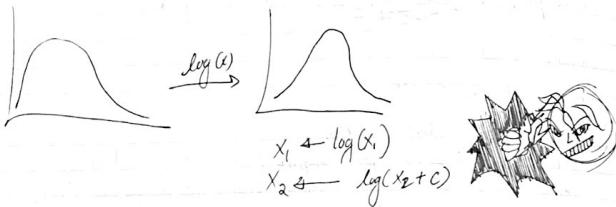
When to use which?

Small number (0-20) of positive example ($y=1$)

Large number of negative number.

Large number of positive and negative number.

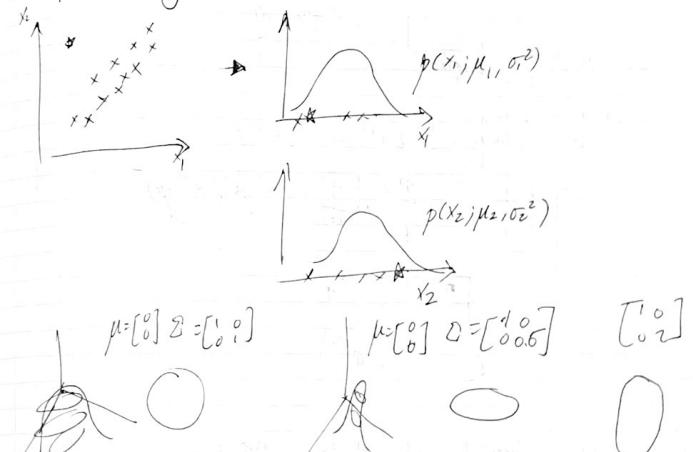
choosing what features to use?



coming up with more/new features

Multivariate Gaussian Distribution

example: monitoring machines in a data center.



$$\mu = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}$$
$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

Multivariate Gaussian Distribution anomaly detection

Parameters $\mu \in \mathbb{R}^n$, $\Sigma \in \mathbb{R}^{nn}$

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

Parameter fitting:
Given training set $\{x^{(1)}, \dots, x^{(m)}\}$

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$



use multivariate if $m > n$, $m \gg n$.

Problem Formulation

predicting movie rating.

N_u = # of users N_m = # of movies

$r(t, j) = 1$ if user j has rated movie i

$y(t, j)$ rating user j gave to movie i .

$\theta^{(j)}$ = parameter vector for user j

$x^{(i)}$ = feature vectors for movie i

rating prediction = $(\theta^{(j)})^T x^{(i)}$

$m^{(j)}$ = # of movies rated by user j

To learn $\theta^{(j)}$:

$$\min_{\theta^{(j)}} \frac{1}{2m^{(j)}} \sum_{i: r(i,j)=1} [(\theta^{(j)})^T x^{(i)} - y^{(i,j)}]^2 + \frac{\lambda}{2} \sum_{k=1}^n (\theta_k^{(j)})^2$$

only for single user j .

for all users: to learn $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(N_u)}$

$$\min_{\theta^{(1)}, \dots, \theta^{(N_u)}} \frac{1}{2} \sum_{j=1}^{N_u} \sum_{i: r(i,j)=1} [(\theta^{(j)})^T x^{(i)} - y^{(i,j)}]^2 + \frac{\lambda}{2} \sum_{j=1}^{N_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

Gradient Descent:

$$\theta_K^{(j)} := \theta_K^{(j)} - \alpha \sum_{i: r(i,j)=1} [(\theta^{(j)})^T x^{(i)} - y^{(i,j)}] \cdot x_k^{(i)} \quad k=0$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left(\sum_{i: r(i,j)=1} [(\theta^{(j)})^T x^{(i)} - y^{(i,j)}] x_k^{(i)} + \lambda \theta_k^{(j)} \right) \quad k \neq 0$$

Collaborative Filtering

Given $\theta^{(1)}, \dots, \theta^{(n)}$ to learn $X^{(i)}$:

$$\min_{X^{(i)}} \frac{1}{2} \sum_{j:r(i,j)=1} [(\theta^{(j)})^T X^{(i)} - y^{(i,j)}]^2 + \frac{\lambda}{2} \sum_{k=1}^n (X_k^{(i)})^2$$

Given $\theta^{(1)}, \dots, \theta^{(n)}$ to learn $X^{(1)}, \dots, X^{(n)}$:

$$\min_{X^{(1)}, \dots, X^{(n)}} \frac{1}{2} \sum_{t=1}^m \sum_{j:r(i,j)=1} [(\theta^{(j)})^T X^{(t)} - y^{(t,j)}]^2 + \frac{\lambda}{2} \sum_{i=1}^n \sum_{k=1}^n (X_k^{(i)})^2$$

Movie	alice (1)	bob (2)	carol (3)	dave (4)	x_1	x_2
Cone at least	5	5	0	0	? 1	? 0
Romance forever	5	?	?	0	?	?
late paprika love	?	4	0	?	?	?
last step car chase	0	0	5	4	?	?
Swords vs karate	0	0	5	?	?	?

$$\theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix} \quad \theta^{(2)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix} \quad \theta^{(3)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix} \quad \theta^{(4)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}$$

$$(\theta^{(1)})^T \cdot X^{(1)} \approx 5 \quad (\theta^{(2)})^T \cdot X^{(1)} \approx 5$$

$$\begin{bmatrix} 0 & 5 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = 0 + 5 + 0 = 5$$

$$\min_{\substack{X^{(1)}, \dots, X^{(n)} \\ \theta^{(1)}, \dots, \theta^{(n)}}} J(X^{(1)}, \dots, X^{(n)}, \theta^{(1)}, \dots, \theta^{(n)}) = \frac{1}{2} \sum_{(i,j) \in r} [(\theta^{(j)})^T X^{(i)} - y^{(i,j)}]^2 + \frac{\lambda}{2} \sum_{i=1}^n \sum_{k=1}^n (X_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^m \sum_{k=1}^n (\theta_k^{(j)})^2$$

$$\begin{aligned} x_k^{(i)} &:= X_k^{(i)} - \alpha \left(\sum_{j:r(i,j)=1} ((\theta^{(j)})^T X^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda X_k^{(i)} \right) \frac{\partial}{\partial X_k^{(i)}} \\ \theta_k^{(j)} &:= \theta_k^{(j)} - \alpha \left(\sum_{i:r(i,j)=1} ((\theta^{(j)})^T X^{(i)} - y^{(i,j)}) X_k^{(i)} + \lambda \theta_k^{(j)} \right) \end{aligned}$$

Vectorization: Low Rank Matrix Factorization

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix} \quad \text{Predicted Rating:} \quad \begin{bmatrix} (\theta^{(1)})^T(X^{(1)}) & (\theta^{(2)})^T(X^{(1)}) & \dots & (\theta^{(n_u)})^T(X^{(1)}) \\ (\theta^{(1)})^T(X^{(2)}) & & & \\ \vdots & & & \\ (\theta^{(1)})^T(X^{(n_u)}) & (\theta^{(2)})^T(X^{(n_u)}) & & (\theta^{(n_u)})^T(X^{(n_u)}) \end{bmatrix}$$

$\theta \rightarrow$ what genre they like $X_i \rightarrow$ movie genre type

$$\rightarrow X \cdot \theta^T \quad \rightarrow \quad X = \begin{bmatrix} -(\theta^{(1)})^T \\ -(\theta^{(2)})^T \\ \vdots \\ -(\theta^{(n_u)})^T \end{bmatrix} \quad \theta = \begin{bmatrix} -(\theta^{(1)})^T \\ \vdots \\ -(\theta^{(n_u)})^T \end{bmatrix}$$

For each product i , we learn a feature vector $x^{(i)} \in \mathbb{R}^n$.
movie genre

Movie i related to movie j

If $\|x^{(i)} - x^{(j)}\|$ is small, movie i & j are similar.

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix} \rightarrow \mu = \begin{bmatrix} 2.5 \\ 0.5 \\ 2 \\ -0.5 \\ 1.25 \end{bmatrix} \quad b = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ -2.5 & -2.5 & 2.5 & 1.25 & ? \\ -1.25 & -1.25 & 3.75 & 7.25 & ? \end{bmatrix}$$

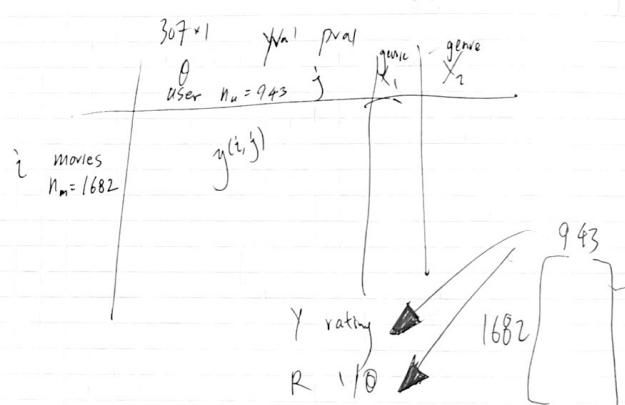
for user j , on movie t , predict
 $(\theta^{(t)})^T(X^{(t)}) + \mu_j$

EX 8 Work Space

$$X = \begin{bmatrix} 307 \\ 2 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\sigma^2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$



collaborative filtering predict rating for non-rated movies

$$X = \begin{bmatrix} 100 \\ 1684 \end{bmatrix} \quad \theta = \begin{bmatrix} 100 \\ 1 \end{bmatrix}$$

predicts rating for movie i by user j .

$$y^{(i,j)} = (\theta^T)^T X^T$$

$$X = 5 \begin{bmatrix} 3 \end{bmatrix} \quad \text{theta} = 4 \begin{bmatrix} 3 \end{bmatrix}$$

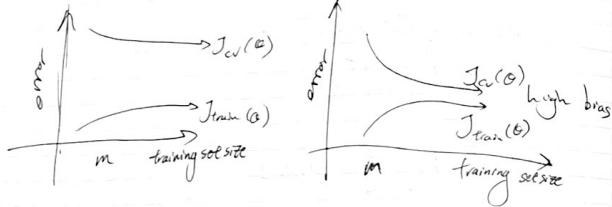
$$Y = 5 \begin{bmatrix} 4 \end{bmatrix} \quad X_{\text{grad}} = 5 \times 3$$

$$R = 5 \begin{bmatrix} 4 \end{bmatrix} \quad \text{Theta-grad} = 4 \times 3$$

$$\begin{matrix} 5 \times 4 & \text{theta} \\ 4 \times 3 & \rightarrow n_u \\ 5 \times 4 & \\ 4 \times 5 & 5 \times 3 \end{matrix}$$

Week 10 Notes

learning with large datasets
 $m = 100,000,000?$ $m = 1,000?$



Stochastic Gradient Descent

batch & it looks at all the examples
 $300,000,000$
 issue with normal gradient descent is if m is large,
 then computation is expensive.

Batch Gradient Descent

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

Repeat

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

$$\frac{\partial}{\partial \theta_j} J_{train}(\theta)$$

}

Stochastic Gradient Descent

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

$$J_{train}(\theta) = \frac{1}{m} \sum_{i=1}^m cost(\theta, (x^{(i)}, y^{(i)}))$$

Batch Gradient Descent $\rightarrow M$ examples each iteration
 Stochastic Gradient Descent $\rightarrow 1$ example each iteration
 Mini-Batch Gradient Descent $\rightarrow B$ examples each iteration

Stochastic Gradient Descent convergence & learning rate α .