

\* علم الیست ML

• توانایی های جدید برای کامپیوتر

• آنچه بزرگ شده است که از AI خارج شده است

\* مثالهای کاربردی

• Database Mining ← ① Web ② داده های بزرگ حاصل از آنالیز  
• برنامه هایی که به دست نمی توان نوشت ← هارت ، دانش همزه ، دانش محنتی در کاربرد

Self-customizing Programs

• در یادگیری در نظر

\* یادگیری ماشین چیست ؟

• آرنولد ساموئل : زمینه مطالعاتی که به کامپیوتر این قابلیت را می دهد که کاری را انجام دهد که صراحتاً برای آن برنامه ریزی نشده است

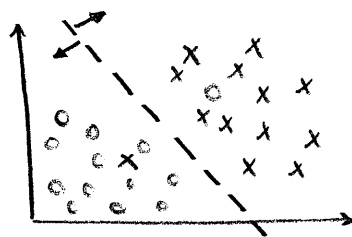
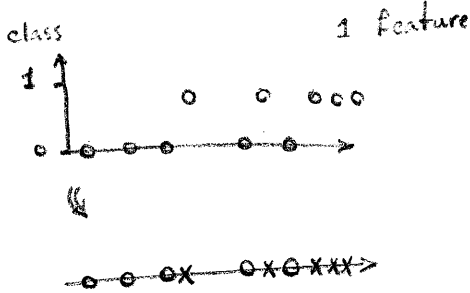
• نام اصلی : یادگیری به صورت یاد گرفتن توسط ماشین  
• Performance بر اساس معیار برای یادگیری حاصل  
P ← E ← T

\* یادگیری با نظار (Supervised learning)

• به مجموعه ای از داده ها اعمال می شود که جوابهای درست آن درجود است

Regression : پیش بینی خروجی های دارای مقادیر پیوسته

Classification : پیش بینی کلاس داده ها با خروجی گسسته



\* یادگیری بدون نظار (Unsupervised learning)

• در مجموعه داده ها به دنبال ساختار می گردیم

social network anal. , organizing large clusters , Genomics , news.google.com

astronomical data analysis , market segmentation

• ساختار داده ها را برای دسته بندی آنها استفاده می کنیم : clustering ← دسته ها مهم است

• ساختار داده ها را برای تاثیر آنها استفاده می کنیم : ICA , PCA ← داده های ورودی اصلی مهم است  
cocktail party problem

# Linear Regression

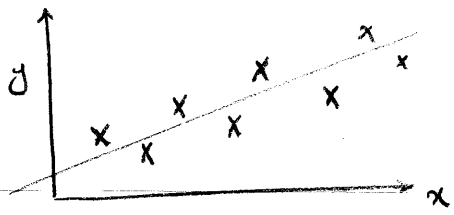
تقریب خط به داده ها، با داشتن جواب درست برای هر داده ورودی برای پیش بینی خروجی برای داده های جدید

$m = \#$  of inputs  
 $x =$  input vars/features  
 $y =$  output vars/features

$(x, y)$  : training example  
 $(x^{(i)}, y^{(i)})$  :  $i^{th}$  training example

Training Set  $\Rightarrow$  Learning Alg.  $\Rightarrow$  Hypothesis: mapping from  $x$  to  $y$

1D Feature:  $h_{\theta}(x) = \theta_0 + \theta_1 x$   $\rightarrow$  Univariate



$\theta$ : parameters  $\rightarrow$  goal: find  $\theta$ 's so  $h_{\theta}(x)$  is close to  $y$  for training examples  $(x, y)$

find me the  $\theta_0$  and  $\theta_1$  so that the average of error of my real data to fitted data is minimized.

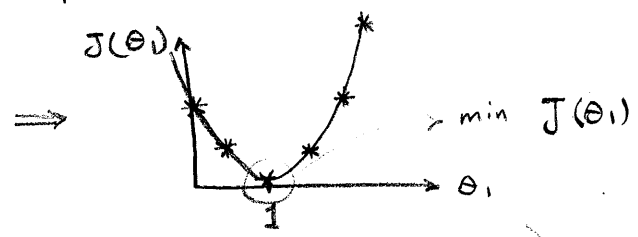
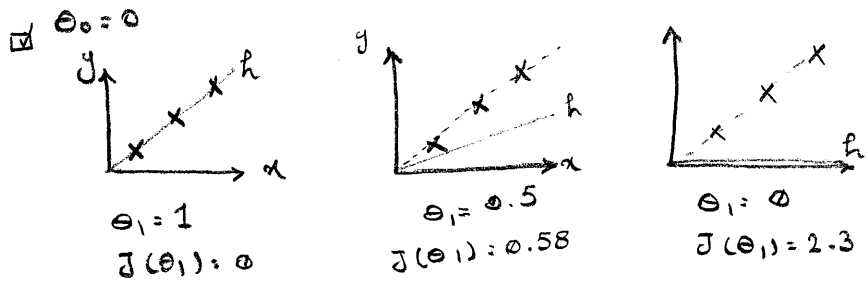
$$\min_{\theta_0, \theta_1} (h_{\theta}(x) - y)^2$$

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

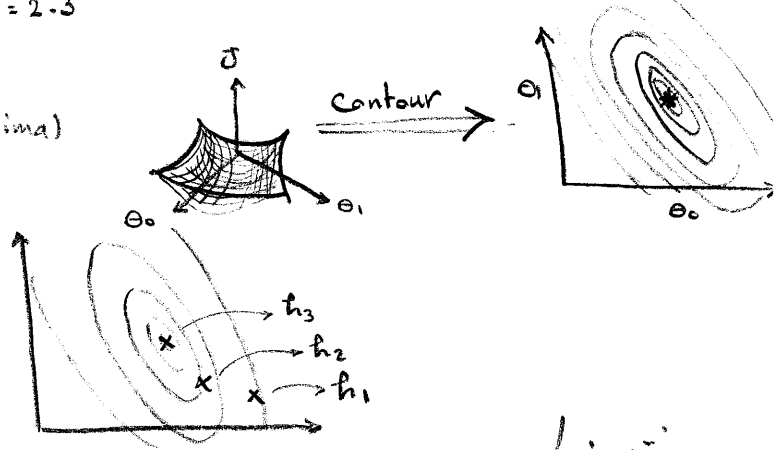
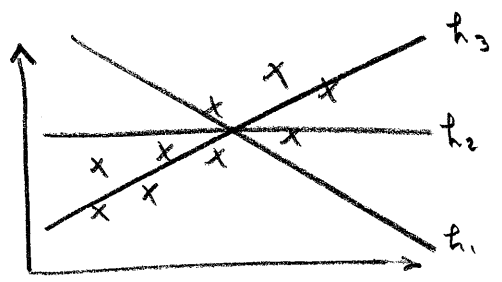
$$\min \frac{1}{2m} \sum_{i=1}^m (h_{\theta}^{(i)}(x) - y^{(i)})^2$$

- hypothesis:  $h_{\theta}(x) = \theta_0 + \theta_1 x$
- params:  $\theta_0, \theta_1$
- cost func:  $J(\theta_0, \theta_1) = \frac{1}{2m} \sum (h_{\theta}(x) - y)^2$
- goal: minimize  $J(\theta_0, \theta_1)$

$J(\theta_0, \theta_1)$ : cost function square error function



$J$  always have bowl shape (no local minima)



پیشن خودکار  $\theta_0$  و  $\theta_1$  نیست  
 - شروع از یک مقدار  $\theta_0, \theta_1$   
 - اندر  $\theta_0, \theta_1$  را عوض می کنیم تا  $J$  کاهش پیدا کند و به کمینه برسد

در مجرایم از تپه پایین میایم، به چه جهت نگاه کنیم؟ به آن سمتی که شیب بیشتری دارد. (steepest descend)

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

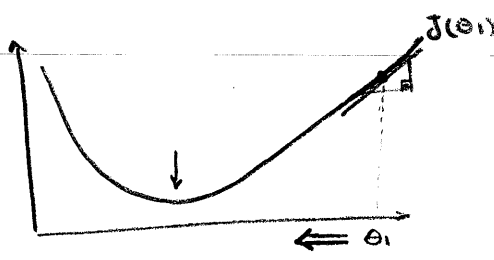
temp0 :=  $\theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$   
 temp1 :=  $\theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$   
 $\theta_0 := \text{temp0}$   
 $\theta_1 := \text{temp1}$

پارامترها باید نزول بکنند.

temp0 :=  $\theta_0 - \alpha \dots$   
 $\theta_0 := \text{temp0}$   
 temp1 :=  $\theta_1 - \alpha \dots$   
 $\theta_1 := \text{temp1}$

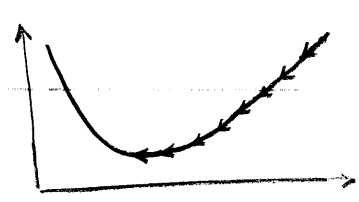
نیجی دیگری دارد. درست →

derivative term

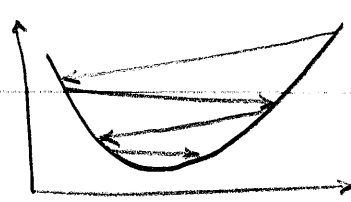


$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$   
 slope  $\gg 0$   
 $\theta_1 := \theta_1 - \alpha (+)$   
 $\theta_1 \downarrow$

د بالعکس



لاصغیر  
 الوریتم لند



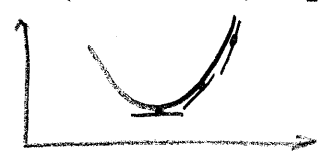
learning rate  
 بزرگ  $\alpha$   
 احتمال عبور از min → overshoot  
 احتمال عدم سبایی



$\theta_1 := \theta_1 - \alpha (0)$

local minima

هر چه به کمینه محلی نزدیک می شویم به طر حذر دار اندازه نگاه لایج می شود (سبایی به افزایش  $\alpha$  نیست)



اعمال GD به LR

مدم اصلی، محاسبه مشتق

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \left( \frac{1}{2m} \sum (h_{\theta} - y)^2 \right)$$

$$= \frac{\partial}{\partial \theta_j} \left( \frac{1}{2m} \sum (\theta_0 + \theta_1 x - y)^2 \right)$$

$\left\{ \begin{array}{l} j=0 : \frac{1}{m} \sum (h_{\theta}(x) - y) \\ j=1 : \frac{1}{m} \sum (h_{\theta}(x) - y) \cdot x \end{array} \right.$

۳.۴ - ترکیب ← جاذب‌های مشتق، حلقه همزمانی

از آنجایی که  $J(\theta_0, \theta_1)$  همواره convex است الگوریتم همیشه به مینیمم محلی می‌رسد.  
 معمولاً  $(\theta_0, \theta_1) = (0, 0)$  برای آغاز انتخاب می‌شود.

- این الگوریتم به گونه‌ای آموزش می‌دهد که در  $(\sum_{i=1}^m)$  نابینان Batch نام دارد.

Extensions

1. حذف  $\alpha$
2. تعداد بیشتر feature

Multiple Features

$n = \#$  of features

$x_j^{(i)}$  = value of feature  $j$  in  $i$ th training example

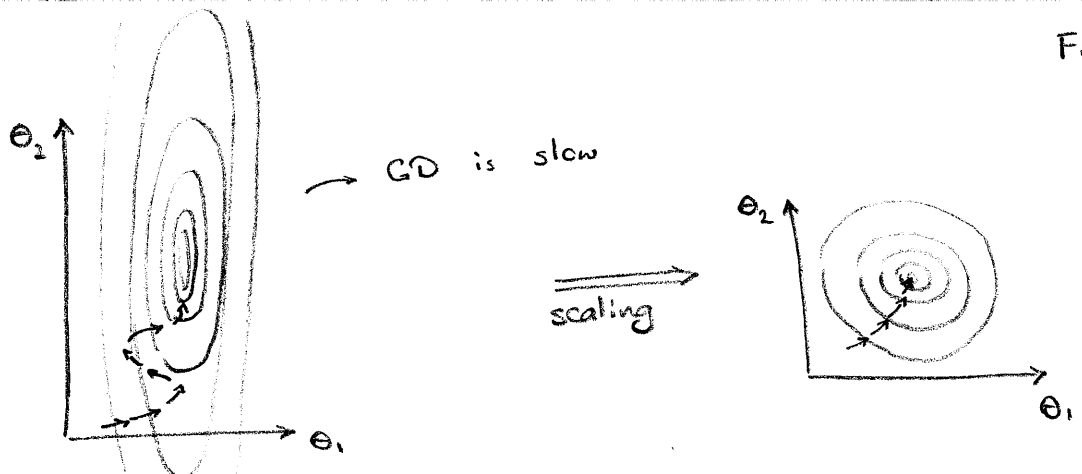
$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n = \sum_{i=0}^n \theta_i x_i \quad \alpha_0 = 1 \rightarrow x = \begin{bmatrix} \alpha_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$= \theta^T x \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

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

Gradient Descent:  $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad j = 0, \dots, n$

Feature Scaling



get every feature approx. into  $[-1, +1]$

- $0 \leq x_1 \leq 3 \quad \checkmark$
- $-2 \leq x_2 \leq 0.5 \quad \checkmark$
- $-100 \leq x_3 \leq 100 \quad \times$
- $0.001 \leq x_4 \leq 0.01 \quad \times$

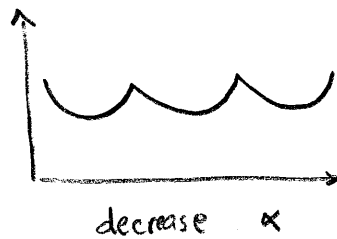
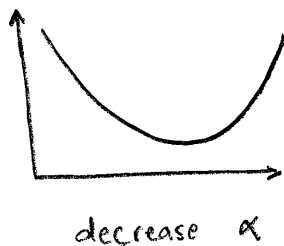
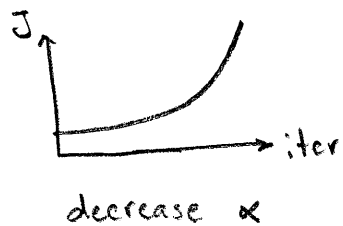
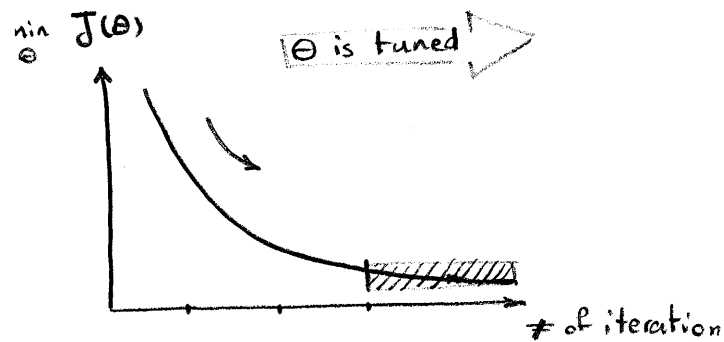
1. Divide By Maximum  $\rightarrow x_i / x_{max}$
2. Mean Normalization  $\rightarrow (x_i - \bar{x}) / x_{max}$
3.  $(x_i - \bar{x}) / \text{range}(x)$
4.  $(x_i - \bar{x}) / \sigma_x \rightarrow$  standard deviation

### Choosing learning Rate

- \*  $J(\theta)$  should decrease after every iter.
- \* Show the convergence

↳ Automatic Convergence Test

↳  $\alpha$  کبھی  $10^{-3}$  تک ہونا چاہیے



- $\alpha$  too small  $\rightarrow$  slow convergence
- $\alpha$  too big  $\rightarrow$  may not decrease on every iter
- $\alpha$  too big  $\rightarrow$  may not converge / slow convergence

- ✓  $\alpha$ : 0.001, 0.01, 0.1, 1, 10, ...
- ✓  $\alpha$ : 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, ...

### Complex Features

✓  $x_1 = \text{frontage}$   
 $x_2 = \text{depth} \Rightarrow x = \text{area} = x_1 \times x_2 \rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x$

### Polynomial Regression

✓  $x_1 = \text{size} \Rightarrow z_1 = x_1, z_2 = x_1^2, z_3 = x_1^3$   
 $h_{\theta}(z) = \theta_0 + \theta_1 z_1 + \theta_2 z_2 + \theta_3 z_3$   
 feature scaling become more important

Cubic

✓  $x_1 = \text{size} \Rightarrow z_1 = x_1, z_2 = \sqrt{x_1}$   
 $h_{\theta}(z) = \theta_0 + \theta_1 z_1 + \theta_2 z_2 = \theta_0 + \theta_1 x_1 + \theta_2 \sqrt{x_1}$

### Normal Equation

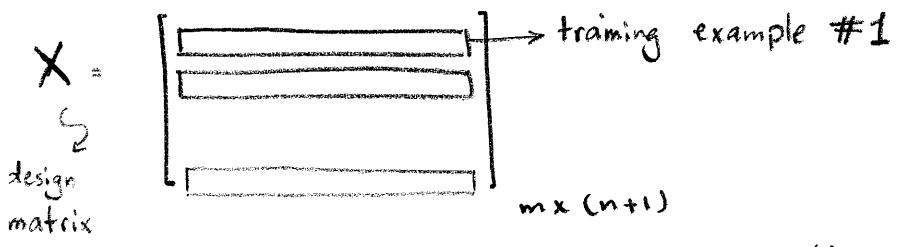
Method to solve for  $\theta$  analytically.

✓  $J(\theta) = a\theta^2 + b\theta + c$   
 $\frac{d}{d\theta} J(\theta) = 2a\theta + b = 0 \rightarrow \theta^*$



Generally:

$J(\theta) = \frac{1}{2m} \sum (h_{\theta}(x) - y)^2 \rightarrow \frac{\partial}{\partial \theta_j} J(\theta) = 0 \quad \forall j \rightarrow \theta_0^*, \theta_1^*, \dots, \theta_n^*$



$$X = \begin{bmatrix} \text{---} (X^{(1)})^T \text{---} \\ \text{---} (X^{(2)})^T \text{---} \\ \text{---} (X^{(m)})^T \text{---} \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & x_1^{(2)} & \dots \\ 1 & x_2^{(1)} & & \dots \\ & & & \dots \end{bmatrix}$$

$\Theta = (X^T X)^{-1} X^T y$   $\rightarrow$  pinv  $(X' * X) * X' * y$

does not need feature scaling

Gradient Desc.	Need $\alpha$	Many Iter	Works with Large $n$ ✓
Normal Eq.	No Need $\alpha$ ✓	No Iteration ✓	Slow with Large $n$

$(X^T X)^{-1} \rightarrow O(n^3)$

- ☑  $n = 100 \rightarrow$  NE ✓
- $= 1000 \rightarrow$  NE ✓
- $= 10000 \rightarrow$  GD ✓
- $= 10^6 \rightarrow$  GD ✓✓

Normal Eq. & Non-invertability.

$X^T X \rightarrow$  non-invertible / singular / degenerate

Causes:  $\rightarrow$  linearly dependant  $\rightarrow$  check features first

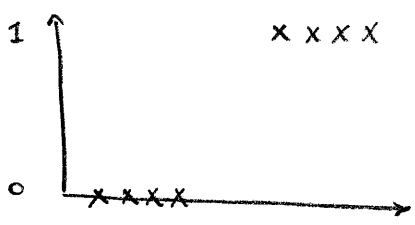
- Redundant features
- ☑  $x_1 = \text{size in ft.}^2 \rightarrow x_1 = (3.28)^2 x_2$
- $x_2 = \text{size in m}^2$

- Too many features
- ☑  $m \leq n \rightarrow$  delete some features
- $\searrow$  use regularization

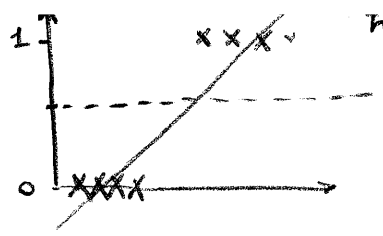
Classification \*  
Two Class .

$y \in \{0, 1\}$   
 $0$ : Negative Class  $\rightarrow$  Absence  
 $1$ : Positive Class  $\rightarrow$  Presence

☑  $0$ : benign tumor ,  $1$ : malignant

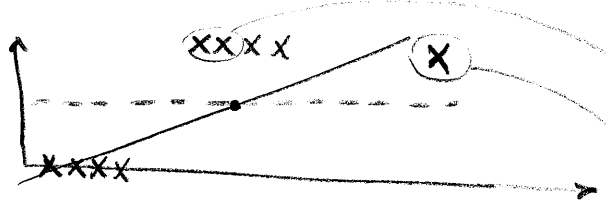


use method we know



threshold  $h_{\theta}(x) \geq 0.5$   
 $\downarrow$   
 1

bad thing



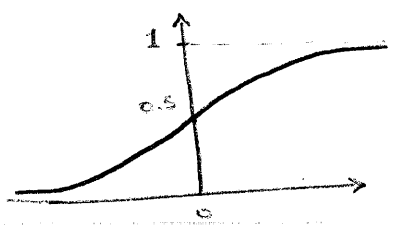
misclassify } reason 1  
 data sample caused it } reason 2  
 $h_{\theta}(x)$  can be  $> 1$  or  $< 0$

جیسی کہ اس شکل میں خواہم، اب یوں کہ  $0 \leq h_{\theta}(x) \leq 1$  ہر وقت ہر طرح Logistic Regression

Logistic Regression

$h_{\theta}(x) = \theta^T x$   $\xrightarrow{0 \leq h \leq 1}$   $h_{\theta}(x) = g(\theta^T x)$

$g(z) = \frac{1}{1 + e^{-z}}$

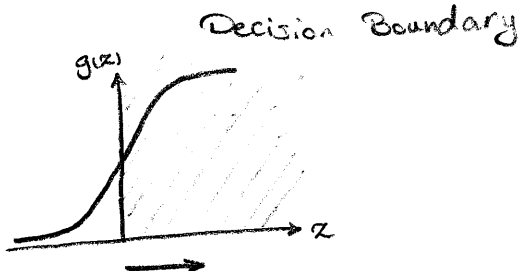


sigmoid / logistic function

$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$   $\rightarrow$  estimated probability that  $y=1$  on input  $x$

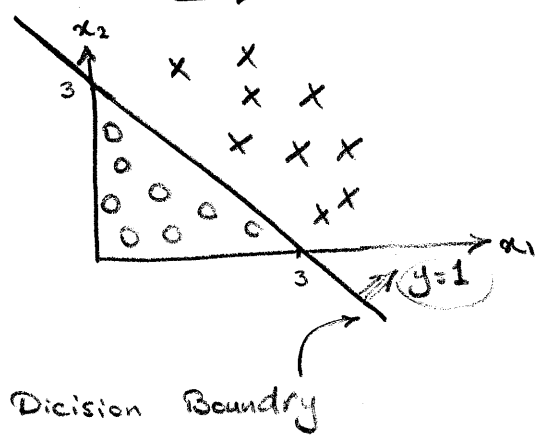
$h_{\theta}(x) = P(y=1 | x; \theta)$   
 $P(y=0 | x; \theta) = 1 - h_{\theta}(x)$

threshold = 0.5  $\rightarrow h_{\theta}(x) \geq 0.5 \Rightarrow y=1$   
 $g(z) \geq 0.5 \Rightarrow z \geq 0$   
 $\theta^T x \geq 0$



$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$   
 $\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$

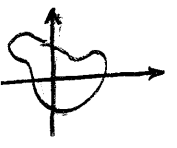
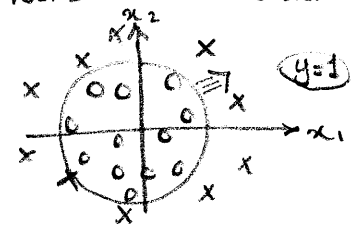
predict "y=1" if  $-3 + x_1 + x_2 \geq 0$   
 $\hookrightarrow$  کب کب  
 $\hookrightarrow h_{\theta}(x) = 0.5$



Decision Boundary

Non-linear Decision Boundary

$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$   
 $\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \rightarrow \text{predict } "y=1" \text{ if } x_1^2 + x_2^2 \geq 1$



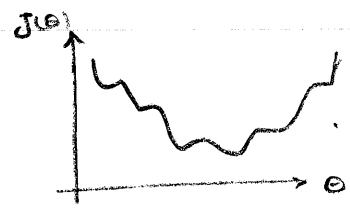
نمونه فرضیه، تابعی که Decision Boundary تعیین کند. (مثلاً دایره)

Training Set :  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

Cost Function

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

$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_{\theta}(x^{(i)}), y^{(i)}) \rightarrow \text{linear regression cost}$   
 $\frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$



برای تابع convex است ولی برای تابع دیگری non convex است

$\rightarrow$  logistic regression cost

برای تابع convex است ولی برای تابع دیگری non convex است

آبیت کجایی این تابع، تخمین ML است

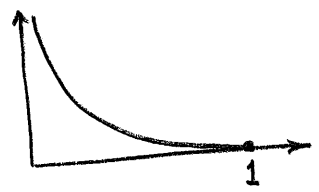
$$\begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

$y=1 \Rightarrow h_{\theta}(x)=1 \Rightarrow \text{cost} = 0$

خاص:

$h_{\theta}(x) \rightarrow 0 \Rightarrow \text{cost} \rightarrow \infty$

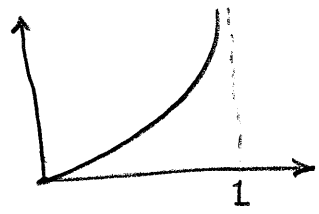
if  $h_{\theta}(x)=0 \rightarrow \text{predict } p(y=1 | x; \theta) = 0$  but  $y=1 \rightarrow \text{large penalty}$



$y=0 \Rightarrow h_{\theta}(x)=0 \Rightarrow \text{cost} = 0$

$h_{\theta}(x) \rightarrow 1 \Rightarrow \text{cost} \rightarrow \infty$

if  $h_{\theta}(x)=1 \rightarrow \text{predict } p(y=0 | x; \theta) = 0$  but  $y=0 \rightarrow \text{large penalty}$





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

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

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

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i^{(j)}) - y_i^{(j)}) x_j^{(i)} \rightarrow \text{! linear regression شبر}$$

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

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

کنترل  $\alpha$  ، feature scaling ، این الگوریتم هم اعمال می شود.

Advanced Optimization .

cost function  $J(\theta) \rightarrow$  want  $\min_{\theta} J(\theta)$

optimization  $\rightarrow$  input  $J(\theta), \frac{\partial}{\partial \theta_j} J(\theta)$

- Gradient Descend
- Conjugate Gradient
- BFGS
- L-BFGS

- ☺ No need to pick  $\alpha$
- ☺ Faster than GD
- ☹ More Complex

function [j\_val, gradient] = costFunction(theta)

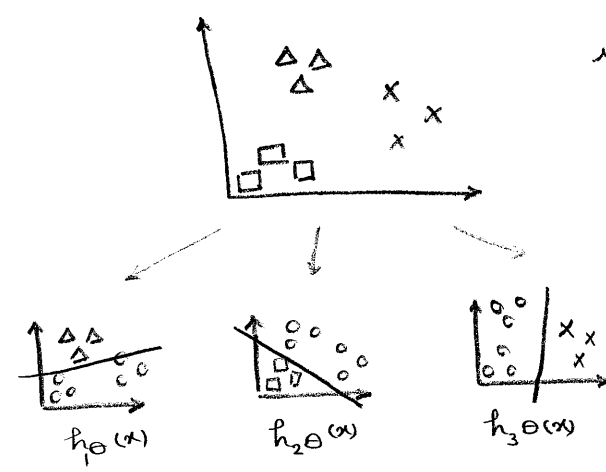
options = optimset('GradObj','on','MaxIter','100')

initialTheta = ...

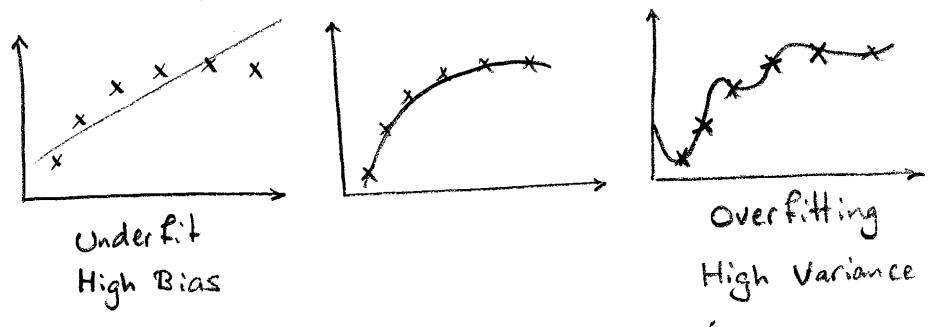
[...] = fminunc(@costFunction, initialTheta, options);

Multi Class Classification : one-vs-all .

- ساختن training set منفی با مقدار دادن target روی کلاس سلف
- آموزش classifier روی هر training set
- برای هر داده ورودی  $x$  ، هر سلف  $h_{\theta}(x)$  را حساب می کنیم که در واقع احتمال عضویت در آن کلاس هستند. بیشینه مقدار آنها کلاس را مشخص می کند.



Regularization \*  
 Overfitting •



↳ not enough data, function goes through all data

Too many features → fit well  
 ↘ fail to generalize (for new data)

Solution •

1. Reduce # of Features
  - Manually
  - Model Selection Alg.

2. Regularization
  - keep all features, but reduce magnitude/values of  $\theta$  → <sup>بسی کم feature</sup> <sup>مقدار</sup> <sup>کو</sup> <sup>کاهش</sup> <sup>دهد</sup> <sup>و</sup> <sup>بسی</sup> <sup>کم</sup> <sup>مقدار</sup> <sup>دهد</sup>

مقدار  $\theta$  ؟ کوچک ← فرضیه ساده تر ← خطر overfitting کمتر

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

↳  $\lambda$ : regularization parameter → balances fitting & generalization → automation needed  
 ↳  $\theta_0$  is not penalized conventionally.

$\lambda$  too large →  $h_{\theta}(x) = \theta_0$  → underfitting

Regularized Linear Regression •

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

Gradient Descent version:

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

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

$0.99 \approx 1$  <sup>عملیاتی</sup> <sup>مقدار</sup> <sup>کو</sup> <sup>کاهش</sup> <sup>دهد</sup>

$$\hookrightarrow 1 - \alpha \frac{\lambda}{m} < 1 \rightarrow \alpha \frac{\lambda}{m} > 2$$

Normal Equation version:

$$\frac{\partial}{\partial \theta_j} J(\theta) = 0 \rightarrow \theta_j^* \Rightarrow \theta = (X^T X + \lambda \underbrace{\begin{bmatrix} 1 & & & 0 \\ & 1 & & \\ & & \dots & \\ 0 & & & 1 \end{bmatrix}}_{(n+1) \times (n+1)})^{-1} X^T y$$

$$\text{if } n=2 \rightarrow \theta = (X^T X + \begin{bmatrix} \lambda & 0 & 0 \\ 0 & \lambda & 0 \\ 0 & 0 & \lambda \end{bmatrix})^{-1} X^T y$$

if  $m \leq n \rightarrow \# \text{ examples} < \# \text{ features} \rightarrow$  w/o reg.: singular matrix  
 w reg.: invertible

### Regularized Logistic Regression

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

Gradient Descent Version: Like Linear Regression with  $h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$

Advanced Optimization

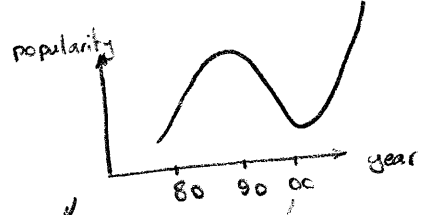
Cost Function  $\rightarrow$  change  $J(\theta)$   
 $\rightarrow$  change  $\frac{\partial}{\partial \theta_j} J(\theta)$   $j = 2, \dots, n+1 \rightarrow \theta_1, \dots, \theta_n$

### Neural Networks \*

logistic regression with full set of polynomial features  $\rightarrow$  برای تعداد بیشتر Feature  
 تعداد ترکیبهای نمایی آینه‌زایی می‌شود.

logistic regression with a subset of polynomial  $f \rightarrow$  کمترین میزان / دشواری انتخاب  $O(n^2)$

high number of simple features  $\rightarrow$  intensity of pixels



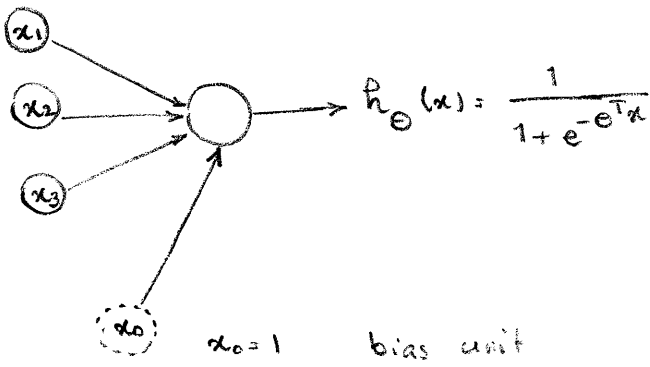
Background  
 - الگوریتم‌هایی که از ستر تولید می‌کنند

auditory cortex  
 somatosensory cortex  
 می‌تواند می‌تواند برخی از وظایف بینایی را انجام دهد. همین طور می‌تواند

"one learning hypothesis"

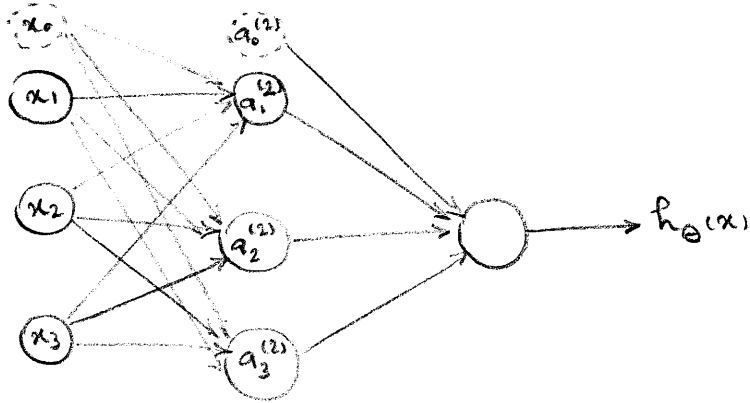
Brain Port

- Seeing with tongue
- Human Ecolocation (sonar)
- Haptic belt: Direction Sense
- Implanting a 3rd Eye



Sigmoid activation function

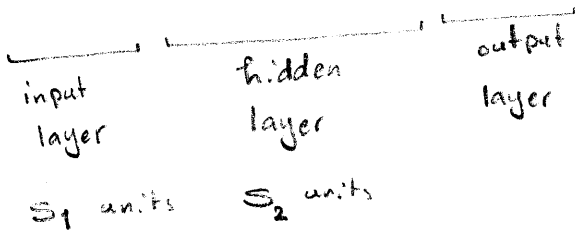
$\Theta$ : parameters / weights



Neural Network

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

$\Theta^{(j)}$ : weights between layer  $j$  and  $j+1$



$$a_1^{(2)} = g(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \dots) = g(z_1^{(2)})$$

$$h_\theta(x) = a_1^{(3)} = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \dots)$$

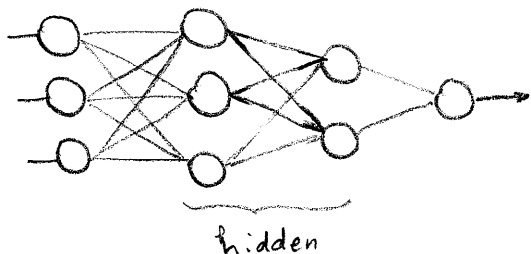
$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} \Rightarrow z^{(2)} = \Theta^{(1)} x \Rightarrow a^{(2)} = g(z^{(2)}) \xrightarrow{\text{Add } a_0} a = \begin{bmatrix} 1 \\ g(z_1^{(2)}) \\ g(z_2^{(2)}) \\ g(z_3^{(2)}) \end{bmatrix}$$

$$\Rightarrow z^{(3)} = \Theta^{(2)} a^{(2)} \Rightarrow a^{(3)} = g(z^{(3)}) = h_\theta(x)$$

کاملاً لایه به لایه است در هر لایه همین کار را forward propagation می‌کنند

در هر لایه یک logistic regres. انجام می‌شود (داده‌ها) Feature های پیچیده تر تبدیل می‌شوند

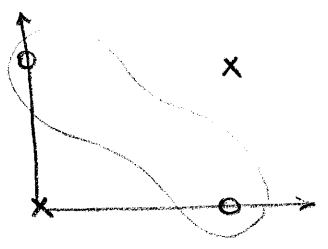
$$\begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} \xrightarrow{\Theta^{(1)}} \begin{matrix} a_1 \\ a_2 \\ a_3 \end{matrix} \xrightarrow{\Theta^{(2)}} h_\theta(x)$$



نکته مهم نودی و کاملاً شبیه Architecture می‌کنند

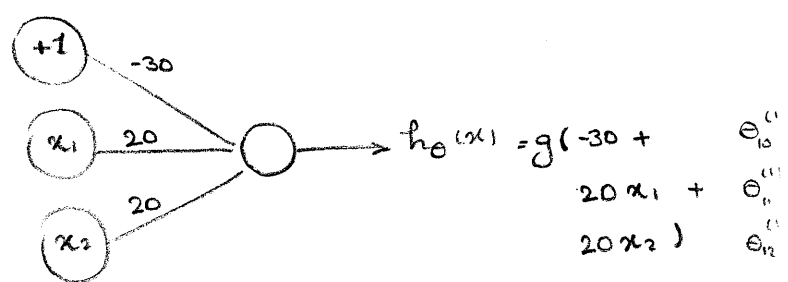
☑ XNOR example

$x: y = 1$   
 $0: y = 0$



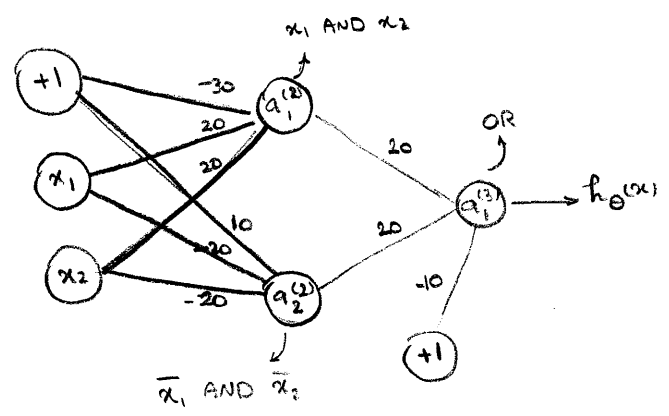
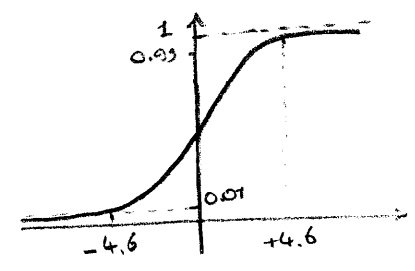
☑ AND example

$y = x_1 \text{ AND } x_2$

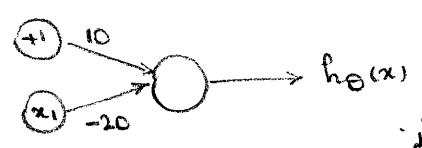


$h_{\theta}(x) = g(-30 + 20x_1 + 20x_2)$

$x_1$	$x_2$	
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(10) \approx 0$
1	1	$g(30) \approx 1$



☑ NOT example



$x_1$	$h_{\theta}(x) = g(10 - 20x_1)$
0	$g(10) \approx 1$
1	$g(-10) \approx 0$

یخراهم  $x_1$ ، NOT لینے  
 سبکی آل مدنی نزل یلایم

One-vs-all  $\leftrightarrow$  Multiple Output Units

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

$y^{(i)} \in \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \right\}$

Cost Function

$L$ : # of layers

$S_{\ell}$ : # of units in layer  $\ell$  (not counting bias units)

Binary Classification:  $y = 0$  or  $1$   
 $h_{\theta}(x) \in \mathbb{R}$   
 $S_{\ell} = 1$   
 $k = 1$

Multi Class:  $y \in \mathbb{R}^k$   
 $h_{\theta}(x) \in \mathbb{R}^k$   
 $S_{\ell} = k$

Cost function of NN is generalization of Logistic Regression.

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

$$+ \sum_{\ell=1}^{L-1} \sum_{i=1}^{S_{\ell}} \sum_{j=1}^{S_{\ell+1}} (\theta_{ji}^{(\ell)})^2$$

$(h_{\theta}(x))_i = i^{\text{th}}$  output

$$a^{(1)} = x$$

$$z^{(2)} = \theta^{(1)} a^{(1)} \Rightarrow a^{(2)} = g(z^{(2)}) \Rightarrow \text{add } a_0^{(2)} = 1 \quad \left. \vphantom{z^{(2)}} \right\} \text{forward propagation}$$

$$z^{(3)} = \dots$$

$\delta_j^{(l)}$  = error of node  $j$  in layer  $l$

$$\delta^{(4)} = a^{(4)} - y$$

$$\delta^{(3)} = (\theta^{(3)})^T \delta^{(4)} \cdot * g'(z^{(3)}) \quad \left. \vphantom{\delta^{(3)}} \right\} \text{back propagation}$$

$$\delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} \cdot * g'(z^{(2)})$$

No  $\delta^{(1)}$   $\rightarrow$  we can't change inputs

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

Back Propagation Alg. •

Training Set  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

Set  $\Delta_{ij}^{(e)} = 0$

For  $i=1$  to  $m$

Set  $a^{(1)} = x^{(i)}$

perform forward propagation  $\Rightarrow a^{(e)}$   $e=2, 3, \dots, L$

$$\delta^{(L)} = a^{(L)} - y^{(i)}$$

Calculate  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

$$\Delta_{ij}^{(e)} += \delta^{(e+1)} (a^{(e)})^T$$

$$\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i)$$

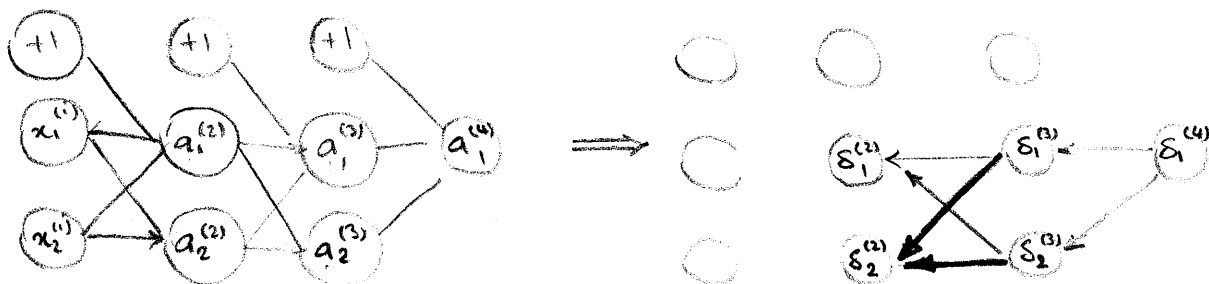
$$\frac{\partial}{\partial \theta_{ij}^{(e)}} J(\theta) = \Delta_{ij}^{(e)}$$

$$\Delta_{ij}^{(e)} = \begin{cases} \frac{1}{m} \Delta_{ij}^{(e)} + \lambda \theta_{ij}^{(e)} & \text{if } j \neq 0 \\ \frac{1}{m} \Delta_{ij}^{(e)} & \text{if } j = 0 \end{cases}$$

$$\delta_2^{(2)} = \theta_{12}^{(1)} \delta_1^{(3)} + \theta_{22}^{(2)} \delta_2^{(3)}$$

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

$$\delta_2^{(3)} = \theta_{12}^{(3)} \delta_1^{(4)}$$



function [jVal, gradient]: costFunction(theta)

Implementation Detail.

optTheta = fminunc (@costFunction, initialTheta, options)

should be vector

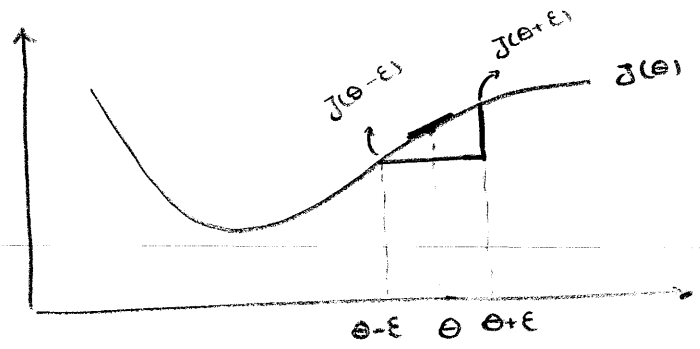
Unrolling the Matrices: thetaVec = [theta1(:); theta2(:); theta3(:)]

similar for D -> gradientVec

Bringing Back to matrix: reshape(thetaVec(0:mxn), m, n)  
reshape(thetaVec(mxn+1:2\*man), m, n)

Gradient Checking

debugging BP



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

$$\epsilon = 10^{-4}$$

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

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

one sided difference NOT GOOD

for i = 1:n

thetaPlus = theta;

thetaPlus(i) = thetaPlus(i) + EPSILON;

thetaMinus = theta;

thetaMinus(i) = thetaMinus(i) - EPSILON;

gradApprox(i) = (J(thetaPlus) - J(thetaMinus)) / (2 \* EPSILON);

Random Initialization

Zero Initialization

نوروزی کی ب لایه بر خود اپرا  
Symmetry  
Redundant Representation

$$a_1^{(2)} = a_2^{(2)}$$
$$\epsilon_1^{(2)} = \epsilon_2^{(2)}$$
$$\frac{\partial}{\partial \theta_{o1}^{(1)}} J(\theta) = \frac{\partial}{\partial \theta_{o2}^{(1)}} J(\theta)$$
$$\theta_{o1}^{(1)} = \theta_{o2}^{(1)}$$

$$\theta_{ij}^{(l)} = 0$$

$$\theta_{ij}^{(l)} \in [-\epsilon, +\epsilon]$$

Symmetry Breaking

Random Initialization

Put all together.

Choose Network Arc

inputs = # of features

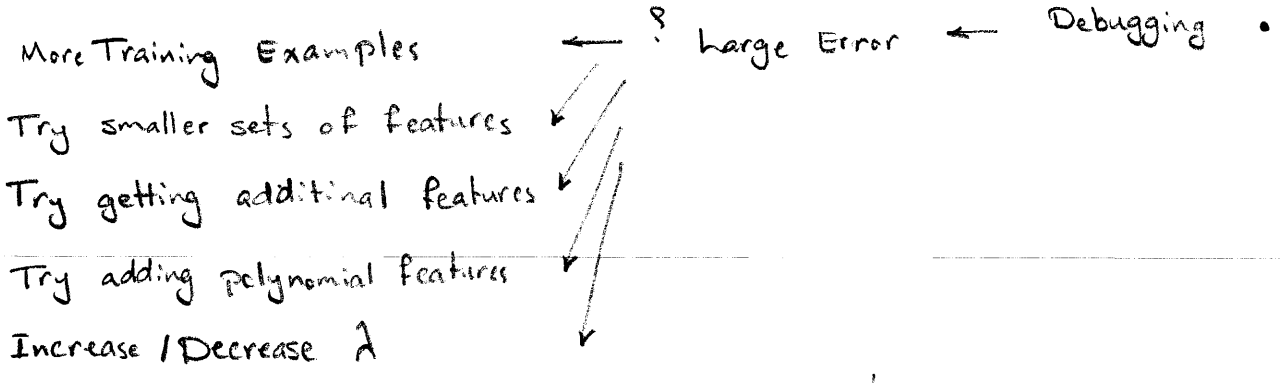
outputs = # of classes

نوروزی کی ب لایه بر خود اپرا در بسته از ب لایه اشتیم تعداد نوروزی ای ای ب لایه نوروزی

1. Randomly initialize weights  $\rightarrow \Theta$
2. Forward Propagation  $\rightarrow h_{\Theta}(x) \rightarrow a^{(i)}$
3. Implement Cost Function  $\rightarrow J(\Theta)$
4. Implement Back Propagation  $\rightarrow \frac{\partial}{\partial \Theta} J(\Theta) \rightarrow \delta^{(i)}$
5. Gradient Checking  $\rightarrow$  BP vs Numerical Estimate
6. Use Gradient Descent to minimize  $J(\Theta) \rightarrow J(\Theta)$  is non-convex  
 $\hookrightarrow$  Advanced optimization

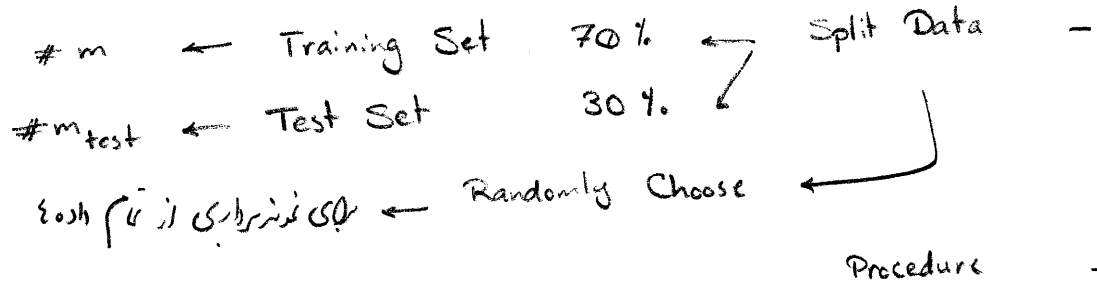
Training -

Advice on Applying ML \*



performance  $\leftarrow$  Diagnostic •  
 (فهمیدن دلیل خطای مدل)

Hypothesis Evaluating •



1. Learn Parameters  $\Theta \rightarrow J(\Theta)$  minimized on training
2. Compute Test Error  $\rightarrow J_{\text{test}}(\Theta)$  on test set

Misclassification Error -

$$\text{err}(h_{\Theta}(x), y) = \begin{cases} 1 & h_{\Theta}(x) \geq 0.5 \\ 0 & h_{\Theta}(x) < 0.5 \end{cases}$$

$$\text{Test\_Err} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(\cdot)$$

Model Selection •

training err  $<$  generalization err  $\leftarrow$  overfitting -

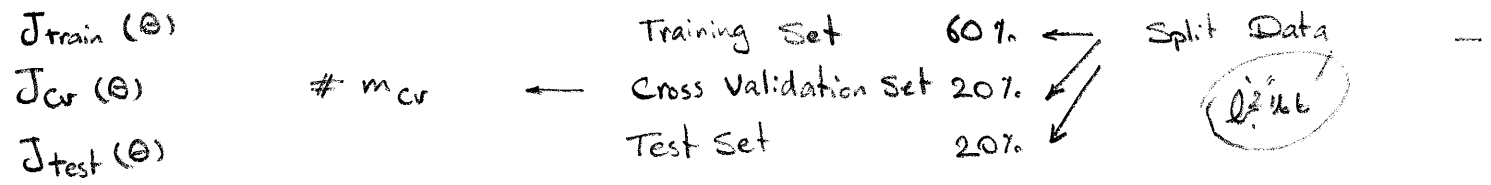
d=1  $h_{\Theta}(x) = \Theta_0 + \Theta_1 x \rightarrow \Theta^{(1)} \rightarrow J_{\text{test}}(\Theta^{(1)})$

d=2  $h_{\Theta}(x) = \Theta_0 + \Theta_1 x + \Theta_2 x^2 \rightarrow \Theta^{(2)} \rightarrow J_{\text{test}}(\Theta^{(2)})$

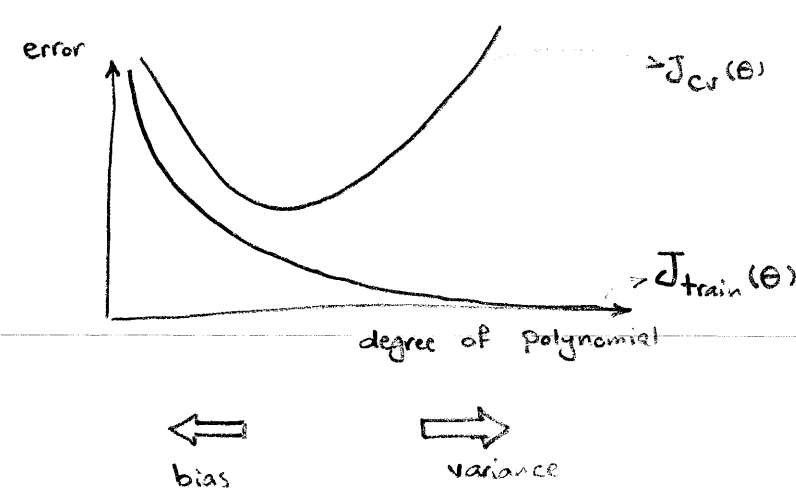
d=3  $h_{\Theta}(x) = \Theta_0 + \Theta_1 x + \Theta_2 x^2 + \Theta_3 x^3 \rightarrow \Theta^{(3)} \rightarrow J_{\text{test}}(\Theta^{(3)}) \checkmark \rightarrow$  new examples ???

d = degree of polynomial





- ✓  $d=1 \rightarrow h_{\theta}(x) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
  - $d=2 \rightarrow h_{\theta}(x) \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$  ✓
  - $d=3 \rightarrow h_{\theta}(x) \rightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$
- overfit است  $d$  بزرگتر

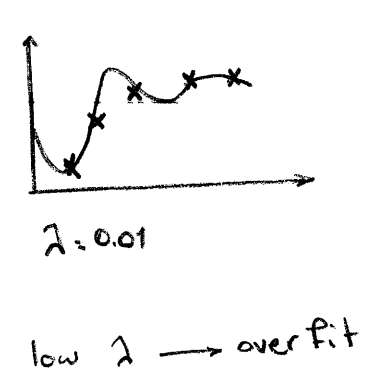
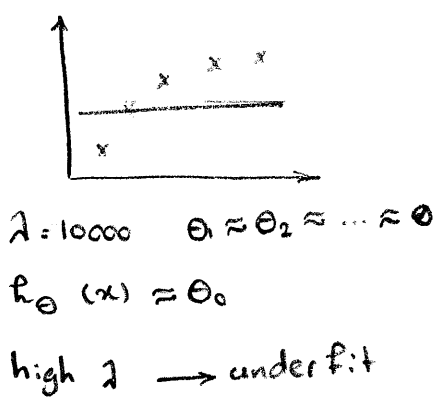


Diagnosing Bias vs. Variance

High Bias:  $J_{train}(\theta)$  will be high  
 $J_{cv}(\theta) \approx J_{train}(\theta)$

High Variance:  $J_{train}(\theta)$  will be low  
 $J_{cv}(\theta) \gg J_{train}(\theta)$

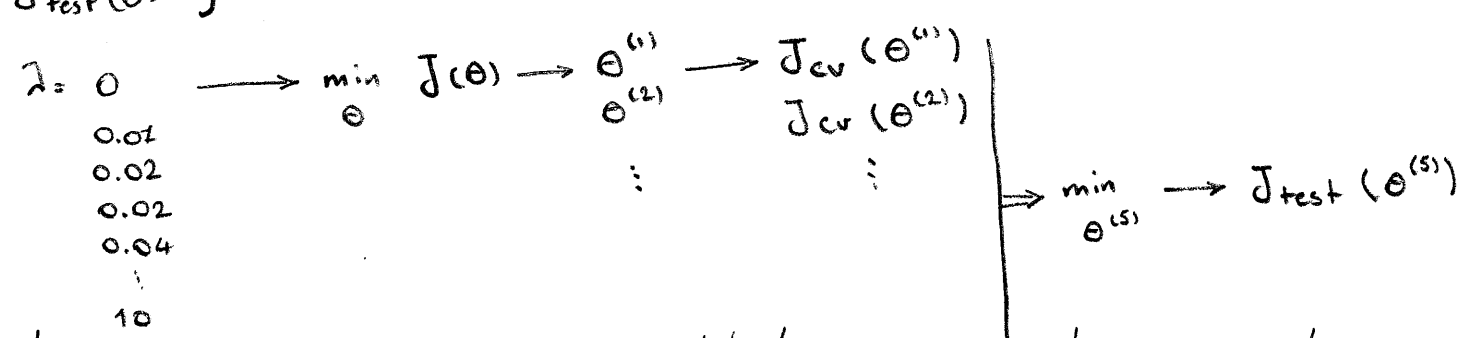
Regularization and Bias / Var



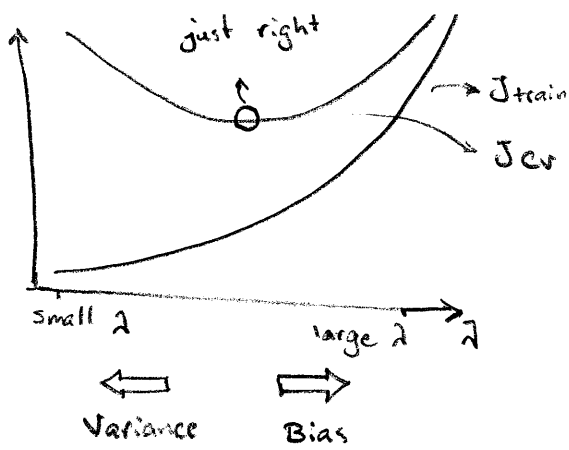
$$J(\theta) = \frac{1}{2m} [cost + \lambda \text{ reg.}]$$

$J_{train}(\theta)$   
 $J_{cv}(\theta)$   
 $J_{test}(\theta)$

} w/o regularization terms



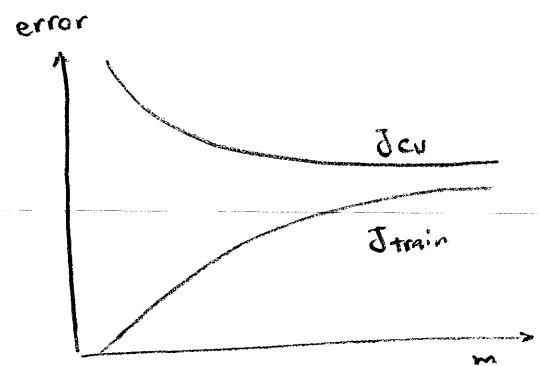
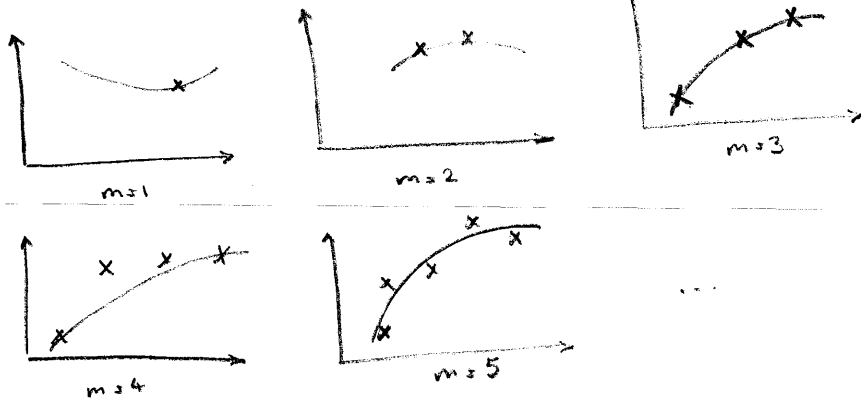
$J_{test}$  کمترین باشد  $\theta$  مربوط به این  $J_{cv}$  را مقایسه می کنیم و کوچکترین را می بینیم  $\rightarrow$  بهترین را می بینیم  $\rightarrow$  بررسی اد بر این می کنیم



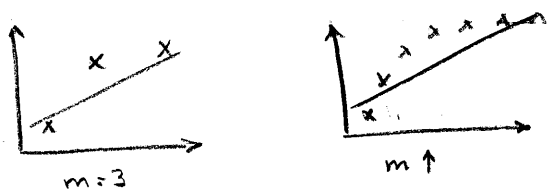
Learning Curves

$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$

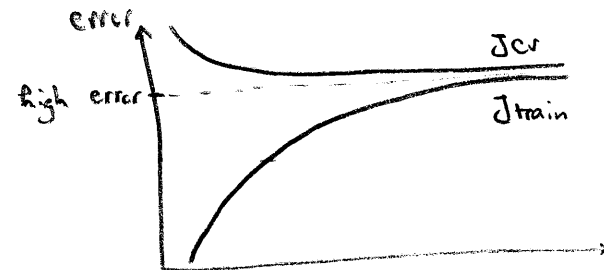
plot error - # of examples -



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

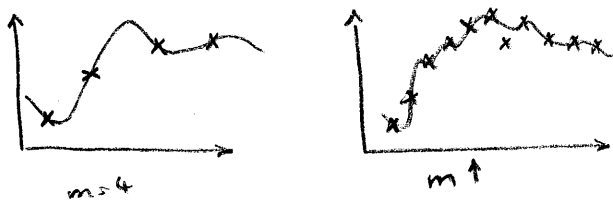


High Bias

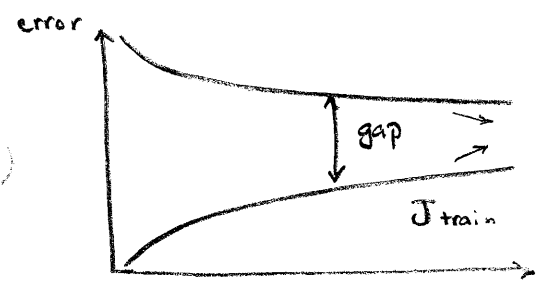


\* more data cannot help

$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_{100} x^{100}$  (and small  $\lambda$ )



High Variance



\* more data helps sometimes

what to do next?

- More training example  $\rightarrow$  fix high var
- Smaller sets of features  $\rightarrow$  fix high var
- getting additional features  $\rightarrow$  fix high bias
- adding polynomial features  $\rightarrow$  fix high bias
- decreasing  $\lambda$   $\rightarrow$  fix high bias
- increasing  $\lambda$   $\rightarrow$  fix high var

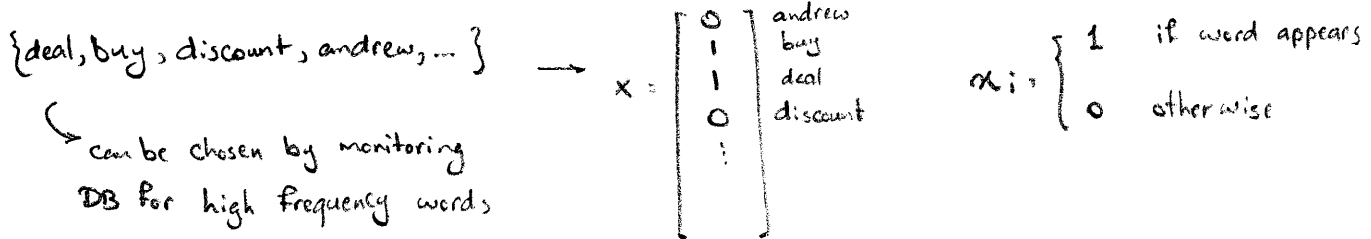
Tips for NN

- Small NN → fewer parameters
  - prone to under fitting
  - computationally cheaper
- Large NN → more parameters → neurons / hidden layers ↑ → use model selection
  - prone to overfitting
  - computationally more complex
  - use regularization

Designing Mh Systems

Example: Spam

- $x$  = features of email
- $y$  = Spam (1) / non Spam (2)
- get 100 words → sort alphabetically → make training vector



- Collect lots of data →  honeypots
- Sophisticated Features →  email routing information → email headers
  - stemming? Punctuation? → email body
  - detect misspelling

Recommended Approach

- Start: Quick Prototype → Test on CV. data
- Plot learning Curves → Decide: more data, more features
- Error Analysis → Manually check False Positives / False Negatives → Trends, Patterns, ...
  - ↳ what type of example it is? → Numeric Evaluation is important!!!
  - ↳ what features are needed to classify them correctly

- Pharma : 12
- Replica / Fake : 4
- Steal Pass word : 53
- Other : 31
- Deliberate misspelling : 5
- Unusual Email Routing : 16
- Unusual Punctuation : 31

↳ Mistakenly Classified By Alg.

- Error Analysis may not be helpful to decide about improving performance.

↳ Implement & Decide

☑ Stemming, Case Sensitivity for spam example → with stemming: 3% error  
w/o stemming: 5% error

Error Metrics for Skewed Classes

☑ Logistic Regression for Cancer  
1% error on test set  
0.5% of patients have cancer

این روش (لایه می‌گردد) که 0.5% خطا داشتیم!  
حالا اگر نت از 99.2% به 99.5% رفت، آیا این بهتر شده است؟

-  $y = 1$  → presence of rare class that we want to detect

		Actual Class	
		1	0
Predicted Class	1	True Positive	False Positive
	0	False Negative	True Negative

Precision =  $\frac{\text{True Positives}}{\text{True Pos} + \text{False Pos}}$  → Predicted +

چقدر از پیش‌بینی‌های ما درباره کلاس نادر درست بوده است؟

Recall =  $\frac{\text{True Positive}}{\text{True Pos} + \text{False Neg}}$  → Actual +

چقدر از اعضای کلاس نادر را درست تشخیص داده‌ایم؟

☑ در مثال بالا کس  $y = 0$  سبب  $\text{Recall} = 0$  است.

- Predict  $y = 1$  only if very confident

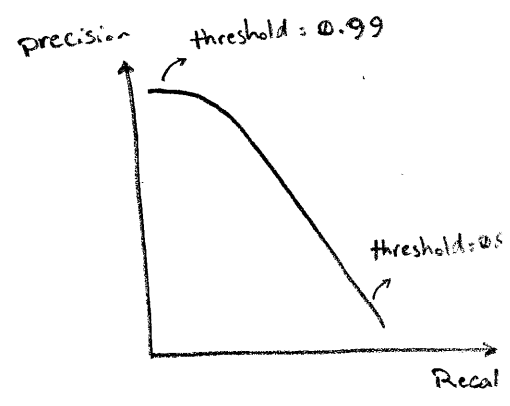
$$\begin{cases} h_{\theta}(x) \geq 0.5 \rightarrow 0.7 \\ h_{\theta}(x) < 0.5 \rightarrow 0.7 \end{cases}$$

precision ↑  
Recall ↓

- Predict  $y = 1$  avoiding missing too many cancer

$$\begin{cases} h_{\theta}(x) \geq 0.5 \rightarrow 0.3 \\ h_{\theta}(x) < 0.5 \rightarrow 0.3 \end{cases}$$

Precision ↓  
Recall ↑



other possibilities

- F1 Score

	Precision	Recall	Average	F1 Score
Alg 1	0.5	0.4	0.45	0.444 ✓
Alg 2	0.7	0.1	0.4	0.175
BAD Alg 3	0.02	1.0	0.51	0.0392
Alg 4	1.0	0.02	0.51	0.0392

$F_1 \text{ Score} = 2 \frac{PR}{P+R}$

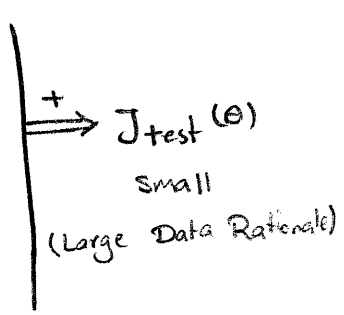
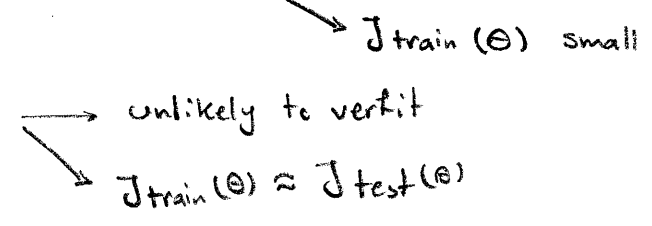
$P=0 \text{ or } R=0 \Rightarrow F_1=0$

$P=1 \text{ and } R=1 \Rightarrow F_1=1$

Data for Machine Learning

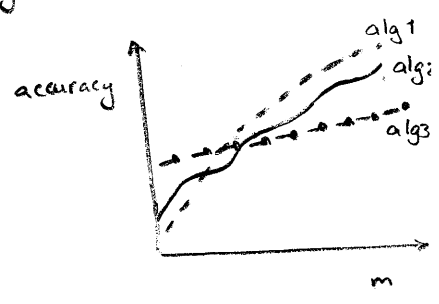
Learning Algorithms with many parameters → low bias algorithms

Use a very large training set



test: can a human expert predict results with that info? is data enough?

$m \uparrow \Rightarrow$  accuracy  $\uparrow$  (monotonically)  
 if alg is low bias

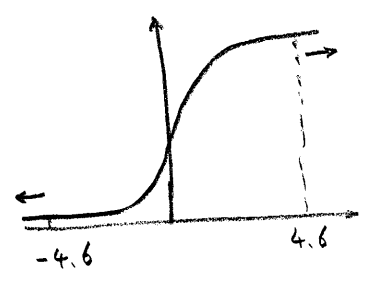


Support Vector Machines \*

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

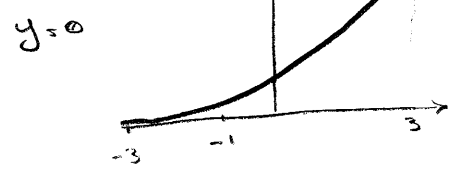
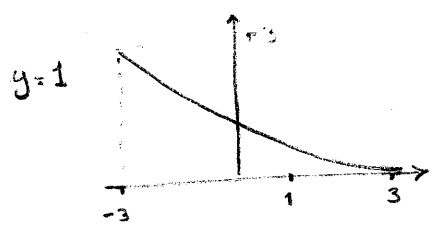
if  $y=1 \rightarrow$  we want  $h_{\theta}(x) \approx 1 \rightarrow \theta^T x \gg 0$

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



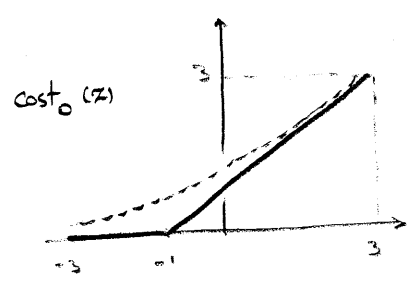
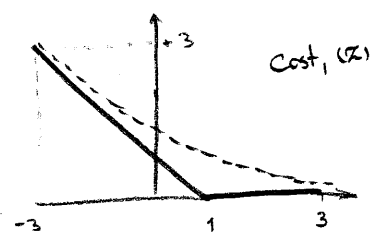
Optimization Objective.

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



Cost logistic regression  
 $= A + \lambda B$   
 err ← A      reg ← B

approx.



Cost SVM

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_2(\theta^T x^{(i)}) + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

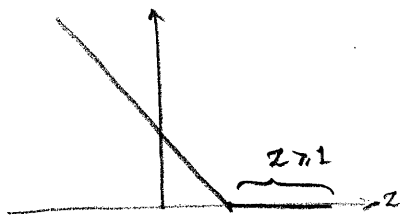
\* omit  $\frac{1}{m}$  in SVM  
 \* change from  $A + \lambda B$  in Logistic Regres. →  $C A + B$  in SVM

$$\hookrightarrow \text{Cost}_{SVM} = \min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_2(\theta^T x^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$$\hookrightarrow h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

↪ *نقطه جدایی*

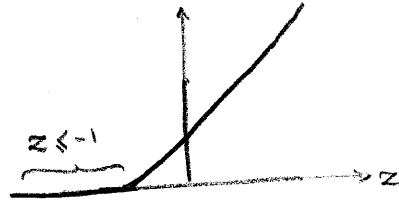
Large Margine Intuition



cost<sub>1</sub>(z)

y=1

want  $\theta^T x \geq 1$



cost<sub>0</sub>(z)

y=0

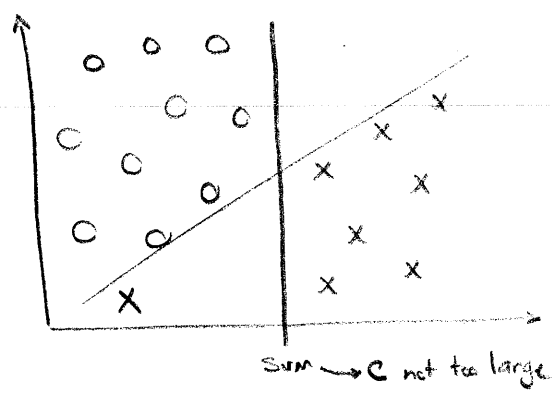
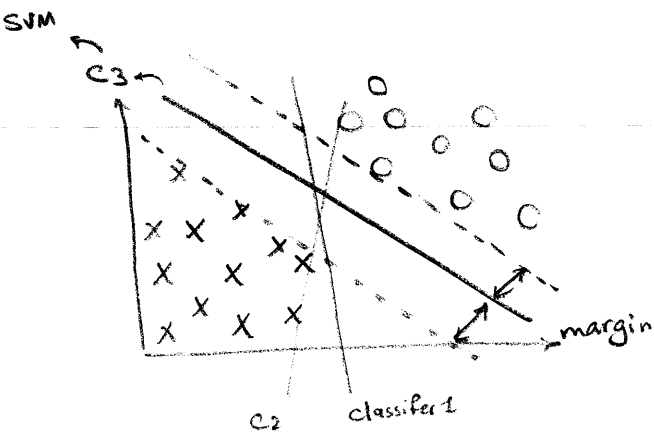
want  $\theta^T x \leq -1$

if  $c \gg 0$

$$\min c \sum [y \text{cost}_1(\cdot) + (1-y) \text{cost}_0(\cdot)] + \frac{1}{2} \sum \theta_j^2$$

$y^{(i)} = 1 \rightarrow \min c x_0 + \frac{1}{2} \sum \theta_j^2$   
s.t.  $\theta^T x^{(i)} \geq 1$

$y^{(i)} = 0 \rightarrow \theta^T x^{(i)} \leq -1$



Large Margine Classifier

Outliers? Non linear separable case?

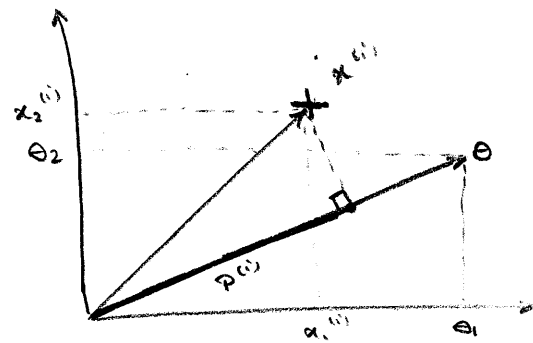
$$\min \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} (\sqrt{\theta_1^2 + \theta_2^2})^2 = \frac{1}{2} \|\theta\|^2$$

Simplification:  $\theta_0 = 0, n=2$

s.t.  $\theta^T x^{(i)} \geq 1$  if  $y^{(i)} = 1$

$\theta^T x^{(i)} \leq -1$  if  $y^{(i)} = 0$

$$\left. \begin{aligned} \theta^T x^{(i)} &= p^{(i)} \|\theta\| \\ \theta^T x^{(i)} &= \theta_1 x_1 + \theta_2 x_2 \end{aligned} \right\}$$

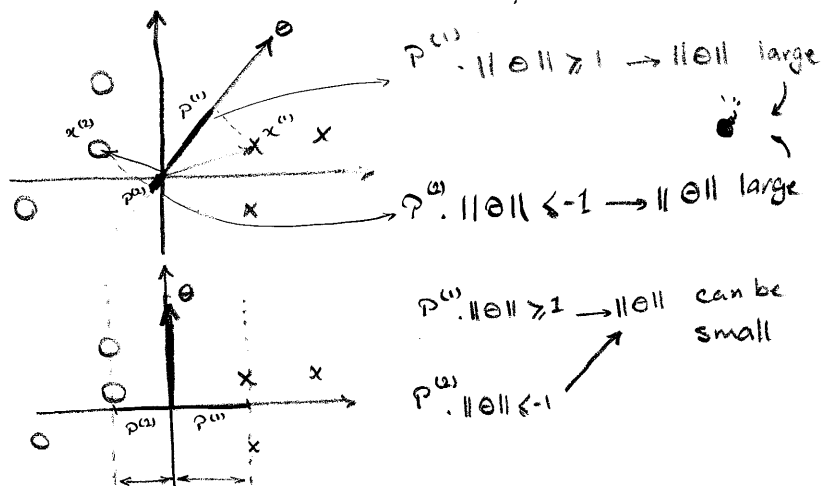


$$\min \frac{1}{2} \sum \theta_j^2 = \frac{1}{2} \|\theta\|^2$$

s.t.  $p^{(i)} \cdot \|\theta\| \geq 1$  if  $y^{(i)} = 1$

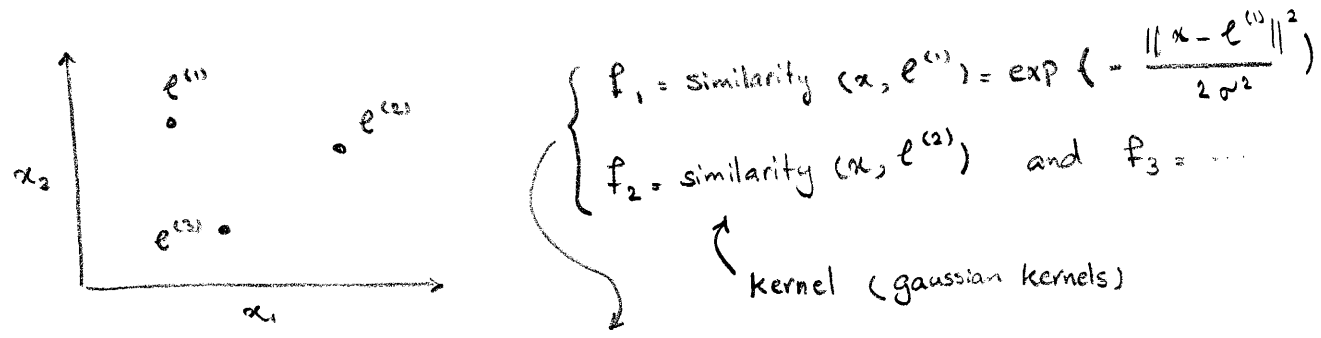
$p^{(i)} \cdot \|\theta\| \leq -1$  if  $y^{(i)} = 0$

simplification:  $\theta_0 = 0$



Non-Linear Decision Boundary.

predict  $y=1$  if  $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \dots \geq 0$   
 $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \theta_4 f_4 + \dots \geq 0$

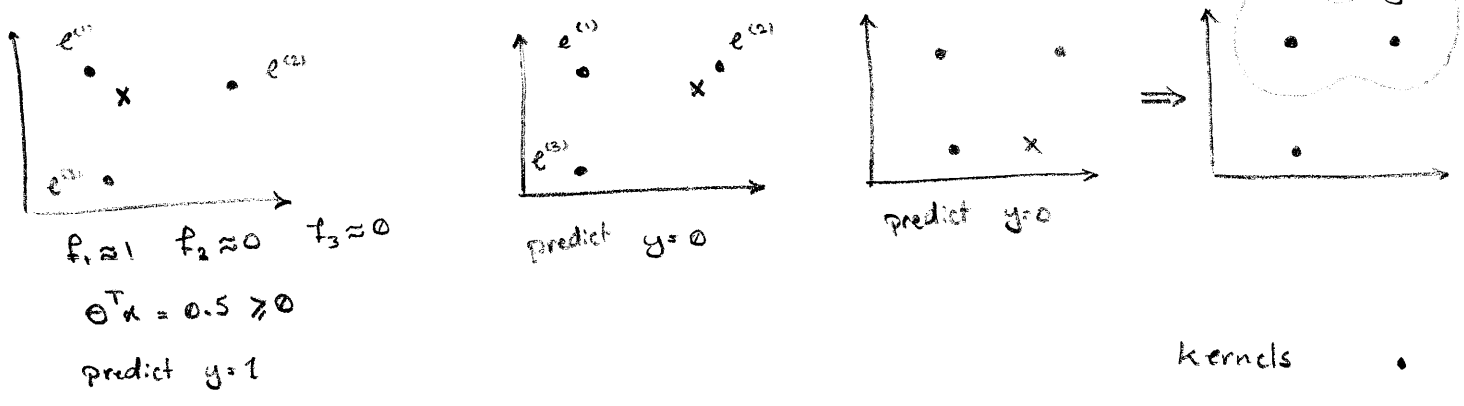


features are proximity to landmarks

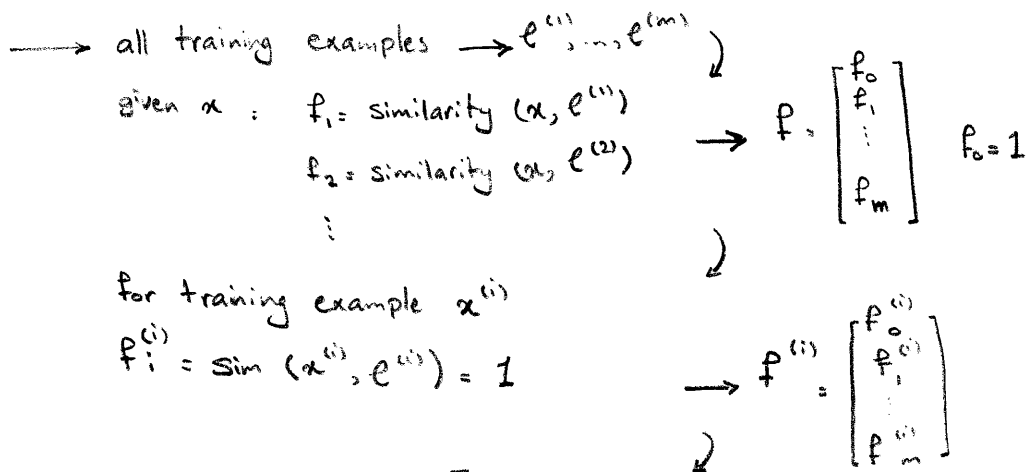
$$f_1 = \exp\left(-\frac{\sum_{j=1}^n (x_j - e_j^{(1)})^2}{2\sigma^2}\right) \begin{cases} \rightarrow x \approx e^{(1)} \Rightarrow f_1 \approx 1 \\ \rightarrow x \text{ far } e^{(2)} \Rightarrow f_1 \approx 0 \end{cases}$$

$$\text{SVM} \begin{cases} h_{\theta}(x) = 1 & \text{if } \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

□  $\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$



where to get landmarks?



predict  $y=1$  if  $\theta^T f \geq 0$

$\theta$  from training:  $\min_{\theta} c \sum \sim + \frac{1}{2} \sum_{j=1}^m \theta_j^2$

$m = n$

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

-  $\theta^T \theta$      $\theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$

-  $\theta^T M \theta$      $\rightarrow$  scaling for large databases ( $m = 10,000$ )

$c (= \frac{1}{\lambda})$      $\rightarrow$  large  $c$  : lower bias, high variance

$c$  پارامتر

$\searrow$  small  $c$  : high bias, low variance

پارامتر  $\sigma^2$  کنترل نویز

$\sigma^2$      $\rightarrow$  large :  $f_i$  vary smoothly  $\rightarrow$  high bias, low var

$\searrow$  small :  $f_i$  vary less smooth  $\rightarrow$  low bias, high var

Using SVM

- Use Library: liblinear, libsvm

- Choice of  $c$

- Choice of kernel     $\rightarrow$  linear kernel (no kernel)     $\rightarrow$   $y = \pm 1$  if  $\theta^T x \geq 0$

$\rightarrow$   $n$  large,  $m$  small ( $n \gg m$ )

$\rightarrow$  gaussian kernel

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

$\rightarrow$  use feature scaling

$\rightarrow$  need  $\sigma^2$

$\rightarrow$   $n$  large,  $m$  large

$\rightarrow$  other  $\rightarrow$  convergence condition: kernel satisfy "Mercer's Theorem"

$\rightarrow$  polynomial kernel

$\rightarrow k(x, \ell) = (x^T \ell)^2$   
 $= (x^T \ell)^3$   
 $= (x^T \ell + 1)^3$   
 $= (x^T \ell + \alpha)^\gamma$

$\rightarrow$  string kernel, chisquare kernel, histogram intersection kernel, ...

- Multi class     $\rightarrow$  built in in packages

$\rightarrow$  one-vs-all     $\rightarrow$  pick class  $i$  with largest  $(\theta^{(i)})^T x$

-  $n$  vs  $m$      $\rightarrow$   $n \gg m$      $\rightarrow$  logistic regression, SVM w/o kernel

$n = 10,000$   
 $m = 10 \sim 1,000$

$\rightarrow$   $n$  small,  $m$  intermediate     $\rightarrow$  SVM with gaussian kernel

$n = 1 \sim 1,000$   
 $m = 10 \sim 10,000$

$\rightarrow$   $n$  small,  $m$  large     $\rightarrow$  add more feature and logistic regression/linear SVM

$n = 1 \sim 1,000$   
 $m = 50,000 +$

NN works for all but must be slower to train (but local optima is a disadvantage)



# Unsupervised learning \*

Find some structure in data ← unlabeled data •

## k-means Clustering \*

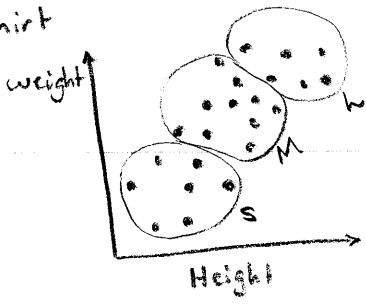
algorithm •

randomly initialize cluster centroids → number = clusters  
assignment step → each data is closer to what centroids  
move centroids step → move centroids to average of groups

input: k # of clusters, m # of training examples

non-separated clusters.

### ✓ T-shirt



k-means can do that!

→ Market Segmentation

Optimization Objective •

$c^{(i)}$  = index of closest cluster to input  $x^{(i)}$

$\mu_k$  = cluster centroid k

$\mu_{c^{(i)}}$  = cluster centroid of cluster to which  $x^{(i)}$  is assigned

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2 \quad \rightarrow \text{cost function, distortion}$$

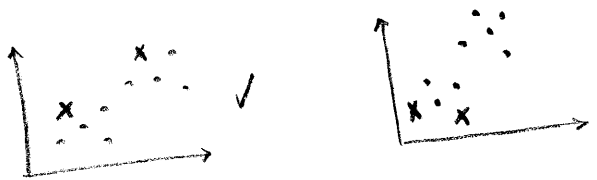
cluster assignment :  $\min J(w)$  w.r.t  $c^{(1)}, c^{(2)}, \dots, c^{(m)}$  holding  $\mu_1, \dots, \mu_k$

move centroid :  $\min J(w)$  w.r.t  $\mu_1, \mu_2, \dots, \mu_k$  holding  $c^{(1)}, \dots, c^{(m)}$

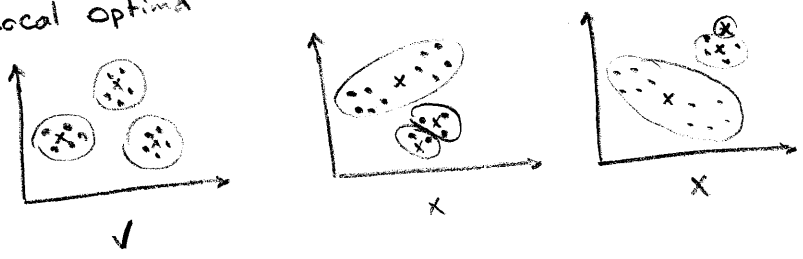
Random Initialization.

-  $k < m$

- Randomly pick k examples → set  $\mu_1, \dots, \mu_k$  to these examples.



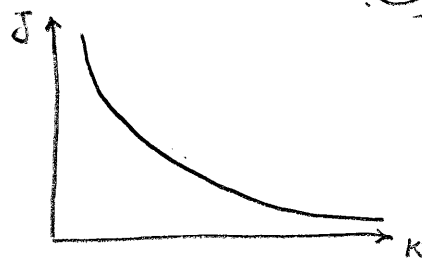
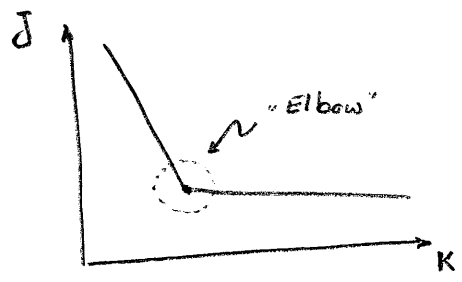
- Local optima



- Run Random Initializations Several times → Pick one giving lowest  $J(w)$

Choosing the number of clusters

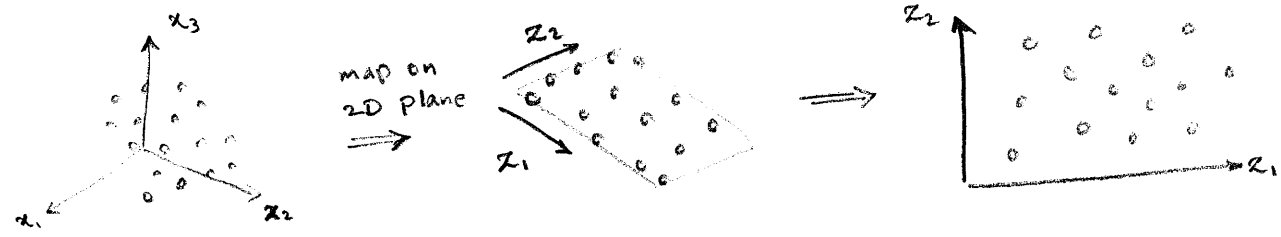
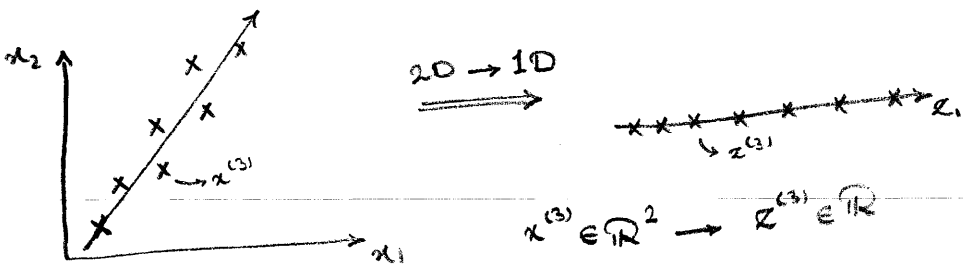
تعداد کلاسترهای موجود در داده معمولاً داخلی نیست.



اگر elbow وجود داشت آن انتخاب می‌کنیم.

اگر نمودار clustering کارایی دیگری روی داده انجام می‌دهیم، می‌توان K را با توجه به feedback این تغییرات.

Dimensionality Reduction



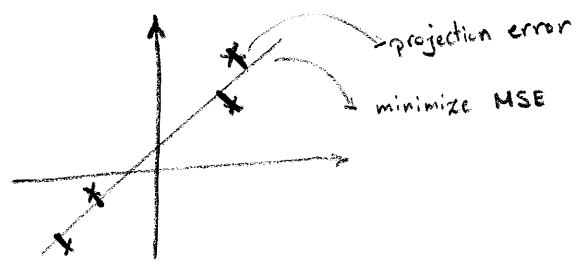
Speed up learning  
 Retain 99% of variance

Data Visualization, Data Compression  
 $nD \Rightarrow 2D$   
 $nD \Rightarrow 3D$

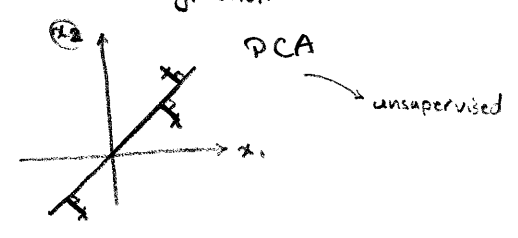
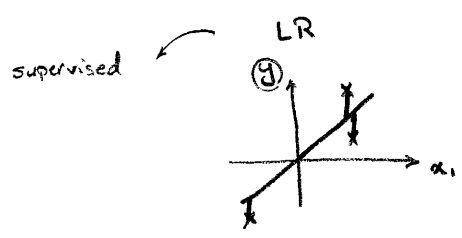
PCA Formulation

problem: find a direction onto which to project data so as to minimize the projection error

$mD \Rightarrow kD \rightarrow$  find  $k$  vectors  $u^{(1)}, u^{(2)}, \dots, u^{(k)}$



linear regression :  $\text{نشان}$   
 PCA :  $\text{نشان}$   
 supervised vs unsupervised



PCA Algorithm

- Training Set :  $x^{(1)}, \dots, x^{(m)}$
- Preprocessing  $\rightarrow$  Feature scaling / mean normalization  $\rightarrow$  Zero Mean is Mandatory

- PCA :  $nD \rightarrow kD$

compute covariance matrix  $\rightarrow \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)})(x^{(i)})^T$

compute eigenvectors of  $\Sigma \rightarrow$  Singular Value Decomposition  $\rightarrow$  svd (matlab)

Take the  $k$  first eigen vectors

$$U = \begin{bmatrix} | & | & & | & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} & \dots & u^{(m)} \\ | & | & & | & | \end{bmatrix} \in \mathbb{R}^{n \times n} = \text{svd}(\Sigma)$$

$$U_{\text{reduce}} = \begin{bmatrix} | & | & & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ | & | & & | \end{bmatrix}_{n \times k} \Rightarrow Z_{k \times 1} = U_{\text{reduce}}^T \times X_{n \times 1}$$

(matlab)

Sigma =  $X' * X / m$  ;

[U,S,v]=svd(sigma);

Ured = U(:,1:k);

Z = Ured' \* X ;

Choosing the number of PCs

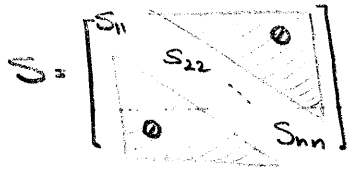
Average Square Projection Err =  $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$

Total variation in data =  $\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$

Choose  $k$  to be smallest value so that

$$\frac{\sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01 \rightarrow 1\%$$

\* 99% of variance is retained \*



- start  $k=1 \rightarrow$  check  $\oplus \rightarrow k++$

inefficient

$$1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \leq 0.01$$

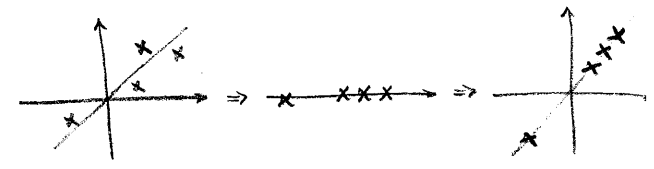
- start  $k=1 \rightarrow$  check  $\oplus \oplus \rightarrow k++$

highly efficient

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99$$

Reconstruction

$$x = U_{\text{reduce}}^T x \rightarrow x_{\text{approx}} = U_{\text{reduce}} \cdot Z$$



$(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$

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

$\{x^{(1)}, \dots, x^{(m)}\} \xrightarrow[nD]{\text{PCA}} \{z^{(1)}, \dots, z^{(m)}\} \xrightarrow[kD]{\text{PCA}} \{(z^{(1)}, y^{(1)}), \dots, (z^{(m)}, y^{(m)})\} \rightarrow$  training

Cross Validation & Test set should undergo the same PCA, PCA defined on training set.

Supervised learning Speed up

Not a solution to prevent overfitting •

Use regularization instead ←

Use no information about y  
 Throws away some information

? Part of ML Strategy •

Try your best w/o PCA → improve algorithm for original data

Use PCA ← if it does not what you want

Anomaly Detection \* Problem •

▣ Aircraft Engine

$x_1$  = heat generated

$x_2$  = vibration intensity



- Problem: Is  $x_{test}$  anomalous?

↳ Make Model  $p(x)$  →  $P(x_{test}) < \epsilon$  → anomaly  
 =  $P(x_{test}) \geq \epsilon$  → normal

▣ Fraud Detection

$x_1$  = Typing Speed of User

$x_2$  = Login Frequency

$x_3$  = Transaction Type

Identify unusual users  
 $p(x) < \epsilon$

▣ Monitoring Computers in data center

$x_1$  = memory use

$x_2$  = # disk access / sec

$x_3$  = CPU load

$x_4$  = CPU load / Net. traffic

Gaussian Distribution •

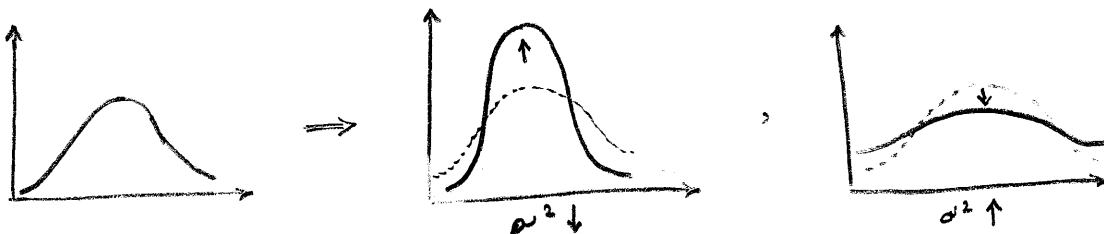
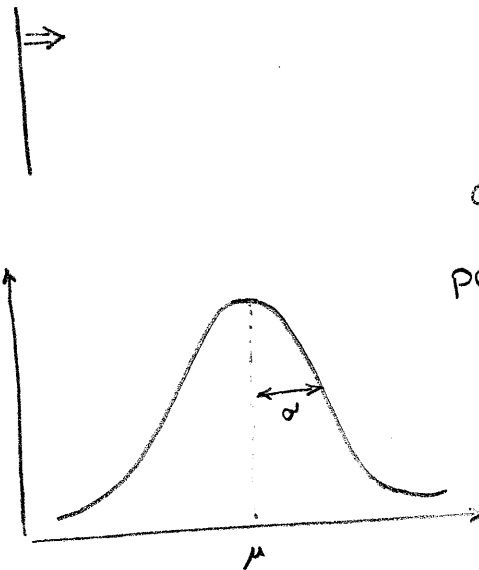
$p(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

- Gaussian Distribution

mean →  $\mu$

variance →  $\sigma^2$

$x \sim \mathcal{N}(\mu, \sigma^2)$   
 distributed as



- parameter estimation:

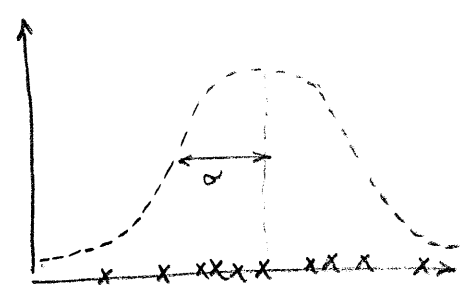
data  $\{x^{(1)}, \dots, x^{(m)}\}$

$x^{(i)} \sim \mathcal{N}(\mu, \sigma^2)$

$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$

$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$

in ML  $\frac{1}{m}$   
 in Statistics  $\frac{1}{m-1}$  → comes from Maximum Likelihood



Anomaly Detection Alg •

- Density Estimation Problem

Independence Assumption between  $x_i$  (holds in real world application)

$x_1 = \mathcal{N}(\mu_1, \sigma_1^2)$

$x_2 = \mathcal{N}(\mu_2, \sigma_2^2)$

⋮

$x_n = \mathcal{N}(\mu_n, \sigma_n^2)$

$P(x) = P(x_1; \mu_1, \sigma_1^2) \dots P(x_n; \mu_n, \sigma_n^2)$

$= \prod_{j=1}^n P(x_j; \mu_j, \sigma_j^2)$

- Algorithm

1. Choose features  $x_i$  might be indicative of anomalous examples
2. Fit parameters  $\mu_i$  and  $\sigma_i^2$  →  $\mu = \frac{1}{m} \sum x^{(i)}, \sigma^2 = \frac{1}{m} \sum (x - \mu)^2$
3. Given new example compute  $P(x)$  →  $P(x) = \prod \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu)^2}{2\sigma^2})$
4. Anomaly if  $P(x) < \epsilon$

Evaluating Algorithm •

Better Result? ← | New Features —  
 | Threshold Change —  
 | Labeled Data —

$y=0 \rightarrow$  normal  
 $y=1 \rightarrow$  anomaly

Training Set  $x^{(1)}, \dots, x^{(m)} \rightarrow$  normal examples

CV Set  $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m)}, y_{cv}^{(m)})$

Test Set  $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m)}, y_{test}^{(m)})$

19000 good (normal) engine  
 20 Flawed engine → { training: 6000 good ( $y=0$ ) →  $\mu, \sigma^2$  (60%)  
 CV: 2000 good, 10 anomalous (20%)  
 test: 2000 good, 10 anomalous (20%)

1. Fit  $p(x)$  on training set  $x^{(1)}, \dots, x^{(m)}$

2. On Cross validation / test set

$$y = \begin{cases} 1 & p(x) < \epsilon \\ 0 & p(x) \geq \epsilon \end{cases}$$

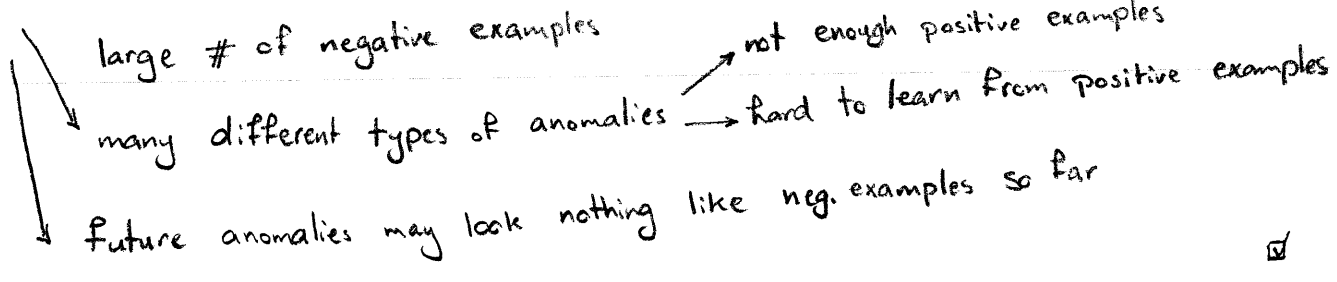
3. Possible Evaluation Metric (Skew classes):  $\rightarrow$  normal examples are dominating

- TP, FP, TN, False Negative
- Precision / Recall
- $F_1$  Score

4. Use Cross Validation to Choose  $\epsilon$  and new Features

Anomaly Det vs. Sup Learning -

anomaly det  $\rightarrow$  very small # of positive examples



Anomaly  $\rightarrow$  Fraud Detection, Manufacturing, Monitoring

Supervised  $\rightarrow$  Email Spam, Weather Prediction, Cancer Classification

Choosing what Features to Use -

1. Plot Histogram of Features

$\rightarrow$  gaussian  $\checkmark$

$\rightarrow$  non gaussian  $\rightarrow$  algorithm usually handles  $\rightarrow$  do nothing

$\searrow$  skewed hist  $\rightarrow$  use log transform



$$\checkmark \log(x_1), \log(x_2 + \epsilon), \sqrt{x_3}, x_4^{\frac{1}{3}}$$

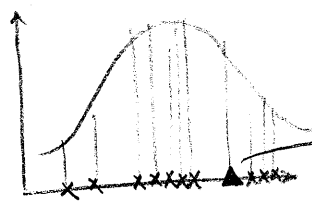
$$\checkmark \text{hist}(x, 50) \rightarrow \text{hist}(x^{0.5}, 50) \rightarrow \text{hist}(x^{0.4}, 50)$$

$$\rightarrow \text{hist}(-\log(x), 50)$$

آزادگی نه غوطه لوسی نمود.

2. Error Analysis

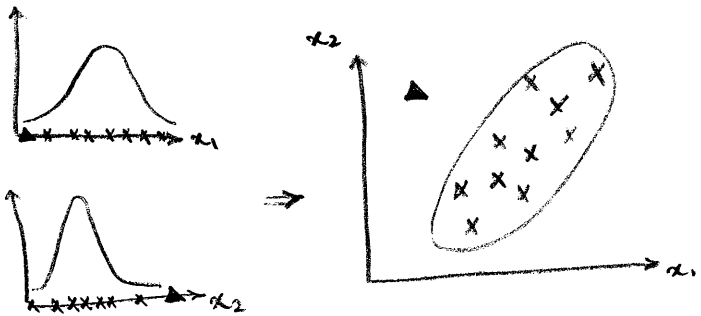
$p(x) \rightarrow$  large for normal  
 $\searrow$  small for anomaly



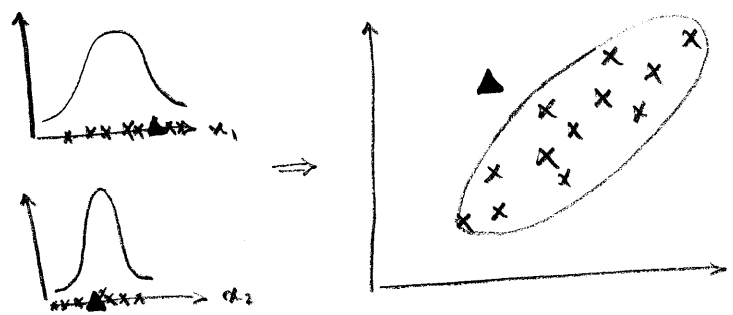
inspect anomalous example for suspicious values

Choose Features that might take on unusually large or small values in the event of anomaly

$$\checkmark x_3 = \text{CPU load} \quad x_4 = \text{Net traffic} \rightarrow x_5 = \text{CPU load/Net traffic} \rightarrow x_6 = (\text{CPU load})^2 / \text{Net traffic}$$



OK



Normal Anomaly Detection Fails

↳ Don't Model  $p(x_i)$  separately

$$P(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

$$\mu \in \mathbb{R}^n \rightarrow \mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\Sigma \in \mathbb{R}^{n \times n} \rightarrow \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

Algorithm -

1. Fit model  $p(x) \rightarrow \mu, \Sigma$
2. Calculate  $P(x^*)$  for example  $x^* \rightarrow P(x^*)$
3.  $P(x^*) < \epsilon \rightarrow$  Anomaly

$$\Sigma = \begin{bmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & \\ & & \ddots & \\ & & & \sigma_n^2 \end{bmatrix} \rightarrow \text{no correlation between features}$$

Original Model -

Original Anomaly Detection

Multi variate

- \* Manually create features to capture when  $x_1, x_2$  take unusual combination
- \* Computationally Cheaper  $\rightarrow$  Scale better
- \* Good for small  $m$

- \* Automatically captures correlations between features
  - \* Computationally More Expensive  $\rightarrow \Sigma^{-1}$
  - \*  $m > n \rightarrow \Sigma$  is noninvertible for  $m < n$   
usually used when  $m \gg 10n$
  - \* Redundant Features  $\rightarrow$  Rank of  $\Sigma \downarrow \rightarrow \Sigma$  maybe non-invertible
- $x_1 = x_2$   
 $x_3 = x_4 + x_5$

Predicting Movie Ratings

	User 1	U2	U3	U4
Romance 1	5	5	0	0
R 2	5	?	?	0
R 3	?	4	0	?
Action 1	0	0	5	4
Action 2	0	0	5	?

↓

Romance  
Lover

Action  
Lover

$n_u = \# \text{ user}$   
 $n_m = \# \text{ movies}$   
 $r(i,j) = \text{if user } j \text{ rated movie } i$   
 $y(i,j) = \text{rating user } j \text{ rated movie } i$   
 Problem: guess question marks?  
 ↳ recommend to user if they predicted to like a movie

Content-Based Recom.

	U1	U2	U3	U4	$x_1$	$x_2$
$x^{(1)} \rightarrow R1$	5				0.9	0
$x^{(2)} \rightarrow R2$	5				1.0	0.01
$x^{(3)} \rightarrow R3$	?				0.99	0
$x^{(4)} \rightarrow A1$	0				0.1	1.0
$x^{(5)} \rightarrow A2$	0				0	0.9

$x^{(1)} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix}$   
 $x^{(4)} = \begin{bmatrix} 1 \\ 0.1 \\ 1.0 \end{bmatrix}$  (with  $x_1=1$ )  
 $n_u = 4$   
 $n_m = 5$   
 $n = 2 \rightarrow \# \text{ features}$   
 $m = 5 \rightarrow \# \text{ examples}$

for each user (j) → learn  $\theta^{(j)} \in \mathbb{R} \rightarrow$  predict user j as rate movie (i) with  $(\theta^{(j)})^T x^{(i)}$

$x^{(3)} = \begin{bmatrix} 1 \\ 0.99 \\ 0 \end{bmatrix}$      $\theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}$      $(\theta^{(1)})^T x^{(3)} = 4.95$

Problem -

- $r(i,j)$
- $y(i,j) = \text{rating of user } j \text{ for movie } i$
- $\theta^{(j)} = \text{params for user } j$
- $x^{(i)} = \text{features for movie } i$
- $m^{(j)} = \# \text{ of movies rated by user } i$

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

user j : 
$$\min_{\theta^{(j)}} \frac{1}{2} \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2}$$

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

$$J(\theta^{(1)}, \dots, \theta^{(n_u)})$$



Collaborative Filtering

- Given  $y^{(i,j)}$  &  $\theta^{(j)}$
- Want  $x$

□

	$u_1$	$u_2$	$u_3$	$u_4$	$x_1, x_2$
$x^{(1)}$	5	5	0	0	? ?
$x^{(2)}$	5				
$x^{(3)}$	?				
$x^{(4)}$	0				
$x^{(5)}$	0				

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

$$(\theta^{(1)})^T x^{(1)} \approx 5$$

$$\approx 5$$

$$\approx 0$$

$$\approx 0$$

probably a love movie because  $u_1$  and  $u_2$  that loves romance movie loved it!

learn  $x^{(i)}$

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

optimization alg -

learn  $\alpha$

$$\min_{\alpha^{(1)} \dots \alpha^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T \alpha^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (\alpha_k^{(i)})^2$$

- { Given  $x$  &  $y \Rightarrow \theta$
- { Given  $\theta \Rightarrow x$

all together -

Guess  $\theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \theta \rightarrow \dots$

improve features & parameters by going back & forth

To Minimize Simultaneously:

$$J(x^{(1)} \dots x^{(n_m)}, \theta^{(1)} \dots \theta^{(n_u)}) = \frac{1}{2} \sum_{(i,j):r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

$$\min_{x, \theta} J(x, \theta)$$

no need to  $x_0 = 1$

1. Initialize  $\alpha^{(1)} \dots \alpha^{(n)}$ ,  $\theta^{(1)} \dots \theta^{(n)}$  to small random values

2. Minimize  $J(\omega)$  using gradient descend

$$\alpha_k^{(i)} := \alpha_k^{(i)} - \alpha \left( \sum (\theta_k^{(j)})^T \alpha^{(i)} - y^{(i,j)} \right) \theta_k^{(j)} + \lambda \alpha_k^{(i)}$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left( \sum (\theta_k^{(j)})^T \alpha^{(i)} - y^{(i,j)} \right) \alpha^{(i)} + \lambda \theta_k^{(j)}$$

3. For a user with param  $\theta$  and movie with (learned features  $\alpha$ ) predict a rating

$$\theta^T \alpha \rightarrow (\theta^{(j)})^T \alpha^{(i)}$$

user

5	5	0	0
5	?	?	0
?	4	0	?
0	0	5	4
0	0	5	0

movie

$$(\theta^{(j)})^T \alpha^{(i)} = \begin{bmatrix} (\theta^{(j)})^T \alpha^{(i)} \\ \vdots \\ (\theta^{(m)})^T \alpha^{(m)} \end{bmatrix}$$

Vectorized Form -

$y^{(i,j)}$

$$X = \begin{bmatrix} -(\alpha^{(1)})^T \\ \vdots \\ -(\alpha^{(m)})^T \end{bmatrix}$$

$$\Theta = \begin{bmatrix} (\theta^{(1)})^T \\ \vdots \\ (\theta^{(m)})^T \end{bmatrix}$$

Low Rank Matrix Factorization

vectorize version of collaborative filtering

How to find movie  $j$  related to movie  $i$ ?

$$\| \alpha^{(i)} - \alpha^{(j)} \| \text{ small}$$

$$\min_j \| \alpha^{(i)} - \alpha^{(j)} \| \rightarrow \text{most similar movie to movie } i$$

Implementation -

new user 5 who haven't rated any movie

5	5	0	0	?
5	?	?	0	?
?	4	0	?	?
0	0	5	4	?
0	0	5	5	?

$$\xrightarrow{n=2} \min J(\omega) = \dots + \frac{\lambda}{2} [(\theta_1^{(5)})^2 + (\theta_2^{(5)})^2] \rightarrow \theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

همه چیز را صفر از برای می کند

Mean Normalization

$$\mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix} \rightarrow \text{new } Y = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$$

train

for user  $j$  on movie  $i$  predict  $(\theta^{(j)})^T \alpha^{(i)} + \mu_i$

$$\text{us} \rightarrow \theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \rightarrow (\theta^{(5)})^T \alpha^{(i)} + \mu_i = \mu_i$$

Mini-Batch Gradient Descent •

examples in each iter: Batch  $\rightarrow$  Mini Batch  $\rightarrow$  Stochastic  
 $m \quad b \quad 1$

$\hookrightarrow$  mini-batch size 2~100

□  $b=10$

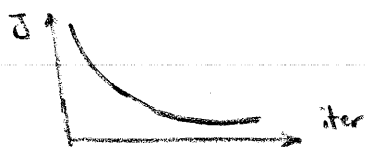
Get  $b=10$  examples

$$\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=1}^{i+9} (h_{\theta}(x^{(k)}) - y^{(k)}) x_j^{(k)} \quad \forall j=0, \dots, n$$

- Mini Batch GD  $\rightarrow$  Use Vectorization  $\rightarrow$  Partial Parallelism  
 $\searrow$  free parameter  $b$  ☹️

- Batch GD  $\rightarrow$  Plot  $J_{train}$  vs iter

Convergence •

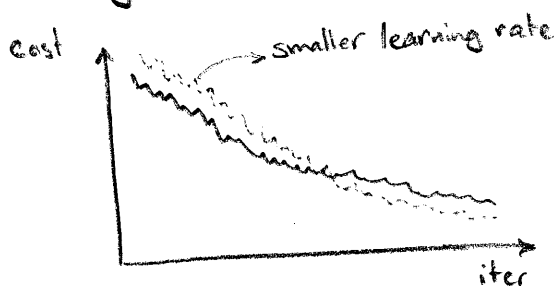


- Stochastic GD

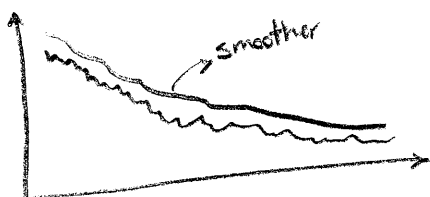
$$\text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

during learning, compute cost before updating  $\theta$

every 1000 iter plot averaged 1000 costs  $\rightarrow t=1000$



$\Rightarrow$  in SGD the algorithm oscillate around min so  $\alpha \downarrow \Rightarrow$  better result

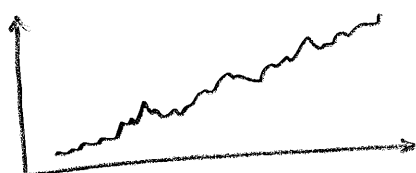


$\Rightarrow$  every 5000 iter

$\underline{t} \uparrow \rightarrow$  smoother curve



$\Rightarrow$  more appropriate bigger  $\underline{t}$  shows trend better



$\Rightarrow$  diverging  $\rightarrow$  use smaller  $\alpha$

- learning rate  $\rightarrow$  usually constant

$\searrow$  slowly decrease  $\rightarrow \alpha = \frac{c1}{iter\# + c2}$

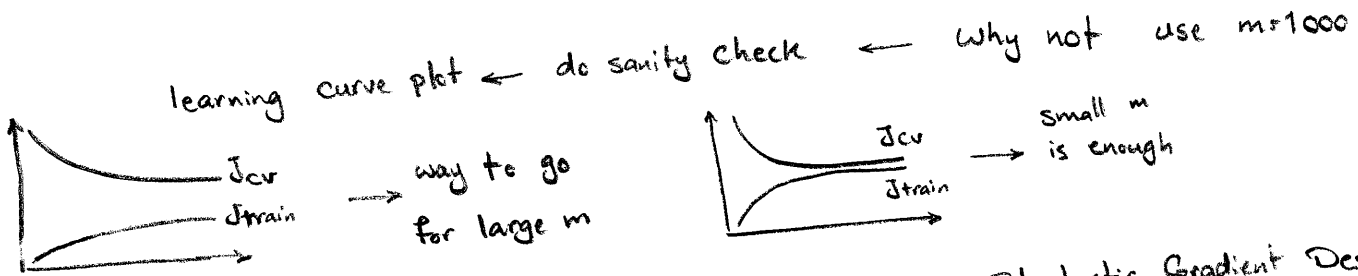
$\nearrow$  ☺️ convergence

$\nearrow$  ☹️  $c1$  &  $c2$  params

Large Dataset \*

data  $\uparrow$   $\rightarrow$  accuracy  $\uparrow$   $\leftarrow$  low bias ML algorithm.

gradient descent for linear regression  $\rightarrow \sum_{i=1}^m \text{circle}$   $m = 100,000,000$



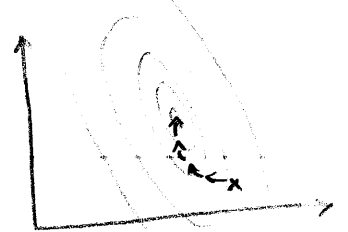
Stochastic Gradient Descent.  
here: for linear regression  $\downarrow$

-  $h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$

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

Repeat  $\left\{ \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad \forall j=0, \dots, n \right\}$

Batch GD  $\rightarrow m = 300,000,000$   $\rightarrow$  all in memory  $\rightarrow \sum$  over all



- Stochastic GD:

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

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

- Algorithm

1. Randomly Shuffle Dataset

2. Repeat  $\left\{ \right.$

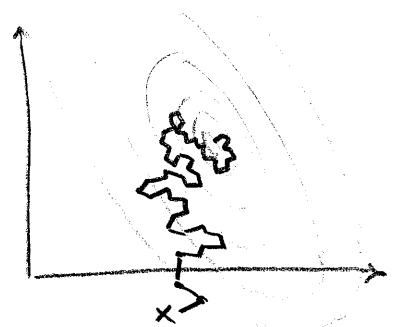
for  $i=1, \dots, m$   $\left\{ \right.$

$\theta_j := \theta_j - \alpha (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \quad \forall j=0, \dots, n$

$\left. \right\}$

$\frac{\partial}{\partial \theta_j} \text{cost}(\theta, (x^{(i)}, y^{(i)}))$

SGD  $\rightarrow$  every iter is really fast  
 $\rightarrow$  do not converge completely to global minima  
 $\rightarrow$  not necessarily improves  $\theta$  in each iter



### Shipping Service Price

offer a price  
 {  $y=1$  user accept  
    $y=0$  user goes away

learn  $\Rightarrow P(y=1 | x; \theta)$   
 to optimize price

continuous stream of data

use data once then discard it.  
 adapt to change user preferences.

Repeat forever {

Get  $(x, y)$  corresponding <sup>to</sup> user.

Update  $\theta$  using  $(x, y)$

$$\theta_j := \theta_j - \alpha (h_{\theta}(x) - y) x_j \quad j=0, \dots, n$$

### Product Search

user search a query  
 suggest a number of results  
 {  $y=1$  user clicks  
    $y=0$  otherwise

learn  $\Rightarrow P(y=1 | x; \theta)$   
 to predict

Click-Through-Rate (CTR)

e.g. suggest 10 phones

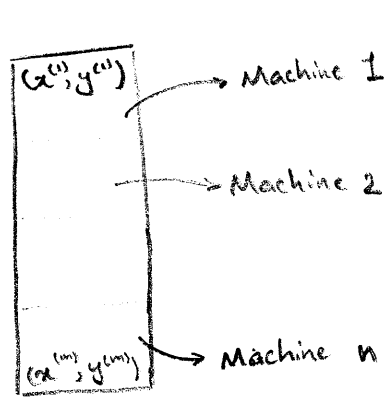
10 training sets  $(x, y)$   $\rightarrow$  10 iter of learning  $\rightarrow$  discard data

special offers, news articles, product recommendations, ...

Map-Reduced Data Parallelism

### Batch Gradient Descent

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



Machine 1: Use  $(x^{(1)}, y^{(1)}) \dots (x^{(100)}, y^{(100)})$   
 $temp_j^{(1)} = \sum_{i=1}^{100} (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$

Combine:  
 $\theta_j := \theta_j - \alpha \frac{1}{m} (temp_j^{(1)} + temp_j^{(2)} + \dots)$

- network latency, overhead, combination time  $\rightarrow$  speed is slightly slower than  $n \times$
- works for functions expressed as computing sum of function over training set.
- advanced optimization with logistic regression
- work for multi-core processing  $\rightarrow$  no net. latency