

# Machine Learning

Examples:

- Database Mining:

- Applications that can't be programmed
- Self-customizing programs
- Understanding human learning.

## Supervised Learning

- Data set with right answers
- Regression problem - continuous data output
- Classification problem - discrete valued output

## Unsupervised Learning

- Data has no labels or right answers  
e.g. (clustering algorithms)  
(associative)



# Linear Regression with one variable

- Also known as "univariate linear regression".

Hypothesis function:

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

- Cost Function: accuracy of  $h_{\theta}(x)$  related to data

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

$\frac{1}{2}$  the mean of the squares      mean of

- Also known as: squared error function <sup>the squares.</sup>  
- Mean squared error.

- Gradient Descent

- We take the derivative of our cost function, it will give us the direction to move towards.

$\alpha$  = learning rate.

- We make steps down by  $\alpha$ .

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

- For a Linear Regression

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

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

- if uses all the data in the sum it is also called "Batch" gradient Descent.

## Multiple Features

SON MENTIRAS!!  
↓ TODO!!!

②

- $n$ : number of features
- $x^{(i)}$ : input (features) of  $i^{\text{th}}$  training example. Vector of dimension  $n$ .
- $x_j^{(i)}$ : value of feature  $j$  in  $i^{\text{th}}$  training example.

- Now linear Regression looks like

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

- For convenience we define  $x_0^{(i)} = 1$

$$x = \begin{pmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{pmatrix} \quad \theta = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{pmatrix} \quad h_{\theta}(x) = \theta^T x$$

- Cost Function

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

- Gradient Descent

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

- Feature Scaling

- Make sure features are on a similar scale

- Divide given features by maximum value.  $x^{(i)}_j = \frac{x^{(i)}_j}{\max(x^{(i)}_j)}$

- Mean normalization

$$x_i \rightarrow \frac{x_i - \bar{x}_i}{x_{\max} - x_{\min}} \quad \text{range}$$

- How to choose learning rate  $\alpha$

The left graph shows a slow decrease in  $J(\theta)$  over iterations, indicating a small learning rate  $\alpha$ . The right graph shows a very large initial drop followed by oscillations, indicating an overly large learning rate  $\alpha$ . Problem is solve by choosing smaller  $\alpha$ .

Try  $\alpha$  by 4 diff. of magnitudes:  
1, .1, 0.01, 0.001

- Polynomial fitting
 
$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 \dots$$

NAAH... ASÍ NO ES  
③

take  $x, x^2, x^3$  as features, then

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

- Exact Minimization of Cost Function:

$X$  is the data with an added column of ones:  $X = \begin{pmatrix} 1 & \text{data} \end{pmatrix}$   $y = \begin{pmatrix} \dots \end{pmatrix}$

$x_0$	size	# bedrooms	Price
1	-	-	-
1	-	-	-
⋮	⋮	⋮	⋮

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

this are the values of  $\theta$  that minimize  $J(\theta)$ .

$$X = \begin{pmatrix} x^{(1)T} \\ x^{(2)T} \\ \vdots \\ x^{(m)T} \end{pmatrix}$$

- Gradient Descent vs Normal Equation

For  $n$  big # of features normal equation scales as  $O(n^3)$   $n \approx 10,000$  ~~too slow~~  
 somewhat the threshold for implementing gradient descent.

(4)

- Normal Equation:

what if  $X^T X$  is non-invertible?

- This can be because of redundant features. (linearly dependent)

e.g.  $x_1 = \text{size in feet}$

$x_2 = \text{size in m}^2$

- Also if there are more features than data point.  $m \leq n$

- Delete some features or use regularization.

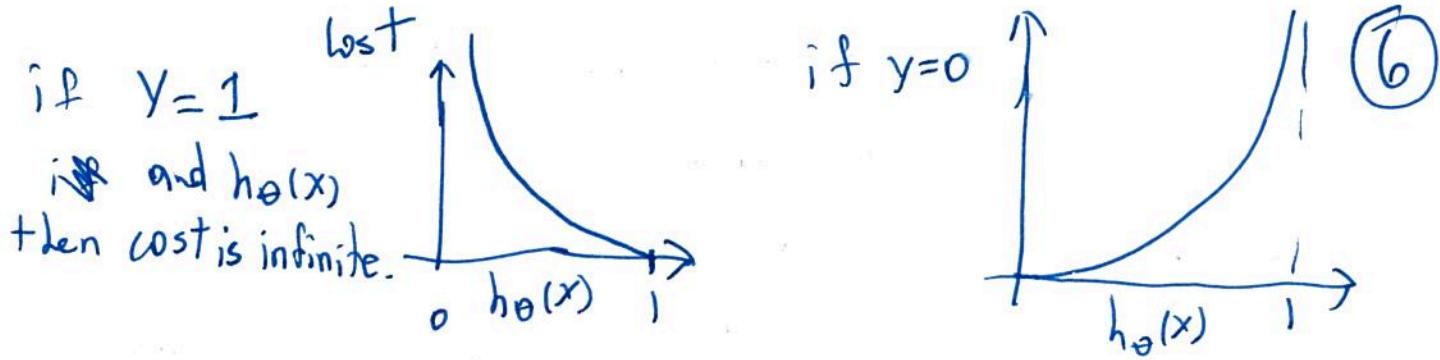
## Classification:

-  $y \in \{0, 1\}$  Two classes e.g. 0: Negative class  
1: Positive class

- Linear regression not a good idea.

$h_\theta(x)$  can be  $> 1$  and  $< 0$ !!

- Logistic Regression:  $0 \leq h_\theta(x) \leq 1$



- We write cost( $h_\theta(x)$ ,  $y$ ) as

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

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

- We want to minimize  $J(\theta)$

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

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

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

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

- Optimization Algorithm

- Gradient Descent

- Conjugate Gradient

- BFGS

- L-BFGS

• Function is `fminunc(theta, options);`

`[Jval, gradient] = function(t, theta)`

`function [Jval, gradient] = function(t, theta)`

`Jval = J(theta)`

`gradient(1) = 2 * J(theta)`

`gradient(n+1) = 2 * sum(J(theta))`

# Logistic Regression Model

(5)

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

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

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

Logistic function - sigmoid function

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

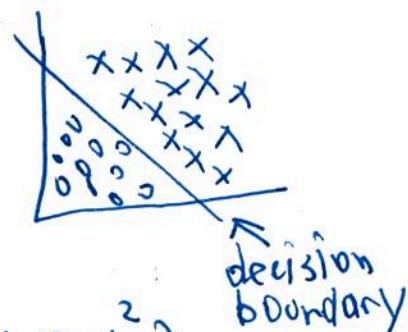
Interpretation:  $h_{\theta}(x) = 0.7 \Rightarrow 70\% \text{ chance}$

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

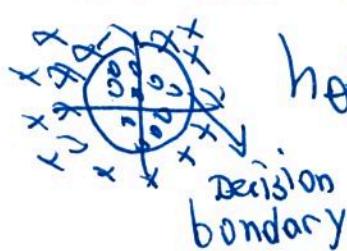
$g(z) \geq 0.5 \text{ when } z \geq 0$

$h_{\theta}(x) \geq 0.5 \text{ when } \theta^T x \geq 0$

- Decision Boundary:  $h_{\theta}(x) = 0.5$   
 $\theta^T x = 0$



- Non-linear decision boundaries



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

Adding polynomial terms can change shape of boundary.

- Higher order polynomial terms: can have more complex decision boundaries.

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2 + \theta_5 x_1 x_2^2 + \dots)$$

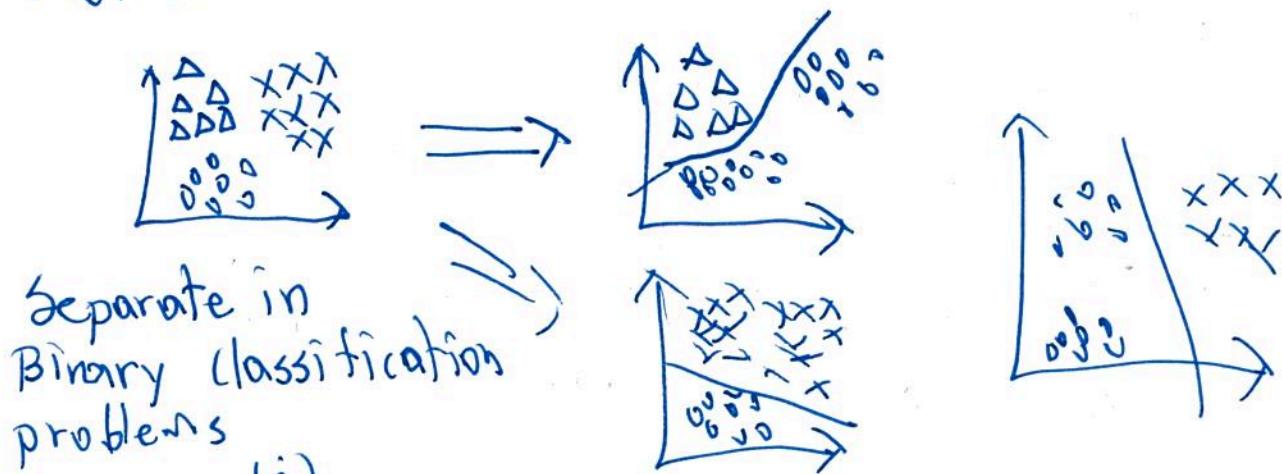
- Cost function

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

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

## - Multi-class classification: One vs All

⑦

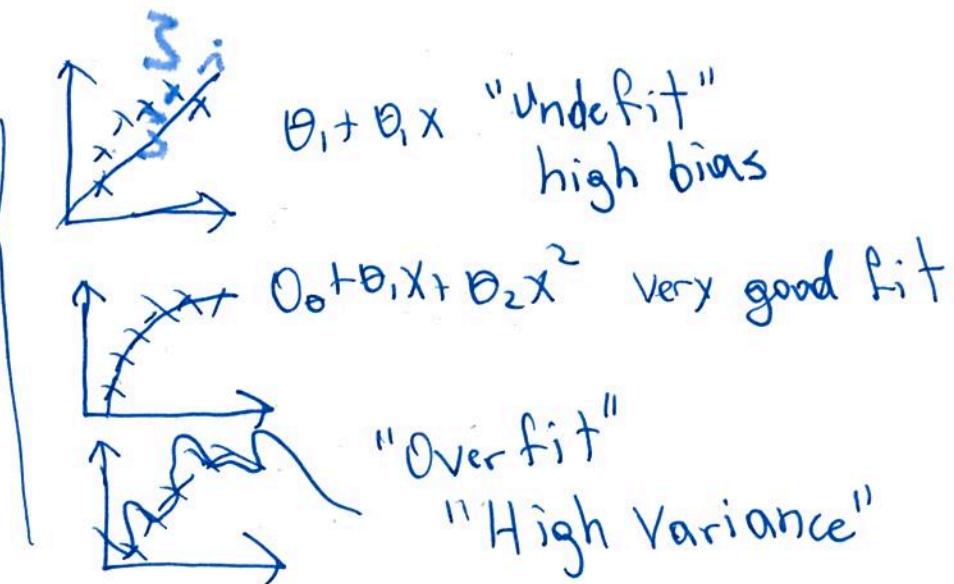


$$^{(i)} h_{\theta}(x) = P(y=i | \theta)$$

$\max_i h_{\theta}^{(i)}(x)$  To predict we choose the highest classification output.

## - Overfitting

Linear  
Regression



- It can also happen for logistic Regression,

## - Addressing Overfitting:

- Reduce # of features
  - Manually select which features to keep.
  - Model selection algorithm
- Regularization
  - Reduce magnitude/values of parameters  $\theta_j$ .
  - Works well for a lot of features

- We penalize high order terms ③

- "Simpler" hypothesis

- Less prone to overfitting

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

$\lambda$ : regularization parameter

if  $\lambda$  is big "under fitting"

if  $\lambda$  is small "over fitting"

## Regularized Linear Regression

- Gradient Descent

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

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

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

- Normal Equation

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

- Non-Invertible

if  $m \leq n$   $X^T X$  is singular

if  $\lambda > 0$  the matrix  $X^T X + \lambda I_n$  is invertible.

## Regularized Logistic Regression

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

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

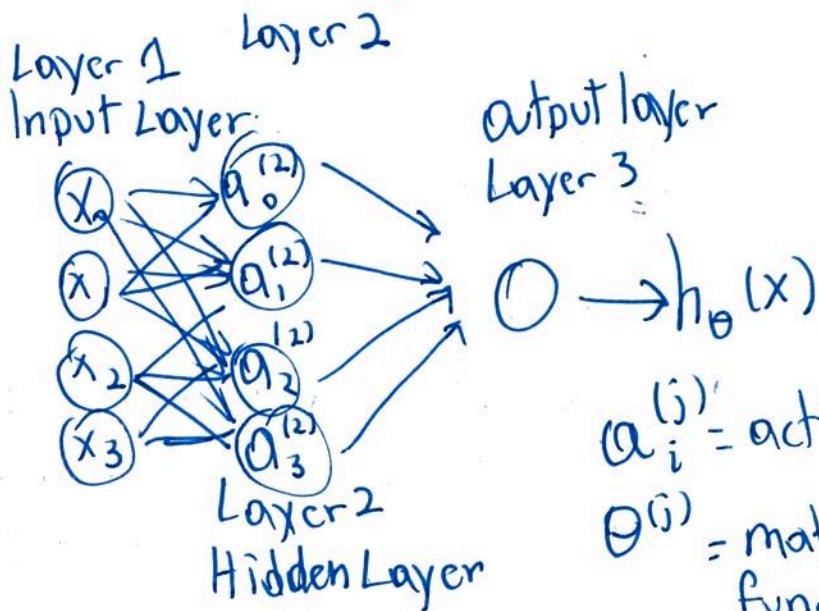
# Neural Networks

⑨

- Non-linear hypotheses.

- With  $n$  large just too many features  
e.g. Pixel examples, recognize object from image.

Neural network  $\Rightarrow$



$a_i^{(j)}$  = activation of unit  $i$  layer  $j$   
 $\Theta^{(j)}$  = matrix weights controlling function mapping from layer  $j$  to layer  $j+1$ .

3 input units

3 hidden units

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

- if network has  $s_j$  units in layer  $j$ ,  $s_{j+1}$  units in layer  $j+1$ , then  $\Theta^{(j)} \in \mathbb{R}^{s_j \times s_{j+1}}$

$$a_1^{(2)} = g[\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3]$$

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

# ⑩

## Forward Propagation: Vectorized Implementation

$$a_1^{(2)} = g(z_1^{(2)}) \quad z_1^{(2)} = \theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3$$

$$X = \begin{pmatrix} x_0 \\ x_1 \\ \vdots \\ x_3 \end{pmatrix} \quad z^{(2)} = \begin{pmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{pmatrix} \text{ then } a^{(2)} = g(z^{(2)})$$

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

Use  $a^{(1)}$  to define activations then  $z^{(2)} = \theta^{(1)} a^{(1)}$

Extra bias unit  $a_0^{(2)} \rightarrow a^{(2)} \in \mathbb{R}^4$

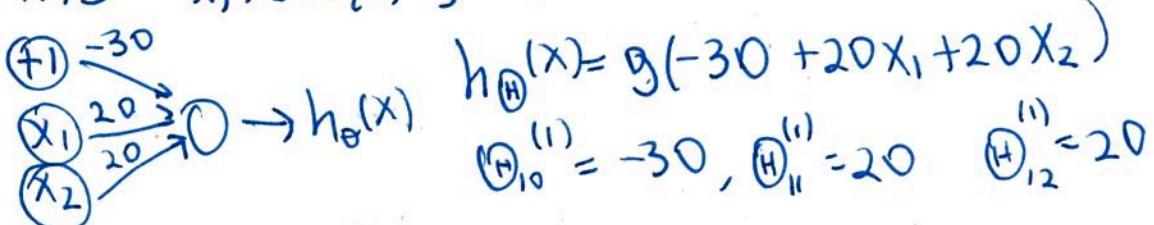
$$\text{Add } a_0^{(2)} = 1 \quad z^{(3)} = \theta^{(2)} a^{(2)}$$

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

- same as logistic regression but using updated  $a_i^{(j)}$  features, the  $a$  are the learned features from the initial features  $x_i$ .

Example:

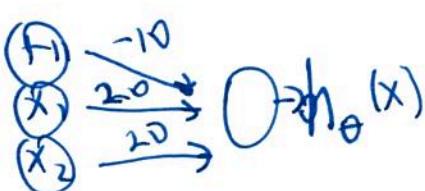
AND  $x_1, x_2 \in \{0, 1\}$   $y = x_1 \text{ AND } x_2$



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

$h_\theta(x) \approx x_1 \text{ AND } x_2$

OR function



$x_1$	$x_2$	$h_\theta(x)$
0	0	0
0	1	1
1	0	1
1	1	1

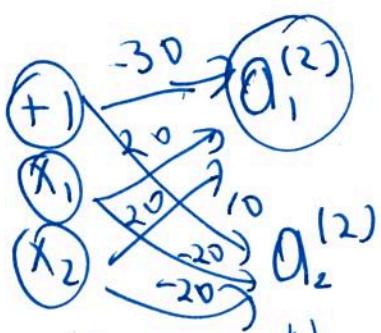
NOT

$$\begin{array}{c} (+) \\ (x_1) \end{array} \xrightarrow{\begin{array}{l} 10 \\ -20 \end{array}} \rightarrow h_{\theta}(x) \quad \begin{array}{c|cc} x & h_{\theta}(x) = g(10 - 20x_1) \\ 0 & g(10) = 1 \\ 1 & g(-10) = 0 \end{array}$$

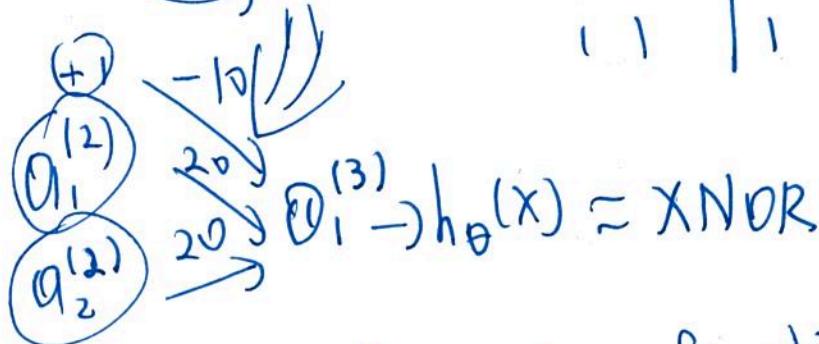
(11)

XNOR

$x_1$	$x_2$	XNOR
0	0	1
0	1	0
1	0	0
1	1	1



$x_1$	$x_2$	$O_1^{(2)}$	$O_2^{(2)}$	$h_{\theta}(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1



Multi-Class classification

Example: Recognize a number 10 classes, each number a class.  $h_{\theta}(x) \approx \begin{pmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$  first class

$$h_{\theta}(x) \in \mathbb{R}^4 \quad 4 \text{ classes} \quad h_{\theta}(x) \approx \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} \text{ second class}$$

Multiple output units: One vs all

Training set:  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$

$y^{(i)}$  one of  $\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$

$(x^{(i)}, y^{(i)})$   $h_{\theta}(x^i) \approx y^{(i)}$   
Image object Vector representing a class.

(12)

# Neural Network Classification

- Binary classification: Training set

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

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

$$h_{\theta}(x) \in \mathbb{R}$$

$$S_L = 1$$

$$K = 1$$

no. of layers  $\uparrow$   
of output units

- Multi-class classification ( $K$  classes)

$$y \in \mathbb{R}^K \quad \text{e.g. } \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

$$K \text{ output units} \quad h_{\theta}(x) \in \mathbb{R}^K \quad S_L = K \quad K \geq 3$$

- Cost Function: generalization of logistic regression.

Logistic regression  $\Rightarrow$

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

Neural Network  $\Rightarrow$

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

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

(13)

## Gradient Computation

- we want to minimize  $J(\theta)$
- we need to compute

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

## Backpropagation algorithm

$\delta_j^{(l)}$ : "error" of node  $j$  in layer  $l$ .

activation  $a_j^{(l)}$  e.g.  $\delta_j^{(4)} = a_j^{(4)} - y_j$

Vectorization:  $\delta_j^{(4)} = a^{(4)} - y$   $\nabla h_\theta(x)_j$

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

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

\* No  $\delta^{(1)}$ \*

$$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_j^{(l)} \delta_i^{(l+1)} \quad \begin{matrix} \text{(ignoring } \lambda; \text{ if } \\ \lambda=0 \end{matrix}$$

- For many training sets ( $m$ )

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

for  $i = 1$  to  $m$

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

Perform forward propagation to compute  $a^{(l)}$

for  $l = 2, 3, \dots, L$ .

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

Then back propagation  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)} \quad \Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} / (a^{(l)})^T$$

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

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \quad \begin{matrix} j = 0 \\ i = 0 \end{matrix}$$

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

$$\text{Formally } \hat{\theta}_j^{(1)} = \frac{\partial}{\partial z_j^{(1)}} \text{Cost}(i) \quad j \geq 0$$

(14)

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

## ~~Advanced Optimization~~ Advanced Optimization $\mathbb{R}^{n+1}$

need  $\mathbb{R}^{n+1}$

function [val, gradient] = costFunction(theta)

optTheta = fminunc(@costFunction, initialTheta, options)

Neural Network ( $L=4$ )

$\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$  matrices

$D^{(1)}, D^{(2)}, D^{(3)}$  matrices

"Unroll into vectors"

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

DVec = [D1(:,1), D2(:,1); D3(:,1)];

Algorithm ↗

- Initial parameters  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$

- Unroll to get initialTheta

- Use fminunc(@costFunction, initialTheta, options)

- function [val, gradient] = costFunction(thetaVec)

- from thetaVec reshape to get  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$

- Use forward propagation/back propagation to compute  $D^{(1)}, D^{(2)}, D^{(3)}$  and  $J(\theta)$

- Unroll  $D^{(1)}$  to get gradientVec.

- Gradient Checking

$$\frac{d}{d\theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} \quad \epsilon = 10^{-4}$$

$\theta \in \mathbb{R}^n$  unrolled version of  $\theta$

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

for  $i = 1:n$

$$\theta_+ = \theta$$

$$\theta_+(i) = \theta_+ + \epsilon;$$

$$\theta_- = \theta$$

$$\theta_- = \theta_- - \epsilon;$$

$$\bullet \text{grad}(i) = \frac{J(\theta_+) - J(\theta_-)}{2\epsilon}$$

end

\* Once you check gradient turn off!! when training

- Initial value of  $\theta$ :

- if initial  $\theta$  are zeros then don't work as ~~you can't calculate~~  $\theta_{ij}^{(1)}$

Even after one gradient decent update.

$$\theta_{01}^{(1)} = \theta_{02}^{(1)}$$

The hidden units are identical.

$$\begin{aligned}\theta_{ij}^{(2)} &= \theta_{ij}^{(1)} & f_i^{(2)} &= f_i^{(1)} \\ \frac{\partial}{\partial \theta_{01}} J(\theta) &= \frac{\partial}{\partial \theta_{02}} J(\theta)\end{aligned}$$

- Use random initialization

$$\Theta_{01} = \text{rand}(10,11) * 12\epsilon - \epsilon$$

$$\theta_{ij}^{(1)} \in [-\epsilon, \epsilon]$$

Symmetry Breaking

- How to choose architecture: Default 1 hidden layer. If more hidden some no. of hidden units in every layer.

-  $J(\theta)$  non-convex

Check that

~~grad~~ grad  $\approx$  Dvec

from back propagation

(15)

# Debugging a learning algorithm:

(16)

- Suppose the solution makes errors large in its predictions.

- Get more training examples
- Try smaller sets of features (prevent overfitting)
- Try getting additional features.
- Try polynomial features ( $x_1^2, x_2^2, x_1 x_2$ , etc.)
- Try decreasing  $\lambda$
- Try increasing  $\lambda$

## Machine learning Diagnostic:

test to gain insight what is/isn't working and gain guidance as to how to improve.

- Evaluate your hypothesis

- Split Data set and use 70% as training set and 30% as test set.

$$m_{\text{test}} = \# \text{ of test examples.}$$
$$J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

• Compute test error

For Logistic Regression  $J_{\text{test}}(\theta) = -\frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} x^{(i)}_t \log(h_{\theta}(x^{(i)}_t)) + (1 - h_{\theta}(x^{(i)}_t))^2$

Linear Regression

Another way is to check the 0/1 misclassification error:

$$\text{err}(h_{\theta}(x)) = \begin{cases} 1 & \text{if } h_{\theta}(x) \geq 0.5, y=0 \\ 0 & \text{if } h_{\theta}(x) \leq 0.5, y=1 \\ 0 & \text{otherwise} \end{cases}$$

- Model Selection

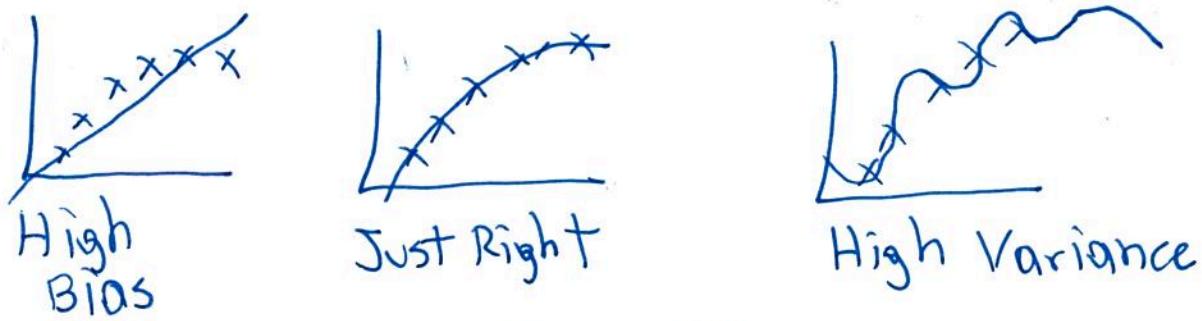
d: degree of polynomial

- check performance of each polynomial  $J_{\text{test}}(\theta)$  and choose the minimum.

- Separate Data set even more. Training set 60%  
Cross Validation set (CV) 20%  
Test set 20%

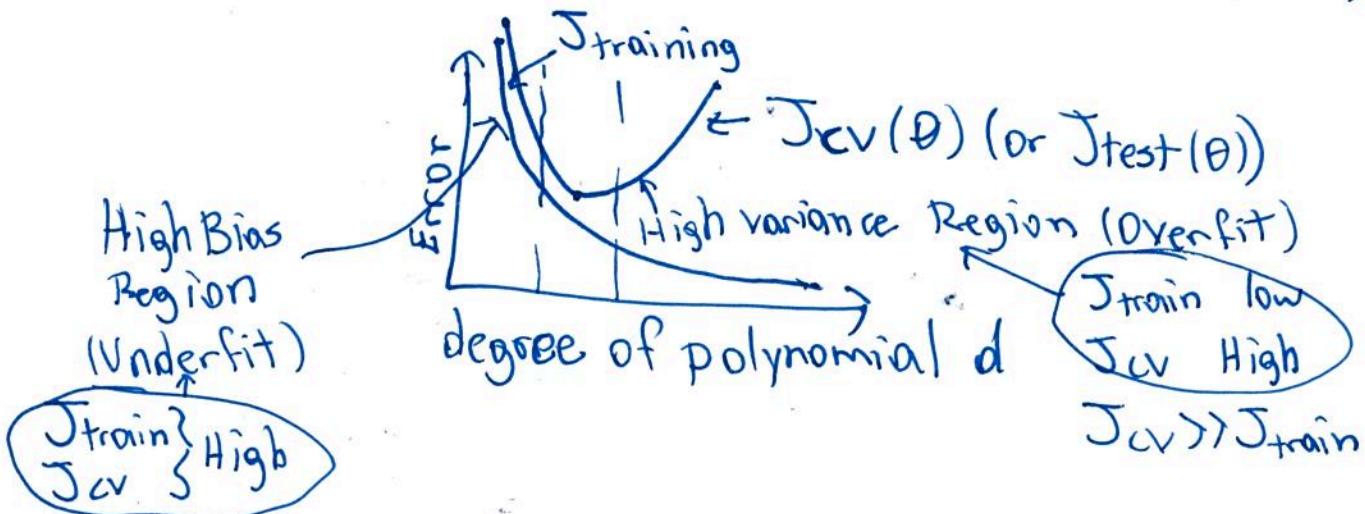
# Diagnostic Bias vs Variance

(17)



Training Error:  $J_{\text{train}}(\theta)$

Cross Validation Error:  $J_{\text{cv}}(\theta)$  (or  $J_{\text{test}}(\theta)$ )



Regularization effect to High Bias and High Variance

- Large  $\lambda \rightarrow$  High Bias (underfit)
- Small  $\lambda \rightarrow$  High Variance (overfit)
- choosing  $\lambda$ :  $J_{\text{train}}(\theta)$ ,  $J_{\text{cv}}(\theta)$ ,  $J_{\text{test}}(\theta)$  have no regularization parameter ( $\lambda=0$ )
- Range of values for  $\lambda$  e.g.  $[0, 10]$

$$\lambda_1 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^1 \rightarrow J_{\text{cv}}(\theta^1)$$

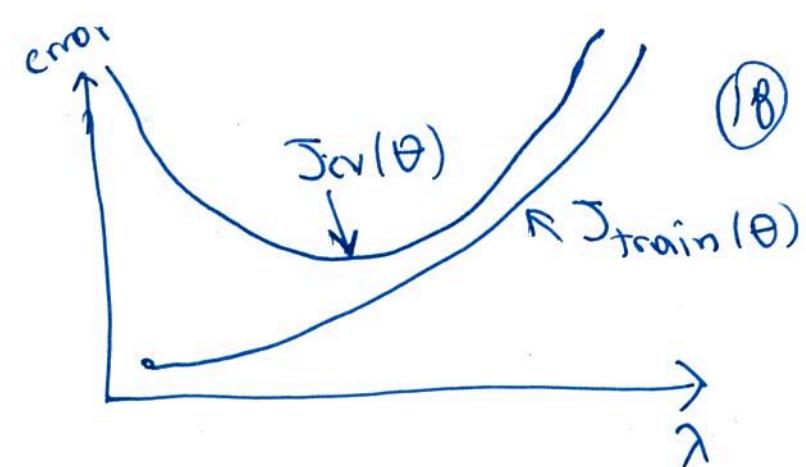
$$\lambda_2 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^2 \rightarrow J_{\text{cv}}(\theta^2)$$

$$\lambda_3 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^3 \rightarrow J_{\text{cv}}(\theta^3)$$

choose  $\lambda$ ; that gives lowest  $J_{\text{cv}}(\theta^i)$

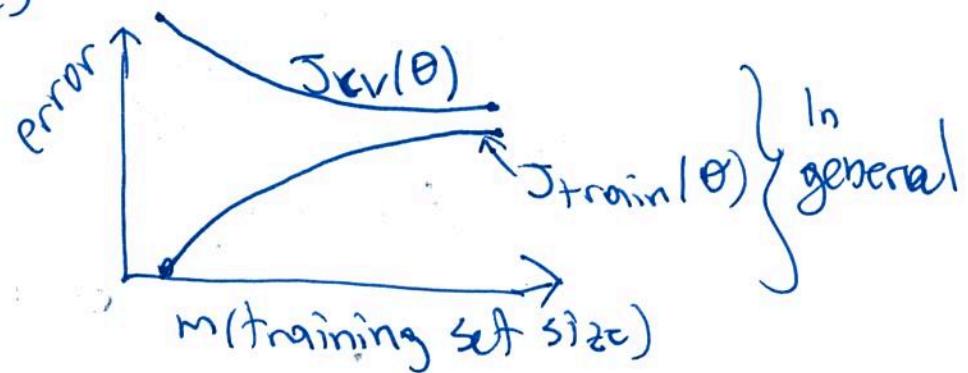
then Test error:  $J_{\text{test}}(\theta^{i*})$

$J(\theta) \leftarrow$  regularization  $\lambda$   
 $J_{\text{train}}(\theta) \rightarrow \lambda = 0$  to get  
 $J_{\text{CV}}(\theta) \rightarrow \lambda = 0$  to get

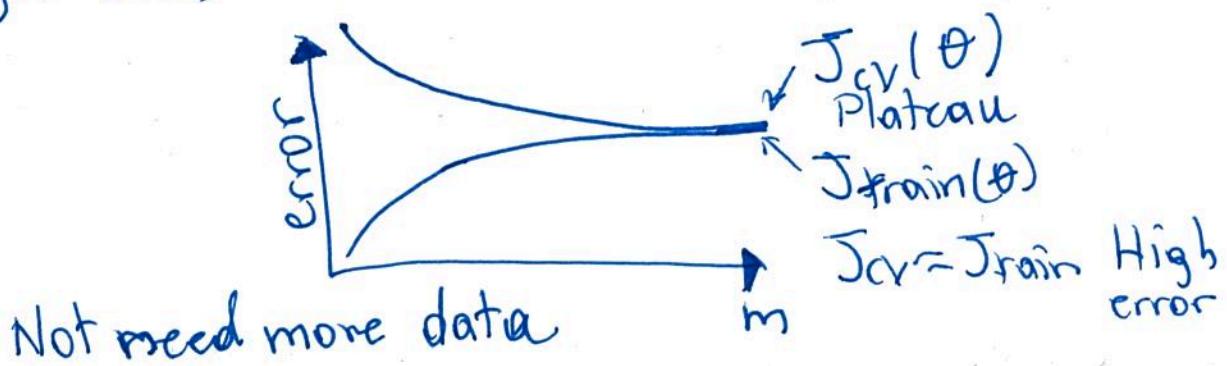


## Learning Curves

$J_{\text{train}}$   
 $J_{\text{CV}}$

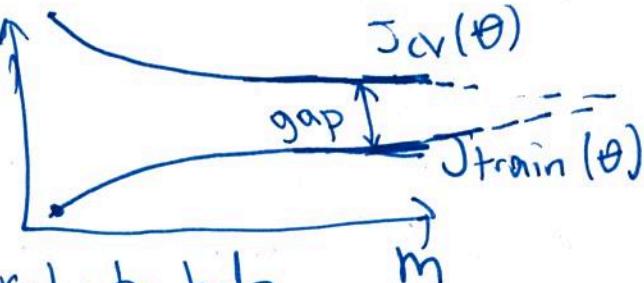


### - High Bias



### - High Variance

Not need more data  
getting more training data is likely to help.



- Get more training examples } fixes High Variance
- Try smaller set of features } fixes High Variance
- Try get additional features } fixes High Bias
- Try polynomial features
- Decreasing  $\lambda \rightarrow$  fixes High Bias
- Increasing  $\lambda \rightarrow$  fixes High Variance

# Neural Networks and overfitting

(19)

- Small NN underfits
- Large NN overfits, change  $\lambda$ .

Example  $\rightarrow$  Spam Classifier

Features are certain words e.g.  
Suppose email contains

words  $w_1, w_2, w_4$   
 $x_j = \begin{cases} 1 & \text{if word } j \\ 0 & \text{otherwise} \end{cases}$

$$x = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

buy  
deal  
discount  
now

- what is the best way to spend time?

- Collect lots of data
- Develop sophisticated features (routing information)
- Email message body and develop sophisticated features.
- Develop sophisticated algorithms to detect misspellings.

## Error Analysis

- Start with simple algorithm and test in cross validation data
- Plot learning curves  Error analysis
- Manually examine the examples (in cross validation set) that your algorithm made errors.

## Error Analysis Example

$M_{cv} = 500$  cross validation set

Algorithm misclassifies 100 emails

Manually examine 100 errors and categorize them

- i) what type of email is it? pharma, replicas, fishing

~~Pharma~~: 12

Replia: 4

fishing: 54  $\rightarrow$  Doing Poorly

Other: 31

- ii) what cues (features) you think would have helped the algorithm classify them

Deliberate Misspelling: 5  $\rightarrow$  rare not worth the time

Unusual Routing: 16

Unusual Punctuation: 32  $\rightarrow$  strong indication for more features for punctuation.

- iii) Numerical evaluation - single row number

e.g. Treat discount/discounts/discounted) ---  
be treated as same word?

(can use "stemming" program software)

Negative: universe/university

- Look CV error for with and without stemming

e.g. Distinguish upper vs lower case (Mom/man)

# Error Metrics for Skewed Classes

(21)

e.g. cancer classification

$h_\theta(x)$

Find 1% error but only 0.5% have cancer

function  $y = \text{predict}(\text{cancer}(x))$  ↗ this predicts better  
 $x=0$  ↗ 0.5% error.  
return

- When we have much more examples from one class than the other.

- Precision / Recall

$y=1$  in presence of rare class.

Compute:

		Actual Class	
		Predicted Class	
Predicted Class	True Positive	False Positive	True Negative
	False Negative	True Positive	

- Precision (of all patients where we predicted  $y=1$  what fraction has cancer?)

$$\frac{\# \text{True Positives}}{\# \text{predicted Positive}} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

- Recall (of patients that have cancer, what fraction did we correctly detect as having cancer?)

$$\frac{\text{True Positives}}{\# \text{actual Positives}} = \frac{\text{True Positives}}{\text{True Positives} - \text{False Negatives}}$$

- Then if algorithm predicts  $y=0$  then recall = 0 so not good.

# Trading Off precision and recall

(22)

Logistic regression:  $0 \leq h_{\theta}(x) \leq 1$

- Suppose we want to predict cancer only if very confident.

- change threshold  $h_{\theta}(x) \begin{cases} 1 & h_{\theta}(x) \geq 0.7 \\ 0 & h_{\theta}(x) < 0.7 \end{cases}$

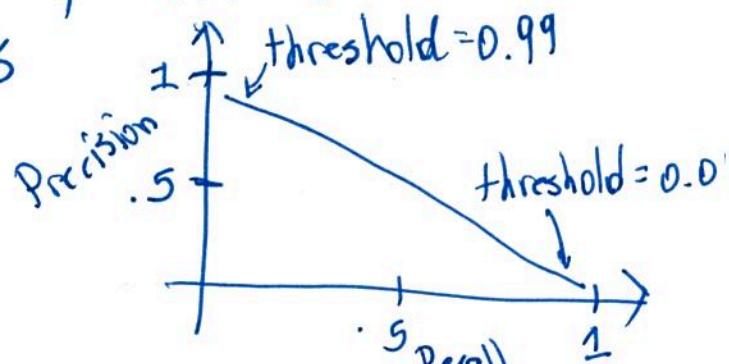
\* Higher Precision, lower recall

- Avoid false negatives, missing too many cases of cancer.

$h_{\theta}(x) \begin{cases} 1 & h_{\theta}(x) \geq .3 \\ 0 & h_{\theta}(x) < 0.3 \end{cases}$

\* Higher Recall, lower precision.

for most classifiers



- How to decide best threshold

	Precision	Recall	Which is the best?
Alg 1	.5	.4	
Alg 2	.7	.1	
Alg 3	.02	.1	

- Average  $\{.45, .4, .51\}$  \* Not Good

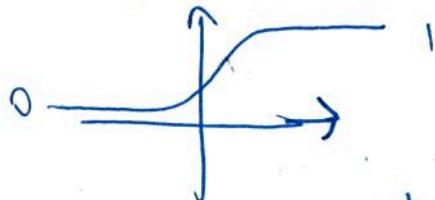
- F<sub>1</sub> Score  $\frac{2PR}{P+R}$   $\{0.444, 0.175, 0.0392\}$   
(check in cross Validation set)

# Support Vector Machine (SVM)

(23)

## - Optimization Objective

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

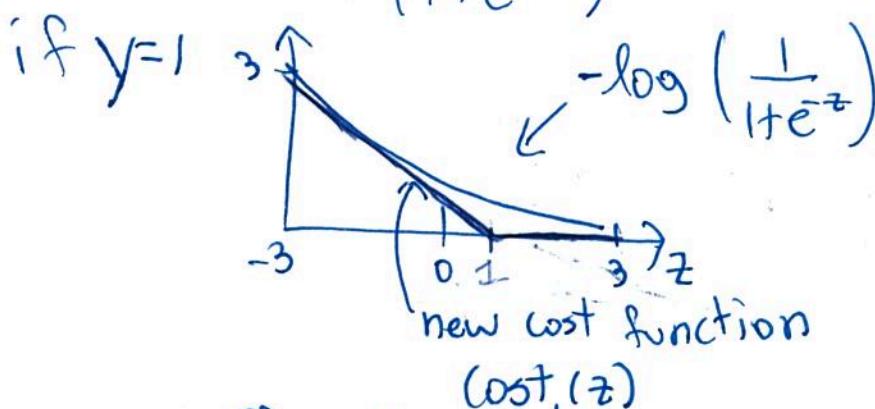


$$z = \theta^T x$$

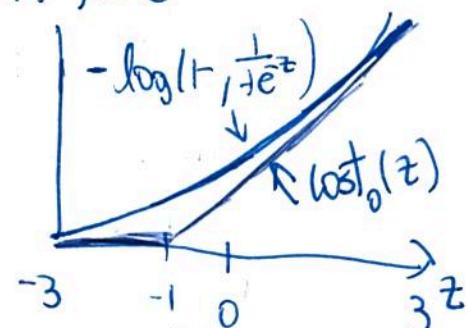
if  $y=1$ , we want  $h_{\theta}(x) \approx 1$ ,  $\theta^T x \gg 0$   
 if  $y=0$ , we want  $h_{\theta}(x) \approx 0$ ,  $\theta^T x \ll 0$

Cost of example  $(x, y)$

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



$$\text{if } y=0$$



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

By convention get rid of  $1/m$

Normally  $A + \lambda B$  (logistic regression) } same solution  $\theta$   
 now  $C + \lambda B$  for SVM } if  $C = 1/\lambda$   
 ↑ weight

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

- Hypothesis:

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

(24)

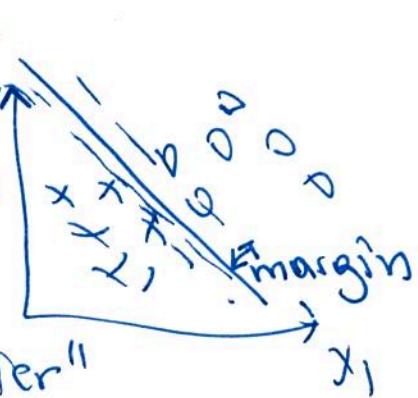
- if  $y=1$  we want  $\theta^T x \geq 1$  (not just  $\geq 0$ )
  - if  $y=0$  we want  $\theta^T x \leq -1$  (not just  $< 0$ )
- adds safety margin factor

$\cdot C$  very large:

e.g. linearly separable case

Tries to separate from the data as big as possible

"Large margin classifier"



## SVM Decision Boundary

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

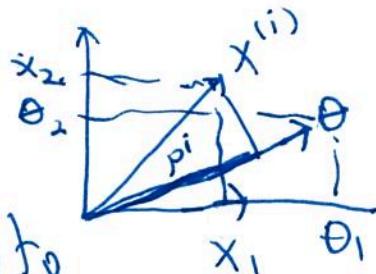
$$\theta^T x^i = ?$$

$$= p_i \|\theta\|$$

and minimizing  $p_i$

$x^i$  and  $\theta$  are perpendicular, or try to

be. So SVM tries to minimize the projections of  $x^{(i)}$  into  $\theta$



SVM tries that decision boundary is as perpendicular as possible to  $x^i$  vectors.

The separation margin are the projections so if projections are big  $\|\theta\|$  can be smaller.

# Kernels

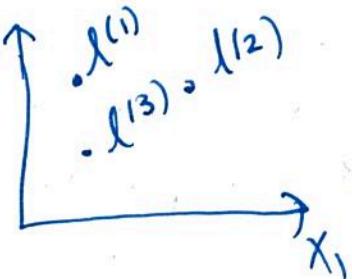
(25)

## Non-Linear Decision Boundary

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

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \dots \text{ Polynomial}$$

Is there a better choice?  $x_2 \uparrow$



Given  $x$ :

new features  $f_1 = \text{similarity}(x, l^{(1)})$   
 $= \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$

$f_2 = \text{similarity}(x, l^{(2)})$

Kernel (Gaussian kernel)  
 function

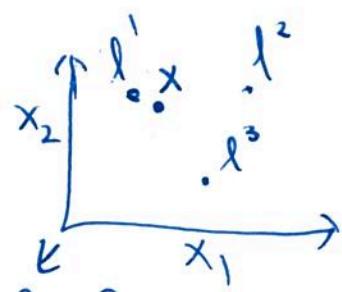
$\text{if } x \approx l^{(1)} : f_1 \approx 1$

$\text{if } x \text{ is far from } l^{(1)} : f_1 \approx 0$

- Predict  $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$

$\text{e.g. } \theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$

$\theta_0 + \theta_1 \approx -0.5 \geq 0 \quad \leftarrow f_1 \approx 1 \quad f_2 \approx 0 \quad f_3 \approx 0$



- For point close to  $l^{(1)}$  and  $l^{(2)}$  predict "1" and far predict "0"  $\sigma$  defines how far or closer.

- Choosing landmarks: where?

SVM with Kernels:

Location of landmarks where examples are.

$$(x^i, y^i) : x \rightarrow f^i = \begin{bmatrix} \downarrow & \uparrow \\ \vdots & \vdots \\ \downarrow & \uparrow \end{bmatrix} \quad x^{(i)} \in \mathbb{R}^{n+1} \quad f^{(i)} = \begin{bmatrix} f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \in \mathbb{R}^{m+1}$$

Given  $X$  compute  $f \in \mathbb{R}^{m+1}$  (n landmarks)

(26)

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

$\theta \in \mathbb{R}^{m+1}$

$$\min_{\theta} C \sum_{i=1}^n y^{(i)} \text{cost}_1(\theta^T f^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T f^{(i)})$$

$$\theta^T X \rightarrow \theta^T f^{(i)}$$

$$+ \frac{1}{2} \sum_{j=1}^m \theta_j^2$$

\* Kernels with logistic regression very slow.

\* Use packages.

- $C$  Large  $c$ : Lower Bias, high variance
- $C$  small  $c$ : High bias, low variance.
- Large  $\sigma^2$  High bias features vary smoothly
- Small  $\sigma^2$  High variance features less smoothly
- ~~Large  $\sigma^2$~~

Use SVM

Software package (e.g. liblinear, libsvm, ...) to solve for  $\theta$

→ choose parameter  $C$ .

→ choice Kernel (similarity function)

e.g. no kernel ("Linear Kernel")

→ In case Gaussian Kernel need to choose  $\sigma^2$

\* Perform feature scaling before!!

- Similarity function needs to satisfy Mercer's Theorem to make sure SVM not diverges.

- Polynomial Kernel:  $k(x, l) = (x^T l + \text{constant})^{\text{degree}}$   
use when data is always positive.

- String Kernel, histogram Kernel

# SVM for Multi-class Classification

(27)

- Package
- One vs All

Logistic Regression vs SVM ↴

$n = \# \text{ of features}$

$m = \# \text{ of training examples}$

- if  $n$  is large (relative to  $m$ ):  $n \gg m$

↳ Use logistic Regression or SVM with linear kernel

- if  $n$  is small,  $m$  is intermediate

↳ ~~SVM~~ SVM with Gaussian kernel

- if  $n$  is small,  $m$  is large

SVM will struggle. ~~(SVM)~~

Create/add more features, then use logistic Regression or SVM without Kernel.

- Neural network will work well for most but

may be slower to train.

- SVM (convex function) finds ~~global~~ local minima.



# Unsupervised Learning

(28)

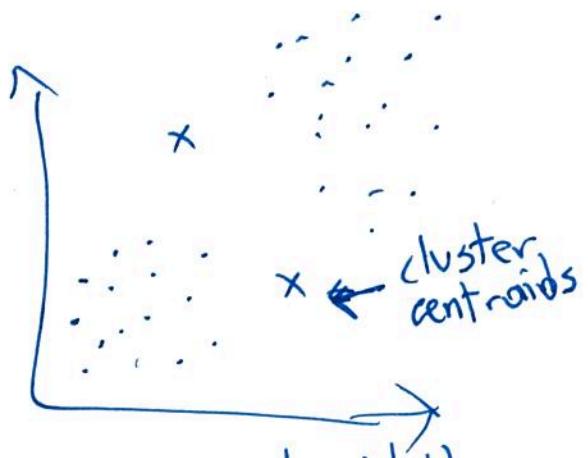
- Clustering: Market Segmentation  
e.g. Social Network analysis  
Astronomical data

## K-means Algorithm

- Iterative algorithm
- Separates data to closest centroids.
- Calculates mean position of separated cluster and moves centroid there.

repeat

    Converges to a centroid point.



- Input:
  - $K$  (numbers of clusters)
  - Training set (unlabeled)  
 $x^{(i)} \in \mathbb{R}^n$
- Random  $K$  cluster centroids  $M_1, M_2, \dots, M_K \in \mathbb{R}^n$
- Repeat {
  - for  $i=1$  to  $n$  } \* Cluster assignment step
    - $c^{(i)} := \underset{\substack{\text{cluster centroid closest to } x^{(i)}}}{\text{index}}$
    - $\underset{\substack{(1 \text{ to } K) \\ \text{index}}}{\min_k} \underset{\substack{\text{index} \\ k}}{\|x^{(i)} - M_k\|} = c^{(i)}$
  - Move centroid } for  $k=1$  to  $K$ 
    - $M_k := \text{average of points assigned to cluster } k.$

\* If no element was assigned to a cluster centroid  
eliminate or randomly assign.

$$M_2 = \frac{1}{3} [x^{(1)} + x^{(5)} + x^{(6)}] \in \mathbb{R}^n$$

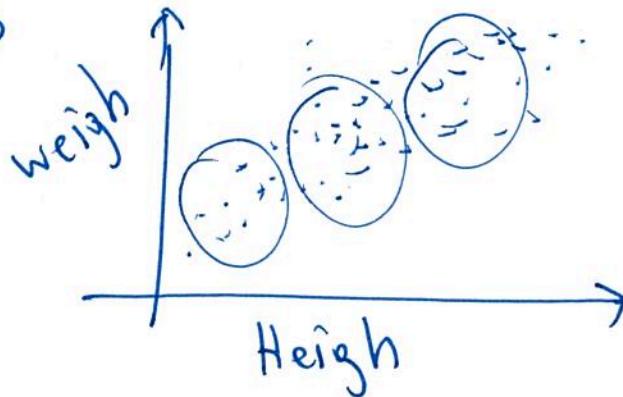
e.g. if  $x^{(1)}, x^{(5)}, x^{(6)}$  are in  $K=2$

## K-means for non-separated clusters

(29)

e.g. shirt sizes

- will kind of separate try to data.
- Market segmentation example.



## Optimization Objective

$c^{(i)}$  Index of cluster for example  $x^{(i)}$

$M_k$  cluster centroid k

$M_{c^{(i)}}$  cluster centroid assign to example  $x^{(i)}$

$$J(c^{(1)}, \dots, c^{(m)}; M_1, \dots, M_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - M_{c^{(i)}}\|^2$$

↑  
Distortion  
function       $\min_{\substack{(c^{(1)}, \dots, c^{(m)}) \\ M_1, \dots, M_K}} J(c^{(1)}, \dots, c^{(m)}, M_1, \dots, M_K)$

## Random Initialization

- should  $k < m$
- pick K random training examples
- set  $M_1 \dots M_k$  to these k examples

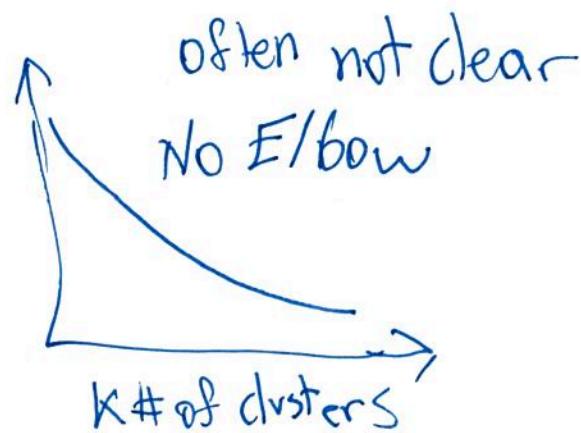
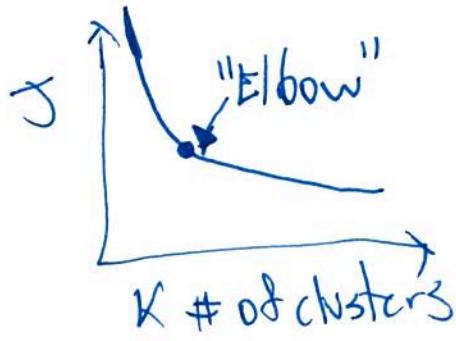
## Local Optima

- Do many Random initializations
- ~~keep minimum val~~
- Keep values that give minimum cost function (lowest distortion) for a few clusters.
- For many clusters makes little difference.

# Choosing number of clusters

(30)

## Elbow Method



## - Data Compression

Reduce data from 2D to 1D e.g. inches and cm.  
Can also be in 3D to 2D  
e.g. if Data is in a plane.

## - Data Visualization

Reduce features

## Dimensionality Reduction

PCA: Principal Component Analysis.

Tries to find surface to minimize projections of Data.

## - Preprocessing (feature scaling/mean normalization)

$$M_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

$$x_j^{(i)} = x_j - M_j$$
$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - M_j}{S_j}$$

Reduce from  $n$  to  $K$  dimensions

"covariance" matrix

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x^{(i)}) (x^{(i)})^T \quad \leftarrow \text{sigma}$$

Compute "eigenvectors" of matrix  $\Sigma$ :

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

$$U = \begin{bmatrix} u_1^{(1)} & u_1^{(2)} & \dots & u_1^{(n)} \\ \vdots & \vdots & & \vdots \end{bmatrix} \quad U \in \mathbb{R}^{n \times n}$$

$$\underbrace{U}_{\text{K vectors}} \Rightarrow \underbrace{U_{\text{reduce}} = \begin{bmatrix} u_1^{(1)} & \dots & u_1^{(K)} \\ \vdots & & \vdots \end{bmatrix}}_{\in n \times K}$$

$$\Rightarrow z = U_{\text{reduce}}^T X = t \in \mathbb{R}^K \quad ; \quad x = \begin{bmatrix} x^{(1)T} \\ \vdots \\ x^{(n)T} \end{bmatrix}$$

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

$$z = U_{\text{reduce}}^T x$$

$$\Sigma = \frac{1}{n} X^T X;$$

- Reconstruction from compressed Representation

$$\text{given } z \in \mathbb{R}^K \quad X_{\text{approx}} = U_{\text{reduce}} \cdot z$$

- How to choose  $K$  dimensional features

$$\text{Average squared projection error: } \frac{1}{n} \sum_{i=1}^n \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$$

$$\text{Total variation: } \frac{1}{n} \sum_{i=1}^n \|x^{(i)}\|^2$$

(choose  $K$  to be smallest value so that

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

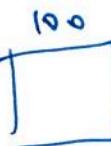
"99% of variance is retained"

(32)

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

$S = \begin{bmatrix} s_{11} & & 0 \\ & s_{22} & \\ 0 & & s_{mm} \end{bmatrix}$  for given  $k$  check can be done by  $1 - \frac{\sum_{i=1}^k s_{ii}}{\sum_{i=1}^n s_{ii}} \leq 0.01$

$$\text{or } \frac{\sum_{i=1}^k s_{ii}}{\sum_{i=1}^n s_{ii}} \geq 0.99$$

E.g.  10,000 Pixels  $x^{(i)} \in \mathbb{R}^{10,000}$

Unlabeled dataset  $x^{(1)}, x^{(2)}, \dots, x^{(m)} \in \mathbb{R}^{10,000}$

$$\downarrow \text{PCA}$$

$$z^{(1)}, z^{(2)}, \dots, z^{(m)} \in \mathbb{R}^{100}$$

New training set  $(z^1, y^1), (z^2, y^2), \dots, (z^m, y^m)$

- PCA only on training set, mapping  $x^{(i)} \rightarrow z^{(i)}$  to find parameters. Then mapping  $\hat{x}_{cv}^{(i)}$  and  $\hat{x}_{test}^{(i)}$ .

- Application of PCA

+ compression  
+ Visualization

\* Bad Use of PCA: To prevent Overfitting

Idea is  $\Rightarrow$  less features less likely to overfit.

Use regularization!!!

(33)

# Anomaly Detection

(34)

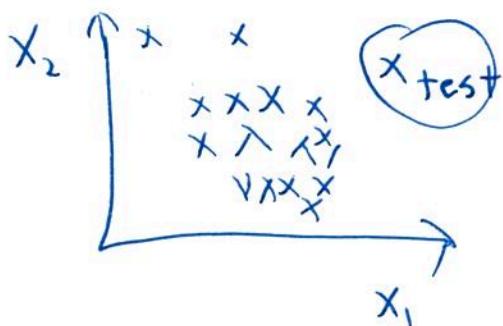
E.g. Aircraft engine features

$x_1$  = heat generated.

$x_2$  = vibration intensity.

New engine:  $x_{\text{test}}$

Is  $x_{\text{test}}$  anomalous?



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

- Model  $p(x)$

if  $p(x_{\text{test}}) < \varepsilon$  flag anomaly

if  $p(x_{\text{test}}) \geq \varepsilon \rightarrow \text{OK}$

E.g. Fraud Detection

- $x^{(i)}$  = features of user  $i$ 's activities.
- Model  $p(x)$  from data.
- Identify unusual users.

- Gaussian-Normal Distribution

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

- Density Estimation

- Training set:  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

$$x_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$$

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

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

- Anomaly Detection algorithm

1. choose features  $x_i$  that might indicate anomalous examples.

2. fit parameters  $\mu_1, \dots, \mu_n$ ,  $\sigma_1^2, \dots, \sigma_n^2$

$$\cancel{\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}} \rightarrow \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

$$3. p(x) = \prod_{j=1}^n \frac{1}{2\pi\sigma_j} \exp \left[ -\frac{(x_j - \mu_j)^2}{2\sigma_j^2} \right]$$

Anomaly if  $p(x) < \epsilon$

- Evaluating an anomaly detection system.

Assume labeled data, of anomalous and non-anomalous examples. ( $y=0$  if normal) ( $y=1$  if anomalous)

- Take training set and calculate  $\mu$  and  $\sigma$ ,  $p(x)$  only non-anomalous

- CV set and test set ~~use~~ go anomalous and non-anomalous.

- Set will be very skewed

- Evaluation metrics:

- True positive, false positive, false negative

- Precision / Recall

- F<sub>1</sub>-Score.

- Can also use CV to set  $\epsilon$ .

# Anomaly Detection vs Supervised Learning

(36)

- |   |   |
|---|---|
| <ul style="list-style-type: none"><li>- Very small # of positive examples (<math>y=1</math>). (0-20 is common)</li><li>- Large # of negative (<math>y=0</math>) examples.</li><li>- Many different "types" of anomalies.</li><li>- Future anomalies may look nothing like any of the anomalous examples.</li><li>- Fraud Detection</li><li>- Manufacturing</li><li>- Monitoring machines in a data center</li></ul> | <ul style="list-style-type: none"><li>- Large # of positive and negative examples.</li><li>- Enough positive examples for algorithm to get a sense of what positive examples look like.</li><li>- Future + examples likely to be similar to ones in training set.</li><li>- Email Spam classification</li><li>- Weather Prediction (Sunny /rainy/etc)</li><li>- Cancer Classification</li></ul> |
|---|---|

## Non-gaussian features

- Plot histogram of features.
- Play with transformations of data e.g.  $\log(x)$ ,  $\log(x+c)$ ,  $\sqrt{x}$

# Multivariate Gaussian Distribution

(37)

$\mu \in \mathbb{R}^n$ ,  $\Sigma \in \mathbb{R}^{n \times n}$   
(covariance matrix)

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

↑ determinant

- Captures correlations between features.
- Original scales better, computationally cheaper.
- Original even if  $m$  (training size) is small.
- Multi needs  $m > n$  or  $\Sigma$  is non-invertible.  
↳ vs. if  $m \gg n$
- $\Sigma$  may be non-invertible because of redundant features.

# Recommender Systems

(38)

-Content-Based

Movie	$\theta^{(1)}$ Alice (1)	$\theta^{(2)}$ Bob (2)	$\theta^{(3)}$ Carol (3)	$\theta^{(4)}$ Dave (4)	$x_1$ (romance)	$x_2$ (action)	
once at last	5	5	0	0	.9	0	$\leftarrow x^{(1)}$
romance forever	5			0	1	.01	$\leftarrow x^{(2)}$
ate Puppies		4	0		.99	0	$\leftarrow x^{(3)}$
Non Stop (or	0	0	5	4	.1	1	$\leftarrow x^{(4)}$
swords vs Karate	0	0	5		0	.9	

-Each movie represented by feature vector  $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$$-n_u = 4, n_m = 5$$

$$-n = 2 \quad x^{(1)} = \begin{pmatrix} .9 \\ 0 \end{pmatrix}$$

-For each user  $j$ , learn parameter  $\theta^{(j)} \in \mathbb{R}^3$ . Predict user  $j$  as rating movie  $i$   $(\theta^{(j)})^T x^{(i)}$  stars

-Problem Formulation

$$r(i,j) = 1 \text{ if user } j \text{ has rated movie } i \text{ (0 otherwise)}$$

$y^{(i,j)}$  = rating by user  $j$  on movie  $i$  (if defined)

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

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

For user  $j$ , movie  $i$ , predicted rating:  $(\theta^{(j)})^T x^{(i)}$

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

$$\text{To learn } \theta^{(j)}: \min_{\theta^{(j)}} \sum_{i:r(i,j)=1} [(\theta^{(j)})^T x^{(i)} - y^{(i,j)}]^2 \frac{1}{2m^{(j)}}$$

$$\theta^{(j)} \in \mathbb{R}^{n+1}$$

For recommended system convention  
get rid of  $m^{(j)}$

$$+ \frac{1}{2m^{(j)}} \sum_{k=1}^n \theta_k^{(j)}$$

## Optimization Objective

$$\min_{\theta^{(j)}} \frac{1}{2} \sum_{i: r(i,j)=1} [(\theta^{(j)})^T X^{(i)} - y^{(i,j)}]^2 + \frac{\lambda}{2} \sum_{k=1}^{n_u} (\theta_k^{(j)})^2 \quad (30)$$

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

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

Gradient Descent Update:

$$\theta_K^{(j)} = \theta_K^{(j)} - \alpha \sum_{i: r(i,j)=1} [(\theta^{(j)})^T X^{(i)} - y^{(i,j)}] X_K^{(i)} \quad (\text{for } k=0)$$

$$\theta_K^{(j)} = \theta_K^{(j)} - \alpha \left( \sum_{i: r(i,j)=1} [(\theta^{(j)})^T X^{(i)} - y^{(i,j)}] X_K^{(i)} + \lambda \theta_K^{(j)} \right)$$

\* Collaborative Filtering / feature learning

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

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

- Randomly guess  $\theta$  then  $x$  then  $\theta$

↳ Iterate  $\theta \rightarrow x \rightarrow \theta \rightarrow x \dots$

End at  $\theta$  and  $x$  quite good, converge.

Minimize  $X^{(1)}, \dots, X^{(n_u)}$  and  $\theta^{(1)}, \dots, \theta^{(n_u)}$  simultaneously

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

we don't need  $x_0$  and  $\theta_0$

$$\theta \in \mathbb{R}^n \quad X \in \mathbb{R}^{n_u \times n}$$

(40)

1- Choose Initialize  $X^{(1)} \dots X^{(n_m)}$   
 $\theta^{(1)} \dots \theta^{(n_u)}$  to small random values.

2- Minimize  $J(X^{(1)}, \dots, X^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$  using gradient descent or an advanced optimization algorithm. For every  $j = 1, \dots, n_u$  and  $i = 1, \dots, n_m$ :

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

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

3- Giving a user with  $\theta$  and a movie with (learned) features  $X$ , predict a star rating of  $\theta^T X$ .  $(\theta^j)^T (X^{(j)})$ .

### Vectorization Implementation

Movie	A	B	C	D
Love last	5	5	0	0
Romance 4E	5	?	2	0
Intepuppers	?	4	0	?
Nonstop	0	0	5	4
Karate	0	0	5	?

$$y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}_{5 \times 4}$$

$\uparrow y^{(i,j)}$

Predicted ratings:

$$\begin{bmatrix} (\theta^1)^T X^{(1)} & \dots & (\theta^{n_u})^T X^{(1)} \\ \vdots & & \vdots \\ (\theta^1)^T X^{n_m} & \dots & (\theta^{n_u})^T X^{n_m} \end{bmatrix}$$

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

"Low rank matrix factorization"

(41)

## Finding related movies

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

- How to find movie  $j$  related to movie  $i$ ?  
 $\frac{x^{(j)}}{x^{(i)}}$

If small  $\|x^{(i)} - x^{(j)}\| \rightarrow$  movie  $i$  and  $j$  are similar.

- Mean normalization

$$Y = \begin{bmatrix} 55 & 0 & ? \\ 52 & 0 & ? \\ 0 & 0 & 4 & ? \\ 0 & 5 & 0 & ? \end{bmatrix} \quad M = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ -1.25 \end{bmatrix}$$

$$Y \rightarrow \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 \\ - & - & - & - \\ - & - & - & - \\ -1.25 \end{bmatrix} \quad Y = Y - M \text{ subtract per row.}$$

$Y$  is normalized

For user  $j$  non movie  $i$  predict

$$(\theta^{(j)})^T (X^{(i)}) + M_i$$

Then for a user that hasn't rated anything

$$\theta = (0) \quad (\theta^{(s)})^T (X^i) + M_i = M_i$$

# Large Datasets

(42)

## \* Stochastic Gradient descent

Note: Batch gradient descent (all data needed)

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

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

1: Randomly shuffle dataset

2: Repeat { for  $i = 1, \dots, m$  }

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

- Small steps for each example.

- Ends up around the global minimum.

## \* Mini-Batch Gradient Descent

~~Batch gradient~~

Use  $b$  examples in each iterations

$$\text{e.g. if } b = 10 \quad \theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=1}^{10} (h_{\theta}(x^{(k)}) - y^{(k)}) x_j^{(k)}$$

Repeat { for  $i = 1, 11, 21, 31, \dots, 991$  }

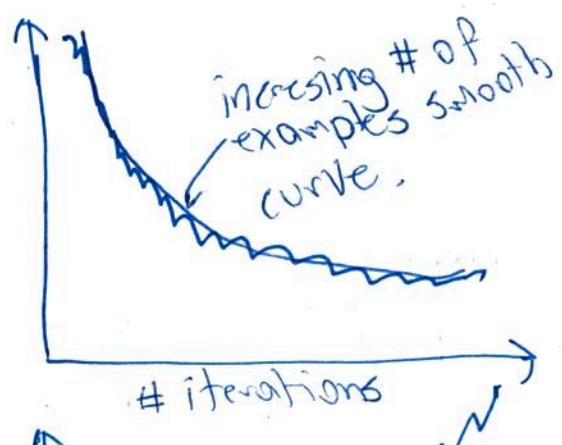
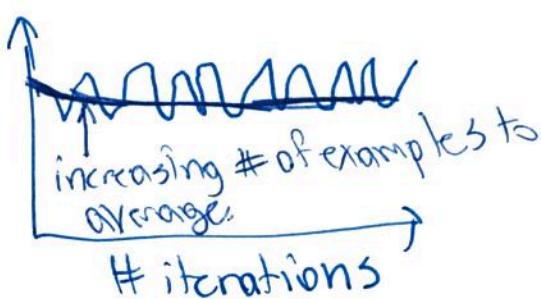
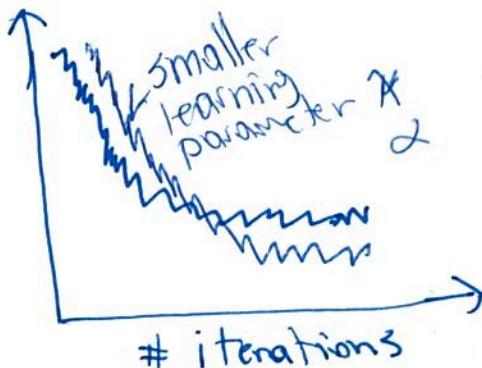
$$\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+10} (h_{\theta}(x^{(k)}) - y^{(k)}) x_j^{(k)}$$

- Use vectorize way  $\nearrow$  sometimes faster than both.

- Convergence

- Compute  $\text{cost}(\theta, (x^i, y^i))$  before updating  $\theta$  using  $(x^{(i)}, y^{(i)})$
- Every 1000 iterations plot  $\text{cost}(\theta, x^i, y^i)$  averaged over the last 1000 examples processed.

Examples:



- For stochastic to converge to global minima

$$\alpha = \frac{C_1}{\# \text{ iteration} + C_2}$$

# Online Learning

(42)

- continuous stream of data.

e.g. shipping website, offer ship their package package and users use it ( $y=1$ ) or not ( $y=0$ ).

- features  $x$ , price, user properties etc.

Probability for use to optimize price.

$$p(y=1|x; \theta)$$

\* Logistic Regression

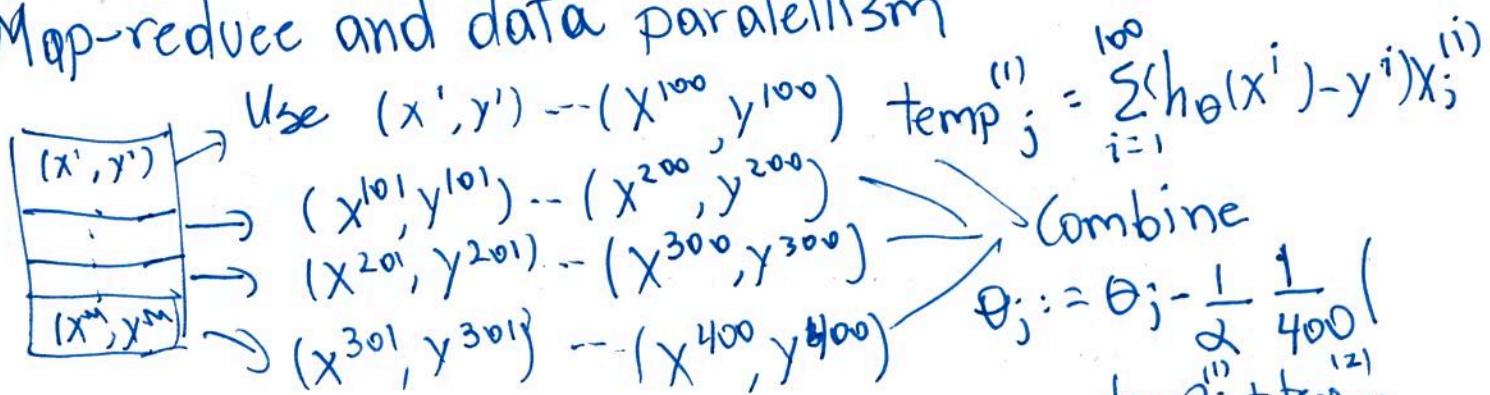
Repeat forever {

Get  $(x, y)$  correspond to user

Update  $\theta$  using  $(x, y)$ :  $\theta_j = \theta_j - \alpha (h_\theta(x) - y)x_j$

(can adapt to changing user preference.)

Map-reduce and data parallelism



Many algorithms can be expressed as computing sums of function over the training set.

e.g. logistic regression

- Can algorithm be expressed as sum over training set?

# Photo OCR

(optical character recognition)

1- text detection

2- character segmentation

3- character classification

\* System like this is called a machine learning pipeline.

- Sliding Windows classifier

- Pedestrian classifier: Supervised learning.

Move window in image to check over the photo.

Aspect ratio is ~~similar~~ equal always.

- Artificial Data Synthesis

- In the case OCR use fonts with background.

- Use different examples <sup>computer</sup> and apply transformations or distortions to amplify the data set.

- Make sure you have a low bias before getting more data.

- How much would it take to get 10x as much data?

Ways of getting more data:

- Artificial data

- Collect/label it yourself: calculate # of hours

- Crowd Source (e.g. Amazon Mechanical Turk)

# Estimating the errors due to each component - Ceiling Analysis

(46)

Component	Accuracy	To check every part of the pipeline feed
Overall System	72% ↓ 17%	100% accuracy results.
Text Detection	89%	
Character Segmentation	90% ↓ 1%	
Character Recognition	100% ↓ 10%	Manually give exact results and see changes.

Summary:

Supervised:

- Linear regression
- Logistic regression
- Neural Networks
- SVMs

Unsupervised

- k-means
- PCA
- Anomaly Detection

Special Algo applications

- Recommender Systems
- Large Scale machine learning

\* Advice on building a machine learning system

- Bias/Variance
- Regularization
- Evaluation of Algorithms
- Learning Curves
- Error Analysis
- Ceiling Analysis.