

What is this ?

Monday, May 29, 2017
3:22 PM

- Notes of Coursera course on Machine learning by Andrew Ng after discussion

Week 1

Monday, July 11, 2016
8:48 PM

Introduction:

What is Machine Learning ?

- a. Science of getting computers to Learn without being explicitly programmed
- b. Well Posed Learning Problem: Learn from
 - i.
- ii. $P(T)$ improves with E
- c. Grew out of AI

Experience	E
Task	T
Performance Measure	P

Examples:

- a. Database Mining
 - i. Medical Records
 - ii. Biology
- b. Applications that can't be programmed by hand
 - i. Autonomous Helicopter
 - ii. Handwriting recognition
- c. Self-Customizing Programs
 - i. Customized Recommendations
- d. Understanding Human Learning

Supervised Learning (Teach computer how to do something)

- a. We have the correct answer to the problems.
- b. Regression: Continuous valued output
 - i. Ex: House price prediction
 - ii. Predict the house price, having list of house and prices
 - iii. Predict future sales from past data
- c. Classification: Discrete valued output.
 - i. Ex; Classification of tumors
 - ii. Predict a tumor is malignant or benign
 - iii. Predict if the account is hacked or not

Unsupervised Learning (Let the computer learn by itself)

- a. No labels in data
- b. Given a data set => find the structure in data
 - i. Find the clusters in data
- c. Example:
 - i. Clustering news from different channels on same topic into groups
 - ii. Clustering the genes into groups.
- d. Used to organize Computers in Cluster
 - i. Computers which are interacting frequently can be placed together
- e. Grouping Social Network friends into groups
- f. Market Segmentation.

Linear Regression in One Variable:

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

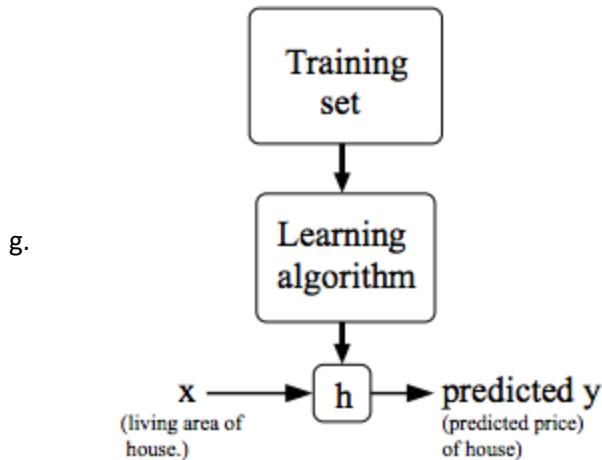
Parameters: θ_0, θ_1

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

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

Model Representation

- a. Supervised learning
- b. Notations:
 - i. x -> Input Variable
 - ii. y -> Output Variable
 - iii. X -> Input Space
 - iv. Y -> Output Space
- c. Training Set -> Learning Algorithm -> Hypothesis
- d. Size of House (x) -> Hypothesis (h) -> Estimated value of y
- e. h maps from x's to y's
- f. How to represent h
 - i. $h(x) = \theta_0 + \theta_1 * x$
 - ii. Linear regression with one variable x.
 - iii. Univariate linear regression
 - iv. θ_0 and θ_1 are parameters



Cost Function

- a. Used to measure the accuracy of our hypothesis.
- b. Choose theta₀ and theta₁ so that h(x) is closer to y for our training examples
- c. Minimize $(1/2m) * \text{Sum} (h(x) - y)^2$ over theta₀ and theta₁
- d. Squared error function or Mean squared error

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

e. $\underset{\theta_0, \theta_1}{\text{Minimize}} \ J(\theta_0, \theta_1)$
 $\underbrace{\quad}_{\text{Cost function}}$

- f. Simplified hypothesis function => theta₀ = 0

Gradient Descent

- a. Take a step in the direction of the deepest descent
- b. Different starting points -> Different paths
- c. Can go to the local minimum
- d. Assignment -> := , Assertion -> =

Gradient descent algorithm

```
repeat until convergence {
     $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$  (for  $j = 0$  and  $j = 1$ )
}
```

e.

Correct: Simultaneous update

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

Incorrect:

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

- f. If alpha is too small gradient descent can be slow
- g. If alpha is too large gradient descent can overshoot the minimum. It may fail to converge or even diverge
- h. At local minimum the gradient descent will leave the theta values unchanged as derivative term is zero.
- i. No need to decrease the alpha value over time. As the derivative term gets smaller when nearing to local minima => smaller steps when nearing the local minimum.
- j. FOR LINEAR REGRESSION:

repeat until convergence: {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x_i) - y_i)$$

k.

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m ((h_\theta(x_i) - y_i)x_i)$$

}

- l. Cost function for linear regression will always be a convex function.
 - i. Local optimum == Global optimum
- m. Batch Gradient Descent is another name for the above algorithm
 - i. Each Step of gradient descent uses all the training examples

Linear Algebra Review:

- a. Matrix: Rectangular Array of Numbers
- b. Dimension: No of Rows X No of Columns
- c. A_{ij} => ith row and jth column
- d. Vector: An $n \times 1$ matrix -> n dimensional vector.
- e. Implicitly 1 indexed arrays unless mentioned
- f. Scalar Multiplication.

- g. Matrix Vector Multiplication
- h. Matrix - Matrix Multiplication
- i. Matrix Multiplication properties
 - i. Not Commutative
 1. $A \times B \neq B \times A$
 2. $A \rightarrow m \times n, B \rightarrow n \times m \Rightarrow A \times B \rightarrow m \times m \text{ & } B \times A \rightarrow n \times n$
 - ii. Is Associative
 1. $A \times (B \times C) = (A \times B) \times C$
 - iii. Identity Matrix
 1. $A \times I = I \times A = A$
- j. Matrix Inverse
 - i. $A \times \text{Inv}(A) = \text{Inv}(A) \times A = I$
 - ii. Only square matrices have inverse
- k. Matrix Transpose

Questions:

q1: Does regression and classification fall under unsupervised also?

A: The classifications are orthogonal

q2: Learning in unsupervised way ?

A:

q3: what is the inference if the cost function is zero

A:

q4: Diff between linear regression and polynomial regression

A:

q5: Why this cost function for Linear regression ?

A:

q6: Is multivariate cost function convex ?

A:

Notes:

1. Note: YouTube video lectures has derivations

Action Items:

Learn about Maximum likelihood

Week 2

Sunday, September 18, 2016
4:08 PM

Question:

q1: Do we need both mean normalization and Feature scaling?

What is advantage of mean normalization after feature scaling ?

A:

q2: Practical case of alpha varying like WWWW

A:

q3: Need to know more pseudo inverse

Relation between x and $x^T x$ invertibility

q4: [TODO][Check] Normal equation only works for Regression but not for classification

A:

Linear Regression with Multiple Variables

Environment Setup Instructions

Multivariate Linear Regression

a. Notation

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

$x^{(i)}$ = the column vector of all the feature inputs of the i^{th} training example

m = the number of training examples

$n = |x^{(i)}|$; (the number of features)

b. Hypothesis

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

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

c. Cost Function

Cost function:

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

d. Gradient Descent

repeat until convergence: {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \quad \text{for } j := 0 \dots n$$

}

e. Feature Scaling:

Make sure that the features on the same scale

- i. If variables vary on different scale the contour plot can be more elliptical => More skew => More time to global minimum
- ii. Get every feature into approximately $-1 \leq x \leq 1$
Not strictly 1. Ranges near -1 to 1 should be good
Scale up the features whose range is very less
Scale down the features whose range is very large
Circular contour plot => Less time to global minimum

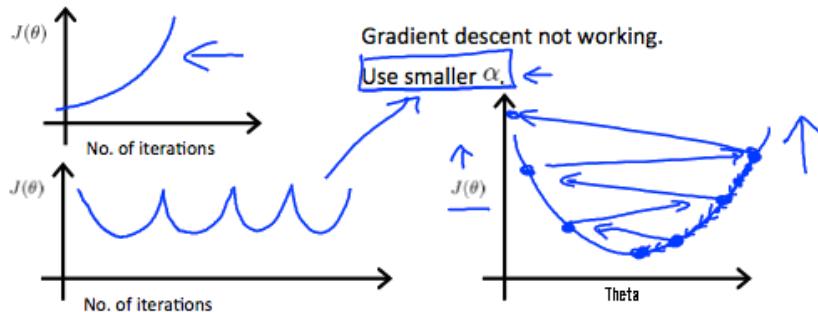
f. Mean Normalization

- i. Replace x with $(x - \mu)$ to make features have approximately zero mean (Do not apply for $x_0 = 1$)
- ii. More generally replace
 $x \Rightarrow (x - \mu)/s$
 $\mu \rightarrow \text{Average}$, $s \rightarrow \text{Standard deviation or range i.e. } (\max - \min)$

g. Make Sure Gradient descent is working properly

- i. Curve of $J(\theta)$ vs No of iterations should decrease after every iteration and converge over number of iterations
- ii. Automatic convergence test:
Declare convergence if $J(\theta)$ is decreasing less than epsilon
- iii. Use smaller value of alpha if $J(\theta)$ is increasing over number of iterations or if $J(\theta)$ is decreasing and increasing continuously.

Making sure gradient descent is working correctly.



- For sufficiently small α , $J(\theta)$ should decrease on every iteration.
- But if α is too small, gradient descent can be slow to converge.

iv. Try with alpha as below:
0.001, 0.003, 0.01, 0.003, 0.1, 0.3, 1

Plot $J(\theta)$ vs #of iterations

Pick alpha which is decreasing the $J(\theta)$ rapidly

Features and Polynomial regression

- We can make new features combining the existing features.
- You can fit a quadratic function or higher order polynomials to our data
 - Feature scaling is needed when you add higher order terms
 - Features can be var, var^2, var^(1/2)

Computing Parameters Analytically

- Normal Equation: Method to solve for theta analytically

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

- Feature scaling is not needed.

Gradient Descent

- Need to choose α .
- Needs many iterations.
- Works well even when n is large.

$$n = 10^6$$

Normal Equation

- No need to choose α .
- Don't need to iterate.
- Need to compute $(X^T X)^{-1}$ $\frac{n \times n}{O(n^3)}$
- Slow if n is very large.

$$n = 100$$

$$n = 1000$$

$$\dots \dots n = 10000$$

- What if $\text{inv}(X^T X)$ is not invertible
 - Redundant Features \rightarrow Linearly dependents

2. Too many features
 3. Use pinv to calculate inv even if the matrix is non invertible
- iv. [TODO][Check] Normal equation only works for Regression but not for classification

Octave Tutorial

%	Comment
<code>~=</code>	Not equal
<code>PS1('>>');</code>	change prompt
<code>;</code>	Suppresses output
<code>disp(a);</code>	Prints a
<code>disp(sprintf('2 decimals: %.2f', a))</code>	2 decimals: 3.14
In Matrix => ;	Separates rows
[1 2 3]	Matrix 1 x 3
[1; 2; 3]	Column Vector
1:0.1:2	From 1 to 2 in steps of 0.1
1:6	1 2 3 4 5 6
<code>ones(2,3)</code> <code>2*ones(2,3)</code>	2 2 2 2 2 2
<code>zeros(1,3)</code>	0 0 0
<code>rand(3, 3)</code>	random numbers from uniform distribution
<code>randn(1,3)</code>	gaussian random distribution
<code>hist(w, no_of_bins)</code>	histogram
<code>eye(4)</code>	Identity matrix of 4 x 4
<code>help</code>	documentation

<code>size(A)</code>	#row #columns
<code>length(v)</code>	length of the longest dimension
<code>pwd</code>	current directory
<code>ls</code>	list directories as in Linux

load(features.dat)	loads files
who	shows variable in current scope
whos	gives the detailed view of who
clear varname	removes varname
save hello.mat v	save variable v in hello.mat
load hello.mat	loads
save hello.txt	save as text ASCII
A(3,2)	Element in 3rd row and 2nd column
A(2, :)	Get everything in second column
A(:, 2)	Get everything in second row
A([1,3], :)	Get everything from first and third rows and all columns
A(:, 2) = [10; 11; 12]	Edit 2nd col
A = [A, [a; b; c]]	append another column
A(:)	Put all elements of A in single vector
C = [A B] or [A, B] C = [A; B]	concatenate A and B

A * C	Matrix multiplication
A .* B	Element wise multiplication
A.^2	Element wise squared
1 ./ V	Element wise reciprocal of V
log(v)	Element wise log
exp(v)	Element wise exp
abs(v)	Element wise abs
-v	-1 * v
A'	Transpose of A
[val, ind] = max(a)	Max value and index
val = max(a)	Column wise max

A < 3	element wise comparison
find(a < 3)	
magic(n)	n x n all rows cols and diags add up to same number
sum(a), prod(a), floor(a)	
max(A, [], 1)	Max of matrix col wise 1 -> first dimension
max(max(A)) or max(A(:))	Max of all the matrix
sum(A, 1) sum(A, 2)	sum of elements along dimensions 1 or 2
sum(sum(A .* eye(9))) sum(sum(A .* flipup(eye(9)))) // flipup, flipud	Sum of diagonal elements
pinv(a)	pseudo inverse

plot(t)	plot function
hold on;	plot new functions on the old ones
xlabel, ylabel, legend, title	params of plot
print =dpng 'myplot.png'	save plot
close	get away from the figure displayed during plot
figure(1) figure(2)	Multiple figures
subplot(1,2,1)	first element of 1,2 grid
axis([a b c d])	a b for x range c d for y range
clf;	clears a figure
imagesc(A) imagesc(A), colorbar, colormap gray;	visualize a matrix A
,	chaining multiple commands in same line

<pre>>> for i=indices, > disp(i); > end;</pre>	For
<pre>>> i = 1; >> while i <= 5, > v(i) = 100; > i = i+1; > end;</pre>	While
break and continue	
<pre>function y = squareThisNumber(x) y = x^2;</pre>	Function
<pre>function [y1,y2] = squareAndCubeThisNumber(x) y1 = x^2; y2 = x^3;</pre>	
addpath(dir)	find the functions
<pre>J = 1/(2*m) * sum(sqrErrors); </pre>	Cost functions

Vectorized implementation:

$$\Theta := \Theta - \alpha \delta$$

where $\delta = \frac{1}{m} \sum_{i=1}^m (h_\Theta(x^{(i)}) - y^{(i)}) x^{(i)}$

Week 3

Monday, May 29, 2017
9:33 PM

Classification and Representation

- a. Examples:
 - i. Email -> Spam / Not Spam
 - ii. Online Transaction -> Fraud /Not Fraud

- iii. Tumor \rightarrow Malignant / Benign
- b. Binary classification

$$\underline{y \in \{0, 1\}}$$

0: "Negative Class" (e.g., benign tumor)
 1: "Positive Class" (e.g., malignant tumor)

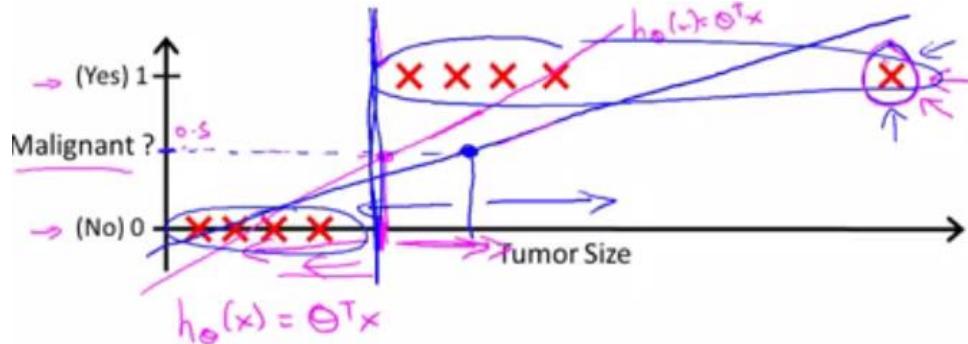
- i. Yes/No Type of problems
- ii. Positive/Negative assignment for 1/0 is arbitrary.
- c. Threshold Classifier

Threshold classifier output $h_\theta(x)$ at 0.5:

i. \rightarrow If $h_\theta(x) \geq 0.5$, predict "y = 1"

\cdot If $h_\theta(x) < 0.5$, predict "y = 0"

- ii. Outliers in data can change the line that we are fitting.



- iii. The output can be $>> 1$ or $<< 0$ even though our labels are 0 and 1 only

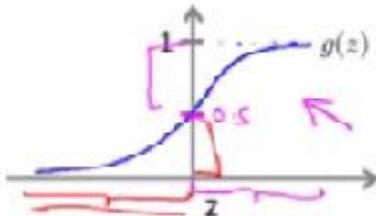
Logistic Regression Model

- a. Hypothesis
 $g(z)$ is sigmoid or logistic function

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

$$z = \theta^T x$$

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



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

$$P(y = 0|x; \theta) + P(y = 1|x; \theta) = 1$$

- b. Decision Boundary
 \cdot We predict

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

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

ii. But

$$g(z) \geq 0.5$$

when $z \geq 0$

iii. Implies

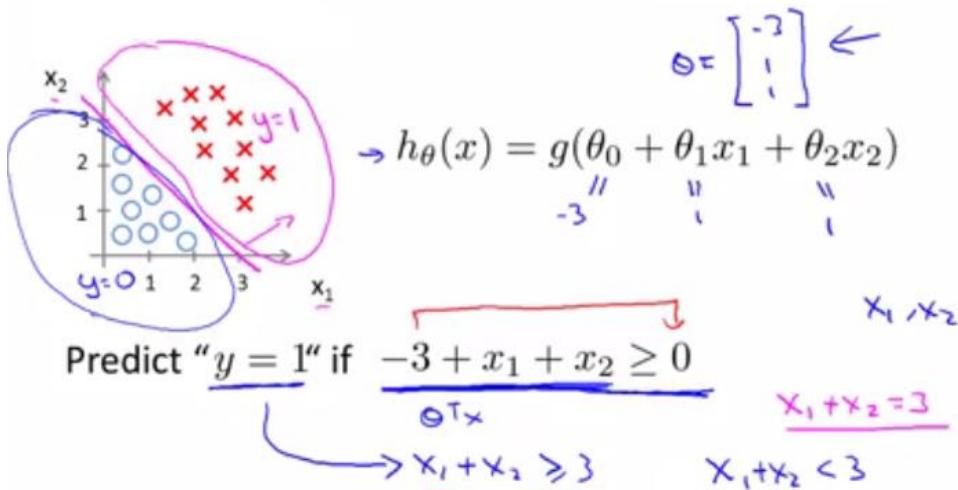
$$h_{\theta}(x) = g(\theta^T x) \geq 0.5$$

when $\theta^T x \geq 0$

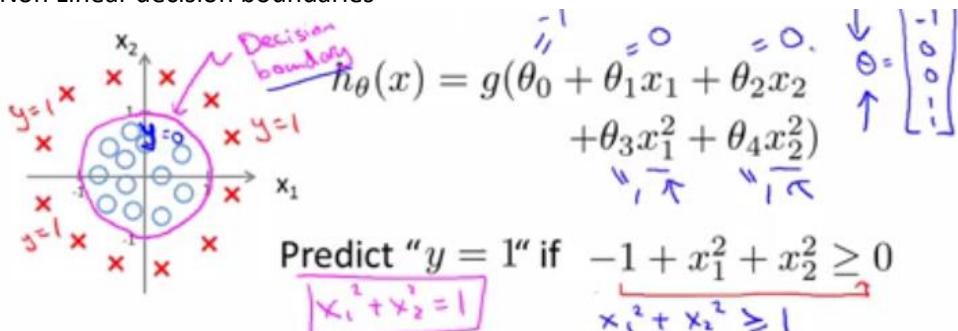
iv. Decision boundary is a property of the parameters but not the property of the Data.

Once we have Theta that completely defines the decision boundary.

c. Linear Decision Boundary



d. Non Linear decision boundaries



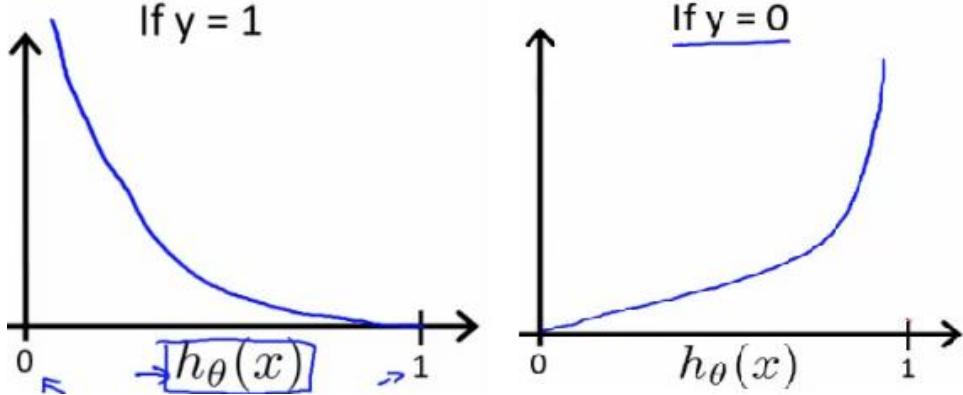
e. Cost Function

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

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

$$\text{Cost}(h_{\theta}(x), y) = -\log(1 - h_{\theta}(x)) \quad \text{if } y = 0$$

$\text{Cost}(h_\theta(x), y) = 0$ if $h_\theta(x) = y$
 $\text{Cost}(h_\theta(x), y) \rightarrow \infty$ if $y = 0$ and $h_\theta(x) \rightarrow 1$
 $\text{Cost}(h_\theta(x), y) \rightarrow \infty$ if $y = 1$ and $h_\theta(x) \rightarrow 0$



f. Simplified Cost Function

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

A vectorized implementation is:

$$\begin{aligned} h &= g(X\theta) \\ J(\theta) &= \frac{1}{m} \cdot \left(-y^T \log(h) - (1 - y)^T \log(1 - h) \right) \end{aligned}$$

This cost function can be derived from principle of maximum likelihood. And this has nice property of being convex.

g. Gradient Descent

Gradient Descent

Remember that the general form of gradient descent is:

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

We can work out the derivative part using calculus to get:

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

Notice that this algorithm is identical to the one we used in linear regression. We still have to simultaneously update all values in theta.

A vectorized implementation is:

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

h. Check Gradient Descent

Plot

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

as a function of the number of iterations and make sure

$J(\theta)$ is decreasing on every iteration.

i. Optimization Algorithm for Gradient Descent

- i. Conjugate Gradient
- ii. BFGS
- iii. L-BFGS

Advantages:

- i. No need to manually pick alpha
- ii. Often faster than gradient descent

Disadvantages:

- i. More complex

Octave:

- i. Write costFunction taking input theta and giving output jVal and gradient Vector

```
function [jVal, gradient] = costFunction(theta)
    jVal = [...code to compute J(theta)...];
    gradient = [...code to compute derivative of J(theta)...];
end
```

- ii. options = optimset('Gradobj','on','MaxIter','100');

- iii. initialTheta = zeros(<dim1>,<dim2>);

```

iv. [optTheta, functionVal, exitVal] = fminunc(@costFunction, initialTheta, options);
    options = optimset('GradObj', 'on', 'MaxIter', 100);
    initialTheta = zeros(2,1);
    [optTheta, functionVal, exitFlag] = fminunc(@costFunction, initialTheta,
        options);
    ...

```

Multiclass Classification

- a. Examples:
 - i. Email Tagging
 - ii. Medical Diagnosis -> Not ill, Cold, Flu
 - iii. Weather -> Sunny, Cloudy, Rainy, Snow
- b. One Vs All

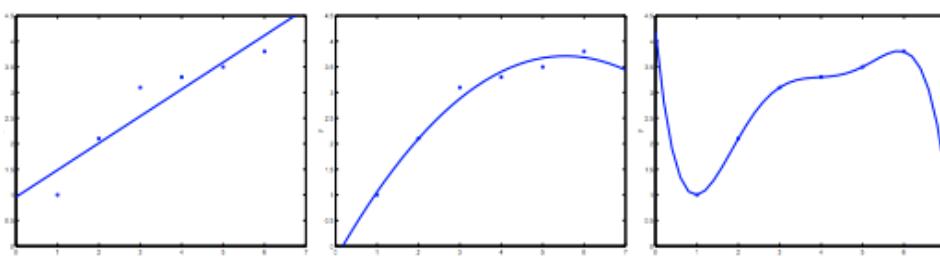
Train a logistic regression classifier $h_{\theta}^{(i)}(x)$ for each class i to predict the probability that $y = i$.

On a new input x , to make a prediction, pick the class i that maximizes

$$\max_i h_{\theta}^{(i)}(x)$$

$$\begin{aligned}
 y &\in \{0, 1 \dots n\} \\
 h_{\theta}^{(0)}(x) &= P(y = 0|x; \theta) \\
 h_{\theta}^{(1)}(x) &= P(y = 1|x; \theta) \\
 &\dots \\
 h_{\theta}^{(n)}(x) &= P(y = n|x; \theta) \\
 \text{prediction} &= \max_i(h_{\theta}^{(i)}(x))
 \end{aligned}$$

Solve the Problem of Overfitting:



- a. Underfitting -> High Bias
Hypothesis not even fitting the training set

Overfitting -> High Variance
Hypothesis fits training data well but fails to generalize

- b. Addressing Overfitting:
 - i. Reduce number of features
 - Manually select features
 - Model selection algorithm
 - ii. Regularization
 - Keep all features but reduce magnitude values of theta

- c. Regularization

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

Smaller values of theta

- i. Shrink the Parameters by adding additional theta terms.
- ii. Simpler hypothesis
- iii. Less prone to overfitting
- iv. We can't select the parameters for shrinking unless we know in prior. So we shrink all the parameters.
- v. Large value of lambda => underfitting
- vi. Note: theta0 is not penalized

- d. Regularized Linear Regression:

Cost Function:

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

Gradient Descent

Repeat {

$$\begin{aligned} \theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_j &:= \theta_j - \alpha \left[\left(\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \right) + \frac{\lambda}{m} \theta_j \right] \quad j \in \{1, 2, \dots, n\} \\ &\} \end{aligned}$$

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

Normal Equation

$$\theta = \left(X^T X + \lambda \cdot L \right)^{-1} X^T y$$

where $L = \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & 1 & \\ & & & \ddots \\ & & & & 1 \end{bmatrix}$

Note: Below matrix is always invertible

$$\left(X^T X + \lambda \cdot L \right)^{-1}$$

e. Regularized Logistic Regression:

Cost Function:

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

Gradient Descent:

Repeat {

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

$$\theta_j := \theta_j - \alpha \left[\underbrace{\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}}_{(j = 1, 2, 3, \dots, n)} + \underbrace{\frac{\lambda}{m} \theta_j}_{\theta_1, \dots, \theta_n} \right] \leftarrow$$

}

$$\frac{\partial J(\theta)}{\partial \theta_j}$$

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

Week 4

Sunday, June 04, 2017

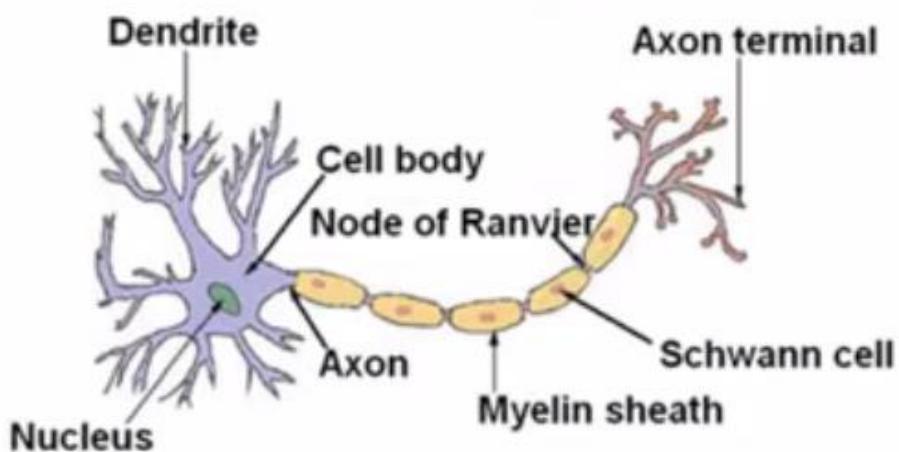
9:54 AM

Neural Networks

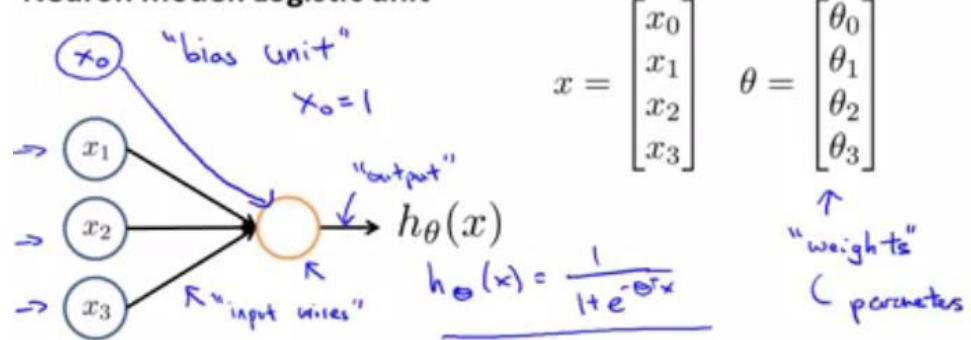
1. Motivations:

- a. Non Linear Classification
 - i. Number of quadratic features increase as $O(n^2)$
 - ii. Also Cubic features increase as $O(n^3)$
 - iii. So the number of features to consider is large

- b. Computer Vision
 - i. Features are not linearly separable
 - c. Neurons and Brain: Mimics Brain
 - i. Neurons learn to do the respective work.
 - ii. Auditory part of brain can learn to see.
2. Model Representation:



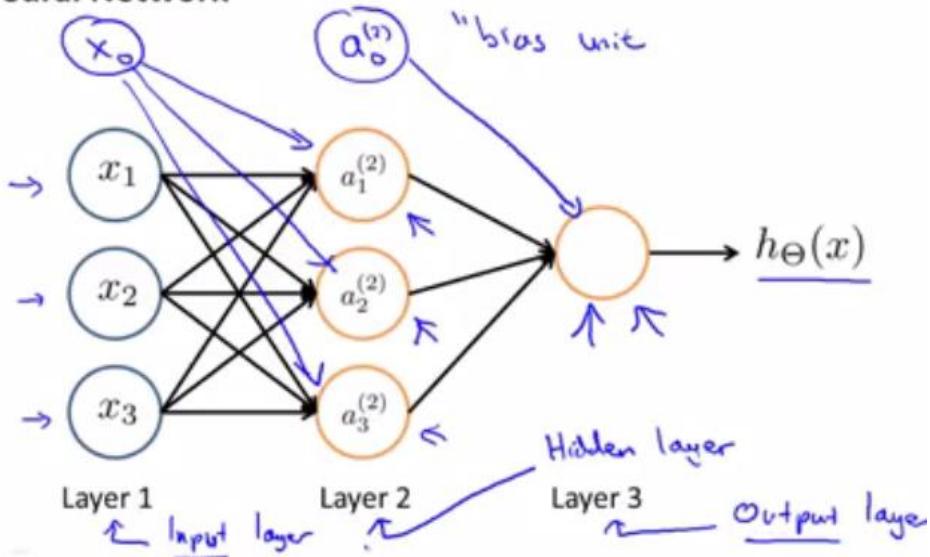
Neuron model: Logistic unit



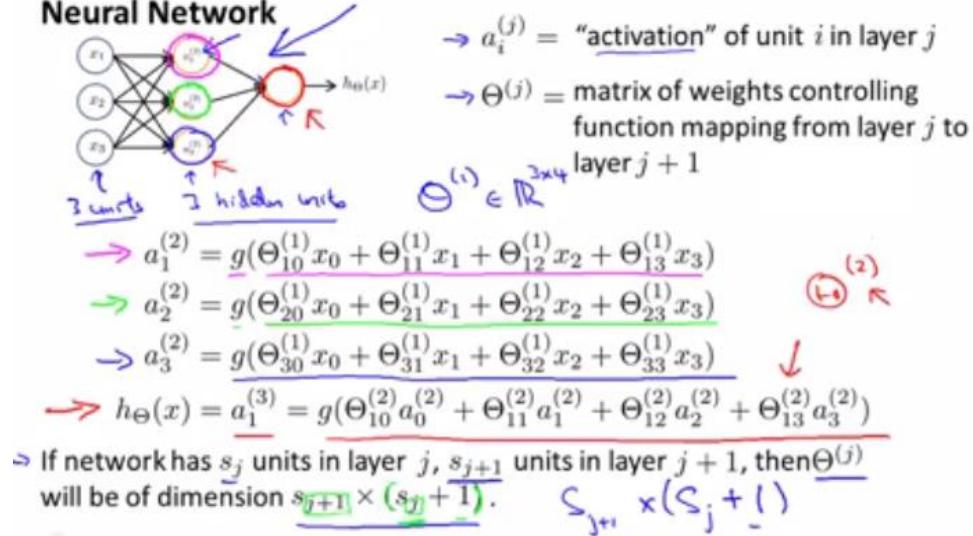
Sigmoid (logistic) activation function.

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

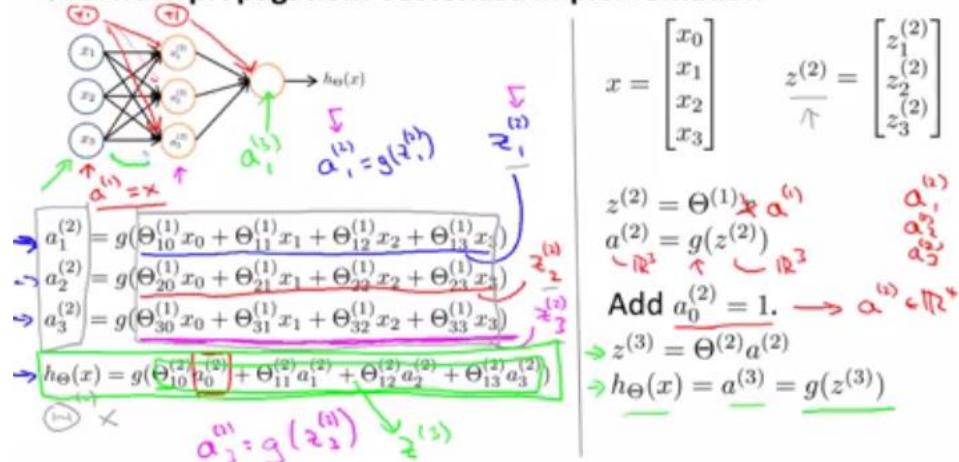
Neural Network



Neural Network



Forward propagation: Vectorized implementation



Note:

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

b. Neural networks will learn its own features. These features are formed from the input features.

c. We can have neurons connected differently i.e. different architectures.

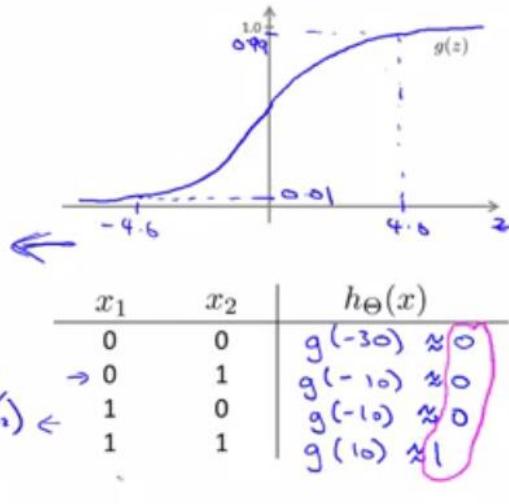
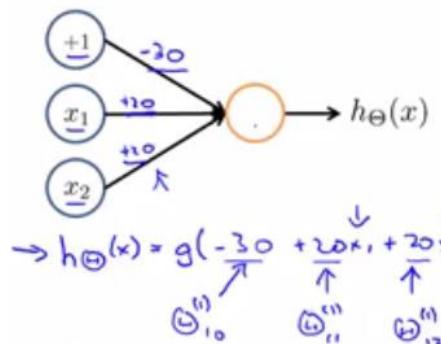
3. Examples and Intuitions:

a. AND

Simple example: AND

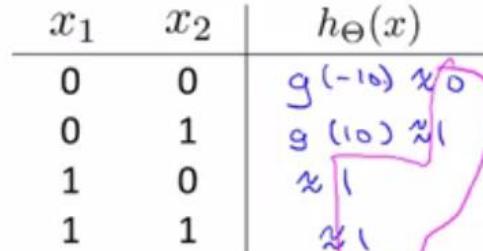
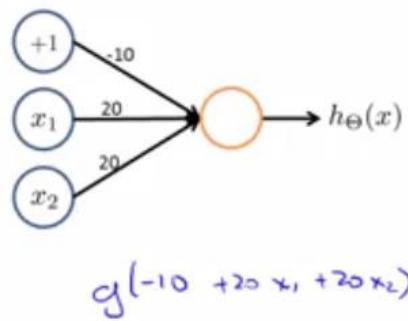
$\Rightarrow x_1, x_2 \in \{0, 1\}$

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



b. OR

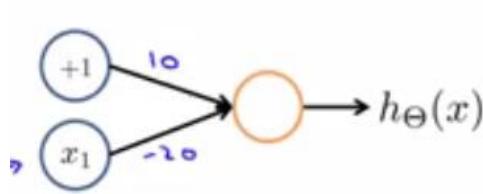
Example: OR function



c. NOT

Negation:

$\text{NOT } x_1$

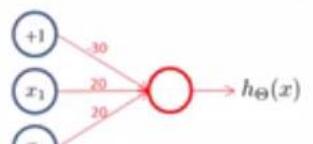


x_1	$h_\Theta(x)$
0	$g(10) \approx 1$
1	$g(-10) \approx 0$

$$h_\Theta(x) = g(10 - 20x_1)$$

d. XNOR

Putting it together: $x_1 \text{ XNOR } x_2$



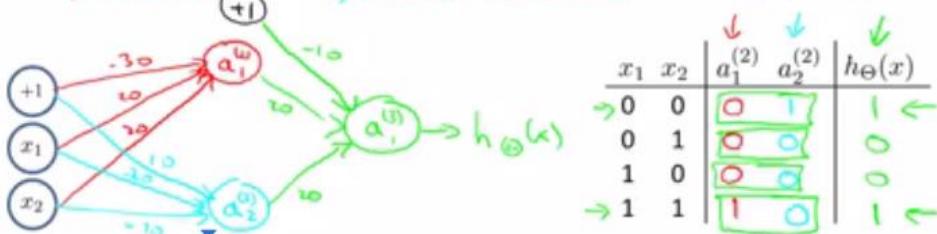
$\rightarrow x_1 \text{ AND } x_2$



$\rightarrow (\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$



$\rightarrow x_1 \text{ OR } x_2$



x_1	x_2	$a_1^{(2)}$	$a_2^{(2)}$	$h_\Theta(x)$
0	0	0	1	1 ↪
0	1	0	0	0 ↪
1	0	0	0	0 ↪
1	1	1	0	1 ↪

4. Multiclass Classification in Neural networks

Multiple output units: One-vs-all.

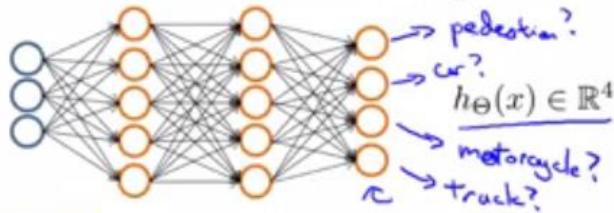


Pedestrian

Car

Motorcycle

Truck



Want $h_\Theta(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, $h_\Theta(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$, $h_\Theta(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$, etc.
 when pedestrian when car when motorcycle

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

pedestrian car motorcycle truck

Week 5

Sunday, June 04, 2017
6:42 PM

Cost Function and Backpropagation.

Cost Function:

1. Definitions:

- L = total number of layers in the network
- s_l = number of units (not counting bias unit) in layer l
- K = number of output units/classes

$$\alpha_j^{(l)}$$

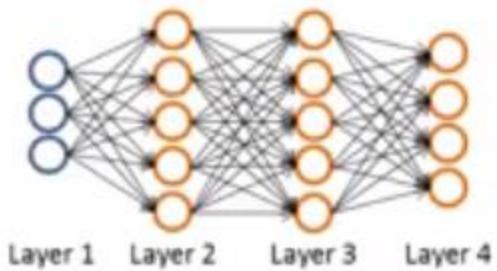
$$\delta_j^{(l)}$$

2. Cost Function

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

$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[y_k^{(i)} \log((h_\Theta(x^{(i)}))_k) + (1 - y_k^{(i)}) \log(1 - (h_\Theta(x^{(i)}))_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{j,i}^{(l)})^2$$

3. Backpropagation Algorithm



Forward propagation:

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

Gradient computation: Backpropagation algorithm

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

For each output unit (layer $L = 4$)

$$\delta_j^{(4)} = \underline{a_j^{(4)} - y_j} \quad (h_{\Theta}(x))_j \quad \underline{\delta^{(4)}} = \underline{a^{(4)} - y}$$

$$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} \cdot * g'(z^{(3)})$$

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

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

$$\frac{\partial J(\Theta)}{\partial \Theta_{ij}^{(k)}} = a_j^{(k)} \delta_i^{(k+1)} \quad \begin{array}{l} \text{(ignoring } \lambda \text{ if } \\ \lambda = 0) \end{array}$$

Backpropagation algorithm

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

Set $\Delta_{ij}^{(l)} = 0$ (for all l, i, j). (use to compute $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$)

For $i = 1$ to $m \leftarrow (x^{(i)}, y^{(i)})$.

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

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

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

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

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

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

$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} \text{ if } j \neq 0$

$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \text{ if } j = 0$

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$$

4. Implementation Note: Unrolling Parameters

In order to use optimizing functions such as "fminunc()", we will want to "unroll" all the elements and put them into one long vector:

```
1 thetaVector = [ Theta1(); Theta2(); Theta3(); ]
2 deltaVector = [ D1(); D2(); D3() ]
```

If the dimensions of Theta1 is 10x11, Theta2 is 10x11 and Theta3 is 1x11, then we can get back our original matrices from the "unrolled" versions as follows:

```
1 Theta1 = reshape(thetaVector(1:110),10,11)
2 Theta2 = reshape(thetaVector(111:220),10,11)
3 Theta3 = reshape(thetaVector(221:231),1,11)
4
```

Learning Algorithm

→ Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$.

→ Unroll to get `initialTheta` to pass to

`fminunc(@costFunction, initialTheta, options)`

```
function [jval, gradientVec] = costFunction(thetaVec)
    From thetaVec, get  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ .
    Use forward prop/back prop to compute  $D^{(1)}, D^{(2)}, D^{(3)}$  and  $J(\Theta)$ 
    Unroll  $D^{(1)}, D^{(2)}, D^{(3)}$  to get gradientVec.
```

5. Gradient Checking

- a. Very slow. Use only to verify once.

b.
$$\frac{\partial}{\partial \theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

```

epsilon = 1e-4;
for i = 1:n,
    thetaPlus = theta;
    thetaPlus(i) += epsilon;
    thetaMinus = theta;
    thetaMinus(i) -= epsilon;
    gradApprox(i) = (J(thetaPlus) - J(thetaMinus))/(2*epsilon)
end;

```

6. Random Initialization:

Random initialization: Symmetry breaking

➤ Initialize each $\Theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$
 (i.e. $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$)

E.g.

↗ Random 10x11 matrix (b/w. 0
and 1)

→ **Theta1 =** rand(10,11)* (2*INIT_EPSILON)
- INIT_EPSILON; $[-\epsilon, \epsilon]$

→ **Theta2 =** rand(1,11)* (2*INIT_EPSILON)
- INIT_EPSILON;

7. Putting it Together

Training a Neural Network

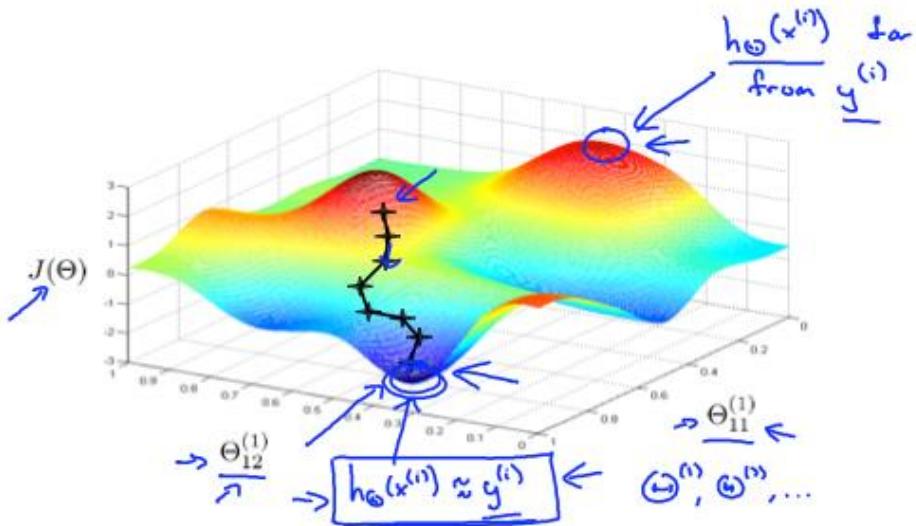
1. Randomly initialize the weights
2. Implement forward propagation to get $h_\theta(x^{(i)})$ for any $x^{(i)}$
3. Implement the cost function
4. Implement backpropagation to compute partial derivatives
5. Use gradient checking to confirm that your backpropagation works. Then disable gradient checking.
6. Use gradient descent or a built-in optimization function to minimize the cost function with the weights in theta.

When we perform forward and back propagation, we loop on every training example:

```

1 for i = 1:m,
2     Perform forward propagation and backpropagation using example (x(i),y(i))
3     (Get activations a(l) and delta terms d(l) for l = 2,...,L)

```



$J(\theta)$ is not convex we can get stuck in local minima.
 -> Activation of node j in layer l

-> Error due to node j in layer l

Week 6

Wednesday, June 28, 2017
 1:56 AM

Evaluating the Learning Algorithm:

1. Deciding what to try next:

Suppose you have implemented regularized linear regression to predict housing prices.

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

However, when you test your hypothesis on a new set of houses, you find that it makes unacceptably large errors in its predictions. What should you try next?

- a. Get more training examples
- b. Try smaller set of features
- c. Try getting additional features
- d. Try adding polynomial features
- e. Try decreasing lambda -> Regularization
- f. Try increasing lambda

2. Machine Learning Diagnostic

Diagnostic: A test that you can run to gain insight what is/isn't working with a learning algorithm, and gain guidance as to how best to improve its performance.

3. Evaluating a Hypothesis:

- a. Low training error does not imply good hypothesis.
- b. We can plot the curve and check for overfitting
 - i. Not possible when there are multiple variables
- c. How to Evaluate ?
 - i. Split the data into Training Set and Test Set in 7:3 Ratio.
 - ii. m_{test} = number of test examples
 - iii. If the data is overfitting train error will be low and test error will be high.
- d. Procedure
 - i. Learn parameter theta from training data minimizing training error $J(\theta)$
 - ii. Compute test set error:

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} \left(h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)} \right)^2$$

$$J_{test}(\theta) = -\frac{1}{m_{test}} \sum_{i=1}^{m_{test}} y_{test}^{(i)} \log h_{\theta}(x_{test}^{(i)}) + (1 - y_{test}^{(i)}) \log (1 - h_{\theta}(x_{test}^{(i)}))$$

In case of logistic regression you can use:

Misclassification Error (0/1 misclassification error)

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

$$\text{Test error} = \frac{1}{m_{test}} \sum_{i=1}^{m_{test}} \text{err}(h_{\theta}(x_{test}^{(i)}), y_{test}^{(i)}).$$

4. Model Selection and Train/Validation/Test Sets

- a. Once parameters $\theta_0, \theta_1, \dots, \theta_4$ were fit to some set of data (training set), the error of the parameters as measured on that data (the training error $J(\theta)$) is likely to be lower than the actual generalization error.
- b. Model Selection

- $\lambda=1$ 1. $\rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x$
- $\lambda=2$ 2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$
- $\lambda=3$ 3. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3$
- $\vdots \quad \vdots$
- $\lambda=10$ 10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10}$

Take theta values for all the models.

Measure their performance in the test set.

How well does this model generalize ?

Report test set error.

c. Problem:

We chose the parameter d to test set. The hypothesis will overfit for the test set while selecting d.

Solution:

Split it into three pieces Test Set/ Validation Set/ Test in ration say 6 : 2 : 2

Training/ validation / test error:

Training error:

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

Cross Validation error:

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

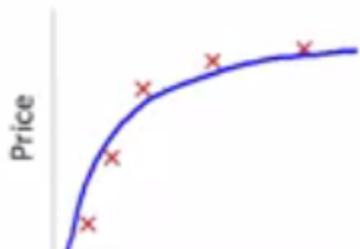
Test error:

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

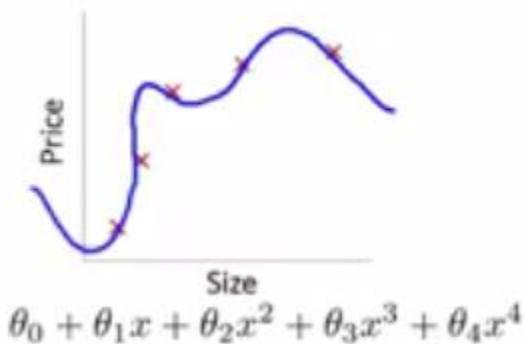
- d. Now use Validation test to get more hyper parameters.
 Test the models on validation test.
 Pick up the model with lowest Cross validation error.
 And use Test set to evaluate.
 Use different data for all the Sets above.
- e. Procedure
 - Optimize the parameters in Θ using the training set for each polynomial degree.
 - Find the polynomial degree d with the least error using the cross validation set.
 - Estimate the generalization error using the test set with $J_{test}(\Theta(d))$, ($d = \text{theta from polynomial with lower error}$);
5. Diagnosing Bias vs Variance
 - High Bias:



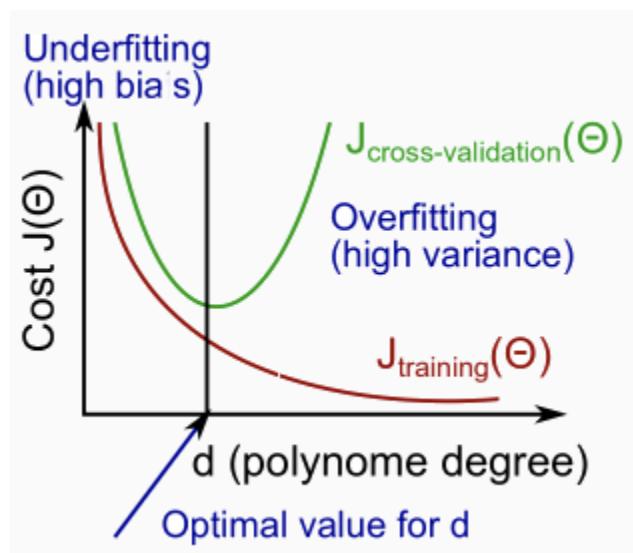
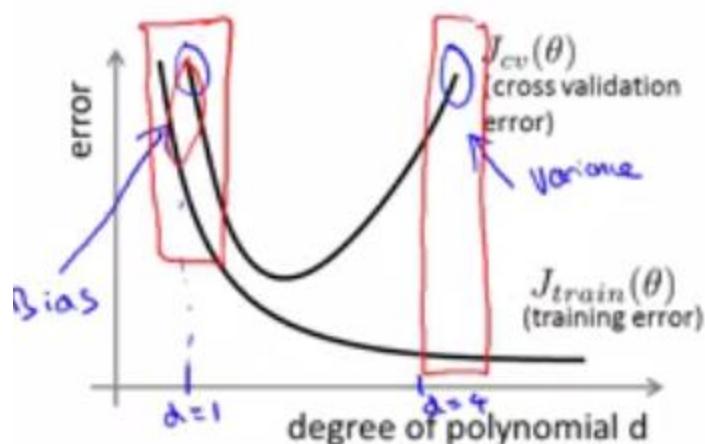
b. Just Right



c. High Variance:



d. Plot Error vs Degree of polynomial d.



Bias (underfit):

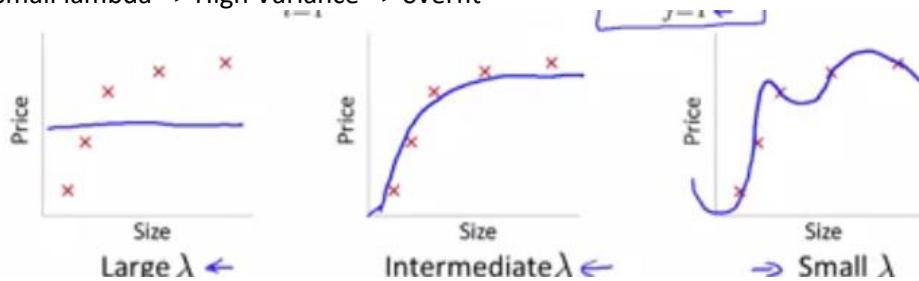
$\rightarrow J_{\text{train}}(\theta)$ will be high }
 $J_{\text{cv}}(\theta) \approx J_{\text{train}}(\theta)$ }

Variance (overfit):

$J_{\text{train}}(\theta)$ will be low }
 $J_{\text{cv}}(\theta) \gg J_{\text{train}}(\theta)$ }

6. Regularization and Bias/Variance
- Large lambda => High Bias => Underfit
 Intermediate lambda => Just Right

Small lambda => High Variance => overfit



Model selection for Lambda

$$\text{Model: } h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

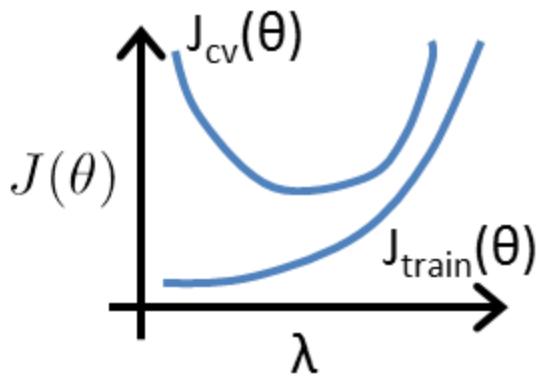
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$

1. Try $\lambda = 0 \leftarrow$ $\min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
2. Try $\lambda = 0.01 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
3. Try $\lambda = 0.02 \rightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$
4. Try $\lambda = 0.04 \vdots$
5. Try $\lambda = 0.08 \vdots$
12. Try $\lambda = 10 \uparrow 10 \cdot 24 \rightarrow \theta^{(12)} \rightarrow J_{cv}(\theta^{(12)})$
- Pick (say) $\theta^{(5)}$. Test error: $J_{test}(\theta^{(5)})$

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

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

$$\Rightarrow \boxed{J_{cv}(\theta)} = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

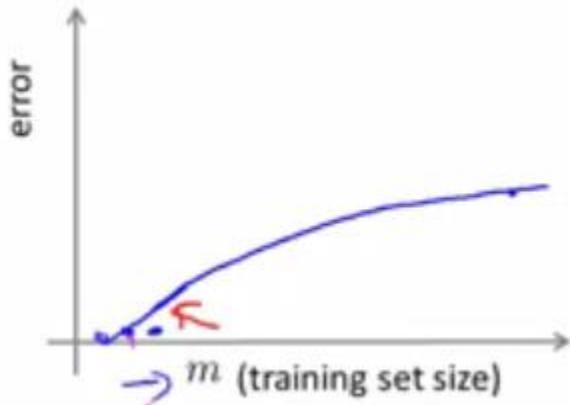


Procedure:

- Create a list of lambdas (i.e. $\lambda \in \{0, 0.01, 0.02, 0.04, 0.08, 0.16, 0.32, 0.64, 1.28, 2.56, 5.12, 10.24\}$);
- Create a set of models with different degrees or any other variants.
- Iterate through the λ s and for each λ go through all the models to learn some Θ .
- Compute the cross validation error using the learned Θ (computed with λ) on the $J_{CV}(\Theta)$ **without regularization or $\lambda = 0$** .
- Select the best combo that produces the lowest error on the cross validation set.
- Using the best combo Θ and λ , apply it on $J_{test}(\Theta)$ to see if it has a good generalization of the problem.

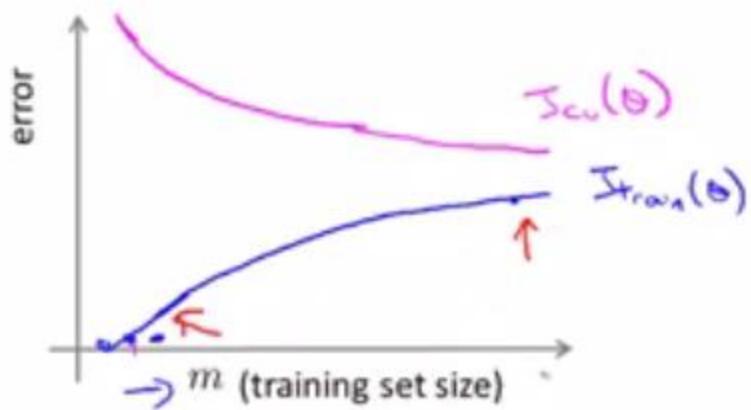
7. Learning Curves:

- Training set error:

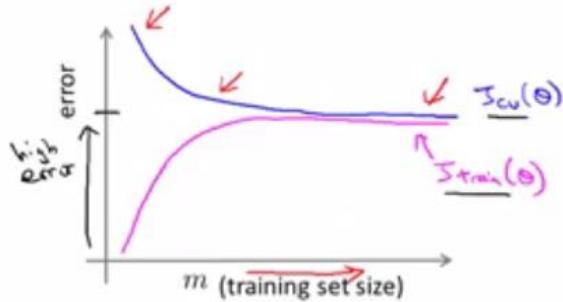


Easy to fit smaller data.

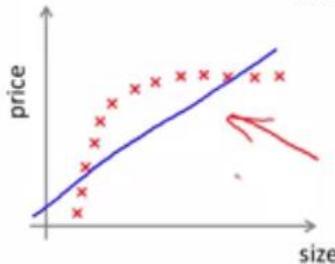
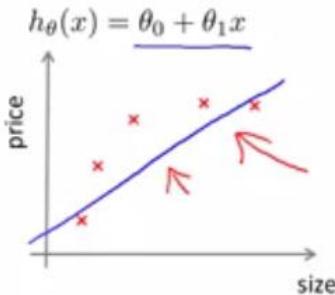
- Validation or Test error;



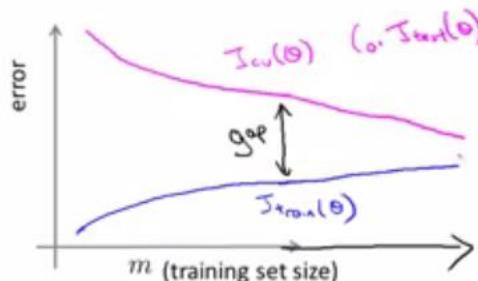
High bias



If a learning algorithm is suffering from high bias, getting more training data will not (by itself) help much.

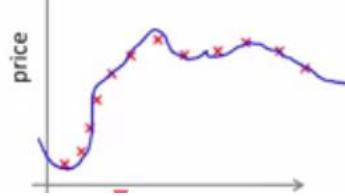
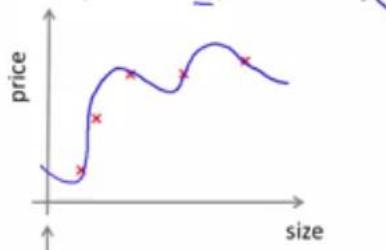


High variance

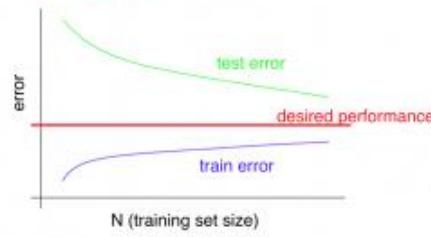
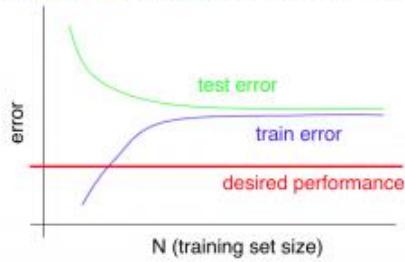


If a learning algorithm is suffering from high variance, getting more training data is likely to help.

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



Typical learning curve for high bias(at fixed model complexity): Typical learning curve for high variance(at fixed model complexity):

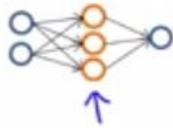


8. Deciding What to Do Next:

- Get more training examples → fixes high variance
- Try smaller sets of features → fixes high variance
- Try getting additional features → fixes high bias
- Try adding polynomial features (x_1^2, x_2^2, x_1x_2 , etc) → fixes high bias.
- Try decreasing λ → fixes high bias
- Try increasing λ → fixes high variance

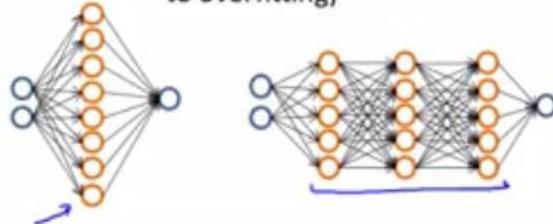
Neural Networks and Overfitting:

"Small" neural network
(fewer parameters; more prone to underfitting)



Computationally cheaper

→ "Large" neural network
(more parameters; more prone to overfitting)



Computationally more expensive.

Use regularization (λ) to address overfitting.

Diagnosing Neural Networks

- A neural network with fewer parameters is **prone to underfitting**. It is also **computationally cheaper**.
- A large neural network with more parameters is **prone to overfitting**. It is also **computationally expensive**. In this case you can use regularization (increase λ) to address the overfitting.

Using a single hidden layer is a good starting default. You can train your neural network on a number of hidden layers using your cross validation set. You can then select the one that performs best.

Model Complexity Effects:

- Lower-order polynomials (low model complexity) have high bias and low variance. In this case, the model fits poorly consistently.
- Higher-order polynomials (high model complexity) fit the training data extremely well and the test data extremely poorly. These have low bias on the training data, but very high variance.
- In reality, we would want to choose a model somewhere in between, that can generalize well but also fits the data reasonably well.

Machine Learning System Design

1. Building a spam classifier

- a. Supervised Learning:
 $X \Rightarrow$ features of email
 $Y \Rightarrow$ spam (1) or not spam (0)
- b. How to spend time:
 - i. Collect lots of data
 - ii. Develop sophisticated features based on email routing
 - iii. Develop sophisticated features based on message body
 - iv. Develop sophisticated algorithm to detect misspelling

It is difficult to tell which is more useful

2. Recommended approach:

- a. Start with a simple algorithm that you can implement quickly. Implement it and test it on your cross validation data.
- b. Plot learning curves to decide if more data, more features etc. are likely to help
- c. Error Analysis: Manually examine the examples in cross validation set that your algorithm made errors on. See if you spot any systematic trend in what type of examples it is making errors on.

3. Error Metrics for Skewed Classes:

- a. Consider cancer classification example

Train logistic regression model $y = 1$ if cancer or $y = 0$ otherwise

Find that got 1% Error on test set

Only 0.50% of patients have cancer

$Y = 0$ always $\Rightarrow 0.5\%$ Error.

- b. Precision/Recall

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

		Actual class		\rightarrow Precision (Of all patients where we predicted $y = 1$, what fraction actually has cancer?)
		1	0	
Predicted 1 class	1	True positive	False positive	$\frac{\text{True positives}}{\# \text{predicted positive}} = \frac{\text{True positive}}{\text{True pos} + \text{False pos}}$
	0	False negative	True negative	\rightarrow Recall (Of all patients that actually have cancer, what fraction did we correctly detect as having cancer?)
	$y = 0$			$\frac{\text{True positives}}{\text{actual positives}} = \frac{\text{True positives}}{\text{True pos} + \text{False neg}}$
		$\text{recall} = 0$		

- c. Trading Off Precision and Recall.

Trading off precision and recall

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

Predict 1 if $h_\theta(x) \geq 0.8$ ~~0.7~~ ~~0.9~~ 0.3 \leftarrow

Predict 0 if $h_\theta(x) < 0.8$ ~~0.7~~ ~~0.9~~ 0.3

\Rightarrow Suppose we want to predict $y = 1$ (cancer) only if very confident.

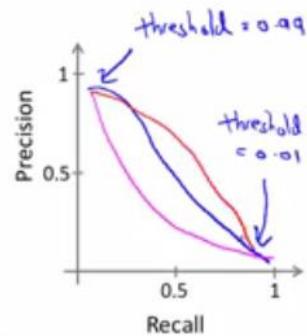
\rightarrow Higher precision, lower recall.

\Rightarrow Suppose we want to avoid missing too many cases of cancer (avoid false negatives).

\rightarrow Higher recall, lower precision.

$$\rightarrow \text{precision} = \frac{\text{true positives}}{\text{no. of predicted positive}}$$

$$\rightarrow \text{recall} = \frac{\text{true positives}}{\text{no. of actual positive}}$$



More generally: Predict 1 if $h_\theta(x) \geq \text{threshold}$.

- d. How to compare precision/recall numbers?

	Precision(P)	Recall (R)	Average	F ₁ Score
Algorithm 1	0.5	0.4	0.45	0.444 ←
Algorithm 2	0.7	0.1	0.4	0.175 ←
Algorithm 3	0.02	1.0	0.51	0.0392 ←

Average: $\frac{P+R}{2}$

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

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

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

Predict $y=1$ all the time

4. Data for High Accuracy learning system

Classify between confusable words {to, two, too}, {then, than}

Algorithms Used:

Perceptron, Winnow, Memory-based, Naïve Bayes

[Banko and Brill, 2001]: Increasing data => Increase in Perf

Features must represent data and many features => Low bias

Large training set => Low variance

Week 7

Thursday, June 29, 2017

6:28 PM

Support Vector Machines:

1. Optimization Objective (Mathematical Definition of SVMs)

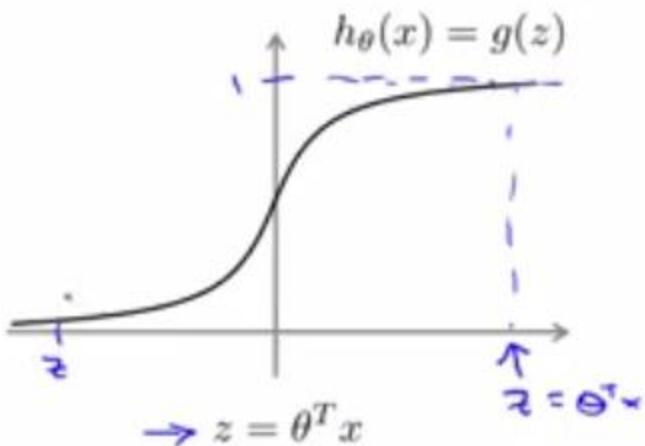
a. Alternative view of logistic regression

Sigmoid function:

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

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

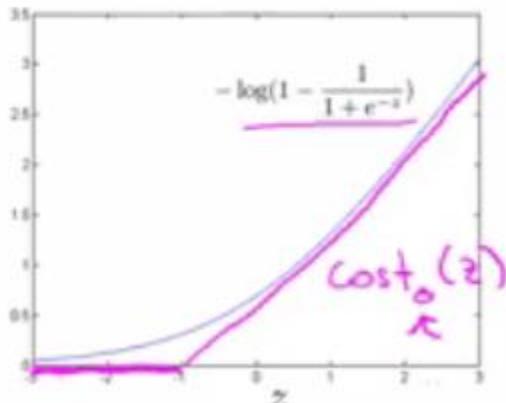
If $y = 0$, we want $h_{\theta}(x) \approx 0$, $\theta^T x \ll 0$



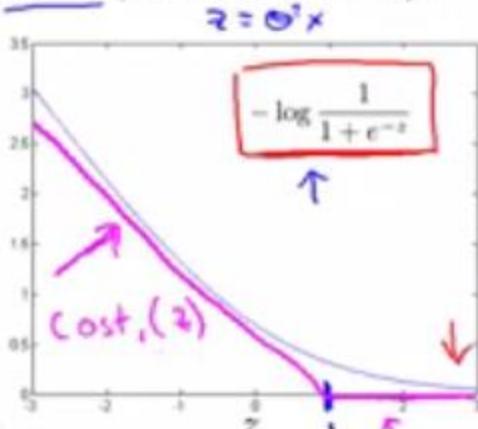
Cost of example:

$$\begin{aligned}
 & -(y \log h_\theta(x) + (1 - y) \log(1 - h_\theta(x))) \\
 &= -y \log \frac{1}{1 + e^{-\theta^T x}} - (1 - y) \log(1 - \frac{1}{1 + e^{-\theta^T x}})
 \end{aligned}$$

If $y = 0$ (want $\theta^T x \ll 0$):



If $y = 1$ (want $\theta^T x \gg 0$):



Logistic regression:

$$\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \underbrace{\left(-\log h_{\theta}(x^{(i)}) \right)}_{\text{cost}_1(\theta^T x^{(i)})} + (1-y^{(i)}) \underbrace{\left(-\log(1-h_{\theta}(x^{(i)})) \right)}_{\text{cost}_0(\theta^T x^{(i)})} \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Support vector machine:

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) + \frac{\lambda}{2m} \sum_{j=0}^n \theta_j^2$$

We can assume Logistic regression as $A + (\lambda) * B$

SVMs use a different way of regularization which is like $C * A + B$

$$\begin{array}{l} \cancel{A} + \cancel{\lambda} \cancel{B} \\ \rightarrow C \cancel{A} + \cancel{B} \end{array}$$

b. Hypothesis:

SVM Does not predict the Probability but gives directly 0 or 1

SVM hypothesis

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

Hypothesis:

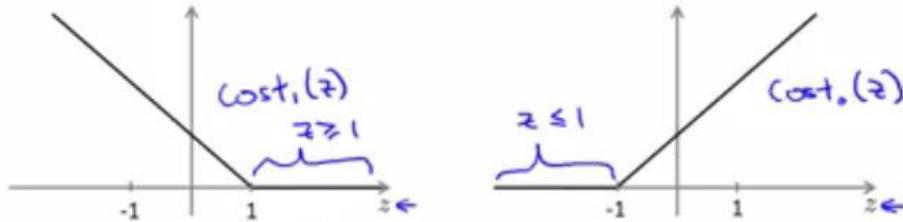
$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

2. Large Margin Intuition:

We need the prediction to be confident i.e. $>$ or < 0 is not enough we need > 1 or < -1

Support Vector Machine

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



- If $\underline{y = 1}$, we want $\underline{\theta^T x \geq 1}$ (not just ≥ 0)
- If $\underline{y = 0}$, we want $\underline{\theta^T x \leq -1}$ (not just < 0)

SVM Decision Boundary

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

Whenever $y^{(i)} = 1$:

$$\theta^T x^{(i)} \geq 1$$

$$\min_{\theta} C \rightarrow 0 + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

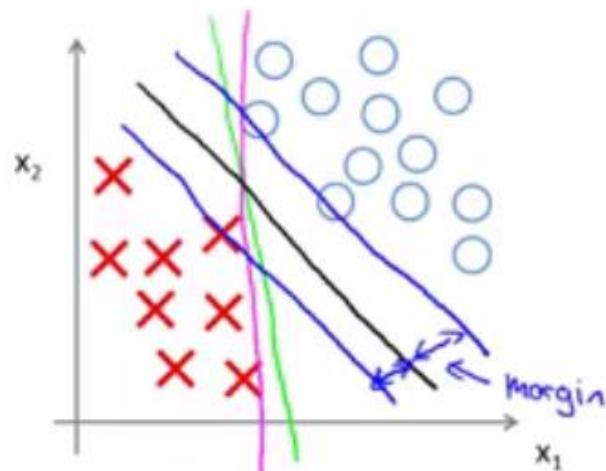
$$\text{st. } \theta^T x^{(i)} \geq 1 \quad \text{if } y^{(i)} = 1 \\ \theta^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0$$

Whenever $y^{(i)} = 0$:

$$\theta^T x^{(i)} \leq -1$$

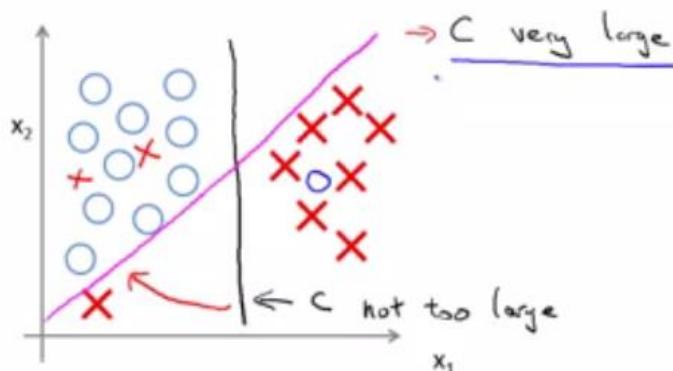
This Optimization of the cost function will lead to a large margin classifier as below
SVMs select the Black Decision Boundary but not the Magenta or the Green ones.

SVM Decision Boundary: Linearly separable case



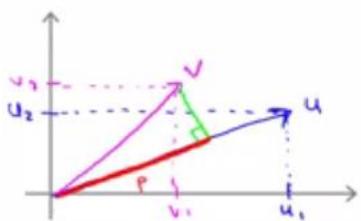
Large margin classifier

Large margin classifier in presence of outliers



3. Mathematics Behind Large Margin Classifier

Vector Inner Product



$$\Rightarrow u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad \Rightarrow v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$u^T v = ? \quad [u_1 \ u_2] \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$\|u\| = \text{length of vector } u$$

$$= \sqrt{u_1^2 + u_2^2} \in \mathbb{R}$$

p = length of projection of v onto u.

$$u^T v = \frac{p}{\|u\|} \cdot \|u\| \leftarrow = v^T u$$

Signed

$$= u_1 v_1 + u_2 v_2 \leftarrow p \in \mathbb{R}$$

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

$$p < 0$$

In the below diagram while computing the decision boundary we are trying to get the length of the theta vector.

SVM Decision Boundary

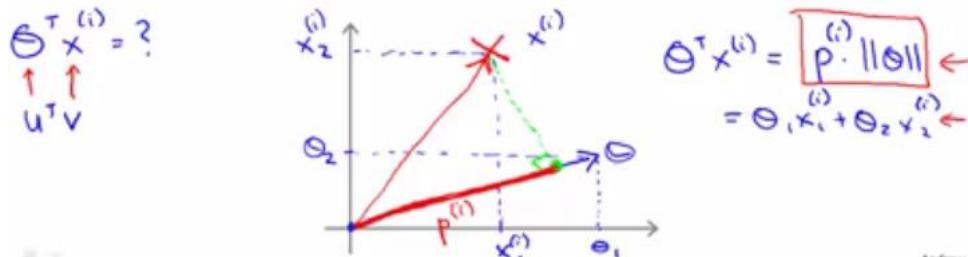
$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\Theta_1^2 + \Theta_2^2) = \frac{1}{2} \left(\sqrt{\Theta_1^2 + \Theta_2^2} \right)^2 = \frac{1}{2} \|\theta\|^2$$

s.t. $\boxed{\theta^T x^{(i)} \geq 1}$ if $y^{(i)} = 1$
 $\rightarrow \theta^T x^{(i)} \leq -1$ if $y^{(i)} = 0$

Simplification: $\Theta_0 = 0$, $n=2$

$$\|\theta\| = \sqrt{\Theta_1^2 + \Theta_2^2}$$

$$\begin{bmatrix} \Theta_1 \\ \Theta_2 \end{bmatrix}, \Theta_0 = 0$$



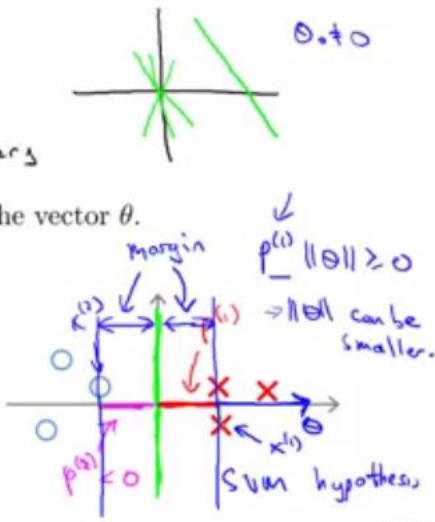
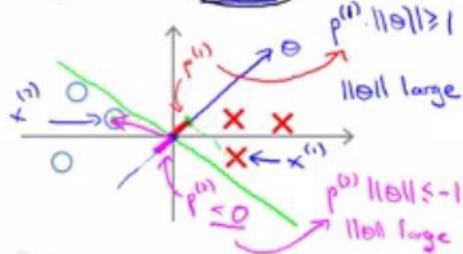
SVM Decision Boundary

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

s.t. $\begin{cases} p^{(i)} \cdot \|\theta\| \geq 1 & \text{if } y^{(i)} = 1 \\ p^{(i)} \cdot \|\theta\| \leq -1 & \text{if } y^{(i)} = -1 \end{cases}$

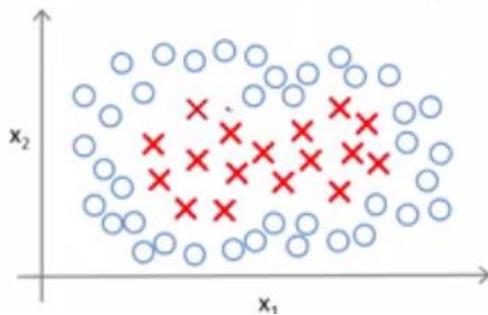
where $p^{(i)}$ is the projection of $x^{(i)}$ onto the vector θ .

Simplification: $\theta_0 = 0$



Kernels:

1. Non Linear decision boundary
 - a. We can add high order polynomial



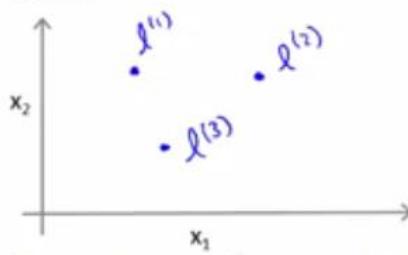
Using high order polynomial makes things computationally expensive.

$$\begin{aligned} \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 \\ + \theta_4 x_1^2 + \theta_5 x_2^2 + \dots \geq 0 \end{aligned}$$

- b. Instead of selecting new polynomial features

We can choose few points calling them landmarks now compute the new features based on the proximity to the landmarks.

Kernel



Given x , compute new feature depending on proximity to landmarks $l^{(1)}, l^{(2)}, l^{(3)}$

Given x :

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

$$f_2 = \text{similarity}(x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$$

$$f_3 = \text{similarity}(x, l^{(3)}) = \exp(\dots)$$

\curvearrowleft kernel (Gaussian kernels) $k(x, l^{(1)})$

In the below the new feature i.e. Similarity is the gaussian distribution peaking at landmarks

Kernels and Similarity

$$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^{(1)})^2}{2\sigma^2}\right)$$

If $x \approx l^{(1)}$:

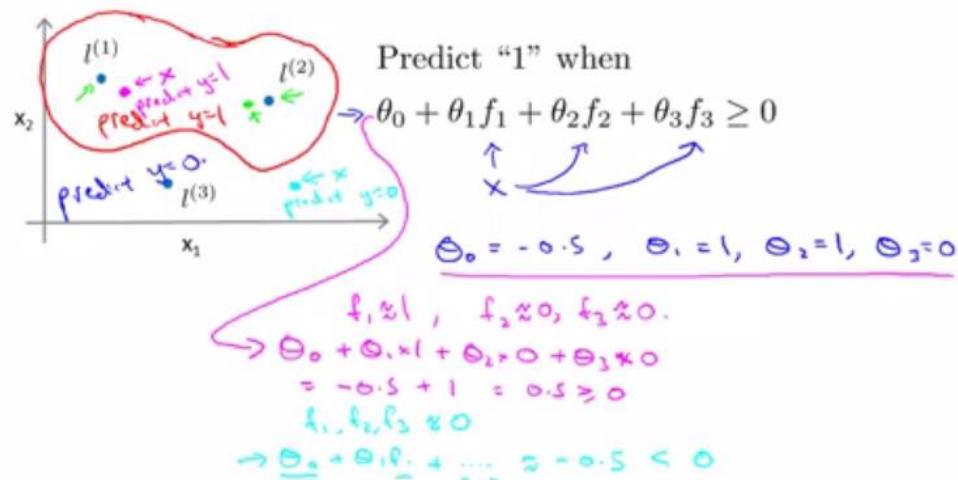
$$f_1 \underset{\uparrow}{\approx} \exp\left(-\frac{0}{2\sigma^2}\right) \approx 1$$

$l^{(1)} \rightarrow f_1$
 $l^{(2)} \rightarrow f_2$
 $l^{(3)} \rightarrow f_3$.

If x if far from $l^{(1)}$:

$$f_1 = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0.$$

Example:



c. Where do we get the landmarks from:

- i. Choose the landmarks from the training examples.

SVM with Kernels

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

Given example \underline{x} :

$$\begin{aligned} \rightarrow f_1 &= \text{similarity}(x, l^{(1)}) \\ \rightarrow f_2 &= \text{similarity}(x, l^{(2)}) \\ \dots \end{aligned}$$

$$f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} \quad f_0 = 1$$

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

$$\begin{aligned} \underline{x}^{(i)} \rightarrow f_1^{(i)} &= \sin(x^{(i)}, l^{(1)}) \\ f_2^{(i)} &= \sin(x^{(i)}, l^{(2)}) \\ \vdots & \\ f_m^{(i)} &= \sin(x^{(i)}, l^{(m)}) \end{aligned}$$

$$\begin{aligned} x^{(i)} \in \mathbb{R}^{n+1} &\quad (\text{or } \mathbb{R}^n) \\ f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \\ f_0^{(i)} &= 1 \end{aligned}$$

SVM with Kernels

Hypothesis: Given \underline{x} , compute features $f \in \mathbb{R}^{m+1}$ $\theta \in \mathbb{R}^{m+1}$

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

$$\begin{aligned} \text{Training:} \\ \rightarrow \min_{\theta} C \sum_{i=1}^m y^{(i)} \underbrace{\text{cost}_1(\theta^T f^{(i)})}_{\cancel{\theta^T f^{(i)}}} + (1 - y^{(i)}) \underbrace{\text{cost}_0(\theta^T f^{(i)})}_{\cancel{\theta^T f^{(i)}}} + \frac{1}{2} \sum_{j=1}^m \theta_j^2 \quad \begin{array}{l} \cancel{\theta^T f^{(i)}} \\ \cancel{\theta^T f^{(i)}} \end{array} \quad \begin{array}{l} n=m \\ \cancel{\theta^T f^{(i)}} \end{array} \quad \rightarrow \theta_0 \end{aligned}$$

$$\begin{aligned} - \sum_j \theta_j^2 &= \theta^T \theta \quad \leftarrow \theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_m \end{bmatrix} \quad (\text{ignoring } \theta_0) \\ - \theta^T M \theta \end{aligned}$$

Kernels can be applied to Logistic regression also.

But the SVMs are Numerically simple and this simplicity cannot be achieved in Logistic Regression.

The M above is used for Numerical optimization.

- d. SVM Parameters:

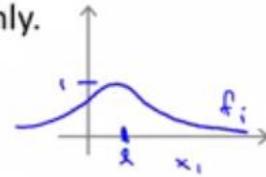
Larger sigma \Rightarrow Smoother curves \Rightarrow Less Variance

Smaller sigma \Rightarrow Sharp curves \Rightarrow High Variance

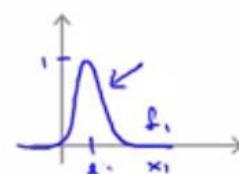
$C \left(= \frac{1}{\lambda} \right)$. \rightarrow Large C: Lower bias, high variance. (small λ)
 \rightarrow Small C: Higher bias, low variance. (large λ)

σ^2 Large σ^2 : Features f_i vary more smoothly.
 \rightarrow Higher bias, lower variance.

$$\exp \left(- \frac{\|x - l^{(i)}\|^2}{2\sigma^2} \right)$$



Small σ^2 : Features f_i vary less smoothly.
Lower bias, higher variance.



2. Using and SVM

a. Few things:

- i. Use SVM Software packages e.g. liblinear, libsvm
- ii. Need to specify
 - 1. Parameter C
 - 2. Kernel function
 - 3. Sigma square value

b. How to choose a kernel

i. No Kernel:

No kernel ("linear kernel")
Predict "y = 1" if $\theta^T x \geq 0$

$$\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \geq 0 \quad \rightarrow \underline{n} \text{ large}, \underline{m} \text{ small} \quad x \in \mathbb{R}^{n+1}$$

ii. Gaussian Kernel:

Gaussian kernel:

$$f_i = \exp \left(- \frac{\|x - l^{(i)}\|^2}{2\sigma^2} \right), \text{ where } l^{(i)} = x^{(i)}. \quad x \in \mathbb{R}^n, n \text{ small}$$

Need to choose σ^2 . and m large

c. Kernel Similarity functions

Kernel (similarity) functions:

function $f = \text{kernel}(x_1, x_2)$

$$f = \exp\left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right)$$

return

Note: Do perform feature scaling before using Gaussian kernel.

$$\begin{aligned} \|x - l\|^2 &\rightarrow v = x - l \\ \|v\|^2 &= v_1^2 + v_2^2 + \dots + v_n^2 \\ &= (x_1 - l_1)^2 + (x_2 - l_2)^2 + \dots + (x_n - l_n)^2 \\ &\quad \text{1000 feet}^2 \quad 1-5 \text{ bedrooms} \end{aligned}$$

One features magnitude will dominate other ones with smaller magnitude.

Note: Not all similarity functions $\text{similarity}(x, l)$ make valid kernels.
 (Need to satisfy technical condition called "Mercer's Theorem" to make sure SVM packages' optimizations run correctly, and do not diverge).

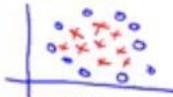
- d. Few Off Shelf kernels:
 - i. Polynomial kernel
 - ii. String kernel
 - iii. Chi-square kernel
 - iv. Histogram intersection kernel
- e. Multi-class classification
 - i. Use one vs all method
 - ii. Train K SVMs for K Class classification
- f. Logistic Regression vs SVMs:

n = number of features ($x \in \mathbb{R}^{n+1}$), m = number of training examples
If n is large (relative to m): (e.g. $n \geq m$, $n = 10,000$, $m = 10 \dots 1000$)
Use logistic regression, or SVM without a kernel ("linear kernel")

If n is small, m is intermediate: ($n = 1 \dots 1000$, $m = 10 \dots 10,000$) ←
→ Use SVM with Gaussian kernel

If n is small, m is large: ($n = 1 \dots 1000$, $m = 50,000+$)
→ Create/add more features, then use logistic regression or SVM without a kernel

Neural network likely to work well for most of these settings, but may be slower to train.



Optimization problem that SVM has is a Convex optimization problem and so the good SVM from software packages will find a global minimum.

We need not worry about local optimum in SVM.

Week 8

Wednesday, July 12, 2017
1:24 AM

Unsupervised Learning:

1. Few points
 - a. No Labels in the training set.
 - b. Find the structure in the training set.
 - c. Clustering Algorithm.
 - d. Applications of clustering
 - i. Market segmentation
 - ii. Social Network Analysis
 - iii. Organize computing clusters
 - iv. Astronomical data analysis

2. K-Means Clustering Algorithm:

Input:

- K (number of clusters) ←
- Training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ ←

$x^{(i)} \in \mathbb{R}^n$ (drop $x_0 = 1$ convention)

K-means algorithm

$$\mu_1 \quad \mu_2$$

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

Repeat {

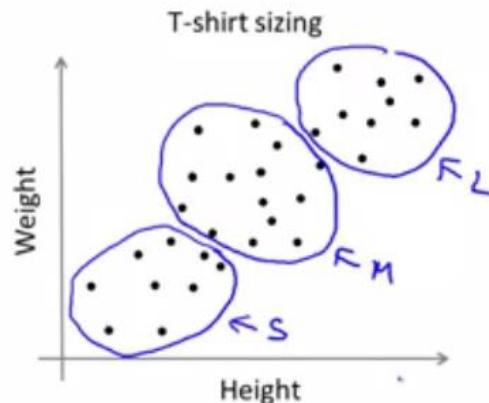
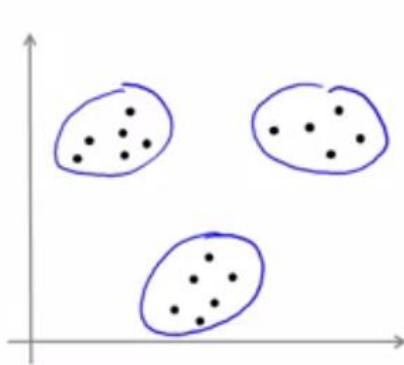
Cluster
Assignment
step

```

for i = 1 to m
    c(i) := index (from 1 to K) of cluster centroid
    closest to x(i)
        min ||x(i) - μk||2
        ↪ c(i)
for k = 1 to K
    → μk := average (mean) of points assigned to cluster k
        x(1), x(2), x(3), x(4) → c(1)=2, c(2)=2, c(3)=2, c(4)=2
    }
    μ2 =  $\frac{1}{4} [x^{(1)} + x^{(2)} + x^{(3)} + x^{(4)}] \in \mathbb{R}^n$ 
}

```

- a. If you get a cluster centroid with no values assigned to it
 - i. You can eliminate that class and use $K-1$ classification
 - ii. Else you can randomly initialize the
- b. K Means for not well separated Clusters



3. Optimization Objective:

K-means optimization objective

- $c^{(i)}$ = index of cluster ($1, 2, \dots, K$) to which example $x^{(i)}$ is currently assigned
- μ_k = cluster centroid k ($\mu_k \in \mathbb{R}^n$) $\leftarrow k \in \{1, 2, \dots, K\}$
- $\mu_{c^{(i)}}$ = cluster centroid of cluster to which example $x^{(i)}$ has been assigned $x^{(i)} \rightarrow S \quad c^{(i)} = S \quad \underline{\mu_{c^{(i)}}} = \mu_S$

Optimization objective:

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

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

Distortion

K-means algorithm

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

```

Repeat {
    Cluster assignment step
    Minimize  $J(\dots)$  w.r.t.  $(c^{(1)}, c^{(2)}, \dots, c^{(m)}) \leftarrow$ 
    (holding  $\mu_1, \dots, \mu_K$  fixed)
    for  $i = 1$  to  $m$ 
         $c^{(i)}$  := index (from 1 to  $K$ ) of cluster centroid
        closest to  $x^{(i)}$ 
    for  $k = 1$  to  $K$ 
         $\mu_k$  := average (mean) of points assigned to cluster  $k$ 
    }   minimize  $J(\dots)$  w.r.t.  $\mu_1, \dots, \mu_K$ 

```

4. Random Initialization:

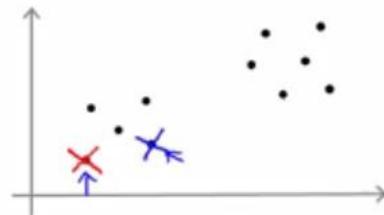
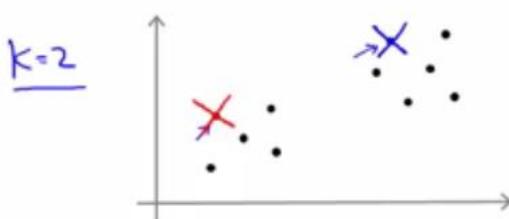
Random initialization

Should have $K < m$

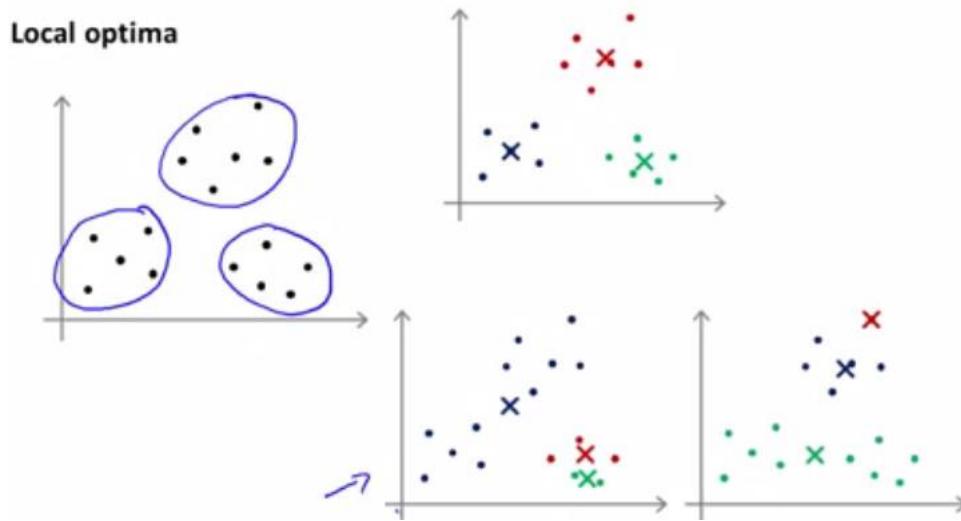
Randomly pick K training examples.

Set μ_1, \dots, μ_K equal to these K examples.

$$\begin{aligned} \mu_1 &= x^{(i)} \\ \mu_2 &= x^{(j)} \end{aligned}$$



Local Optima problem: This can happen when we are get unlucky random initialization. This is the local optima of the J Function i.e. distortion function.



When we K to be small then there are more chances that the random initialization will make a difference else when our K is very large it is highly probable that the first randomly selected values will give us a good optimization.

Random initialization

```
For i = 1 to 100 { s0 - 1000
    → Randomly initialize K-means.
    Run K-means. Get  $c^{(1)}, \dots, c^{(m)}$ ,  $\mu_1, \dots, \mu_K$ .
    Compute cost function (distortion)
    →  $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$ 
}
```

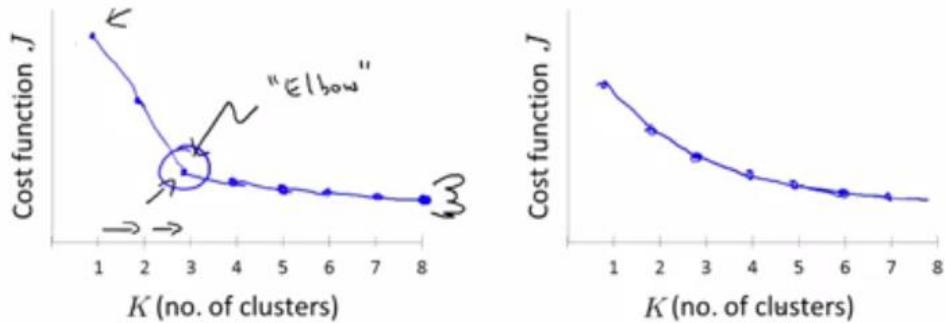
Pick clustering that gave lowest cost $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$
K=2-10

5. Choosing the number of Clusters:

- What is the value of K ?
 - We can group in many ways
- Elbow method:
 - Run K - means with multiple values and then calculate the cost function
 - Plot the Cost function vs K graph

- iii. This finds an elbow in the Graph i.e. the point where the distortion goes down rapidly and goes down slowly thereafter select that value of K
- iv. It's harder to location elbow in the graph on right.

Elbow method:



c. Choosing value of K:

Evaluate K-means based on a metric for how well it performs for later purpose.

On every iteration of K-means, the cost function $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$

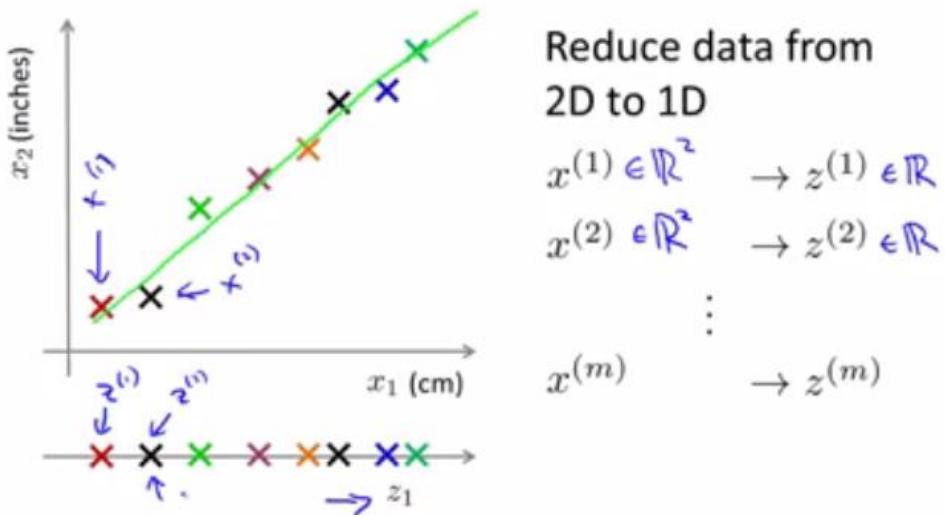
- d. (the distortion function) should either stay the same or decrease; in particular, it should not increase.

Principal Component Analysis:

1. Dimensionality Reduction:

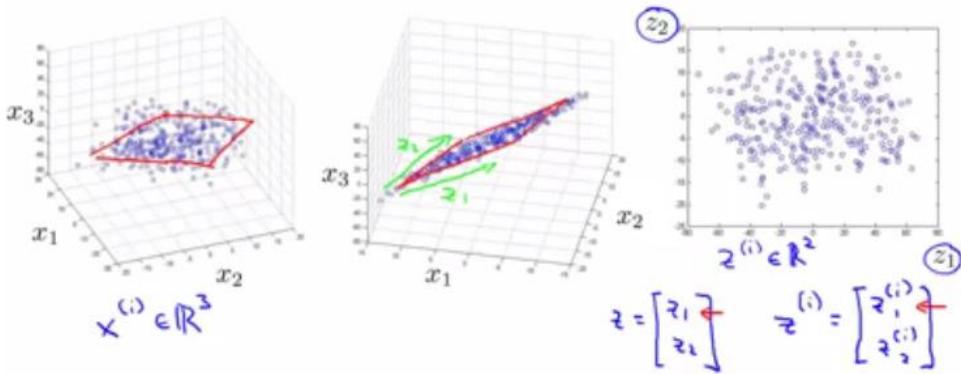
a. Data Compression

- i. We have redundant features.
- ii. Features need not be exactly linearly dependent.
- iii. Memory Requirement can be reduced.
- iv. Project data from 2D Plane to 1D Line



- i. Project data from 3D Space to 2D Plane.

Reduce data from 3D to 2D



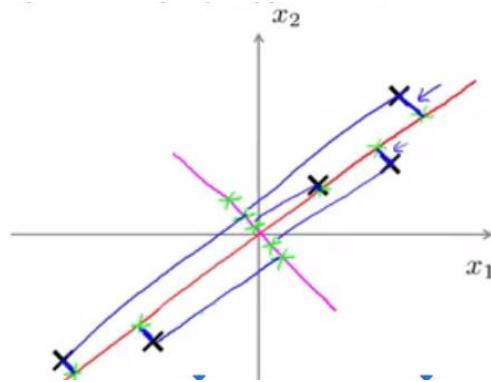
- b. Data Visualization:

- i. Often we convert data in high dimension space to 2D or 3D so that we can plot the data to visualize them

2. Principal Component Analysis Problem Formulation:

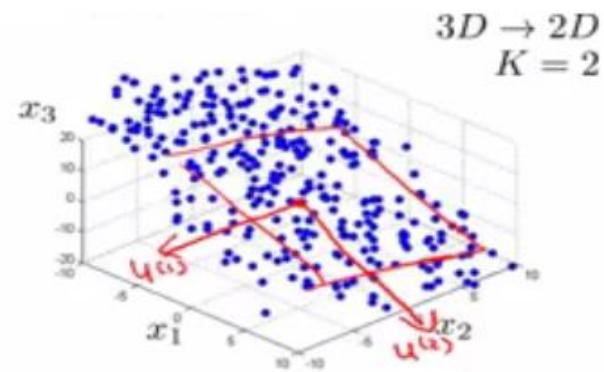
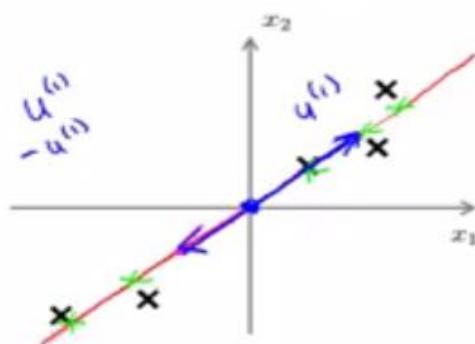
- a. Intro

- i. Minimize Projection Error



ii. Mean normalization and Feature Scaling is needed

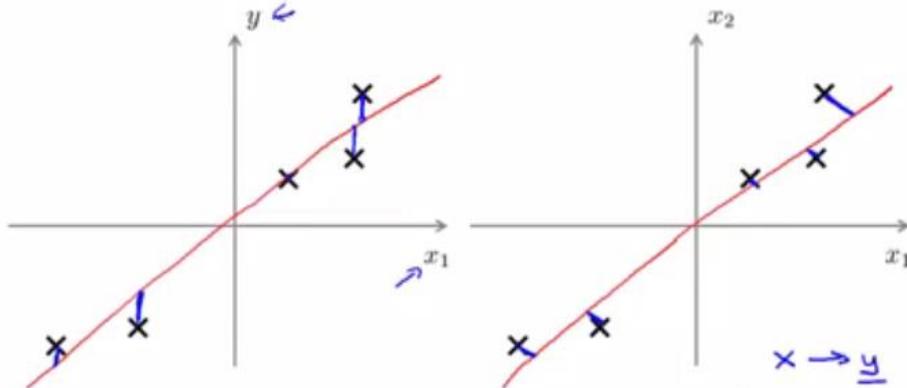
b. PCA Problem Formulation:



Reduce from 2-dimension to 1-dimension: Find a direction (a vector $u^{(1)} \in \mathbb{R}^n$) onto which to project the data so as to minimize the projection error.

Reduce from n-dimension to k-dimension: Find k vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project the data, so as to minimize the projection error.

c. PCA is not linear Regression:



d. PCA Algorithm:

i. Data Preprocessing:

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

Preprocessing (feature scaling/mean normalization):

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

Replace each $x_j^{(i)}$ with $x_j^{(i)} - \mu_j$.

If different features on different scales (e.g., x_1 = size of house, x_2 = number of bedrooms), scale features to have comparable range of values.

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{s_j}$$

ii. PCA

Reduce data from n -dimensions to k -dimensions

Compute "covariance matrix":

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

nxn Sigma

Compute "eigenvectors" of matrix Σ :

$$\rightarrow [U, S, V] = svd(\Sigma);$$

nxn matrix. → Singular value decomposition
eig (Sigma)

$$U = \begin{bmatrix} | & | & | & \dots & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(n)} \\ | & | & | & \dots & | \end{bmatrix}$$

U ∈ ℝ^{n × n}
u⁽¹⁾, ..., u⁽ⁿ⁾

From $[U, S, V] = svd(\Sigma)$, we get:

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

$$x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^k$$

$$z = \begin{bmatrix} | & | & | & \dots & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(n)} \\ | & | & | & \dots & | \end{bmatrix}^T x = \underbrace{\begin{bmatrix} -(u^{(1)})^T \\ \vdots \\ -(u^{(n)})^T \end{bmatrix}}_{k \times n} \underbrace{x}_{k \times 1}$$

U reduce

iii. Summary

Principal Component Analysis (PCA) algorithm summary

→ After mean normalization (ensure every feature has zero mean) and optionally feature scaling:

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

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

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

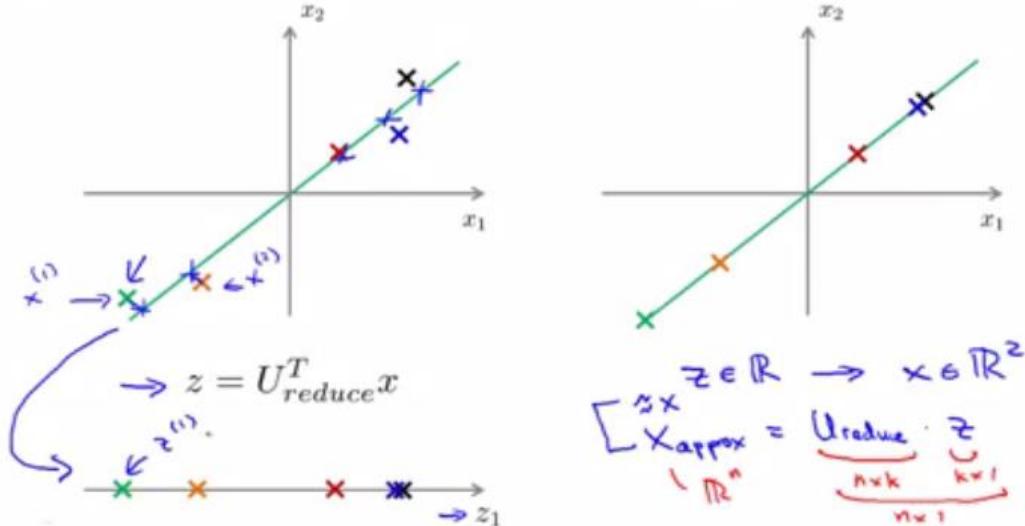
$$\rightarrow z = U_{\text{reduce}}' * x;$$

$$\uparrow \quad \uparrow \quad x \in \mathbb{R}^n \quad \cancel{x \in \mathbb{R}}$$

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

$$\rightarrow \text{Sigma} = (1/m) * X' * X;$$

e. Reconstruction from compressed Representation:



f. Choosing the number of Principal components:

Choosing \$k\$ (number of principal components)

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

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

Typically, choose \$k\$ to be smallest value so that

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

› “99% of variance is retained”

How to choose the value of K which gives the above variance.
Use the S Matrix that was computed with svd function

Choosing k (number of principal components)

Algorithm:

Try PCA with $k=1$ ~~$k \leq k-1$~~

Compute $U_{reduce}, z^{(1)}, z^{(2)}, \dots, z^{(m)}, x_{approx}^{(1)}, \dots, x_{approx}^{(m)}$

Check if

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

$k=17$

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

$$\rightarrow S = \begin{matrix} S_{11} & & \\ & S_{22} & \\ & & S_{33} \\ & \ddots & \ddots & \ddots & S_{nn} \end{matrix}$$

For given k

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

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

Andrew

Summary:

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

Pick smallest value of k for which

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

(99% of variance retained)

g. Advice for applying PCA:

i. Supervised Learning Speedup

Lower memory and computation power.

$$\underline{x^{(1)}, x^{(2)}, \dots, x^{(m)}} \in \underline{\mathbb{R}^{10000}} \leftarrow \downarrow \text{PCA}$$

$$\underline{z^{(1)}, z^{(2)}, \dots, z^{(m)}} \in \underline{\mathbb{R}^{1000}} \leftarrow$$

Note: Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA $\xrightarrow{x \rightarrow z}$ only on the training set. This mapping can be applied as well to the examples $x_{cv}^{(i)}$ and $x_{test}^{(i)}$ in the cross validation and test sets.

ii. BAD Use of PCA: To Prevent Overfitting

1. Use regularization instead of PCA when overfitting.

2. PCA does not use labels hence throws away some data when you are compressing
Without seeing the labels
- iii. PCA Should not be used to Design ML System
 - Design of ML system:**
 - - Get training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
 - - ~~Run PCA to reduce $x^{(i)}$ in dimension to get $z^{(i)}$~~
 - - Train logistic regression on $\{(z^{(1)}, y^{(1)}), \dots, (z^{(m)}, y^{(m)})\}$
 - - Test on test set: Map $x_{test}^{(i)}$ to $z_{test}^{(i)}$. Run $h_\theta(z)$ on $\{(z_{test}^{(1)}, y_{test}^{(1)}), \dots, (z_{test}^{(m)}, y_{test}^{(m)})\}$
 - › How about doing the whole thing without using PCA?
 - Before implementing PCA, first try running whatever you want to do with the original/raw data $x^{(i)}$. Only if that doesn't do what you want, then implement PCA and consider using $z^{(i)}$.

Week 9

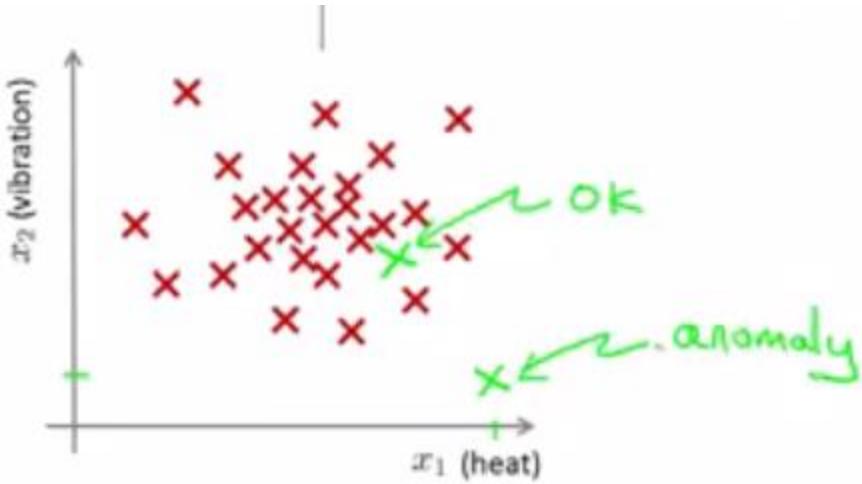
Monday, July 24, 2017
12:19 AM

Anomaly Detection:

Problem Motivation:

As a part of QA Testing of Aircraft engine we use features like
 $x_1 \Rightarrow$ heat detected
 $x_2 \Rightarrow$ vibration intensity

The problem here is that if the new engine with parameters should be sent for further testing does it look different than normal engines.

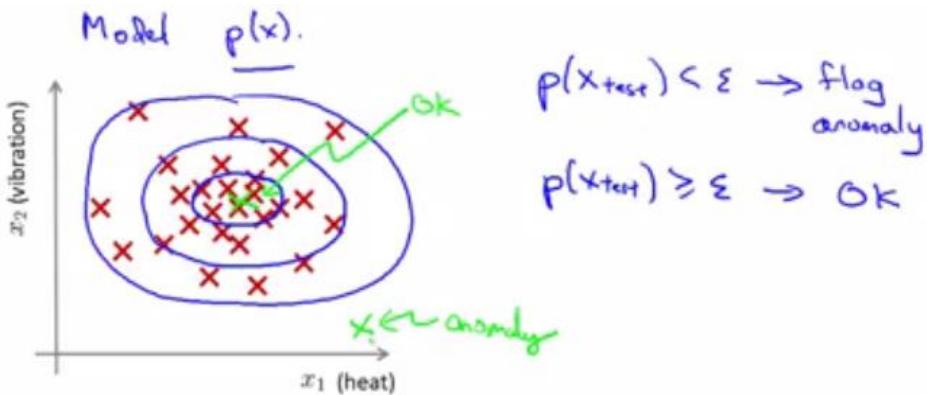


Given:

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

We need determine

Is x_{test} anomalous?



Example:

- i. Fraud Detection
- ii. Manufacturing QA
- iii. Monitoring Computers in Data Center

Gaussian Distribution: (Or Normal Distribution)

sigma	Standard deviation
u	Mean

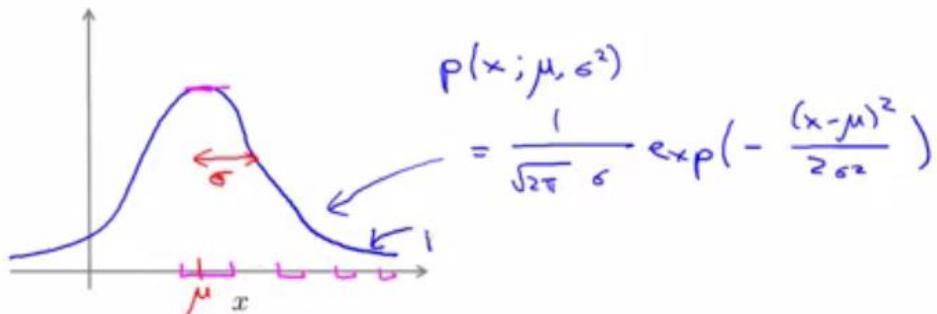
sig square	Variance
------------	----------

Probability of x taking the values away from mean reduces.

Say $x \in \mathbb{R}$. If x is distributed Gaussian with mean μ , variance σ^2 .

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

"distributed as"

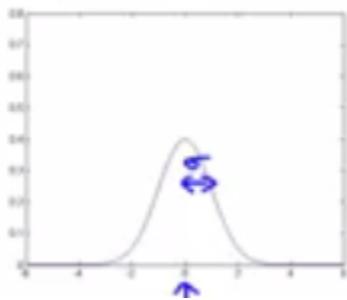


Example:

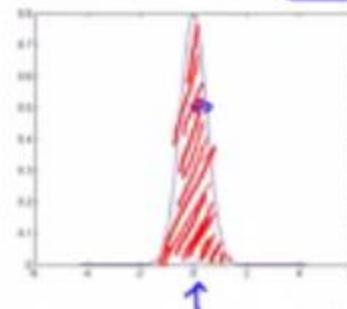
Area of curve must integrate to 1.

This is the property of the PDF.

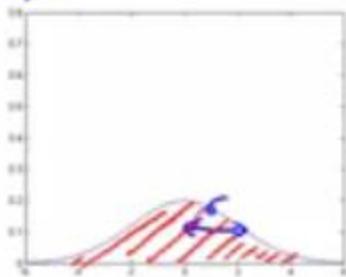
$$\Rightarrow \mu = 0, \sigma = 1$$



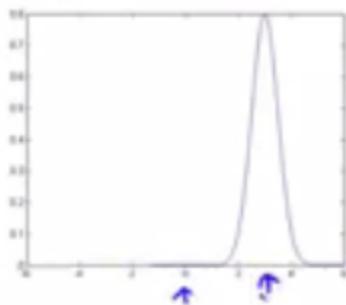
$$\Rightarrow \mu = 0, \sigma = 0.5$$



$\rightarrow \mu = 0, \sigma = 2$



$\rightarrow \mu = 3, \sigma = 0.5$



Parameter Estimation:

If we have to predict the μ and σ^2 from the dataset given. We have to estimate the gaussian distribution that the data came from.

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

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

These estimated are the maximum likelihood estimates of these parameters.

In Variance sometimes people use $1/(M-1)$. in ML people use M and M being large this does not matter much practically

Algorithm for Anomaly detection:

Density estimation

- Training set: $\{x^{(1)}, \dots, x^{(m)}\}$
- Each example is $x \in \mathbb{R}^n$

We assume that features are distributed with gaussian distributions with some mean and variance.

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

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

$$x_3 \sim \mathcal{N}(\mu_3, \sigma_3^2)$$

$$\begin{aligned} p(x) &= p(x_1; \mu_1, \sigma_1^2) p(x_2; \mu_2, \sigma_2^2) p(x_3; \mu_3, \sigma_3^2) \dots p(x_n; \mu_n, \sigma_n^2) \\ &= \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) \end{aligned}$$

$\sum_{i=1}^n i = 1+2+3+\dots+n$
 $\prod_{i=1}^n i = 1 \times 2 \times 3 \times \dots \times n$

Putting it all together:

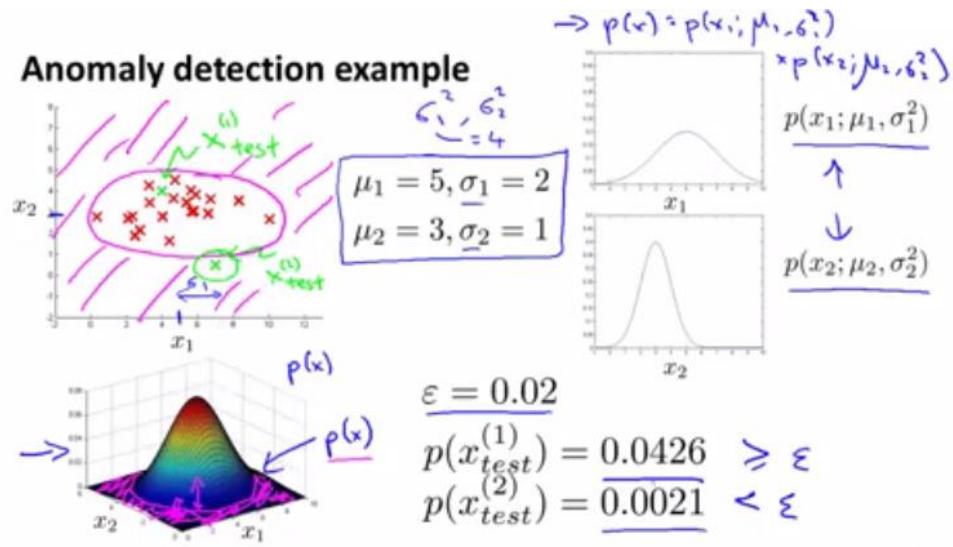
Anomaly detection algorithm

- Choose features x_i that you think might be indicative of anomalous examples. $\{x^{(1)}, \dots, x^{(n)}\}$
- Fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$
 - $\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$
 - $\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$
- Given new example x , compute $p(x)$:

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

Anomaly if $p(x) < \epsilon$

Example:



Developing and Evaluating the algorithm:

Algorithm evaluation

- Fit model $p(x)$ on training set $\{x^{(1)}, \dots, x^{(m)}\}$ $(x_{test}^{(i)}, y_{test}^{(i)})$
 - On a cross validation/test example x , predict
- $$y = \begin{cases} 1 & \text{if } p(x) < \epsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \epsilon \text{ (normal)} \end{cases} \quad y=0$$

Possible evaluation metrics:

- - True positive, false positive, false negative, true negative
- - Precision/Recall
- - F_1 -score

Can also use cross validation set to choose parameter ϵ

Anomaly Detection vs Supervised Learning:

Anomaly detection	vs.	Supervised learning
<ul style="list-style-type: none"> → Very small number of positive examples ($y = 1$). (<u>0-20</u> is common). → Large number of negative ($y = 0$) examples. $p(x)$ ← → Many different "types" of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like; future anomalies may look nothing like any of the anomalous examples we've seen so far. 		<p>Large number of positive and negative examples. ←</p> <p>Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set.</p>

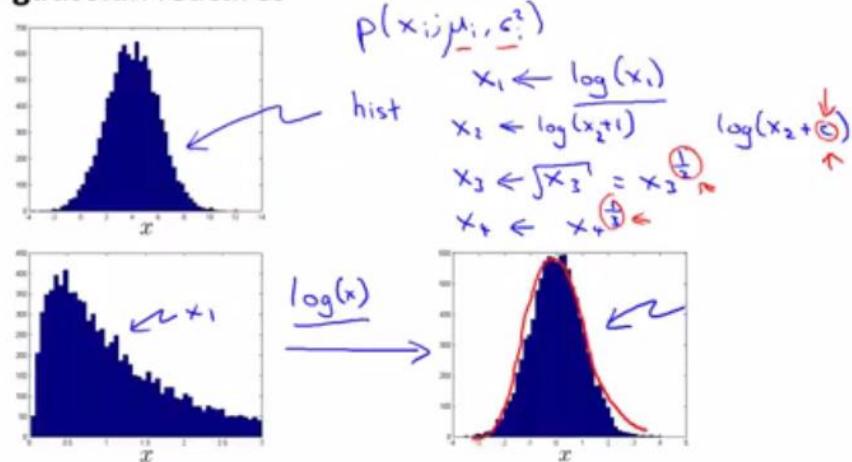
Anomaly detection	vs.	Supervised learning
<ul style="list-style-type: none"> → • Fraud detection $y=1$ → • Manufacturing (e.g. aircraft engines) → • Monitoring machines in a data center ⋮ 	→	<ul style="list-style-type: none"> • Email spam classification • Weather prediction (sunny/rainy/etc). • Cancer classification ⋮

Choosing what features to Use:

Even in case the data looks not gaussian the algorithms work fine.
Its better though to convert the data to look like gaussian distribution.

You can use some transformations to make the data into gaussian.

Non-gaussian features



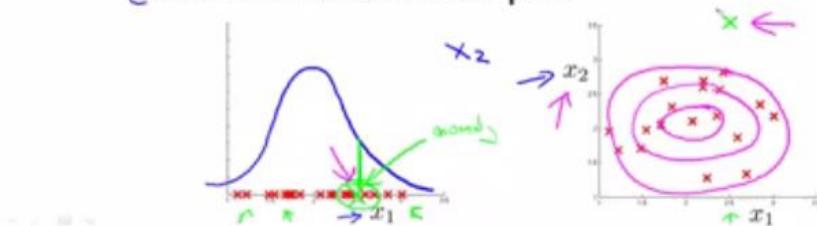
When you are getting an anomaly in the $p(x) > \epsilon p$. There you might try to more features that will get the anomaly.

→ Error analysis for anomaly detection

- [Want $p(x)$ large for normal examples x .
- [$p(x)$ small for anomalous examples x .

Most common problem:

- [$p(x)$ is comparable (say, both large) for normal and anomalous examples



We will have to select the features that might help us distinguish the anomalies from the rest.

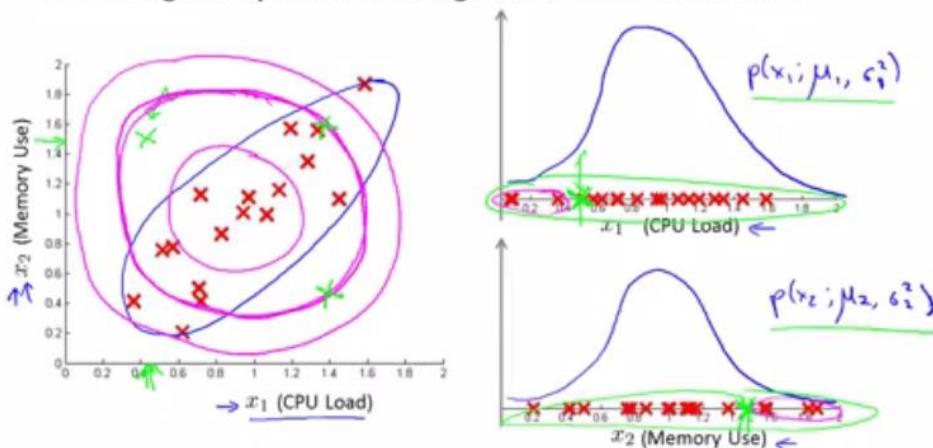
- Monitoring computers in a data center
- Choose features that might take on unusually large or small values in the event of an anomaly.
- x_1 = memory use of computer
- x_2 = number of disk accesses/sec
- x_3 = CPU load ←
- x_4 = network traffic ←

$$x_5 = \frac{\text{CPU load}}{\text{network traffic}}$$

$$x_6 = \frac{(\text{CPU load})^2}{\text{network traffic}}$$

Multivariate Gaussian Distribution:

Motivating example: Monitoring machines in a data center



Multivariate Gaussian (Normal) distribution

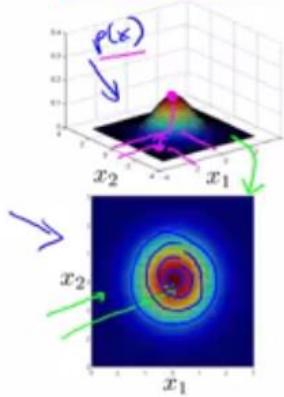
- $x \in \mathbb{R}^n$. Don't model $p(x_1), p(x_2), \dots$, etc. separately.
Model $p(x)$ all in one go.
- Parameters: $\mu \in \mathbb{R}^n, \Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix)

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

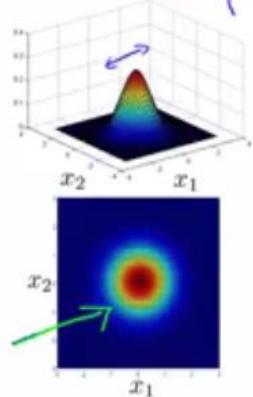
$|\Sigma| = \text{determinant of } \Sigma \quad | \det(\text{Signal})$

Multivariate Gaussian (Normal) examples

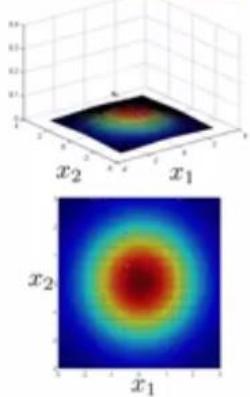
$$\rightarrow \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



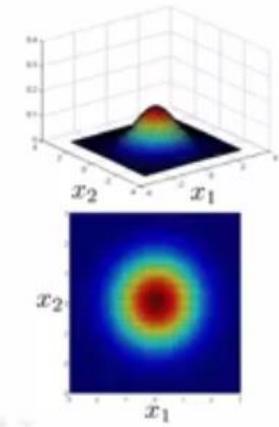
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 0.6 \end{bmatrix}$$



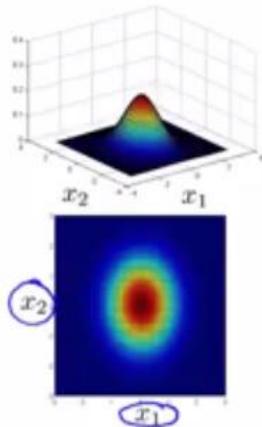
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$



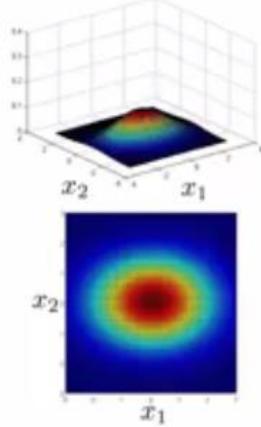
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



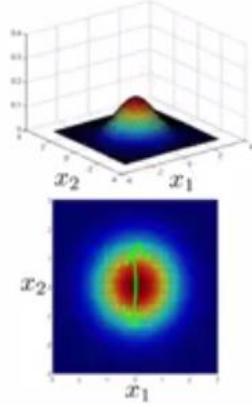
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 1 \end{bmatrix}$$



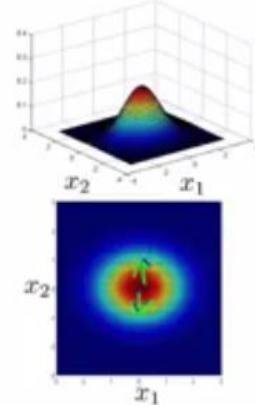
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$



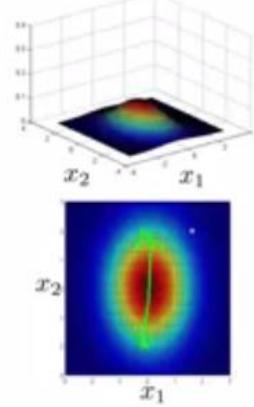
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 0.6 \end{bmatrix}$$

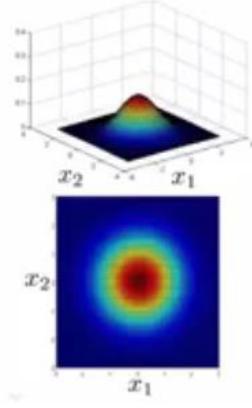


$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$

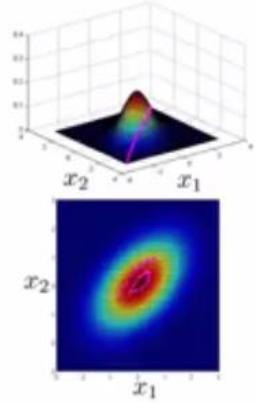


Off diagonal elements allow us to model the correlation

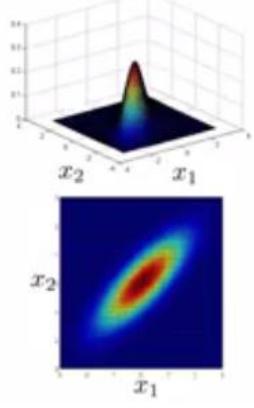
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



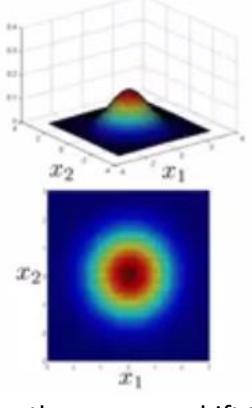
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$



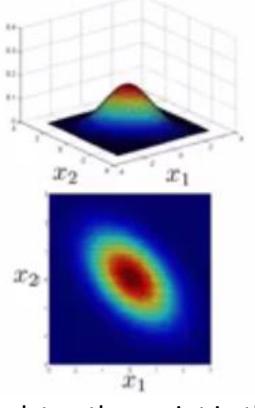
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$



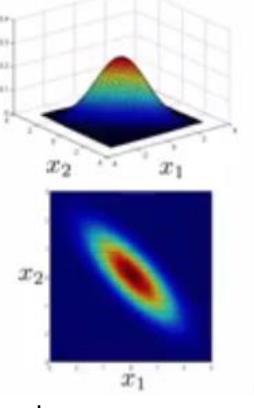
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



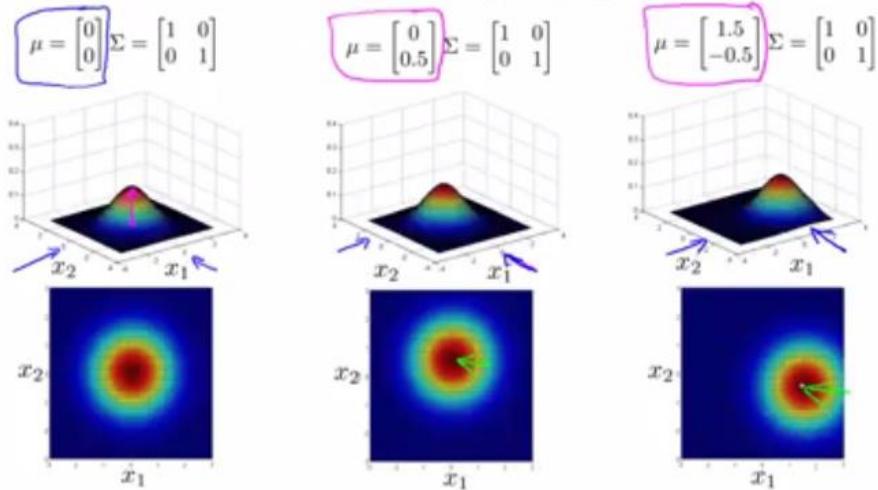
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & -0.5 \\ 0.5 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}$$



Varying the Σ we can shift the peak to other point in the graph

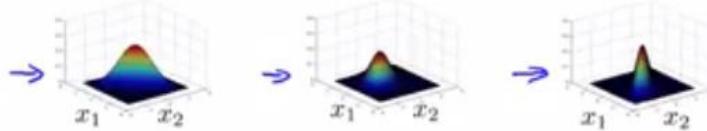


Anomaly Detection with Multi variate gaussian distribution:

Multivariate Gaussian (Normal) distribution

Parameters μ, Σ $\mu \in \mathbb{R}^n$ $\Sigma \in \mathbb{R}^{n \times n}$

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



Parameter fitting:

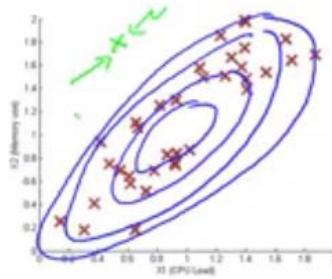
Given training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \leftarrow \times \in \mathbb{R}^n$

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

Anomaly detection with the multivariate Gaussian

1. Fit model $p(x)$ by setting

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



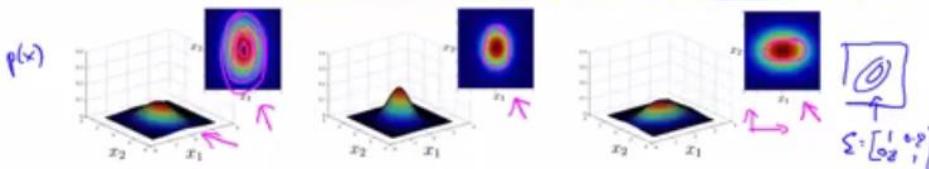
2. Given a new example x , compute

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

Flag an anomaly if $\underline{p(x) < \varepsilon}$

Relationship to original model

Original model: $p(x) = p(x_1; \mu_1, \sigma_1^2) \times p(x_2; \mu_2, \sigma_2^2) \times \cdots \times p(x_n; \mu_n, \sigma_n^2)$



Corresponds to multivariate Gaussian

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

where $\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix}$



\rightarrow Original model

vs. \rightarrow Multivariate Gaussian

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

Manually create features to capture anomalies where x_1, x_2 take unusual combinations of values.

$$\rightarrow x_3 = \frac{x_1}{x_2} = \frac{\text{CPU load}}{\text{memory}}$$

\rightarrow Computationally cheaper (alternatively, scales better to large n) $n=10,000, m=100,000$

OK even if m (training set size) is small

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

\rightarrow Automatically captures correlations between features

$$\Sigma \in \mathbb{R}^{n \times n} \quad \Sigma^{-1}$$

Computationally more expensive

$$\Sigma$$

Must have $\underline{m > n}$ or else Σ is non-invertible. $m \geq 10n$

People generally capture extra features and use the original model instead of Multivariate Gaussian Model.

If you find the sigma matrix to be singular or non-invertible

- a. Check $m \geq n$ or say $m \geq 10n$
- b. You might have redundant features
 - i. $x_1 = x_2$
 - ii. $x_3 = x_1 + x_2$

Features are Linearly independent.

Recommender Systems:

Problem Formulation;

Example: Predicting movie rating:

Given the Matrices Below r and y we have to find the missing data i.e. here the data of the movies not watched.

Example: Predicting movie ratings

→ User rates movies using one to five stars
zero

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)
Love at last	5	5	0	0
Romance forever	5	?	0	0
Cute puppies of love	?	5	0	?
Nonstop car chases	0	0	5	4
Swords vs. karate	0	0	5	?

$n_u = 4$ $n_m = 5$

→ n_u = no. users
→ n_m = no. movies
 $r(i, j) = 1$ if user j has rated movie i
 $y^{(i,j)}$ = rating given by user j to movie i (defined only if $r(i, j) = 1$)

Here for this analysis we use user rating from 0 to all starts. Generally ppl user from 1 - all

Content Based Recommendations:

Content-based recommender systems

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	$n_u = 4, n_m = 5$	$x^{(1)} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix}$
Love at last 1	5	5	0	0	x_1 (romance)	0.9 → 0
Romance forever 2	5	?	?	0	x_2 (action)	1.0 → 0.01
Cute puppies of love 3	?	4	0	?		0.99 → 0
Nonstop car chases 4	0	0	5	4		0.1 → 1.0
Swords vs. karate 5	0	0	5	?		0 → 0.9

$\theta^{(j)} \in \mathbb{R}^3$

For each user j , learn a parameter $\theta^{(j)} \in \mathbb{R}^3$. Predict user j 's rating for movie i with $(\theta^{(j)})^T x^{(i)}$ stars.

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

Problem formulation

→ $r(i, j) = 1$ if user j has rated movie i (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: $\underline{(\theta^{(j)})^T x^{(i)}} \quad \theta^{(j)} \in \mathbb{R}^{n+1}$

→ $m^{(j)}$ = no. of movies rated by user j

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

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

We can ignore m_j as optimization does not change.

Optimization objective:

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

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

To learn $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n_u)}$:

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

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

Optimization algorithm:

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

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

Gradient descent update:

$$\begin{aligned} \theta_k^{(j)} &:= \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} \quad (\text{for } k = 0) \\ \theta_k^{(j)} &:= \theta_k^{(j)} - \alpha \left(\sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right) \quad (\text{for } k \neq 0) \end{aligned}$$

$\frac{\partial}{\partial \theta_k^{(j)}} \mathcal{J}(\theta^{(1)}, \dots, \theta^{(n_u)})$

Collaborative Filtering:

Problem Motivation:

It is difficult to get the feature vector for movies.

Instead we get the taste of the user from him.

We now infer the values of the feature vectors of movies.

Movie	Alice (1) $\theta^{(1)}$	Bob (2) $\theta^{(2)}$	Carol (3) $\theta^{(3)}$	Dave (4) $\theta^{(4)}$	x_1 (romance)	x_2 (action)
Love at last	5	5	0	0	?	?
Romance forever	5	?	?	0	?	?
Cute puppies of love	?	4	0	?	?	?
Nonstop car chases	0	0	5	4	?	?
Swords vs. karate	0	0	5	?	?	?

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

Now that we know the values of theta we predict the values of x using the rating that the user has given to the movie.

$x^{(i)}$

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

$$(\theta^{(2)})^T x^{(i)} \approx 5$$

$$(\theta^{(3)})^T x^{(i)} \approx 0$$

$$\theta^{(4)})^T x^{(i)} \approx 0$$

Implies

$$x^{(i)} = \begin{bmatrix} 1 \\ 1.0 \\ 0.0 \\ 0.0 \end{bmatrix}$$

Optimization Algorithm:

Optimization algorithm

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

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

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

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

Collaborative filtering

Given $x^{(1)}, \dots, x^{(n_m)}$ (and movie ratings),
can estimate $\theta^{(1)}, \dots, \theta^{(n_u)}$

Given $\theta^{(1)}, \dots, \theta^{(n_u)}$,
can estimate $x^{(1)}, \dots, x^{(n_m)}$

Guess $\Theta \rightarrow x \rightarrow \Theta \rightarrow x \rightarrow \Theta \rightarrow x \rightarrow \dots$

Collaborative Filtering Optimization Objective:

Collaborative filtering optimization objective $(i_j) : r(i_j) \in \mathbb{R}$

Given $x^{(1)}, \dots, x^{(n_m)}$, estimate $\theta^{(1)}, \dots, \theta^{(n_u)}$:

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

Given $\theta^{(1)}, \dots, \theta^{(n_u)}$, estimate $x^{(1)}, \dots, x^{(n_m)}$:

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

Minimizing $x^{(1)}, \dots, x^{(n_m)}$ and $\theta^{(1)}, \dots, \theta^{(n_u)}$ simultaneously:

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

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

$\Theta \rightarrow x \rightarrow \Theta \rightarrow x \rightarrow \dots$

Andrew Ng

Collaborative filtering algorithm

- 1. 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). E.g. for every $j = 1, \dots, n_u, i = 1, \dots, n_m$:
- $$x_k^{(i)} := x_k^{(i)} - \alpha \left(\sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)} \right) \quad \frac{\partial J}{\partial x_k^{(i)}}$$
- $$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right) \quad \frac{\partial J}{\partial \theta_k^{(j)}}$$
3. For a user with parameters θ and a movie with (learned) features x , predict a star rating of $\theta^T x$.

Vectorization: Low Rank Matrix Factorization

Collaborative filtering

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	
Love at last	5	5	0	0	
Romance forever	5	?	?	0	
Cute puppies of love	?	4	0	?	
Nonstop car chases	0	0	5	4	
Swords vs. karate	0	0	5	?	
	↑	↑	↑	↑	

$n_m = 5$
 $n_u = 4$

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$$

$y^{(i,j)}$

Collaborative filtering

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$$

Predicted ratings: $(\theta^{(i)})^T (x^{(j)})$

$$\begin{bmatrix} (\theta^{(1)})^T (x^{(1)}) & (\theta^{(2)})^T (x^{(1)}) & \dots & (\theta^{(n_u)})^T (x^{(1)}) \\ (\theta^{(1)})^T (x^{(2)}) & (\theta^{(2)})^T (x^{(2)}) & \dots & (\theta^{(n_u)})^T (x^{(2)}) \\ \vdots & \vdots & \vdots & \vdots \\ (\theta^{(1)})^T (x^{(n_m)}) & (\theta^{(2)})^T (x^{(n_m)}) & \dots & (\theta^{(n_u)})^T (x^{(n_m)}) \end{bmatrix}$$

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

Low rank matrix factorization

Finding related movies

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

$\rightarrow x_1 = \text{romance}, x_2 = \text{action}, x_3 = \text{comedy}, x_4 = \dots$

How to find $\underline{x}^{(j)}$ related to $\underline{x}^{(i)}$?

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

5 most similar movies to movie i :

\rightarrow Find the 5 movies j with the smallest $\|\underline{x}^{(i)} - \underline{x}^{(j)}\|$.

Implementation Detail: Mean Normalization

Users who have not rated any movies

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	Eve (5)
Love at last	5	5	0	0	?
Romance forever	5	?	?	0	?
Cute puppies of love	?	4	0	?	?
Nonstop car chases	0	0	5	4	?
Swords vs. karate	0	0	5	?	?

\downarrow

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix}$$

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

$n=2 \quad \underline{\theta}^{(s)} \in \mathbb{R}^2 \quad \underline{\theta}^{(s)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \frac{\lambda}{2} [\underline{\theta}_1^{(s)}]^2 + [\underline{\theta}_2^{(s)}]^2 \leq$

$(\underline{\theta}^{(s)})^T \underline{x}^{(i)} = 0$

Mean Normalization:

Mean Normalization:

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? & 2.5 \\ 5 & ? & ? & 0 & ? & 2.5 \\ ? & 4 & 0 & ? & ? & 2 \\ 0 & 0 & 5 & 4 & ? & ? \\ 0 & 0 & 5 & 0 & ? & ? \end{bmatrix}$$

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

For user j , on movie i predict:

$$\rightarrow (\theta^{(s)})^T (x^{(i)}) + \mu_i$$

\downarrow
learn $\theta^{(s)}, \mu_i$

User 5 (Eve):

$$\underline{\theta^{(s)}} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\underbrace{(\theta^{(s)})^T (x^{(i)})}_{\sim} + \boxed{\mu_i}$$

Week 10

Monday, July 24, 2017
12:27 AM

Gradient Descent with Large Datasets

“It’s not who has the best algorithm that wins.
It’s who has the most data.”

Large data => More computation.

How can you say that using more data will improve performance

Plot a learning curve for a range of values of m and verify that the algorithm has high variance when m is small.

Stochastic Gradient Descent:

Repeat {
 $\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$
 (for every $j = 0, \dots, n$)
 }
 $M = 300,000,000$

The highlighted term above will be very computer intensive.

This is called Batch Gradient Descent.

One step of the gradient descent needs us to go through all the data.

We need not go through all the data for making a gradient descent update. We do this with every example.

Every iteration will be faster and each iteration can move towards or away from the global minimum.

We end up going close to global min and keep moving there.

Batch gradient descent

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

Stochastic gradient descent

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

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

1. Randomly shuffle dataset.

2. Repeat {

for $i=1, \dots, m$ {

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

} (for $j=0, \dots, n$)

$$\hookrightarrow \frac{\partial}{\partial \Theta_j} cost(\theta, (x^{(i)}, y^{(i)}))$$

Mini Batch Gradient Descent:

Batch gradient descent: Use all m examples in each iteration

Stochastic gradient descent: Use 1 example in each iteration

Mini-batch gradient descent: Use b examples in each iteration

$b = \text{mini-batch size.}$ $b = 10.$ $2-100$

Mini-batch gradient descent

Say $b = 10, m = 1000$.

Repeat {

- for $i = 1, 11, 21, 31, \dots, 991$ {
$$\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_\theta(x^{(k)}) - y^{(k)}) x_j^{(k)}$$

(for every $j = 0, \dots, n$)

}

}

Makes more progress over the batch gradient decent.

Why b why not 1 ? Ans: Vectorization will be faster. => Parallelization.

Takes advantage for the vectorization.

Needs to tune one more parameter b.

Stochastic Gradient Descent Convergence:

Week 11

Sunday, August 06, 2017

2:22 PM

Application Example:

Problem Description and Pipeline:

Photo OCR - Photo Optical Character Recognition.

- i. Given an Image photo. The algorithm should detect and read the text in the image.
- ii. This could be helpful for the build people
- iii. Also for the Car Navigation.

Phot OCR Pipeline:

Photo OCR pipeline

1. Text detection

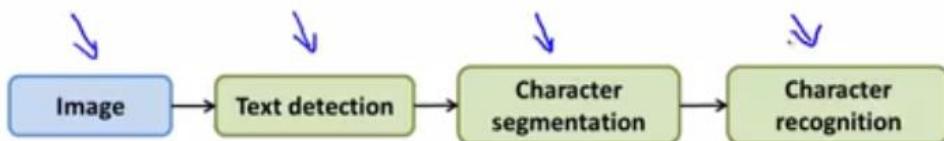


2. Character segmentation

ANTIQUE MALL

3. Character classification

A N T



Sliding Windows:

Text detection:

- i. Text aspect ratio varies.
- ii. We can relate this to pedestrian detection in the image.
 - 1. This is simpler we can use fixed aspect ratio for different pedestrians.
 - 2. They are at different distances but the aspect ratio i.e. height to width ratio will be same

Supervised learning can be used for pedestrian detection. For checking if the image has the pedestrian or not.

Sliding Window is traversed over the image to check if the window has the pedestrian or not

Sliding window detection



We generally take a larger patch and resize it to detect the pedestrians.

Similar to pedestrian detection we can get positive examples for images with and without test.

- i. We now train model a supervised learning.
- ii. Do sliding window detection.

We need to draw rectangle around the text regions. We use the Expanders by checking the proximities of the pixels and make them in one rectangles.

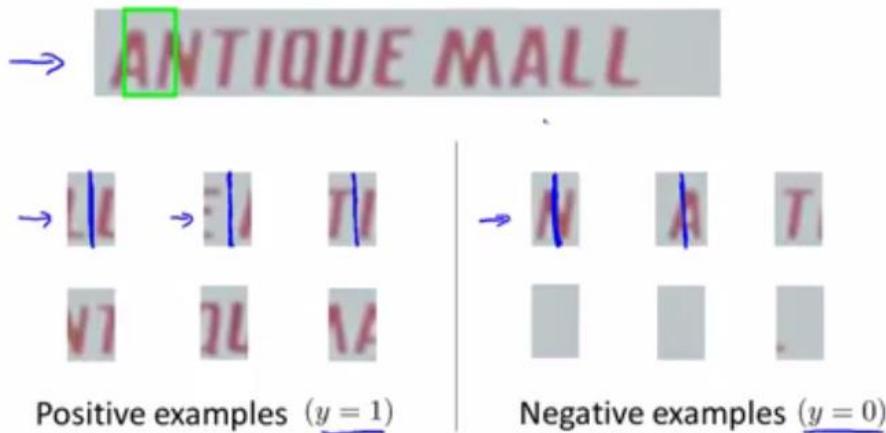


We can discard few rectangles based on the aspect ratios that does not look normal for the test boxes.

Character Segmentation:

Get positive and negative examples for detecting the split between the characters. Now we can train another supervised learning algorithm

1D Sliding window for character segmentation



Character Classification:

We can apply Neural Nets or etc. to classify the images to alphabets.

Getting lots of Data and Artificial Data:

We have two main variations:

- i. Generate new data from scratch
- ii. Get data from already existing data.

For images: (Photo OCR)

- i. Color does not matter for text detection here.
- ii. We can get different font styles here.
- iii. We can change the background of the character.
- iv. We can use Scaling or Cropping or Rotations.
- v. We can introduce distortions into the already existing images.

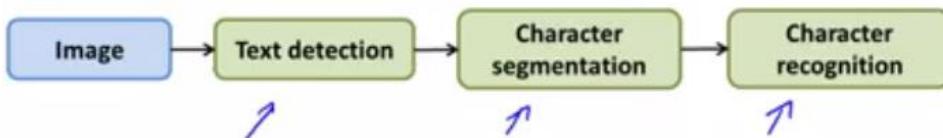
Discussion on getting more data

1. Make sure you have a low bias classifier before expending the effort. (Plot learning curves). E.g. keep increasing the number of features/number of hidden units in neural network until you have a low bias classifier.
2. "How much work would it be to get 10x as much data as we currently have?"
 - Artificial data synthesis
 - Collect/label it yourself
 - "Crowd source" (E.g. Amazon Mechanical Turk)

Ceiling Analysis: What part of Pipeline to Work on Next

We must find the part of the work that will give maximum performance improvement to focus the efforts on.

Estimating the errors due to each component (ceiling analysis)



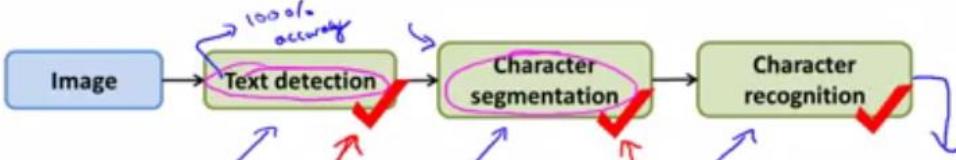
What part of the pipeline should you spend the most time trying to improve?

We will find the accuracy of the system by going through each block at a time.

We will assume that all the blocks earlier in the pipeline are working with 100% accuracy while I am checking the current block.

We can simply use the labels as predictions to act 100% accuracy of the previous blocks.

This way we will find the accuracy for the entire system every time by focusing on the individual block.



What part of the pipeline should you spend the most time trying to improve?

Component	Accuracy
Overall system	72%
Text detection	89%
Character segmentation	90%
Character recognition	100%

Annotations in blue ink show arrows pointing from the overall system accuracy down to each component's accuracy, with a downward arrow between the overall system and text detection.

Summary:

- Supervised Learning $(x^{(i)}, y^{(i)})$
 - Linear regression, logistic regression, neural networks, SVMs
- Unsupervised Learning $x^{(i)}$
 - K-means, PCA, Anomaly detection
- Special applications/special topics
 - Recommender systems, large scale machine learning.
- Advice on building a machine learning system
 - Bias/variance, regularization; deciding what to work on next: evaluation of learning algorithms, learning curves, error analysis, ceiling analysis.

Further reading

Monday, May 29, 2017

3:43 PM

- Andrew Ng on [L1 vs L2 regularization](#)
 - **Summary:** L2 is the "popular" technique for regularization. However, in this paper, Andrew Ng argues that L1 is often more efficient, taking logistic regression as an example; given (among other factors) that it offers "automatic" feature selection.
 - **Note:** The main difference is that in L1, we use a parameter lambda to minimize $|\Theta|$ rather than Θ^2 ; L1 has a property called [sparsity](#) due to which many parameters are not just reduced but made 0, thus offering built-in feature selection
- Back Propagation.

<http://neuralnetworksanddeeplearning.com/chap2.html>