

# Week 1

# What is Machine Learning

## Introduction

### What's ML

· Early: Arthur Samuel

### Supervised / Unsupervised

· Late: Tom Mitchell

Example: playing checkers

def: E = the experience of playing many games of checkers

T = the task of playing checkers

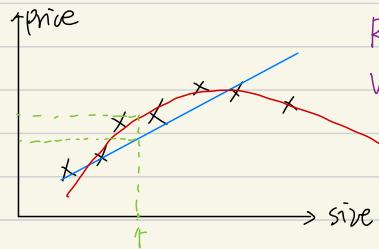
P = the probability that the program will win the next game

(task - T, experience - E, performance - P)

### Supervised Learning

监督学习

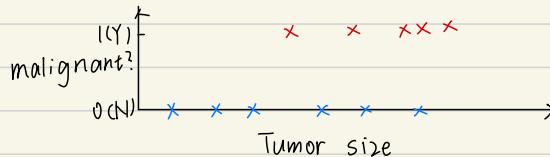
### Regression



Regression: Predict continuous valued output

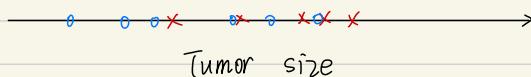
classification

- e.g. Breast cancer (malignant, benign)

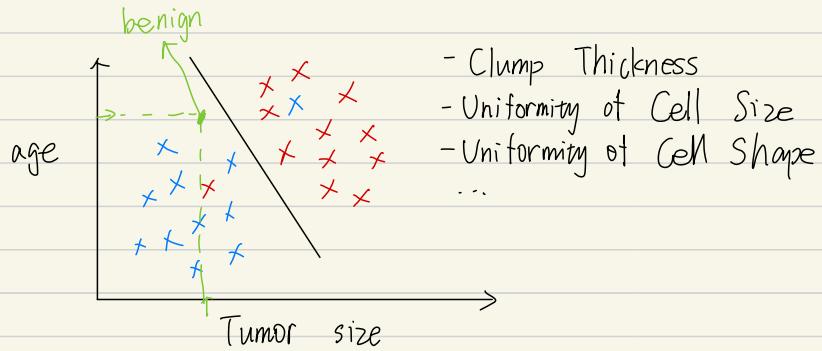


Classification: Discrete valued output (0 or 1)  
(can also have more than 2 values)

Another way to plot:

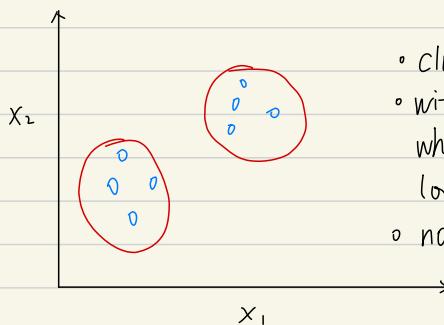


with many  
features



## Unsupervised Learning

无监督学习

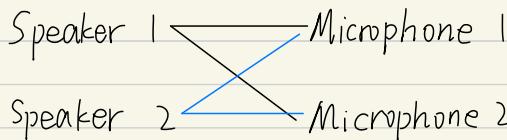


- clustering algorithm
- with little or no idea what the result should look like.
- no feedback for prediction

e.g. google news , genes

- organize computing clusters
- social network analysis
- market segmentation
- astronomical data analysis

Cocktail party problem:



Octave

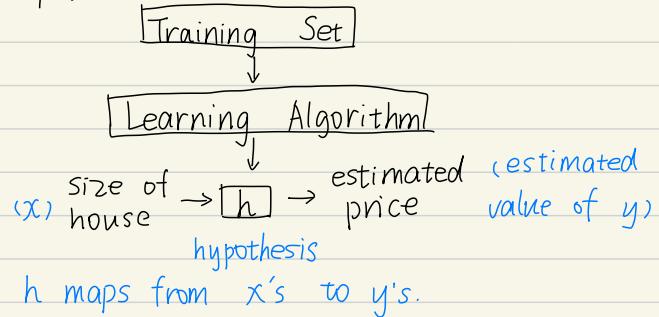
$[U, S, V] = svd(\text{repmat}(\text{sum}(x \cdot x, 1), \text{size}(x, 1), 1) \cdot x \cdot x')$

# Model and Cost Function

## Model Representation

Notation :  
 $m$  = number of training examples  
 $x$ 's = "input" variable / features  
 $y$ 's = "output" variable / "target" variable  
 $(x, y)$  = one training example  
 $(x^{(i)}, y^{(i)})$  =  $i^{\text{th}}$  training example

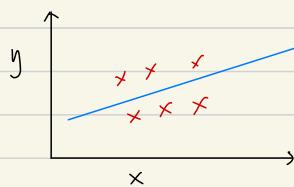
e.g. house price



how to represent  $h$ ?

$$h_0(x) = \theta_0 + \theta_1 x \quad (\text{linear})$$

shorthand :  $h(x)$



• Linear regression with one variable

• Univariate linear regression.  
one variable

## Cost Function

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

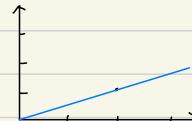
$\theta_i$ 's: parameters

- o How to choose  $\theta_i$ 's?



$$\theta_0 = 1.5$$

$$\theta_1 = 0$$



$$\theta_0 = 0$$

$$\theta_1 = 0.5$$

Idea: Choose  $\theta_0, \theta_1$  so that  $h_{\theta}(x)$  is close to  $y$  for our training examples  $(x, y)$

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

why  $\frac{1}{2}$ ?  
convenience for the gradient descent.

cost function

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

also called: squared error function  
mean squared error

- o find  $\theta_0, \theta_1$  to minimize  $J(\theta_0, \theta_1)$

## Cost Function Intuition I

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

Parameters:  $\theta_0, \theta_1$

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

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

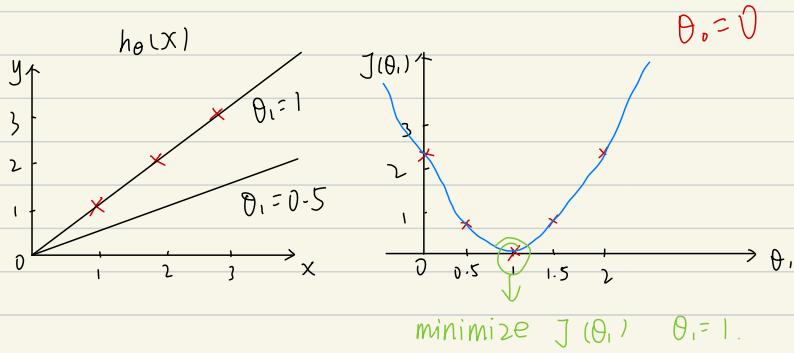
o Simplified:  $h_{\theta}(x) = \theta_1 x$  (Let  $\theta_0 = 0$ )

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

minimize  $J(\theta_1)$

$h_{\theta}(x)$ : for fixed  $\theta_1$ , this is a function of  $x$

$J(\theta_1)$ : function of the parameter  $\theta_1$



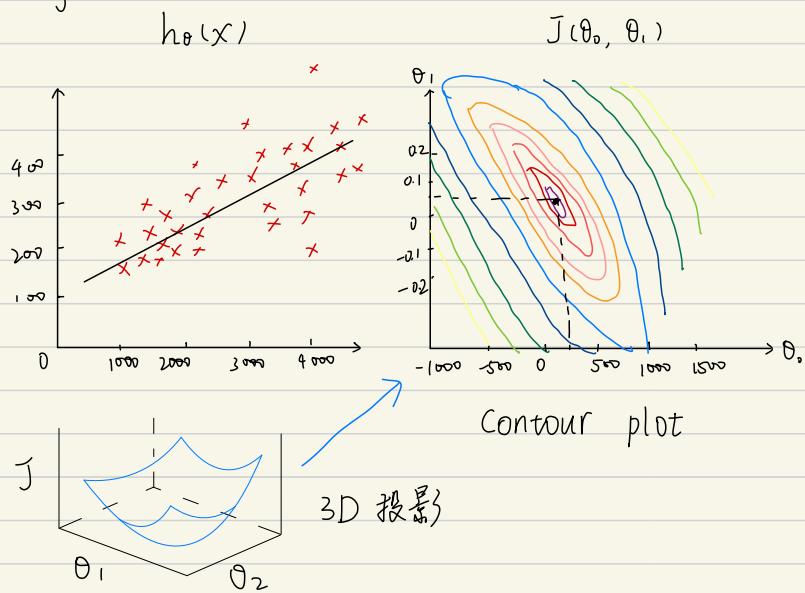
## Cost Function Intuition II

- Contour plot: a graph that contains many contour lines.

- Contour line: (of a two variable function) has a constant value at all points of the same line.

e.g.

Contour  
plot



# Parameter Learning

## Gradient Descent

Have some function  $J(\theta_0, \theta_1)$   $J(\theta_0, \theta_1, \theta_2, \dots, \theta_n)$

Want  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

$\min_{\theta_0, \dots, \theta_n} J(\theta_0, \theta_1, \dots, \theta_n)$

Outline:

- o start with some  $\theta_0, \theta_1$  (say  $\theta_0 = 0, \theta_1 = 0$ )
- o keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$  until we hopefully end up at a minimum

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  $j=1$ )  
}

learning rate (the size of each step)

Note:

- o Assignment  $a := b$

- o Truth assertion  $a = b$

$a = a + 1 \times$

o Correct: Simultaneous update

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

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

$\theta_0 := \text{temp0}$

$\theta_1 := \text{temp1}$

o Incorrect

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

$\theta_0 := \text{temp0}$

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

$\theta_1 := \text{temp1}$

Question: Why simultaneous update?

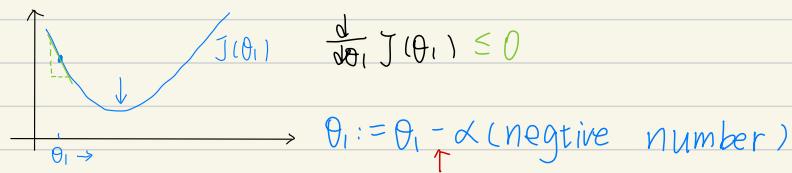
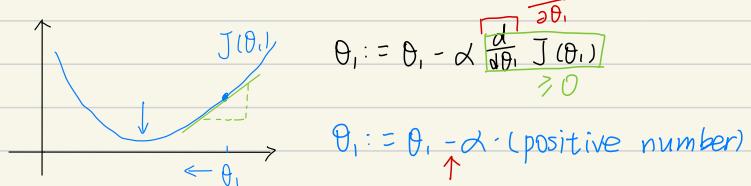
## Gradient descent Intuition

Gradient descent algorithm :

repeat until convergence }

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad \begin{matrix} \text{learning rate} \\ \downarrow \\ \text{derivative} \end{matrix} \quad \begin{matrix} (\text{simultaneously update} \\ j=0 \text{ and } j=1) \end{matrix}$$

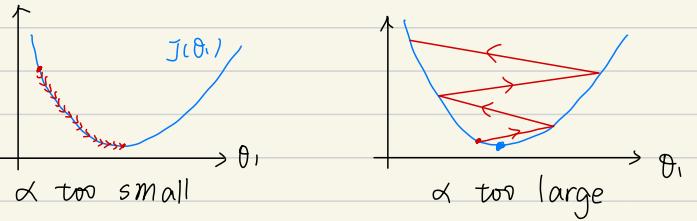
- $\min_{\theta_1} J(\theta_1) \quad \theta_1 \in \mathbb{R}$



- adjust  $\alpha \rightarrow$  converge in a reasonable time

- if  $\alpha$  too small, gradient descent can be slow

- if  $\alpha$  too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.

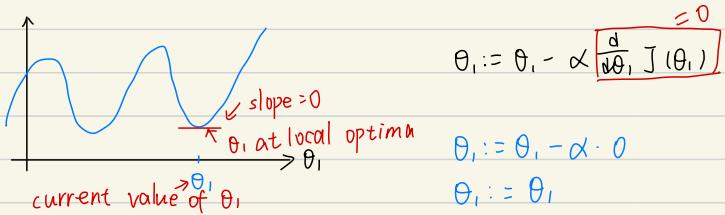


- How does gradient descent converge with a fixed step size  $\alpha$ ?

approach the bottom of the convex function  
 $\rightarrow \frac{d}{d\theta_1} J(\theta_1)$  approaches 0

$$\theta_1 := \theta_1 - \alpha \cdot 0$$

- Gradient descent can converge to a local minimum, even with the learning rate  $\alpha$  fixed.



- As we approach a local minimum gradient descent will automatically take smaller steps. So, no need to decrease  $\alpha$  over time.

## Gradient descent for linear regression

- Gradient descent algorithm:

repeat until convergence {

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

(for  $j=1$  and  $j=0$ )

}

- Linear regression Model:

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

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

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{\partial}{\partial \theta_j} \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{\partial}{\partial \theta_j} \cdot \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2 \end{aligned}$$

$$j=0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

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

- Gradient descent algorithm:

repeat until convergence {

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

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

}

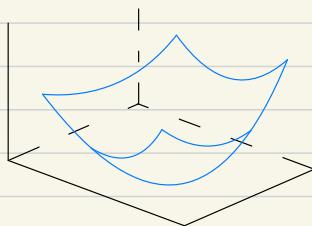
$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

} update  $\theta_0$  and  $\theta_1$

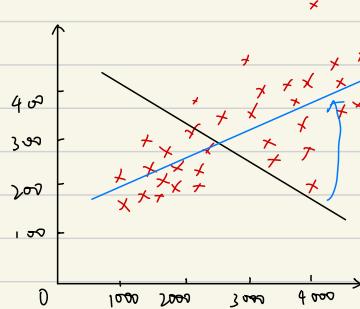
} simultaneously

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

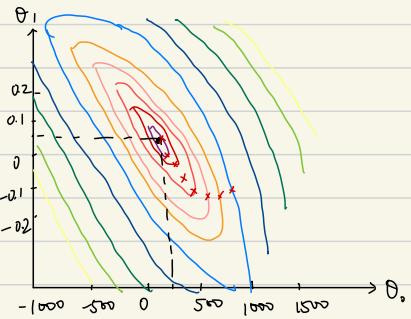
- Convex function: Bowl-shaped



$h_\theta(x)$



$J(\theta_0, \theta_1)$



"Batch" Gradient Descent

- "Batch": Each step of gradient descent uses all the training examples

# Linear Algebra Review

## Matrices and Vectors

- Matrix: Rectangular array of numbers (2-dimensional arrays)

$$\begin{array}{ccc} \rightarrow & \begin{bmatrix} 1402 & 191 \\ 1371 & 821 \\ 949 & 1437 \\ \rightarrow & \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ \uparrow & \uparrow & \uparrow \end{bmatrix} \\ 147 & 1448 \end{bmatrix} \\ \uparrow & \uparrow \end{array}$$

4x2 matrix

$$\boxed{\mathbb{R}^{4 \times 2}}$$

$$\begin{array}{ccc} \rightarrow & \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ \uparrow & \uparrow & \uparrow \end{bmatrix} \\ 2 \times 3 \end{array}$$

$$\boxed{\mathbb{R}^{2 \times 3}}$$

- Dimension of matrix: number of rows  $\times$  number of cols

- Matrix Elements (entries of matrix)

$$A = \begin{bmatrix} 1402 & 191 \\ 1371 & 821 \\ 949 & 1437 \\ 147 & 1448 \end{bmatrix}$$

$A_{ij}$  = "i, j entry" in the  $i^{\text{th}}$  row,  $j^{\text{th}}$  column.

$$A_{11} = 1402 \quad A_{12} = 191 \quad A_{32} = 1437 \quad A_{41} = 147$$

$A_{43}$  = undefined (error)

- Vector: An  $n \times 1$  matrix

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix} \leftarrow 4\text{-dimensional vector } \mathbb{R}^4$$

$y_i = i^{\text{th}}$  element  $y_1 = 460, y_2 = 232, \dots$

1-indexed vs 0-indexed

$$y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}$$

(math)

$$y = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

(machine learning)

- A, B, C, X — matrix (uppercase)  
a, b, x, y — number, scalar, vector (lowercase)

## Addition and Scalar Multiplication

- Matrix Addition

$$\begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} + \begin{bmatrix} 4 & 0.5 \\ 2 & 5 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 5 & 0.5 \\ 4 & 10 \\ 3 & 2 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} + \begin{bmatrix} 4 & 0.5 \\ 2 & 5 \end{bmatrix} = \text{error}$$

- Scalar Multiplication

$$3 \times \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 6 & 15 \\ 9 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} \times 3$$

$$\begin{bmatrix} 4 & 0 \\ 6 & 3 \end{bmatrix} / 4 = \frac{1}{4} \begin{bmatrix} 4 & 0 \\ 6 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ \frac{3}{2} & \frac{3}{4} \end{bmatrix}$$

o Combination of Operands

$$3 \times \begin{bmatrix} 1 \\ 4 \\ 2 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix} - \begin{bmatrix} 3 \\ 0 \\ 2 \end{bmatrix} / 3$$

scalar multiplication                                    scalar division

Matrix Vector  
Multiplication

Example:  $\begin{bmatrix} 1 & 3 \\ 4 & 0 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 5 \end{bmatrix} = \begin{bmatrix} 16 \\ 4 \\ 7 \end{bmatrix}$

Details:

$$\begin{bmatrix} \quad \end{bmatrix} \times \begin{bmatrix} x \\ \end{bmatrix} = \begin{bmatrix} y \\ \end{bmatrix}$$

$m \times n$  matrix       $n \times 1$  matrix       $m$ -dimensional  
( $m$  rows,                ( $n$ -dimensional      vector  
 $n$  columns)                vector)

To get  $y_i$ , multiply  $A$ 's  $i^{th}$  row with elements of vector  $x$ , and add them up

o House sizes:

2104

1416

1534

852

$$h_0(x) = -40 + 0.25x$$

matrix  $X$ :  $\downarrow 4 \times 2$       Vector:  $\downarrow 2 \times 1$

$$\begin{bmatrix} 1 & 2 \\ 1 & 1 \\ 1 & 15 \\ 1 & 8 \end{bmatrix} \times \begin{bmatrix} -40 \\ 0.25 \end{bmatrix} = \begin{bmatrix} -40 \times 1 + 0.25 \times 2 \\ -40 \times 1 + 0.25 \times 1 \\ -40 \times 1 + 0.25 \times 15 \\ -40 \times 1 + 0.25 \times 8 \end{bmatrix}$$

Prediction = DataMatrix  $\times$  Parameters

(Simplify the code, more efficient)

for  $i = 1 : 4$ ,

$\text{prediction}(i) = \dots$

## Matrix Matrix

### Multiplication

Example

$$\begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ 0 & 1 \\ 5 & 2 \end{bmatrix} = \begin{bmatrix} 11 & 10 \\ 9 & 14 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 5 \end{bmatrix} = \begin{bmatrix} 11 \\ 9 \end{bmatrix} \quad \begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 10 \\ 14 \end{bmatrix}$$

Details:  $A$      $\times$      $B$      $=$      $C$

$$\begin{bmatrix} \textcolor{red}{m} \\ \textcolor{blue}{n} \end{bmatrix} \times \begin{bmatrix} \textcolor{green}{n} \\ \textcolor{red}{l} \end{bmatrix} = \begin{bmatrix} \textcolor{green}{m} \times \textcolor{red}{l} \end{bmatrix}$$

House sizes :

2104

1416

1534

852

Have 3 competing hypotheses:

$$1. h_0(x) = -40 + 0.25x$$

$$2. h_0(x) = 200 + 0.1x$$

$$3. h_0(x) = -150 + 0.4x$$

$$\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix} \times \begin{bmatrix} -40 \\ 0.25 \\ 200 \\ 0.1 \\ -150 \\ 0.4 \end{bmatrix} = \begin{bmatrix} 486 \\ 314 \\ 344 \\ 173 \\ 342 \\ 285 \\ 464 \\ 191 \end{bmatrix}$$

## Matrix Multiplication Properties

- $3 \times 5 = 5 \times 3$  "Commutative"
- Let A and B be matrices. Then in general,  
 $A \times B \neq B \times A$ . (not commutative)
- $3 \times 5 \times 2 = (3 \times 5) \times 2 = 3 \times (5 \times 2)$  "Associative"
- $A \times B \times C = (A \times B) \times C = A \times (B \times C)$

• Identity Matrix. Denoted I (or  $I_{n \times n}$ )

I is identity.  $I \times Z = Z \times I = Z$  for any Z.

Examples of identity matrices

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & \ddots \end{bmatrix}$$

For any matrix  $A$ ,  $\underset{m \times n}{A} \cdot \underset{n \times n}{I} = \underset{n \times m}{I} \cdot \underset{m \times n}{A} = A$

## Inverse and Transpose

$$I = \text{"identity"} \quad 3 \times \underbrace{(3^{-1})}_{\frac{1}{3}} = I \quad 12 \times \underbrace{(12^{-1})}_{\frac{1}{12}} = I \quad 0 \times (0^{-1}) = \text{undefined}$$

- Matrix inverse: if  $A$  is an  $m \times m$  matrix, and if it has an inverse,  $\underset{\text{square matrix} (\# \text{ rows} = \# \text{ columns})}{A A^{-1} = A^{-1} A = I}$

e.g.  $\begin{bmatrix} 3 & 4 \\ 2 & 16 \end{bmatrix} \begin{bmatrix} 0.4 & -0.1 \\ -0.05 & 0.075 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I_{2 \times 2}$

code: `pinv(A)`

- Matrices that don't have an inverse are "singular" or "degenerate"

### Matrix Transpose

Example:  $A = \begin{bmatrix} 1 & 2 & 0 \\ 3 & 5 & 9 \end{bmatrix} \quad A^T = \begin{bmatrix} 1 & 3 \\ 2 & 5 \\ 0 & 9 \end{bmatrix}$

Let  $A$  be an  $m \times n$  matrix, and let  $B = A^T$ .  
Then  $B$  is an  $n \times m$  matrix, and  $B_{ij} = A_{ji}$

# Week 2 Multiple Features

## Multivariate Linear Regression

### Multiple Features

	Size (feet <sup>2</sup> )	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
	$x_1$	$x_2$	$x_3$	$x_4$	$y$
	2104	5	1	45	460
m	1416	3	2	40	232
	1534	3	2	40	315
	852	2	1	36	178
	...	...	...	...	...

Notation:  $m$  = number of training examples  $n=4$

$n$  = number of features

$x^{(i)}$  = input (features) of  $i^{\text{th}}$  training example

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

$$x^{(i)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix} \in \mathbb{R}^4 \quad x_3^{(i)} = 2$$

- Hypothesis:  $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$

E.g.  $h_{\theta}(x) = 80 + 0.1 x_1 + 0.01 x_2 + 3 x_3 - 2 x_4$

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

For convenience of notation, define  $x_0 = 1$ . ( $x_0^{(i)} = 1$ )

$$X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta^T = [\theta_0, \theta_1, \dots, \theta_n]$$

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n = \theta^T x$$

## Gradient Descent for Multiple Variables

- Hypothesis:  $h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$
- Parameters:  $\theta$  (n+1)-dimensional vector
- Cost function:  $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$
- Gradient descent:
  - Repeat until convergence!
  - $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$
  - } (simultaneously update for every  $j=0, \dots, n$ )
- $n=1$ : Repeat}
  - $\theta_0 := \theta_0 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})}_{\frac{\partial}{\partial \theta_0} J(\theta)}$
  - $\theta_1 := \theta_1 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}}_{\text{(simultaneously update } \theta_0, \theta_1)}$
  - }
- New algorithm ( $n \geq 1$ ):

Repeat {

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

(simultaneously update  $\theta_j$  for  $j = 0, \dots, n$ )

}

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

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

⋮

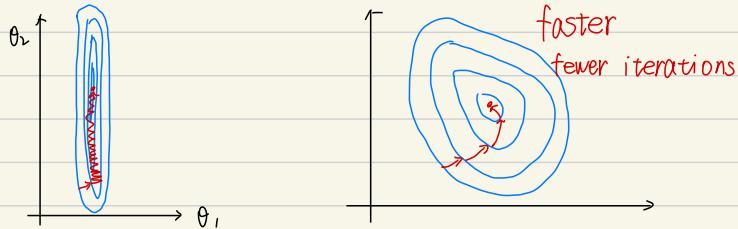
## Gradient Descent in Practice 1 - Feature Scaling

### Feature Scaling

Idea: Make sure features are on a similar scale.

E.g.  $x_1 = \text{size (0-2000 feet)}$

$x_2 = \text{number of bedrooms (1-5)}$



$$x_1 = \frac{\text{size (feet)}}{2000}$$

$$0 \leq x_1 \leq 1$$

$$x_2 = \frac{\text{number of bedrooms}}{5}$$

$$0 \leq x_2 \leq 1$$

Feature Scaling Get every feature into approximately a  $-1 \leq x_i \leq 1$  range

$$\begin{array}{ll}
 x_0 = 1 & \\
 0 \leq x_1 \leq 3 & \checkmark \\
 -2 \leq x_2 \leq 0.5 & \checkmark \\
 -100 \leq x_3 \leq 100 & \cancel{x} \quad \rightarrow \text{to } 3 \quad \checkmark \\
 -0.0001 \leq x_4 \leq 0.0001 & \cancel{x} \quad \rightarrow \frac{1}{3} \text{ to } \frac{1}{3} \quad \checkmark
 \end{array}$$

o Mean normalization: Replace  $x_i$  with  $x_i - \bar{x}_i$  to make features have approximately zero mean (Do not apply to  $x_0 = 1$ )

o Apply feature scaling and mean normalization:  
 E.g.  $x_1 = \frac{\text{size} - 1000}{2000}$        $x_2 = \frac{\# \text{bedrooms} - 2}{5}$   
 $-0.5 \leq x_1 \leq 0.5$        $-0.5 \leq x_2 \leq 0.5$

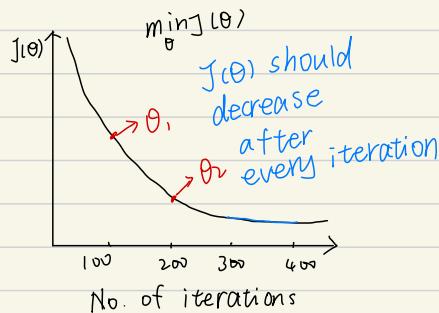
o  $x_1 := \frac{x_1 - \bar{x}_1}{s_1}$  (avg value of  $x_1$  in training set  
 range (max - min)  
 cor standard deviation)

## Gradient descent in Practice II - Learning Rate

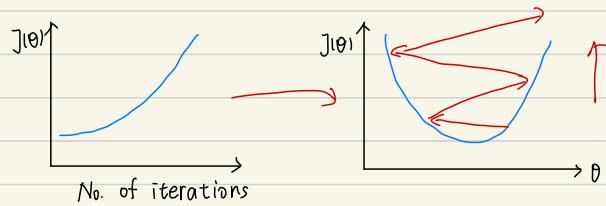
Gradient descent:  $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$

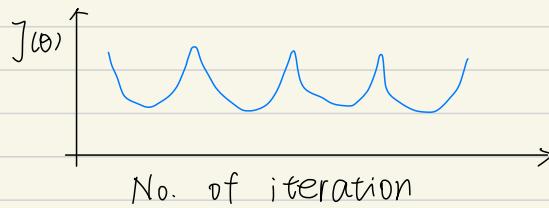
- o "Debugging": How to make sure gradient descent is working correctly
- o How to choose learning rate  $\alpha$ .

- Making sure gradient descent is working correctly.



- Example automatic convergence test: Declare convergence if  $J(\theta)$  decreases by less than  $\epsilon$  in one iteration. ( $\epsilon$  might be  $10^{-3}$ , but is hard to decide)
- Gradient descent not working. Use smaller  $\alpha$ :





- For sufficiently small  $\alpha$ ,  $J(\theta)$  should decrease on every iteration.
- But if  $\alpha$  is too small, gradient descent can be slow to converge.
- Summary:
  - If  $\alpha$  is too small: slow converge
  - If  $\alpha$  is too large:  $J(\theta)$  may not decrease on every iteration; may not decrease. (slow converge also possible)
- To choose  $\alpha$ , try
 
$$\dots, 0.001, \underbrace{0.003}_{\times 3}, \underbrace{0.01}_{\approx 3}, \underbrace{0.03}_{\times 3}, 0.1, \underbrace{0.3}_{\approx 3}, 1, \dots$$

## Features and Polynomial Regression

- House prices prediction

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \underbrace{\text{frontage}}_{x_1} + \theta_2 \times \underbrace{\text{depth}}_{x_2}$$

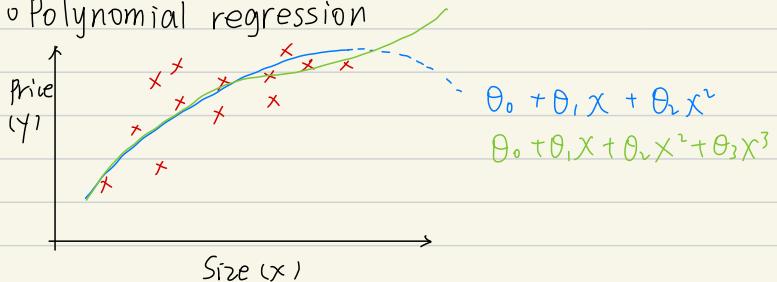


- We can combine multiple features into one. For example, we can combine  $x_1$  and  $x_2$  into a new feature  $x_3$  by taking  $x_1 \cdot x_2$

$$\text{Area } x = \text{frontage} \times \text{depth}$$

$$h_\theta(x) = \theta_0 + \theta_1 x_1 \leftarrow \text{land area}$$

- Polynomial regression



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

$$= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3$$

$$x_1 = (\text{size})$$

$$\text{size : } 1 - 1000$$

$$x_2 = (\text{size})^2$$

$$\text{size}^2 : 1 - 1000,000$$

$$x_3 = (\text{size})^3$$

$$\text{size}^3 : 1 - 10^9$$

- Choice of features

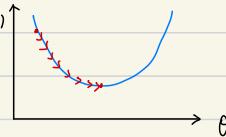
$$h_\theta(x) = \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2$$

$$h_\theta(x) = \theta_0 + \theta_1 (\text{size}) + \theta_2 \sqrt{\text{size}}$$

## Computing Parameters Analytically

### Normal Equation

- Gradient Descent

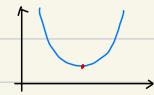


- Normal equation: Method to solve for  $\theta$  analytically.

- Intuition: If 1D ( $\theta \in \mathbb{R}$ )

$$J(\theta) = a\theta^2 + b\theta + c$$

$\frac{\partial}{\partial \theta} J(\theta) = \dots \stackrel{\text{set}}{=} 0$ , solve for  $\theta$ .



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

$\frac{\partial}{\partial \theta_j} J(\theta) = \dots \stackrel{\text{set}}{=} 0$  (for every  $j$ )      solve for  $\theta_0, \theta_1, \dots, \theta_n$

- Example:  $m=4$

	Size (feet <sup>2</sup> )	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
$x_0$	$x_1$	$x_2$	$x_3$	$x_4$	$y$
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	40	315
1	852	2	1	36	178

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$$

$m \times (n+1)$

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

$m$ -dimensional vector

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

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

$n$  features

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1}$$

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

$m \times (n+1)$

$$\text{E.g. If } x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \end{bmatrix} \quad X = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ \vdots & \vdots \\ 1 & x_1^{(m)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$m \times 2$

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

$(X^T X)^{-1}$  is inverse of matrix  $X^T X$

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

$0 \leq x_1 \leq 1$

~~Feature Scaling~~

$0 \leq x_2 \leq 1000$  ✓

$0 \leq x_3 \leq 10^{-5}$  ✓

m training examples, n features

Gradient Descent

Need to choose  $\alpha$

Need many iterations

Works well even when n  
is large

$O(kn^2)$

$n = 10^6$

$\sqrt{n} > 10^4$

Normal Equation

No need to choose  $\alpha$

Don't need to iterate

Need to compute  $(X^T X)^{-1}$   
 $(n \times n)^{-1}$   $O(n^3)$

Slow if n is very large

$n = 100, 1000, 10000$  ]

Normal Equation

Noninvertibility

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

What if  $X^T X$  is non-invertible? (singular/  
degenerate)

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

{  $\text{pinv}$  even when  
 $\text{inv}$  non-invertible }

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

- Redundant features (linearly dependent)

E.g.  $x_1 = \text{size in feet}$   $1\text{m} = 3.28 \text{ feet}$

$$x_2 = \text{size in m} \quad x_1 = (3.28)^{-1} x_2$$

Delete a feature that is linearly dependent with another

- Too many features (e.g.  $m \leq n$ )  $m=10, n=100$

$$\theta \in \mathbb{R}^{101}$$

Delete some features, or use regularization.

# Week 3 Logistic Regression

## Classification and Representation

### Classification

- Classification:

- Email: Spam / Not Spam?
- Online Transactions: Fraudulent (Yes / No)?
- Tumor: Malignant / Benign?

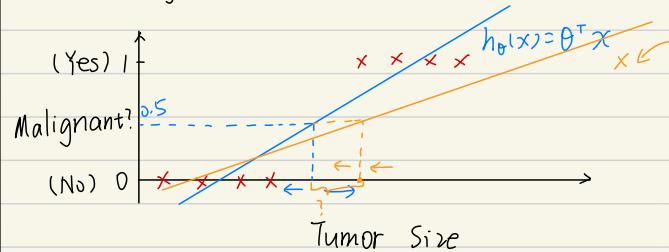
$$y \in \{0, 1\}$$

$$y \in \{0, 1, 2, 3\} \dots$$

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

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

- Linear Regression (not applicable)



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

If  $h_\theta(x) > 0.5$ , predict "y=1"

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

- Classification:  $y=0$  or  $1$

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

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

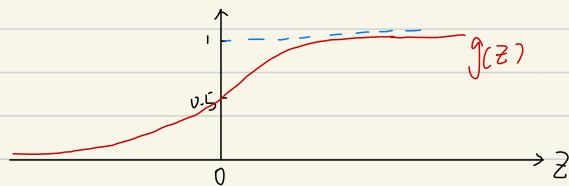
## Hypothesis Representation

- Logistic Regression Model: Want  $0 \leq h_{\theta}(x) \leq 1$

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

$g(z) = \frac{1}{1 + e^{-z}}$  (Sigmoid function / Logistic function)

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



- Interpretation of Hypothesis Output

$h_{\theta}(x)$  = estimated probability that  $y=1$  on input  $x$

Example: If  $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$

$h_{\theta}(x) = 0.7$ : Tell patient that 70% chance of tumor being malignant.

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

"probability that  $y=1$ , given  $x$ , parameterized by  $\theta$ "

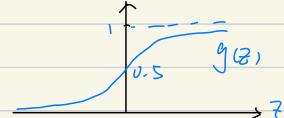
$$y=0 \text{ or } 1 : P(y=0 | x; \theta) + P(y=1 | x; \theta) = 1$$

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

## Decision Boundary

- Logistic regression:  $h_\theta(x) = g(\theta^T x) = P(y=1|x; \theta)$

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



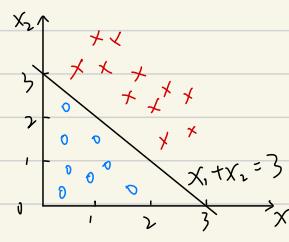
- Suppose predict "y=1" if  $h_\theta(x) \geq 0.5$   
predict "y=0" if  $h_\theta(x) < 0.5$

$g(z) \geq 0.5$  when  $z \geq 0$

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

- Decision Boundary: the line that separates the area where  $y=0$  and where  $y=1$ . Created by the hypothesis function.

- E.g. linear



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

$$\theta_0 = -3 \quad \theta_1 = 1 \quad \theta_2 = 1 \quad \theta^T x$$

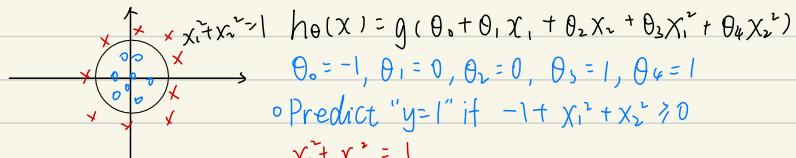
- Predict "y=1" if  $-3 + x_1 + x_2 \geq 0$

$$x_1 + x_2 \geq 3$$

$$\text{Boundary: } x_1 + x_2 = 3$$

$$y=0: x_1 + x_2 < 3$$

- Non-linear decision boundaries



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

$$\theta_0 = -1, \theta_1 = 0, \theta_2 = 0, \theta_3 = 1, \theta_4 = 1$$

- Predict "y=1" if  $-1 + x_1^2 + x_2^2 \geq 0$

$$x_1^2 + x_2^2 = 1$$

## Logistic Regression Model

### Cost Function

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

$m$  examples  $x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$   $x_0 = 1$   $y \in \{0, 1\}$

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

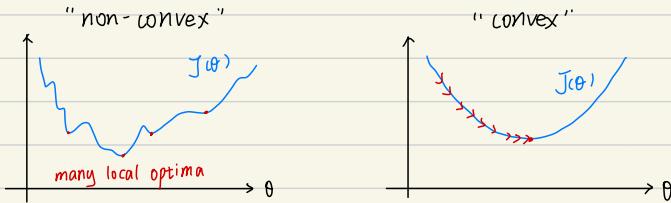
- How to choose parameters  $\theta$ ?

- Cost Function

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

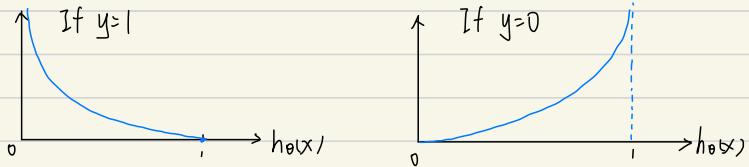
Logistic regression:  $J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_{\theta}(x^{(i)}) - y^{(i)})$

- If choose  $\text{cost}(h_{\theta}(x), y) = \frac{1}{2} (h_{\theta}(x) - y)^2$ : non-convex



- Logistic regression cost function

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



- If  $y=1$ :  $h_{\theta}(x)=1 \Rightarrow \text{Cost}=0$   
 $h_{\theta}(x) \rightarrow 0 \Rightarrow \text{Cost} \rightarrow \infty$

Captures intuition that if  $h_{\theta}(x)=0$  (predict  $P(y=1|x;\theta)=0$ ), but  $y=1$ , we will penalize learning algorithm by a very large cost

## Simplified Cost Function and Gradient Descent

Logistic regression cost function:

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

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

Note:  $y=0$  or  $1$  always

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

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

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

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

- To fit parameters  $\theta$ :  $\min_{\theta} J(\theta)$

- To make a prediction given new  $x$ : Output  $h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$

- Gradient Descent: Want  $\min_{\theta} J(\theta)$ :

Repeat:

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

} (simultaneously update all  $\theta_j$ )

- Repeat {
 
$$\theta_j := \theta_j - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$
}
 
$$h_\theta = g(X\theta)$$
- Algorithm looks identical to linear regression
- Vectorized implementation:  $\theta := \theta - \frac{\alpha}{m} X^T (g(X\theta) - \vec{y})$

## Advanced Optimization

- Optimization Algorithm  
Cost function  $J(\theta)$ . Want  $\min J(\theta)$ .  
Given  $\theta$ , we have code that compute:  $J(\theta)$ ,  $\frac{\partial}{\partial \theta_j} J(\theta)$  (for  $j=0, 1, \dots, n$ )
- Gradient Descent: Repeat {  $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$  }
- Optimization algorithm: Given  $\theta$ , we have code that can compute:  $J(\theta)$  and  $\frac{\partial}{\partial \theta_j} J(\theta)$  (for  $j=0, 1, \dots, n$ )
- Optimization algorithm: Gradient descent
  - Conjugate gradient
  - BFGS
  - L-BFGS

Advantages: No need to manually pick  $\alpha$   
Often faster than gradient descent  
Disadvantages: More complex

Example:  $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$

$$J(\theta) = (\theta_0 - 2)^2 + (\theta_1 - 5)^2 \quad \min_{\theta} J(\theta) : \theta_0 = \theta_1 = 5$$

$$\frac{\partial}{\partial \theta_0} J(\theta) = 2(\theta_0 - 2)$$

$$\frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5)$$

Code: function [jVal, gradient] = costFunction(theta)

$$jVal = (\theta(1) - 2)^2 + (\theta(2) - 5)^2;$$

gradient = zeros(2,1);

$$\text{gradient}(1) = 2 * (\theta(1) - 2)$$

$$\text{gradient}(2) = 2 * (\theta(2) - 5)$$

users will provide  
a gradient the maximum number  
of iteration

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

initialTheta = zeros(2,1); initial guess for theta

[optTheta, functionVal, exitFlag] ...

= fminunc(@costFunction, initialTheta, options);  
advanced a pointer to function  $\theta \in \mathbb{R}^d, d \geq 2$

theta =  $\begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$

$\theta_0$   $\theta_1$

$\theta_2$

$\theta_{n+1}$

function [jVal, gradient] = costFunction(theta)

jVal = [code to compute  $\frac{\partial}{\partial \theta_0} J(\theta)$ ];

gradient(1) = [code to compute  $\frac{\partial}{\partial \theta_1} J(\theta)$ ];

gradient(2) = [code to compute  $\frac{\partial}{\partial \theta_2} J(\theta)$ ];

## Multiclass Classification

gradient( $\theta$ ) = [code to compute  $\frac{\partial}{\partial \theta} J(\theta)$ ];

### Multiclass Classification:

One - vs - all

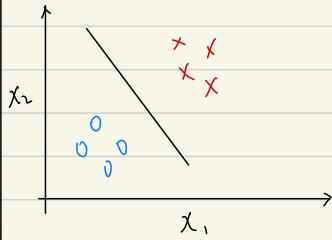
- Multiclass classification

Email foldering/tagging: Work, Friends, Family, Hobby  
 $y=1$        $y=2$        $y=3$        $y=4$

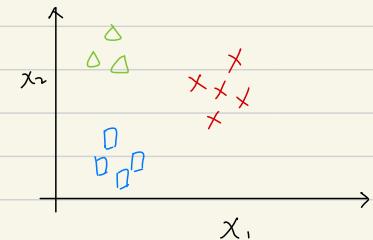
Medical diagrams: Not ill, Cold, Flu  
 $y=1$        $y=2$        $y=3$

Weather: Sunny, Cloudy, Rain, Snow  
 $y=1$        $y=2$        $y=3$        $y=4$

### Binary classification



### Multi-class classification

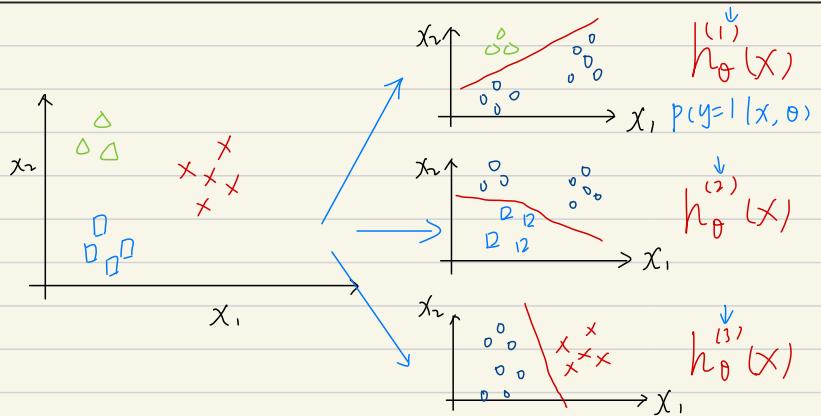


- One - vs - all (one - vs - rest)

Class 1: ▲

Class 2: □

Class 3: ✕



- $h_\theta^{(i)}(x) = P(y=i|x; \theta) \quad (i=1, 2, 3)$

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

- $y \in \{0, 1, \dots, n\}$

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

⋮

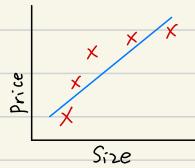
$$h_\theta^{(n)}(x) = P(y=n|x; \theta)$$

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

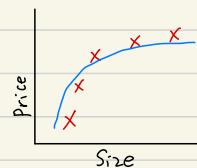
## Solving the Problem of Overfitting

### The Problem of Overfitting

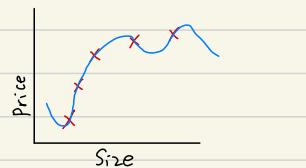
- Example: Linear regression (housing prices)



"underfitting"  
"high bias"



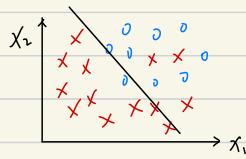
"just right"



"overfitting"  
"high variance"

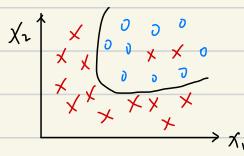
- Overfitting: If we have too many features, the learned hypothesis may fit the training set very well ( $J(\theta) = \frac{1}{m} \sum_{i=1}^m (\text{h}_\theta(x^{(i)}) - y^{(i)})^2 \approx 0$ ), but fail to generalize to new examples (predict prices on new examples)

- Example: Logistic regression  $g$  = sigmoid function



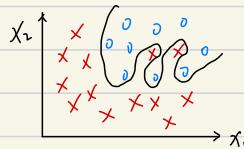
$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

"underfitting"  
"high bias"



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$

"just right"



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3 x_2 + \theta_4 x_1^2 x_2^2 + \dots)$$

"overfitting" "high variance"

### Addressing overfitting:

$x_1$  = size of house

$x_2$  = no. of bedrooms

$x_3$  = no. of floors

$x_4$  = age of house

:

$x_{100}$



### Options:

1. Reduce number of features.

- Manually select which features to keep

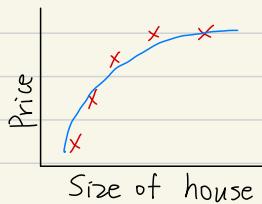
- Model selection algorithm (later in course)

## 2. Regularization

- Keep all the features, but reduce magnitude/values of parameters  $\theta_j$ .
- Works well when we have a lot of features, each of which contributes a bit to predicting  $y$ .

### Cost Function

#### Intuition



$$\theta_0 + \theta_1 x + \theta_2 x^2$$



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

- o Suppose we penalize and make  $\theta_3, \theta_4$  really small:

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + 1000 \theta_3^2 + 1000 \theta_4^2$$

$\downarrow \theta_3 \approx 0 \quad \theta_4 \approx 0$

- o Regularization: Small values for parameters  $\theta_0, \theta_1, \dots, \theta_n$ 
  - "Simpler" hypothesis
  - Less prone to overfitting

- o Housing :

- Features:  $x_1, x_2, \dots, x_{100}$

- Parameters:  $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

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

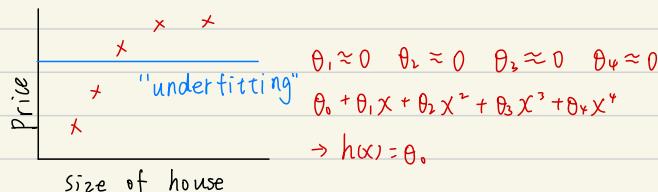
$\theta_0, \theta_1, \theta_2, \dots, \theta_n$

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

regularization parameter

o In regularized linear regression, we choose  $\theta$  to minimize  $J(\theta)$ .

What if  $\lambda$  is set to an extremely large value (perhaps for too large for our problem, say  $\lambda = 10^9$ )?



## Regularized Linear Regression

o Regularized linear regression

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

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

o Gradient descent

Repeat {

$$\frac{\partial}{\partial \theta_0} J(\theta)$$

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

$$\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad (j = 0, 1, 2, \dots, n)$$

}

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

regularized

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

Same  
 $1 - \alpha \frac{\lambda}{m} < 1$ : reducing  $\theta_j$  on every update

◦ Normal equation

$$X = \begin{bmatrix} (x^{(1)})^\top \\ \vdots \\ (x^{(m)})^\top \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \in \mathbb{R}^m \quad \rightarrow \min_{\theta} J(\theta)$$

$$\frac{\partial}{\partial \theta_j} J(\theta) \stackrel{\text{set}}{=} 0 : \theta = (X^\top X + \lambda I)^{-1} X^\top y$$

$$I = \begin{bmatrix} 0 & & & \\ 1 & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}_{(m+1) \times (m+1)}$$

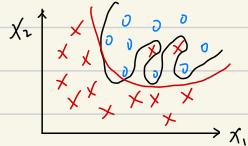
◦ Non-invertibility (optional / advanced)

Suppose  $m < n$ ,  $\theta = (X^\top X)^{-1} X^\top y$   
non-invertible / singular      pinv, inv

If  $\lambda > 0$ ,  $\theta = (X^\top X + \lambda I)^{-1} X^\top y$   
invertible

## Regularized Logistic Regression

◦ Regularized logistic regression



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

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

◦ Gradient descent

Repeat {

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

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

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

}

◦ Advanced optimization

$fminunc(@costFunction, \dots)$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

function [jVal, gradient] = costFunction(theta)

jVal = [code to compute  $J(\theta)$ ];

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

gradient(1) = [code to compute  $\frac{\partial}{\partial \theta_0} J(\theta)$ ];

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

gradient(2) = [code to compute  $\frac{\partial}{\partial \theta_1} J(\theta)$ ];

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

gradient(3) = [code to compute  $\frac{\partial}{\partial \theta_2} J(\theta)$ ];

$$\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)} + \frac{\lambda}{m} \theta_2$$

:

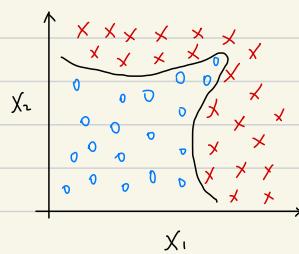
gradient(n+1) = [code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$ ];

# Week 4 Neural Networks: Representation

## Motivations

### Non-linear Hypothesis

#### Non-linear Classification



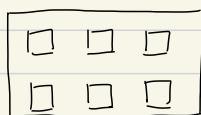
$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^3 x_2 + \theta_6 x_1 x_2^2 + \dots)$$

$n=100:$

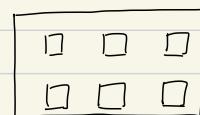
$$\text{If: } x_1^2, x_1 x_2, x_1 x_3, \dots, x_1 x_{100}, x_2^2, x_2 x_3, \dots \approx 5000 \text{ features} \approx \frac{1}{2} n^2 = O(n^2)$$

$$\text{If: } x_1 x_2 x_3, x_1^2 x_2, \dots, x_{10} x_{11} x_{12}, \dots \approx 170,000 \text{ features} = O(n^3)$$

#### Computer Vision: Car detection



Cars

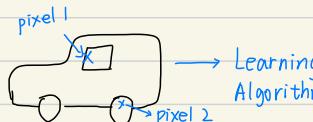


Not a car

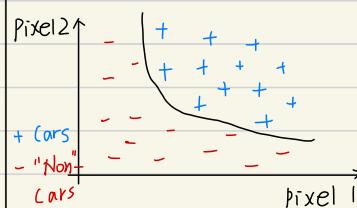
Testing:



What is this?



50x50 pixel images  $\rightarrow 2500$  pixels  
 $n=2500$  ( $7500$  if RGB)



$$x = \begin{bmatrix} \text{pixel 1 intensity} \\ \vdots \\ \text{pixel 2500 intensity} \end{bmatrix}$$

Quadratic features ( $x_i \times x_j$ ):  
 $\approx 3$  million features  
expensive computation!

## Neurons and the Brain

### Neural Networks

- Origins: Algorithm that try to mimic the brain
- Was very widely used in 80s and early 90s; popularity diminished in late 90s.
- Recent resurgence: State-of-the-art technique for many application

### The "one learning algorithm" hypothesis

- Auditory cortex learns to see (unable to listen)
- Somatosensory cortex learns to see (unable to touch)

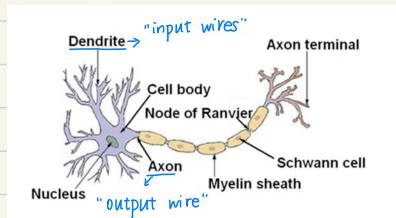
### Sensor representations in the brain

- Seeing without your tongue
- Human echolocation (sonar)
- Haptic belt: Direction sense
- Implanting a 3rd eye

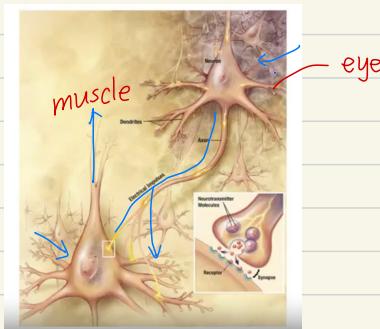
## Neural Networks

### Model Representation I

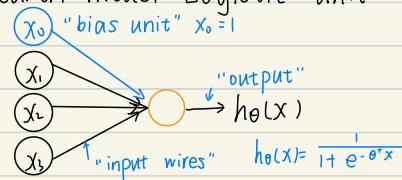
#### Neuron in the brain



#### Neurons in the brain



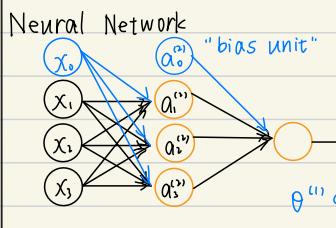
#### Neuron model: Logistic unit



$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

"weights"  
(parameters)

Sigmoid (logistic) activation function,  $g(z) = \frac{1}{1 + e^{-z}}$



Layer 1      Layer 2      Layer 3  
 ↑                  ↓                  ↓  
 input layer    hidden layer    output layer

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

$\theta^{(j)} \in \mathbb{R}^{3 \times 4}$  = matrix of weights controlling function

mapping from layer  $j$  to layer  $j+1$

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

$$a_2^{(2)} = g(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3)$$

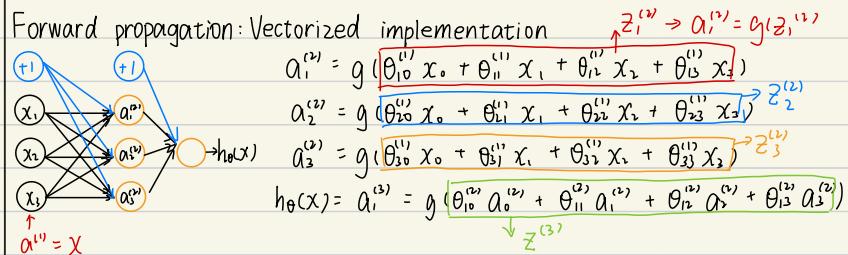
$$a_3^{(2)} = g(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3)$$

$$h_\theta(x) = a^{(3)} = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)})$$

- If network has  $s_j$  units in layer  $j$ ,  $s_{j+1}$  units in layer  $j+1$ , then  $\theta^{(j)}$  will be of dimension  $s_{j+1} \times (s_j + 1)$  bias nodes

## Model Representation II

Forward propagation: Vectorized implementation



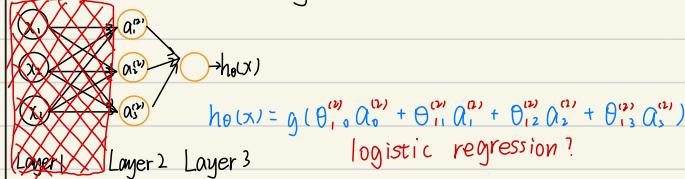
$$a^{(1)} = X$$

$$X = \begin{bmatrix} X_0 \\ X_1 \\ X_2 \\ X_3 \end{bmatrix} \quad Z^{(1)} = \begin{bmatrix} Z_1^{(1)} \\ Z_2^{(1)} \\ Z_3^{(1)} \end{bmatrix}$$

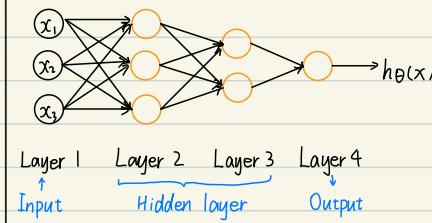
$\alpha_{1,1}^{(1)}$   
 $\alpha_{1,2}^{(1)}$   
 $\alpha_{1,3}^{(1)}$

$Z^{(1)} = \Theta^{(1)} X = \Theta^{(1)} a^{(1)}$ ,  $a^{(1)} = g(Z^{(1)}) \in \mathbb{R}^3$   
 $Z^{(2)} = \Theta^{(2)} a^{(1)}$ ,  $h_\theta(x) = a^{(2)} = g(Z^{(2)})$

Neural Network learning its own features



Other network architectures

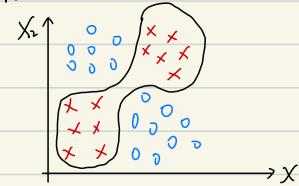
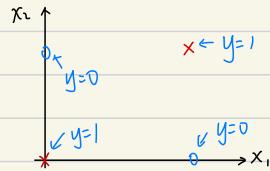


## Applications

### Examples and Intuitions I

Non-linear classification example: XOR / XNOR

$x_1$  and  $x_2$  are binary (0 or 1)



$$y = x_1 \text{ XOR } x_2$$

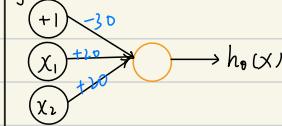
$$x_1 \text{ XNOR } x_2$$

$$\text{Not } (x_1 \text{ XOR } x_2)$$

Simple example: AND

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

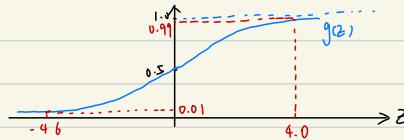
$$y = x_1 \text{ AND } x_2$$



$$h_\theta(x) = g(-3x_1 + 2x_2)$$

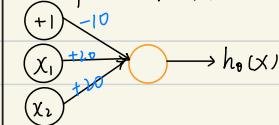
$$\theta_{01}^{(i)}, \theta_{11}^{(i)}, \theta_{12}^{(i)}$$

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



$x_1$	$x_2$	$h_\theta(x)$
0	0	$g(-3 \cdot 0 + 2 \cdot 0) \approx 0$
0	1	$g(-3 \cdot 0 + 2 \cdot 1) \approx 0$
1	0	$g(-3 \cdot 1 + 2 \cdot 0) \approx 0$
1	1	$g(-3 \cdot 1 + 2 \cdot 1) \approx 1$

Example: OR function

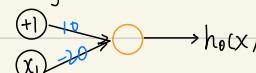


$$h_0(x) = g(-10 + 20x_1 + 20x_2)$$

$x_1$	$x_2$	$h_0(x)$
0	0	$g(-10) \approx 0$
0	1	$g(10) \approx 1$
1	0	$g(10) \approx 1$
1	1	$g(30) \approx 1$

## Examples and Intuitions II

Negation: NOT  $x_1$ ,

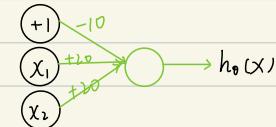
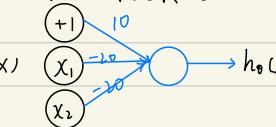
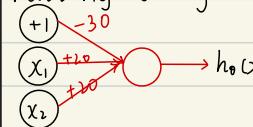


$$h_0(x) = g(10 - 20x_1)$$

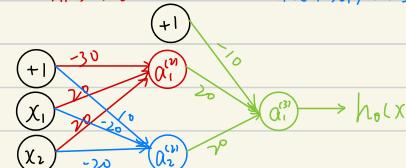
$x_1$	$h_0(x)$
0	$g(10) \approx 1$
1	$g(-10) \approx 0$

$$(\text{NOT } x_1) \text{ AND } (\text{NOT } x_2) = 1 \text{ iff } x_1 = x_2 = 0$$

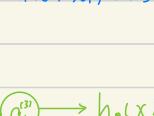
Putting it together:  $x_1 \text{ XNOR } x_2$



$x_1 \text{ AND } x_2$



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



$x_1 \text{ OR } x_2$

$x_1$	$x_2$	$a_1^{(v)}$	$a_2^{(v)}$	$h_0(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

$$\text{AND: } \theta^{(1)} = [-30 \ 20 \ 20] \quad > \text{AND and NOR: } \theta^{(1)} = \begin{bmatrix} -30 & 20 & 20 \\ 10 & -20 & -20 \end{bmatrix}$$

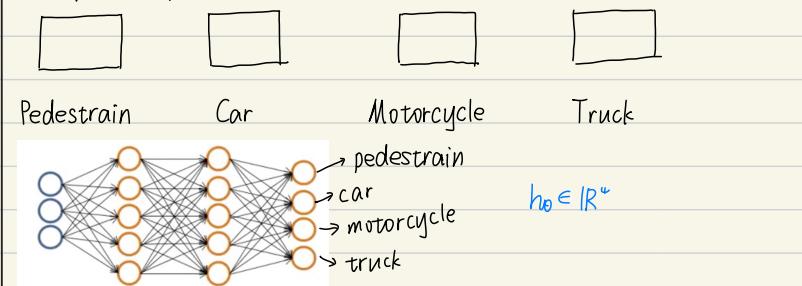
$$\text{NOR: } \theta^{(1)} = [10 \ -20 \ -20] \quad \rightarrow \text{OR: } \theta^{(2)} = [-10 \ 20 \ 20]$$

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{bmatrix} a_0^{(1)} \\ a_1^{(1)} \\ a_2^{(1)} \end{bmatrix} \rightarrow [a^{(1)}] \rightarrow h_{\theta}(x) \quad \begin{cases} a^{(2)} = g(\theta^{(1)} \cdot x) \\ a^{(3)} = g(\theta^{(2)} \cdot a^{(1)}) \\ h_{\theta}(x) = a^{(3)} \end{cases}$$

◦ Video: Handwritten digit classification

## Multiclass Classification

Multiple output units: One-vs-all



$$\text{Want } h_{\theta}(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \quad h_{\theta}(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \quad h_{\theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \quad \text{etc.}$$

when pedestrian      when car      when motorcycle

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

$$y^{(i)} \text{ one of } \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \quad \text{previously } y \in \{1, 2, 3, 4\}$$

pedestrian      car      motorcycle      truck

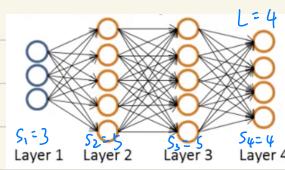
$$h_{\theta}(x^{(i)}) \approx y^{(i)} \quad \mathbb{R}^4$$

# Week 5 Neural Networks: Learning

## Cost Function and Backpropagation

### Cost Function

#### Neural Network (Classification)



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

$L$  = total no. of layers in network

$S_l$  = no. of units (not counting bias unit) in layer  $l$ .

- Binary classification

$$y = 0 \text{ or } 1, \text{ 1 output unit, } h_\theta(x) \in \mathbb{R}, S_L = 1, k = 1.$$

- Multi-class classification ( $k$  classes)

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

E.g.  $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

### Cost function:

- Logistic regression:

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

- Neural Network:

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log (h_\theta(x^{(i)}))_k + (1-y_k^{(i)}) \log (1-(h_\theta(x^{(i)}))_k) \right] - \frac{\lambda}{2m} \sum_{i=1}^m \sum_{j=1}^n \sum_{k=1}^K (\theta_{jk}^{(i)})^2 \rightarrow S_{kl} \times S_l \text{ matrix}$$

the double sum simply adds up the logistic regression costs calculated for each cell in the output layer.

the triple sum simply adds up the square of all the individual  $\theta$ s in the entire network

## Backpropagation Algorithm

Gradient computation

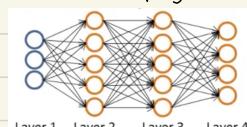
$$J(\theta) = -\frac{1}{m} \left[ \sum_{k=1}^K y_k^{(i)} \log h_\theta(x^{(i)})_k + (1 - y_k^{(i)}) \log (1 - h_\theta(x^{(i)})_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{j=1}^{S_l} \sum_{j'=1}^{S_{l+1}} (\theta_{j,j'}^{(l)})^2$$

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

Need code to compute:  $J(\theta)$  &  $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) \quad \theta_{ij}^{(l)} \in \mathbb{R}$

Given one training example  $(x, y)$ :

Forward propagation:  $a^{(0)} = x$



$$a^{(0)} \quad a^{(1)} \quad a^{(2)} \quad a^{(3)} \quad a^{(4)}$$

$$z^{(0)} = \theta^{(0)} a^{(0)}$$

$$a^{(1)} = g(z^{(0)}) \quad (\text{add } a_0^{(1)})$$

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

$$a^{(2)} = g(z^{(1)}) \quad (\text{add } a_0^{(2)})$$

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

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

Backpropagation algorithm

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

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

$$\delta_j^{(4)} = [a_j^{(4)} - y_j] \quad \text{vec: } \delta^{(4)} = a^{(4)} - y$$

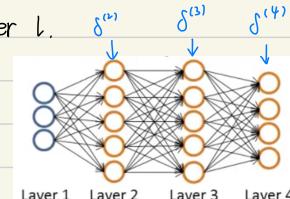
formula?

$$\delta^{(3)} = (\theta^{(3)})^\top \delta^{(4)} * \boxed{g'(z^{(3)})} \rightarrow a^{(3)} * (1 - a^{(3)})$$

$$\delta^{(2)} = (\theta^{(2)})^\top \delta^{(3)} * \boxed{g'(z^{(2)})} \rightarrow a^{(2)} * (1 - a^{(2)})$$

NO  $\delta^{(1)}$

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



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

$$g'(z) = g(z) \cdot (1 - g(z))$$

## Backpropagation algorithm

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

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

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

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

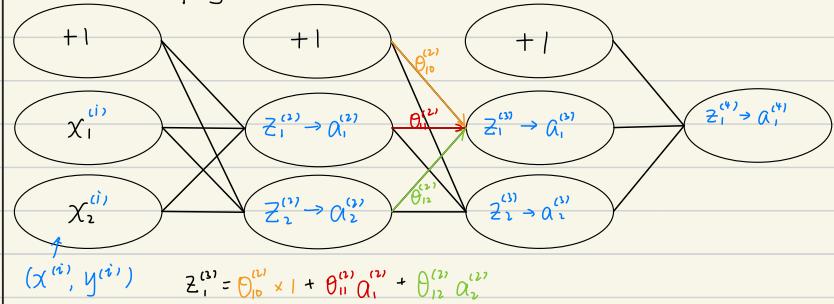
$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$  Vec:  $\Delta^{(l)} = \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$

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

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

## Backpropagation Intuition

### Forward Propagation



What is backpropagation doing?  $\rightarrow k=1$

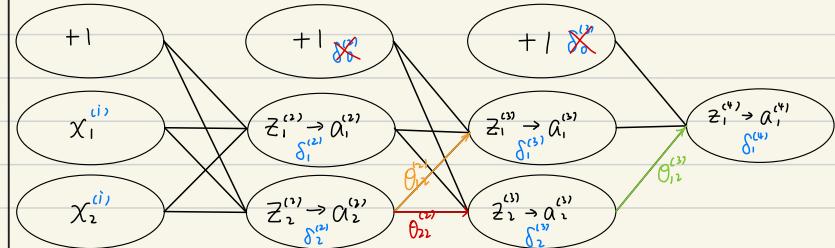
$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^m \sum_{j=1}^{S_{l+1}} (\theta_{ij}^{(l)})^2$$

Focusing on a single example  $x^{(i)}, y^{(i)}$ , the case of 1 output unit, and ignoring regularization ( $\lambda=0$ ):  $\text{cost}(i) = y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))$

(Think of  $\text{cost}(i) \approx (h_\theta(x^{(i)}) - y^{(i)})^2$ )

$\rightarrow$  how well is the network doing on example  $i$ ?

Backpropagation



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

Formally,  $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i)$  (for  $j \geq 0$ ),

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

$$\delta_1^{(4)} = a_1^{(4)} - y^{(i)} \quad \delta_2^{(4)} = \theta_{12}^{(3)} \cdot \delta_1^{(3)} \quad \delta_1^{(3)} = \theta_{11}^{(2)} \delta_1^{(2)} + \theta_{21}^{(2)} \delta_2^{(2)}$$

## Backpropagation in Practice

Implementation Note: Unrolling

Parameters

Advanced optimization

function  $[jVal, gradient] = \text{costFunction}(\theta)$   
...  $\in \mathbb{R}^{m+1}$   $\in \mathbb{R}^{m+1}$  (vectors)

$\text{optTheta} = \text{fminunc}(@\text{costFunction}, \text{initialTheta}, \text{options})$

Neural Network ( $L = 4$ ):

$\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$  - matrices ( $\text{Theta1}, \text{Theta2}, \text{Theta3}$ )

$D^{(1)}, D^{(2)}, D^{(3)}$  - matrices ( $\text{D1}, \text{D2}, \text{D3}$ )

"Unroll" into vectors

Example

$$S_1 = 10, S_2 = 10, S_3 = 1$$

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

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

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

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

$$\text{Theta1} = \text{reshape}(\text{thetaVec}(1:110), 10, 11);$$

$$\text{Theta2} = \text{reshape}(\text{thetaVec}(111:220), 10, 11);$$

$$\text{Theta3} = \text{reshape}(\text{thetaVec}(221:231), 1, 11);$$

Learning Algorithm

Have initial parameters  $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$

Unroll to get  $\text{initialTheta}$  to pass to:

$\text{fminunc}(@\text{costFunction}, \text{initialTheta}, \text{options})$

function [jval, gradientVec] = costFunction(thetaVec)

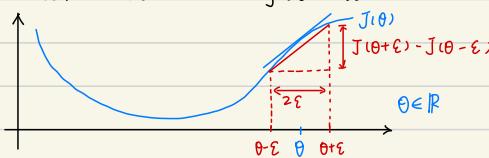
From thetaVec, get  $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$  reshape

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.

## Gradient Checking

Numerical estimation of gradients



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

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

Implement:  $\text{gradApprox} = (J(\theta + \text{EPSILON}) - J(\theta - \text{EPSILON})) / (2 * \text{EPSILON})$

Parameter vector  $\theta$

$\theta \in \mathbb{R}^n$  (E.g.  $\theta$  is "unrolled" version of  $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$ )

$$\theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_n]$$

$$\frac{\partial}{\partial \theta_1} J(\theta) \approx \frac{J(\theta_1 + \epsilon, \theta_2, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \dots, \theta_n)}{2\epsilon}$$

$$\frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta_1, \theta_2 + \epsilon, \theta_3, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \theta_3, \dots, \theta_n)}{2\epsilon}$$

⋮

$$\frac{\partial}{\partial \theta_n} J(\theta) \approx \frac{J(\theta_1, \theta_2, \theta_3, \dots, \theta_{n-1}, \theta_n + \epsilon) - J(\theta_1, \theta_2, \theta_3, \dots, \theta_{n-1}, \theta_n - \epsilon)}{2\epsilon}$$

```

for i = 1:n,
    thetaPlus = theta;
    thetaPlus(i) = thetaPlus(i) + EPSILON;
    thetaMinus = theta;
    thetaMinus(i) = thetaMinus(i) - EPSILON;
    gradApprox(i) = (J(thetaPlus) - J(thetaMinus)) / (2 * EPSILON)
end;
Check that gradApprox ≈ DVec ← From backprop

```

$\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_i + \epsilon \\ \vdots \\ \theta_n \end{bmatrix}$      $\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_i - \epsilon \\ \vdots \\ \theta_n \end{bmatrix}$

## Random Initialization

Initial value of  $\theta$

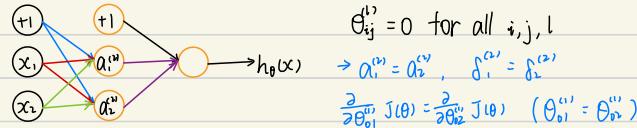
For gradient descent and advanced optimization method, need initial value for  $\theta$ .

`optTheta = fminunc (@costFunction, initialTheta, options)`

Consider gradient descent:

Set `initialTheta = zeros(n, 1)`?

Zero initialization



After each update, parameters corresponding to inputs going into each of two hidden units are identical.

Random initialization: Symmetry breaking

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

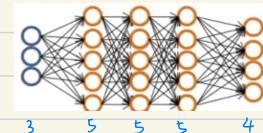
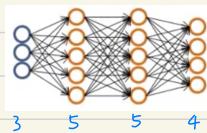
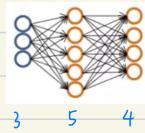
E.g.  $\Theta_{\text{theta1}} = \text{rand}(10, 11) * (2 * \text{INIT\_EPSILON}) - \text{INIT\_EPSILON};$   
 $[\epsilon, \epsilon] \quad \downarrow$   
random 10x11 matrix (between 0 and 1)

$\Theta_{\text{theta2}} = \text{rand}(1, 11) * (2 * \text{INIT\_EPSILON}) - \text{INIT\_EPSILON};$

## Putting it together

### • Training a neural network

Pick a network architecture (connectivity pattern between neurons)



No. of input units: Dimension of features  $x^{(i)}$

No. of output units: Number of classes

$$y \in \{1, 2, 3, \dots, 10\}, y = 5 \quad y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

Reasonable default: 1 hidden layer, or if > 1 hidden layer, have some no. of hidden units in every layer (usually the more the better)

### • Training a neural network

1. Randomly initialize weights
2. Implement forward propagation to get  $h_{\Theta}(x^{(i)})$  for any  $x^{(i)}$
3. Implement code to compute cost function  $J(\theta)$
4. Implement backprop to compute partial derivatives  $\frac{\partial J}{\partial \theta_j}$ ,  $J(\theta)$

for  $i=1:m$

Perform forward propagation and backpropagation using example  $(x^{(i)}, y^{(i)})$

(Get activations  $a^{(l)}$  and delta terms  $\delta^{(l)}$  for  $l=2, \dots, L$ )

$$\Delta^{(l)} = \Delta^{(l)} + \delta^{(l+1)}(a^{(l)})^T$$

5. Use gradient checking to compare  $\frac{\partial}{\partial \theta} J(\theta)$  computed using backpropagation vs. using numerical estimate of gradient of  $J(\theta)$ .

Then disable gradient checking code

- b. Use gradient descent or advanced optimization method with backpropagation to try to minimize  $J(\theta)$  as a function of parameters  $\theta$ .

$J(\theta)$  - non-convex

# Week 6 Advice for Applying Machine Learning

## Evaluating a Learning Algorithm

### Deciding what to try next

Debugging a learning algorithm:

Having implemented regularized linear regression to predict housing prices:  $J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$

What if it makes unacceptably large errors in its prediction when we test the hypothesis on a new set of houses?

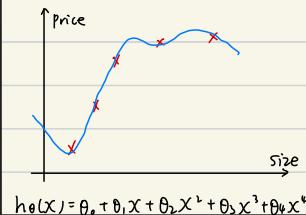
- Get more training examples
- Try smaller sets of features
- Try getting additional features
- Try adding polynomial features ( $x_0^2, x_1^2, x_1x_2$ , etc.)
- Try decreasing / increasing  $\lambda$

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.

Diagnostics can take time to implement, but doing so can be a very good use of your time.

## Evaluating a hypothesis



Fails to generalize to new examples not in training set.

Dataset:

	Size	Price	random order
Training set	2104	400	
	1600	330	$(x^{(1)}, y^{(1)})$
	2400	369	$(x^{(2)}, y^{(2)})$
	70%	232	:
	1416	232	
	3000	540	$(x^{(m)}, y^{(m)})$
	1985	300	
	1534	315	
	1427	199	$(x_{\text{test}}^{(1)}, y_{\text{test}}^{(1)})$
Test set	30%	1380	$(x_{\text{test}}^{(2)}, y_{\text{test}}^{(2)})$
	1494	212	:
		243	
			$(x_{\text{test}}^{(m_{\text{test}})}, y_{\text{test}}^{(m_{\text{test}})})$

$m_{\text{test}}$  = no. of test examples

Training/testing procedure for linear regression

- Learn parameter  $\theta$  from training data (minimizing training error  $J(\theta)$ )
- Compute test set error:  $J_{\text{test}}(\theta) = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2$

Training/testing procedure for logistic regression

- Learn parameter  $\theta$  from training data
- Compute test set error:

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

- Misclassification error (0/1 misclassification error):

$$\text{err}(h_0(x), y) = \begin{cases} 1 & \text{if } h_0(x) \geq 0.5 \text{ when } y=0, \text{ or } h_0(x) < 0.5 \text{ when } y=1 \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Test error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_0(x_{\text{test}}^{(i)}), y_{\text{test}}^{(i)})$$

## Model Selection and Train / Validation / Test Sets

Once parameters  $\theta_0, \theta_1, \dots, \theta_d$  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.

Model selection

$d$  = degree of polynomial

$$d=1 \quad 1. \quad h_0(x) = \theta_0 + \theta_1 x \rightarrow \theta^{(1)} \rightarrow J_{\text{test}}(\theta^{(1)})$$

$$d=2 \quad 2. \quad h_0(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^{(2)} \rightarrow J_{\text{test}}(\theta^{(2)})$$

$\vdots \quad \vdots \quad \vdots$

$$d=103. \quad h_0(x) = \theta_0 + \theta_1 x + \dots + \theta_{103} x^{103} \rightarrow \theta^{(103)} \rightarrow J_{\text{test}}(\theta^{(103)})$$

If choose  $\theta_0 + \dots + \theta_5 x^5$ , how well does the model generalize? Report test set error  $J_{\text{test}}(\theta^{(5)})$ ?

Problem:  $J_{\text{test}}(\theta^{(5)})$  is likely to be an optimistic estimate of generalization error.

I.e. our extra parameter ( $d$  = degree of polynomial) is fit to test set.

Dataset:

Size	Price	random order
Training set 60%	2104	$(x^{(1)}, y^{(1)})$
	1600	$\vdots$
	2400	$(x^{(m)}, y^{(m)})$
	1416	$(x_{cv}^{(1)}, y_{cv}^{(1)})$
	3000	$\vdots$
	1985	$(x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})}) \rightarrow m_{cv} = \text{no. of CV examples}$
cross validation set (CV) 20%	1534	$(x^{(1)}, y^{(1)})$
	1427	$(x^{(2)}, y^{(2)})$
test set 20%	1380	$\vdots$
	1494	$(x^{(m_{test})}, y^{(m_{test})})$
	243	

Train/validation/test error

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

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

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

Model selection

$$d=1 \quad h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{\text{cv}}(\theta^{(1)})$$

$\vdots$

$$d=10 \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \theta^{(10)} \rightarrow J_{\text{cv}}(\theta^{(10)})$$

If  $\min = J_{\text{cv}}(\theta^{(4)})$ , then choose  $d=4$ .

Pick  $\theta_0 + \theta_1 x + \dots + \theta_4 x^4$

→ Estimate generalization error for test set  $J_{\text{test}}(\theta^{(4)})$ .

## Bias vs. Variance

### Diagnosing Bias vs. Variance