

# Week 1

Supervised Learning: (Given "right answer" for each example in the data)

Mainly 2 types of problems:

- (1) Classification: discrete-valued output
- (2) Regression: real-valued output (continous)

Unsupervised Learning:

(1) Clustering

(2) Non-clustering

## Linear Regression with 1 variable

### (1) The Model

training set examples:

housing data:

size in feet<sup>2</sup> ( $X$ ) | price (\$) in 1000's ( $y$ )

2104
1416
1534
852
:

460
232
315
178
?

$\brace{m}$

Notations:

$m$  = no. of training examples

$X$ 's = "input" variable / features

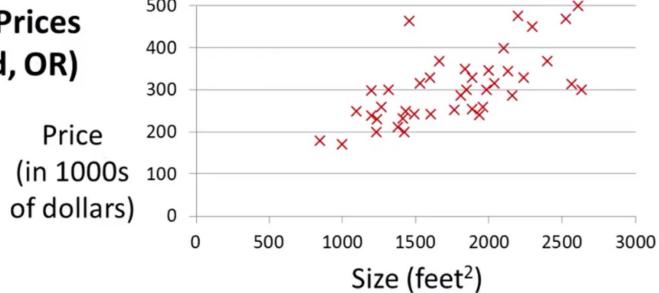
$y$ 's = "output" variable / target

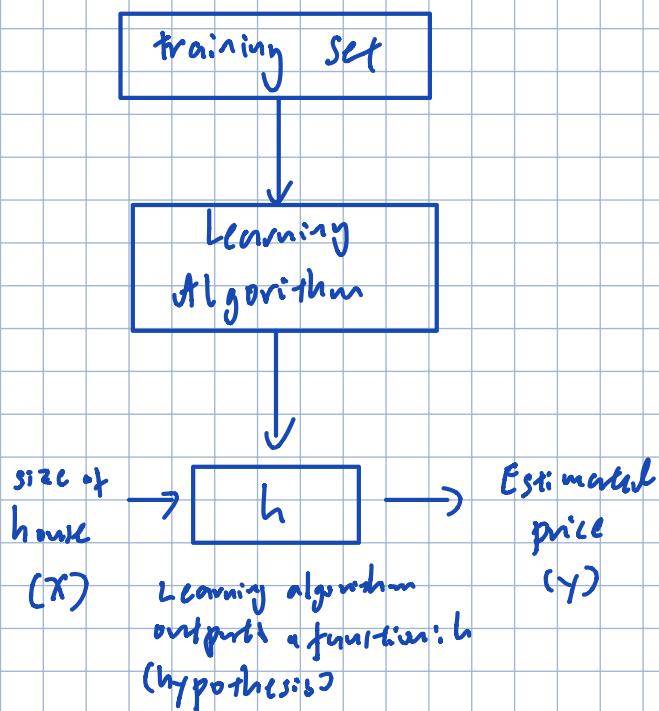
variable

Housing Prices  
(Portland, OR)

$(x, y)$  = a single training example

$(x^{(i)}, y^{(i)})$  =  $i^{\text{th}}$  training example





so  $h$  is a function  
maps from  $X \rightarrow Y$

$\rightarrow$  Supervised Learning Goal:  
given a training set, to learn a function  
 $h: X \rightarrow Y$  so that  $h(x)$  is a "good" predictor

$\uparrow$  Space of input value       $\downarrow$  Space of output value

## ② Cost Function

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

$\theta_i$ 's: Parameters of the model

How to choose  $\theta_i$ 's?

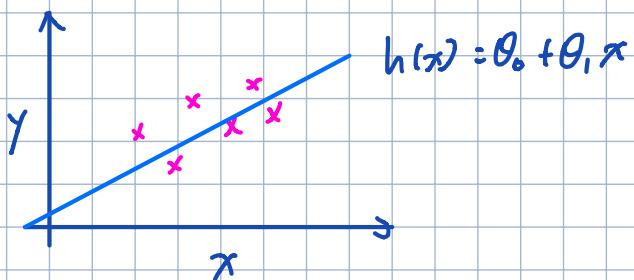
Idea:

choose  $\theta_0, \theta_1$  so that  $h_\theta(x)$   
is close to  $\textcircled{Y}$  for our  
training examples  $(x, \textcircled{y})$

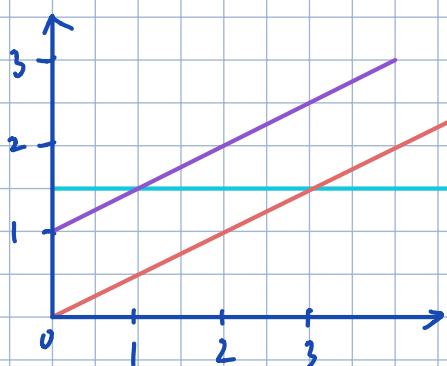
How do we represent  $h$ ?

we choose:

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



This is linear regression with 1 variable  
OR Univariate linear regression.



$$\textcircled{1} \quad \theta_0 = 1.5, \theta_1 = 0 \quad \textcircled{2} \quad \theta_0 = 0, \theta_1 = 0.5$$

$$h(x) = 1.5 + 0 \cdot x$$

$$h(x) = 0 + 0.5x$$

$$\textcircled{3} \quad \theta_0 = 1, \theta_1 = 0.5$$

$$h(x) = 1 + 0.5x$$

We want to solve a minimisation problem:

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

↑  
 find values of  
 $\theta_0, \theta_1$  to  
 minimise  
 ← marks  
 rough  
 for easier  
 minimisation

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

By convention, we define a cost function:

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

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

↓  
 $\hat{y}^{(i)}$   
 mean squared error

squared error function  
 or

i.e. it is  $\frac{1}{2} \bar{x}$ , where  $\bar{x}$ : mean of the square of  $h_\theta(x^{(i)}) - y^{(i)}$

\* mean is halved as a convenience for the computation of the gradient descent.

We can measure the accuracy of hypothesis function by using a cost function.

### A Simplified Model

$$h_\theta(x) = \theta_1 x$$

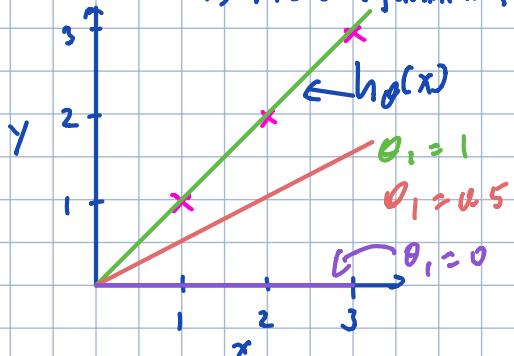
$$\text{i.e. } \theta_0 = 0$$

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

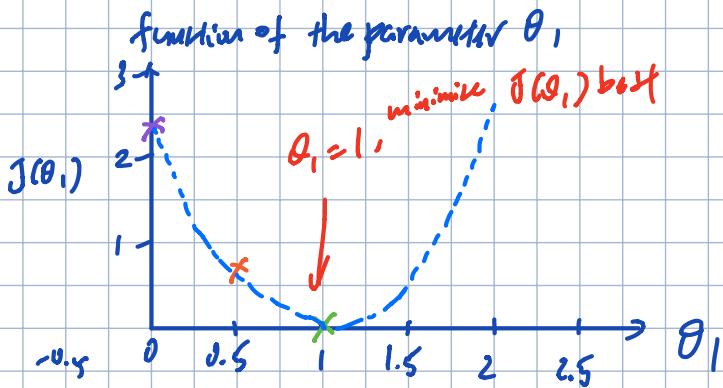
↓  
 $\theta_1 x^{(i)}$

hypothesis function  
 $h_\theta(x)$

for fixed  $\theta_1$ , this is a function of  $x$



cost function  
 $J(\theta_1)$



$$\begin{aligned}
 J(\theta_0) &= \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\
 \theta_0 &= 1 \\
 &= \frac{1}{2m} \sum_{i=1}^m (\theta_0 x^{(i)} - y^{(i)})^2 \\
 &= \frac{1}{2m} (0^2 + 0^2 + 0^2) \\
 &= 0^2
 \end{aligned}$$

i.e.  $J(1) = 0$

$$\begin{aligned}
 J(0) &= \frac{1}{2 \times 3} [1^2 + 2^2 + 3^2] \\
 &= \frac{1}{6} [14] \\
 &= 7/3
 \end{aligned}$$

$$\begin{aligned}
 J(\theta_0=0.5) &= \frac{1}{2m} [(0.5-1)^2 + (1-2)^2 \\
 &\quad + (1.5-3)^2] \\
 &= \frac{1}{2 \times 3} [3.5] \\
 &\approx 0.58
 \end{aligned}$$

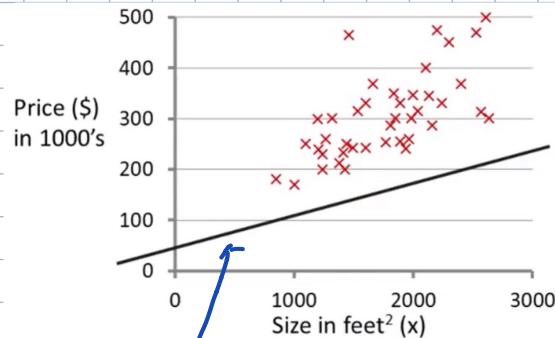
i.e.  $J(0.5) \approx 0.58$

\* Each value of  $\theta_0$  correspond to a different hypothesis function

Now with  $\theta_0$  and  $\theta_1$

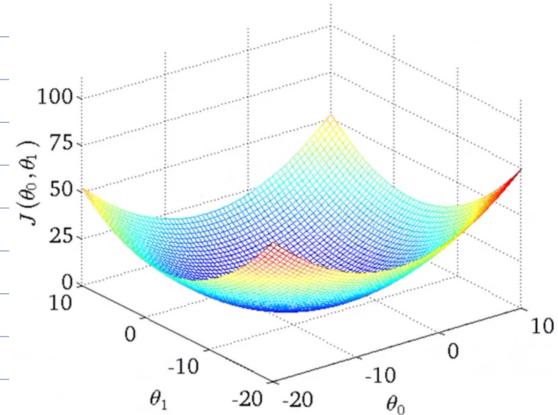
$h_{\theta}(x)$

for fixed  $\theta_0, \theta_1$ , this is  
a function of  $x$

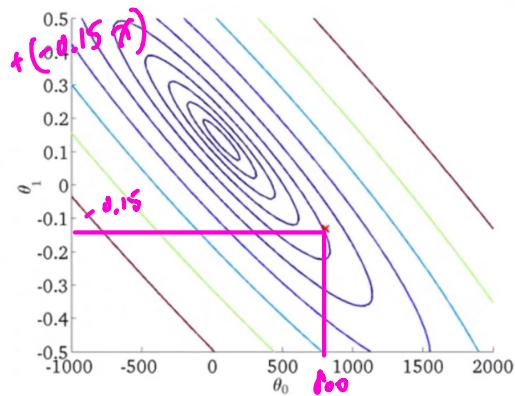
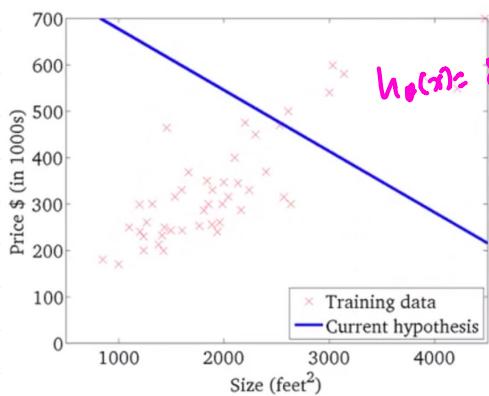


$$h_{\theta}(x) = 50 + 0.06x$$

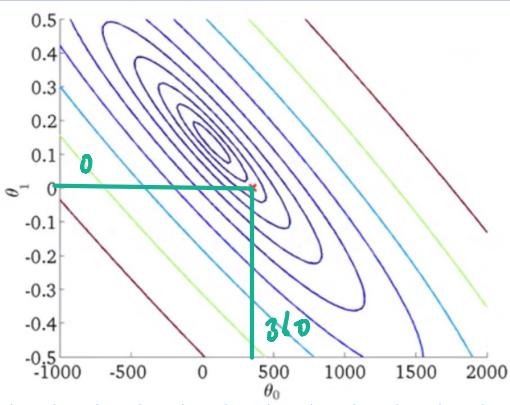
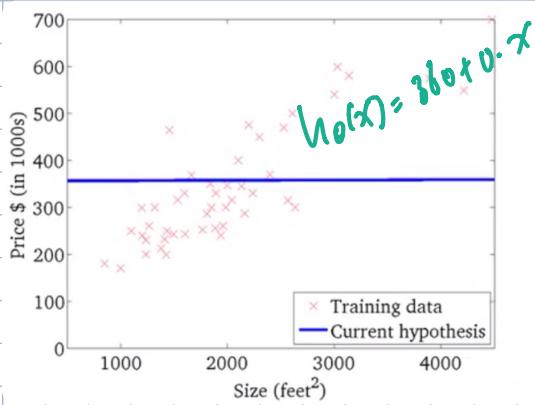
$J(\theta_0, \theta_1)$   
function of parameters  $\theta_0, \theta_1$



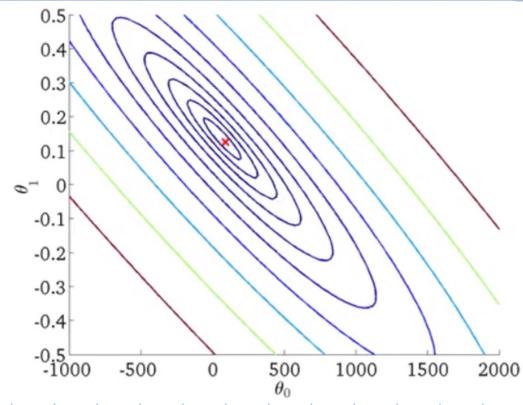
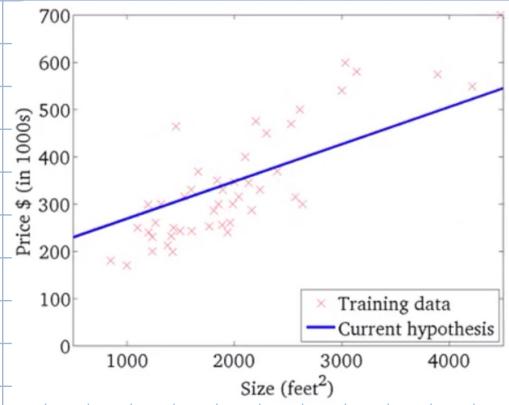
(-7)



e.g.



e.g.



We want an algorithm that finds  $\theta_0, \theta_1$  that minimise  $J$ .

### ③ Gradient Descent

It is a general technique used in all ML, can be used to minimize cost function  $J$  for linear regression.

Problem setup:

Have some function  $J(\theta_0, \theta_1)$

more generally:  $J(\theta_0, \theta_1, \theta_2, \dots, \theta_n)$

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

Want to come up with an algorithm for  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

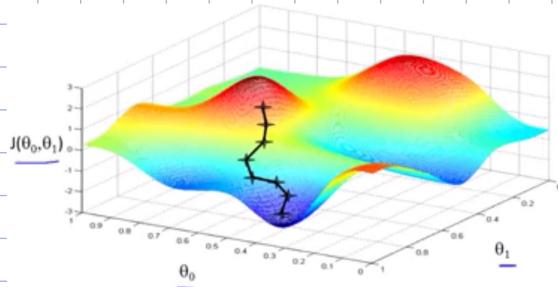
Out line:

1. start with some  $\theta_0, \theta_1$  (say, initialise to  $\theta_0 = \theta_1 = 0$ )

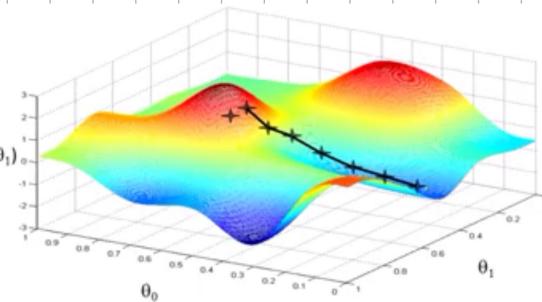
2. keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$ , until end up a minimum

Graphically:

initial point 1



initial point 2



### Gradient Descent Algorithm

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j=0 \text{ and } j=1)$$

}

assignment learning  
operator rate

taking derivatives to find function (tangent)

→ slope of the tangent is the derivative of that point → gives us a direction

correct: simultaneous update

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

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

$$\theta_0 := \text{temp0}$$

$$\theta_1 := \text{temp1}$$

incorrect

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

~~$$\theta_0 := \text{temp0}$$~~

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

~~$$\theta_1 := \text{temp1}$$~~

$$\text{e.g. } \theta_0 = 1, \theta_1 = 2, \theta_j := \theta_j + \sqrt{\theta_0 \cdot \theta_1},$$

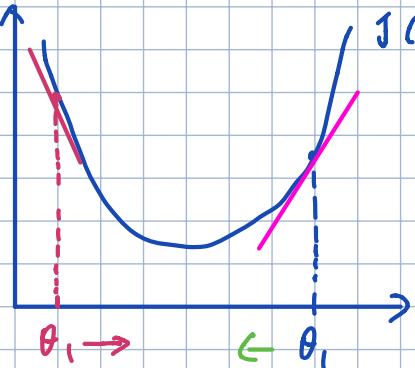
$$\theta_0 := \theta_0 + \sqrt{\theta_0 \cdot \theta_1} = 1 + \sqrt{2}$$

$$\theta_1 := \theta_1 + \sqrt{\theta_0 \cdot \theta_1} = 2 + \sqrt{2}$$

### Simplified example

1 parameter

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



$J(\theta_1)$  ( $\theta_1 \in \mathbb{R}$ )

only 1 variable  $\rightarrow$  total derivative

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

$$\theta_1 := \theta_1 - \alpha \cdot (\text{positive number})$$

$$\theta_1 := \theta_1 - \alpha \cdot (\text{negative number/slope})$$

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

If  $\alpha$  is too large, gradient descent can overshoot the minimum

$\hookrightarrow$  might fail to converge, even diverge.

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

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1) \quad \therefore \frac{d}{d\theta_1} J(\theta_1) \text{ decrease when approaching min.}$$

## ④ Use gradient descent back to linear regression problem

Gradient Descent algorithm

repeat until convergence of

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

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

}

Linear Regression model

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

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

$$\min J(\theta_0, \theta_1)$$

$\theta_0, \theta_1$

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

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

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

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

repeat until convergence {

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

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

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

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

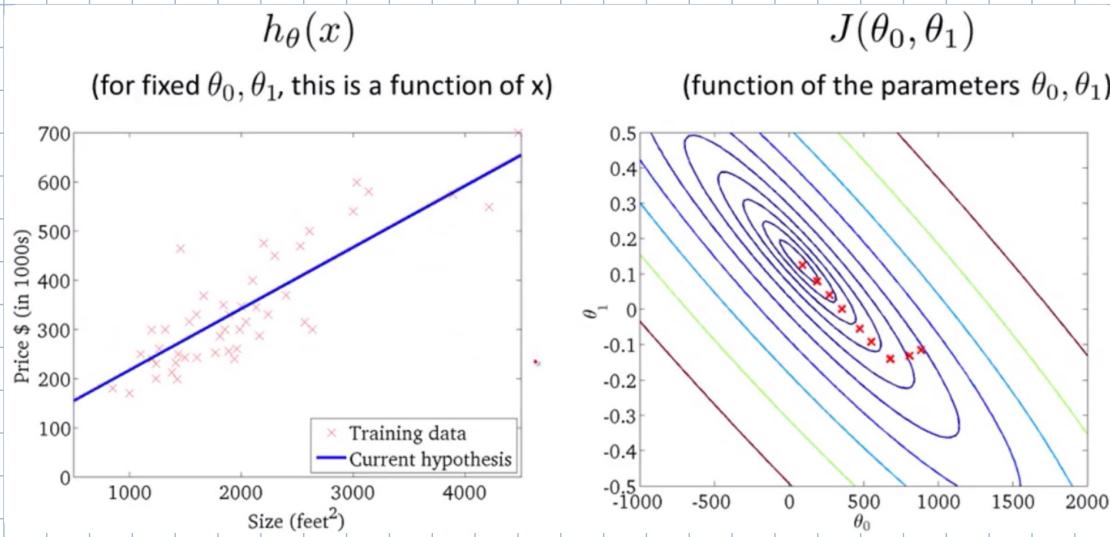
}

\* Cost function for linear Regression is always a bowl-shaped function

↳ Convex function

↳ no local optima, only global optima.

↳ Gradient descent always converge to global optima



\* "Batch" Gradient Descent

Batch : each step of gradient descent uses all the training examples.

(i.e.  $\sum_{i=1}^m$ )

There exist a method to numerically solving the min. of  $J$  w/o iteration.

↳ normal equations method

↳ but gradient descent scales better with large data sets.

## Linear Algebra Review

### Matrix

Rectangular array of numbers

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

$4 \times 2$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$$

$2 \times 3$

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

Dimension of matrix: no. of rows  $\times$  no. of columns

### matrix elements

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

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

e.g.  $A_{1,1} = 1402$

### Vector

Special case of matrix:  $n \times 1$  matrix

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

e.g.  $n=4 \Rightarrow 4 \times 1$   
also 4-dimensional vector or  $\mathbb{R}^4$

$y_i = i^{\text{th}}$  element

1-indexed vs 0-indexed

e.g.  $y_1 = 460$

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

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

## Matrix Addition

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

$3 \times 2 \quad 3 \times 2 \quad 3 \times 2$

## Scalar Multiplication

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

## Matrix - Vector Multiplication

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

$3 \times 2 \quad 2 \times 1 \quad 3 \times 1$

## Matrix - Matrix Multiplication

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

$2 \times 3 \quad 3 \times 2 \quad 2 \times 2$

## Matrix Inverse

If  $A$  is square matrix, and if it has an inverse:

$$AA^{-1} = A^{-1}A = I$$

$$\text{e.g. } \begin{bmatrix} 3 & 4 \\ 2 & 16 \end{bmatrix} \begin{bmatrix} 0.4 & -0.1 \\ -0.05 & 0.075 \end{bmatrix} = I$$

## Matrix Transpose

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

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

### Multivariate Linear Regression

Size (feet <sup>2</sup> ) $x_1$	Number of bedrooms $x_2$	Number of floors $x_3$	Age of home (years) $x_4$	Price (\$1000) $y$
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...	...	...	...	...

$\brace{m}$

$n$  = number of features

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

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

$$\text{e.g. } x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix} \in \mathbb{R}^4$$

$$x_3^{(2)} = 2$$

### Hypothesis

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

In general:

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

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

so now  $X = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_m \end{bmatrix} \in \mathbb{R}^{n+1}$        $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$

Now:

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

### Gradient Descent for multiple variables

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

Parameters:  $\theta_0, \theta_1, \dots, \theta_n$

$\theta$  is  $n+1$  dimensional vector

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

$\hookrightarrow J(\theta)$

Gradient descent:

Repeat {

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

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

}

Repeat {

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

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

}

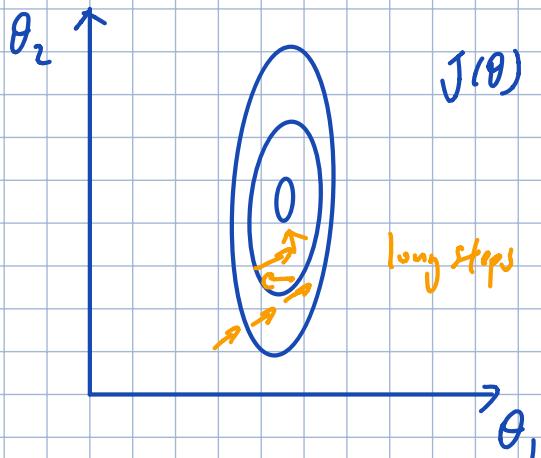
## Gradient Descent in practice - Feature Scaling

Idea: make sure features are on a similar scale

If allows gradient descent to converge more quickly

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

$x_2 = \text{no. of bedroom}$  (1-5)

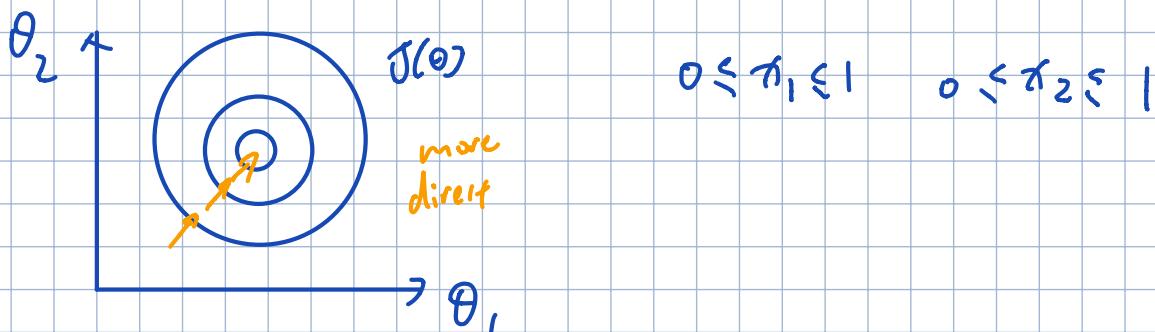


very tall, skinny, skewed

contour of  $J(\theta)$

$\therefore \theta_2$  can take on large value

$$\text{If } x_1 = \frac{\text{size(feet)}}{2000}, \quad x_2 = \frac{\text{no. of bedroom}}{5}$$



## Feature Scaling

Get every feature to approximately  $-1 \leq x_i \leq 1$  range

### Mean Normalization

Using  $x_i - \bar{x}_i$  to make features have approximately zero mean

(Do NOT apply to  $\theta_0 = 1$ )

$$\text{e.g. } x_1 = \frac{\text{size} - 1000}{2000}, \quad x_2 = \frac{\text{no. of bedrooms} - 2}{5}$$

$$-0.5 \leq x_1 \leq 0.5, \quad -0.5 \leq x_2 \leq 0.5$$

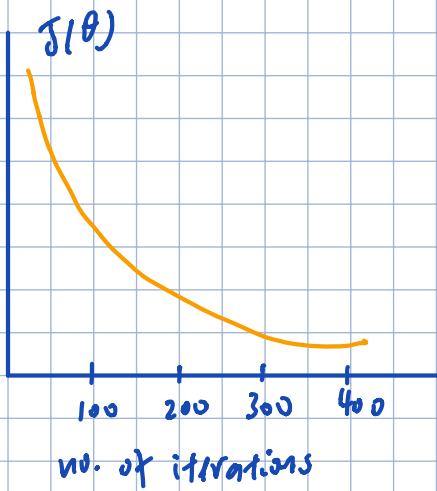
\*  $x_i \leftarrow \frac{x_i - \bar{x}_i}{s_i}$

!  $\bar{x}_i = \text{avg value of } x_i \text{ in training set.}$   
 $s_i = \text{range or std deviation}$   
 $(\max - \min)$

$$\bar{x}_i = \frac{x_{i\text{max}} - x_{i\text{min}}}{s_i} \quad (\text{for 1 feature})$$

## Gradient Descent in practice — Learning Rate

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



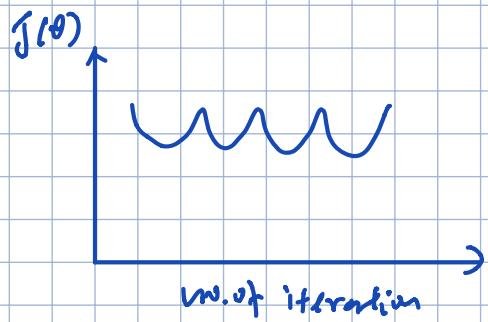
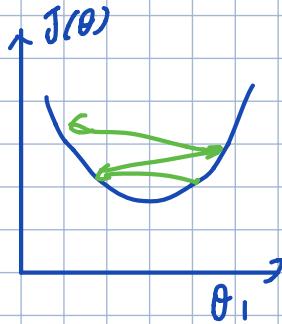
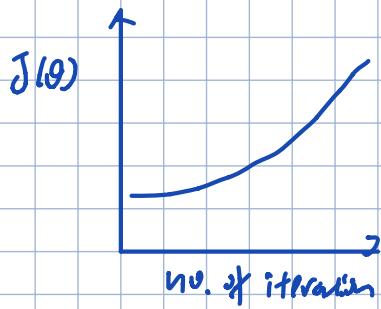
× plot  $J(\theta)$  as gradient descent runs.

$J(\theta)$  should decrease after every iteration

Example of automatic convergence test:

Declare convergence if  $J(\theta)$  decreases by less than  $10^{-3}$  in 1 iteration.

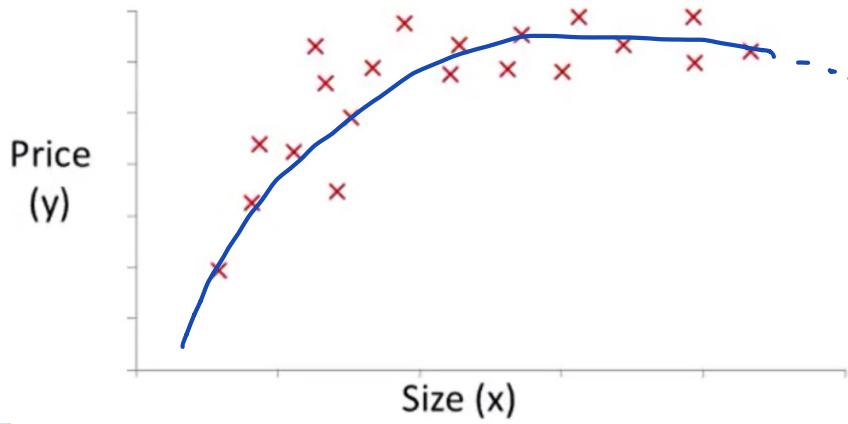
If learning rate is too big



\* For sufficiently small  $\alpha$ ,  $J(\theta)$  should decrease on every iteration

But if  $\alpha$  is too small  $\rightarrow$  slow to converge

## Polynomial Regression



quadratic  
will  
deviate  
→ not ideal

$$\begin{aligned} \text{we can try: } h_{\theta}(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\ &= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3 \end{aligned}$$

$$\left. \begin{aligned} x_1 &= (\text{size}) \\ x_2 &= (\text{size})^2 \\ x_3 &= (\text{size})^3 \end{aligned} \right\} \text{feature scaling!}$$

$$\text{we can also try: } h_{\theta}(x) = \theta_0 + \theta_1 (\text{size}) + \theta_2 \sqrt{(\text{size})}$$

## Normal Equation

A method to solve for  $\theta$  analytically.

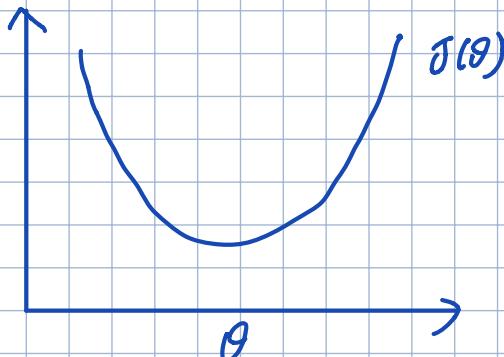
i.e. instead of iterative finding the optima, 1 step will do.

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

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

for  $\theta$  is real numbers

$$\frac{d}{d\theta} J(\theta) = 0, \text{ solve for } \theta$$



for regression problem

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

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

$$\frac{\partial}{\partial \theta_j} J(\theta) = 0$$

solve for  $\theta_0, \dots, \theta_n$

e.g.  $m=4$  (4 training examples)

new feature $\downarrow$	Size (feet <sup>2</sup> )	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
$\theta_0$	$x_1$	$x_2$	$x_3$	$x_4$	$y$
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

Construct matrix  $X$ , contains all features

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

$4 \times 5$

Construct vector  $y$ , contains target

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

$4 \times 1$   
 $m$ -dimensional vector

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

In the general case, say we have  $m$  examples:

$$\text{i.e. } (\mathbf{x}^{(1)}, \mathbf{y}^{(1)}) \dots (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})$$

and  $n$  features:

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

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(1)\top} \\ \mathbf{x}^{(2)\top} \\ \vdots \\ \mathbf{x}^{(m)\top} \end{bmatrix} \quad m \times (n+1)$$

(design matrix)

e.g. only 1 feature

$$\mathbf{x}^{(i)} = \begin{bmatrix} 1 \\ x_i^{(i)} \end{bmatrix}$$

$$\mathbf{X} = \begin{bmatrix} 1 & \mathbf{x}^{(1)} \\ 1 & \mathbf{x}^{(2)} \\ \vdots & \vdots \\ 1 & \mathbf{x}^{(m)} \end{bmatrix} \quad m \times 2$$

$$\mathbf{Y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

Normal Equations

this gives the optimum value of  $\theta$  that minimize  $J(\theta)$

$$\Rightarrow \boxed{\theta = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}}$$

OCTAVE:  $\text{pinv}(\mathbf{X}' \times \mathbf{X}) \times \mathbf{X}' \times \mathbf{y}$

\* no need for feature scaling

e.g.  $\theta_0 \approx 0.001 \approx 1$

$\theta_1 \approx 1000$

Comparison: m training examples, n features

### Gradient Descent

- need to choose  $\alpha$
- need many iteration
- works well even when n is large
- $O(kn^2)$

### Normal Equation

- no need to choose  $\alpha$
- don't need to iterate
- need to compute  $(X^T X)^{-1}$  memory  $\rightarrow O(n^3)$
- slow if n is very large  
 $n > 10^6$

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

Octave:  $\text{pinv}$      $\text{pseudo-inv}$  ← works even if  $X^T X$  is not invertible  
             $\text{inv}$          $\text{inv}$

2 most common causes for  $X^T X$  not invertible

① Redundant features (linearly dependent)

e.g.  $X_1 = \text{size in feet}^2$

$X_2 = \text{size in m}^2$

② too many features (e.g.  $m \leq n$ )

↳ delete some features

↳ use regularisation

## WEEK 3

### Classification problems

discrete prediction! Email: Spam/not spam

Online transaction: Fraudulent (Y/N)

Tumor: Malignant / Benign

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

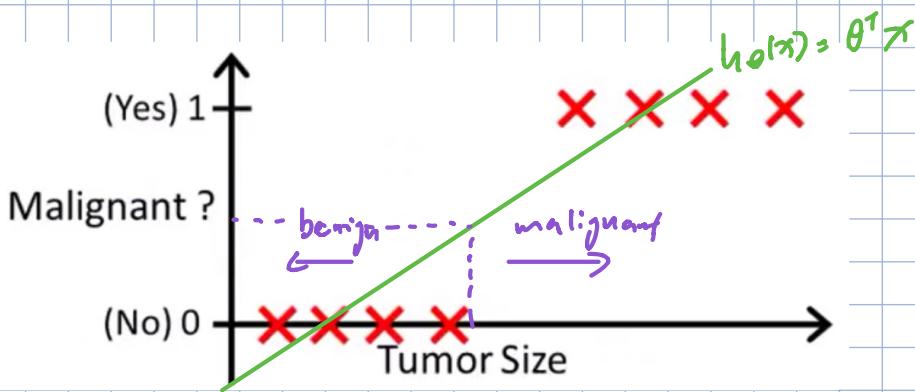
$0$ : "-ve class" (e.g. benign tumor)

$1$ : "+ve class" (e.g. malignant tumor)

→ binary classification problem

→ multi-class classification problem:  $y \in \{0, 1, 2, 3\}$

Intuition for logistic regression:



linear regression to this dataset:

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

If we want to make prediction using the linear regression model:

threshold classifier output  $h_{\theta}(x)$  at 0.5:

- { if  $h_{\theta}(x) \geq 0.5$ , predict  $y=1$
- { if  $h_{\theta}(x) < 0.5$ , predict  $y=0$

Linear regression is NOT good for classification problem

e.g.  $y=0$  or  $y=1$

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

↳ Logistic Regression:  $0 \leq h_{\theta}(x) \leq 1$  (probability)

## Logistic Regression

① hypothesis representation  $h_{\theta}(x)$ ?

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

recall linear regression:  $h_{\theta}(x) = \theta^T x$

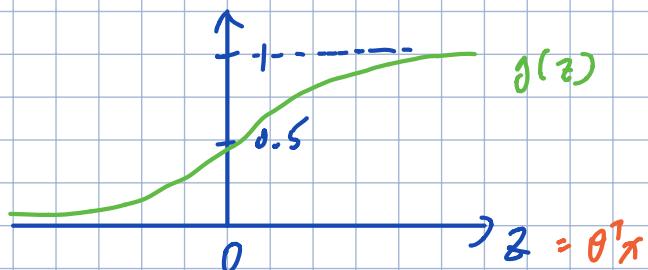
↳ for logistic regression:  $h_{\theta}(x) = g(\theta^T x)$

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

$$\left. \begin{aligned} h_{\theta}(x) \\ = \end{aligned} \right\} \frac{1}{1 + e^{-\theta^T x}}$$



Sigmoid function or logistic function



$g(z)$  asymptotes at 0 and 1  
as  $z \rightarrow \pm\infty$   
 $\begin{cases} z \rightarrow 0, e^{-z} \rightarrow 0, g(z) \rightarrow 1 \\ z \rightarrow -\infty, e^{-z} \rightarrow \infty, g(z) \rightarrow 0 \end{cases}$

so given that hypothesis is

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

⇒ given a training set, we need to pick parameters  $\theta$

⇒ then this hypothesis let us make predictions

② hypothesis output interpretation

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

e.g. if  $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumor size} \end{bmatrix}$  one of the features with

1 new sample with some  $x_0, x_1 \Rightarrow h_{\theta}(x) = 0.7$

after training  $\theta \Rightarrow 70\%$  chance of tumor being malignant

Formally:  $h_{\theta}(x) = P(Y=1 | x; \theta)$  this probability is parameterized by  $\theta$

for classification task:  $y=0$  or  $1$

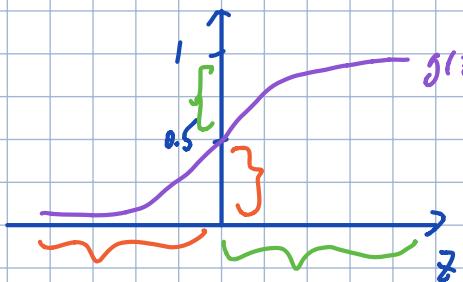
$$\hookrightarrow P(Y=0 | x; \theta) + P(Y=1 | x; \theta) = 1 \quad (\text{for the same patient})$$

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

### (3) Decision boundary

Suppose predict " $y=1$ " if  $h_{\theta}(x) \geq 0.5$

predict " $y=0$ " if  $h_{\theta}(x) < 0.5$



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

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

when  $\theta^T x \geq 0$

$\hookrightarrow$   
predict  $y=1$  when  $\theta^T x \geq 0$

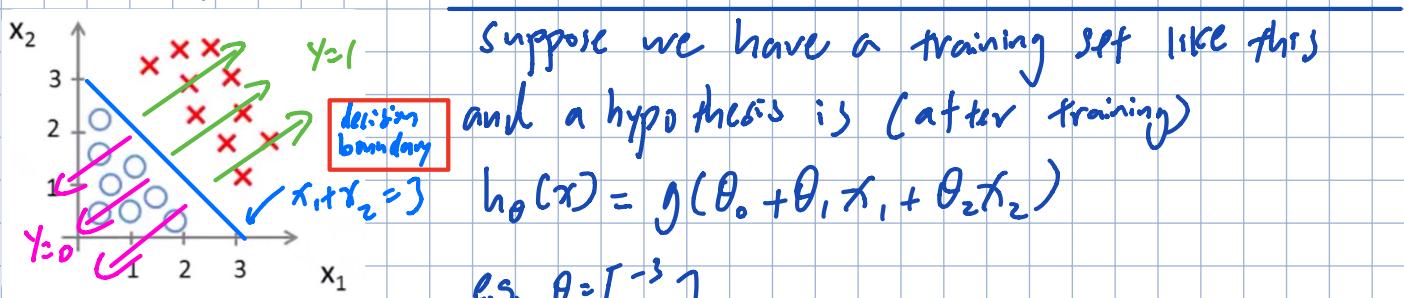
$g(z) < 0.5$  when  $z < 0$

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

when  $\theta^T x < 0$

$\hookrightarrow$   
predict  $y=0$  when  $\theta^T x < 0$

feature space



Suppose we have a training set like this

and a hypothesis is (after training)

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

$$\text{e.g. } \theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

we know: predict  $y=1$  if  $-3 + x_1 + x_2 \geq 0$

$$(-3 + x_1 + x_2) \geq 0$$

$$\hookrightarrow x_1 + x_2 \geq 3$$

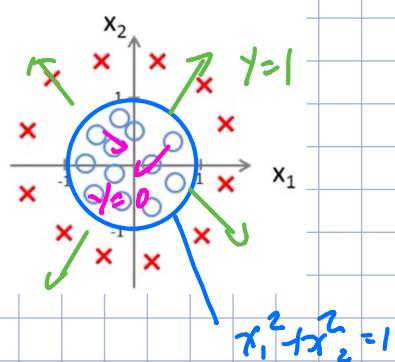
predict  $y=0$  if  $-3 + \theta_1 x_1 + \theta_2 x_2 < 0$

$$\hookrightarrow \theta_1 + \theta_2 < 3$$

here  $\theta_1 + \theta_2 = 3 \Rightarrow h_\theta(x) = 0.5$

- \* Decision boundary is a property of hypothesis  $h_\theta(x)$ , including  $\theta_0, \theta_1, \theta_2$ .
  - i.e. not a property of the data set.
    - $\hookrightarrow$  once trained,  $\theta_0, \theta_1, \theta_2$  will completely define the decision boundary.

### Non-linear boundaries



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

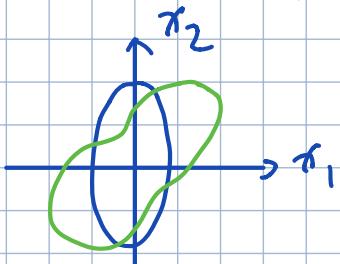
add higher order polynomial terms

obtained through training  $\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$

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

$$\hookrightarrow x_1^2 + x_2^2 \geq 1$$

even higher polynomial terms:



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

## ④ Cost Function

Recall the logistic regression problem:

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

$m$  example(s)

$n$  features

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$x_0 = 1$$

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

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

how to choose parameters  $\theta$

Recall for linear regression:

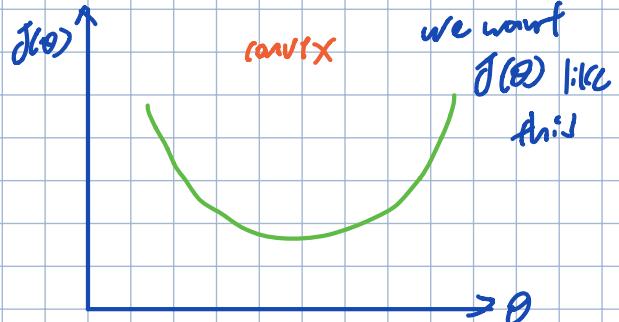
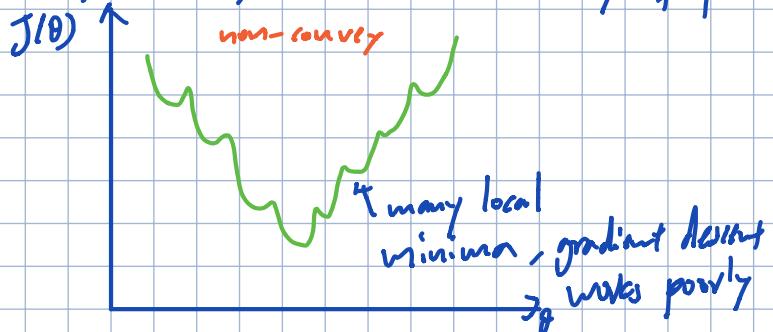
$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{m} \sum_{i=1}^m \text{cost}(h_{\theta}(x^{(i)}), y^{(i)}) \end{aligned}$$

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

a cost for learning algorithm have to pay, if output  $h_{\theta}(x^{(i)})$  where actual label is  $y^{(i)}$ .

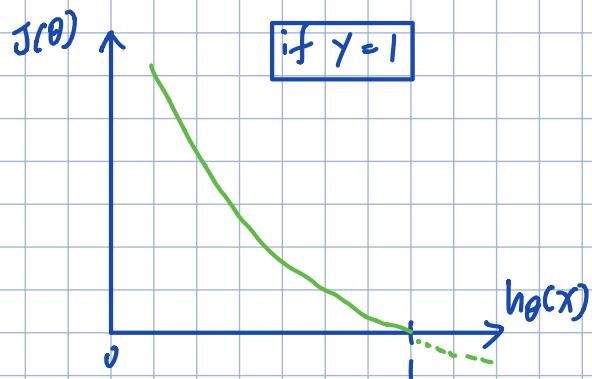
We cannot use  $\frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$  as cost function in logistic regression.

If we plug in  $h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$  into  $J(\theta)$ , we get a "non-convex" function, due to non-linearity of function,



For logistic regression

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



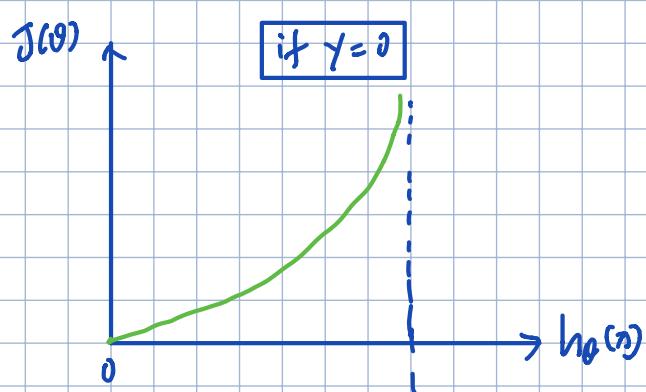
$\star$  Cost = 0 if  $y=1, h_{\theta}(x)=1$

but as  $h_{\theta}(x) \rightarrow 0, \text{cost} \rightarrow \infty$

captures intuition that if  $h_{\theta}(x)=0$

(i.e. predict  $P(y=1|x; \theta)=0$ ), but  $y=1$ ,

we penalise learning algorithm by a large cost.



Vectorized implementation:

$$h = g(X\theta)$$

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

Because  $y=0 \text{ or } 1$  always, we can simplify the cost function

one-line logistic regression cost function:

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

Logistic regression cost function

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m \text{cost}(h_{\theta}(x^{(i)}), y^{(i)}) \\ &= -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)})) \right] \end{aligned}$$

can be derived from principle of maximum likelihood estimation

To fit parameters  $\theta$ :

$$\min_{\theta} J(\theta) \Rightarrow \text{get } \theta$$

## ⑤ Gradient Descent to find $\theta$

To make a prediction given new  $x$ :

$$\text{output } h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} \Leftrightarrow p(y=1|x; \theta)$$

Using gradient descent to minimise cost function:

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

Repeat {

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

}

(Simultaneously update all  $\theta_j$ )

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

vectorised implementation

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

\* looks exactly same to linear regression!

But  $h_{\theta}(x)$  is different

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

\* feature scaling also applied to gradient descent for logistic regression.

## ⑥ Advanced Optimisation

Goal: allows logistic regression run more quickly with gradient descent,  
also allows very large ML problems, with large no. of features  
to be run.

optimisation algorithm so far:

cost function  $J(\theta)$ , want min.  $J(\theta)$

Given  $\theta$ , we have code that can compute:

$$J(\theta) \quad \leftarrow \text{to monitor convergence}$$
$$\boxed{\frac{\partial}{\partial \theta_j} J(\theta)} \quad \leftarrow \text{needed for gradient descent}$$

(for  $j = 0, 1, \dots, n$ )

what gradient descent does is:

Repeat {

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

}

\* Gradient descent is not the only algorithm for optimisation  $J(\theta)$

↳ Conjugate gradient

↳ BFGS

↳ L-BFGS

Advantages:

- \* no need to manually pick  $\alpha$
- \* often faster than gradient descent

Disadvantages:

- \* more complex

e.g.  $\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$

$$J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

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

$$\frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$$

If we want to minimise  $\min_{\theta} J(\theta)$

then  $\theta_1 = 5, \theta_2 = 5$

function [jVal, gradient] = costFunction(theta)

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

$$\text{gradient} = \text{zeros}(2, 1); \quad \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\frac{\partial}{\partial \theta_1} J(\theta) \rightarrow \text{gradient}(1) = 2 * (\theta(1) - 5); \quad \left. \begin{array}{l} \text{gradient}(2) = 2 * (\theta(2) - 5); \\ \text{element of gradient} \end{array} \right\}$$

% call advanced optimisation

options = optimset('GradObj', 'on', 'MaxIter', '100');  
data structure that stores the options we want

initialTheta = zeros(2, 1);

[optTheta, functionVal, exitFlag]

$\rightarrow \theta GR^d, d=2$

= fminunc(@(costFunction, initialTheta, options));

$\uparrow$  pointer to cost function

function optimisation unconstrained

## Logistic Regression

theta =  $\begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \rightarrow \theta(1) \text{ in octave}$

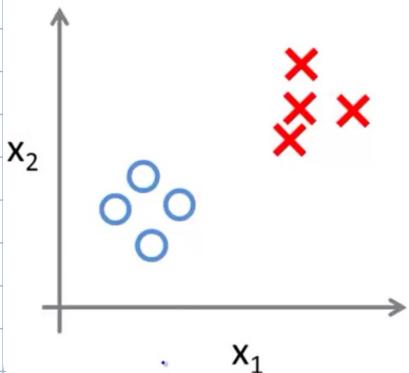
```
function [jVal, gradient] = costFunction(theta)

    jVal = [code to compute  $J(\theta)$ ];
    gradient(1) = [code to compute  $\frac{\partial}{\partial \theta_0} J(\theta)$ ];
    gradient(2) = [code to compute  $\frac{\partial}{\partial \theta_1} J(\theta)$ ];
    :
    gradient(n+1) = [code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$ ];
```

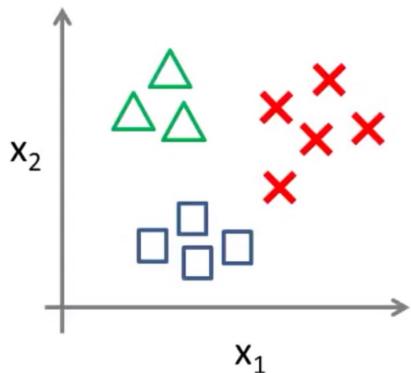
### ③ Multi-class classification:

one-vs-all

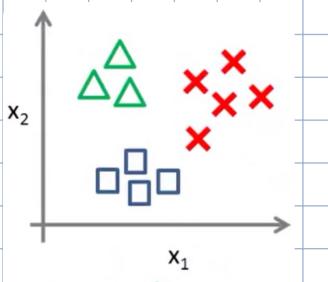
Binary classification:



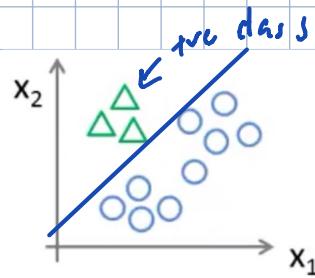
Multi-class classification:



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



Class 1:  $\triangle$   
Class 2:  $\square$   
Class 3:  $\times$



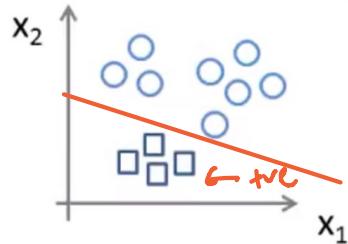
$$h_{\theta}^{(1)}(\mathbf{x}) \leftarrow \text{for 1st class}$$

$$h_{\theta}^{(1)}(\mathbf{x})$$

$$\Delta = 1$$

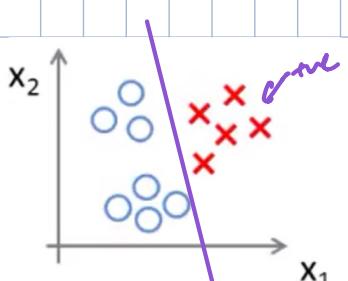
$$\circ = 0$$

$$P(Y=1 | \mathbf{x}; \theta)$$



$$h_{\theta}^{(2)}(\mathbf{x})$$

$$P(Y=2 | \mathbf{x}; \theta)$$



$$h_{\theta}^{(3)}(\mathbf{x})$$

$$P(Y=3 | \mathbf{x}; \theta)$$

In Summary:

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

$\Rightarrow$  trained 3 classifiers

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

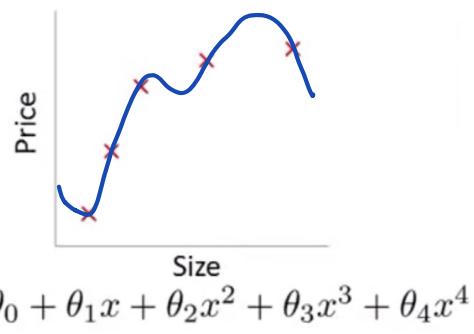
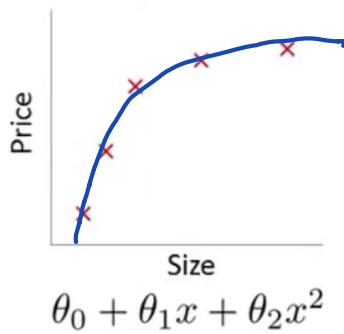
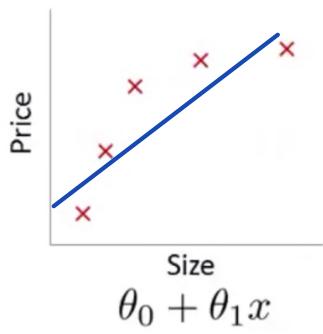
\* On a new input  $x$ , to make a prediction, pick the class  $i$  that maximises

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

i.e. run all 3 classifier, pick the class  $i$  that maximises the 3.  
(most confident)

## Overfitting

e.g. linear regression

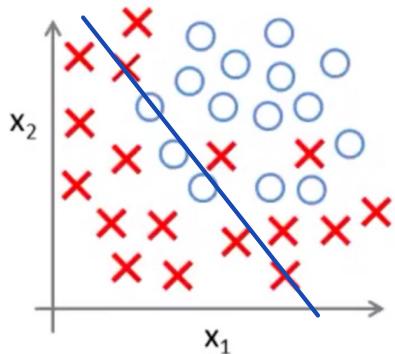


under fit  
high bias  
pre-conceptions

over fit  
high variance

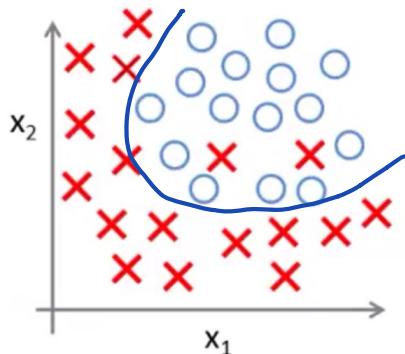
Overfitting: if we have too many features, the learned hypothesis may fit the training set very well (i.e.  $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \approx 0$ ), but fail to generalise to new examples (i.e. prediction on new examples)

## e.g. logistic regression

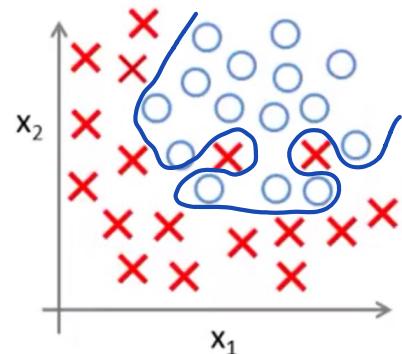


$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

( $g$  = sigmoid function)



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



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

under fit

over fit

## Addressing overfitting

2 main options:

### ① Reduce number of features

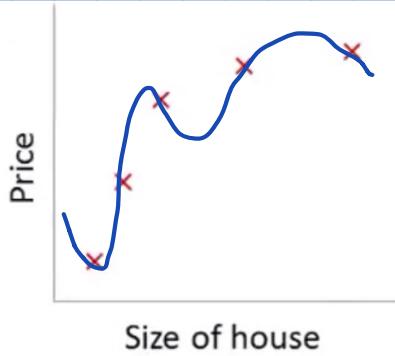
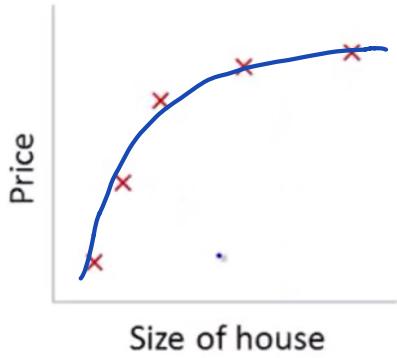
- ↳ manually select which features to keep
- ↳ model selection algorithms

### ② Regularisation

- ↳ keep all features, but reduce magnitude/value of  $\theta_j$
- ↳ works well when we have a lot of features, each of which contributes a bit to predicting  $y$ .

## Regularisation : Cost function

### Intuition



Suppose we penalise and make  $\theta_3, \theta_4$  really small

the usual squared error cost function + some extra terms on  $\theta_3, \theta_4$

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

If we were to minimise this function, the only way is let  $\theta_3$  and  $\theta_4$  to be small.  $\Rightarrow \theta_3 \approx \theta_4 \approx 0$  likely.

$\hookrightarrow$  similar to get rid of the  $\theta_3, \theta_4$  term

### Regularisation:

Small values for parameters  $\theta_0, \theta_1, \dots, \theta_n$

$\hookrightarrow$  simpler hypothesis

$\hookrightarrow$  less prone to overfitting

## Regularized Linear Regression

e.g. housing:

- features:  $x_1, x_2, \dots, x_{100}$
- parameters:  $\theta_0, \theta_1, \dots, \theta_{100}$

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

hard to pick in advance the relevant features  
Shrink all parameters  
regularisation parameter

$\theta_1, \dots, \theta_{100}$  are all shrink as a result

↳ by convention,  $\theta_0$  (i.e. the intercept) is NOT penalised.

↳ in practice, include  $\theta_0$  or not makes little difference

$\lambda$  (regularisation parameter) controls a trade off between 2 different goals:

(1)  $\sum_{i=1}^m (\hat{h}_\theta(x^{(i)}) - y^{(i)})^2$ : fit training data well

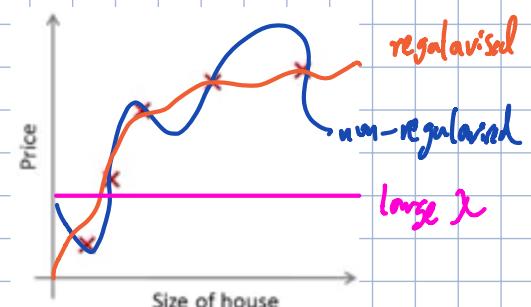
(2)  $\sum_{j=1}^n \theta_j^2$ : keeps the parameter ( $\theta$ ) small

If  $\lambda$  is set to an extremely large value:

then all  $\theta \approx 0$

↳  $\hat{h}_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \dots$

$\hat{h}_\theta(x) = \theta_0 \Rightarrow$  flat horizontal line  $\Rightarrow$  underfit  $\Rightarrow$  high bias



## Regularised Linear Regression

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

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

### Gradient Descent

Repeat {

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

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

( $j = 1, 2, \dots, n$ )

treat  $\theta_0$  separately  
 $\because$  no regularization on  $\theta_0$

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

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

regularized

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

$$1 - \alpha \frac{\lambda}{m} \text{ usually } < 1 \quad \Rightarrow \quad \theta_j \rightarrow \theta_j \times 0.99 \quad \Rightarrow \quad \theta_j^2 \text{ smaller}$$

\* This is just a mathematical representation, in computation, we just perform gradient descent on the regularised  $J(\theta)$ .

## Normal equation

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

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad R^m$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = 0$$

$$\min_{\theta} J(\theta) \Rightarrow \theta = (X^T X + \lambda \begin{bmatrix} 0 & 1 & \dots & 1 \end{bmatrix})^{-1} X^T y$$

$\underbrace{\hspace{10em}}$   
 $(n+1) \times (n+1)$

$\leftarrow$  get  $\theta$

sample feature  
Suppose  $m \leq n$

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

$\uparrow$  non-invertible

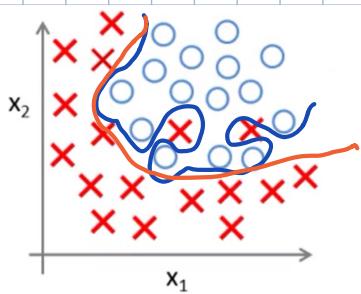
$m \leq n$ non-invertible
$m = n$ may be non-invertible

If  $\lambda > 0$

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

$\underbrace{\hspace{10em}}$   
invertible

## Regularised Logistic Regression



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \dots)$$

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 }

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

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

(j = 1, 2, ..., n)

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

## Advanced Optimisation

## Advanced optimization

function [jVal, gradient] = costFunction(theta) theta(h+1)

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

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

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

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

→ **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  $\frac{\partial}{\partial \theta_2} J(\theta)$ ];

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

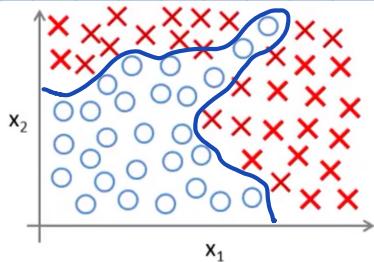
**gradient(n+1)** = [ code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$  ];

## Week 4

### Neural Networks

#### Representation: Non-linear hypotheses

Motivation:



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

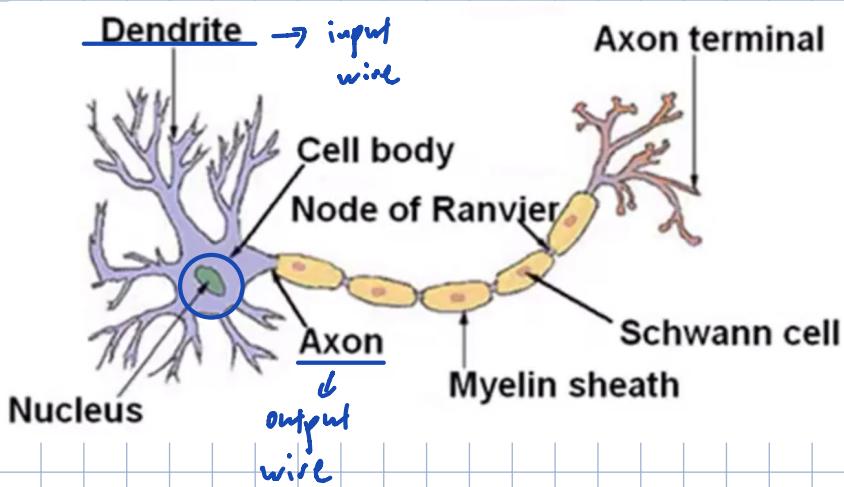
If we want to apply logistic regression for the above classification, we can include lots of polynomial terms

↳ works well for 2 features

↳ will overfit for a lot of features

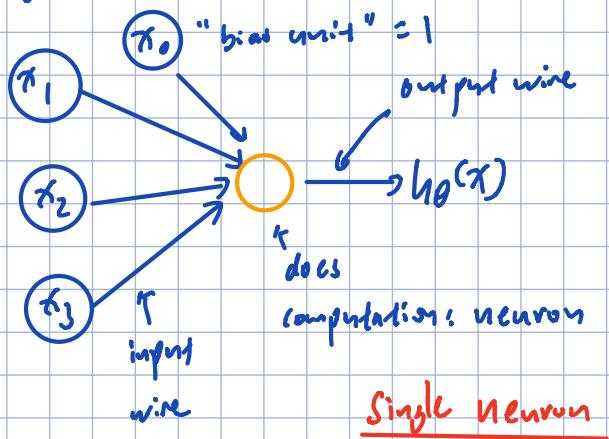
↳ can do feature selection, but will not get nice decision boundary

#### Model Representation



fun fact:  
Neuron in the  
brain

Logistic unit : a simple model



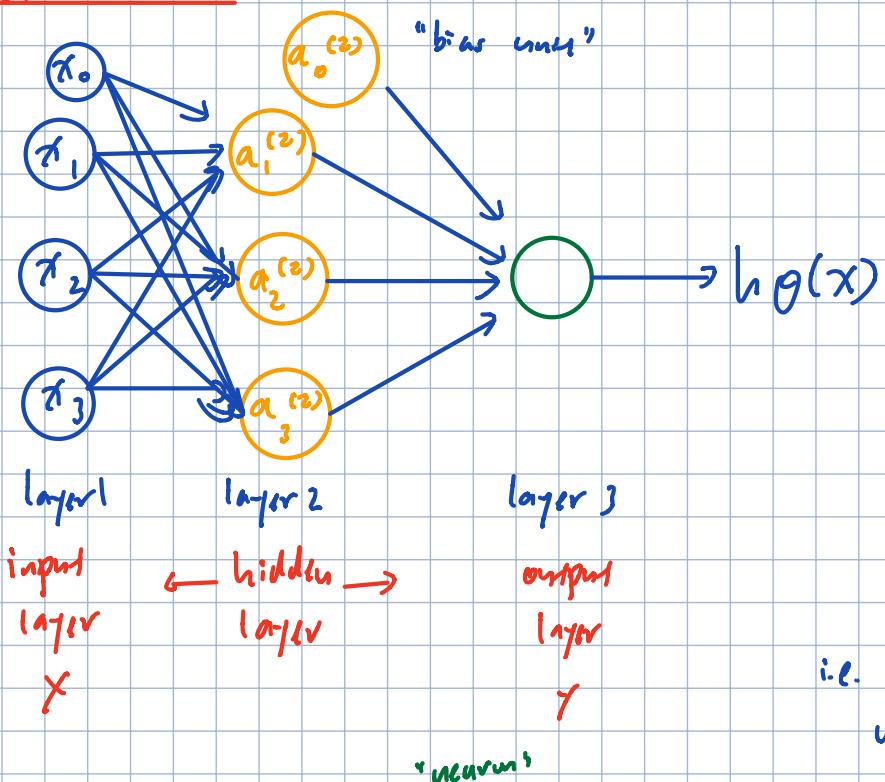
$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

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

T  
weights

↳ an artificial neuron with a Sigmoid / logistic activation function:  $g(z) = \frac{1}{1 + e^{-z}}$

### Neural Network



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

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

The computation that are represented by the diagram

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

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

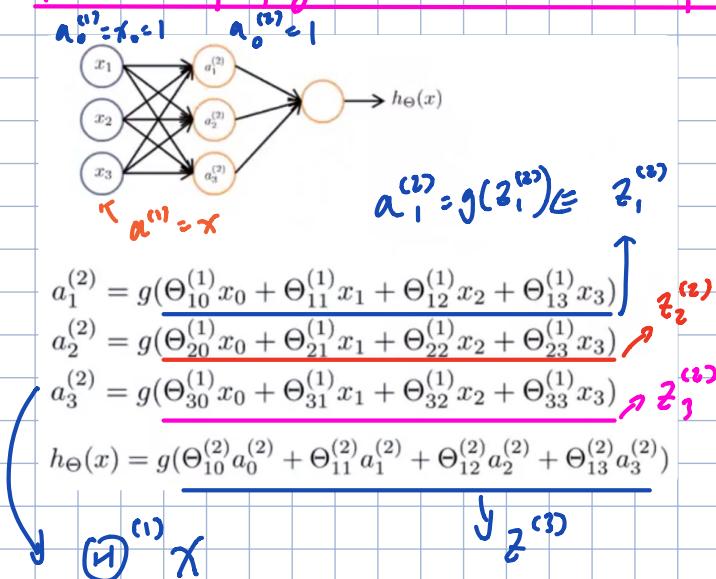
$$a_3^{(2)} = g(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3)$$

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

Here, 3 input, 3 hidden layer  $\Rightarrow \theta^{(2)} \in \mathbb{R}^{3 \times 4}$

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

### Forward propagation: Vectorized implementation



$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

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

$$z^{(2)} = \Theta^{(1)} x = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

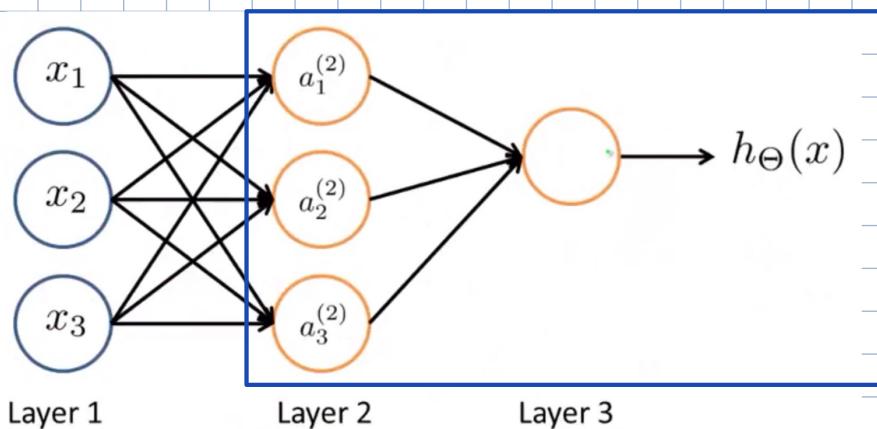
$\Downarrow \mathbb{R}^3$                    $\Downarrow \mathbb{R}^3$       element-wise application

$$\hookrightarrow \text{Add } a_0^{(2)} = 1 \quad \Rightarrow a^{(2)} \in \mathbb{R}^4$$

$$\hookrightarrow z^{(3)} = \Theta^{(2)} a^{(2)}$$

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

Neural Network learning its own features!



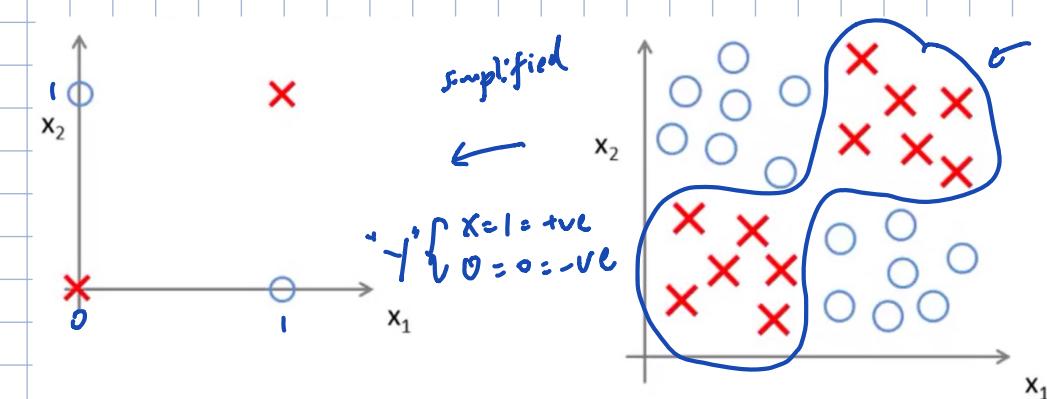
$$h_{\Theta}(x) = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})$$

↳ logistic regression  
with 3 "features":  
 $a_0^{(2)}, a_1^{(2)}, a_2^{(2)}, a_3^{(2)}$

but these are learned  
from actual feature  
 $x_1, x_2, x_3$   
through  $\Theta^{(2)}$

### Examples and Intuition

Consider we have features  $x_1$  and  $x_2$  that are binary values. (0 or 1)



With to learn a  
non-linear  
decision boundary

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

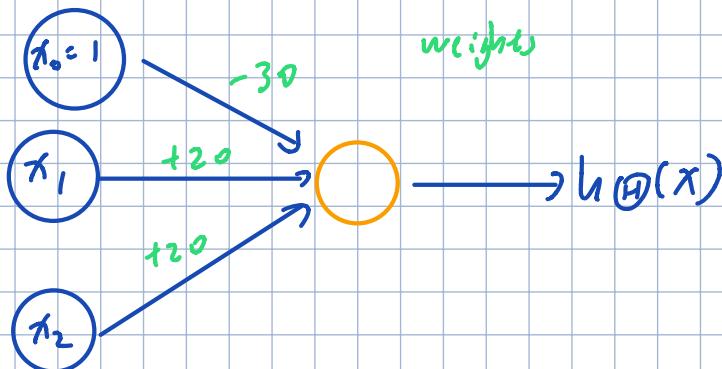
$$x_1 \text{ XNOR } x_2$$

A	B	A XOR B	A XNOR B
0	0	0	1
0	1	1	0
1	0	1	0
1	1	0	1

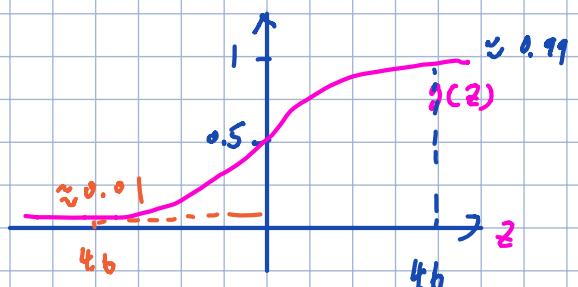
In order to build up to a network that fits the XNOR example, we start with a slightly simpler one and show a network that fits the AND function:

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

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



A	B	$\sigma$ AND B
0	0	0
0	1	0
1	0	0
1	1	1



mathematically, the hypothesis:

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

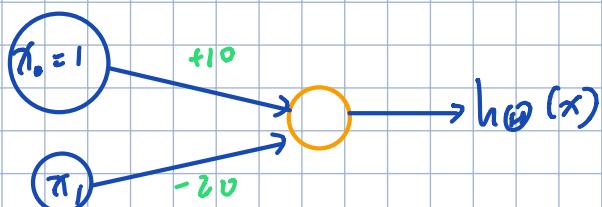
$\Theta_{00}^{(1)}$      $\Theta_{11}^{(1)}$      $\Theta_{12}^{(1)}$

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

AND

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

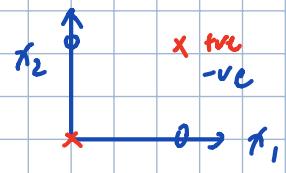
Negation: NOT  $x_1$



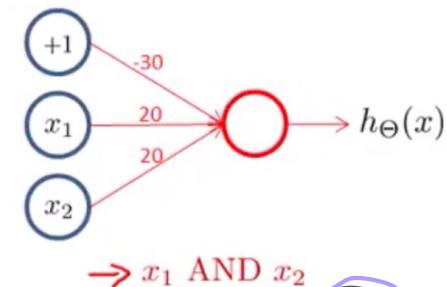
$$h_{\theta}(x) = g(10 - 20x_1)$$

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

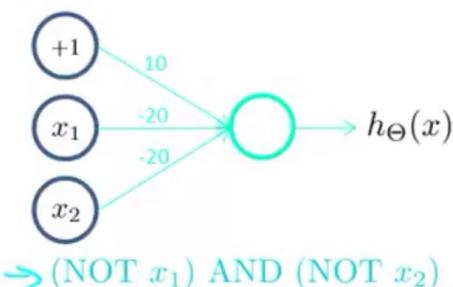
$\pi_1, XNOR \pi_2$



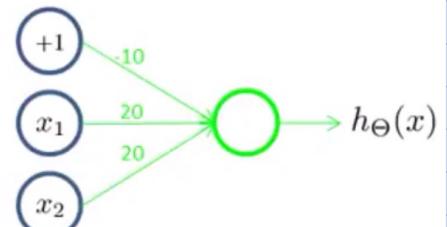
Putting it together:  $x_1$  XNOR  $x_2$



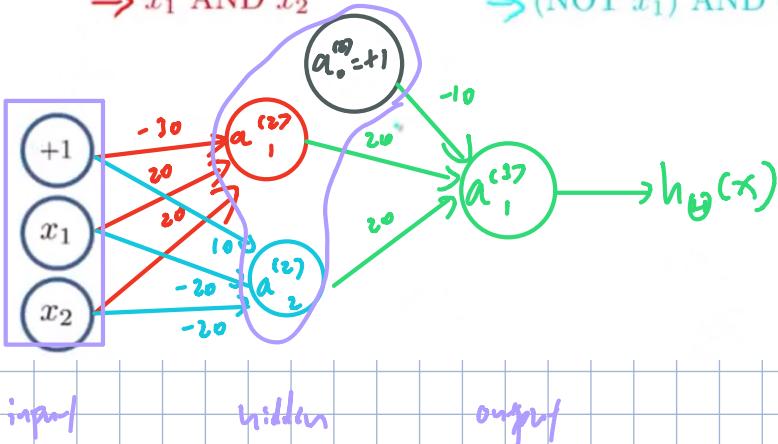
$\rightarrow x_1 \text{ AND } x_2$



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



$x_1 \text{ OR } x_2$



input

hidden

output

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

Neural Network: multi-class classification

Multiple output units: One-vs-all.



Pedestrian



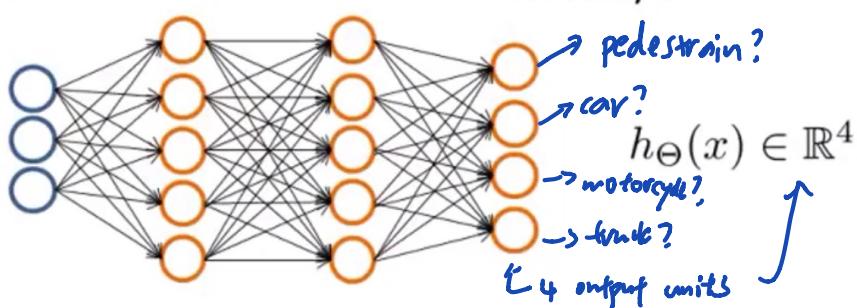
Car



Motorcycle



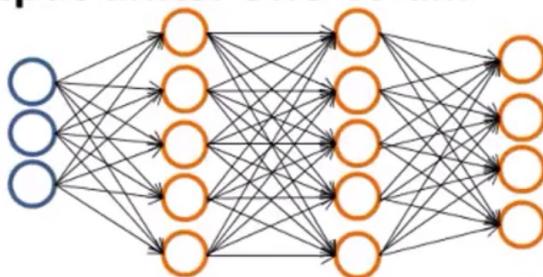
Truck



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

an extension of 1-vs-all method

## Multiple output units: One-vs-all.



$$h_{\Theta}(x) \in \mathbb{R}^4$$

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

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

$y^{(i)}$  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

↑ previously  
↑  $\{1, 2, 3, 4\}$

$(x^{(i)}, y^{(i)})$   
↑  
image

network gives  $h_{\Theta}(x^{(i)}) \approx y^{(i)}$   
↑  
 $\mathbb{R}^4$

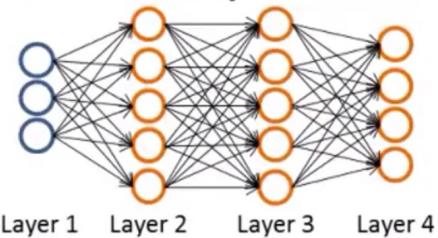
## Week 5

### Neural Networks: Cost function

cost function for fitting the parameters of the network.

#### Classification problems

#### **Neural Network (Classification)**



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

$L$  = total no. of layers in network e.g.  $L=4$

$s_l$  = no. of units (not counting bias unit) in layer  $l$  e.g.  $S_1=3, S_2=5, S_3=3, S_4=4$

2 types of classification problems:

##### ① Binary classification

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

1 output unit

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

$$S_L = 1 \quad \text{or } K=1$$

##### ② Multi-class classification ( $K$ classes)

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

$K$  output units

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

$$S_L = K \quad (K \geq 3)$$

The cost function we use for the neural network is going to be a generalisation of the one we use for logistic regression.

## Logistic Regression Cost Function

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

Instead of having just 1 compression output unit, we have K of them

## Neural Network cost function

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

↑ sum over K outputs for each training example

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

$K=1$  for binary classification

$$\theta: s_{l+1} \times s_l$$

don't sum over bias value:

i.e. don't sum over the terms where  $i=0$

↪  $\theta_{j,0}^{(l)} x_0 + \theta_{j,1}^{(l)} x_1 + \dots$  1st hidden layer

$a_0$   
not included when  $i=0$

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

\* triple sum simply adds up the square of all the individual  $\theta_j$  in the entire network

## Back propagation algorithm

An algorithm that is trying to minimise the cost function,

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

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

Need code to compute: either gradient descent or other advanced algorithm

- $J(\Theta)$
- $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$   $\textcircled{H}_{ij}^{(l)} \in \mathbb{R}$

### Gradient computation

consider the case with one training

example:  $(x, y)$

① Apply forward propagation:

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

$$z^{(2)} = \textcircled{H}^{(1)} a^{(1)}$$

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

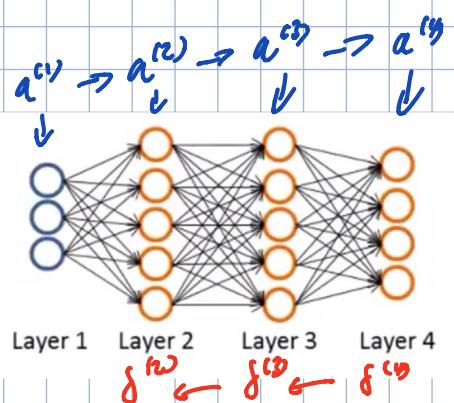
$$z^{(3)} = \textcircled{H}^{(2)} a^{(2)}$$

$$a^{(3)} = g(z^{(3)}) \quad (\text{add } a_0^{(3)})$$

$$z^{(4)} = \textcircled{H}^{(3)} a^{(3)}$$

$$a^{(4)} = h_\Theta(x) = g(z^{(4)})$$

compute what a hypothesis actually outputs given this input  $x$ .



(2) Back propagation:  $\rightarrow$  compute  $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$

Intuition:  $\delta_j^{(l)}$  = "error" of node  $j$  in layer  $L$   
 $t$  for each node

recall  $a_j^{(l)}$ : activation of the  $j$ th unit in layer  $L$ .

Here, for each output unit (layer  $L=4$ )

$$\delta_j^{(4)} = \underbrace{a_j^{(4)} - y_j}_{\text{for each node}} \rightarrow (h_\theta(x))_j$$

earlier layers:

$$\delta^{(3)} = (\Theta^{(3)})^T f^{(4)} \cdot \underbrace{g'(z^{(3)})}_{\text{vector element wise multiplication}}$$

$$\delta^{(2)} = (\Theta^{(2)})^T f^{(3)} \cdot \underbrace{o'(z^{(2)})}_{\text{vector element wise multiplication}}$$

no  $\delta^{(1)}$ : just features

vector element wise multiplication  
 vector

derivatives of activation functions  
 $o'$  evaluated on input value  
 given by  $z^{(3)}$

$$= a^{(3)} \cdot (1 - a^{(3)})$$

vectors vectors vectors  
 of  $l_3$

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

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} f_i^{(l+1)} \quad \left( \begin{array}{l} \text{ignoring } \lambda; \text{ regularization} \\ \lambda=0 \end{array} \right)$$

Putting together:

$$\rightarrow \text{to get } \frac{\partial}{\partial \Theta_{ij}} J(\Theta)$$

Backpropagation algorithm (for the entire training set)

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

→ set  $\Delta_{ij}^{(l)} = 0$  (for all  $l, i, j$ ), eventually used to compute

$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$ , in the process, they are used as accumulator

that will slowly add things in order to get  $\frac{\partial}{\partial \Theta_{ij}} J(\Theta)$ .

→ loop through training set

For  $i=1$  to  $m$  for  $i^{\text{th}}$  iteration, works with  $(x^{(i)}, y^{(i)})$

↳ set  $a^{(1)} = x^{(i)}$  set input  $x$  as  $a^{(1)}$ , first layer

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

↳ using  $y^{(i)}$ , compute  $\delta^{(L)} = a^{(L)} - y^{(i)}$  hypo - actual

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

↳  $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$  vectorise:  $\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$   
not run from  $1 \rightarrow m$ .

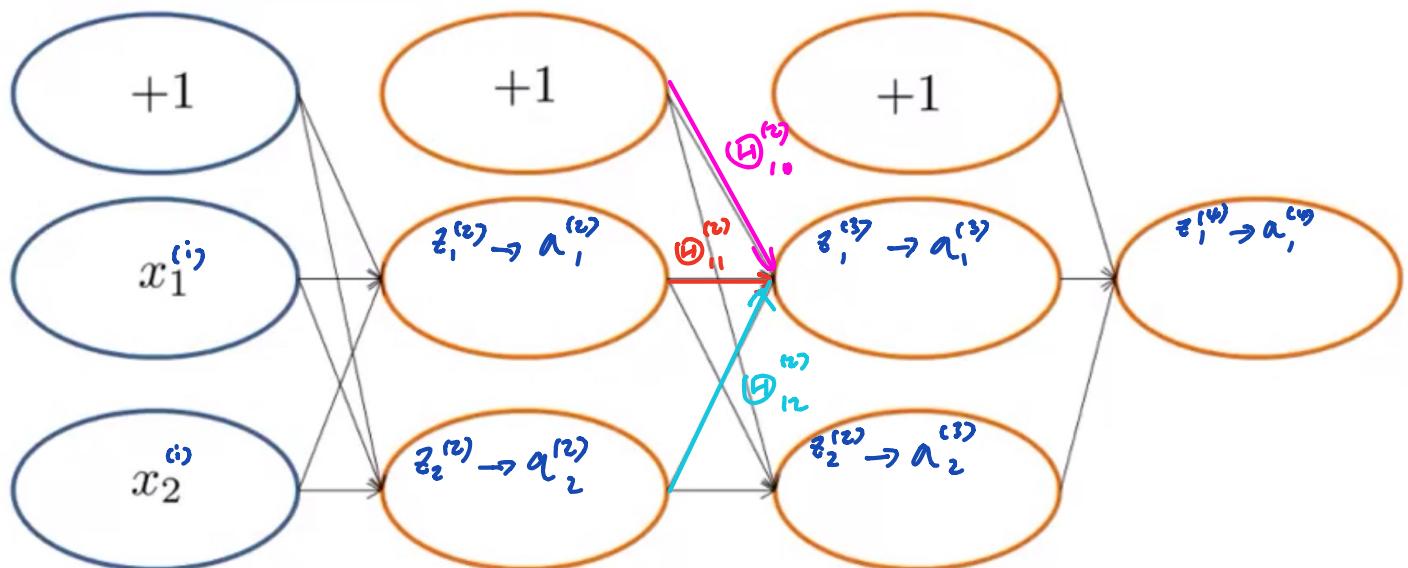
→  $D_{ij}^{(l)} = \frac{1}{m} (\Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)})$  if  $j \neq 0$  should be no. of nodes in  $l+1$  layer

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

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

↳ this gradient term can be used for gradient descent etc.

## Back propagation : intuition



$$(x^{(i)}, y^{(i)}) \uparrow_{\text{feed}} \quad z_1^{(3)} = \theta_{10}^{(2)} \times 1 + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)}$$

### What is backpropagation doing?

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

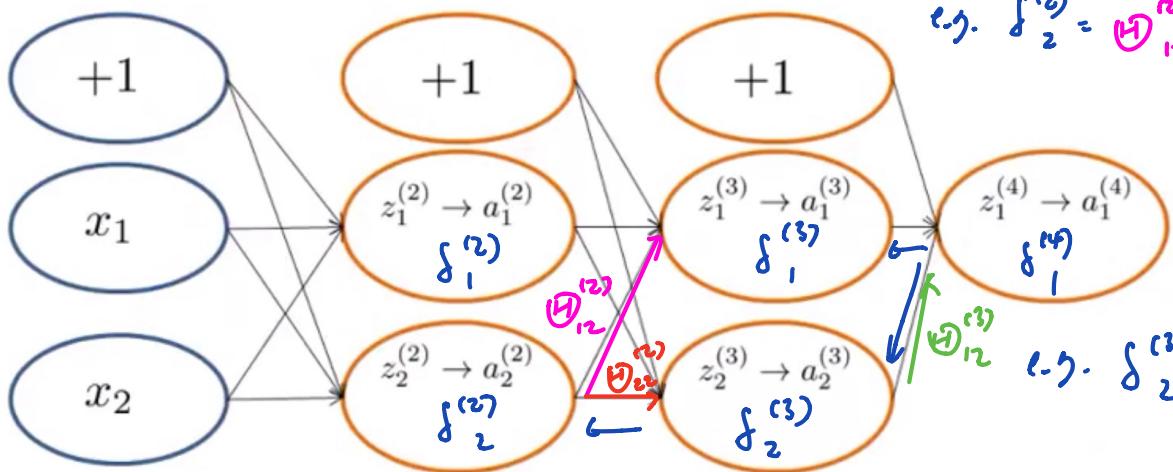
*k=1*

Focusing on a single example  $x^{(i)}$ ,  $y^{(i)}$ , the case of 1 output unit, and ignoring regularization ( $\lambda = 0$ ),

$$\text{cost}(i) = y^{(i)} \log(h_\Theta(x^{(i)})) + (1 - y^{(i)}) \log \overline{h_\Theta(x^{(i)})} \quad \leftarrow \text{for } (x^{(i)}, y^{(i)})$$

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

I.e. how well is the network doing on example i?



$$\delta_1^{(4)} = a_1^{(4)} - y^{(4)}$$

e.g.  $\delta_2^{(2)} = \Theta_{12}^{(2)} \delta_1^{(3)} + \Theta_{22}^{(2)} \delta_2^{(3)}$

e.g.  $\delta_2^{(3)} = \Theta_{12}^{(3)} \delta_1^{(4)} + \Theta_{22}^{(3)} \delta_2^{(4)}$

$\delta_j^{(l)}$  = "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  $j \geq 0$ ), where  
 $\text{cost}(i) = y^{(i)} \log h_\Theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\Theta(x^{(i)}))$

varying  $z_j^{(l)}$  to change the value

of intermediate output  
 $\rightarrow$  change cost function

$(1 - h_\Theta(x^{(i)}))$

$\delta_j$  so  $\delta$  are measure of how much would we like to change the neural network's weights, in order to affect those intermediate values ( $z_j$ ).  $\Rightarrow$  affect final output  $h_\Theta(x)$   $\Rightarrow$  affect overall cost.

### Implementation note: unrolling parameters

unrolling parameters from matrices into vectors, which we need in order to use the advanced optimisation routines.

#### Advanced optimisation

function [JVal, gradient] = cost Function(theta)

input  $\theta$

$\rightarrow$  output  $J(\theta) - \frac{\partial}{\partial \theta} J(\theta)$

...

$\text{optTheta} = \text{fminunc}(@\text{cost Function}, \text{initial Theta}, \text{options})$

Both routines assume that theta and initial theta are parameters vectors (e.g.  $\in \mathbb{R}^{n+1}$ ). Also assumes that the gradient (2nd return value) is also a vector  $\in \mathbb{R}^{n+1}$ .

This works fine when we were using logistic regression, but now for neural networks, our parameters are no longer vectors.

e.g. for a full neural network ( $L=4$ ), we have

$\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)} \rightarrow$  matrices (Theta1, Theta2, Theta3)

$\Theta^{(l)}:$    
 no. of  
 nodes  
 in l  
 layer no. of feature + 1

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

gradients that are returned

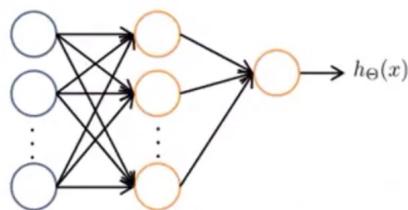
↪ "unroll" those matrices into vectors, so they end up being suitable to be used in the octave code above (for fminunc)

### Example

$$s_1 = 10, s_2 = 10, s_3 = 1$$

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

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



In Octave, convert these matrices  $\rightarrow$  vectors

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

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

Vector  $\rightarrow$  matrix

pull the 1st 110 terms

Theta 1 = reshape(thetaVec(1:110), 10, 11);

Theta 2 = reshape(thetaVec(111:220), 10, 11);

Theta 3 = reshape(thetaVec(221:330), 1, 11);

e.g. Learning Algorithm

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

$\rightarrow$  unroll to get initial Theta to pass to  
 $fminunc(@costFunction, initialTheta, options)$

implementation of the cost function:

function [Jval, gradientVec] = costFunction(thetaVec)

From thetaVec, get  $(\Theta^{(1)})$ ,  $(\Theta^{(2)})$ ,  $(\Theta^{(3)})$ . (redage)

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.

rolled into 1 vector



## Gradient Checking

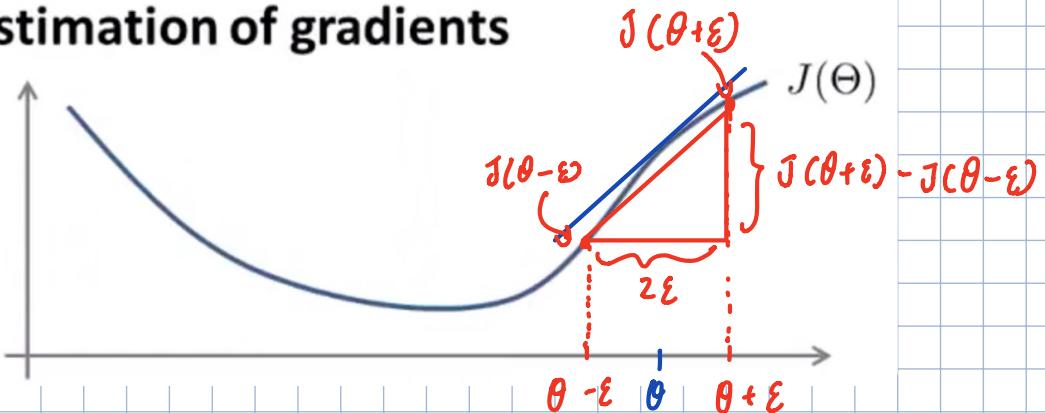
With gradient descent implementation, sometimes the cost function  $J(\Theta)$  of the neural network decreases with each iteration, the implementation can be buggy.

$\Rightarrow$  gradient checking can be used to eliminate these problems.

$\Rightarrow$  make sure that forward prop and backward prop is correct.

e.g. assume  $\theta \in \mathbb{R}$  (i.e.  $\theta$  is just a real number)

## Numerical estimation of gradients



procedure for numerically approximating the derivatives

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

Octave:

$$\text{grad Approx} = C J(\theta_{\text{eta}} + \text{EPSILON}) - J(\theta_{\text{eta}} - \text{EPSILON}) / (2 * \text{EPSILON})$$

$\frac{J(\theta + \epsilon) - J(\theta)}{\epsilon} \Rightarrow$  1-sided differences  
 $\frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} \Rightarrow$  2-sided differences

We considered here  $\theta$  is a rolled number,

more generally,  $\theta$  can be a vector parameter

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

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

similar idea

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

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

:

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

## Octave implementation:

```
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;
```

unrolled version: a long list of all  $\theta$  in the neural network

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

$\nabla \approx \frac{\partial}{\partial \theta_i} J(\theta)$

↳ check that  $\text{gradApprox} \approx \text{DVec}$

[regularized gradients averaged]

over the training set

↳ if true, much more confident that  
the implementation of back prop is correct.

↳ then plug these DVec vectors into gradient descent or other advanced algorithm.

## Numerical Gradient Checking

- implement backprop to compute DVec (unrolled  $D^{(1)}, D^{(2)}, D^{(3)}$ )
- implement numerical gradient check to compute gradApprox.
- make sure they give similar values.
- turn off gradient checking. Using backprop code for learning.

## Random Initialisation

when running gradient descent and advanced optimisation method,  
we need initial value for  $\theta$ .

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

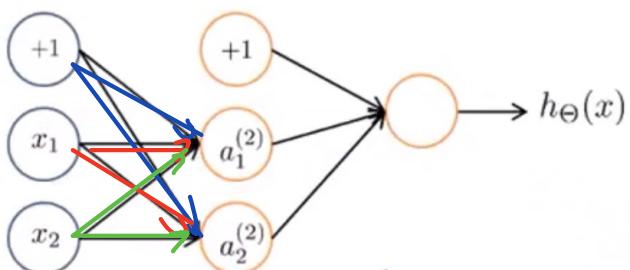
consider gradient descent

$\hookrightarrow$  done set initialTheta = zeros(n, 1) ?

$\hookrightarrow$  works OK for logistic regression

$\hookrightarrow$  DOES NOT work for neural network

why? consider training the following neural network



$$\Theta_{ij}^{(l)} = 0 \text{ for all } i, j, l.$$

$$\text{this gives } a_1^{(2)} = a_2^{(2)} \Rightarrow \delta_1^{(1)} = \delta_2^{(1)}$$

$$\frac{\partial}{\partial \Theta_{01}^{(1)}} J(\Theta) = \frac{\partial}{\partial \Theta_{02}^{(1)}} J(\Theta)$$

after 1 gradient descent iteration/update :  $\Theta_{01}^{(1)} = \Theta_{02}^{(1)}$

all will end up some non-zero value

$\Rightarrow$  after each update, parameters corresponding to inputs going into each of 2 hidden units are identical.

$\hookrightarrow$  the neural network really can't compute interesting/complex function. All nodes in hidden layer will compute the same thing.

\* Need random initialisation : symmetry breaking

$\hookrightarrow$  Initialise each  $\Theta_{ij}^{(l)}$  to a random value in  $[-\epsilon, \epsilon]$

$$\text{i.e. } -\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$$

e.g.  $\Theta_{0,1} = \text{rand}(10,11) * (2 * \text{INIT\_EPSILON}) - \text{INIT\_EPSILON};$

$\Theta_{0,2} = \text{rand}(1,11) * (2 * \text{INIT\_EPSILON}) - \text{INIT\_EPSILON};$

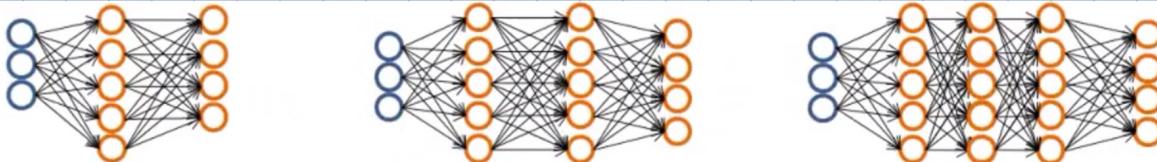
random  $10 \times 11$  matrix (between 0 and 1)

$[-\epsilon, \epsilon]$  this  $\epsilon$  is not the same as the one used in gradient checking.

## Putting together: Neural Network Training

### Choosing architecture

when training a neural network, first we need to pick a network architecture (i.e. connectivity patterns between neurons)



No. of input units : Dimension of feature  $X^{(i)}$

No. of output units : No. of classes (for classification problems)

$$Y \in \{1, 2, \dots, 10\} \Rightarrow Y = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ \vdots \\ 1 \end{bmatrix}$$

if  $Y=5$ , then at output nodes, the vector =  $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

For the hidden layers:

reasonable default: 1 hidden layer, or  $\geq 1$  hidden layer, have same no. of hidden units in every layer (usually the more the better).

$\Rightarrow$  usually the no. of hidden units in each layer will be comparable to the dimension of  $X$  (e.g.  $3 \times, 5 \times$ )

## Training neural network

① Randomly initialise weights

usually randomise the weights to small values near 0.

training example i

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  $\frac{\partial}{\partial \theta_j^{(l)}} J(\Theta)$

usually with a for loop over the training examples

for  $i = 1:m$  {

perform forward prop and back prop using example  $(x^{(i)}, y^{(i)})$

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

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

...

}

..

$\tilde{\Delta}^{(l)}$ : compute  $\frac{\partial}{\partial \theta_j^{(l)}} J(\Theta)$ ; (including regularisation terms)

⑤ Use gradient checking to compute  $\frac{\partial}{\partial \theta_j^{(l)}} J(\Theta)$  computed using

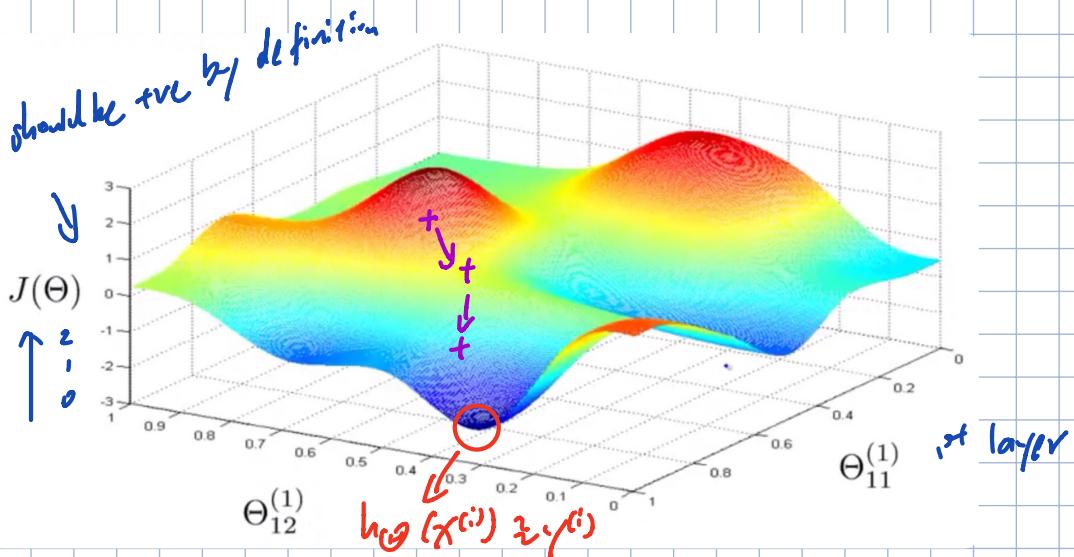
backprop  $\checkmark$  using numerical estimate of gradient of  $J(\Theta)$

$\hookrightarrow$  then disable gradient checking code

⑥ Use gradient descent or advanced optimisation method with  
backprop to try minimise  $J(\Theta)$  as a function of parameters  $(\theta)$ .

$\checkmark$

$\frac{\partial}{\partial \theta_j^{(l)}} J(\Theta)$ ;  $J(\Theta)$  is non-convex for neural network  
 $\hookrightarrow$  might have local minima



$J(\theta)$  measures how well the neural network fits the training data.

back prop computes the direction of the gradient

↳ gradient descent takes little steps down hill.

## Week 6

### Debugging a learning algorithm

Suppose we have implemented regularised linear regression to predict (e.g. house price)

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

However, if we test the hypothesis on a new set of houses, we find that it makes very large errors in its predictions, what should we do next?

- Get more training examples
- try smaller sets of features
- try getting additional features
- try adding polynomial features ( $\mathbf{x}_1^2, \mathbf{x}_2^2, \mathbf{x}_1\mathbf{x}_2$ , etc)
- try decrease/increase  $\lambda$

### ⇒ Machine Learning diagnostic:

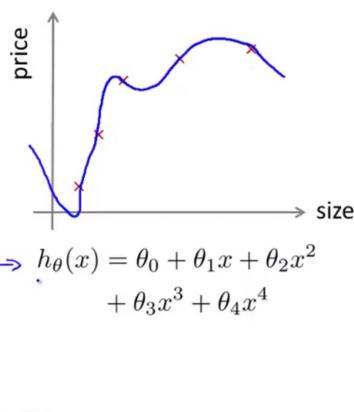
A test that you can run to gain insight what is/isn't working with a learning algorithm, and gain guidance as to how best to improve its performance.

Takes time to implement, but a good use of time.

## Evaluating a hypothesis

Evaluating a hypothesis that has been learned by our algorithm.

### Evaluating your hypothesis



Fails to generalize to new examples not in training set.

$$\begin{aligned} x_1 &= \text{size of house} \\ x_2 &= \text{no. of bedrooms} \\ x_3 &= \text{no. of floors} \\ x_4 &= \text{age of house} \\ x_5 &= \text{average income in neighborhood} \\ x_6 &= \text{kitchen size} \\ &\vdots \\ x_{100} & \end{aligned}$$

Hypothesis with low training error not necessarily a good fit  $\rightarrow$  overfit

The standard way to evaluate a learned hypothesis is as follows:

Dataset:

Size	Price
2104	400
1600	330
2400	369
1416	232
3000	540
1985	300
1534	315
<hr/>	
1427	199
1380	212
1494	243

70%

30%

training set

test set

$(x^{(1)}, y^{(1)})$

:

$(x^{(m)}, y^{(m)})$

$(x_{\text{test}}^{(1)}, y_{\text{test}}^{(1)})$

$(x_{\text{test}}^{(\text{cur}))}, y_{\text{test}}^{(\text{cur}))})$

Wt test.

no. of test examples

$\rightarrow$  better to randomly split train/test set

## Training/testing procedure for linear regression

$\rightarrow$  learn parameter  $\theta$  from training data (minimizing training error  $J(\theta)$ )

or  
70% data

→ Compute test set error:

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

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

↳ sometimes there is an alternative test set that might be easier

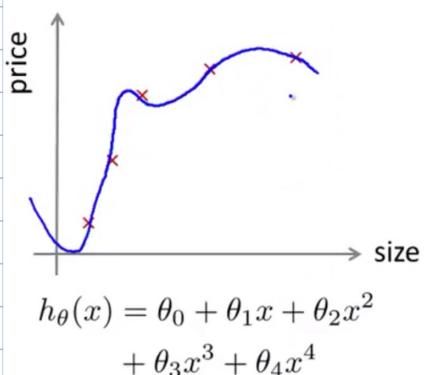
to interpret: Misclassification error (0/1 misclassification error):

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

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

## Model selection

### Overfitting example



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.

## Model selection

after fitting training data

$$d=1 \text{. } h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \theta^{(1)} = \begin{pmatrix} \theta_0 \\ \theta_1 \end{pmatrix} \text{ for } d=1 \rightarrow J_{\text{test}}(\theta^{(1)})$$

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

$$d=3 \text{. } h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \rightarrow \theta^{(3)} \rightarrow J_{\text{test}}(\theta^{(3)})$$

:

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

as it a new parameter in this algorithm

$d$  = degree of polynomial we pick

based on lowest lowest test set error

say choose  $\theta_0 + \dots + \theta_5 x^5$

How well does this model generalise? Report test set error  $J_{\text{test}}(\theta^{(5)})$

It is likely not fair:

Problem:  $J_{\text{test}}(\theta^{(5)})$  is likely to be an optimistic estimate of generalisation error. i.e. our extra parameter ( $d$  = degree of polynomial) is fit to the set.

To address this problem

Dataset:

Size	Price	
2104	400	training set
1600	330	
2400	369	
1416	232	
3000	540	
1985	300	
1534	315	cross-validation (cv) set
1427	199	
1380	212	test set
1494	243	

$(x^{(1)}, y^{(1)})$

:

$(x^{(m)}, y^{(m)})$

$(x_{cv}^{(1)}, y_{cv}^{(1)})$

:

$(x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

$(x_{test}^{(1)}, y_{test}^{(1)})$

:

$(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

$m_{cv}$ : no. of cv example

## Train/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:

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

Instead of using the test set to select the model, we can use validation (or CV) set to select the model.

## Model selection

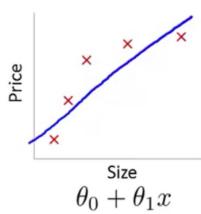
1.  $h_\theta(x) = \theta_0 + \theta_1 x \xrightarrow{\min_{\theta} J(\theta)} \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
2.  $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \xrightarrow{\theta^{(2)}} J_{cv}(\theta^{(2)})$
3.  $h_\theta(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \xrightarrow{\theta^{(3)}} \vdots$
- ⋮
10.  $h_\theta(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \xrightarrow{\theta^{(10)}} J_{cv}(\theta^{(10)})$

Say  $J_{cv}(\theta^{(4)})$  is the lowest

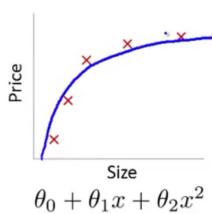
↳ estimate generalization error for test set  $J_{test}(\theta^{(4)})$

## Diagnose: Bias vs Variance

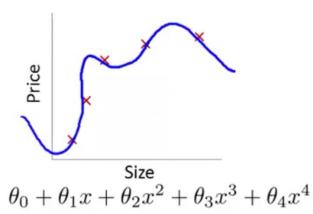
### Bias/variance



High bias  
(underfit)  $d=1$



"Just right"  $d=2$

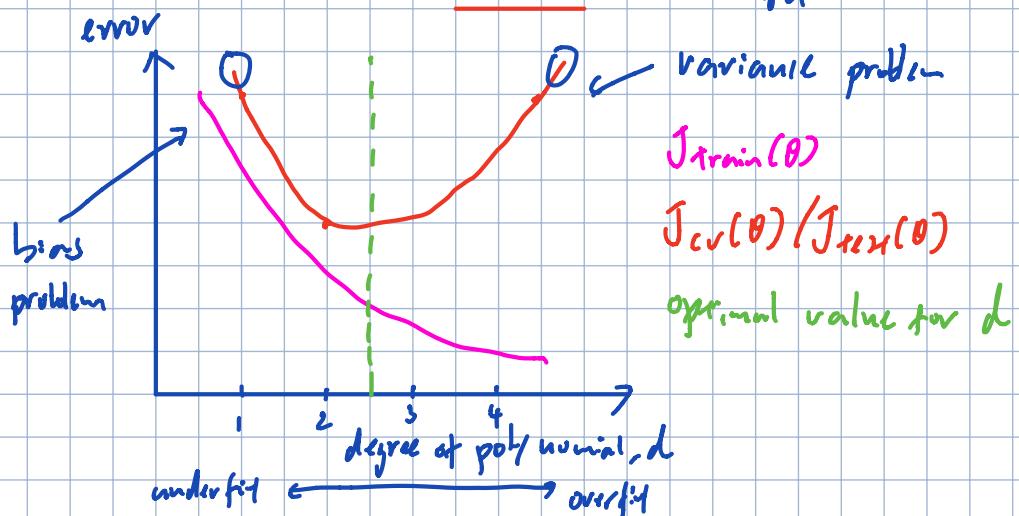


High variance  
(overfit)  $d=4$

## Bias/variance

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

cross validation error:  $J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (\hat{y}_{cv}^{(i)} - y_{cv}^{(i)})^2$



Suppose our learning algorithm is performing less well than we were hoping.

( $J_{\text{cv}}(\theta)$  or  $J_{\text{train}}(\theta)$  is high). Is it bias or variance problem?

Bias (underfit):

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

$J_{\text{cv}}(\theta) \approx J_{\text{train}}(\theta)$

Variance (overfit):

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

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

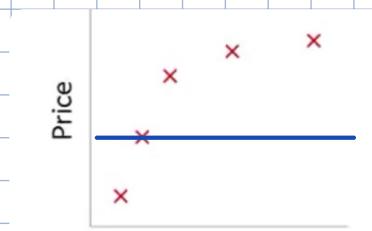
## Regularisation and bias/variance

Suppose we're fitting a high order polynomial, for linear regression with regularisation:

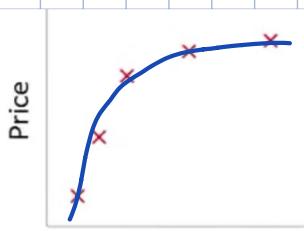
$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

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

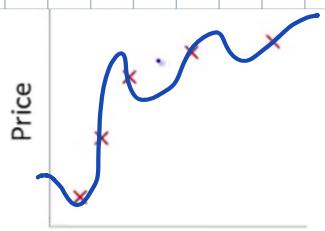
Let's consider 3 cases:



Large  $\lambda$   
High bias (underfit)  
 $\lambda = 10000, \theta_1 \approx 0, \theta_2 \approx 0, \dots$   
 $h_\theta(x) \approx \theta_0$



Intermediate  $\lambda$   
"Just right"



Small  $\lambda \approx 0$   
High variance (overfit)

If we define:

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

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

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

(w/o regularisation)

1. try  $\lambda=0$   $\rightarrow$  min  $J(\theta)$   $\rightarrow \theta^{(1)}$   $\rightarrow J_{\text{cv}}(\theta^{(1)})$

2. try  $\lambda=0.01$

⋮ ⋮ ⋮

3. try  $\lambda=0.02$

⋮ ⋮ ⋮

4. try  $\lambda=0.04$

⋮ ⋮ ⋮

5. try  $\lambda=0.08$

⋮ ⋮ ⋮

⋮

n. try  $\lambda=10$

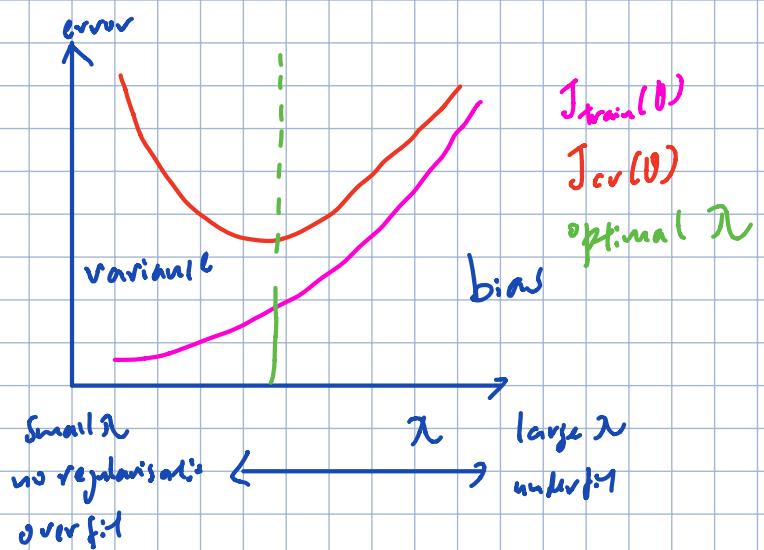
$\theta^{(12)} \rightarrow J_{\text{cv}}(\theta^{(12)})$

Say  $\theta^{(5)}$  has the lowest  $J_{\text{cv}}(\theta^{(5)})$ , we can pick that

Let's compute  $J_{\text{test}}(\theta^{(5)})$  to check generalisation error.

Bias / variance as a function of the regularisation parameter  $\lambda$

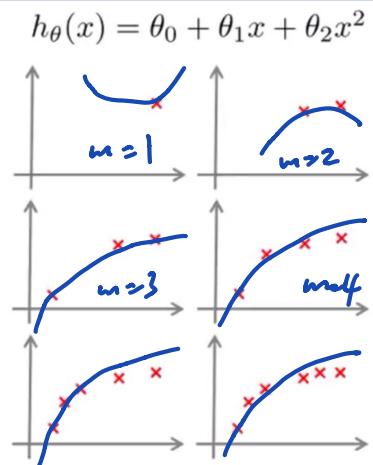
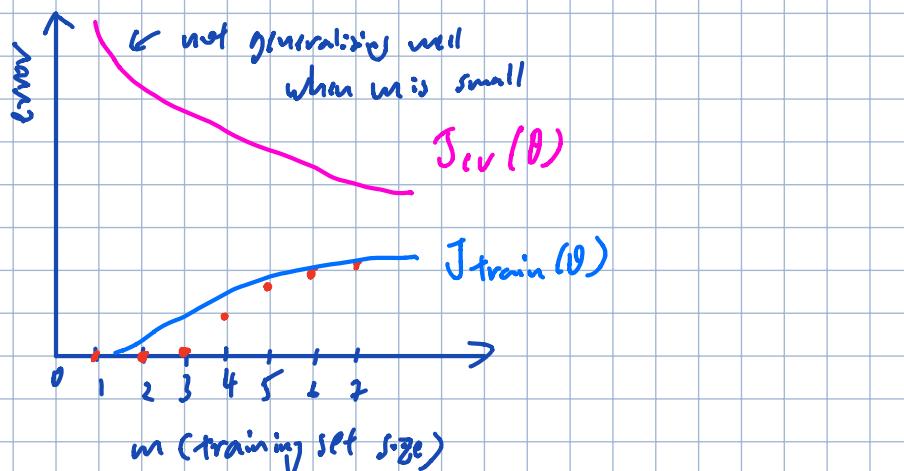
$J_{\text{train}}(\theta) / J_{\text{cv}}(\theta)$  has no regularisation term for this section



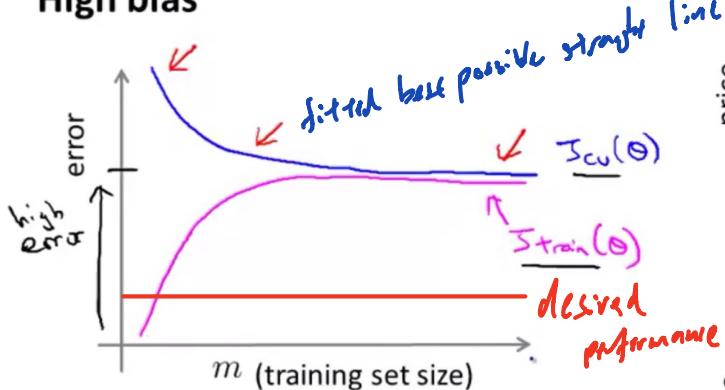
## Learning Curves

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

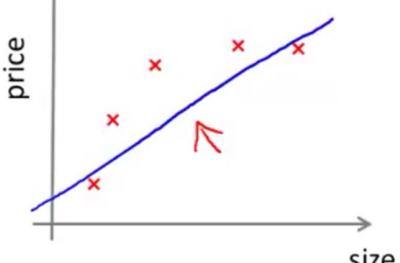
$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$



## High bias

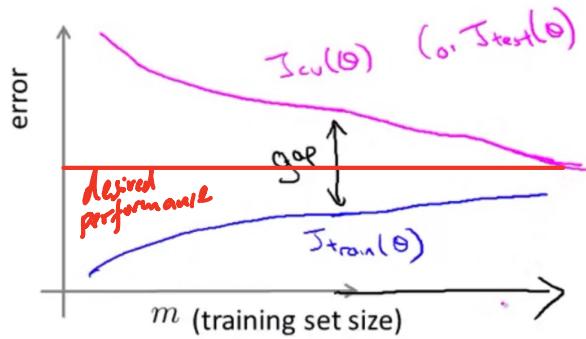


$$h_{\theta}(x) = \underline{\theta_0 + \theta_1 x}$$



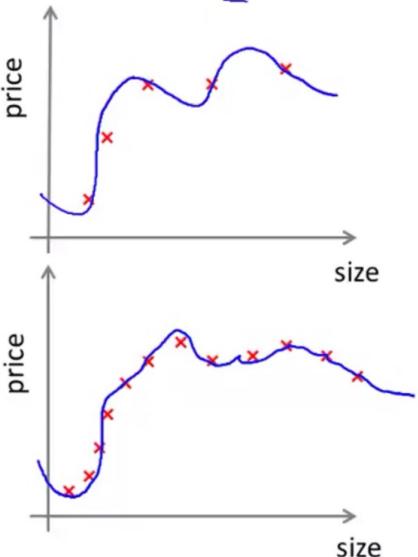
If a learning algorithm is suffering from high bias, getting more training data will not (by itself) help much.

## High variance



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

(and small  $\lambda$ )



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

## Debugging a learning algorithm:

Suppose you have implemented regularized linear regression to predict housing prices. However, when you test your hypothesis in a new set of houses, you find that it makes unacceptably large errors in its prediction. What should you try next?

- Get more training examples → fix high variance
- Try smaller sets of features → fix high variance
- Try getting additional features → fix high bias, current model too simple
- Try adding polynomial features ( $x_1^2, x_2^2, x_1 x_2$ , etc) → fix high bias
- Try decreasing  $\lambda$  → fix high bias
- Try increasing  $\lambda$  → fix high variance

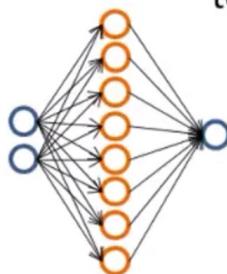
## Neural networks and overfitting

"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 ( $\lambda$ ) to address overfitting.

## Building a spam classifier

supervised learning.

$X$  = features of email     $y$  = Spam(1) or ham(0)

Feature  $x_i$ : choose 100 words indicative of spam/ham.

e.g. deal, buy, discount...

Given an email, we can construct a feature vector

$$X = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ \vdots \end{bmatrix} \begin{array}{l} \text{buy} \\ \text{deal} \\ \vdots \\ \text{now} \\ \vdots \end{array}$$

$$x_j = \begin{cases} 1 & \text{if word } j \text{ appears in email} \\ 0 & \text{otherwise} \end{cases}$$

$$x \in \mathbb{R}^{100}$$

In practice, take most frequently occurring words (10k - 50k)  
in training set, rather than picking 100 words manually.

## Building a spam classifier

How to spend our time to make it have low error

→ collect lots of data

→ develop sophisticated features based on email routing information (from email header)

→ develop sophisticated feature for message body

→ detect misspelling

## Recommended approach

→ Start with a simple algorithm that we can implement quickly

Implement it and test it on our cross-validation data

→ Plot learning curves to decide if more data, more features, etc are likely to help

→ Error analysis: Manually examine the examples (in cross-validation set) that our algorithm made errors on. See if we spot any systematic trend in what type of examples it is making errors on.

e.g.  $M_{cv} = 500$  examples in CV set

Algorithm misclassified 100 emails

↳ manually examine 100 errors, and categorise them based on:

(i) what type of email it is

(ii) what cues (features) we think would have helped the algorithm classify them correctly.

## Error metrics for skewed classes

Cancer classification example:

train logistic regression model  $h_\theta(x)$

$$Y = \begin{cases} 1 & \text{cancer} \\ 0 & \text{otherwise} \end{cases}$$

find that you get 1% error on test set, 99% correct diagnosis

⇒ Only 0.5% of patients have cancer.

⇒ Error percentage is big.

### Precision / Recall

✗  $\gamma=1$  in presence of rare class that we want to detect

		Actual class	
		1	0
predict class	1	TP	FP
	0	FN	TN

### Precision

of all patients where we predict  $y=1$ , what fraction actually has cancer?

$$\frac{\text{TP}}{\text{A predicted positive}} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

### Recall

of all patients that actually have cancer, what fraction did we correctly detect as having cancer?

want high precision

high recall

$$\frac{\text{TP}}{\text{Actual positive}} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

## Trading off precision and recall

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

predict 1 if  $h_\theta(x) \geq 0.5 \rightarrow 0.7$  more confident  $\rightarrow 0.3$   
predict 0 if  $h_\theta(x) < 0.5 \rightarrow 0.7$  less confident  $\rightarrow 0.3$

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

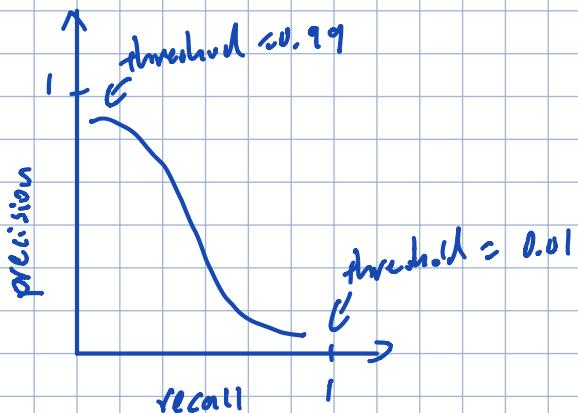
$\rightarrow$  higher precision

$\rightarrow$  lower recall ( $\because$  predict  $y=1$  less often)

$\Rightarrow$  Suppose we want to avoid missing too many cases of cancer  
( $\therefore$  avoid false negative)

$\rightarrow$  low precision

$\rightarrow$  higher recall



## F<sub>1</sub> score (F score)

How to compare precision/recall numbers

	Precision (P)	Recall (R)	Average	F <sub>1</sub> score
Algorithm 1	0.5	0.4	0.45	0.4444
Algorithm 2	0.7	0.1	0.4	0.175
Algorithm 3	0.02	1.0	0.51	0.0392

A single metric to decide (instead of P and K)

$$\text{average : } \frac{P+R}{2} \quad \times$$

Not very good classifier which predicts  $y=1$  all the time, will get very high recall. (Algorithm 3)

classifier that predicts  $y=0$  all the time, (by setting threshold very high),  $\Rightarrow$  very high precision and very low recall

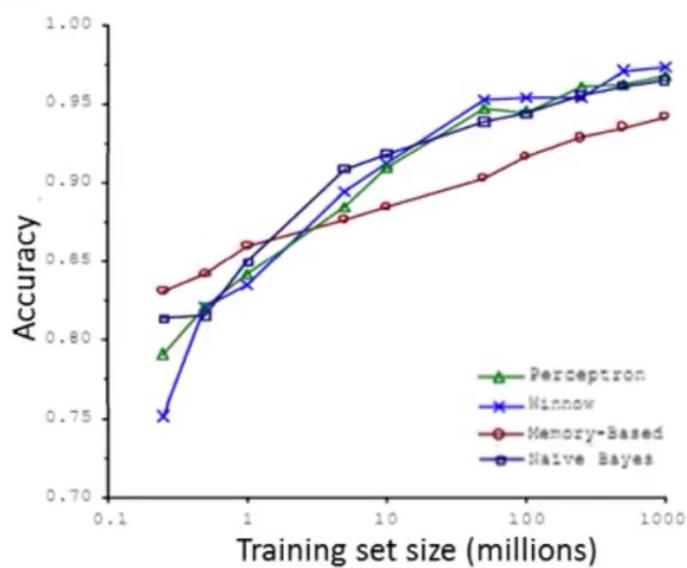
$$F_1 \text{ score} : 2 \frac{P R}{P + R}$$

(good if both  $R$  and  $R'$  are high, punishes if one of them are very low)

\* try a range of thresholds, compute F<sub>1</sub> score on the CV set, pick the threshold that gives the maximum F<sub>1</sub> score

\* Assume feature  $\pi \in \mathcal{R}^{n+1}$  has enough info to predict y accurately

# Parte far ML



## Large data rationale

use a learning algorithm with many parameters (e.g. logistic regression with many features; neural network with hidden units)

## low bias algorithm

use a very large training set (unlikely to overfit) low variance

$$\rightarrow J_{train}(\theta) \approx J_{test}(\theta)$$

hope fully

$J_{\text{try}}(\theta)$  will be small