

# Coursera Machine Learning

## [Welcome lecture]

- ▶ ML as a field grew out of AI, when it became clear that for some problems you can't write a program by hand. (e.g. autonomous helicopters)
- ▶ Examples of everyday ways you use ML
  - search ranking
  - automatic photo tagging
  - spam filtering
  - database mining (e.g. web click data)
  - self-customising programs

## What is Machine Learning?

Def'n according to Arthur Samuel:  
Field of study that gives computers the ability to learn without being explicitly programmed.

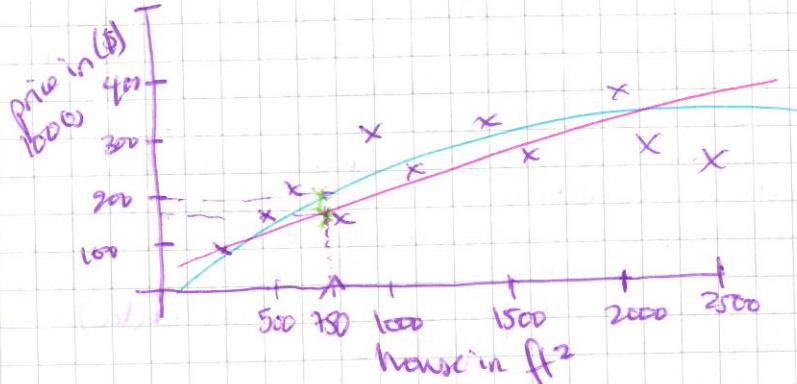
Later def'n:

A well Posed Learning Problem: A computer program is said to learn from experience E wrt some task T and some performance measure P if its performance on T, as measured by P, improves with experience E.

# What is Supervised Learning?

Let's start with an example

e.g. Housing Price Prediction



If given data on how much houses (by size) have sold for in the past, can we predict how much a 750 sq ft house will sell for?

This problem is an example of supervised learning because there are some "right answers" to work with (i.e. existing data)

This is also known as a regression problem because it is predicting continuous valued output (the prices)

What about a ~~more~~ different example?

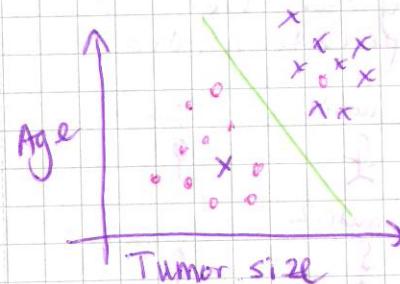
# Breast Cancer Classification



Say you have some past patient data (their tumor size and whether benign or malignant). Now you have a new patient & corresponding tumor size. Can we predict if benign or malignant?

This is an example of classification (predicting discrete valued output)

But many real world problems have many features. We may still want to use classification.



In breast cancer research, also interested in  
- clump thickness  
- uniformity of cell size  
- uniformity of cell shape  
in etc.

Support Vector Machines allow an infinite number of features.

## What is unsupervised learning?

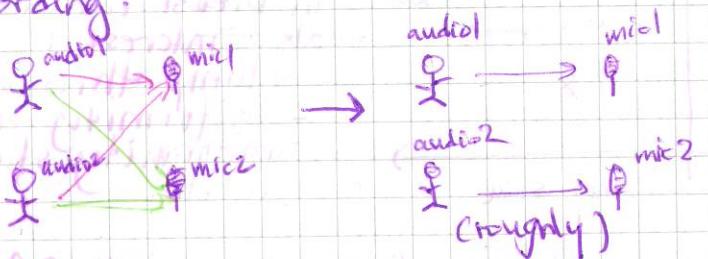
- No historical data with right/wrong classification
- Given unlabelled data, algorithm may seek to find structure (e.g. clustering)
- Every day example: Google News & other news aggregators

~~clustering example~~  
of DNA microarrays  
• social network analysis  
• market segmentation  
• astronomical data analysis  
→ we don't know in advance what the segments are!

Research example:

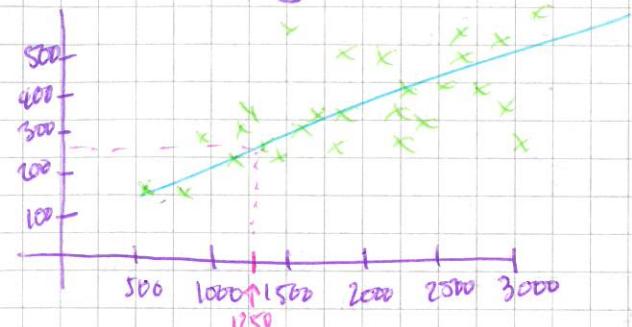
### Cocktail Party Problem

Algorithm that tries to find structure in overlapping sounds/voices  
~~choose~~ so that they can be separated out into separate audio recording.



## Model representation: Linear Regression

Back to housing prices in Portland...



More formally, we have a training set

Notation:

$m$  = number of training examples  
 $x^{(i)}$  = "input" variable/features  
 $y^{(i)}$  = "output" variable/features

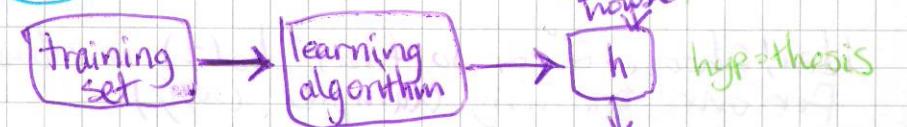
size in $\text{ft}^2$ ( $x^{(i)}$ )	price in 1000s ( $y^{(i)}$ )
2104	460
1416	232
1534	315
802	178
...	...
$x^{(1)}$	$y^{(1)}$
$x^{(2)}$	$y^{(2)}$
$x^{(m)}$	$y^{(m)}$

$m = 4$

$(x, y)$  - one training example

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

Model



$h$  maps from  $x^{(i)}$  to  $y^{(i)}$

How do we represent  $h$ ?

Here, for regression

$$h(x) = \theta_0 + \theta_1 x$$



$$h(x) = \theta_0 + \theta_1 x$$

Later, we will use more complex models!

This one is linear regression with one variable aka univariate linear regression.

How to choose our 2 parameters?

Let's try a few values...

$$\begin{aligned} h(x) &= 1.5 + 0x \\ \theta_0 &= 1.5 \\ \theta_1 &= 0 \end{aligned}$$

$$\begin{aligned} h(x) &= 0 + 0.5x \\ \theta_0 &= 0 \\ \theta_1 &= 0.5 \end{aligned}$$

$$\begin{aligned} h(x) &= 1 + 0.5x \\ \theta_0 &= 1 \\ \theta_1 &= 0.5 \end{aligned}$$

Idea! Choose  $\theta_0, \theta_1$  so that  $h(x)$  is close to  $y$  for our training examples  $(x_i, y_i)$

More formally:

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

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

As a cost function:

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

$$\text{minimize}_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

$J(\theta_0, \theta_1)$  is aka squared error function

Recap: we want to fit a straight line through our data and

Hypothesis:  $h(x) = \theta_0 + \theta_1 x$

Parameters:  $\theta_0, \theta_1$

Cost function:  ~~$J(\theta_0, \theta_1) = (h_\theta(x^{(i)}) - y^{(i)})^2$~~

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

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

Let's consider a simplified version

$$h_\theta(x) = \theta_1 x \quad \theta_0 = 0$$

parameter:  $\theta_1$

$$\min_{\theta_1} J(\theta_1)$$

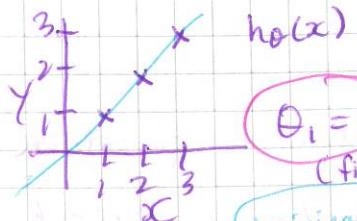
only consider hypothesis functions that pass through the origin



Two functions to understand

$$h_{\theta}(x)$$

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

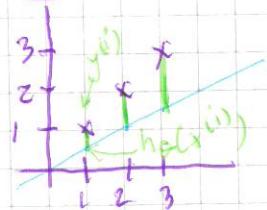


when  $\theta_1 = 1$

$$\begin{aligned} J(\theta_1) &= \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2 \\ &= \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0 \end{aligned}$$

so  $J(1) = 0$

$$\theta_1 = 0.5$$



$$J(\theta_1) = J(0.5)$$

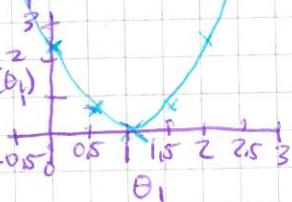
$$= \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2$$

$$= \frac{1}{2m} (0.5^2 + 1^2 + 1.5^2)$$

$$= \frac{3.5}{6} \approx 0.58$$

$$J(\theta_1)$$

(function of the parameter  $\theta_1$ )



Each value of  $\theta_1$  corresponds to a different hypothesis function

We want to pick the  $\theta_1$  where

$$\begin{aligned} &\text{minimize } J(\theta_1) \\ &\theta_1 \end{aligned}$$

## Gradient Descent

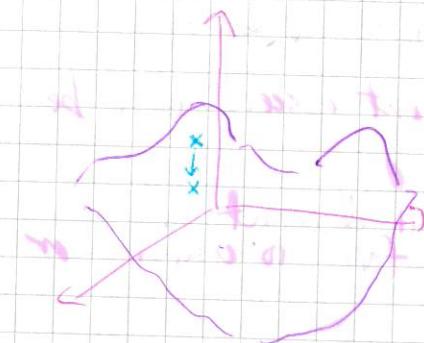
Goal:

We have some function  $J(\theta_0, \theta_1)$

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

Outline:

- start with some  $\theta_0, \theta_1$  (say,  $\theta_0 = 0, \theta_1 = 0$ ) but it doesn't really matter
- keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$  until we hopefully end up at a minimum



Ideas: Imagine you're in a grassy park where you are standing on a hill and want to decrease your elevation ASAP. So you look around 360° and decide where to take your next step. Repeat.

Gradient descent does not guarantee global minimum. Only a local min!

Gradient descent algorithm

repeat until convergence {

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

learning rate  
= size of step

means update both  $\theta_0$  and  $\theta_1$  simultaneously

Correct implementation  
In impl: simultaneous update

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

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

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

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

Picking  $\alpha$ :

• If  $\alpha$  is too small, gradient descent can be slow

• If  $\alpha$  is too large, gradient descent can overshoot the minimum, fail to converge or even diverge.

## Gradient Descent For Linear Regression

To apply to our linear regression model we need to know:

$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{\partial}{\partial \theta_j} \cdot \frac{1}{m} \sum_{i=1}^m (\text{h}_0(x^{(i)}) - y^{(i)})^2 \\ &= \frac{\partial}{\partial \theta_j} \cdot \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2\end{aligned}$$

Need to know for

$$\theta_0: j=0: \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (\text{h}_0(x^{(i)}) - y^{(i)})$$

$$\theta_1: j=1: \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (\text{h}_0(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

Now we have gradient descent for linear regression

repeat until convergence?

$$\theta_0 := \theta_0 - \alpha \left( \frac{1}{m} \sum_{i=1}^m (\text{h}_0(x^{(i)}) - y^{(i)}) \right) \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\theta_1 := \theta_1 - \alpha \left( \frac{1}{m} \sum_{i=1}^m (\text{h}_0(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \right) \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

E

Note: over iteration

Though gradient descent generally doesn't guarantee a global optimum, linear regression always generates a convex (bowl-shaped) function and thus there will always be a global optimum

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

$$w_0 + w_1 x_0 + w_2 x_1 + \dots + w_n x_n = (x) \text{ dot } (w)$$

$$p(x) = p(x) \cdot \exp. + \alpha x_0 + \beta x_1 + \gamma x_2 = (x) \text{ dot } (p)$$

## Linear Regression with Multiple Features

Let's return to our question about housing price prediction.

Size (ft <sup>2</sup> )	# of bedrooms	# of floors	Age of home (yr)	Price (in \$1000)	y
2104	5	1	45	460	
1416	3	2	40	232	
1534	3	2	30	315	
852	2	1	36	178	
...	...	...	...	...	m=4

More notation:

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

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

How do multiple features affect the hypothesis?

Previously:  $h_{\theta}(x) = \theta_0 + \theta_1 x_1$

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

e.g.  $h_{\theta}(x) = 80 + 0.1x_1 + 0.9x_2 + 3x_3 - 2x_4$

More generally:

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

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n \quad (\text{for convenience of notation})$$

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

We can use the fancy term multivariate linear regression.

## 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:  $\theta_0, \theta_1, \dots, \theta_n / \theta$  where  $\theta$  is a  $n+1$  dimensional vector

Cost function:

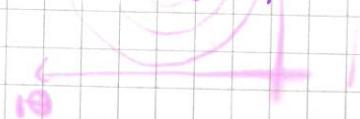
$$J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^{m-1} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

## Gradient Descent

Repeat until convergence {

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

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



## Now Gradient Descent Algorithm ( $n \geq 1$ )

Repeat {

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

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

}

So..

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

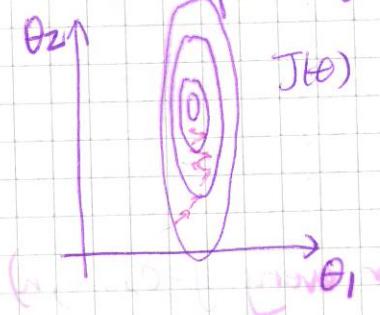
$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (\text{h}_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

## Feature Scaling

Idea: If you make sure features are on a similar scale, gradient descent will converge faster.

e.g.  $x_1 = \text{size (0-2000 ft}^2)$   
 $x_2 = \text{number of bedrooms (1-5)}$

Contour plot might look like



$x_1 = \frac{\text{size (ft}^2)}{2000}$   
 $x_2 = \# \text{ of bedrooms}$

$\theta_2 \uparrow$   
 $J(\theta)$

$\theta_1 \uparrow$   
 $J(\theta)$

so then  
 $0 \leq x_1 \leq 1$   
 $0 \leq x_2 \leq 1$

$\theta_2 \uparrow$   
 $J(\theta)$

$\theta_1 \uparrow$   
 $J(\theta)$

## Goal of feature scaling:

Get every feature into approximately a  $-1 \leq x_i \leq 1$  range but this is arbitrary

$0 \leq x_1 \leq 3$  ✓ different rules of thumb,  
 $-2 \leq x_2 \leq 0.5$  ✓ but  $3 \leq x_3 \leq -3$   
 $-100 \leq x_3 \leq 100$  ✗ def. reasonable  
 $-0.0001 \leq x_4 \leq 0.0001$  ✗

## Mean normalization

Replace  $x_i$  with  $x_i - \bar{x}_i$  to make features approx zero mean

e.g.  $x_1 = \frac{\text{size} - 1000}{2000}$   $x_2 = \frac{\# \text{ bedrooms} - 2}{5}$

$-0.5 \leq x_1 \leq 0.5$   $-0.5 \leq x_2 \leq 0.5$

Generally:

$x_i \leftarrow \frac{x_i - \bar{x}_i}{S_i}$   $\bar{x}_i$  avg value of  $x_i$  in training set

$S_i$  range (max-min) or standard deviation

$\bar{x}_i$  mean for  $x_i$

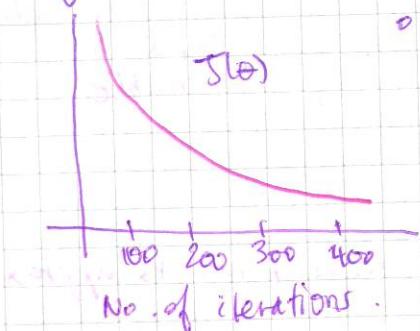
! multiply two  $x_i$  to avoid trouble

! adding  $x_i$  to  $J(\theta)$  to minimize a well minimum with boundaries at previous

## Working Rate

Making sure gradient descent is working

GOOD  
 $\min J(\theta)$



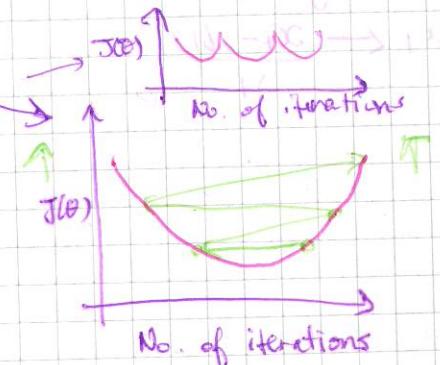
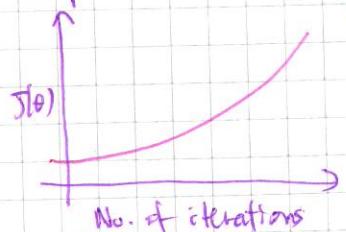
- value of cost function after each iteration of gradient descent

- $J(\theta)$  should decrease after every iteration.
- No. of iterations can vary depending on application but this trend should hold.

Example of automatic convergence test:

Declare convergence if  $J(\theta)$  decreases by less than some  $\epsilon$  (e.g.  $10^{-3}$ ) in one iteration

If plot looks like:



Gradient descent is not working!

Use a smaller  $\alpha$ . Choice of  $\alpha$  is probably causing us to overshoot the minimum.

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

- But if  $\alpha$  is too small,  $\alpha$  vs gradient descent will be slow.

To choose  $\alpha$ , try

...,  $0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, \dots$

until you find which in the range is too large and which is too small and then try to pick largest value possible that will cause convergence in that range.

## Features and Polynomial Regression

Housing Prices Prediction

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

$$+ \theta_2 \times \text{depth}$$

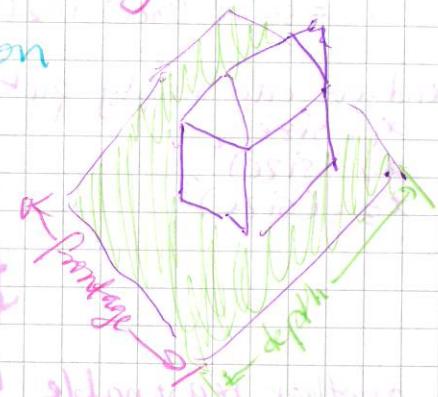
May want to instead define new feature

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

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

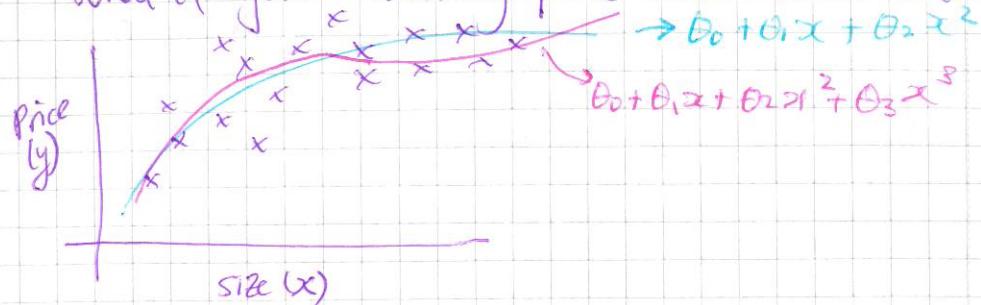
~ land area.

No need to stick with the features you started with!



## Polynomial Regression

What if your housing price data looks like



It may be better fit by a polynomial.

We're familiar with hypotheses of the form

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\ = \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2 + \theta_3(\text{size})^3$$

and we can simply pick

$$\begin{aligned} x_1 &= \text{size} \\ x_2 &= (\text{size})^2 \\ x_3 &= (\text{size})^3 \end{aligned}$$

size:  $1 - 1000$   
 $(\text{size})^2: 1 - 100,000$   
 $(\text{size})^3: 1 - 10^9$

feature scaling even more important in this case!

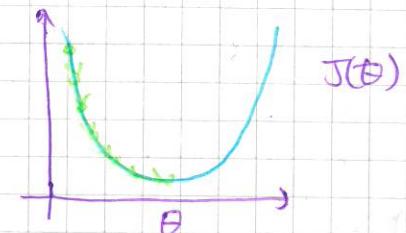
another reasonable hypothesis might be

$$h_0(x) + \theta_0 + \theta_1(\text{size}) + \theta_2 \sqrt{\text{size}}$$

Later, we'll cover algorithms that will automatically decide which features to use from old two

## Normal Equation

Gradient descent



Normal equation: Method to solve for theta analytically.

Some intuition: If 1D ( $\theta \in \mathbb{R}$ )

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

$$\frac{d}{d\theta} J(\theta) = \dots \stackrel{\text{set}}{=} 0$$

Solve for  $\theta$

But in our case,  $\theta$  is not simply a real number.  $\theta \in \mathbb{R}^{n+1}$

$$J(\theta_0, \theta_1, \dots, \theta_m) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots \stackrel{\text{set}}{=} 0 \quad (\text{for every } j)$$

Solve for  $\theta_0, \theta_1, \dots, \theta_n$

But derivation is somewhat involved:

$$\left[ \begin{array}{c} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{array} \right] = \left[ \begin{array}{c} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{array} \right]$$

$$T^{(m)} X =$$

Examples m=4

$x_0$	$x_1$	$x_2$	$x_3$	$x_4$	$y$
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \quad y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

dim: m x (n+1)

m (vector)

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

More generally,

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

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \quad X = \begin{bmatrix} \xrightarrow{(x^{(1)})^T} & & \\ \xrightarrow{(x^{(2)})^T} & \ddots & \\ \vdots & & \\ \xrightarrow{(x^{(m)})^T} & & \end{bmatrix}$$

m x (n+1) dim feature matrix.

simple example: If  $x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \end{bmatrix}$

$$X = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ \vdots & \vdots \\ 1 & x_1^{(m)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

dim: m x 2

Implementing  $\Theta = (X^T X)^{-1} X^T y$ :

$$A = X^T X$$

$$(X^T X)^{-1} = A^{-1}$$

in octave: pinv( $X^T X$ ) \*  $X^T y$

If using normal equation, feature scaling is not necessary.

Comparing to gradient descent:

Gradient descent

- Need to choose  $\alpha$
- Needs many iterations
- Works well even when  $n$  is large

Normal Equation

- No need to choose  $\alpha$
- Don't need to iterate
- Need to compute  $(X^T X)^{-1}$  ( $n \times n$  matrix)
- Slow if  $n$  is very large! Costs  $\Theta(n^3)$

$$n=100$$

$$n=1000$$

$$n=10,000$$

(maybe).

# Normal Equation Noninvertability

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

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

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

if you use  $\text{pinv}(X^T X)$ , you will still

get the right answer

How a non-invertible  $X^T X$  comes about:

- Redundant features (linearly dependent)

e.g.  $x_1 = \text{size in feet}^2$   
 $x_2 = \text{size in m}^2$

$$x_1 = (3, 28)^2 x_2$$

- When you have linearly dependent features, your matrix will be non-invertible.

- Too many features ( $m \leq n$ )

- Delete some features, or use regularization

end of week 1

end of week 2

gradient descent

stochastic gradient descent

(vectorized)  $(X^T X)$

(matrix inversion)

( $X^T X$  is also  $\log(\cdot)$ )

# Octave Tutorial

- size(A) gives  $1 \times 2$ , [rows #cols]

length(v) gives longest dimension

- load data into octave

e.g.  $\text{load('featuresX.dat')}$

$\text{load('featuresX.dat', 1)}$

- who gives variables in current scope.  
Who is more detailed.

- clear removes variable from scope

- save variables into file

e.g.  $\text{save('hello.mat')}$

$\text{save('hello.txt')}$  -ascii

- append column  
human-readable

e.g.  $A = [A, \{ \text{COL VECTOR} \}]$

flatten:  $A(:)$

Multiplication:

$A * B \rightarrow$  matrix multiplication

$A . * B \rightarrow$  element wise multiplication

$\text{diag}, \text{full}, \text{prn}$

$\text{triu}, \text{tril}$

$\text{conj}, \text{real}, \text{imag}$

max [val, ind] = max(A)

returns value that is max & its index  
works for 2D arrays

sum sum(a) sums all elements of vector  
prod(a) as expected

sum(A, 1); sum of each column

sum(A, 2); sum of each row.

pinv pinv(A) gives inverse of A.

Plotting: many different options available

e.g. t = [0:0.1:0.98]

y1 = sin(2 \* pi \* 4 \* t) sum

1000 - v(t+0.1) sum

y2 = cos(2 \* pi \* 4 \* t)

[plot both on same graph, one in red]

plot(t, y1);

hold on

plot(t, y2, 'r');

xlabel('time');

ylabel('value');

legend('sin', 'cos');

title('my plot');

Save: print '-dpng' 'myplot.png'.

hold off  
close (close figure)

can specify figure numbers  
figure(1) ...

subplot (1, 2, 1) % divides plot into  
1x2 grid &  
across first elem

change scale of axes with axis  
eg. axis([0, 5, -1, 1])

clf clear figure

imagesc(A) plots A matrix as a variation  
of colour depending on value

grayscale heatmap:

imagesc(A), colorbar, colormap gray;

"comma chaining commands"  
eg a=1, a=2, a=3 will execute all 3,  
a=1; a=2; a=3

for i=1:10

while i<=5

if v(i) == 1

all need an  
end.

## functions

create a new file e.g. squareThis.m  
takes one param

function  $y = \text{squareThisNumber}(x)$

$$y = x^2$$

will return one value  
saved in y

to access the function, it needs to be in working directory.

to add directory to Octave's search path,  
use searchpath

to return multiple arguments, e.g.-

function [y1, y2] squareAndCube(x)

$$y1 = x^2$$

$$y2 = x^3$$

to call:  
[a, b] = squareAndCube(x)

## Vectorization

e.g.

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

$$= \theta^T x$$

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

## unvectorised implementation

$$\text{prediction} = 0.0$$

for  $j = 1 : n+1$ ,  
 $\text{prediction} = \text{prediction} + \theta_j * x(j)$   
end;

## vectorized implementation

$$\text{prediction} = \theta^T * x$$

## e.g. gradient descent

$$\theta_i := \theta_i - \alpha \delta$$

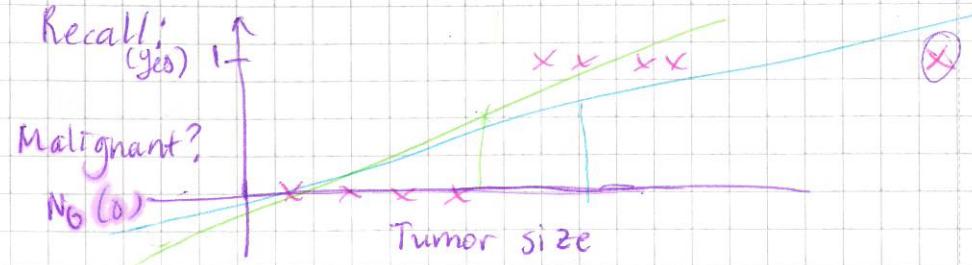
where  $\delta = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$

$$\delta = \begin{bmatrix} \delta_0 \\ \delta_1 \\ \delta_2 \end{bmatrix} \quad \delta_0 = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

# Classification: Logistic Regression

onto the discrete values!

For now,  $y$  can take on 2 values  
 $\rightarrow 0$ : "Negative class" (eg. benign tumor)  
 $y \in \{0, 1\}$   
 $\rightarrow 1$ : "Positive class" (eg. malignant tumor)



We can't use linear regression as data points like one far right are not outliers but still influence fit

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

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

## Hypothesis representation

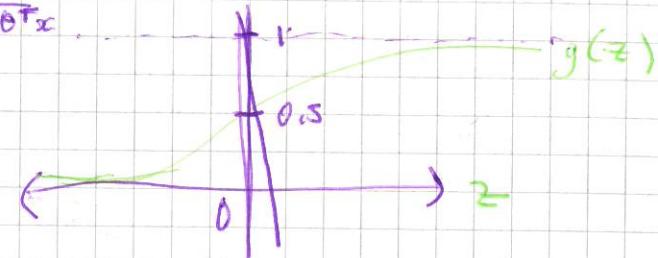
### Logistic Regression Model

Want  $0 \leq h_\theta(x) \leq 1$

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

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

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



As before, will eventually need to fit parameters  $\theta$

## Interpretation of hypothesis output

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

Example: If  $x = [x_0 \ x_1] = [1 \ \text{tumorsize}]$

$$h_\theta(x) = 0.7$$

Tell patient 70% chance of tumor being malignant

Or more formally,  $h_\theta(x) = p(y=1|x; \theta)$

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

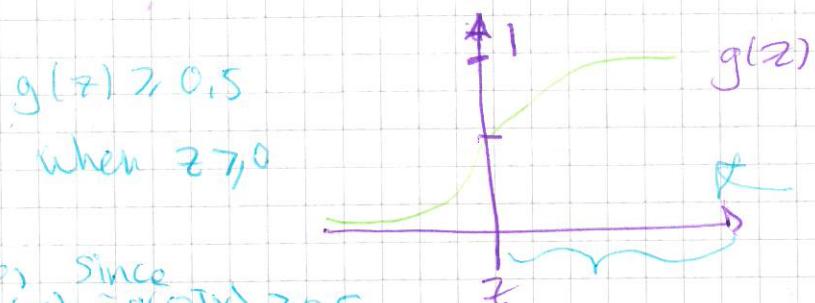
We can compute the probability for  $y=0$  or  $1$ :

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

### The Decision Boundary

but when to pick  $y=0$ ,  $y=1$

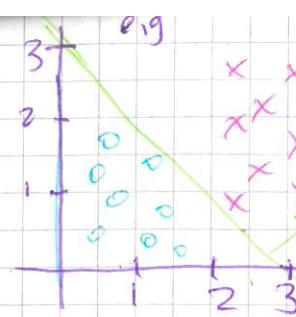
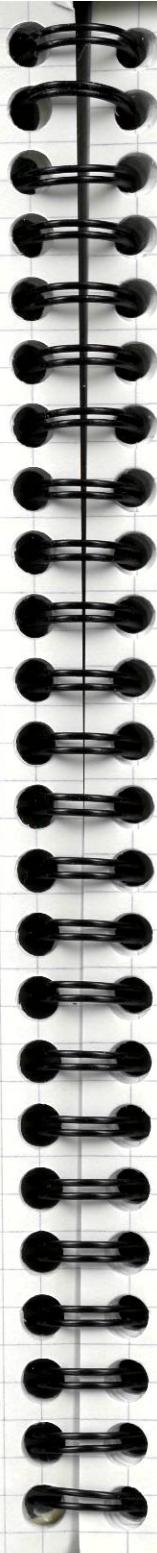
Suppose predict " $y=1$ " if  $h_\theta(x) \geq 0.5$   
 " $y=0$ " if  $h_\theta(x) < 0.5$



So, since  $h_\theta(x) = g(\theta^T x) > 0.5$

whenever  $\theta^T x > 0$

Likewise  $y=0$  predicted when  $\theta^T x < 0$



$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

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

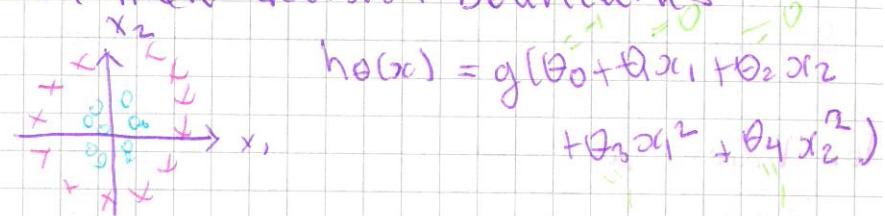
Decision boundary where  $h_\theta(x) = 0$

When would this  $\theta$  predict  $y=1$ ?

$$\text{If } \frac{-3 + \theta_1 x_1 + \theta_2 x_2}{\theta^T x} \geq 0$$

$$x_1 + x_2 \geq 3$$

### Non-linear decision boundaries



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

## Logistic Regression: Cost Function

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

with m examples

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad x_0 = 1 \quad y = \{0, 1\}$$

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

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

so how do we choose parameters  $\theta$

For linear regression

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

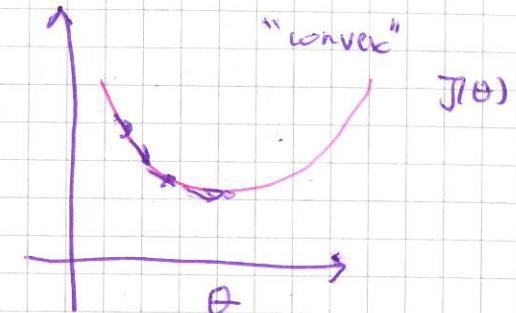
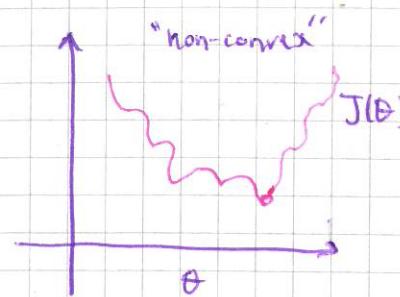
Can change notation

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

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

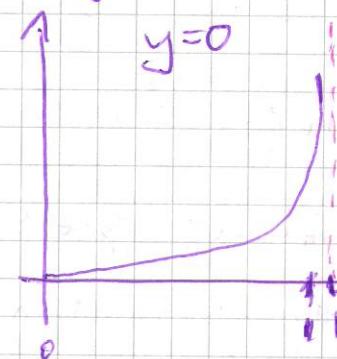
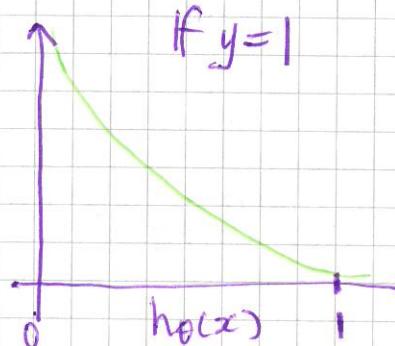
called the cost because that's the price paid by the algorithm (not error) for outputting this value

This can't be directly applied to logistic regression because then  $J(\theta)$  is not convex and it's more complicated to find a global minimum



## Logistic Regression Cost Function

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



recall  $\log z$ :



only care about range between 0 and 1 as  $h_{\theta}(x) \rightarrow 0$  and 1 as  $h_{\theta}(x) \rightarrow \infty$

$\text{cost} = 0$  if  $y=1, h_{\theta}(x)=1$  as  $h_{\theta}(x) \rightarrow 0$   $\text{cost} \rightarrow \infty$   
 if  $h_{\theta}(x)=0$  but  $y=1$  we'll penalize learning alg.

## Logistic Regression: Simplified Cost Function

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

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

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

Can rewrite as

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

## So, Logistic Regression Cost Function

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

Why this cost function  
• convex!

• determined through maximum likelihood estimation

To fit parameters  $\theta$ :

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

To make a prediction given new  $x$ :

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

So how do we get  $\min_{\theta} J(\theta)$ ?

Gradient Descent!

We want  $\min_{\theta} J(\theta)$ :

Repeat {

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

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

} (simultaneously update all  $\theta_j$ )

Algorithm looks identical to linear regression!

The difference is in  $h_\theta(x)$ .

## Logistic Regression: Optimisation

### Optimisation Algorithm

We have cost function  $J(\theta)$ .

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

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

$$\begin{bmatrix} -J(\theta) \\ -\frac{\partial}{\partial \theta_j} J(\theta) \end{bmatrix} \leftarrow \text{for } (j=0, 1, \dots, n)$$

which we then use for gradient descent

Repeat:

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

?

But there are algorithms  
e.g. - Conjugate gradient

- BFGS

- L-BFGS

Advantages

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

Disadvantages

- more complex

But don't try to implement these algorithms yourself! Beyond the scope of this course.  
It's worth investigating and finding a good library.

Example

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

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

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

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

Would implement a cost function:

function [JVal, gradient] = costFunction(theta)

$$\begin{aligned} JVal &= \theta(1) - 5^2 + \dots \\ &\quad + \theta(2) - 5^2 \end{aligned}$$

$$\begin{aligned} \text{gradient}(1) &= 2 * \theta(1) - 5 \\ \text{gradient}(2) &= 2 * \theta(2) - 5 \end{aligned}$$

{code to  
compute  
 $\frac{\partial}{\partial \theta_1} J(\theta)$

{code to  
compute  
 $\frac{\partial}{\partial \theta_2} J(\theta)$

Which then leaves you with what you  
need to use matlab library

options = optimset('gradObj', 'on', 'maxIter', 100);  
initialTheta = zeros(2, 1);

[optTheta, functionVal, exitFlag] ...  
ideally to zero,

= fminunc(@costFunction, initialTheta, options)

function minimization  
unconstrained

pointer to  
costfunction

## Multiclass Classification

examples:

- o autotagging emails
  - work  $y=0$
  - friends  $y=1$
  - family  $y=2$
  - hobby  $y=3$

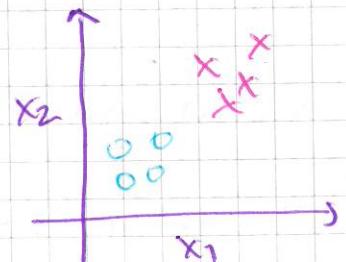
- o medical diagnosis

- not ill
- cold
- flu

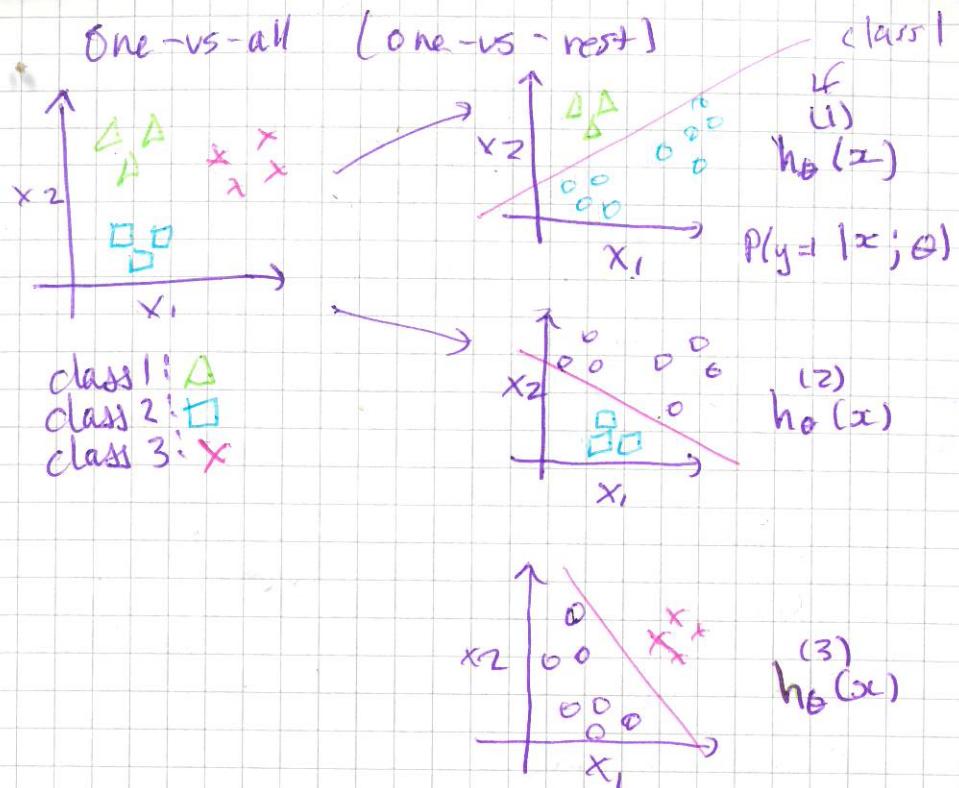
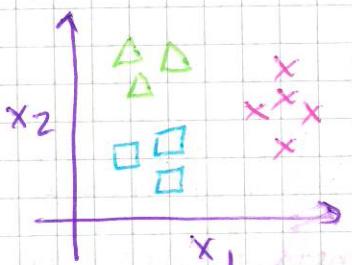
- o weather

- sunny
- cloudy
- rain
- snow

Previously, binary classification



Now, multi-class



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

Approach:

- Train a logistic regression classifier  $h_\theta(x)$  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~~ maximises

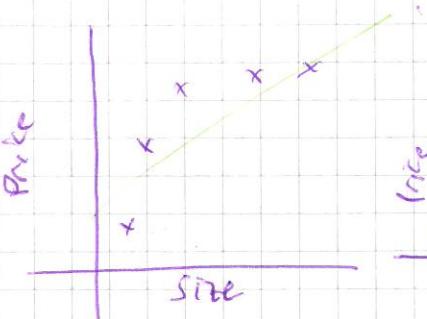
$$\max_i h_\theta^{(i)}(x)$$

(check all 3, pick the one that is most confident)

# The Problem of Overfitting

Regularization

What is overfitting?

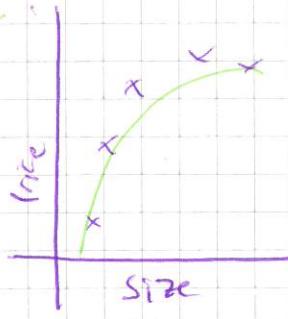


$$\hat{y} = \theta_0 + \theta_1 x$$

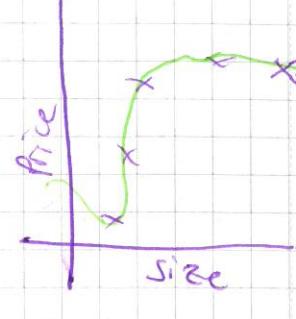
"underfit"

"high bias"

(With already assumed data is linear)



$$\hat{y} = \theta_0 + \theta_1 x + \theta_2 x^2$$



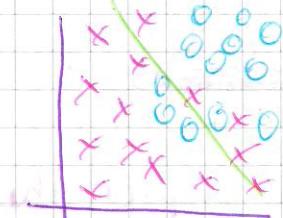
$$\hat{y} = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

"overfit"

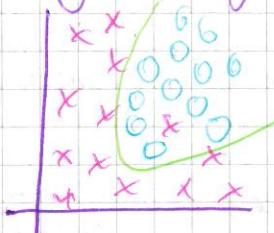
"high variance"

Overfitting: If we have too many features, the learned hypothesis may fit the training set very well, but fail to generalise to new examples

Can overfit for logistic regression as well



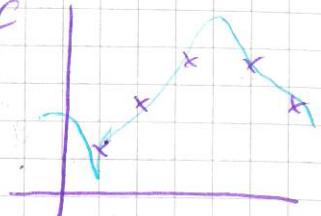
underfit



overfit

How to address overfitting?

- can only plot hypothesis if you have 1-2 features



Options:

1. Reduce number of features

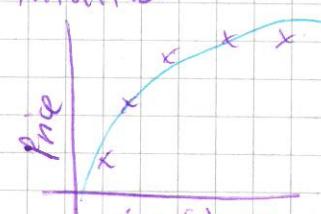
- Manually select which features to keep
- Model selection algorithm (later in course)

2. Regularization

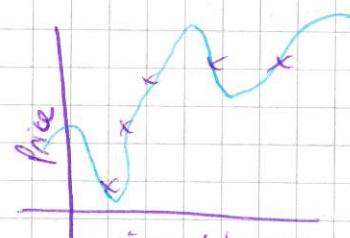
- keep all the features, but reduce magnitude/values of parameters  $\theta_j$

- Works well when we have a lot of features, each of which contributes a bit to predicting  $y$ .

Intuition



$$\hat{y} = \theta_0 + \theta_1 x + \theta_2 x^2$$



$$\hat{y} = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

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

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

then  $\theta_3 \approx 0$   $\theta_4 \approx 0$

## Idea behind Regularisation

If you have small values for parameters  $\theta_0, \theta_1, \dots, \theta_n$

- "simpler" hypothesis

- less prone to overfitting

e.g. housing

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

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

harder to pick which parameters are most relevant!

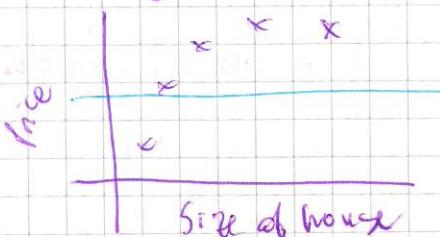
So, we'll modify our cost function

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

But there's a trick to picking  $\lambda$ !

Setting  $\lambda$  to an extremely large value makes  $\hat{h}_\theta$  flat that  $\theta_1, \theta_2, \theta_3, \dots \approx 0$

and suggests  $\hat{h}_\theta(x) = \theta_0$ , a flat line!



underfit!

## Regularised Linear Regression

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

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

### Now: Gradient Descent

Repeat {

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

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

}

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

Can rewrite  $\theta_j$  update as

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

## Normal Equation

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

o Regularisation also takes care of non-invertibility issue.

## Regularised Logistic Regression

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, same as for linear regression but remember  $h_\theta(x)$  has changed!

## Regularised Logistic Regression with more Advanced Optimisation

Recall: you need to provide a cost function to fitline

function [JVal, gradient] = costFunction(theta)

JVal = code to compute  $J(\theta)$

$J(\theta)$  = new cost function w/  $\lambda$  term

gradient(1) = code to compute  $\frac{\delta}{\delta \theta_0} J(\theta)$   
(same as before)

gradient(2) = code to compute  $\frac{\delta}{\delta \theta_1} J(\theta)$   

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

: etc