

Machine Learning

Topics :

1. Introduction
2. Linear Representation with One Variable
4. Linear Regression with Multiple Variables
5. Octave Tutorial
6. Logistic Regression
7. Regularization
8. Neural Network Representation
9. Neural Networks : Learning
10. Advice for Applying Machine Learning
11. Machine Learning System Design
12. Support Vector Machines
13. Clustering
14. Dimensionality Reduction
15. Anomaly Detection
16. Recommender Systems
17. Large Scale Machine Learning
18. Applications Example: Photo OCR

Machine Learning

1. Introduction

Machine Learning:

- grew out of AI
- new capability for computers

Examples:

- data mining
 - web click - better serve users
 - medical records - understand disease
 - biology - understanding DNA
 - engineering -
- applications which can't be programmed by hand
 - autonomous helicopter
 - handwriting recognition
 - computer vision
- self-customizing programs
 - recommendations
- understand human learning (brain, real AI)

1.2 What is Machine Learning?

Definition: Arthur Samuel (1959) checkers playing alg.
Field of study that gives computers the ability to learn without being explicitly programmed

Definition: Tom Mitchell (1998) Rhyming definition

Well posed learning problem: A computer program is said to learn from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with E.

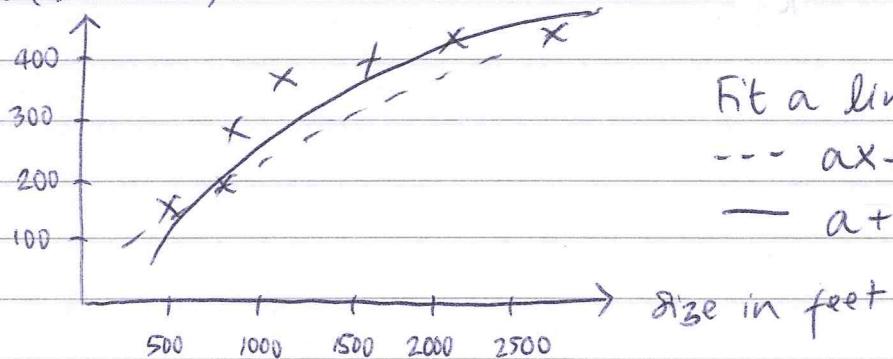
Machine learning algorithms

- supervised learning: we teach the computer
- unsupervised learning: computer learns by itself
- Reinforcement learning
- Recommender systems

1.3 Supervised Learning

1.3.1 House price prediction

price (\$ in 1000's)



Fit a line through the data

$$\dots ax + b = y$$

$$— a + bx^2 = y$$

Supervised learning:
Right answers given

Regression:

Predict continuous valued output (price)

1.3.2 Breast cancer (malignant, benign)

malignant?

*

(Y) 1

(N) 0

X X X X X

: : : : :

Tumour size

: : : : :

0 0 0 * 0 * 0 X X X

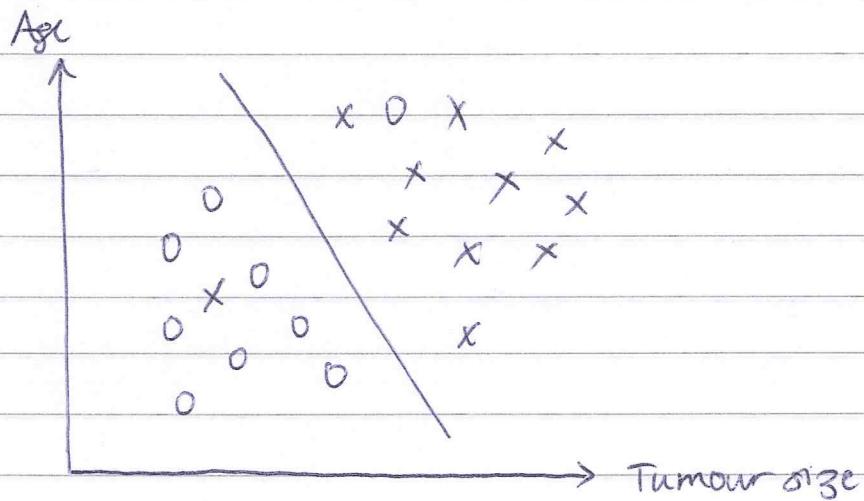
Tumour size

- here we have one feature as input...

Classification:

Discrete valued output ($0 = N$ or $1 = Y$)

*

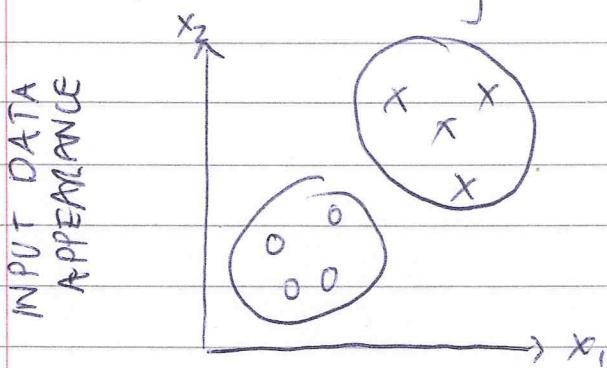


- here two features as input

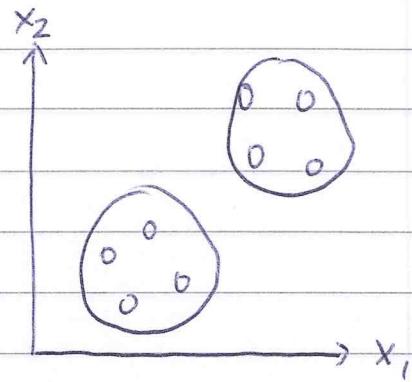
other possible features

- Clump thickness
- uniformity of cell size
- uniformity of cell shape

1.4 Unsupervised Learning



supervised learning



unsupervised learning

input data is provided and the algorithm looks
for structure in it

- Clustering algorithms used in
- Google news for grouping stories on the web
 - Gene data -
 - Organise computer clusters
 - Social network analysis
 - Market segmentation
 - Astronomical data analysis

Cocktail party problem

- two speakers, two microphones - the problem is to separate the sounds from the two people out

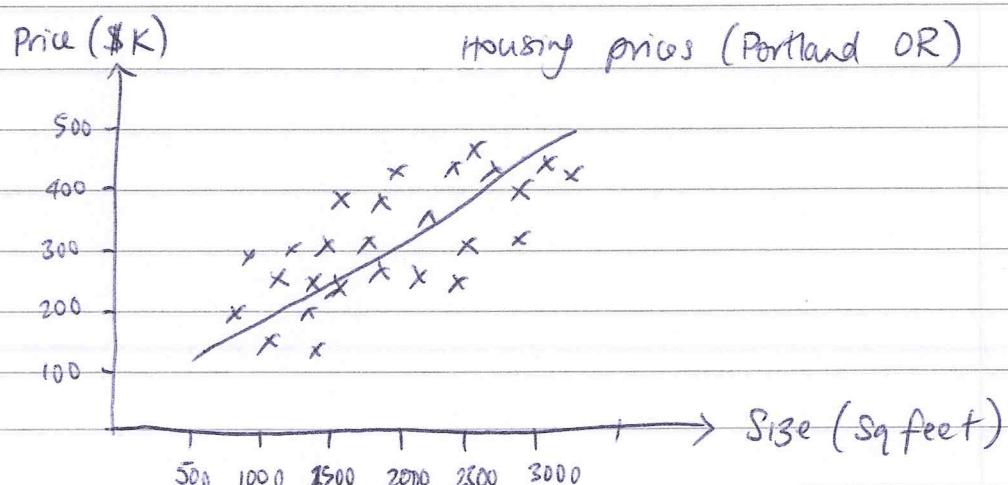
Solution:

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

octave open source free language - fast good for prototyping
 svd = singular value decomposition

2. Linear Regression with one variable

2.1 Model Representation



Supervised learning:

Given the "right answer" for each example in the data

Regression problem:

Predict real valued output
 as opposed to a classification problem with discrete valued output

Training set of housing prices
(Portland, OR)

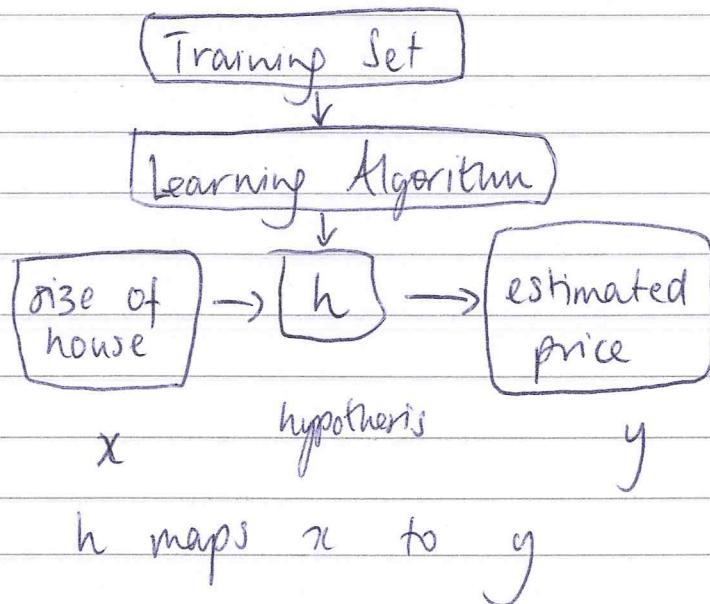
Size in feet (x)	Price (\$K) (y)
2104	460
1416	232
1534	315
852	178
...	...

Notation:

$m = \#$ training examples

x 's = input variable / feature

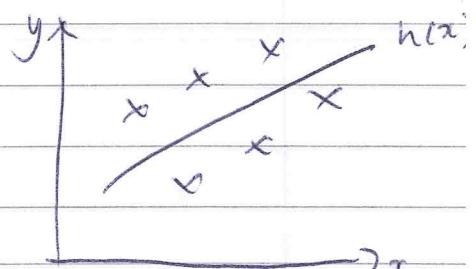
y 's = output variable / target variable
 $(x^{(i)}, y^{(i)})$ = i th training example



how do we represent

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

shorthand $h(x)$



This is called

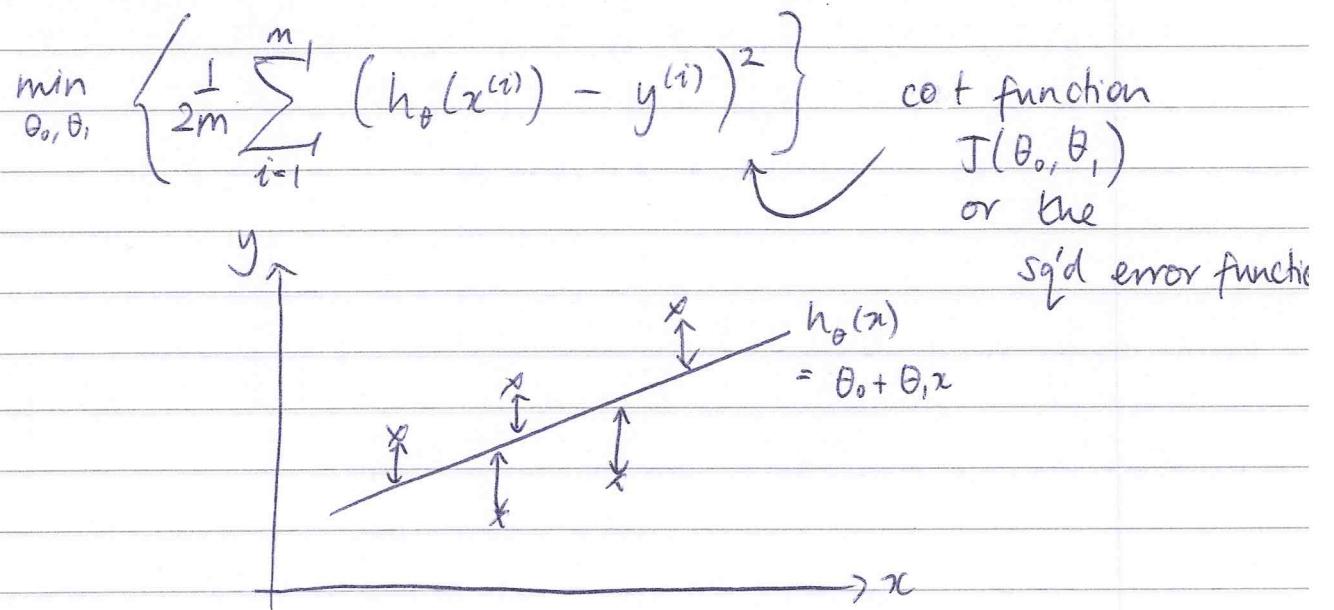
- Linear regression with one variable or
 - Univariate linear regression

2.2 Cost function (Intuition I and II)

Different choices of Θ_0 and Θ_1 parameters give different hypothesis functions

We choose θ_0 and θ_1 so that $h_\theta(x)$ is close to y for our training examples (x, y) . To do

this we solve a minimisation problem i.e.
minimise



2.3 Gradient descent

The idea is that if you plot your cost function in terms of its parameters you work your way (downhill) to a minimum as quickly as possible

Gradient descent algorithm:

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

learning rate = step size downhill

Correct:

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

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

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

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

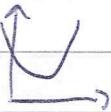
Incorrect: will probably not work is not gradient descent

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

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

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

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

Recall: convex  concave 

4 Linear Regression with Multiple Variables

4.1 Multiple features

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

So now we have

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

for convenience, define $x_0 = 1$

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

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

This is called Multivariate Linear Regression

4.2 Gradient descent for multiple variables

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

Gradient descent:

repeat $\left\{ \theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j} \right\}$ simultaneously update for every $j = 0, \dots, n$

Recall that for one variable we had

Repeat $\left\{ \begin{array}{l} \theta_0 := \theta_0 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})}_{\frac{\partial}{\partial \theta_0} J(\theta)} \\ \theta_1 := \theta_1 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}}_{\frac{\partial}{\partial \theta_1} J(\theta)} \end{array} \right\}$

simultaneously update θ_0 and θ_1

Now we have

Repeat $\left\{ \theta_j := \theta_j - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}}_{\text{update } \theta_j \text{ for } j=0, 1, \dots, n}$ *simultaneously* } *for*

4.3 Gradient descent in practice I : feature scaling

If the various features are orders of magnitude different, it's worth scaling them to similar sizes. This leads to faster convergence.

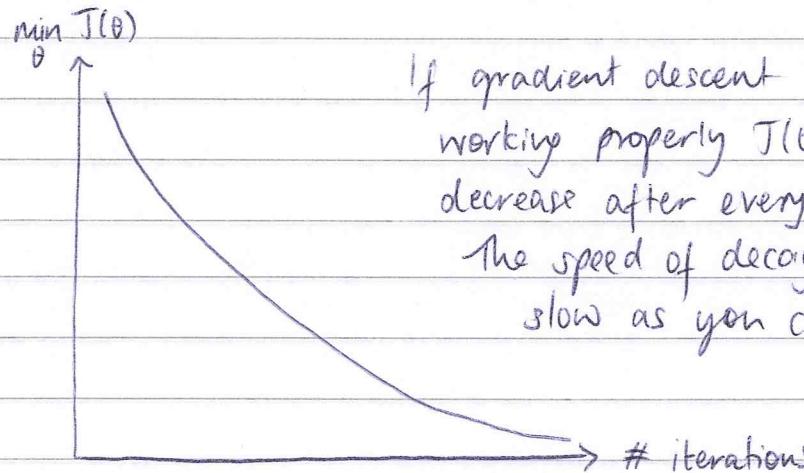
Mean normalisation:

Replace x_i with $(x_i - \mu_i)/\sigma_i$ to make features have approximately zero mean and a reasonable range.

$$\left. \begin{array}{l} \mu_i = \text{mean} \\ \sigma_i = \text{std dev / range} \end{array} \right\} \text{ith feature}$$

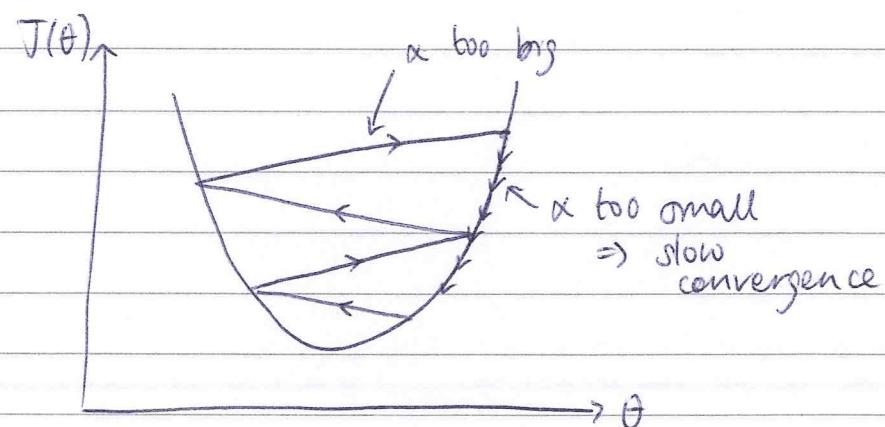
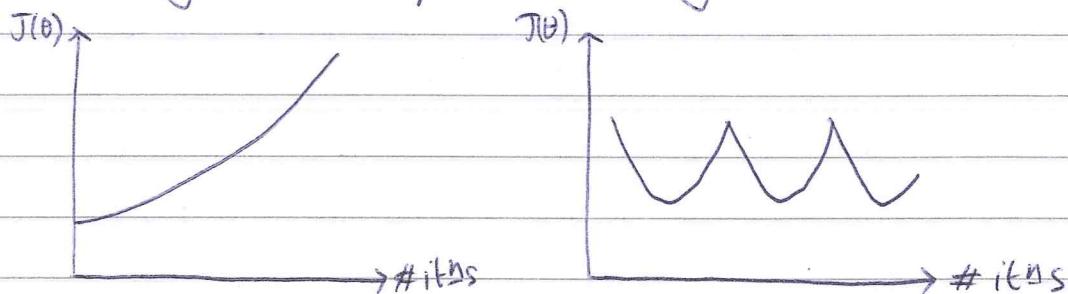
4.4 Gradient descent in practice II : Learning Rate

Making sure gradient descent is working properly.



If gradient descent is working properly $J(\theta)$ will decrease after every iteration. The speed of decay also slows as you converge

If $J(\theta)$ increases then it is likely that you are using too large a learning rate.



4.5 Features & Polynomial Regression

The linear dependence need not be assumed on the features you could for example optimise over

$$h_0(x) = \theta_0 + \theta_1(x, x_2)$$

or

$$h_0(x) = \theta_0 + \theta_1 x^2$$

you can fit any polynomial functions

for example. When using these types of functions feature scaling becomes even more important.

4.6 Normal Equation

This gives us an analytic formula for the optimal θ

Recall: In 1D $J(\theta)$ is just a quadratic in θ .

We can find the minimum analytically.

Clearly this will also work for n features, we just have to solve $n+1$ quadratic equations simultaneously i.e.

$$\theta \in \mathbb{R}^{n+1} \quad J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

$$\text{set } \frac{\partial J(\theta)}{\partial \theta_j} = 0 \quad \forall j \quad (\text{i.e. } j = 0, 1, \dots, n)$$

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

x_0	size ($t + t^2$)	# bedrooms	# floors	age	Rice
1	2104	5	1	95	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$$\underbrace{\quad\quad\quad}_{X} \qquad \qquad \qquad \underbrace{\quad\quad\quad}_{y}$$

$$X = \underbrace{\begin{pmatrix} 1 & X_1^{(1)} & X_2^{(1)} & \cdots & X_n^{(1)} \\ 1 & X_1^{(2)} & X_2^{(2)} & \cdots & X_n^{(2)} \\ 1 & \vdots & \vdots & & \vdots \\ 1 & X_1^{(m)} & X_2^{(m)} & \cdots & X_n^{(m)} \end{pmatrix}}_{X \in \mathbb{R}^{m, n+1}}, \quad \theta = \underbrace{\begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{pmatrix}}_{\theta \in \mathbb{R}^{n+1}}, \quad y = \underbrace{\begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{pmatrix}}_{y \in \mathbb{R}^m}$$

$$x_0^{(j)} = 1 \quad \forall j$$

Want θ such that $\min_{\theta} [(X\theta - y)^T (X\theta - y)]$

Normal equation:

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

In Octave this can be written as

$$\text{pinv}(X' * X) * X' * y$$

m training examples and n features

Gradient Descent

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

\rightarrow If $n > 10,000$ consider using gradient descent rather than.

Normal Equation

- No need to choose α
- No need to iterate
- Need to compute $(X^T X)^{-1}$
- $O(n^3)$
- Slow if n is very large

4.7 Normal Equation Non-invertability

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

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

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

↳ pseudo inverse as oppose to the actual inverse function

Two common causes of $X^T X$ being singular

→ Redundant features (linearly dependent features)

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

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

→ Too many features and not enough samples

e.g. $m \leq n$

Either delete some features or use regularisation

we'll come back to this

5 Octave Tutorial

5.1 Basic Operations

- Logical operations:
 - equal $\rightarrow A == B$
 - not equal $\rightarrow A \sim= B$
 - and $\rightarrow A \& \& B$
 - or $\rightarrow A || B, \text{xor}(A, B)$
- false = 0
- true = 1

- Command line: PS1(' >> ');

- Semicolon suppresses output

- apostrophes for strings
- % for comments
- `disp(sprintf('2 decimals: %.0.6f', a))`

format long 14dp

Short 4dp

- $A = [1 \ 2; 3 \ 4; 5 \ 6]$

- $v = [1 \ 2 \ 3] \quad v = [1; 2; 3]$

- $v = 1: 0.1: 2$

$$\downarrow \\ v = (1, 1.1, 1.2, \dots, 2)$$

- $\text{ones}(2, 3) \Rightarrow \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}$

- $2 * \text{ones}(2, 3) \Rightarrow \begin{pmatrix} 2 & 2 & 2 \\ 2 & 2 & 2 \end{pmatrix}$

- $\text{zeros}(1, 3) \Rightarrow (0 \ 0 \ 0)$

- $w = \text{rand}(1, 3)$

- $w = \text{randn}(1, 3)$

- $w = -6 + \sqrt{10} * (\text{randn}(1, 1000))$

- `hist(w)` prints the histogram
`hist(w, 50)` print with 50 bins

- $\text{eye}(4) \Rightarrow \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$

- help eye

- help rand

- help help

5.2 Moving data around

- $\text{size}(A) \Rightarrow [r, c]$ 1×2 matrix

- $\text{size}(A, 1)$ \Rightarrow number rows

- $\text{size}(A, 2)$ number columns

- $\text{length}(v)$

- $\text{length}(A)$ gives the longest dimension
i.e. $\max(r, c)$

- pwd gives current directory

- cd 'C:\...' to change directory

- ls to list files

- load featuresX.dat or `load('featuresX.dat')`
load priceY.dat

- who tells you which variables you have in memory
whos gives name, size, bytes and class (double, ..., ...)

- clear A
- $v = \text{priceY}(1:10)$ gives v = first 10 elements of priceY
- save hello.mat v;
- clear just clears everything
- save hello.txt v -ascii i.e. save as text (ASCII)

$A(3,2) \Rightarrow A_{3,2}$
 $A(2,:) \Rightarrow$ 2nd row of A
 $A(:,2) \Rightarrow$ 2nd column of A

$A([1,3], :) \Rightarrow$ 1st & 3rd row of A

$A(:,2) = [10; 11; 12] \Rightarrow$ sets 2nd column of A as specified

$A = [A, [100; 101; 102]] \Rightarrow$ appends the matrix with the column vector specified

$A(:) \Rightarrow$ puts all elements of A in a column vector

$C = [A \ B]$ } concatenates

$C = [A; B]$ } matrixies

n.b. $[A \ B] = [A, B]$

5.3 Computing on Data

- $A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{pmatrix} \quad B = \begin{pmatrix} 11 & 12 \\ 13 & 14 \\ 15 & 16 \end{pmatrix} \quad A.*B = \begin{pmatrix} 11 & 24 \\ 39 & 56 \\ 75 & 96 \end{pmatrix}$
- element wise operations e.g. $.*$, $./$, $.+$, etc..

- $\ln(A)$, $\exp(A)$, $\text{abs}(A)$
all element wise operations
- A' transpose A
- $A+1$ adds 1 to all elements
- $\max(A)$ column wise max default
 $\max(\max(A))$ for max of all elements
 $\max(v)$
- $v < 3$ gives true false element wise

$\text{find}(a < 3)$ returns the elements of a which are < 3

$\text{sum}(a)$ $\text{prod}(a)$ sum and product
 $\text{floor}(a)$ $\text{ceil}(a)$ rounding down and up

$\max(A, [], 1)$ max of each column
 $\max(A, [], 2)$ max of each row

$A = \text{magic}(9)$

$\text{sum}(A, 1)$ sums each column

$\text{sum}(A, 2)$ " " row

$A.*\text{eye}(9)$

$\text{sum}(\text{sum}(A.*\text{eye}(9)))$

$\text{flipud}(\text{givens}(9))$ gives 1 on off diag

5.4 Plotting data

$t = [0, 0.01, 1]; y_1 = \sin(2\pi \cdot 4 \cdot t)$
 $\text{plot}(t, y_1)$

hold on; % adds plot to the same graph

$y_2 = \cos(2\pi \times 4 \times t)$

plot(t, y2, 'r')

title('my plot')

xlabel('time')

ylabel('value')

legend('sin', 'cos')

print -dpng 'myplot.png'

close % closes the window with the plot in

figure(1); plot(t, y1)

figure(2); plot(t, y2)

subplot(1, 2, 1); plot(t, y1)

subplot(1, 2, 2); plot(t, y2)

axis([0.5 1 -1 1]) sets axis max and min

imagesc(A), colourbar, colourmap gray;

% comma chaining function call

% semi colons do the same but with
ans suppressed

5.5 Control statements : for, while & if statements and functions

v = zeros(10, 1)

for i=1:10, v(i)=2*pi*i; end;

indices = 1:10

for i=indices, disp(i); end;

break, continue

i=1

while $i \leq 5$, $v(i) = 100$; $i = i + 1$; end;

i=1;

while true,

$v(i) = 999$;

$i = i + 1$;

if $i == 6$,

break;

end;

end;

$v(i) = 2$

if $v(i) == 1$

disp('one');

elseif $v(i) == 2$

disp('two');

else

disp('not one/two');

end;

- Save functions in .m file
open these in WordPad or notepad

function $y = \text{squareThisNumber}(x)$
 $y = x^2;$

of course the file has to be in your working directory or add the place to your search paths

addpath('C:\Users\...')

function $[y_1, y_2] = \text{squareAndCubeThisNumber}(x)$
 $y_1 = x^2;$
 $y_2 = x^3;$

function $J = \text{costFunction}(X, y, \theta)$

- X = "design matrix" containing our training examples
- y is the class labels

$$m = \text{size}(X, 1)$$

$$\text{prediction} = X * \theta^T;$$

$$\text{sqrError} = (\text{prediction} - y) . \wedge 2;$$

$$J = 1/(2*m) * \text{sum}(\text{sqrError});$$

5.6 Vectorization

$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$$
$$= \theta^T x$$
$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$
$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

- Unvectorized implementation:

$$\text{prediction} = 0.0;$$

for $j = 1:n+1$,

$$\text{prediction} = \text{prediction} + \theta(j) * x(j)$$

end;

- Vectorized implementation

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

6. Logistic Regression

6.1 Classification

- Email : spam / not spam
- Online transactions : fraudulent / not fraudulent
- Tumour : Malignant / Benign

$y \in \{0, 1\}$ 0: +ve class
 1: -ve class

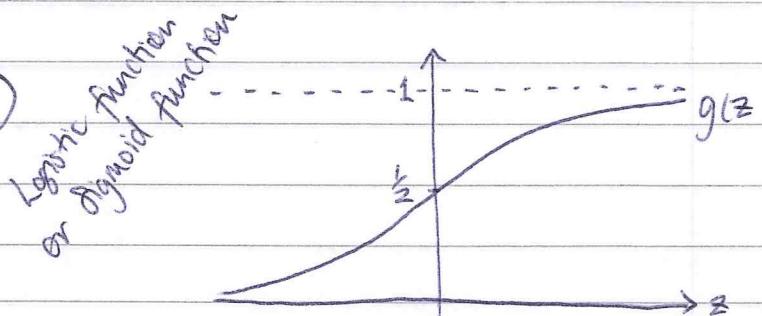
Linear regression not suitable for classification problems

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

6.2 Hypothesis representation

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

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

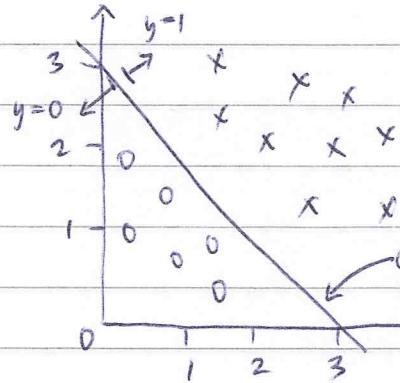


Interpretation:

$h_\theta(x)$ = estimated probability that y
 $= P(y=1 | x; \theta)$

6.3 Decision Boundary

Suppose we predict " $y=1$ " if $h_\theta(x) \geq 0.5$
i.e. $z = \theta^T x \geq 0$
⇒ "y=0" if $h_\theta(x) < 0.5$
i.e. $z = \theta^T x < 0$

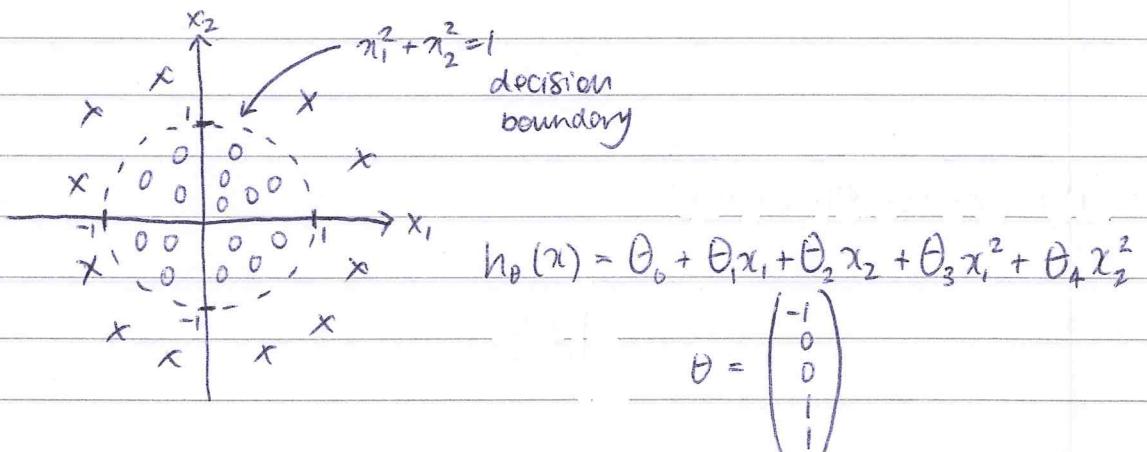


$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

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

predict ' $y=1$ ' if $-3 + \theta_1 x_1 + \theta_2 x_2 \geq 0$
 i.e. $x_1 + x_2 \geq 3$
 i.e.

Non-linear decision boundaries



$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2$$

$$\theta = \begin{pmatrix} -1 \\ 0 \\ 0 \\ 1 \end{pmatrix}$$

6.4 Cost function

Recall:

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

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

$$h_0(x) = \frac{1}{1 + e^{-\theta^T x}}$$

How do we choose our parameters θ ?

In linear regression ...

$$\text{Cost function } J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_0(x^{(i)}) - y^{(i)})^2$$

This cost function is non-convex for

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

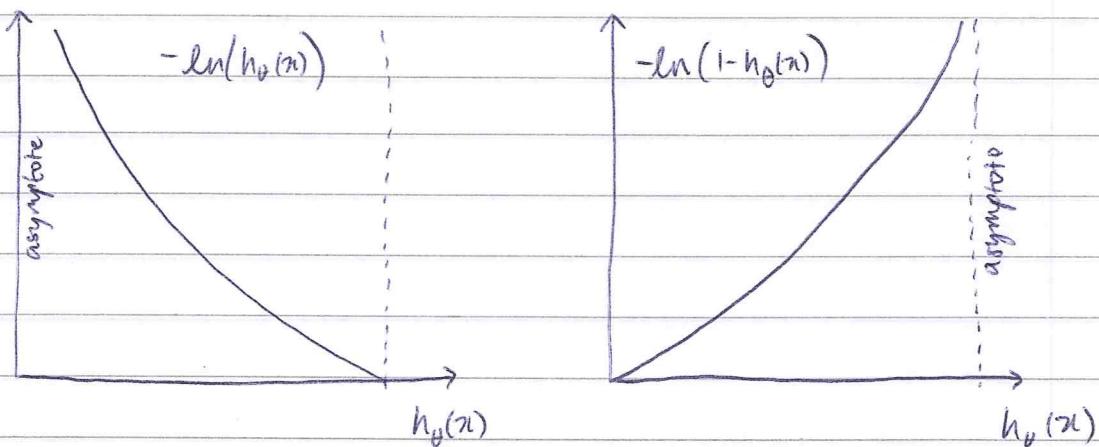
so this is not suitable.

Logistic regression cost function

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

if $y=1$

if $y=0$



6.5 Simplified cost function & gradient descent

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

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

rewrite:

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

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

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

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

Use gradient descent i.e.

$$\text{repeat } \left\{ \theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j} \right\}$$

$$\text{i.e. repeat } \left\{ \theta_j := \theta_j - \alpha \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \right\}$$

simultaneously update all θ_j

6.6 Advanced Optimisation

Optimisation algorithm

→ cost function $J(\theta)$. Want $\min J(\theta)$.

→ given θ , we have code that can compute

- $J(\theta)$ ← to monitor convergence
- $\frac{\partial J(\theta)}{\partial \theta_j}$ for $j = 0, 1, \dots, n$

Optimization algorithms:

- Gradient descent
- Conjugate gradient
- BFGS
- L-BFGS

Advantages

- no need to choose α
- faster than gradient descent

Disadvantages

- more complex

Example:

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad \min_{\theta} J(\theta)$$
$$J(\theta) = (\theta_0 - 5)^2 + (\theta_1 - 5)^2$$
$$\frac{\partial J(\theta)}{\partial \theta_0} = 2(\theta_0 - 5)$$
$$\frac{\partial J(\theta)}{\partial \theta_1} = 2(\theta_1 - 5)$$

function [jVal, gradient] = costFunction(theta)
jVal = (theta(1) - 5)^2 + (theta(2) - 5)^2;
gradient = zeros(2, 1);
gradient(1) = 2 * (theta(1) - 5);
gradient(2) = 2 * (theta(2) - 5);

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

initialTheta = zeros(2, 1);

[optTheta, functionVal, exitFlag] ...

= fminunc (@costFunction, initialTheta, options);

↑ ↑
unconstrained minimisation pointer

} "Gradient descent on steroids"!

$\theta \in \mathbb{R}^d$ $d \geq 2$ for this to work

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

remember octave uses 1 indexing

function [jVal, gradient] = costFunction(theta)

[code to compute ...]

$$\begin{aligned} jVal &= \dots \\ \text{gradient}(1) &= \frac{\partial}{\partial \theta_0} (J(\theta)) \\ \text{gradient}(2) &= \frac{\partial}{\partial \theta_1} (J(\theta)) \\ &\vdots && \vdots \\ \text{gradient}(n+1) &= \frac{\partial}{\partial \theta_n} (J(\theta)) \end{aligned}$$

6.7 Multiclass Classification: One-vs-all

Multiclass classification

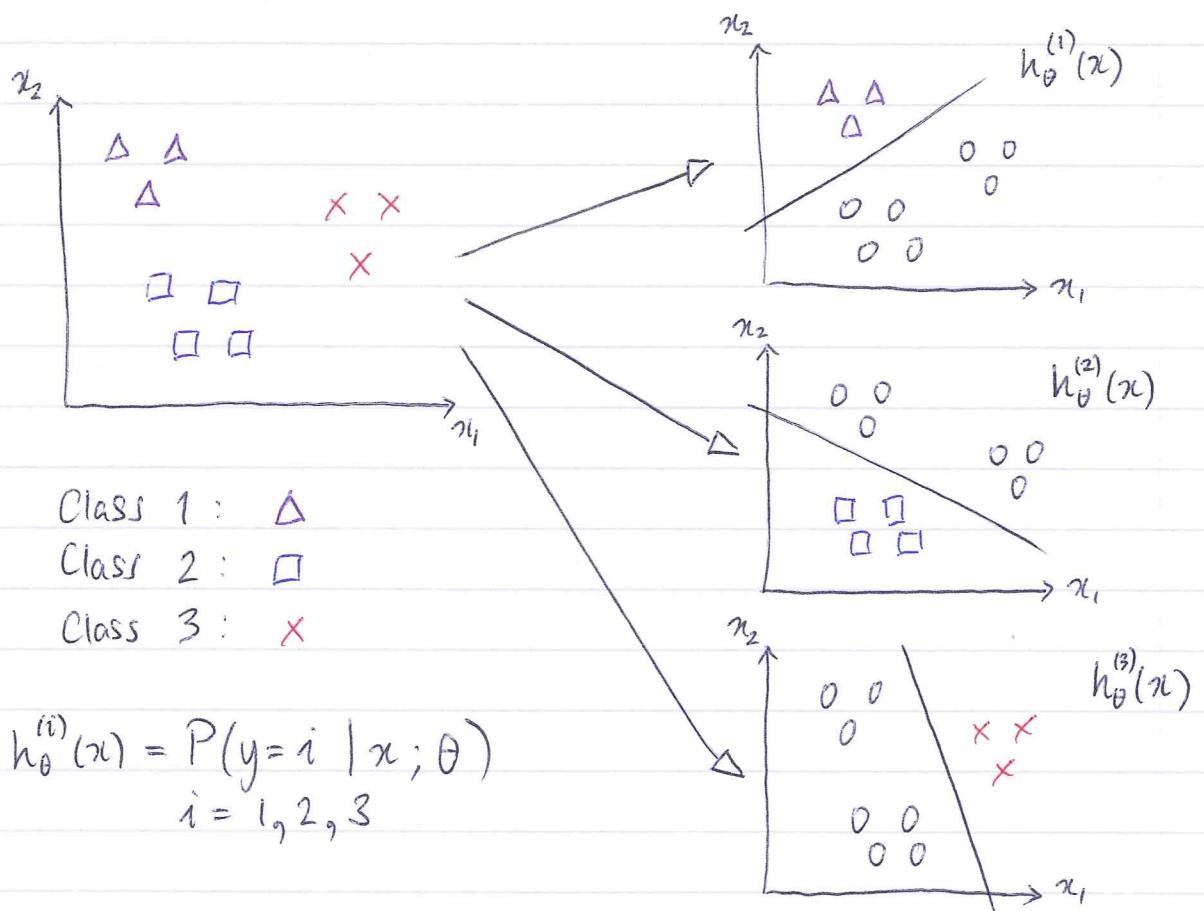
→ Email folder / tagging: Work / friends / family / hobby

→ Medical diagnosis: Not ill / cold / flu

→ Weather: Sunny / cloudy / rain / snow

here $y = 0, 1, 2, \dots, l$, $l = \text{no. of classes}$

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



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

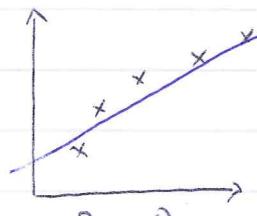
On a new input x , to make a prediction, pick the class i that maximises

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

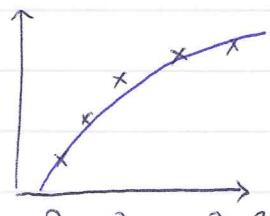
7 Regularization

7.1s The problem of over-fitting

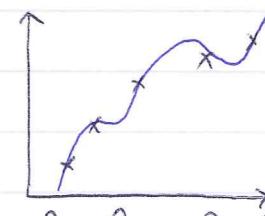
Example: Linear Regression



"underfit" / "high bias"



"just right"



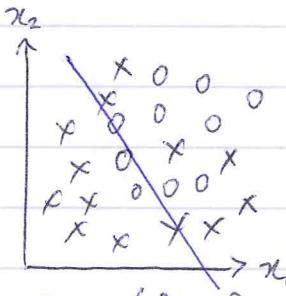
"overfit" / "high variance"

Example: Logistic Regression

Overfitting: If we have too many features the learned hypothesis may fit the training set very well i.e.

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \approx 0,$$

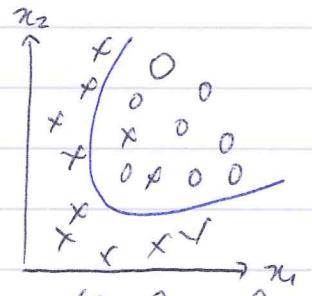
but fail to generalise to new examples (predict prices on new examples).



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

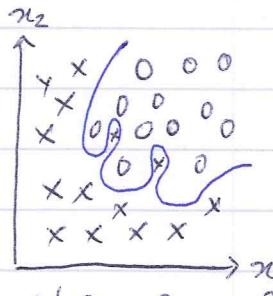
$g = \text{sigmoid function}$

"underfit"



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

"just right"



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

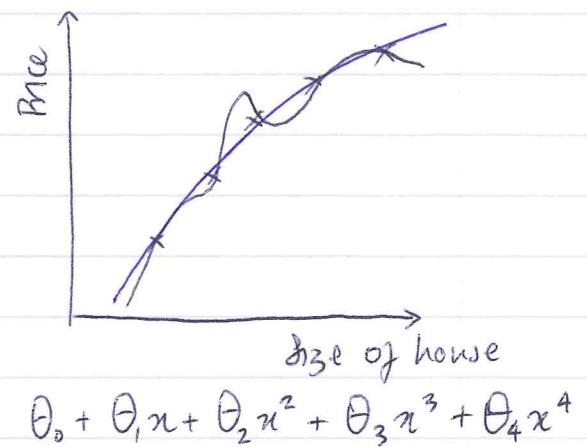
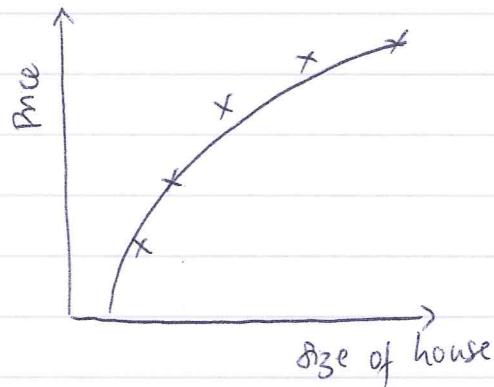
"overfit"

Addressing overfitting

- Common when you have many features and few training examples. Here you need to lose some features
 - manually select which features to keep
 - model selection algorithm to choose
- Regularization
 - keep all the features, but reduce magnitude / values of parameters θ_j
 - works well when we have lots of features, each of which contributes a bit to predicting y .

7.2 Cost function

Intuition



Suppose we penalize and force θ_3 and θ_4 to be small

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

→ the minimisation process will force $\theta_3 \approx 0$ and $\theta_4 \approx 0$

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

Housing:

- Features x_1, x_2, \dots, x_{100}
- Parameters $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

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

regularization parameter

we don't usually regularise θ_0 although in practice it makes little difference if we do.

- too large a λ choice can result in underfitting
- J algorithms which choose λ well for you

7.3 Regularized Linear Regression

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

Gradient descent

$$\text{Repeat } \left\{ \begin{array}{l} \theta_0 := \theta_0 - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_j := \theta_j - \frac{\alpha}{m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} - \lambda \theta_j \right] \end{array} \right\}_{j=1,2,\dots,n}$$

rewrite,

$$\theta_j := \underbrace{\theta_j \left(1 - \frac{\alpha \lambda}{m} \right)}_{< 1 \text{ usually only just e.g. } 1 - \frac{\alpha \lambda}{m} \approx 0.99} - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Normal Equation

$$X = \begin{pmatrix} x_1^{(1)} & x_2^{(1)} & \cdots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \cdots & x_n^{(2)} \\ \vdots & \vdots & & \vdots \\ x_1^{(m)} & x_2^{(m)} & \cdots & x_n^{(m)} \end{pmatrix} \quad y = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{pmatrix}$$

$$\min_{\theta} J(\theta) \quad \frac{\partial J(\theta)}{\partial \theta_j} = 0 \text{ gives min}$$

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

Non-invertability (this happens when $m \leq n$) - it turns out that as long as $\lambda > 0$ we no longer have this problem i.e. you can prove $X^T X + \lambda \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & \ddots \\ 0 & \ddots & 0 \end{bmatrix}$ is invertible

7.4 Regularized Logistic Regression

Here we have

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

Repeat

$$\begin{cases} \theta_0 := \theta_0 - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_j := \theta_j - \frac{\alpha}{m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} - \lambda \theta_j \right] \end{cases}$$

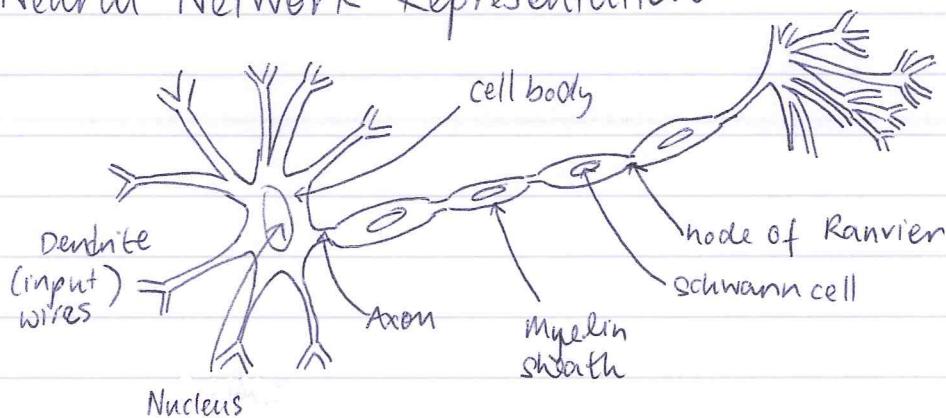
$\frac{\partial J(\theta)}{\partial \theta_j} \quad h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}} \quad j = 1, 2, \dots, n$

function [jVal, gradient] = costFunction(theta)

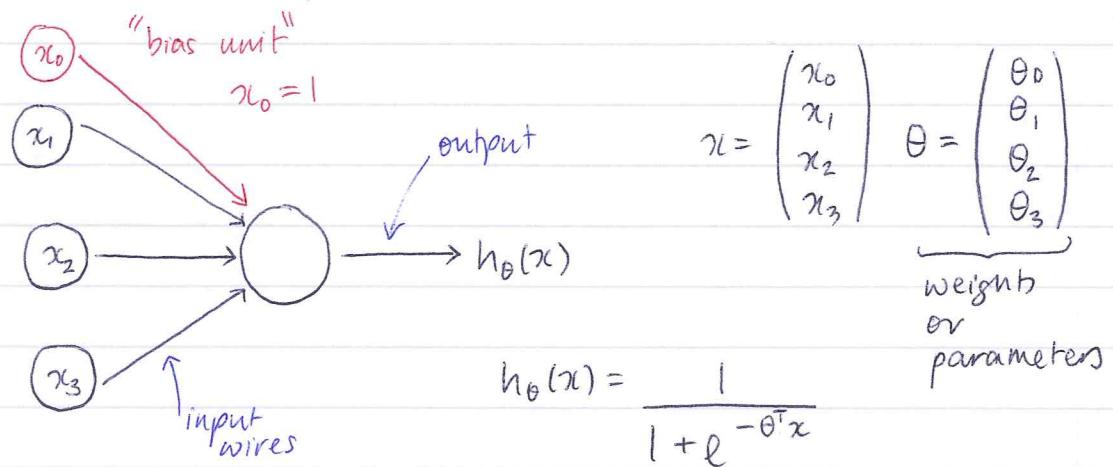
[code to compute ...]

$$\begin{aligned} jVal &= \dots J(\theta) = \text{as above} \\ \text{gradient (1)} &= \dots \frac{\partial}{\partial \theta_0} (J(\theta)) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \text{gradient (2)} &= \dots \frac{\partial}{\partial \theta_1} (J(\theta)) \\ &\vdots \\ \text{gradient (n+1)} &= \dots \frac{\partial}{\partial \theta_n} (J(\theta)) \end{aligned} \quad \left\{ \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} - \frac{\lambda}{m} \theta_j \right\}$$

8 Neural Network Representation

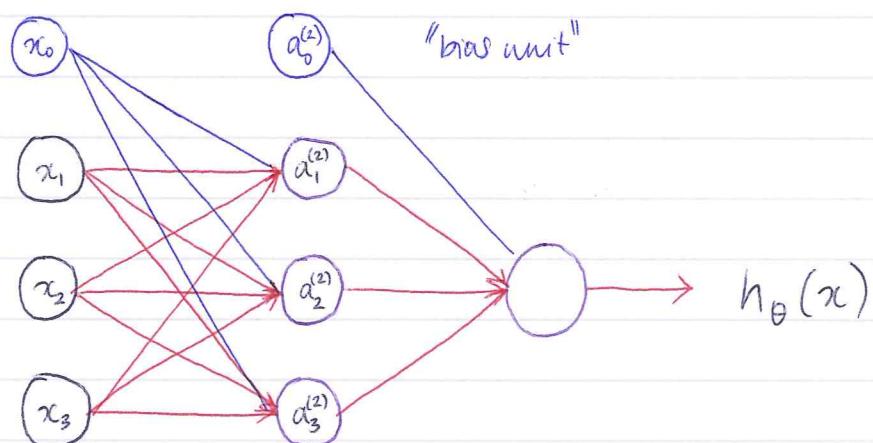


Neuron model : Logistic unit



This is a neuron with a Sigmoid / Logistic activation function

Neutral Network



Layer 1 (input layer)	Layer 2 (hidden layer)	Layer 3 (output layer)
X	$a_i^{(2)}$	y

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

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

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

$$a_2^{(2)} = g(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3) = g(z_2^{(2)})$$

$$a_3^{(2)} = g(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3) = g(z_3^{(2)})$$

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

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

Forward propagation: Vectorised implementation

$$\begin{pmatrix} \theta_{10}^{(1)} & \theta_{11}^{(1)} & \theta_{12}^{(1)} & \theta_{13}^{(1)} \\ \theta_{20}^{(1)} & \theta_{21}^{(1)} & \theta_{22}^{(1)} & \theta_{23}^{(1)} \\ \theta_{30}^{(1)} & \theta_{31}^{(1)} & \theta_{32}^{(1)} & \theta_{33}^{(1)} \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{pmatrix}$$

$$\Rightarrow \theta^{(1)} x = z^{(2)} \text{ and}$$

$$a^{(2)} = \begin{pmatrix} a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \end{pmatrix} = g(z^{(2)})$$

Add $a_0^{(2)}$ corresponding to our bias unit

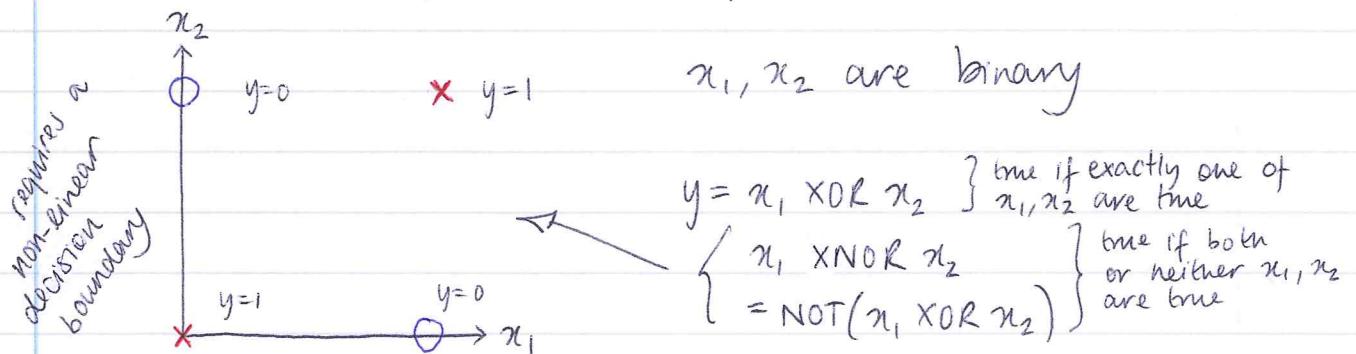
$$a^{(2)} = \begin{pmatrix} a_0^{(2)} \\ a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \end{pmatrix}$$

x can be thought of a $a^{(1)}$

$$\theta^{(2)} a^{(2)} = z^{(3)} \text{ and } h_\theta(x) = a^{(3)} = g(z^{(3)})$$

In neural networks we use logistic regression recursively to learn its own features to feed in $(a^{(2)})$ rather than $(x = a^{(1)})$. The network can have many layers. There is only one input layer and one output layer, all the others are hidden layers.

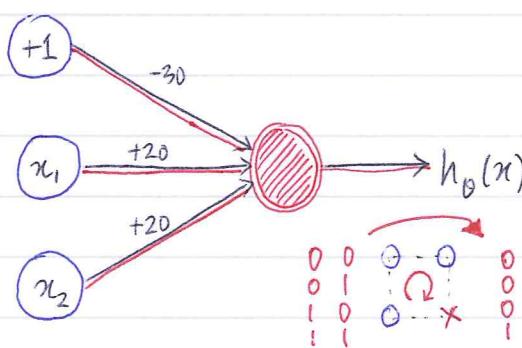
Non-linear classification example : XOR / XNOR



Get a neural network to fit this training set

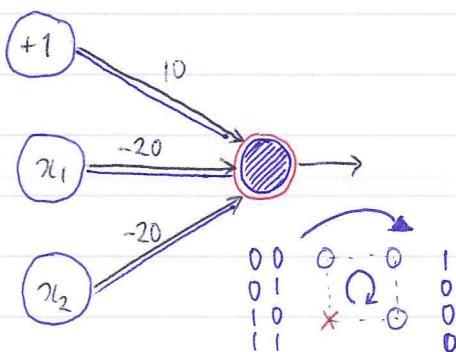
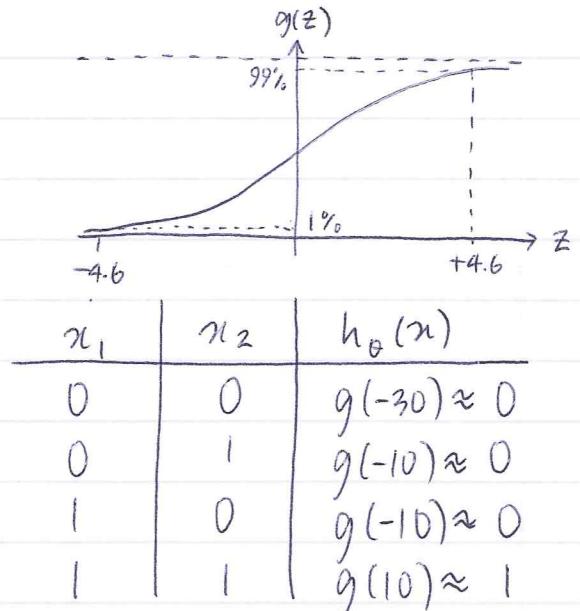
Simple example : AND

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

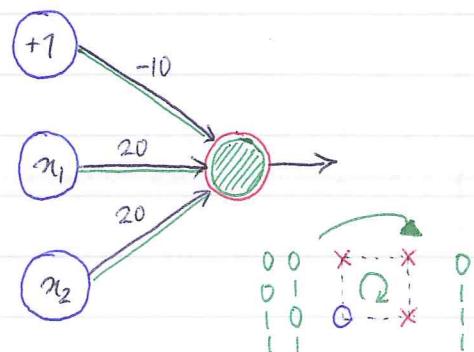


$$h_{\theta}(x) = g(-30 + 20x_1 + 20x_2)$$

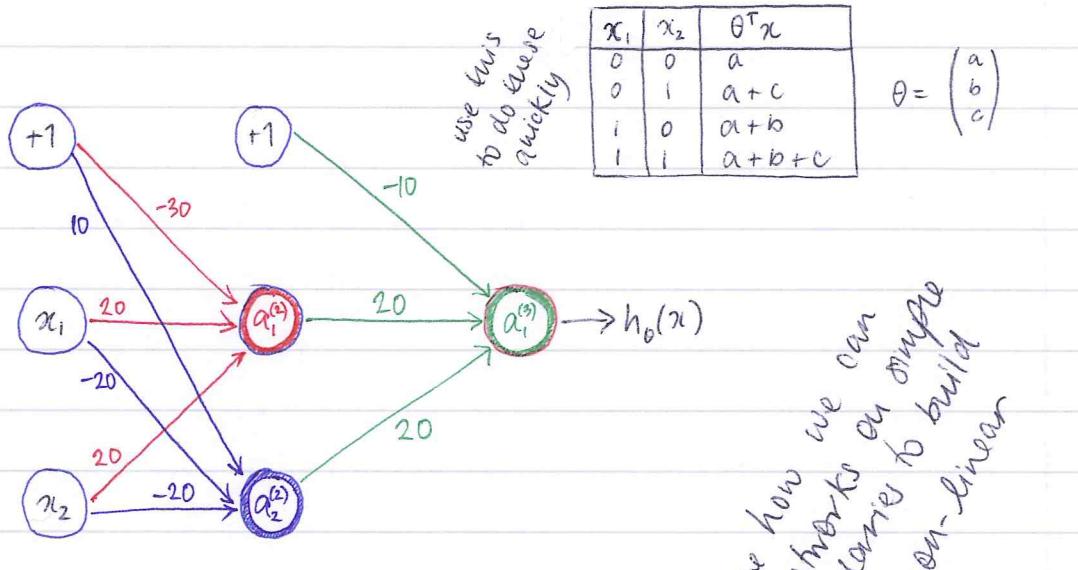
\uparrow \uparrow \uparrow
 $\theta_{10}^{(1)}$ $\theta_{11}^{(1)}$ $\theta_{12}^{(1)}$



(NOT x_1) AND (NOT x_2)



π_1 OR π_2



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

so we see how decision boundaries can be more complex to build non-linear ones!

Multiclass Classification

Multiple output unit: One vs All

If we want to classify data as one of k classes we need to have k elements in our output layer - one for each class. In this case we should predict

$$h_\theta(x) = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \text{ when the example is in our first class}$$

$$h_\theta(x) = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \text{ when the example is in our second class}$$

and so on, i.e. $h_\theta(x) \in \mathbb{R}^k$

As before we have our training set

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

$$\text{but now } y^{(i)} \in \left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \dots, \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \right\}$$

9 Neural Networks: Learning

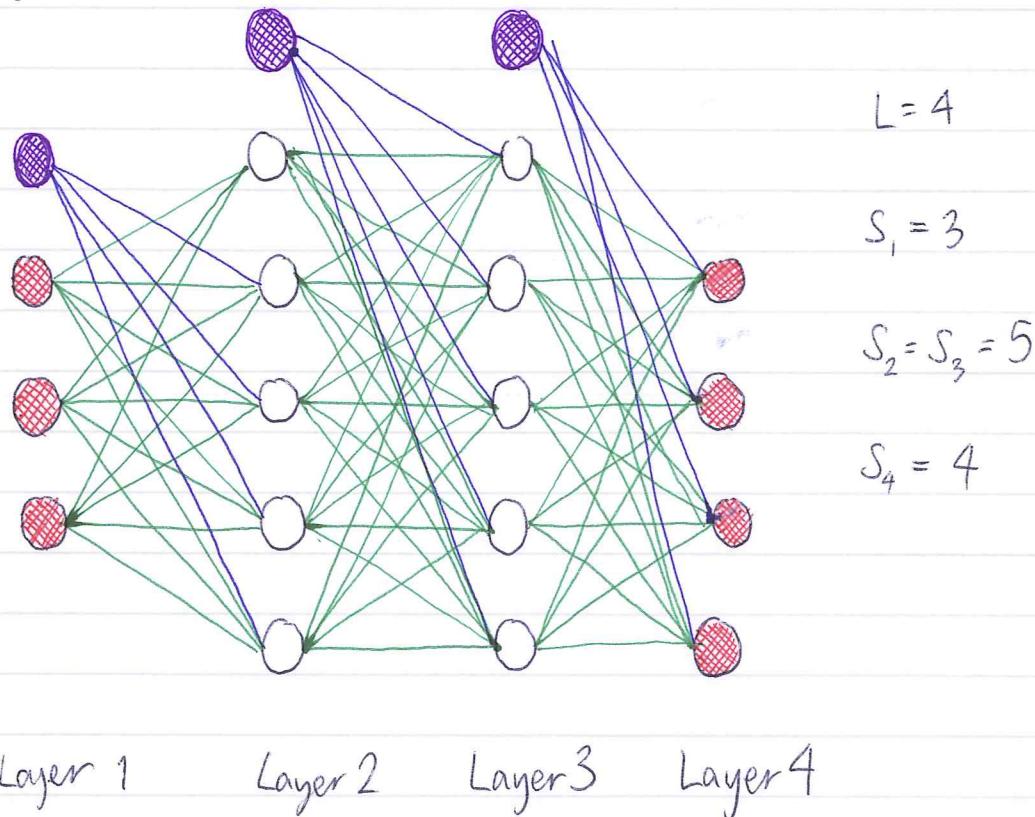
Neural Network (Classification)

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

L = total number of layers in network (inc. input and output layers).

s_L = number of units (not including the bias unit) in Layer L

Example:



Layer 1 Layer 2 Layer 3 Layer 4

Binary classification

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

$$h_\theta(x) \in \mathbb{R} \quad K=1$$

Multi-class classification (K classes)

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

$$h_\theta(x) \in \mathbb{R}^K \quad K \geq 3$$

$$s_L = K$$

9.1 Cost Function

Logistic regression:

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

Neural Network:

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

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

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

9.2 Backpropagation Algorithm

Gradient computation

Given one training example

(x, y) : Forward propagation:

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

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

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

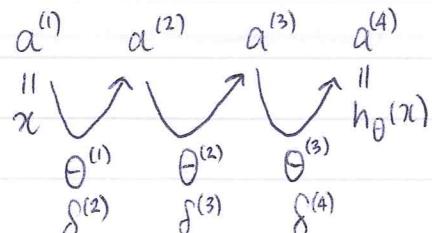
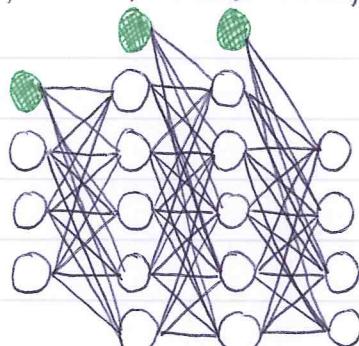
$$z^{(3)} = \theta^{(2)} a^{(2)}$$

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

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

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

Layer 1 Layer 2 Layer 3 Layer 4



Gradient computation: Backpropagation algorithm

Intuition: $\delta_j^{(l)}$ = error for node j in layer l

For each output unit in layer $L=4$

$$\delta_j^{(4)} = a_j^{(4)} - y_j \quad a_j^{(4)} = (h_\theta(x))_j$$

i.e. $\delta_j^{(4)} = a_j^{(4)} - y$

$$\delta^{(3)} = (\Theta^{(3)})^\top \delta^{(4)} \cdot g'(z^{(3)})$$

$$\delta^{(2)} = (\Theta^{(2)})^\top \delta^{(3)} \cdot g'(z^{(2)})$$

no $\delta^{(1)}$

derivatives

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

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

$$\frac{\partial J(\Theta)}{\partial \Theta_{ij}^{(l)}} = a_j^{(l)} \delta_i^{(l+1)} \quad (\text{for } j=0)$$

\uparrow
regularisation term

Now assume we have a large training set

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

Set $\Delta_{ij}^{(l)} = 0 \quad (\forall l, i, j)$

For $i=1$ to m

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

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

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

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

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

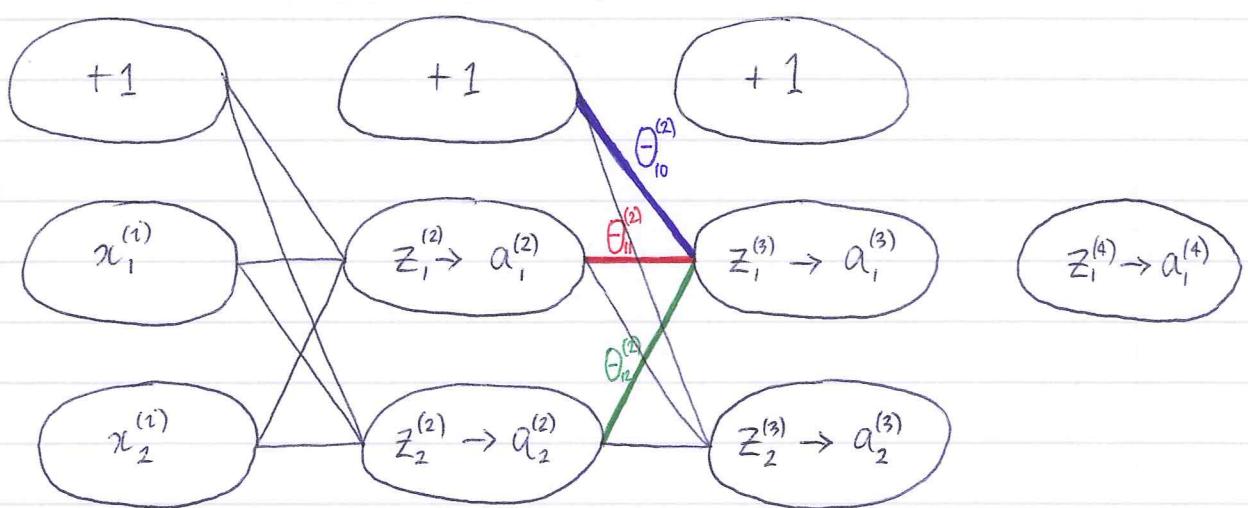
End For

$$D_{ij}^{(l)} := \frac{1}{m} (\Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)}) \quad \text{if } j \neq 0 \quad \left| \begin{array}{l} \frac{\partial J(\Theta)}{\partial \Theta_{ij}^{(l)}} = D_{ij}^{(l)} \\ \frac{\partial J(\Theta)}{\partial \Theta_{0j}^{(l)}} \end{array} \right.$$

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

Propagation Intuition

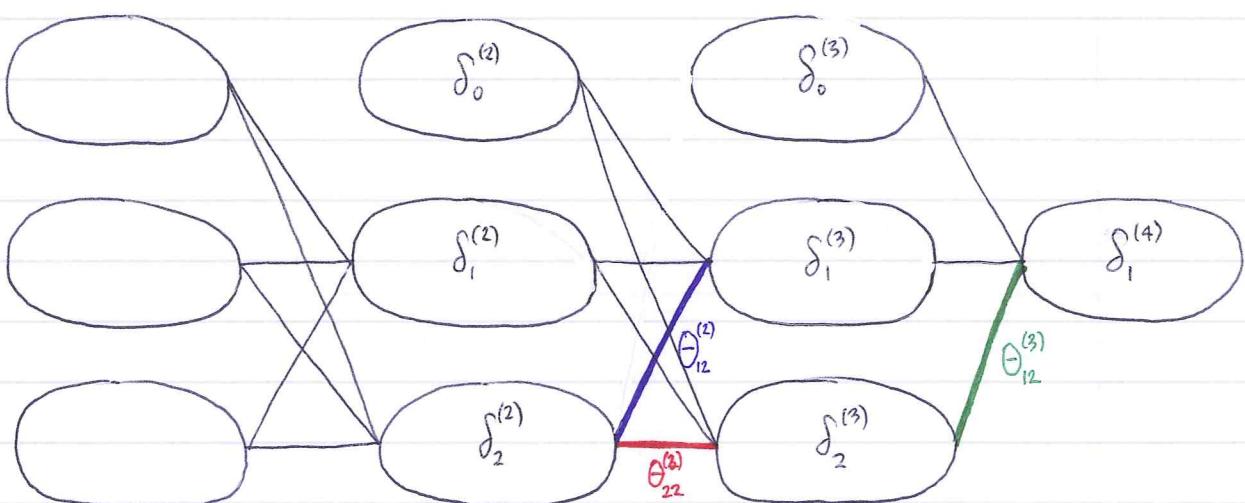
Forward Propagation



$$(\gamma^{(i)}, y^{(i)})$$

$$z_1^{(2)} = \Theta_{10}^{(1)} \times 1 + \Theta_{11}^{(1)} a_1^{(1)} + \Theta_{12}^{(1)} a_2^{(1)}$$

Back Propagation



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

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

$$\begin{aligned} \delta_1^{(4)} &= y^{(i)} - a_1^{(4)} \\ \delta_2^{(2)} &= \Theta_{12}^{(2)} \delta_1^{(3)} + \Theta_{22}^{(2)} \delta_2^{(3)} \end{aligned}$$

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

9.4 Implementation note: unrolling parameters

To use the advanced optimisation functions in Octave we must have our parameters $\Theta_{ij}^{(u)}$ and gradients $D_{ij}^{(u)}$ in vector (rather than matrix) form.

function [jVal, gradient] = costFunction(theta)

$\underbrace{\quad\quad\quad}_{R^{n+1}}$

optTheta = fminunc (@costFunction, initialTheta, options)

Neural network $L=4$

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

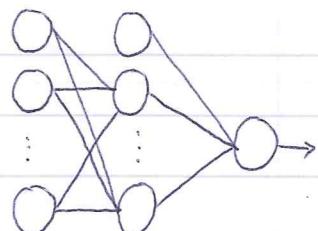
$\rightarrow D^{(1)}, D^{(2)}, D^{(3)}$ - matrices (D1, D2, D3)

"Unroll" into vectors

Example : $S_1 = 10 \quad S_2 = 10 \quad S_3 = 1$

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

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



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

$$\text{DVec} = [\text{D1}(:); \text{D2}(:); \text{D3}(:)];$$

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

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

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

Learning Algorithm

\rightarrow Unroll initial param's $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ to get initialTheta to pass into fminunc (@costFunction, initialTheta, option)

\rightarrow function [jVal, gradient] = costFunction(thetaVec)

From thetaVec get $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ and then use fwd/back to compute $D^{(1)}, D^{(2)}, D^{(3)}, J(\theta)$ then unroll to get gradientVec.

9.5 Gradient checking

Numerical estimation of gradients

To check our code is working we can calculate the gradient by central differencing and then check that our approximation is close to DVec.

i.e.

$\Theta \in \mathbb{R}^n$ (our unrolled version of $\Theta^{(1)}, \Theta^{(2)}, \dots$)

$$\Theta = [\theta_1, \theta_2, \dots, \theta_n]$$

$$\frac{\partial}{\partial \theta_i} J(\theta) = \frac{J(\theta_1, \dots, \theta_{i-1}, \theta_i + \epsilon, \theta_{i+1}, \dots, \theta_n) - J(\theta_1, \dots, \theta_{i-1}, \theta_i - \epsilon, \theta_{i+1}, \dots, \theta_n)}{2\epsilon}$$

for $i=1:n$,

thetaPlus = theta;

thetaPlus(i) = thetaPlus(i) + EPSILON;

thetaMinus = theta

thetaMinus(i) = thetaMinus(i) - EPSILON;

gradApprox(i) = $(J(\theta_{\text{plus}}) - J(\theta_{\text{minus}})) / (2 * \text{EPSILON})$

end;

$\text{EPSILON} = 10^{-8}$
is a good number

Check that $\text{gradApprox} \approx \text{DVec}$

Be sure to disable the gradient checking code before training the classifier as it will slow the code down a lot.

For gradient descent and advanced optimisation we need an initial value for theta, Θ

optTheta = fminunc (@costFunction, initialTheta, options)

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

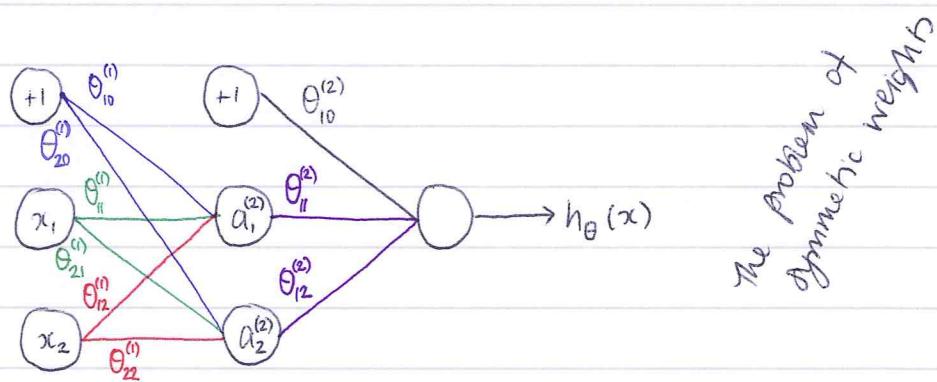
9.6 Random Initialisation

How do we choose our initial Θ ? It turns out that in a neural network we can't choose Θ

$$\Theta_{ij}^{(l)} = x \quad \forall i,j,l.$$

If we do, the parameters corresponding to inputs going into each hidden unit are identical

e.g.



$$\theta_{10}^{(1)} = \theta_{20}^{(1)} \quad \theta_{11}^{(1)} = \theta_{21}^{(1)} \quad \theta_{12}^{(1)} = \theta_{22}^{(1)}$$

$$\theta_{10}^{(2)} = \theta_{12}^{(2)} \Rightarrow a_1^{(2)} = a_2^{(2)} \\ \theta_{11}^{(2)} = \theta_{12}^{(2)} \Rightarrow s_1^{(2)} = s_2^{(2)}$$

also

$$\Rightarrow \frac{\partial J(\theta)}{\partial \theta^{(1)}} = \frac{\partial J(\theta)}{\partial \theta_{20}^{(1)}} \Rightarrow \theta_{10}^{(1)} = \theta_{20}^{(1)},$$

$$\frac{\partial J(\theta)}{\partial \theta_{11}^{(1)}} = \frac{\partial J(\theta)}{\partial \theta_{21}^{(1)}} \Rightarrow \theta_{11}^{(1)} = \theta_{21}^{(1)},$$

$$\text{and } \frac{\partial J(\theta)}{\partial \theta_{12}^{(1)}} = \frac{\partial J(\theta)}{\partial \theta_{22}^{(1)}} \Rightarrow \theta_{12}^{(1)} = \theta_{22}^{(1)}$$

on the next iteration

$$\Rightarrow a_1^{(2)} = a_2^{(2)} \text{ and round and round we go!}$$

This means our neural network becomes degenerate and unable to calculate interesting functions

Random Initialisation: Symmetry Breaking

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

e.g.

$$\Theta_{11}(1) = \text{rand}(10, 1) * (2 * \text{INIT_EPSILON}) - \text{INIT_EPSILON};$$

$$\Theta_{12}(1) = \text{rand}(1, 11) * (2 * \text{INIT_EPSILON}) - \text{INIT_EPSILON};$$

ϵ epsilon small

N.B. $\text{rand}(m, n)$ gives a $m \times n$ matrix whose elements are randomly drawn from the distribution $U(0, 1)$

9.7 Putting it together

Training a neural network:

1. Pick a network architecture (connectivity pattern between neurons)

→ No. input units: Dimension of features x

→ No. output units: Number of classes

→ Reasonable default: 1 hidden layer, or if > 1 , have the same number of hidden units in every layer (usually the more the better)

→ More hidden units than input features is a good idea

2. Randomly initialize weights

3. Implement forward propagation to get $h_\theta(x^{(i)})$ for any $x^{(i)}$

4. Implement code to compute the cost function $J(\theta)$

5. Implement back propagation to compute partial derivatives

$$\frac{\partial J(\theta)}{\partial \Theta_{j,k}^{(l)}}$$

→ for $i=1:m$ {

Perform fwd prop and back prop on example $(x^{(i)}, y^{(i)})$
(Get activations $a^{(l)}$ and delta terms $\delta^{(l)}$ for $l=2, 3, \dots, L$)

$$a^{(l)} := \text{sigmoid}(\Theta^{(l-1)} * a^{(l-1)});$$

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

compute $\frac{\partial J(\Theta)}{\partial \Theta_{jk}^{(l)}}$

5. Use gradient checking to compare $\frac{\partial J(\Theta)}{\partial \Theta_{jk}^{(l)}}$ computed

back prop vs numerical estimate of gradient of $J(\Theta)$

→ Then disable gradient checking code

6. Use gradient descent or advanced optimisation method with back-propagation to try to minimize $J(\Theta)$ as a function of parameters Θ

N.B. $J(\Theta)$, for neural networks, in general is non-convex
so subject to local minima but in practice
this is not in general a massive problem.

10 Advice for applying machine learning

10.1 Decisions what to try next

Debugging a learning algorithm:

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

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

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

ment?

- Get more training examples
- Try a smaller set of features
- Try getting additional features
- Try adopting polynomial features
- Try decreasing λ
- Try increasing λ

Machine Learning diagnostic:

This is a test you can run to gain insight into what is/isn't working with a learning algorithm, and gain guidance as to how best to improve its performance. Diagnostics can take time to implement but can be a good use of time.

10.2 Evaluating a hypothesis

We have seen how if we overfit our training data our hypothesis fails to generalise to new examples not in our training set. For a small number of feature overfitting can be picked up by simply plotting our hypothesis against our training data but for a large number of features it becomes harder to detect.

Training/testing procedure for linear regression

- Take 70% of your data for training and use the remaining 30% for testing
- Learn the parameter Θ from your training set (by minimising $J(\Theta)$)
- Compute your test error:

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

For logistic regression the procedure is the same but the test error is given by

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

or we can use the "misclassification error" or "0/1 misclassification error":

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

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

10.3 Model Selection and Train Validation Test Set

In general, once parameters $\theta_0, \theta_1, \dots, \theta_n$ have been fit to some set of data (the 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 generalisation error.

Model selection:

Suppose you want to decide what degree polynomial to use to fit your data. Let d be the degree

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

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

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

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

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

We could pick d with the lowest test set error $J_{\text{test}}(\theta^{(d)})$ but now

we have no data left to do a generalisation error on since clearly the test data is now biased.

\uparrow
by fitting to training data
 \uparrow
error on test data

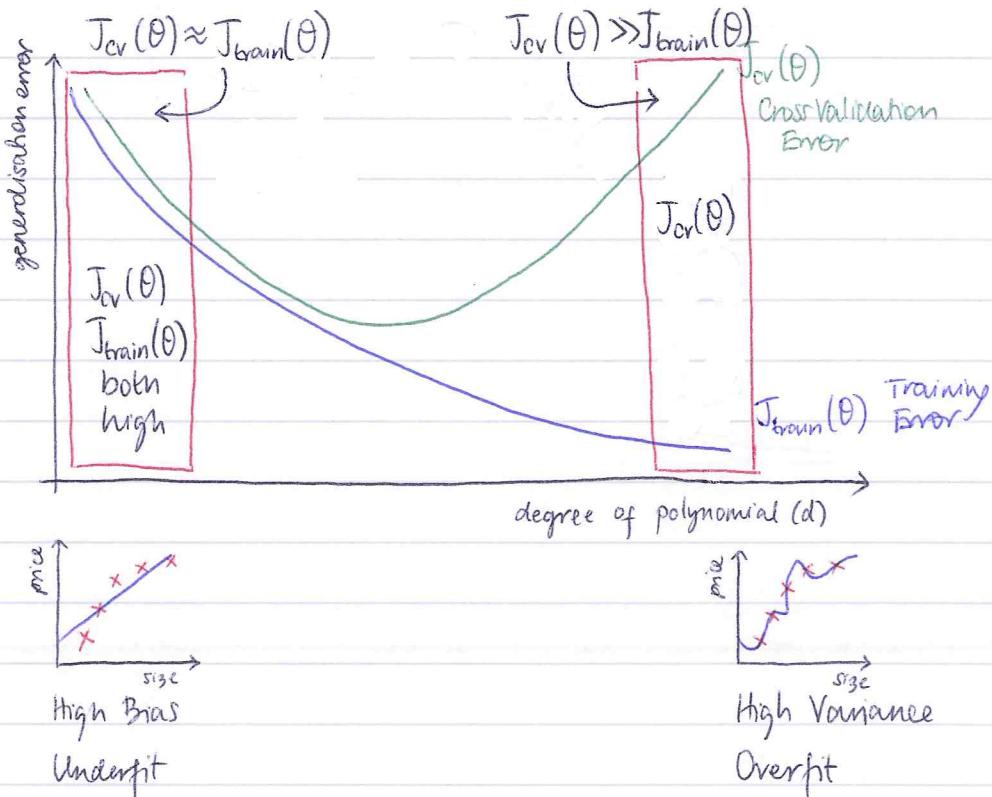
Cross-Validation:

We split our data into 3 sets ...

60% Training Set $(x^{(i)}, y^{(i)})$ m	20% Cross Validation set $(x_{cv}^{(i)}, y_{cv}^{(i)})$ m_{cv}	20% Test set $(x_{test}^{(i)}, y_{test}^{(i)})$ m_{test}
--	---	---

- we minimise $J(\theta)$ on the training set for a range of d
- choose d that gives the minimum $J_{cv}(\theta^{(d)})$
- estimate generalisation error using your test set i.e. $J_{test}(\theta)$

10.4 Diagnosing bias vs variance



10.5 Choosing the regularisation parameter λ

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

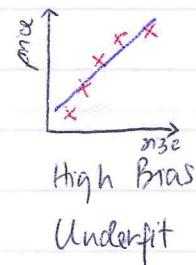
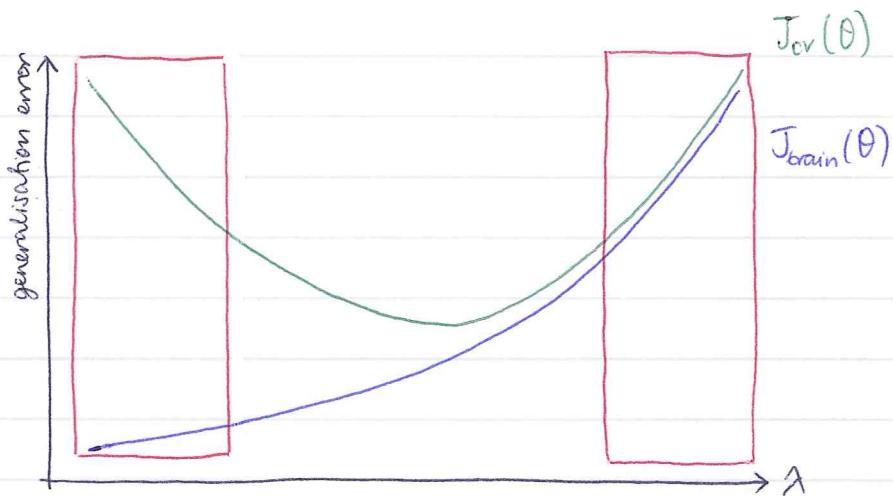
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$

- Try a range of λ and for each minimise $J(\theta)$ on the training set to find $\theta^{(i)}$
- Pick λ such that $J_{\text{cv}}(\theta^{(i)})$ is minimised
- Use the test set to estimate the generalisation error

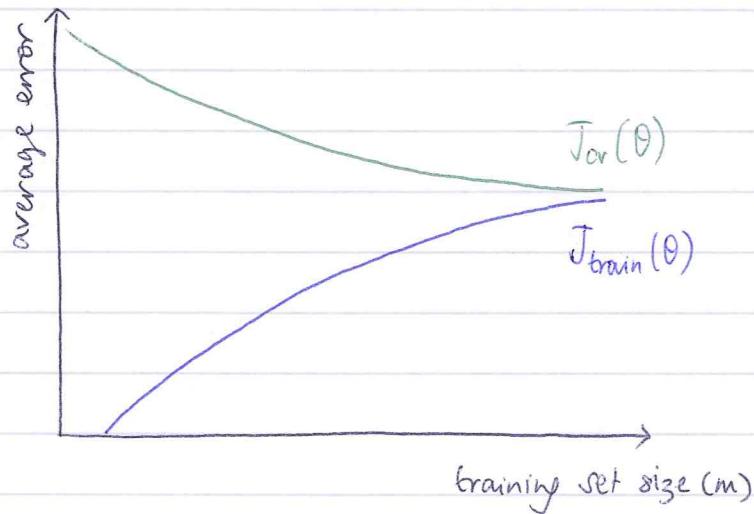
Define $J(\theta)$ as above and...

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

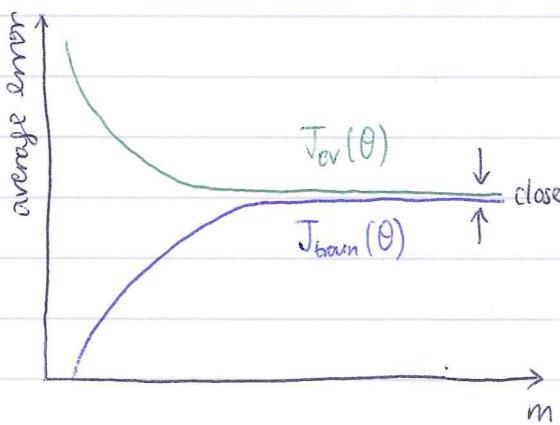
$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_\theta(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$



10.6 Learning Curves

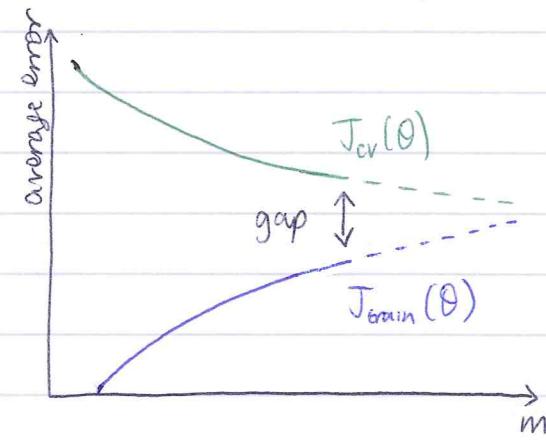


High Bias ↴



Here getting more data will not help

High Variance ↴



Here getting more data will help

10.7 Deciding what to do next revisited ...

- Get more training examples
 - Try a smaller set of features
 - Try increasing λ
 - Try decreasing λ
 - Try getting additional features
 - Try adding polynomial features
- } Fixes high variance
- } Fixes high bias

Neural Networks and Overfitting:

- * "Small" neural networks - fewer parameters
 - prone to underfitting
 - computationally cheaper

- * "Large" neural networks - more parameters
 - prone to overfitting
 - computationally expensive

Regularisation can be used to address overfitting here

- * We can use our cross validation techniques to decide how many units to put in the hidden layer of our neural network (if we have just one hidden layer)

11 Machine Learning System Design

11.1 Prioritising what to work on

Example: Building a spam classifier.

Supervised learning. x = features of email, y = spam(1) / not spam(0)

Features x : choose 100 words indicative of spam/not spam

e.g.

$$x = \begin{pmatrix} 1 \\ 0 \\ 1 \\ 1 \\ \vdots \end{pmatrix} \text{Leena} \quad x \in \mathbb{R}^{100}$$
$$\text{buy} \quad x_j = \begin{cases} 1 & \text{if word } j \text{ appears} \\ 0 & \text{otherwise} \end{cases}$$
$$\text{deal}$$
$$\text{discount}$$

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

How does one best spend their time to make their classifier have a low error?

→ Collect lots of data?

→ e.g. "honeypot" project.

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

→ Develop sophisticated features for message body, e.g. should singular and plural versions be treated as the same word? Punctuation features?

→ Develop sophisticated algorithm to detect misspellings (e.g. mOrtgages, med1cine, wAtches.)

11.2 Error analysis

Recommended approach:

→ Start with a simple algorithm that you can implement quickly. Test it on your cross-validation data

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

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

Example: Spam classifier

$M_{cv} = 500$ examples in the cross-validation set

Suppose the algorithm misclassifies 100 emails.

Manually examine and categorise the errors based on:

(i) what type of email it is

(ii) what features do you think would have helped classify them correctly

e.g. Email Type:		Email Features
Pharma	12	Deliberate Misspellings 5
Replica / fake	14	Unusual routing 16
Phishing	53	Unusual punctuation 32
Other	31	

→ gives an indication of fruitful avenues

The importance of numerical evaluation:

- Should discount / discounts / discounted / discounting be treated as the same word?
- Could use "stemming" software (e.g. "Porter stemmer")
- May have a negative impact by removing differentiation between some words
- Error analysis may not be helpful in this situation.
- Only solution is to try it and see if it works in a quick and dirty implementation
- The cross-validation error of the algorithms performance with and without stemming will allow us to judge its effectiveness
- Single real number evaluation metric

11.3 Error metrics for skewed classes

Cancer classification example

- train logistic regression model $h_0(x) \rightarrow y=1$ if cancer, $y=0$ otherwise
 - suppose you find 1% error on test set (99% correct diagnoses)
 - actually, only 0.5% of patients have cancer
 - this is an example of a SKEWED CLASS i.e. the class occurs on an extreme percentage of the population in this case close to zero
 - in this case it's difficult to judge if our algorithm is performing well by simply looking at our usual error meth
- after all even if we always have $y=0$ our error is just 0.5%!

Precision / Recall Error metrics:

		Actual class	
		1	0
Predicted class	1	True +ve	False +ve
	0	False -ve	True -ve

case: $y=1$ in presence of rare class we wish to detect

Precision: Of all the patients where we predicted $y=1$, what fraction actually has cancer?

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

Recall: Of all the patients which actually have cancer, what fraction did we correctly detect as having it?

$$= \frac{\text{True +ve}}{\# \text{actual +ve}} = \frac{\text{True +ve}}{\text{True +ve} + \text{False -ve}}$$

N.B. • High precision and recall values signify good performance

• Note if we choose $y=0$ always here

11.4 Trading off precision and recall

Logistic regression: $0 \leq h_0(x) \leq 1$

→ predict 1 if $h_0(x) \geq p$? so far we've looked
 0 if $h_0(x) < p$ } at $p=0.5$

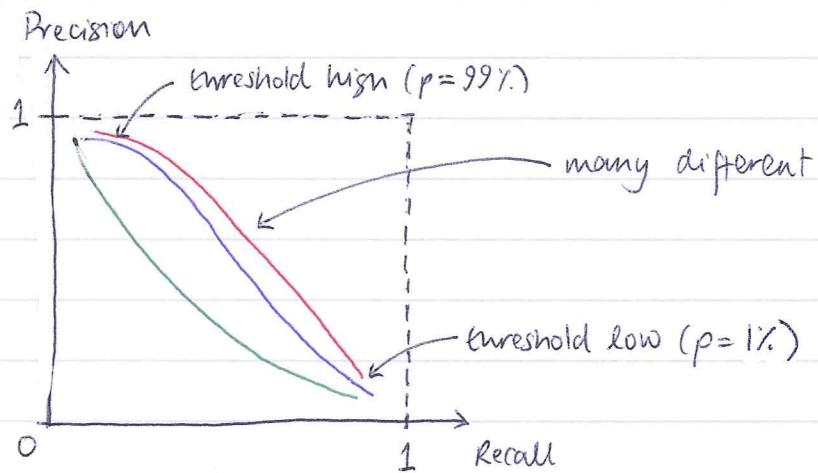
→ suppose we want to predict $y=1$ (cancer) only if we are very confident.

→ we can take a larger p

→ result in higher precision and lower recall

→ suppose we want to avoid missing too many cases

of cancer (avoid false -ves)
 → we can take a lower p
 → results in higher recall, lower precision



F_1 Score (F Score)

- How do we compare precision/recall values for different algorithms
- Ideally we want a single real valued metric

Example

	Precision (P)	Recall (R)	F_1 Score
Algorithm 1	0.5	0.4	0.444
Algorithm 2	0.7	0.1	0.175
Algorithm 3	0.02	1.0	0.0392

- Average is a bad idea as it's likely to have a reasonably consistent value (since the two have a trade off for different values of the threshold)
- F_1 score or F score = $\frac{2PR}{P+R}$

11.5 Data for machine learning

Designing a high accuracy learning system

Studies have shown that inferior classification algorithms will often outperform better algorithms given more data

"It's not who has the best algorithm that wins.
It's who has the most data."

When is this true?

Large data set rationale:

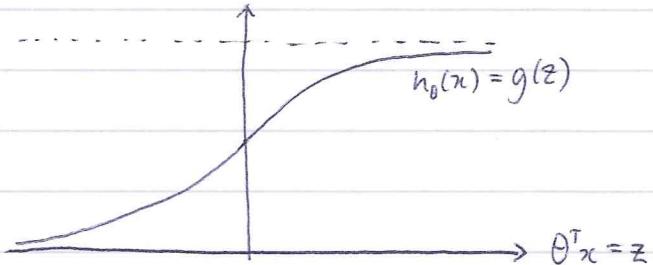
- Assume features $x \in \mathbb{R}^{n+1}$ has sufficient information to predict y accurately
- A useful test is to see if you could classify it with the given data (or a human expert if it's more complicated than classifying spam)
- Use a learning algorithm with many parameters (e.g. logistic regression / linear regression with many features; neural network with many hidden units).
 → LOW BIAS ALGORITHM $J_{\text{train}}(\theta)$ small
- Use a very large training set (unlikely to overfit)
 → LOW VARIANCE ALGORITHM $J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta)$
- LOW BIAS and LOW VARIANCE
 ⇒ $J_{\text{test}}(\theta)$ small

12 Support Vector Machines

12.1 Optimisation objective

Alternative view of logistic regression

$$h_0(x) = \frac{1}{1 + e^{-\theta^T x}}$$



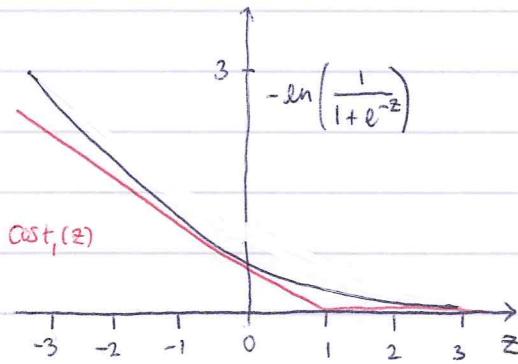
if $y=1$ need $h_0(x)=1$
and $\theta^T x \gg 0$

if $y=0$ need $h_0(x)=0$
and $\theta^T x \ll 0$

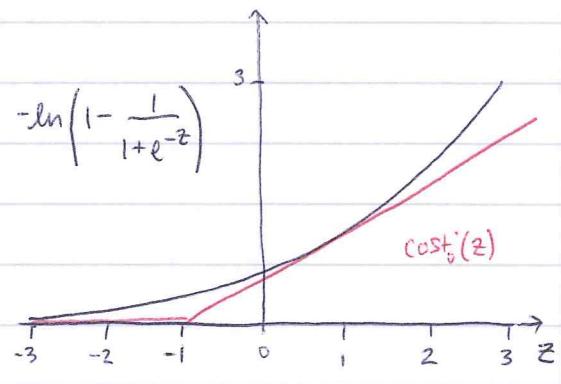
Cost of example : $- (y \ln(h_0(x)) + (1-y) \ln(1-h_0(x)))$

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

if $y=1$ want $\theta^T x \gg 0$



if $y=0$ want $\theta^T x \ll 0$



Logistic Regression :

$$\min_{\theta} \left\{ \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} (-\ln(h_0(x^{(i)})) + (1-y^{(i)}) \ln(1-h_0(x^{(i)}))) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \right\}$$

Support vector machines

$$\min_{\theta} \left\{ C \sum_{i=1}^m [y^{(i)} cost_1(\theta^T x^{(i)}) + (1-y^{(i)}) cost_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2 \right\}$$

regularisation parameter

SVM hypothesis:

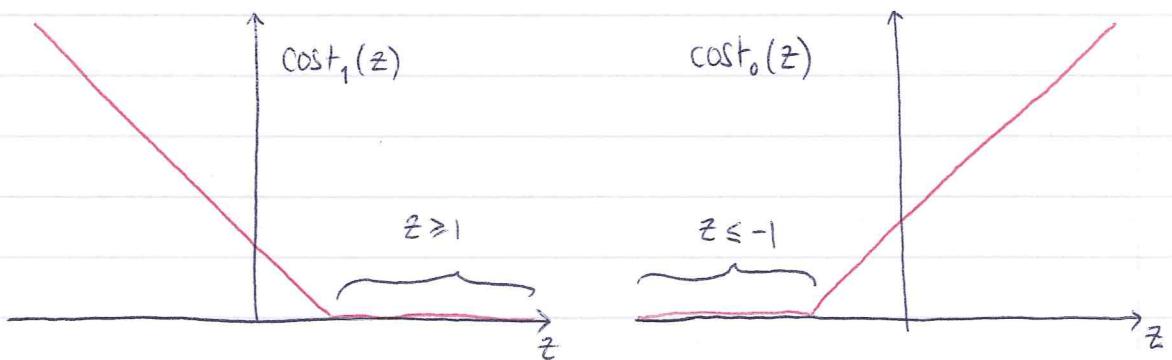
We don't get a probability like with logistic regression:

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

12.2 Large Margin Intuition

Support Vector Machine

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

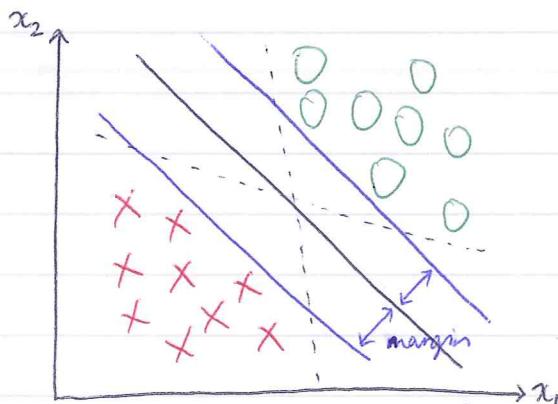


if $y=1$ need $\theta^T x \geq 1$
(not just $\theta^T x \geq 0$)

if $y=0$ need $\theta^T x \leq -1$
(not just $\theta^T x < 0$)

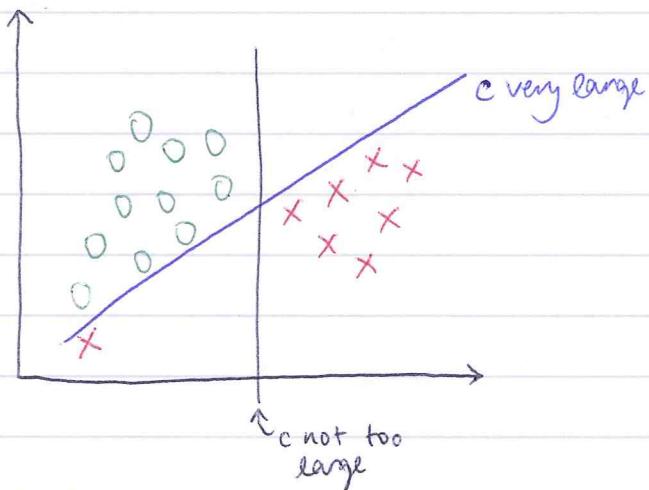
Large margin classifier : take C large

Linearly separable case :



Here there are many linear separators - SVM chooses the one with a big margin

Large margin classifier in the presence of outliers



12.3 The mathematics behind large margin classification

$$\begin{aligned}
 \text{Recall: } \theta^T x &= \theta \cdot x \\
 &= \|\theta\| \|x\| \cos \varphi \quad \leftarrow \varphi = \text{angle between vectors } \theta \text{ and } x \\
 &= (\text{length of projection of } x \text{ onto } \theta) \|\theta\| = p \|\theta\|
 \end{aligned}$$

SVM decision boundary:

$$\min_{\theta} \underbrace{\frac{1}{2} \sum_{j=1}^n \theta_j^2}_{= \|\theta\|^2/2} \quad \text{s.t.} \quad \begin{cases} p^{(i)} \|\theta\| \geq 1 & \text{if } y^{(i)} = 1 \\ p^{(i)} \|\theta\| \leq -1 & \text{if } y^{(i)} = 0 \end{cases}$$

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

simplification: $\theta_0 = 0$

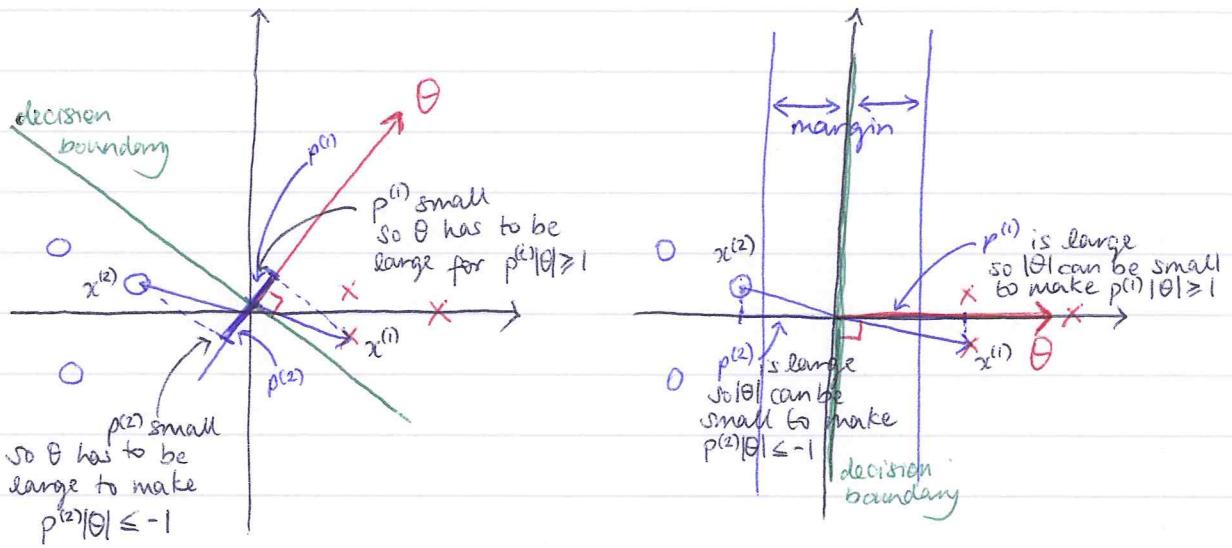
→ in this case the decision boundary looks like

$$\theta_1 x_1 + \theta_2 x_2 = 0$$

→ note that θ is normal to the decision boundary

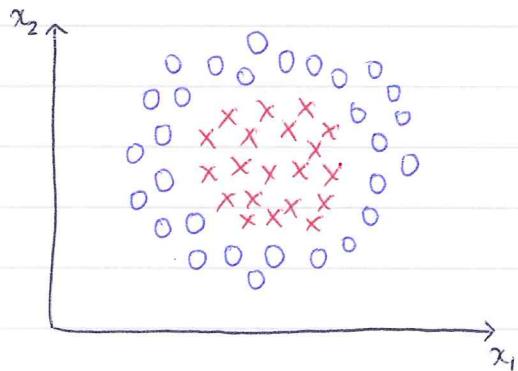
The algorithm tries to enforce } best way to do this is

$$\left. \begin{array}{l} p^{(i)} \|\theta\| \geq 1 \text{ if } y^{(i)} = 1 \\ p^{(i)} \|\theta\| \leq -1 \text{ if } y^{(i)} = 0 \end{array} \right\}$$
 to maximise the $p^{(i)}$ so
 while minimising $\frac{1}{2} \|\theta\|^2$ that θ can be small



12.4 Kernels I

Non-linear decision boundaries



in logistic regression we would add features

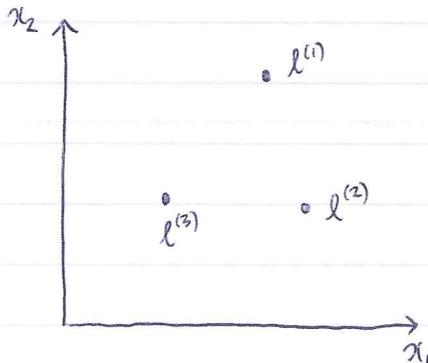
$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2$$

let us write this as

$$h_{\theta}(x) = \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \theta_4 f_4 + \theta_5 f_5$$

$$f_1 = x_1, f_2 = x_2, f_3 = x_1 x_2, f_4 = x_1^2 \text{ and } f_5 = x_2^2$$

Is there a different, better choice of features f_1, f_2, \dots



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

Kernels & similarity

Given x : $f_1 = \text{similarity}(x, l^{(1)}) = \exp(-|x - l^{(1)}|^2 / 2\sigma^2)$

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

$$f_3 = \text{similarity}(x, l^{(3)}) = \exp(-|x - l^{(3)}|^2 / 2\sigma^2)$$

$\underbrace{\hspace{3cm}}$ kernel

$\underbrace{\hspace{3cm}}$ Gaussian Kernel

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

If x_i is close to $l^{(i)}$ then $f_i \approx 1$

If x_i is far from $l^{(i)}$ then $f_i \approx 0$

Predict $y=1$ when $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$

12.5 Kernels II

Choosing the landmarks:

given m training examples

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

choose

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

then

$$\left. \begin{array}{l} f_1 = \text{similarity}(x^{(i)}, l^{(1)}) \\ f_2 = \text{similarity}(x^{(i)}, l^{(2)}) \\ \vdots \\ f_m = \text{similarity}(x^{(i)}, l^{(m)}) \end{array} \right\} f = \begin{pmatrix} f_0 \\ f_1 \\ \vdots \\ f_m \end{pmatrix} \quad f_0 = 1$$

given x , compute $f \in \mathbb{R}^{m+1}$

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

Training :

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

$$\theta = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_m \end{pmatrix} \quad \begin{matrix} \text{implementation} \\ \sum_{j=1}^m \theta_j^2 = \theta^T \theta = |\theta|^2 \end{matrix}$$

$$\theta^T M \theta$$

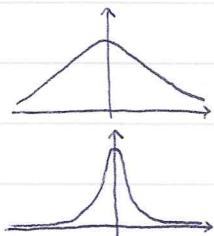
depends on the kernel

SVM parameters

C ($= \frac{1}{\lambda}$) Large C : Lower Bias, High Variance
Small C : High Bias, Low Variance

σ^2 Large σ^2 : Features f_i vary smoothly
High Bias, Low Variance

Small σ^2 : Features f_i vary abruptly
Low Bias, High Variance



12.6 Using an SVM

Use SVM software package (e.g. liblinear, libsvm,...) to solve for parameters θ .

Need to specify

→ Choice of parameter C

→ Choice of kernel (similarity function)

Convention : No kernel ("linear kernel") } suitable when
i.e. predict " $y=1$ " if $\theta^T x \geq 0$ } n large, m small

Gaussian & linear kernels are the most common

Gaussian kernel:

- $f_i = \exp\left(-\frac{|x - l^{(i)}|^2}{2\sigma^2}\right)$, where $l^{(i)} = x^{(i)}$
 - need to choose σ^2
 - can use this to compute complex non-linear decision boundary
- } suitable when
n small and
m large

Some packages require the Kernel (similarity) function to be implemented by the user. E.g.

```
function f = kernel(x1, x2)
f = exp(-\frac{|x1 - x2|^2}{2\sigma^2});
```

return

N.B. Do perform feature scaling before using the Gaussian kernel

Other choices of kernel

Note: Not all similarity functions make valid kernels.

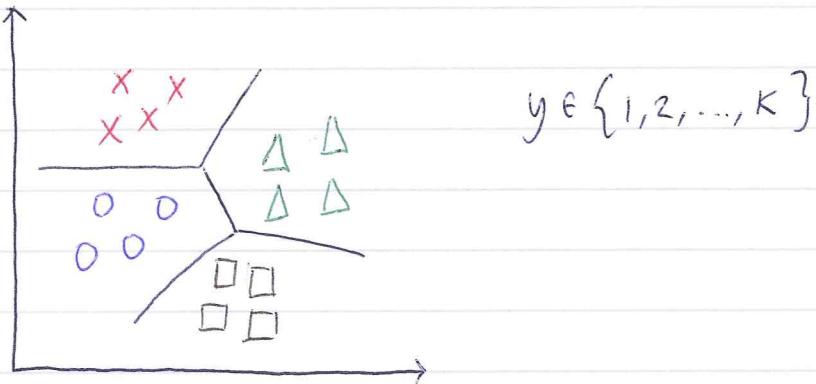
For a function to be a valid kernel it must satisfy "Mercer's Theorem". This ensures an SVM package's optimisers run correctly - doesn't diverge

Many off the shelf kernels are available

- Polynomial Kernel: $k(x, l) = (x^\top l + c)^d$ for some constant c and d
 d is called the degree

- More esoteric: String kernel, chi-square kernel, histogram intersection kernel, ...

Multiclass classification



Many SVM packages already have built-in multi-class classification functionality

Otherwise one-vs-all method works

- train K SVMs, one to distinguish $y=i$ from the rest for $i=1, 2, \dots, K$
- get $\Theta^{(1)}, \Theta^{(2)}, \dots, \Theta^{(K)}$
- for new example predict $y=i$ such that
$$\max_i \{(\Theta^{(i)})^T x\}$$

Logistic Regression vs SVM

n = number of features ($x \in \mathbb{R}^{n+1}$)

m = number of training examples

if n is large relative to m : $m \leq n$

→ use logistic regression or SVM with linear kernel

If n is small, m is intermediate: $m \geq 10n$

→ use Gaussian Kernel

If n is small, m is large: $m \geq 50n$

→ create/add more features, then use logistic regression or SVM without a kernel

* Neural network is likely to work well for all of the above but may be slower to train

B Clustering

Unsupervised learning - given data is unlabelled and the algorithm looks for structure in the data

13.2 K-means Algorithm

Input: $\rightarrow K$ (number of clusters)
 $\rightarrow \text{Training } \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \quad x^{(i)} \in \mathbb{R}^n$

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

Repeat {

cluster assignment step [for $i = 1 : m$
 $c^{(i)} := \text{Index (from } 1 \rightarrow K \text{) of cluster centroid}$
closest to $x^{(i)}$]

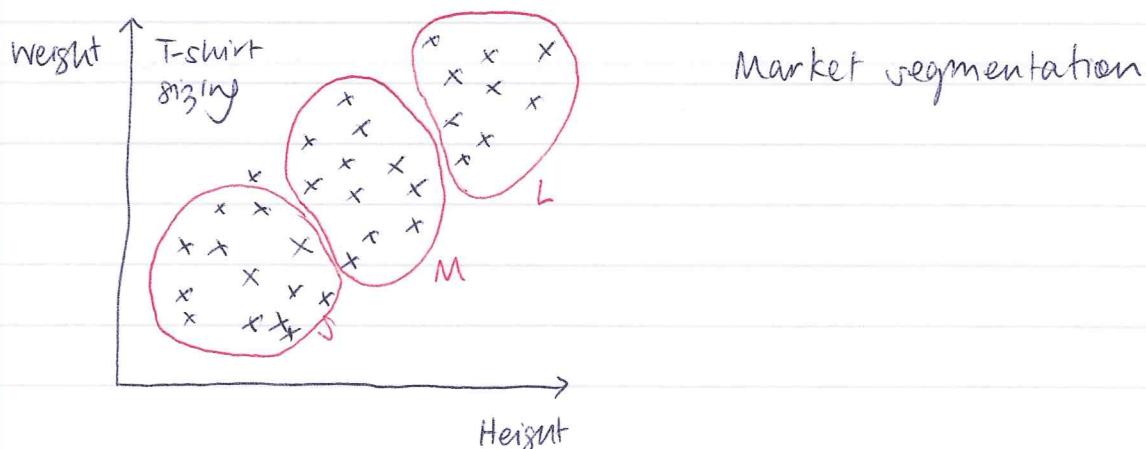
more centroid step [for $k = 1 : K$
 $\mu_k := \text{average (mean) of points assigned}$
to cluster k]

}

If you end up with a cluster with no $x^{(i)}$ assigned to it, you can

- (i) eliminate the cluster
- or (ii) randomly reassign the cluster centre.

K-means for non-separated clusters



13.3 Optimisation Objective

Notation:

$c^{(i)}$ = index of cluster ($1, 2, \dots, K$) to which $x^{(i)}$ is assigned

μ_k = cluster centroid k ($\mu_k \in \mathbb{R}^n$)

$\mu_{c(i)}$ = cluster centroid of cluster to which $x^{(i)}$ is assigned

Optimization objective

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

$$\min_{\substack{c^{(i)} \in \{1, \dots, K\} \\ \mu_k \in \mathbb{R}^n}} \{ J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) \} \quad \text{AKA Distortion}$$

N.B. CLUSTER ASSIGNMENT STEP

→ minimises J over $c^{(i)}$ $i \in [1, m]$

MOVE CENTROID STEP

→ minimises J over μ_k $k \in [1, K]$

13.4 Random Initialisation

Should have $K < m$

Randomly pick K training examples

Set $\mu_1, \mu_2, \dots, \mu_K$ equal to these K examples

→ K-means is vulnerable to LOCAL OPTIMA

for $i = 1 : 100 \{$

Randomly initialize K-means

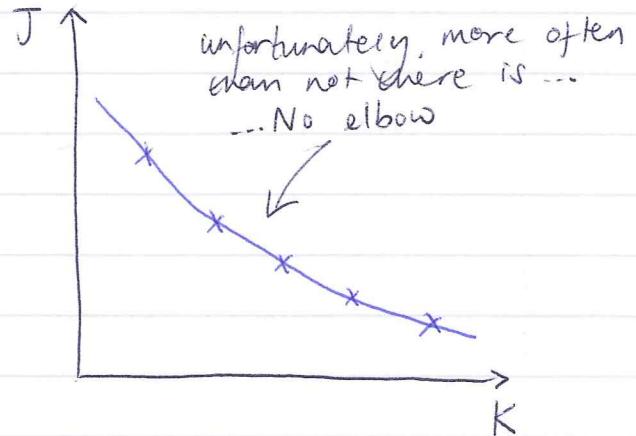
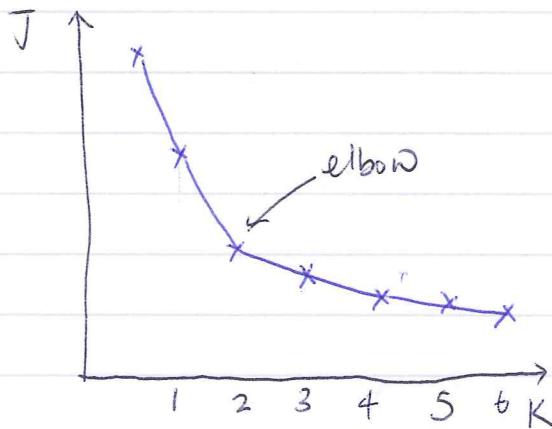
Run K-means. Get $c^{(1)}, c^{(2)}, \dots, c^{(m)}, \mu_1, \mu_2, \dots, \mu_K$

Compute $J(c^{(1)}, c^{(2)}, \dots, c^{(m)}, \mu_1, \mu_2, \dots, \mu_K)$

Pick clustering that gave the lowest J

13.5 Choosing the number of clusters

Elbow method:



N.B. J should always decrease with K

14 Dimensionality Reduction

14.1 Data compression: The data is projected onto the reduced dimension space, i.e. we represent

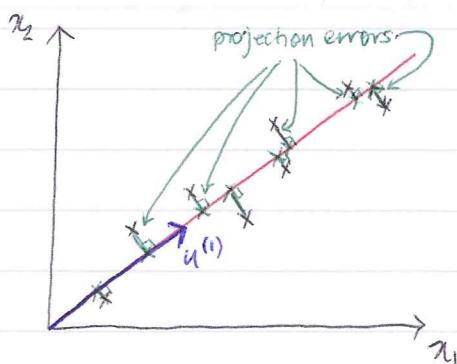
$$\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \rightarrow \{z^{(1)}, z^{(2)}, \dots, z^{(m)}\}$$

$$x^{(i)} \in \mathbb{R}^n \rightarrow z^{(i)} \in \mathbb{R}^k \quad k < n$$

14.2 Data visualization: Dimension reduction to size two/three can allow us to better visualize the data

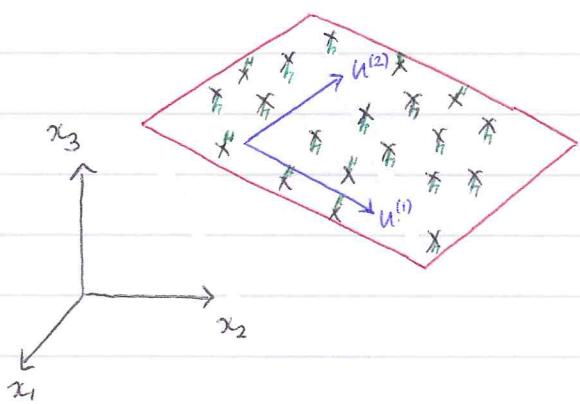
14.3 Principal Component Analysis (PCA)

Problem Formulation: $2D \rightarrow 1D$ example ($n=2$)



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

$3D \rightarrow 2D$ example ($n=3$)

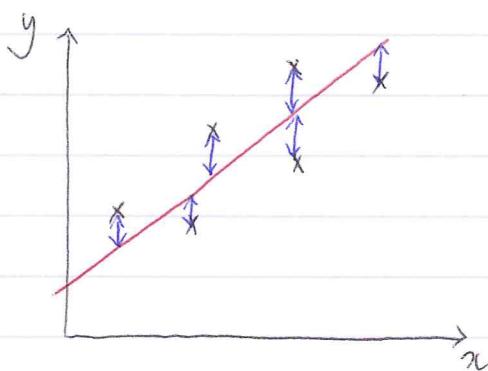


Here we find two direction vectors $u^{(1)}$ and $u^{(2)}$ onto which to project the data so as to minimise the projection error

In general: Reduce from n dimension to k dimension
 Find a set of K vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project the data so as to minimise the projection error

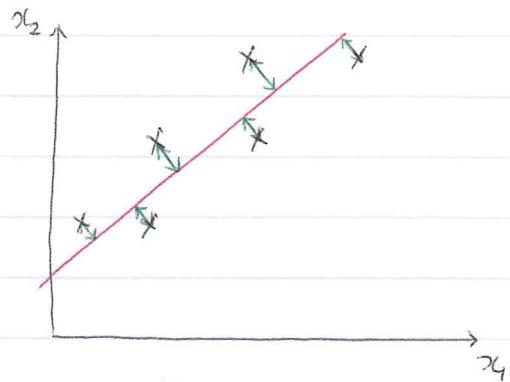
PCA is not linear regression

Linear Regression



y is viewed to be dependent on x i.e. $x^{(i)} \rightarrow y^{(i)}$
 $x^{(i)}, y^{(i)} \in \mathbb{R}^n$

Principal Component Analysis



$x^{(i)} = (x_1, x_2, x_3, \dots, x_n) \in \mathbb{R}^n$
 All features are treated equally

14.4 Principal Component Analysis Algorithm

Data preprocessing

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

Preprocessing (feature scaling / mean normalization):

$$M_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad \text{replace each } x_j^{(i)} \text{ with } x_j^{(i)} - M_j$$

If features are on different scales they should be rescaled to have comparable ranges.

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - M_j}{\sigma_j} \quad \sigma_j \leftarrow \text{range / standard deviation}$$

Algorithm.

We need to find two things:

$$\textcircled{1} \quad z^{(i)} \in \mathbb{R}^k \leftarrow x^{(i)} \in \mathbb{R}^n$$

$$\textcircled{2} \quad u^{(1)}, u^{(2)}, \dots, u^{(k)} \in \mathbb{R}^n \text{ onto which we project our } x^{(i)} \text{ to find our } z^{(i)}$$

To reduce data from n-dimensions to k-dimensions, Compute "covariance matrix":

$$\underbrace{\Sigma}_{n \times n} = \frac{1}{m} \sum_{i=1}^m \underbrace{(x^{(i)})}_{n \times 1} \underbrace{(x^{(i)})^\top}_{1 \times n} \quad = \text{Sigma}$$

Compute "eigenvectors" of the covariance matrix Σ :
 $[U, S, V] = \text{svd}(\text{Sigma})$;

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

take the first k columns

i.e.

$$U_{\text{reduce}} = \begin{bmatrix} \uparrow & \uparrow & \uparrow \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ \downarrow & \downarrow & \downarrow \end{bmatrix} \in \mathbb{R}^{n \times k}$$

$$\mathbf{z} = \underbrace{U_{\text{reduce}}}_{\substack{\uparrow \\ k \times 1}}^T \underbrace{X}_{\substack{\uparrow \\ k \times n}} \underbrace{X}_{\substack{\uparrow \\ n \times 1}}$$

Summary

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

$$\text{Sigma} = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)}) (\mathbf{x}^{(i)})^T = \mathbf{X}' * \mathbf{X} / m ;$$

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

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

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

$\in \mathbb{R}^n$

N.B. as usual

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^{(1)} & \mathbf{x}_2^{(1)} & \dots & \mathbf{x}_n^{(1)} \\ \mathbf{x}_1^{(2)} & \mathbf{x}_2^{(2)} & \dots & \mathbf{x}_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_1^{(m)} & \mathbf{x}_2^{(m)} & \dots & \mathbf{x}_n^{(m)} \end{bmatrix}$$

14.5 Choosing the number of principal components

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

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

Typically choose k to be the smallest value such that

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

i.e. such that "99% of the variance is retained".

Naively we would do the following ...

for $k = 1 : n$

Compute Ureduce, $Z^{(1)}, Z^{(2)}, \dots, Z^{(m)}$, $x_{\text{approx}}^{(1)}, x_{\text{approx}}^{(2)}, \dots, x_{\text{approx}}^{(m)}$

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

If not we choose $k-1$
end;

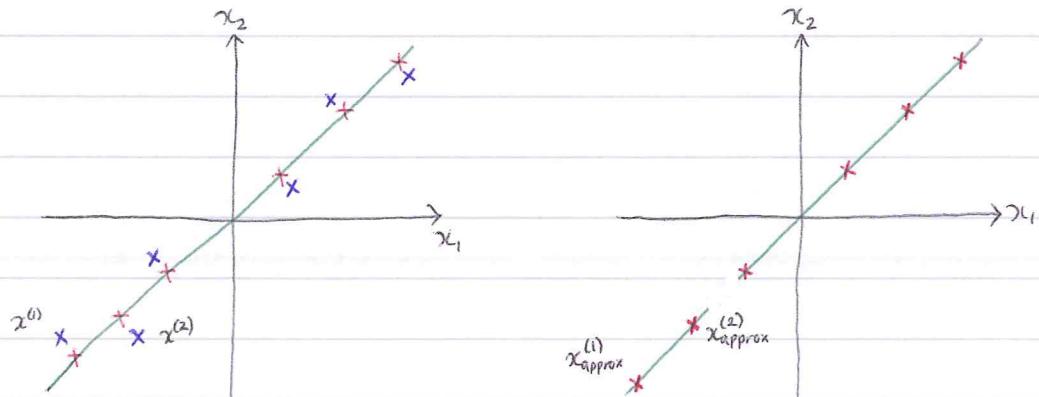
But actually our method uses

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

$$S \text{ is a diagonal matrix} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{For given } k \quad \begin{aligned} 1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} &= \\ \rightarrow \frac{\sum_{i=1}^n S_{ii}}{\sum_{i=1}^n S_{ii}} &\geq 0.99 \end{aligned}$$

i.e. pick the smallest value of k such that

14.6 Reconstruction from Compressed Representation



$$z = \text{Ureduce } x$$

$$* * * + * * \rightarrow z_1$$

$$z \in \mathbb{R} \rightarrow x \in \mathbb{R}^2$$

$$\underbrace{x_{\text{approx}}^{(i)}}_{\in \mathbb{R}^n} = \underbrace{\text{Ureduce}}_{n \times K} \cdot \underbrace{z^{(i)}}_{K \times 1}$$

14.7 Advice for applying PCA

Supervised learning speedup

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

Extract Input:

Unlabelled data set: $x^{(1)}, x^{(2)}, \dots, x^{(n)} \in \mathbb{R}^n$ n large
 \downarrow PCA

$$z^{(1)}, z^{(2)}, \dots, z^{(m)} \in \mathbb{R}^k \quad k < n$$

New training set:

$$(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}), \dots, (z^{(m)}, y^{(m)})$$
$$h_{\theta}(z) = \frac{1}{1 + e^{-\theta^T z}}$$

New example x
 $x \rightarrow z \rightarrow h_{\theta}(z)$
for prediction

Note: Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA only on the training set. h_{θ} is like a parameter.

$n = 10k$ is not unrealistic

Application:

→ Compression

 ↳ Reduce memory needed to store data

 ↳ Speed up learning algorithm

 ↳ choose k by % of variance retained

→ Visualization

 • $k=2$ or 3

Bad use of PCA: To prevent overfitting (might work but bad!)

→ PCA reduces the number of features from n to k

→ This may work but is not a good way to address overfitting as it does not take into account y values

→ Better to just use regularisation

→ PCA should really only be used to speed up your learning algorithm

→ Try doing ML without PCA first and only if the performance is slow implement PCA

15 Anomaly Detection

15.1 Problem Motivation

Example: Fraud detection

→ $x^{(i)}$ = features of user i 's activities

→ Model $p(x)$ from data

→ Identify unusual users by checking which have $p(x) < \epsilon$

Example: Monitoring computers in a data centre

Example: Manufacturing - identifying faulty produce

15.2 Gaussian distribution

$$p(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

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

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

15.3 Algorithm

Density estimation

Training set: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ each $x^{(i)} \in \mathbb{R}^n$

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

Anomaly detection algorithm

1. Choose features x_i that you think might be indicative of anomalous examples

2. Fit parameters $\mu_1, \mu_2, \dots, \mu_n, \sigma_1^2, \sigma_2^2, \dots, \sigma_n^2$

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

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

3. Given a new example x compute $p(x)$:

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

Anomaly if $p(x) < \epsilon$

15.4 Developing and evaluating an anomaly detection system

Real number evaluation metric

\rightarrow When developing a learning algorithm (choosing features, etc.), making decisions is much easier if we have a way of evaluating our learning algorithm

\rightarrow Assume we have some labelled data of anomalous and non-anomalous examples ($y=1$ $y=0$ respectively)
i.e. some of our data is labelled

\rightarrow Assume Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ are non-anomalous examples (exclude examples known to be anomalous)

\rightarrow Cross validation set: $(x_{cv}^{(1)}, y_{cv}^{(1)}), (x_{cv}^{(2)}, y_{cv}^{(2)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

\rightarrow Test set: $(x_{test}^{(1)}, y_{test}^{(1)}), (x_{test}^{(2)}, y_{test}^{(2)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

Example: 10000 good and 20 anomalous

Training set: 6000 good

Cross validation set: 2000 good 10 anomalous

Test set: 2000 good 10 anomalous

Possible evaluation metrics

- True +ve, False +ve, True -ve, False -ve
- Precision / Recall
- F₁ score

Can use cross-validation set to choose ϵ

15.5 Anomaly detection vs. supervised learning

Anomaly Detection

vs.

Supervised Learning

- * Very small number of positive example ($y=1$)
 - * Large number of -ve examples
 - * Many different "types" of anomalies - hard for an algorithm to learn what +ve examples look like; future anomalies are not like past ones
 - * Examples
 - Fraud detection
 - Manufacturing
 - Monitoring machines in a data centre
- in some cases*
- * Large number of positive and negative examples
 - * Enough +ve examples for the algorithm to get a sense of what +ve examples look like; future +ve examples likely to be similar to ones in the training set
 - * Examples
 - Email ~ spam detection
 - Weather prediction
 - Cancer classification

15.6 Choosing what features to use

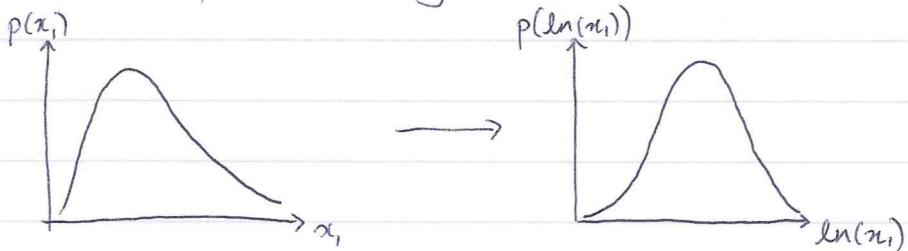
Non-Gaussian features



Plot a histogram

- can use hist function in octave

Examples of transforming data to be more Gaussian



can use other functions too

- $\ln(x + c)$
- x^c

$\text{hist}(x, n)$
number of buckets

Most common problem:

$p(x)$ is comparable for normal and anomalous features

One can try to find functions which will make these stand out rather than

Choose features which take on unusually large or small values in the event of an anomaly

Example: Computer monitoring in a data centre

- | | |
|--|--------------------------------|
| $x_1 = \text{memory use of computer}$ | $x_3 = \text{CPU Load}$ |
| $x_2 = \# \text{ disk accesses/sec}$ | $x_4 = \text{network traffic}$ |
| maybe ... $x_5 = x_3 / x_4$? or $x_6 = x_3^2 / x_4$? | |

15.7 Multivariate Gaussian Distribution

$x \in \mathbb{R}^n$ x_i are not independent

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

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

Recall - diagonal entries give variance

15.8 Anomaly detection using the multivariate Gaussian distribution

As before:

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

Given a new example

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

flag as an anomaly if $p(x) < \epsilon$

Original Model

vs.

Multivariate Gaussian

- * Manually create features to capture where features combined give indicators
- * Computationally cheaper - scales better to large n
- * Okay even if m is small
 $n = 10K, 100K$

- * Automatically captures correlations between features
- * Computationally expensive - requires calculation of Σ^{-1}
- * Must have $m > n$ else Σ is non-invertible. ($m \geq 10n$)

16 Recommender Systems

16.1 Example: Predicting movie ratings

Notation

$$n_u = \# \text{ users}$$

$$n_m = \# \text{ movies}$$

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

$$y^{(i,j)} = \text{rating given by user } j \text{ to movie } i \\ (\text{defined } \Leftrightarrow r(i, j) = 1)$$

Users rate movies with zero-five stars

usually this is one but we use zero for mathematical convenience

movie	users	(1) Alice	(2) Bob	(3) Carol	(4) Dave	Romance	Action
Love at last		1	0	?	5	0.99	0.01
Romance 4ever		0	0	5	4	0.98	0.02
Love puppies		0	?	4	0	0.85	0.05
Car chases		4	5	0	?	0.05	0.95
Big action		5	5	?	0	0.01	0.99

$$\theta^{(j)} = \text{parameter vector for user } j \in \mathbb{R}^{n+1}$$

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

Predicted rating of user j for movie i

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

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

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

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

idea is we want to fill in the gaps so we can make recommendations

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

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

Gradient descent update

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1}^1 ((\theta^{(j)})^\top x^{(i)} - y^{(i,j)}) x_k^{(i)} \quad \text{for } k=0$$

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

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

16.3 Collaborative filtering

Learning features:

* Given $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n_u)}$, to learn $x^{(i)}$:

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

* Given $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n_u)}$ to learn $x^{(1)}, x^{(2)}, \dots, x^{(n_m)}$:

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

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

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

Can do this recursively $\theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \dots$

16.4 Collaborative filtering algorithm

Rather than solving for θ and then x recursively
we do them simultaneously.

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

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

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

Use $x^{(i)} \in \mathbb{R}^n$ and $\theta^{(j)} \in \mathbb{R}^n$ i.e. x_0 and θ_0 not needed

Algorithm

1. Initialize

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

to small random values

2. Minimise $J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$ using gradient descent (or other advanced algorithm). E.g. for every $j=1, \dots, n_u$ and $i=1, \dots, n_m$:

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

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

3. For a user with parameters θ and a movie with (learned) features x , predict a star rating of

$$\theta^T x = (\theta^{(1)})^T x^{(1)}$$

16.5 Vectorization: Low rank matrix factorization

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

○ Finding related movies

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

→ $x_1 = \text{romance}$, $x_2 = \text{action}$, $x_3 = \text{comedy}$, ...

How to find movies j related to movie i

→ small $\|x^{(i)} - x^{(j)}\| \Rightarrow$ movies i and j are "similar"

16.6 Implementation detail: mean normalisation

○ We mean normalise but do not do feature scaling since the features are already on the same scales

Mean normalisation means that our fitted parameters would give a result consistent with predicting the average ratings for a user for which we have no data

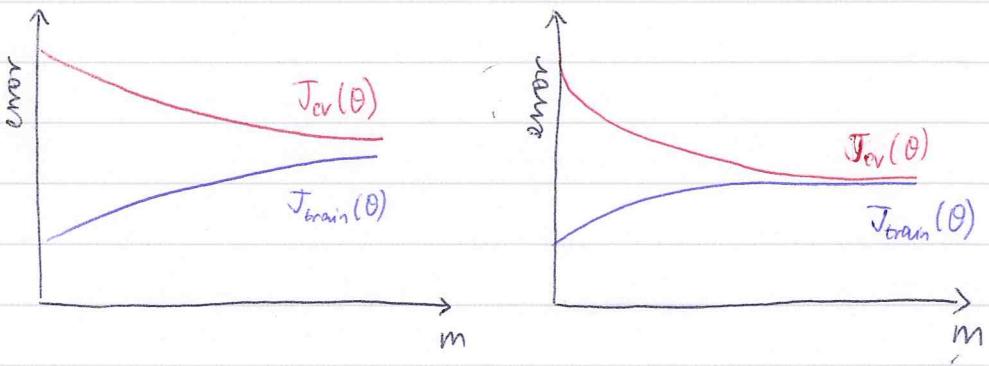
17 Large Scale Machine Learning

17.1 Learning with large data-sets

Recall

"It's not who has the best algorithm who wins. It's who has the most data"

Suppose you have a very large dataset, say $m = 100mn$
it's a good idea to train it on a small randomly chosen subset of the data first to see if more data will help - i.e. $m=1000$ and then plot learning curves of $J_{\text{train}}(\theta)$ and $J_{\text{cv}}(\theta)$ if they indicate high variance for small m



17.2 Stochastic Gradient Descent

When we have large datasets the regression algorithms we have looked at so far become computationally expensive

Recall:

Linear regression with gradient descent

$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j, \quad J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

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

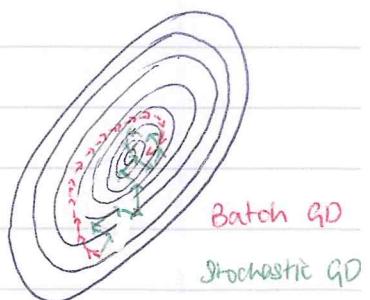
also known as batch gradient descent because we compute the term $\frac{\partial J_{\text{train}}(\theta)}{\partial \theta_j}$ for each $j = 0, 1, \dots, n$

In Stochastic Gradient Descent we instead train on each data pair $(x^{(i)}, y^{(i)})$ at a time i.e. define

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

then

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



1. Randomly shuffle dataset.

2. Repeat }

Unlike batch gradient the algorithm does not monotonically converge to the minimum it can however be much faster for large datasets.

7.3 Mini Batch Gradient Descent

* Batch Gradient Descent - use all m examples in each iteration

* Stochastic Gradient Descent - use one example in each iteration

* Mini-Batch Gradient Descent - use b examples in each iteration

b = mini-batch size typically $2 \leq b \leq 100$

Repeat {

for $i = 1, 1+b, 1+2b, \dots$

$$\theta_j := \theta_j - \frac{\alpha}{b} \sum_{k=i}^{i+b-1} (h_0(x^{(k)}) - y^{(k)}) x_j^{(k)}$$

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

Mini-batch allows us to take advantage of partial vectorization and so can be faster than both Stochastic Gradient Descent and Batch Gradient Descent.

17.4 Stochastic Gradient Descent Convergence

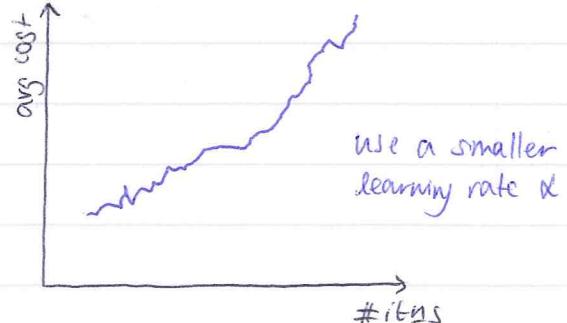
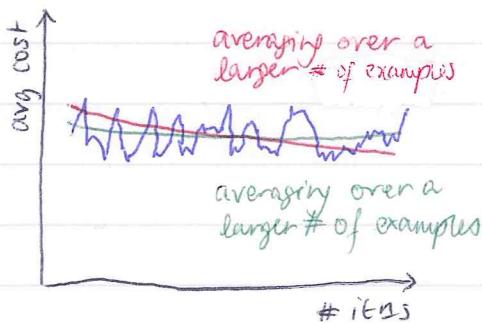
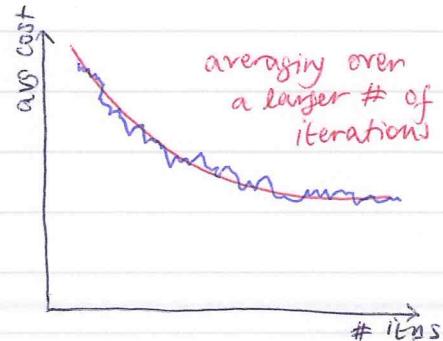
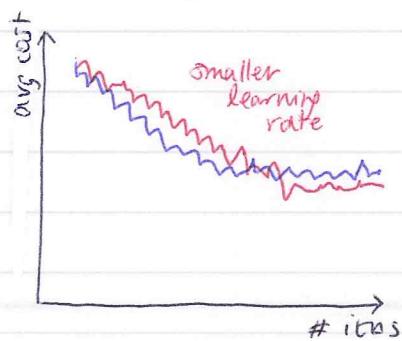
Checking for convergence:

Batch Gradient Descent

- Plot $J_{\text{train}}(\theta)$ as a function of the number of iterations of gradient descent
- $J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Stochastic Gradient Descent

- $\text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$
- During learning compute $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ before updating θ using $(x^{(i)}, y^{(i)})$
- Every 1000 iterations (say), plot $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ averaged over the last 1000 examples processed by the algorithm.



To improve the convergence of the Stochastic Gradient Descent algorithm one can decrease the learning rate as the number of iterations increased.

Typically one might use

$$\alpha = \frac{\text{const 1}}{\text{iteration \#} + \text{const 2}}$$

Of course this gives you an extra two parameters to play with which can be a hassle but for the right choice one might see smoother convergence.

7.5 Online learning

Example: Shipping service where user specifies origin and destination, service quotes a price and then user purchases ($y=1$) or does not purchase ($y=0$) the service.

Feature vector x captures user properties (origin / destination) and our quote and we want to learn $p(y=1|x; \theta)$ to optimise price

Repeat forever {
Get (x, y) corresponding to user
Update θ using (x, y)
 $\rightarrow \theta_j := \theta_j - \alpha (h_\theta(x) - y)x_j \quad j=0, 1, \dots, n$
} each (x, y) is
disregarded once
it is learnt from

Example: Product Search - Click Through Rate (CTR)

Based on the search query we want to know the probability the user will click on a link for a particular product

Show the user the 10 most likely products they will click

- Online learning algorithms can
 - Learn changing user tastes / trends
 - Avoid the need to store huge amounts of data

17.6 Map reduce and data parallelism

Example : Batch Gradient Descent

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

We can split our m training examples into as many equal parts as we have computers - say c . Each computer then calculates the sum

$$\text{temp}_j^{(k)} = \sum_{i=1+(k-1)b}^{kb} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad b = m/c = \text{batch size}$$

and then combine the results at a centralised server

$$\theta_j := \theta_j - \frac{\alpha}{m} \sum_{k=1}^c \text{temp}_j^{(k)} \quad j = 0, 1, \dots, n$$

Map reduce can be used when

- algorithm can be expressed as a sum of functions over the training set.
- works for logistic regression too

Multi-core machines:

- Parallelising over cores rather than separate computers re issue of network latency
- Numerical linear algebra libraries may take advantage of parallelisation automatically

Map reduce can be used to train a neural network.
In this case each core/computer computes forward propagation and back propagation on its portion of the training data followed by the derivative with respect to its portion.

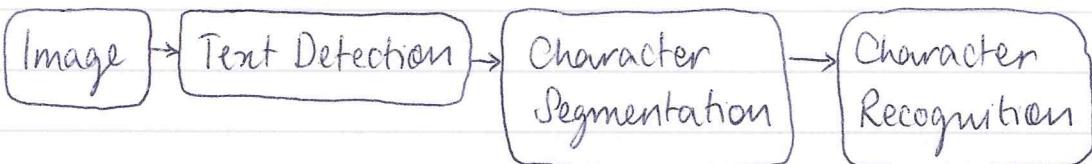
18 Application Example : Photo OCR

18.1 Problem description and pipeline

Photo Optical Character Recognition - Photo OCR

Photo OCR pipeline

1. Text detection
2. Character segmentation
3. Character classification
4. Error correction



18.2 Sliding windows

Learn what text looks like using logistic regression then use a sliding window of a fixed aspect ratio to locate text in the image. Increase the size of the window and repeat. Then use "expansion" on the pixels to identify blocks of text.

Learn what the midpoint between two characters looks like again using logistic regression. Use sliding windows method to perform character segmentation.

Finally we've done character classification using neural networks before!

18.3 Getting Lots of Data and Artificial Data

We've talked about how effective it can be to take a low bias algorithm and train it on masses of data. We can get more data by synthesising it.

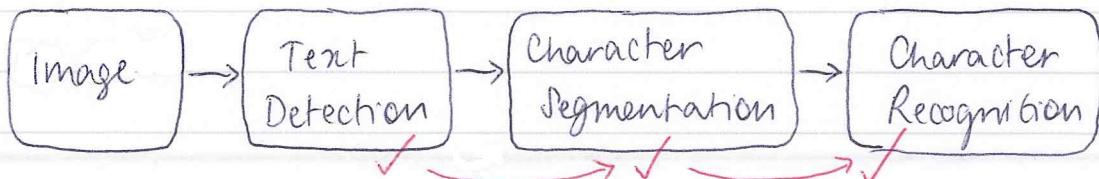
We can synthesis data by taking real data and introducing distortions. In this way we can amplify the size of our data set!

Distortions should be representative of the kinds we would expect to see in our dataset.

Other ways to get data:

- collect and label it yourself
- "crowd sourced" data labelling e.g. "Amazon Mechanical Turk"

18.4 Ceiling Analysis: What part of the pipeline to work on next?



Replace each component in turn (in order) with a perfectly accurate set of results to see the potential gain in overall accuracy

Component	Accuracy
Overall System	72%
Text Detection	89%
Char' Seg'n	90%
Char' Recog'n	100%