

Machine Learning Nanodegree

Contents page

Section	Topic
1	1 Foundations AI 2 Foundations ML 3 Data Science
2	C1 Model Evaluation & Validation P1 Boston House Prices
3	C2 Supervised Learning P2 Student Intervention System
4	C3 Unsupervised Learning P3 Customer Segments
5	C4 Reinforcement Learning P4 Smart Cab

Topic one

Topic two

Topic three

Topic four

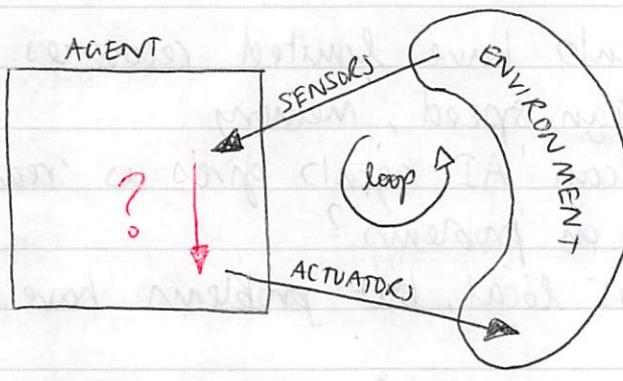
Topic five

Foundations of Artificial Intelligence

* Introduction to Artificial Intelligence

- Outline:
- Agent & Environment
 - Applications of AI
 - AI and Uncertainty

* Intelligent Agent



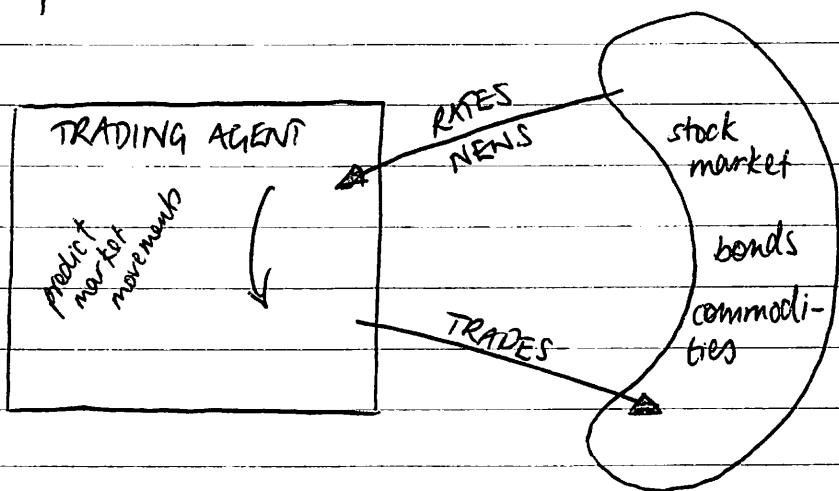
PERCEPTION ACTION CYCLE

* Applications of AI

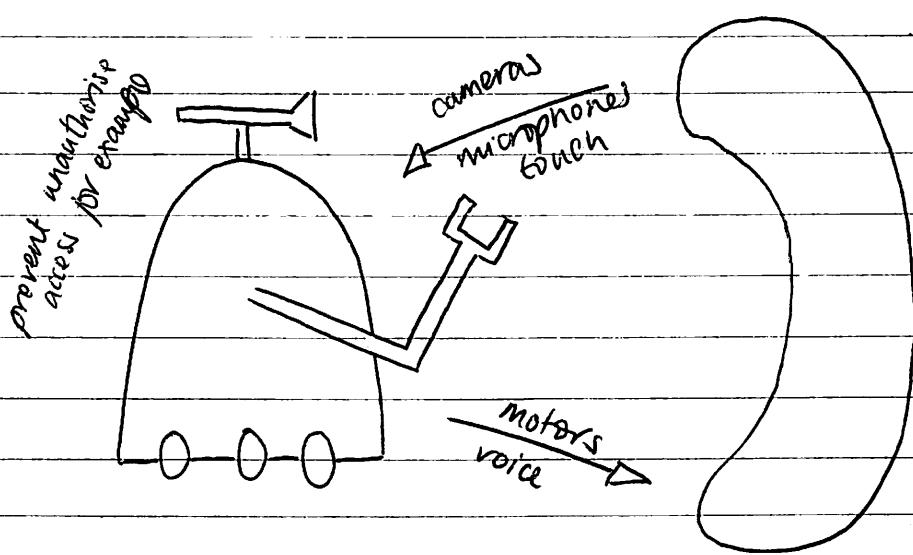
AI has successfully been used in

- Finance
- Robotics
- Games
- Medicine
- The web

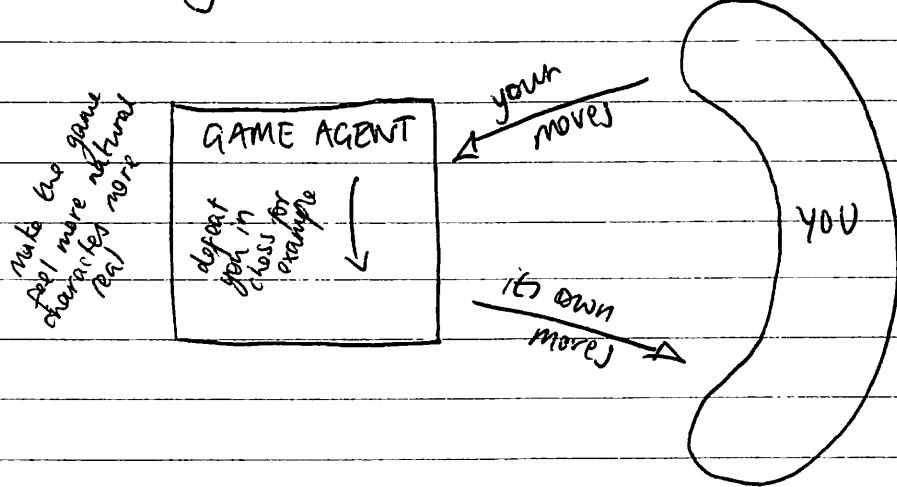
AI in finance:



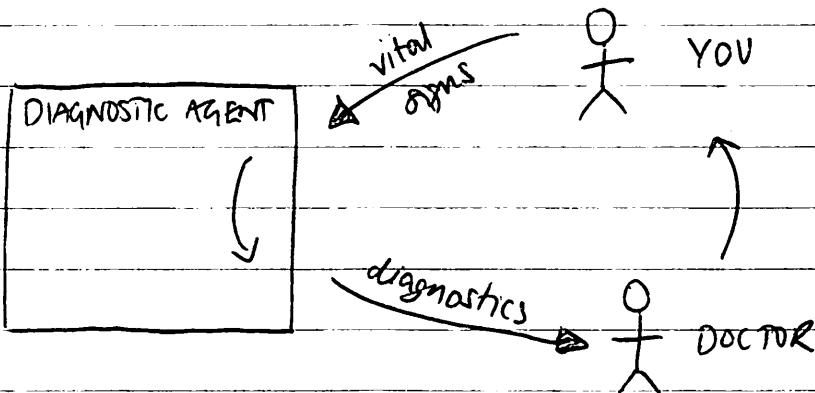
AI in robotics



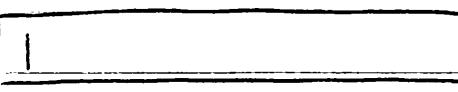
AI in games



AI in medicine

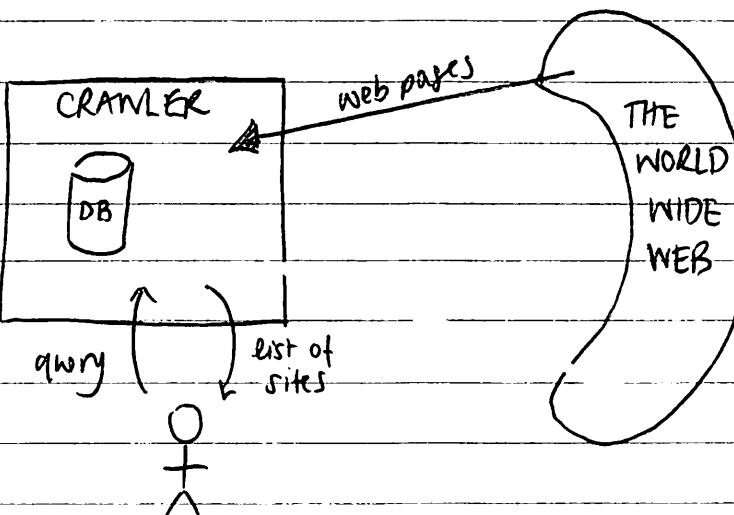


AI and the web



SEARCH **IN FEELING LUCKY**

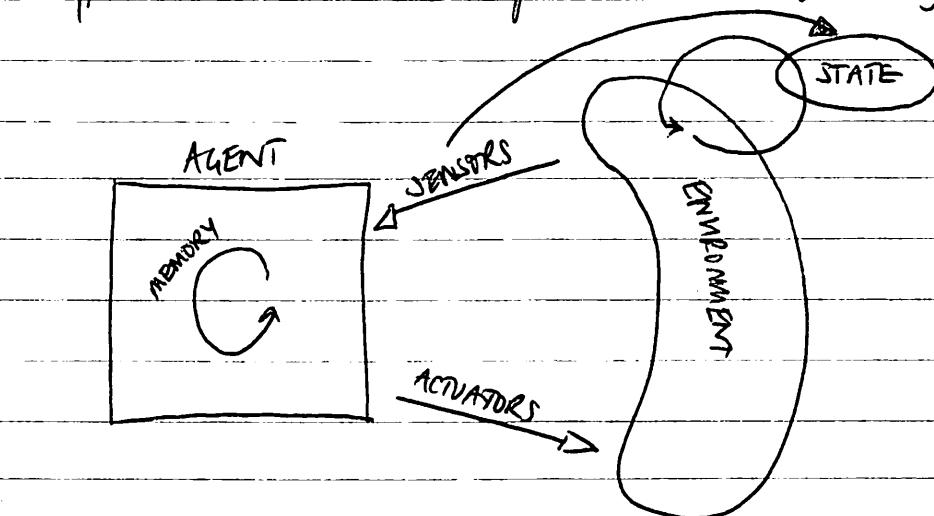
find most relevant web pages for your search



Terminology

* FULLY vs PARTIALLY observable environment

↑
what the agent can sense in any point in time is sufficient to make the optimal choice e.g. seeing all the cards on the table in a game



* DETERMINISTIC vs STOCHASTIC environment

* DISCRETE vs CONTINUOUS environments

* BENIGN vs ADVERSARIAL environments

↑
stochastic
games where there are opponents?

Example:

	PARTIALLY OBSERVABLE	STOCHASTIC	CONTINUOUS	ADVERSARIAL
Checkers				X
Poker	X	X		X
Robot Car	X	X	X	X

momentary sensing is limited

other cars are unpredictable

other drivers

AI as uncertainty management

AI = what to do when you don't know what to do?

Reasons for uncertainty

- Sensor limits
- Adversaries
- Stochastic environments
- Laziness
- Ignorance

AI is a discipline which deals with uncertainty and manages it in decision making

Examples of AI

Language translation e.g. news articles - this is done by examining many articles which are published in both languages. These are used to find the most probable translation.

Example of its use in a Chinese / English menu

Unit 1 summary

- Key application
- Intelligent agents
- Key attributes of environment
- Sources of uncertainty
- Rationality

* Knowledge-Based Artificial Intelligence

- Overview:
- Conundrums in AI
 - Characteristics of AI Problems & Agents
 - AI in practice: Watson
 - What is knowledge-based AI
 - The four schools of AI

* Fundamental Conundrums of Artificial Intelligence

- Intelligent agents have limited resources
 - e.g. processing speed, memory
 - How then can AI agents give us 'real time' performance on problems?
- Computation is local, but problems have global constraints
- Logic is deductive, but many problems are not
 - How can AI agents address inductive and abductive problems?
- The world is dynamic and knowledge is limited
 - How can AI agents address new problems?
- Problem solving, learning and reasoning are complex, explanation and justification even more so
 - How can we get AI Agents to explain or justify its decisions?

* Characteristics of AI Problems

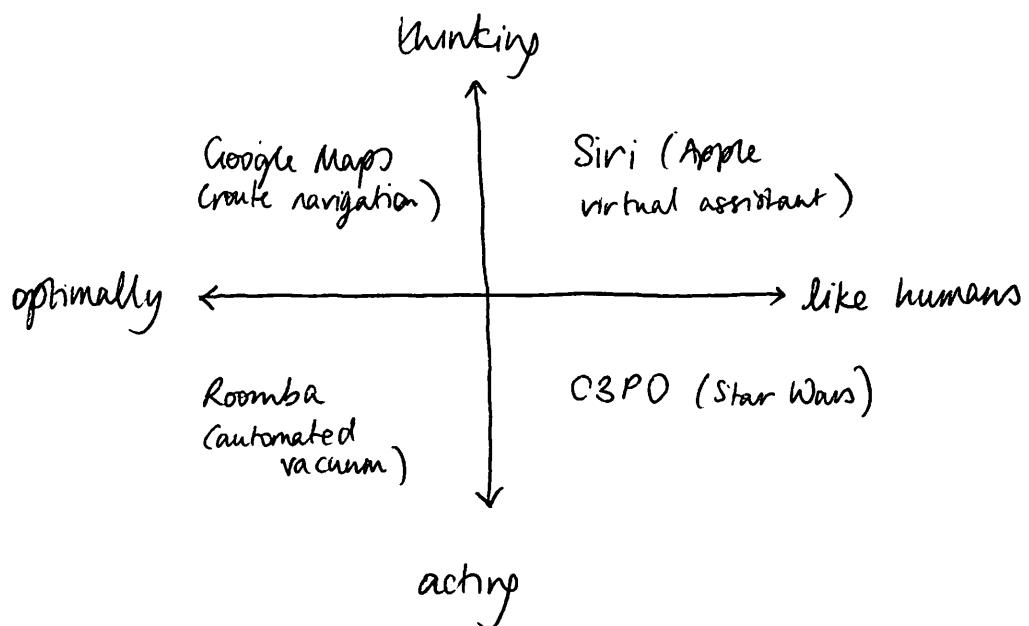
- Knowledge often arrives incrementally
- Problems exhibit recurring patterns
- Problems have multiple levels of granularity
- Many problems are computationally intractable
- The world is dynamic, knowledge is relatively static
- The world is open-ended but knowledge is limited

* Exercise : What is KBAI

Where does an autonomous vehicle fall on the spectrum?

Autonomous vehicles are in the acting optimally quadrant

* Exercise : The four schools of AI



* Problem Solving

- Outline :
- What is a problem?
 - Example: Route finding
 - State spaces

Unit 2 : Problem Solving

Theory and technology of building agents to plan ahead and solve problems

Definition of a problem

- Initial state
- Actions ($s \rightarrow \{a_1, a_2, a_3, \dots\}$) could be state dependent
- Result ($s, a \rightarrow s'$)
- GoalTest ($s \rightarrow \text{true/false}$)
- PathCost ($s \xrightarrow{a} s \xrightarrow{a} s' \dots \rightarrow n$) → n (the cost)
- StepCost ($s, a, s' \rightarrow n$)

Example: Driving from one city to another
convert map of roads to a graph move along edges - you have the frontier, explored and unexplored areas

Function : TREE.SEARCH (problem) :

frontier = {[initial]}

loop :

 if frontier is empty : return FAIL

 path = remove-choice (frontier)

 s = path.end

 if s is a goal : return path

 for a in actions:

 add [path + a \rightarrow Result (s, a)]

 to frontier

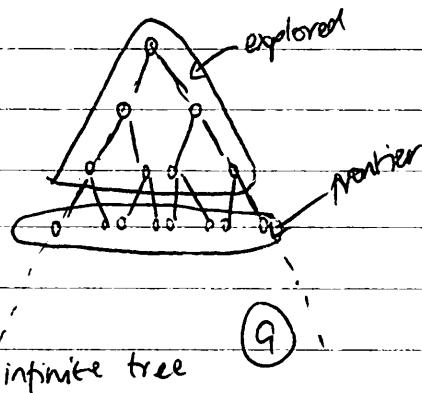
In general the
goal test is applied
when we remove the
path from the frontier
not when it's added to
the frontier

• Breadth first search : explore the shortest paths first (where edges are all 1)

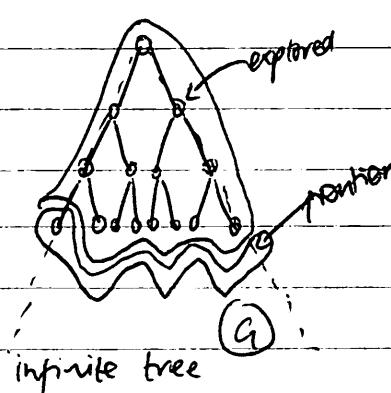
• Cheapest first search : explore the cheapest path first (where the edges have assigned costs)

- Depth first search: explore the longest path first (where edges are all 1)

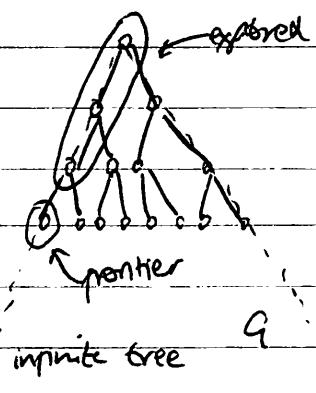
BREADTH FIRST
SEARCH



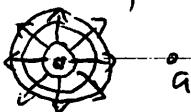
CHEAPEST FIRST
SEARCH



DEPTH FIRST
SEARCH



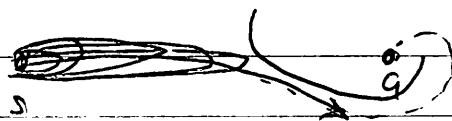
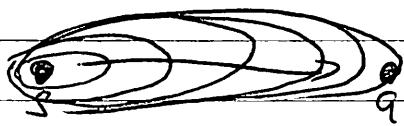
- complete - guaranteed to find the goal if it's at some finite place in an infinite tree
- search area expand uniformly in all directions i.e.



cheaper storage
costs frontier and explored only contain n nodes

- * Greedy-best first search \rightarrow uses estimated distance from goal to decide which direction to expand its search in. This is good but doesn't find the optimal if there are barriers in the search paths

The idea is we want to get the optimal solution faster!



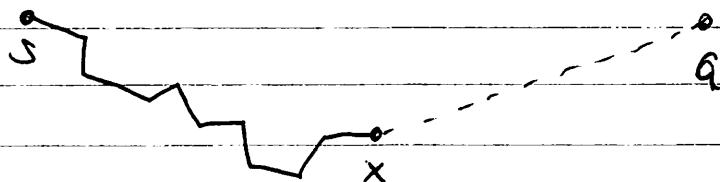
A* search : Best estimated total path cost first

expanding the paths which have the minimum value of the function f

$$f = g + h$$

$g(\text{path})$ = path cost

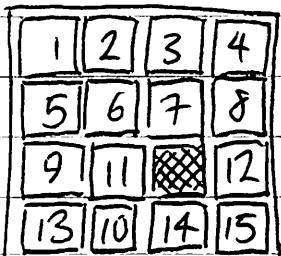
$h(\text{path}) = h(s)$ = estimated distance to goal



A* finds the lowest cost path if :

- $h(s) \leq$ true cost
- i.e. • h never overestimates
- h is optimistic
- h is admissible

* Sliding blocks puzzle / 15 puzzle



heuristics

h_1 = # misplaced blocks

h_2 = sum (distances of blocks)
from correct posn

both these heuristics are admissible

$h_1 \leq$ # moves required fix it

$h_2 \leq$ # " " " "

$h_2 \geq h_1$ so h_2 is a better heuristic

A block can move $A \rightarrow B$
if (A adjacent to B) and (B is blank)

↓

removing this gives h_2

removing both conditions gives h_1

Take $h = \max(h_1, h_2)$

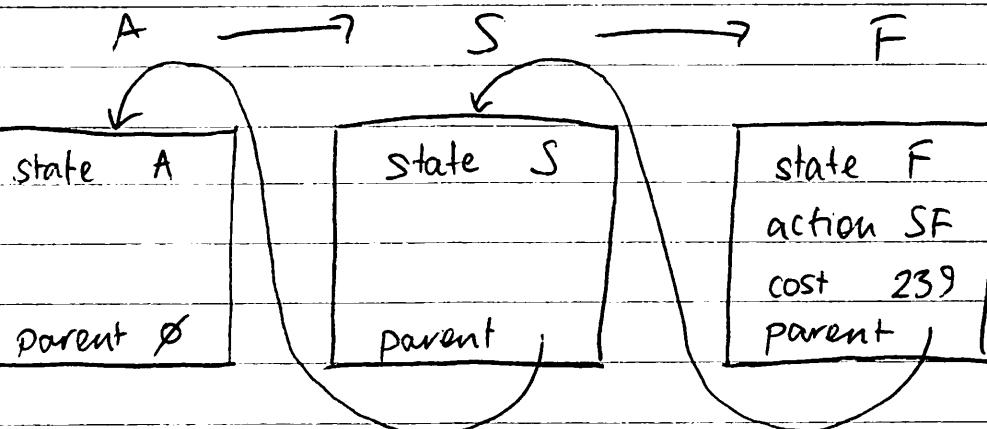
heuristic calculation can have a cost too.

removing conditions only relaxes the problem
so resulting heuristics can only underestimate
the cost

This illustrates how an algorithm can be
devised to come up with the estimated
cost heuristic

Problem-solving technology works when...

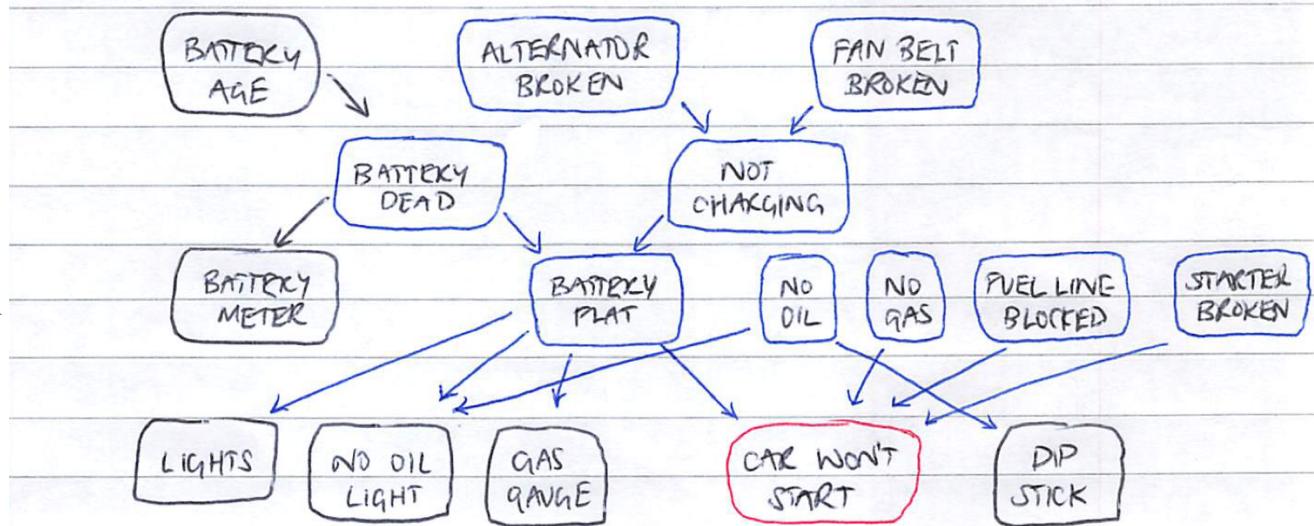
- Fully observable
- Known domain -
- Discrete
- Deterministic
- Static - no external influences on the environment



* Probability in AI

- Outline :
- Example : Bayes Network
 - Conditional Probability
 - Total Probability
 - Joint Probability
 - Bayes Rule

Bayes Networks



- Bayes network assist you in reasoning from observable variables
- Child node is influenced by the Parent in a non-deterministic way

* This course:

- Binary events
- Probability
- Simple Bayes Networks
- Conditional Independence
- Bayesian Networks
- D-separation
- Parameter Count
- Inference

* Bayes Networks are used in

- Diagnosis
- Prediction
- Machine Learning (Finance, Google, Robotics)
- Particle filters
- Hidden Markov Models (HMM)
- MDP / POMDPs
- Kalman Filters

* Probabilities:

$$P(A) = p \quad P(\neg A) = 1-p$$

Independence

$$X \perp Y \Rightarrow \underbrace{P(X)P(Y)}_{\text{marginals}} = \underbrace{P(X,Y)}_{\text{joint}}$$

Dependence

$$P(Y) = \sum_i P(Y|X=i)P(X=i) \quad \text{Total Probability}$$

$$P(\neg X|Y) = 1 - P(X|Y)$$

$$\text{but } P(X|\neg Y) \neq 1 - P(X|Y)$$

* Bayes Rule

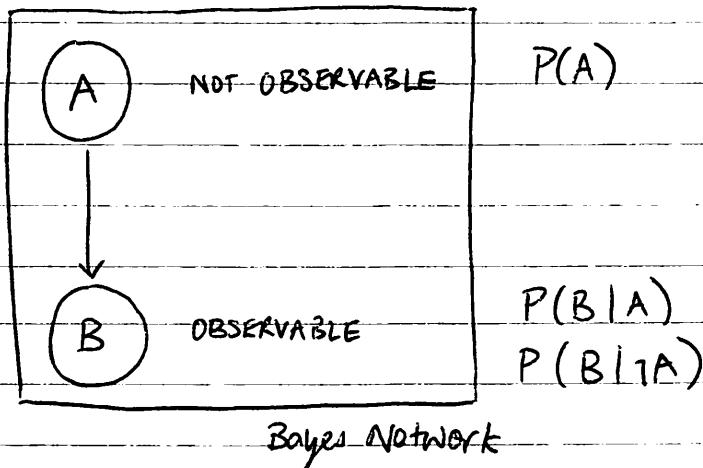
$$P(A|B) = \frac{P(B|A) P(A)}{P(B)} \quad \begin{matrix} \text{LIKELIHOOD} \\ \swarrow \\ \text{POSTERIOR} \end{matrix} \quad \begin{matrix} \leftarrow \text{PRIOR} \\ \searrow \\ \text{MARGINAL LIKELIHOOD} \end{matrix}$$

||

Here B is the evidence
that A has happened or
is the case

$\sum_a P(B|A=a) P(A=a)$

$\underbrace{\qquad\qquad\qquad}_{\text{total probability}}$



Diagnostic Reasoning : $P(A|B)$
 $P(\neg A|\neg B)$

* More complex Bayes Networks

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)} \quad P(\neg A|B) = \frac{P(B|\neg A) P(\neg A)}{P(B)}$$

$$P(A|B) + P(\neg A|B) = 1$$

$$P'(A|B) = P(B|A) P(A)$$

$$P(A|B) = \eta P'(A|B)$$

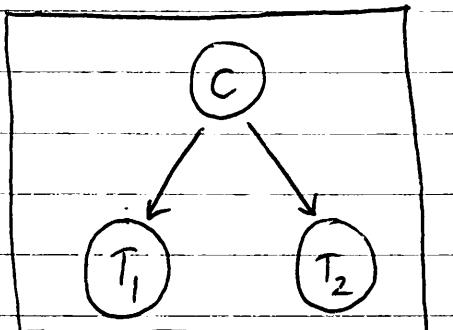
$$P'(\neg A|B) = P(B|\neg A) P(\neg A)$$

$$P(\neg A|B) = \eta P'(\neg A|B)$$

$$\eta = [P'(A|B) + P'(\neg A|B)]^{-1}$$

η is called a 'pseudo probability'

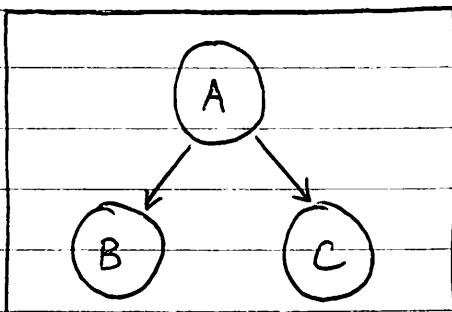
* Conditional Independence



Now let's assume the test for cancer is taken twice

* Conditional Independence

$$P(T_2 | C \cap T_1) = P(T_2 | C)$$



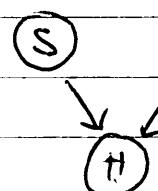
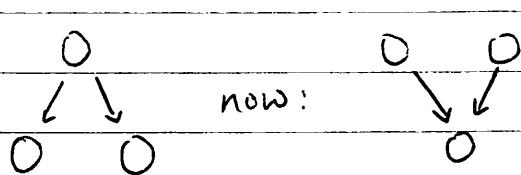
Given A
 $B \perp C$
 i.e. $\underline{B} \perp C | A$

Note: $B \perp C | A \not\Rightarrow B \perp C$

and

$B \perp C \not\Rightarrow B \perp C | A$

Different type of Bayes Network

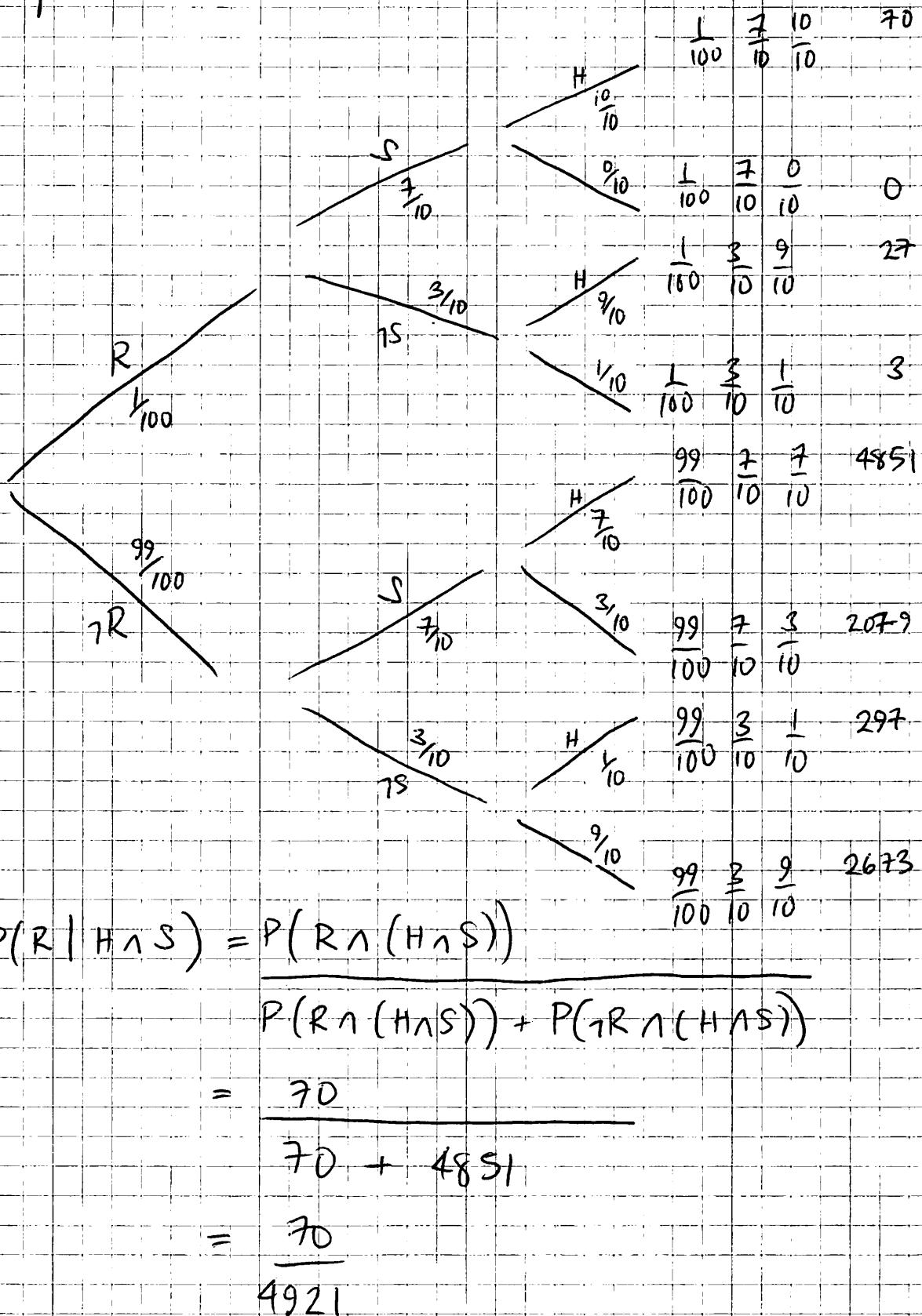


i.e. we can know S|R but S|RH is not
 EXPLAIN THIS CASE

$$\begin{aligned} P(S) &= P(\text{sunny}) = \frac{7}{10} \\ P(R) &= P(\text{raise}) = \frac{5}{100} \\ P(H) &= P(\text{happy}) \end{aligned}$$

$$\left\{ \begin{array}{l} P(H | R \cap S) = 1 \\ P(H | R \cap \neg S) = \frac{2}{10} \\ P(H | \neg R \cap S) = \frac{7}{10} \\ P(H | \neg R \cap \neg S) = \frac{1}{10} \end{array} \right.$$

Example .



$$P(R | H \cap S) = \frac{P(R \cap (H \cap S))}{P(R \cap (H \cap S)) + P(\text{not } R \cap (H \cap S))}$$

$$= \frac{70}{70 + 4851}$$

$$= \frac{70}{4921}$$

$$\text{Bayes Formula} \Rightarrow \frac{P(H \cap S | R) P(R)}{P(H \cap S)}$$

$$= \frac{\frac{70}{100} \times \frac{1}{100}}{\frac{70}{10000} + \frac{4851}{10000}} = \frac{70}{4921} = \frac{1}{703}$$

* Characteristics of AI Agents

- Agents have limited computing power
- Agents have limited sensors
- Agents have limited attention
- Computational logic is fundamentally deductive
- AI agents' knowledge is incomplete relative to the world

* Exercise: What are AI problems?

- Answering questions on Jeopardy
- Configuring the dimensions for the basement of a new house
- Tying Shoelaces
- Deciding on the route to a new destination
- Making sense of a news broadcast
- Designing a robot that walks on water
- Establishing whether a flower pot can be used as a drinking cup
- Deciding whether or not a new animal is a bird

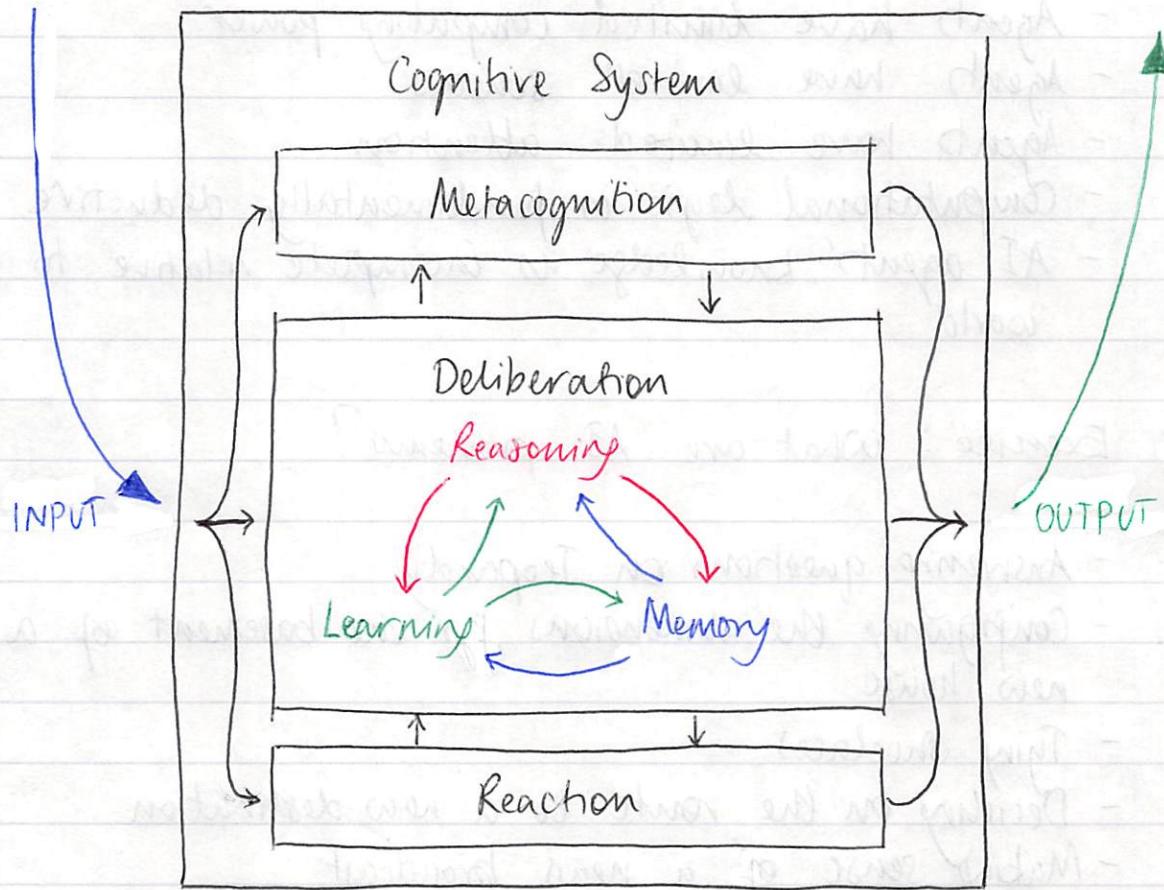
An AI agent could be designed to solve any problem that a human can.

* AI in Practice: Watson (IBM's Jeopardy player)

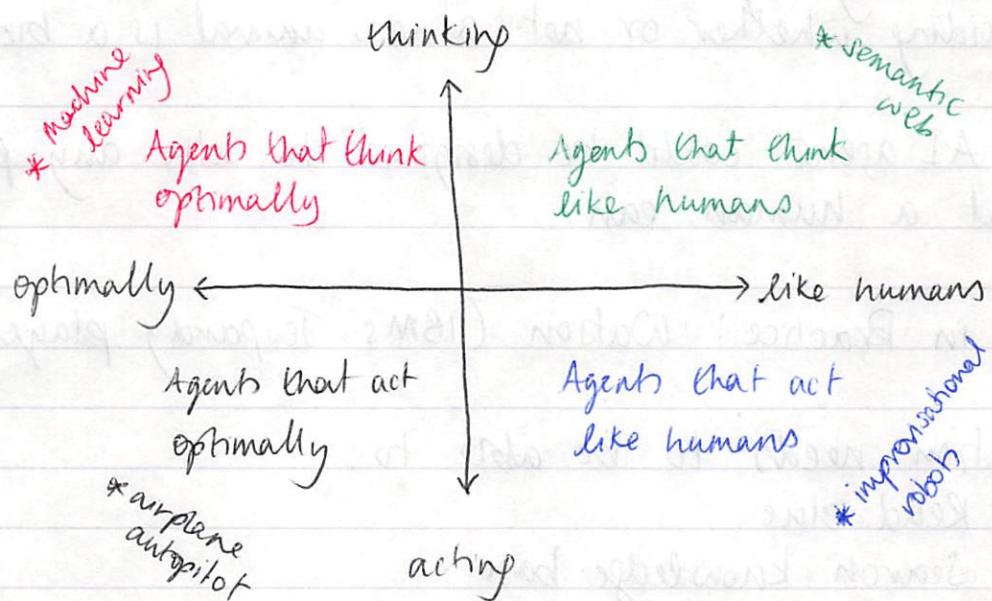
Watson needs to be able to:

- Read clue
- Search knowledge base
- Decide on an answer
- Phrase the answer as a question.

* What is Knowledge-Based AI?



* Foundations: The four schools of AI



Foundations: Machine Learning

* Introduction to Machine Learning

- Outline:
- Introduction to Machine Learning
 - Applications of Machine Learning
 - Useful pre-requisites

* Introduction

* Introduction Part II

* Prerequisites

- Programming experience (esp. Python)
- Statistics
-

* ML is the ROX

* Definition of Machine Learning

Building computational artifacts that learn over time based on experience. Not just building but the mathematics, science, engineering and computing behind the artifact

* Supervised Learning

Using information from labelled datasets to label new datasets.

Function approximation or function induction

* Induction and Deduction

You assume you have a well behaved function which is consistent with your data and you use this to generalise.

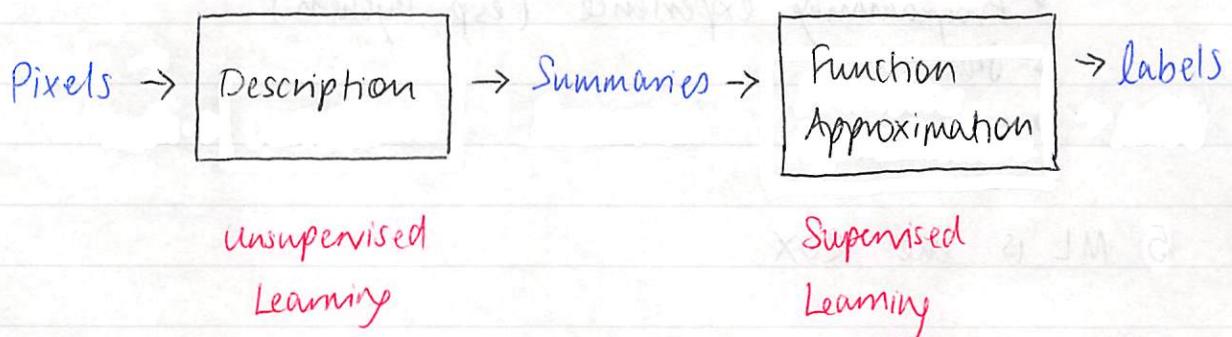
Problem: Inductive bias

Induction:- Example → General Rule

Deduction:- General Rule → Example

* Unsupervised Learning

Derive structure from the data.



* Reinforcement Learning

Learning from delayed reward

* Comparison of These Part of Machine Learning

All these problems can be formulated as some kind of optimisation problem

supervised learning	→ labels data well
reinforcement learning	→ behaviour scores well
unsupervised learning	→ cluster scores well

DATA
IS
CENTRAL!

* Machine Learning Basics

Outline :

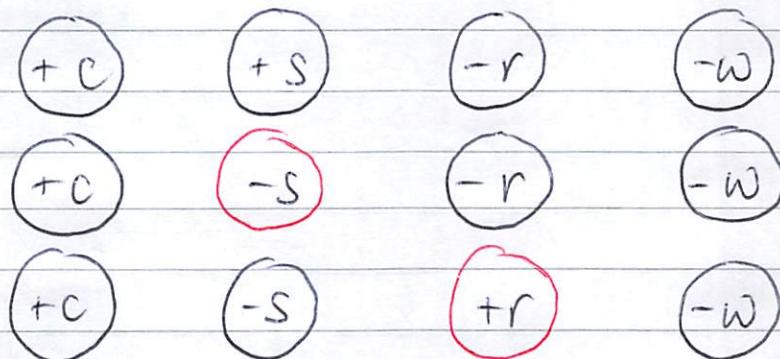
- How modern companies use Machine learning
- Stanley at the DARPA Grand Challenge
- Machine Learning taxonomy
- Overfitting: Occam's Razor
- Example problem: Spam Detection
- Linear methods for supervised learning

* What is Machine Learning

* Gibbs Sampling

... using Markov Chain Monte Carlo (MCMC)

Here we sample one variable at a time



this is consistent

Unit 5: Machine Learning I (Supervised Learning)

- Bayes networks = reason with known models
- Machine learning = learn models from data

Taxonomy of machine learning:

What? parameters, structure, hidden concepts

What from? supervised, unsupervised, reinforcement

What for? prediction, diagnosis, summarisation
How? the learning agent is... passive, active, offline, online while/not while the data is being generated

Outputs? classification, regression

Details? generative, discriminative generate a model distinguish between samples

Supervised Learning

$$\underbrace{\{x_1, x_2, \dots, x_n\}}_{\text{features of the data}} \rightarrow \underbrace{y}_{\text{feature we want to predict}}$$

$$\left[\begin{array}{ccccccc} X_{11} & X_{12} & \cdots & X_{1n} & \rightarrow Y_1 \\ X_{21} & X_{22} & \cdots & X_{2n} & \rightarrow Y_2 \\ \vdots & \vdots & & \vdots & \vdots \\ X_{m1} & X_{m2} & \cdots & X_{mn} & \rightarrow Y_m \end{array} \right] \quad \text{data}$$

$\underbrace{\hspace{10em}}_{X_m}$

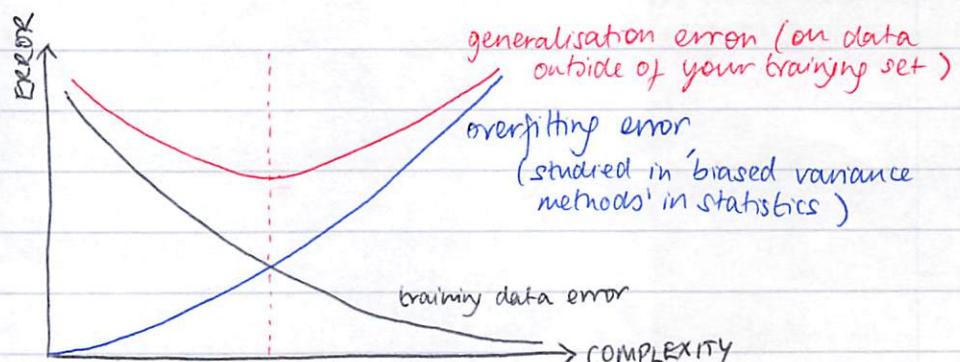
Want to find the function $f(X_m) = Y_m$ so it predicts well Y for X not in our training set

Occam's / Okham's Razor

\Rightarrow Everything else being equal, choose the less complex hypothesis

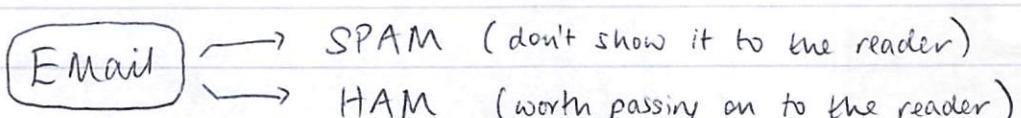
There tends to be a trade off between fit (to the training data) and complexity (of the hypothesis)

(good) fit \longleftrightarrow complexity (low)



$$\text{generalisation error} = \frac{\text{overfitting error}}{\text{error}} + \frac{\text{training data error}}{\text{error}}$$

Spam detection - A classification problem



Y can take one of two values - SPAM / HAM hence a classification problem.

Maximum likelihood

suppose

S S S H H H H H H (8 messages)

let

$$p(S) = \pi$$

we want to find π which maximises the likelihood of the training set assuming that each email is drawn independently from an identical distribution

$$p(y_i) = \begin{cases} \pi & \text{if } y_i = S = 1 \\ 1 - \pi & \text{if } y_i = H = 0 \end{cases}$$

$$p(\text{data}) = \prod_{i=1}^m p(y_i)$$

$$= \pi^{\text{count}(y_i=1)} (1 - \pi)^{\text{count}(y_i=0)}$$

$$\ln(p(\text{data})) = \frac{\text{count}(y_i=1)}{\pi} \ln \pi + \frac{\text{count}(y_i=0)}{1-\pi} \ln(1-\pi)$$

$$\frac{d}{d\pi} (\ln(p(\text{data}))) = \frac{\text{count}(y_i=1)}{\pi} - \frac{\text{count}(y_i=0)}{1-\pi}$$

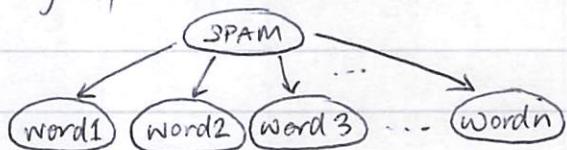
$$\begin{aligned} \frac{d}{d\pi} (\ln(p(\text{data}))) &= 0 \Rightarrow \frac{\text{count}(y_i=1)}{\pi} = \frac{\text{count}(y_i=0)}{1-\pi} \\ &\Rightarrow \text{count}(y_i=0)\pi = \text{count}(y_i=1)(1-\pi) \\ &\Rightarrow [\text{count}(y_i=0) + \text{count}(y_i=1)]\pi = \text{count}(y_i=1) \end{aligned}$$

$$\Rightarrow \pi = \frac{\text{count}(y_i=1)}{\text{count}(y_i=1) + \text{count}(y_i=0)}$$

* Bayes networks examples: Bag of words model

Naïve Bayes Network:

Detecting SPAM emails:



SPAM	HAM
OFFER IS SECRET	PLAY SPORTS TODAY
CLICK SECRET LINK	WENT PLAY SPORTS
SECRET SPORTS LINK	SECRET SPORTS EVENT
	SPORTS IS TODAY
	SPORTS COSTS MONEY

	Dictionary	SPAM	HAM
1	OFFER		
2	IS		
3	SECRET		
4	CLICK		
5	LINK		
6	SPORT		HTT
7	PLAY		
8	TODAY		
9	WENT		
10	EVENT		
11	COSTS		
12	MONEY		
	Total	9	15

Using Maximum likelihood...

$$P(\text{SPAM}) = \frac{3}{8}$$

$$P(\text{"SECRET"} | \text{SPAM}) = \frac{3}{8}$$

$$P(\text{"SECRET"} | \text{HAM}) = \frac{1}{15}$$

$$P(\text{"IS"} | \text{SPAM}) = \frac{1}{9}$$

$$P(\text{"IS"} | \text{HAM}) = \frac{1}{15}$$

$$P(\text{"TODAY"} | \text{SPAM}) = 0$$

$$P(\text{"TODAY"} | \text{HAM}) = \frac{2}{15}$$

Here it's the probability of the word in a message being "..." given ...

Q1 Message M = "SPORTS"

$$P(\text{SPAM} | M) = \frac{\frac{3}{8} \cdot \frac{1}{9}}{\frac{3}{8} \cdot \frac{1}{9} + \frac{5}{8} \cdot \frac{1}{3}} = \frac{3}{3+15} = \frac{1}{6}$$

Q2 M = "SECRET IS SECRET"

$$\begin{aligned} P(\text{SPAM} | M) &= \frac{\frac{3}{8} \cdot \frac{1}{3} \cdot \frac{1}{9} \cdot \frac{1}{3}}{\frac{3}{8} \cdot \frac{1}{3} \cdot \frac{1}{9} \cdot \frac{1}{3} + \frac{5}{8} \cdot \frac{1}{15} \cdot \frac{1}{15} \cdot \frac{1}{15}} = \frac{\frac{3}{8} \times 5 \times \frac{1}{9} \times 5}{\frac{3}{8} \times 5 \times \frac{1}{9} \times 5 + \frac{5}{8} \times 1 \times \frac{1}{15} \times 1} \\ &= \frac{25}{26} \end{aligned}$$

Q3 M = "TODAY IS SECRET"

$$P(\text{SPAM} | M) = \frac{\frac{3}{8} \cdot 0 \cdot \frac{1}{9} \cdot \frac{1}{3}}{\frac{3}{8} \cdot 0 \cdot \frac{1}{9} \cdot \frac{1}{3} + \frac{5}{8} \cdot \frac{2}{15} \cdot \frac{1}{15} \cdot \frac{1}{15}} = 0$$

The third example is a poor estimate clearly. Here we are overfitting.

* Laplace Smoothing

Technique to deal with overfitting.

$$\text{Maximum Likelihood : } P(x) = \frac{\text{count}(x)}{N}$$

$$\text{Laplace Smoothing : } P(x) = \frac{\text{count}(x) + k}{N + k|x|}$$

$\text{count}(x)$ = # of occurrences of this value of x
 $|x|$ is the number of values which x can have
k is the smoothing parameter

This is like adding k dummy emails to each class (SPAM and HAM) each containing all the words in the dictionary.

Using Laplace Smoothing for the Bayes networks example ...

$$k=1 \Rightarrow P(\text{SPAM}) = \frac{3+1}{8+2} = \frac{4}{10} = \frac{2}{5}$$

$$\Rightarrow P(\text{HAM}) = \frac{3}{5}$$

$$P(\text{"TODAY"} | \text{SPAM}) = \frac{0+1}{9+12} = \frac{1}{21}$$

$$P(\text{"TODAY"} | \text{HAM}) = \frac{2+1}{15+12} = \frac{3}{27} = \frac{1}{9}$$

$$P(\text{"IS"} | \text{SPAM}) = \frac{1+1}{9+12} = \frac{2}{21}$$

$$P(\text{"IS"} | \text{HAM}) = \frac{1+1}{15+12} = \frac{2}{27}$$

$$P(\text{"SECRET"} | \text{SPAM}) = \frac{3+1}{9+12} = \frac{4}{21}$$

$$P(\text{"SECRET"} | \text{HAM}) = \frac{1+1}{15+12} = \frac{2}{27}$$

$$P(\text{SPAM} | \text{"TODAY IS SECRET"})$$

$$= \frac{\frac{2}{5} \cdot \frac{1}{21} \cdot \frac{2}{21} \cdot \frac{4}{21}}{\frac{2}{5} \cdot \frac{1}{21} \cdot \frac{2}{21} \cdot \frac{4}{21} + \frac{3}{5} \cdot \frac{1}{9} \cdot \frac{2}{27} \cdot \frac{4}{27}}$$

$$= \frac{4}{4 + \frac{1}{3} \cdot \frac{1}{27} \cdot \frac{1}{27} \cdot 21 \cdot 21 \cdot 21}$$

$$= \frac{4}{4 + \frac{1}{9} \cdot \frac{1}{9} \cdot \frac{4}{7} \cdot \frac{7}{7}} = \frac{4}{4 + \frac{343}{81}}$$

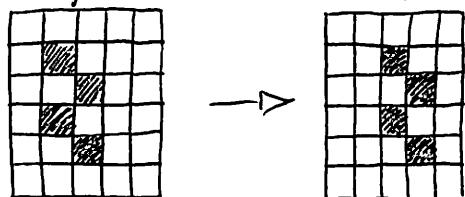
Advanced Spam filters

- known spamming IP?
- have you emailed the person before?
- have 1K other people received the message?
- is the email header and IP consistent?
- is the email all caps
- do inline URLs point to where they say?
- are you addressed by name?

All these can be used in Naive Bayes

* Digital Recognition

Recognising handwritten digits - could try to use Naive-Bayes on 16×16 pixels but this doesn't deal



well with shifting. To deal with this one could use input smoothing in a different way. One could mix pixel counts with those of the neighbouring pixels so if they are shifted we get similar statistics. Here we convolve the input with a Gaussian variable. This may give better results than just using the raw pixel values themselves.

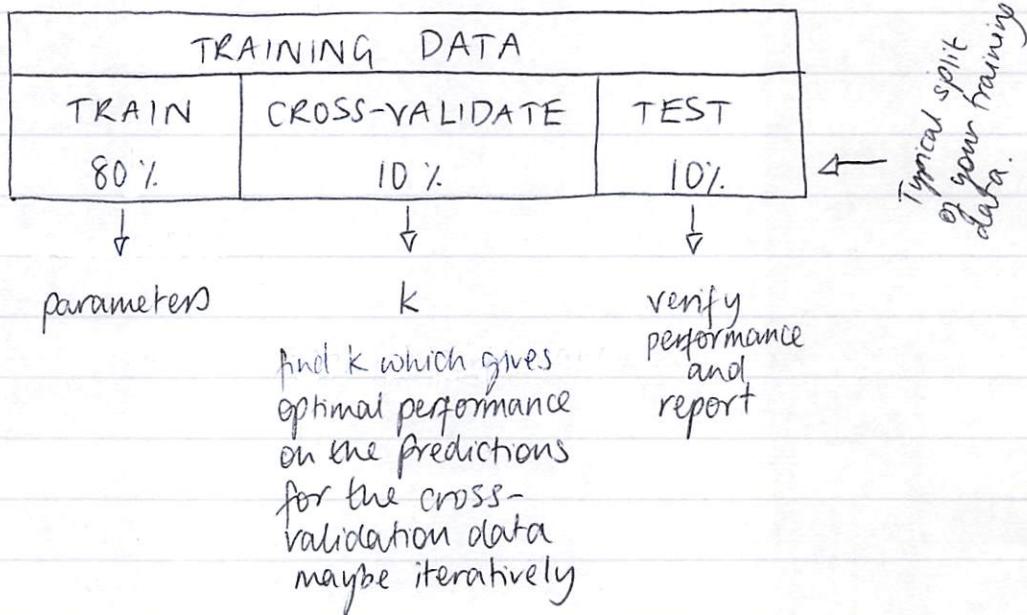
Actually Naive-Bayes is not a great choice for this problem since it turns out that conditional independence of the pixels in this case is too strong an assumption.

* Overfitting Prevention

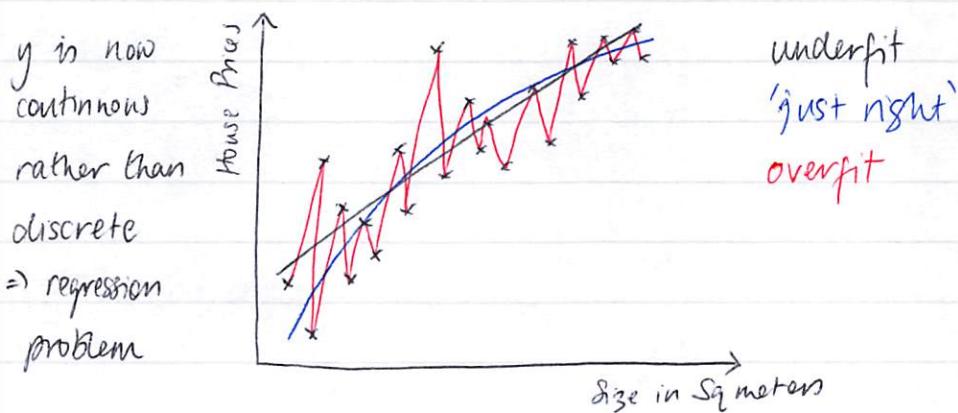
We talked about

- Occam's Razor (trade off between fit / prediction)
 - Laplace smoothing }
 - Input smoothings }
- How do we choose the smoothing parameter?

Cross-validation:



* Classification vs Regression



DATA

$$\begin{bmatrix} X_{11} & X_{12} & \dots & X_{1n} & \rightarrow Y_1 \\ X_{21} & X_{22} & \dots & X_{2n} & \rightarrow Y_2 \\ \vdots & \vdots & & \vdots & \vdots \\ X_{m1} & X_{m2} & \dots & X_{mn} & \rightarrow Y_m \end{bmatrix}$$

$$f(x) = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + w_0$$

$$f(x) = w^T x + w_0$$

We want to find f | $f(x) = y$. We do this by minimising the loss function

$$\text{Loss} = \sum_j^m (y_j - w_1 x_j - w_0)^2 \quad \text{sum square error}$$

$$\text{solution } w^* = \operatorname{arg\min}_w \{ \text{LOSS} \}$$

* Minimising quadratic loss

$$\min_w \sum_{i=1}^m (y_i - w_1 x_i - w_0)^2 = L$$

$$\frac{\partial L}{\partial w_0} = -2 \sum_{i=1}^m (y_i - w_1 x_i - w_0)$$

$$\frac{\partial L}{\partial w_0} = 0 \Rightarrow m w_0 = \sum_{i=1}^m (y_i - w_1 x_i)$$

$$\Rightarrow w_0 = \frac{1}{m} \left(\sum_{i=1}^m y_i - w_1 \sum_{i=1}^m x_i \right)$$

$$\frac{\partial L}{\partial w_1} = -2 \sum_{i=1}^m (y_i - w_1 x_i - w_0) x_i$$

$$\frac{\partial L}{\partial w_1} = 0 \Rightarrow w_1 \sum_{i=1}^m x_i^2 = \sum_{i=1}^m (y_i - w_0) x_i$$

$$\Rightarrow w_1 \sum_{i=1}^m x_i^2 + w_0 \sum_{i=1}^m x_i = \sum_{i=1}^m x_i y_i$$

$$\Rightarrow w_1 \sum_{i=1}^m x_i^2 + \frac{1}{m} \sum_{i=1}^m (y_i - w_0) \sum_{i=1}^m x_i = \sum_{i=1}^m x_i y_i$$

$$\Rightarrow w_1 \sum_{i=1}^m x_i^2 + \frac{1}{m} \sum_{i=1}^m y_i \sum_{i=1}^m x_i - \frac{w_0}{m} \left(\sum_{i=1}^m x_i \right)^2 = \sum_{i=1}^m x_i y_i$$

$$\Rightarrow w_1 = \frac{\sum_{i=1}^m x_i y_i - \frac{1}{m} \sum_{i=1}^m x_i \sum_{i=1}^m y_i}{\sum_{i=1}^m x_i^2 - \frac{1}{m} \left(\sum_{i=1}^m x_i \right)^2}$$

Problems with linear regression

- non-linear data
- outliers
- classification problems

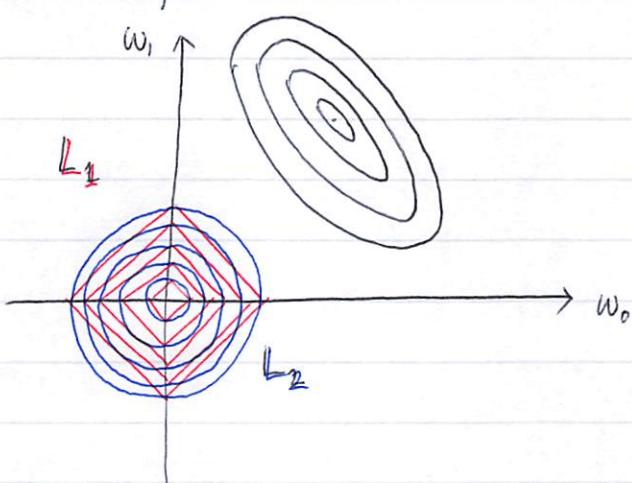
For classification problems we can use logistic regression.

$$\frac{1}{1 + e^{-f(x)}}$$

Regularisation is used for complexity control

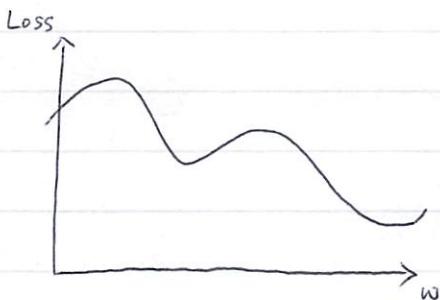
$$\text{Loss} = \text{Loss}(\text{data}) + \text{Loss}(\text{parameters})$$

$$= \sum_j (y_j - w_1 x_j - w_0)^2 + \sum_i |w_i|^p$$



$$p=1 : L_1 \text{ regularisation}$$
$$p=2 : L_2 \text{ regularisation}$$

* Minimisation of more complicated loss functions



Gradient descent

$$w^0 \\ w^{i+1} \leftarrow w^i - \alpha \frac{\partial L}{\partial w_i}$$

* Gradient descent implementation

$$L = \sum_j (y_j - w_1 x_j - w_0)^2 \rightarrow \min$$

$$\frac{\partial L}{\partial w_1} = -2 \sum_j (y_j - w_1 x_j - w_0) x_j$$

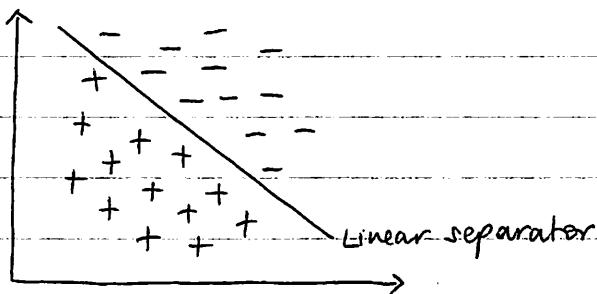
$$\frac{\partial L}{\partial w_0} = -2 \sum_j (y_j - w_1 x_j - w_0)$$

$$w_0 = w_0^0 \text{ and } w_1 = w_1^0$$

$$w_1^i = w_1^{i-1} - \alpha \frac{\partial L}{\partial w_1}(w_1^{i-1})$$

$$w_0^i = w_0^{i-1} - \alpha \frac{\partial L}{\partial w_0}(w_0^{i-1})$$

* Perceptron algorithm

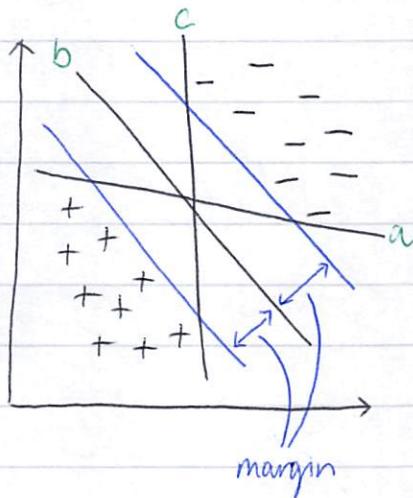


$$f(x) = \begin{cases} 0 & \text{if } w_1 x + w_0 < 0 \\ 1 & \text{if } \underbrace{w_1 x + w_0}_{\text{linear function}} \geq 0 \end{cases}$$

Start with a random guess for w_0 and w_1

$$w_i^k \leftarrow w_i^{k-1} + \alpha (y_i - f(x_i))$$

* Linear separators



separator b is preferable to a and c because of the large margin.

Perceptron only finds a linear separator - not the "best" one.

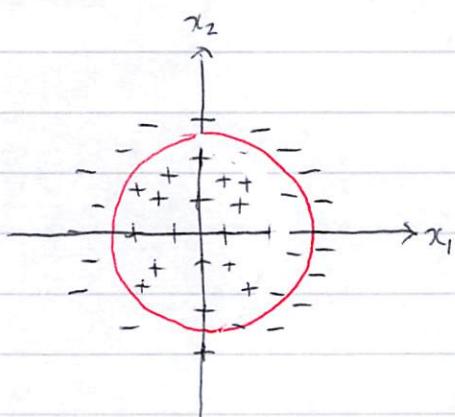
Maximum margin algorithms:

- support vector machines
- boosting

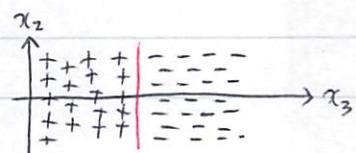
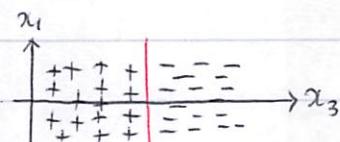
Support vector machines:

These use a "Kernel Trick" to find features which turn complex non-linear decision boundaries into linear ones

Illustration:



$$x_3 = \sqrt{x_1^2 + x_2^2}$$



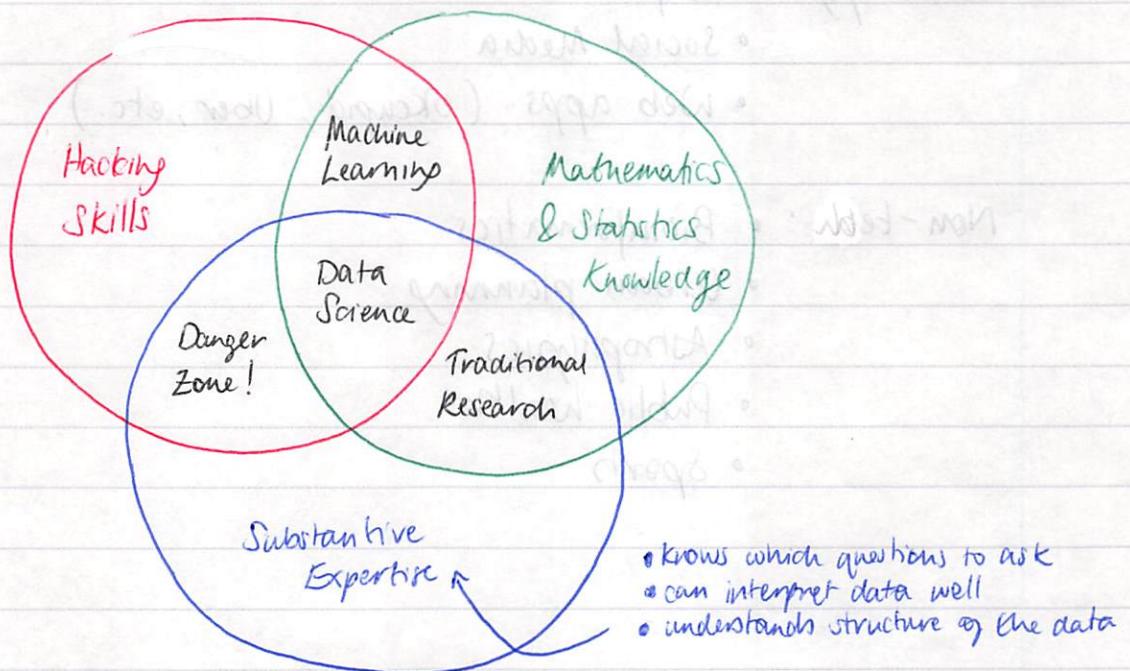
Linear Methods

- Regression vs Classification
- Exact vs Iterative solutions
- Smoothing
- Non-linear problems

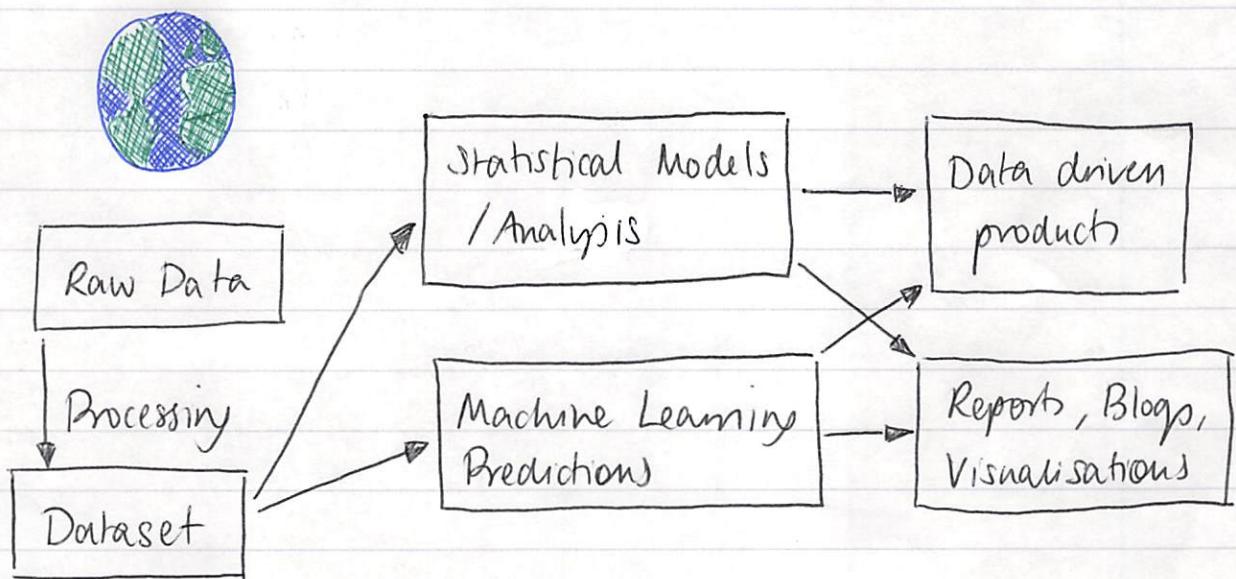
Introductory Material: Data Science

* What is a Data Scientist?

Person who is better at statistics than a software engineer and better at software engineering than a statistician



* What does a Data Scientist do?



* Simpson's Paradox

- See Berkeley admissions data
- vdmlab.com/simpson
- Problems solved by Data Science

Technology : ◦ Netflix
◦ Social Media
◦ Web apps (OKCupid, Uber, etc.)

Non-tech : ◦ Bioinformatics
◦ Urban planning
◦ Astrophysics
◦ Public health
◦ Sports

C1: Model Evaluation & Validation

* Measures of Central Tendency

* Mode - Most frequently occurring

Median - Middle value when the data is ordered
- Half way between the two middle values
when there is an even number of points

Mean - Average = $\sum_{i=1}^n x_i/n$

* Mode

- distribution can have more than one mode (or set of values)
- not all samples of a distⁿ contribute to it
- not affected by outliers

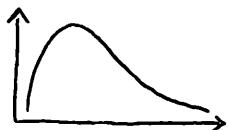
Mean

- one value
- all samples of a distⁿ contribute to it
- affected by outliers

Median

- one value
- all samples of a distⁿ contribute to it
- not affected by outliers

* Order of measures



mode < median < mean



mode = median = mean



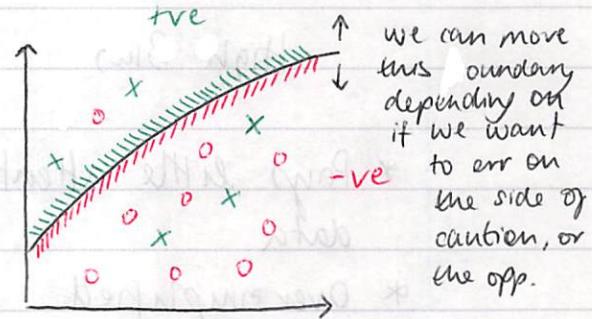
mode > median > mean

Evaluation Metrics

* Confusion Matrix Example:

		actual class	
		+ve	-ve
predicted class	+ve	9	1
	-ve	3	8

Time +ve
False +ve
Time -ve
False -ve



we can move this boundary depending on if we want to err on the side of caution, or the opp.

* Precision & Recall

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

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

Terminology:

$\left\{ \begin{array}{l} \text{True} \\ \text{False} \end{array} \right\}$ $\left\{ \begin{array}{l} \text{+ve} \\ \text{-ve} \end{array} \right\}$
 refers to the correctness of the prediction
 refers to the prediction result

* Example:

ACTUAL			PREDICTED					
			13	4	1	1	0	0
Ariel Sharon								
Colin Powell			0	55	0	8	0	0
Donald Rumsfeld			0	1	25	8	0	0
George W Bush			0	3	0	123	0	1
Gerhard Schroeder			0	1	0	7	14	0
Hugo Chavez			0	3	0	2	1	10
Tony Blair			0	0	1	7	0	26

precision = $\frac{1}{10}$
 recall = $\frac{1}{5} = \frac{1}{5/8}$

Causes of Error

* Bias Variance Dilemma Solution

High Bias

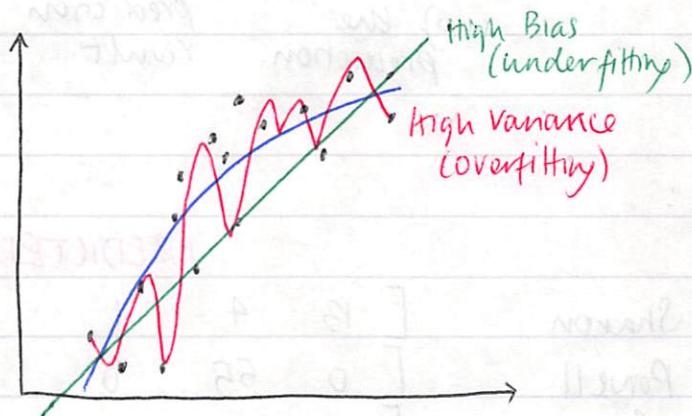
- * Pays little attention to data
- * Oversimplified
- * High error on training set (low r^2 , high SSE)
- * Few features used

High Variance

- * Pays too much attention to data (doesn't generalise well)
- * Overfit
- * Much higher error on test set than training set
- * Carefully minimised error over many features
→ optimised performance on training set

Want to use as few features as possible to

- maximise r^2
- minimise SSE



Data Types

- * Numeric Data

- Discrete
- Continuous

- * Categorical Data

- Can take on numerical values but these don't have mathematical meaning i.e. mean or variance don't make sense
- Ordinal data - Categories with ordering/rank

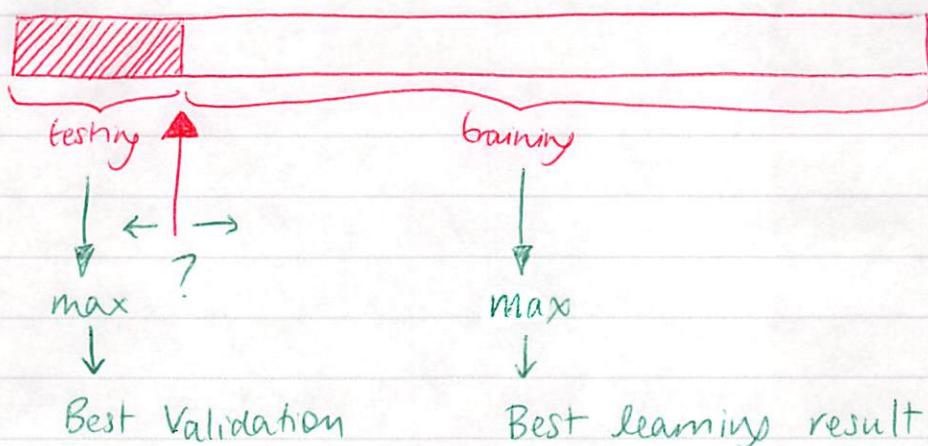
- * Time Series Data

- Data point collected over time at

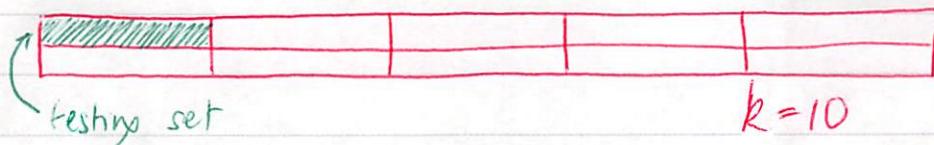
Cross Validation

* k - Fold Cross Validation

Problems with splitting into training and testing data



k fold cross validation



Partition your data into k equal parts and run k separate learning experiments. For each experiment

- pick a different testing set
- train on the remaining data

Average test results from those k experiments

This way we use all our data for both training and testing. Clearly training will take longer in this case than the simple train then test method of validation but k -fold cross validation results in better accuracy in our model.

Representative Power of a Model

* Curse of Dimensionality

As the number of number of features or equivalently dimensions grows, the amount of data we need to generate accurately, grows exponentially

→ You're better off getting more data than you are adding more features!

Topic one

Topic two

Topic three

Topic four

Topic five

C2: Supervised Learning

Decision Trees

* Terminology :

Instances - input

Concept - function: input \rightarrow output

Target concept - function which classifies correctly

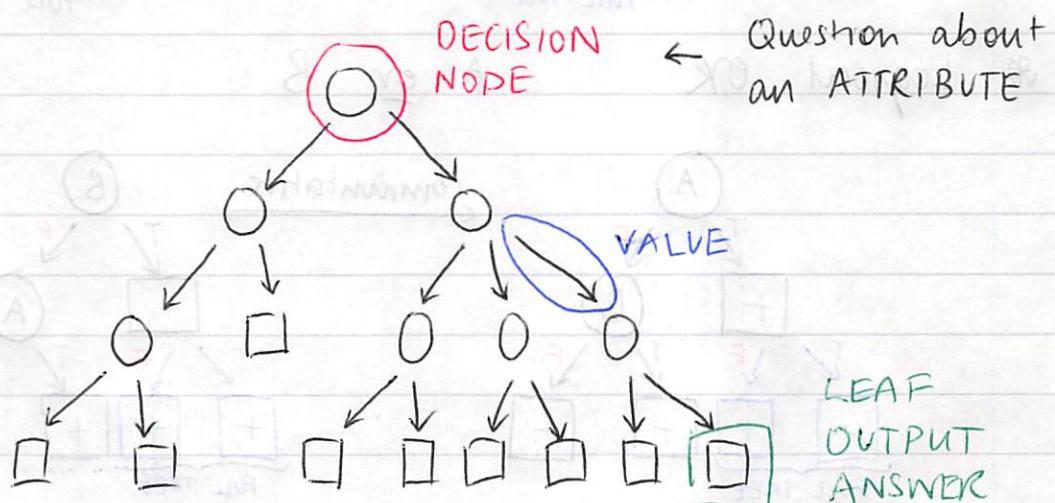
Hypothesis class - class of concepts under consideration

Sample - training set

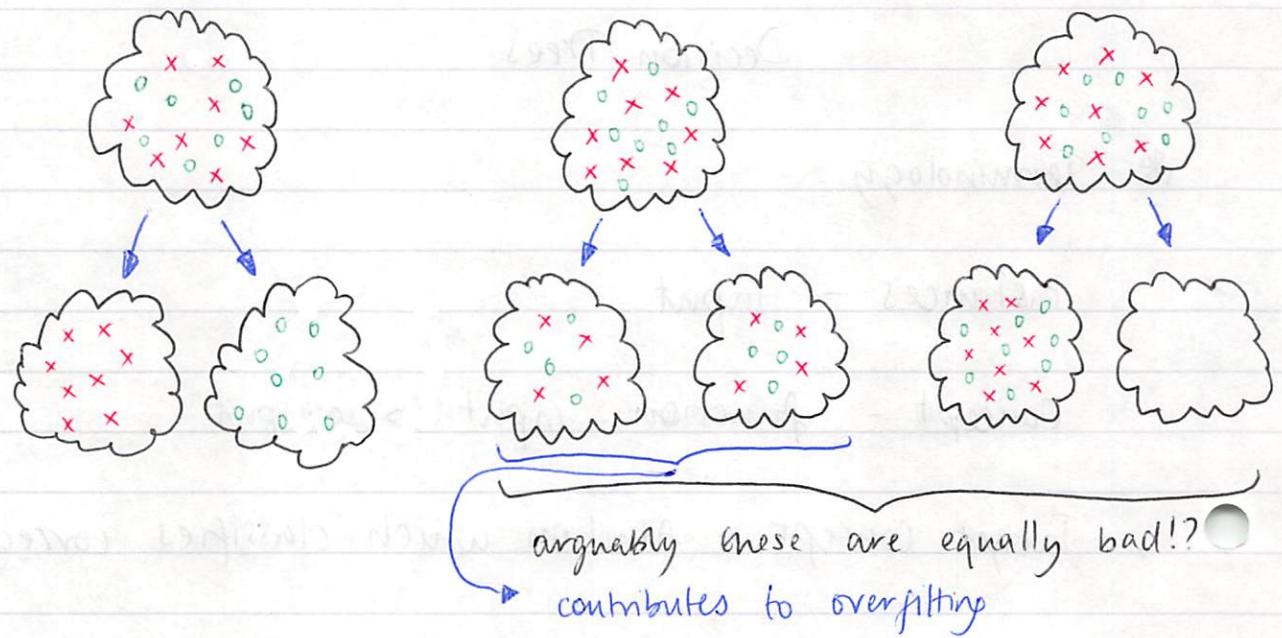
Candidate concept

Testing set - independent of training set

* Representation

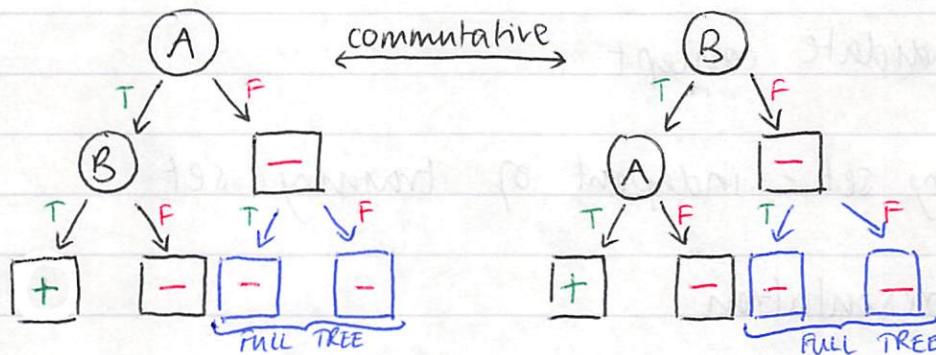


* Best Attribute : In order of preference :

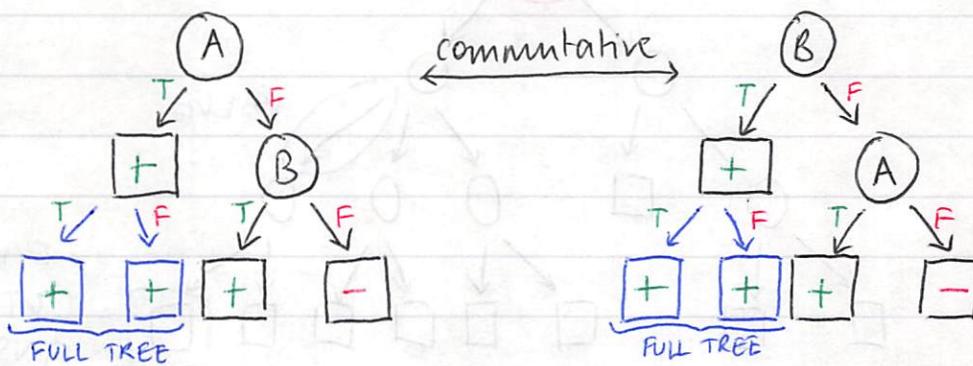


* Decision Trees : Expressiveness

* Logical AND A and B

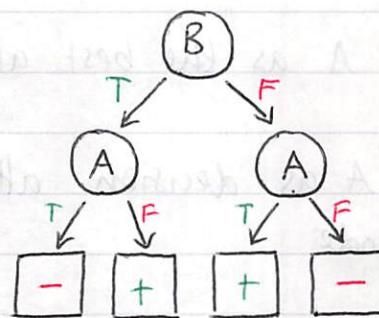
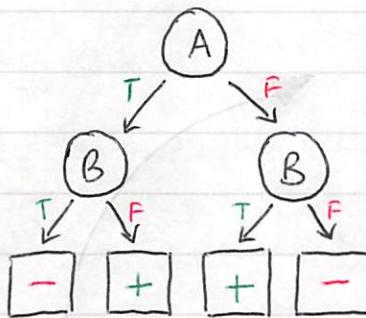


* Logical OR A or B

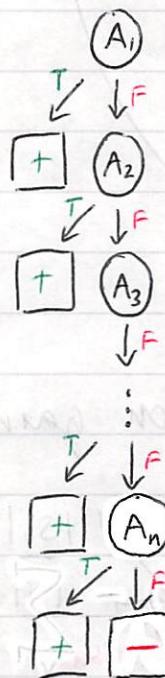


* Logical XOR

$A \text{ xor } B$



* n-OR : ANY

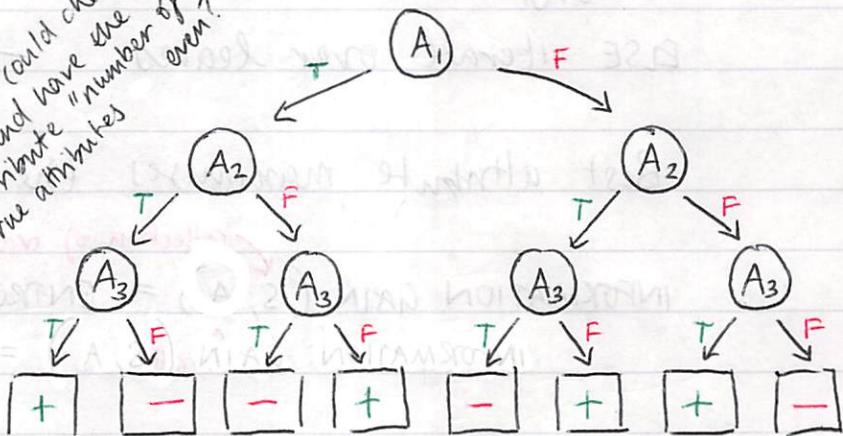


* n-XOR : ODD PARITY

i.e. if number of true attributes is odd then

exponential in n i.e. $2^n - 1$ nodes

(could cheat
and have one
attribute "number of
true attributes
even?"



linear in n
i.e. n nodes

Let's make a truth table ...

A_1	A_2	A_3	...	A_n	output
T	T	T		T	T/F
F	T	T		T	T/F
T	F	T		T	T/F
⋮	⋮	⋮	⋮	⋮	⋮
F	F	F		F	T/F

2^n rows

$2^{(2^n)}$
possible outputs

→ get to one row
from all rows
good history

* ID3

choose A as the best attribute



Assign A as decision attribute
for node



FOR EACH value of A

create a decendent of node



Sort training examples to leaves



| examples perfectly classified
STOP

ELSE iterate over leaves

Best attribute maximises the information gain

$$\text{INFORMATION GAIN } (S, A) = \text{ENTROPY}(S) - \sum_v \frac{|S_v|}{|S|} \text{ENTROPY}(S_v)$$

collection of training samples
attribute

$$\text{ENTROPY} = - \sum_v \underbrace{p(v)}_{\text{probability of seeing value } v} \ln p(v)$$

* ID3 Bias

Restriction Bias : e.g. we restrict our soln to the class of binary trees H

Preference Bias : $h \in H \Leftrightarrow$ Inductive Bias

- Good split at top
- Correct over incorrect
- Shorter trees

* Decision Trees : Other considerations

- Continuous attributes

e.g. Age, weight, distance

→ ask about ranges

→ can ask about the same attribute (diff range)
more than once along the same path

- When do we stop?

→ everything classified correctly

→ no more attributes ← Not smart

→ no overfitting ← tree too deep

← cross validation

← pruning ← Build whole tree and then
look at how cutting it back or
pruning affects the error

← But what if you have bad
data - say the same data point
twice with two different
classification

- Regression

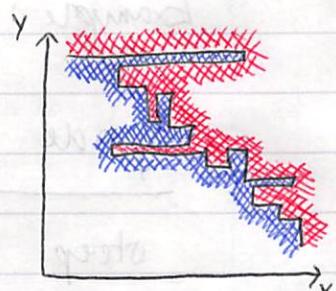
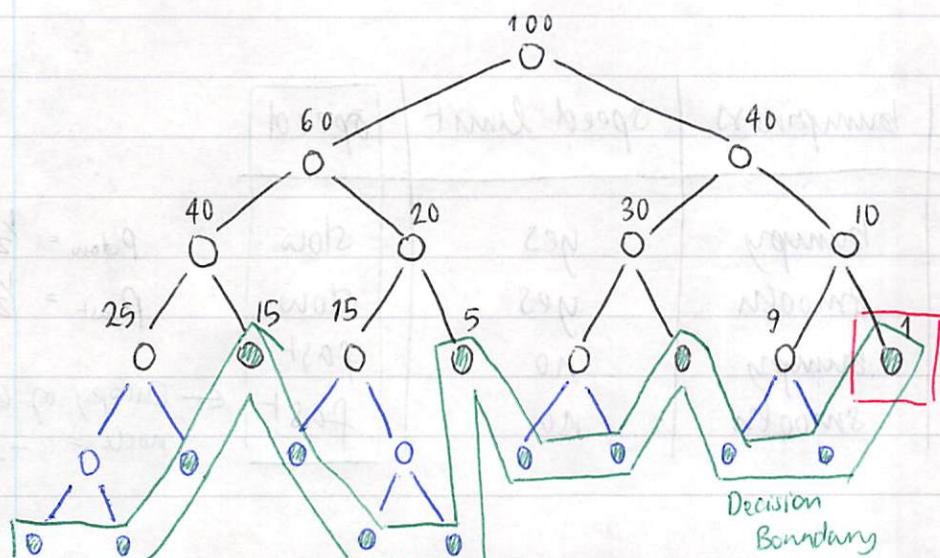
→ Picking the best attribute entropy → variance?

→ Output - average, local linear fit

More Decision Trees

* Decision Tree Parameters

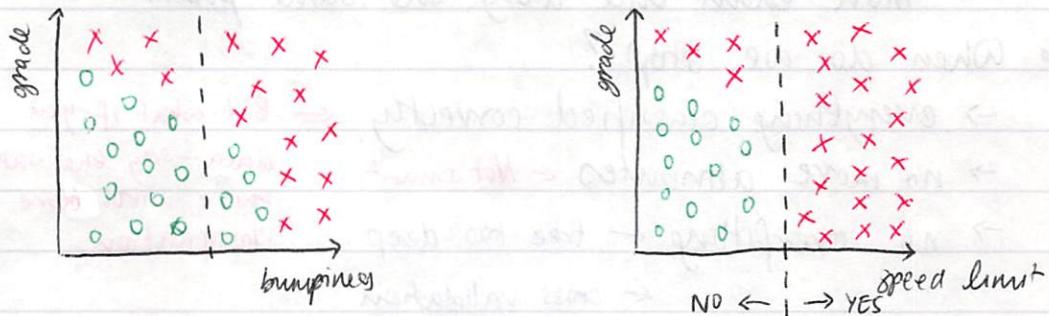
min_samples_split = 2 by default



* Data Impurity & Entropy

Entropy \leftarrow Controls how a decision tree decides where to split the data

Defn: measure of impurity in a bunch of examples



* Formulas

$$\text{entropy} = \sum_i^1 -p_i \log_2(p_i)$$

↑ sum over all classes i
 ↑ fraction of examples in class i

Intuition:

- All examples are one same class \Rightarrow entropy = 0
- Examples are evenly split between classes \Rightarrow entropy = 1

Example:

grade	bumpiness	Speed limit	speed
steep	bumpy	yes	slow
steep	smooth	yes	slow
flat	bumpy	no	fast
steep	smooth	no	fast

$$p_{\text{slow}} = \frac{1}{2}$$

$$p_{\text{fast}} = \frac{1}{2}$$

← Entropy of this node = $-\log_2(\frac{1}{2}) = 1$

* Information Gain

$$\text{Information gain} = \text{entropy}_{\text{(parent)}} - \left[\frac{\text{weighted}}{\text{average}} \right]_{\text{(children)}} \text{entropy}_{\text{(children)}}$$

Decision Tree : Maximise Information Gain

Example:

CHILD	bumpiness	speed limit	PARENT
grade			speed
steep steep	bumpy smooth	yes yes	slow slow
flat steep	bumpy smooth	no no	fast fast

$$IG = 0.3113$$

$$\text{entropy} = 1$$

$$IG = 0$$

$$\text{entropy} = 0$$

$$IG = 1$$

$$\text{entropy} = 1$$

steep

ssf

$$P_{\text{slow}} = 2/3$$

$$P_{\text{fast}} = 1/3$$

$$\text{entropy} = 0.9184$$

ssf

flat

f

← entropy of this node = 1

← entropy of this node = 0

$$\left[\frac{\text{weighted}}{\text{average}} \right]_{\text{(children)}} \text{entropy} = \frac{3}{4} \cdot 0.9184 + \frac{1}{4} \cdot 0$$

$$\text{information gain} = 1 - \frac{3}{4} \cdot 0.9184 = 0.3113$$

SKLearn DecisionTreeClassifier has a "criterion" parameter which specifies the function which measures the split quality

* Bias - Variance Dilemma

Bias algorithm ignores data

High variance algorithm pays too much attention to the data

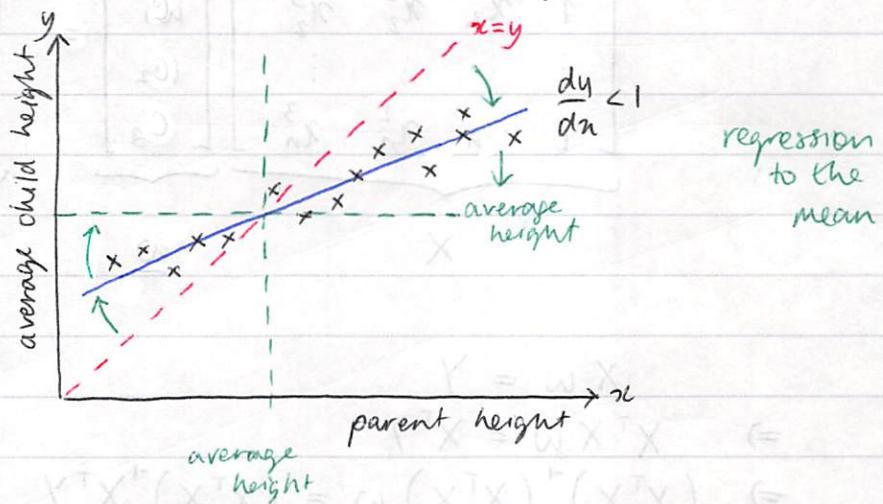
Find the sweet spot.

Regression & Classification

* Regression

Supervised Learning : Take examples of inputs and outputs. Now given a new input predict its output.

Regression - where the name came from :



* Loss / Error function

$$i \quad | \quad x_i \quad \text{Guess} \quad x_i = c$$

$$1 \quad | \quad x_1$$

$$2 \quad | \quad x_2$$

$$\vdots$$

$$n \quad | \quad x_n$$

$$E(c) = \sum_{i=1}^n (x_i - c)^2$$

$$\frac{dE}{dc} = -2 \sum_{i=1}^n (x_i - c)$$

$$\frac{dE}{dc} = 0 \Leftrightarrow \sum_{i=1}^n (x_i - c) = 0 \Leftrightarrow c = \frac{1}{n} \sum_{i=1}^n x_i = \text{mean}$$

* Polynomial Regression

x	y
x_1	y_1
x_2	y_2
\vdots	\vdots
x_n	y_n

cubic regression: fit,

$$w_0 + w_1 x + w_2 x^2 + w_3 x^3 = y$$

$$\begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 & x_n^3 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

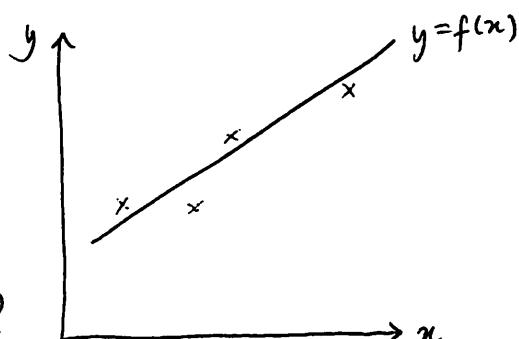
$X \quad w \quad Y$

$$\begin{aligned} Xw &= Y \\ \Rightarrow X^T X w &= X^T Y \\ \Rightarrow (X^T X)^{-1} (X^T X) w &= (X^T X)^{-1} X^T Y \\ \Rightarrow w &= (X^T X)^{-1} X^T Y \end{aligned}$$

* Errors

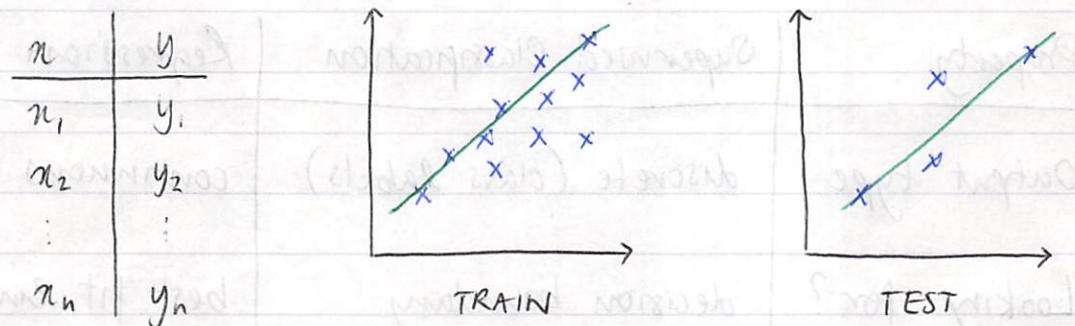
Training data has errors
so we end up not modelling
 f but $f + \epsilon$.

Where do errors come from?

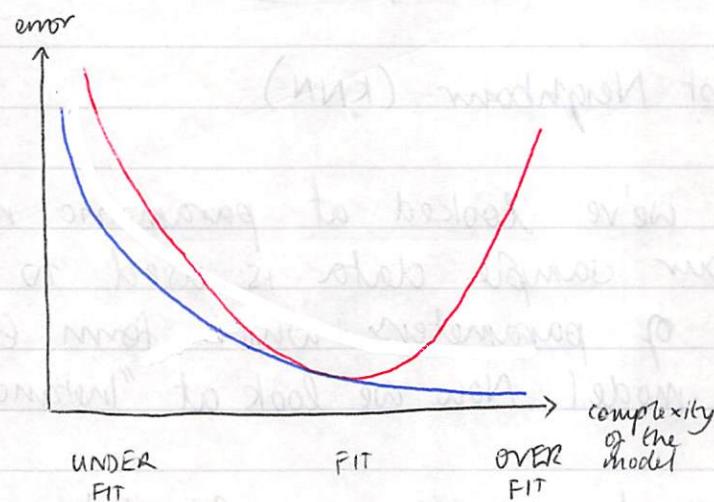


- sensor error
- maliciously - being given bad data
- transcription error
- unmodelled influences

* Cross Validation



We assume the training and testing datasets are 'representative' of the population i.e. IID
This is a fundamental assumption



* r^2 of a regression

r^2 metric describes how much of the change in the output y is explained by change in the input x .

$\rightarrow 0 < r^2 < 1$
line doesn't capture the trend in the data well at all

↑ line describes the relationship between x and y well

* Comparing Classification & Regression

Property	Supervised Classification	Regression
Output type	discrete (class labels)	continuous
Looking for?	decision boundary	best fit line
Evaluation	accuracy	SSE or r^2

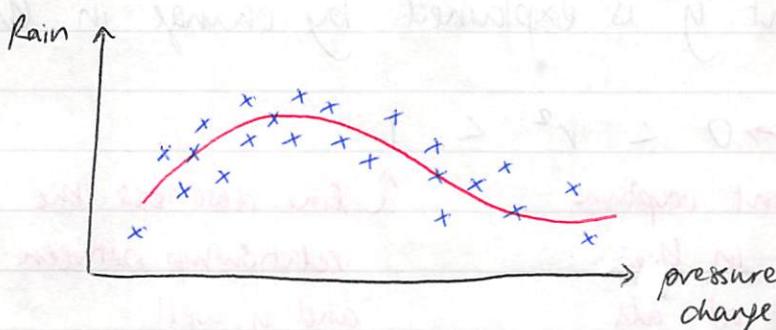
More Regression

* K Nearest Neighbour (KNN)

So far we've looked at parametric regression where our sample data is used to find a number of parameters which form the basis of our model. Now we look at "Instance Based Methods"

KNN instead uses the sample data directly. For a given input we average the output of the k nearest input points.

Eg. 3 NN :



* Kernel Regression

Like KNN but the predicted output for a given input is given by the weighted average of the nearest neighbour's outputs. The weights are proportional to the distance.

In KNN the weights are equal.

* Parametric vs Instance Based Learners

Parametric models are good when we expect a functional relationship between inputs and outputs. The function can be used to make an initial guess for the solution.

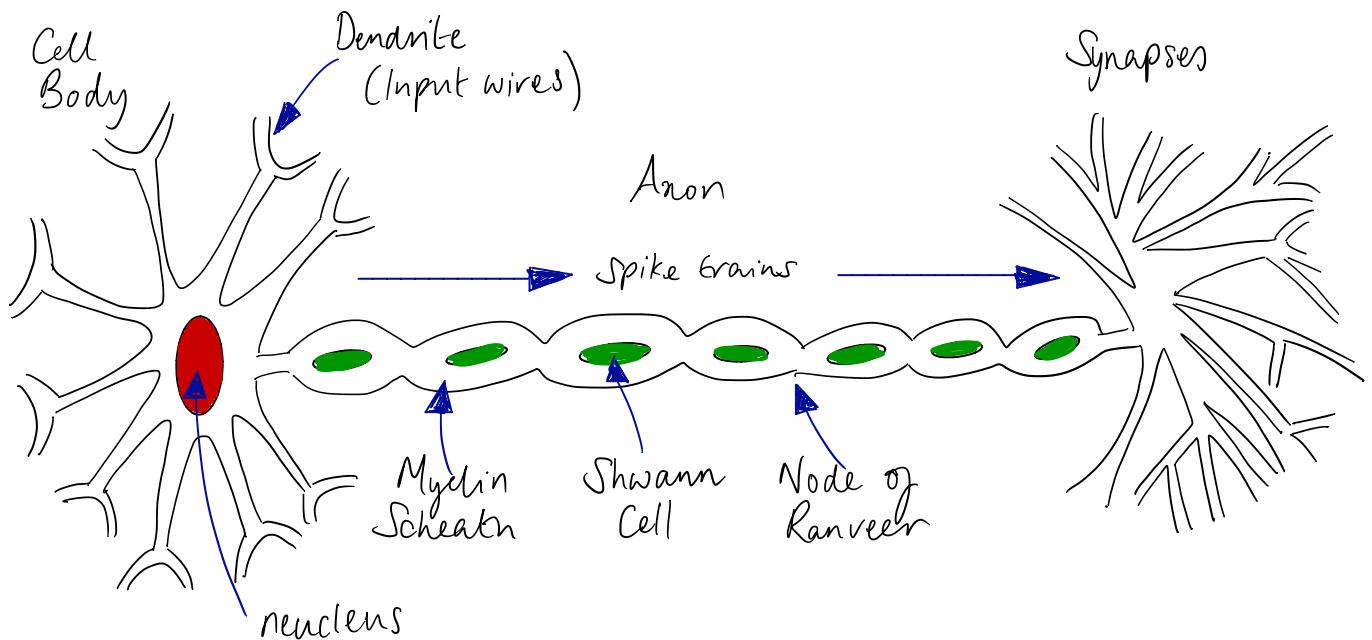
Instance based methods are good when we don't know of what ^{the} relationship between input and output might be.

Parametric models don't require the original data to be stored. However updating the model when new data is added is arduous. So, for parametric approaches training is slow and querying is fast.

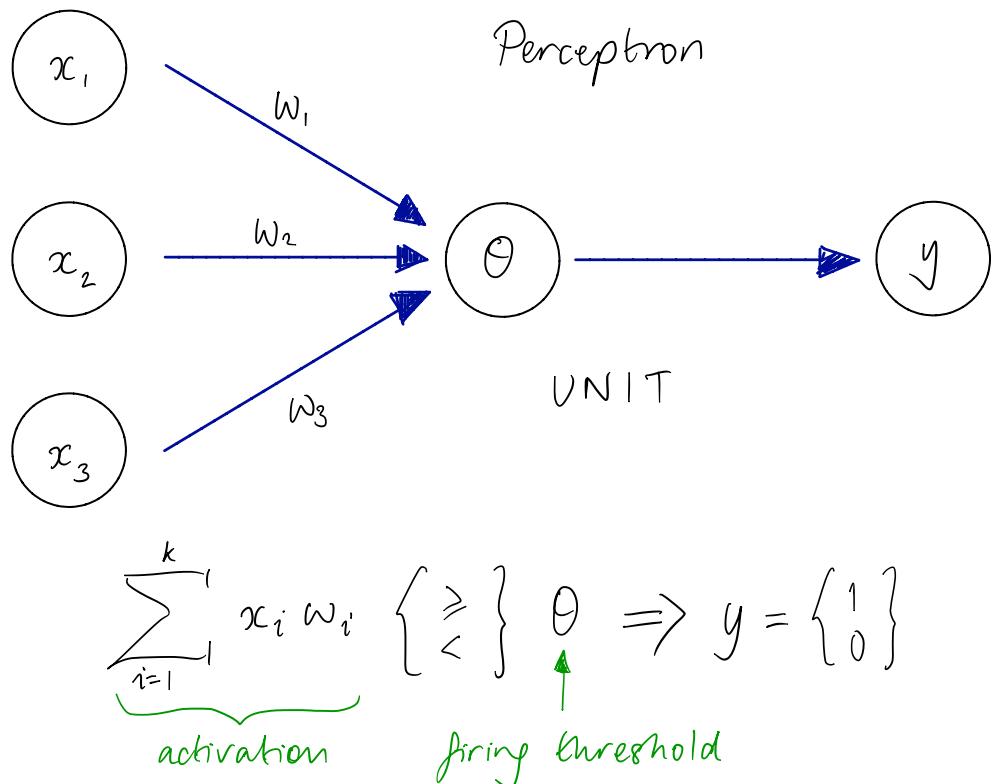
For instance based models training is fast and querying is slow. New data is easily added to update the model.

Neural Networks

* Neuron: Component Unit



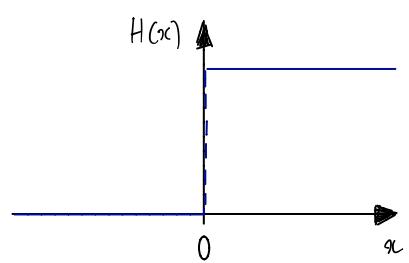
* Artificial Neural Network



We can write this as

$$y = H \left(\sum_{i=1}^k x_i w_i - \theta \right)$$

where H is the heaviside function



* How powerful is a perceptron unit?

Example : $w_1 = \frac{1}{2}, w_2 = \frac{1}{2}, \theta = \frac{3}{4}$

Recall,

$$y = H \left(\underbrace{\sum_{i=1}^k x_i w_i}_{\text{activation}} - \theta \right)$$

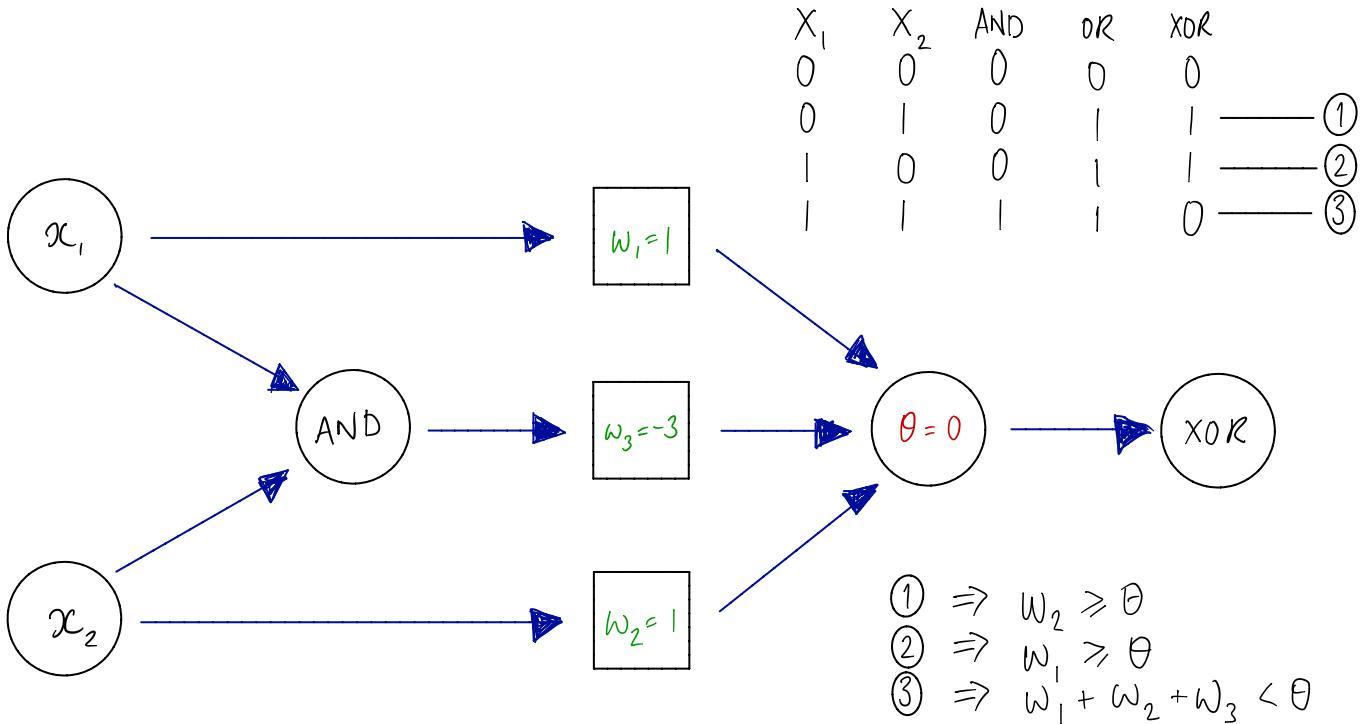
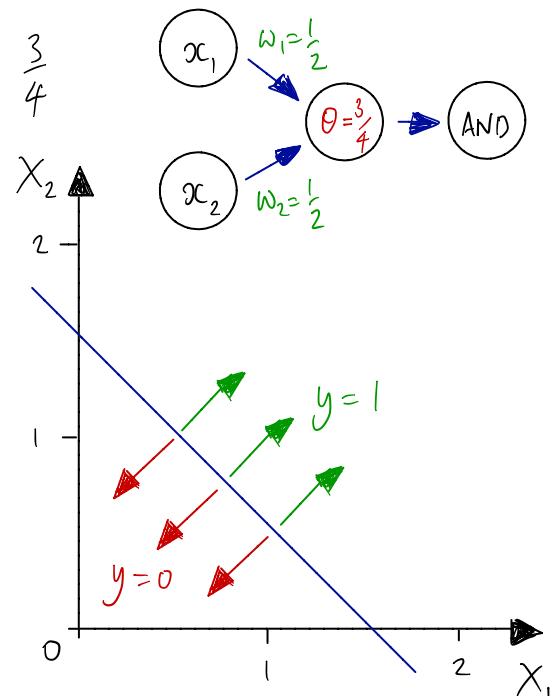
$$\text{activation} = \frac{1}{2} (x_1 + x_2) = \frac{3}{4}$$

$$\Leftrightarrow x_1 + x_2 = \frac{3}{2}$$

If we take $x_1, x_2 \in \{0, 1\}$,

these weights and fixing threshold give the logical AND function. Similarly $w_1 = w_2 = 1$ and $\theta = \frac{1}{2}$ would give logical OR and using just one dimension, i.e.

$w_1 = -1$ and $\theta = -\frac{1}{2}$ would give the logical NOT function. What about XOR?



* Perception Training

Given examples find weights that map inputs to outputs

1. perceptron rule

2. gradient descent / delta rule

threshold

unthresholded

1. Perceptron Rule → Single unit

$$\underline{x} = \begin{pmatrix} 1 & x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \quad \text{"Bias"} \quad \underline{w} = \begin{pmatrix} -\theta & w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix}$$

y = target
 \hat{y} = output
 η = learning rate

\underline{x} = input
 \underline{w} = weights
 n = number of features

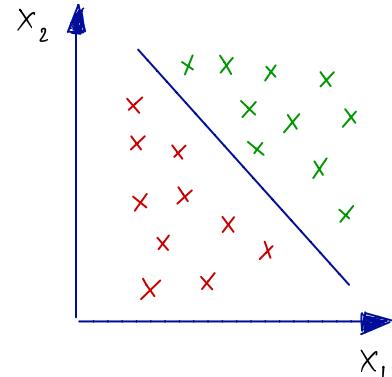
Update Rule : $w_i = w_i + \Delta w_i$

$$\Delta w_i = \eta (y - \hat{y}) x_i$$

$$\hat{y} = H \left(\sum_{i=0}^n w_i x_i \right)$$

$$\left. \begin{array}{l} \hat{y} < y \Rightarrow \Delta w_i < 0 \\ \hat{y} > y \Rightarrow \Delta w_i > 0 \\ \hat{y} = y \Rightarrow \Delta w_i = 0 \end{array} \right\}$$

In the case where our dataset is **LINEARLY SEPARABLE** the perceptron rule algorithm will find the set of weights which map the inputs to outputs in a finite number of iterations... although, we don't know how many iterations it will take.



What about data that is separable by a non-linear boundary?

2. Gradient Descent

$$a = \sum_{i=0}^n x_i w_i$$

Now our output is not thresholded i.e. $\hat{y} = H(a)$ and instead we calculate an error and minimise over the weight

So we have

$$a = \sum_{i=0}^n x_i w_i \quad \text{and} \quad E(w) = \frac{1}{2} \sum_{(x,y) \in D} (y - a)^2$$

$$\begin{aligned} \Rightarrow \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \left\{ \frac{1}{2} \sum_{(x,y) \in D} (y - a)^2 \right\} \\ &= \sum_{(x,y) \in D} (y - a) \frac{\partial}{\partial w_i} \left\{ -\sum x_i w_i \right\} \\ &= \sum_{(x,y) \in D} (y - a) (-x_i) \end{aligned}$$

For gradient descent we have the new update rule

$$\begin{aligned} \Delta w_i &= -\eta \frac{\partial E}{\partial w_i} \\ \Rightarrow \Delta w_i &= \eta (y - a) x_i \quad \text{looks much like Perception Rule!} \end{aligned}$$

* Comparison of learning rules

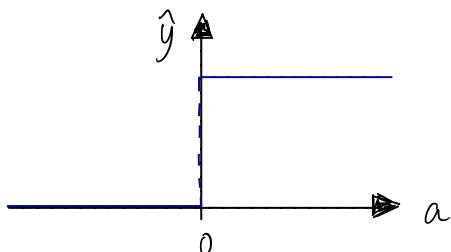
perception : $\Delta w_i = \eta (y - \hat{y}) x_i, \quad \hat{y} = H(a)$

- guarantee of finite convergence in the case of linearly separable data

gradient descent : $\Delta w_i = \eta (y - a) x_i$

- more robust in the case of non-linearly separable datasets
- converges in the limit to a local optimum

Qn: Why not do gradient descent on \hat{y} ?



Because \hat{y} is a non-differentiable function of the activation.

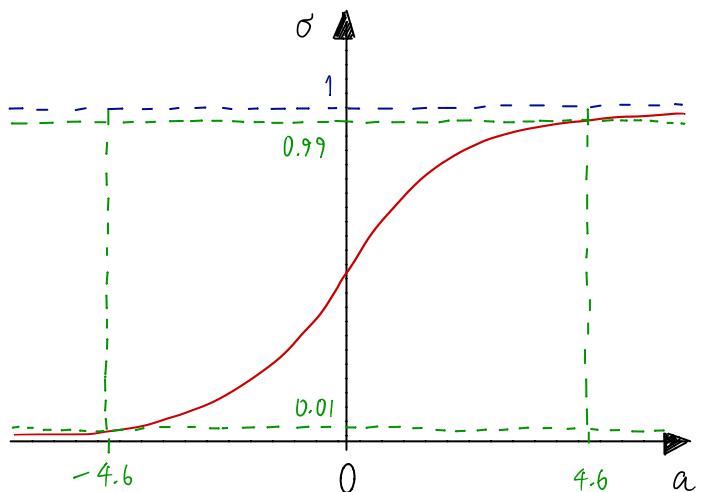
* Sigmoid - Differentiable threshold

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

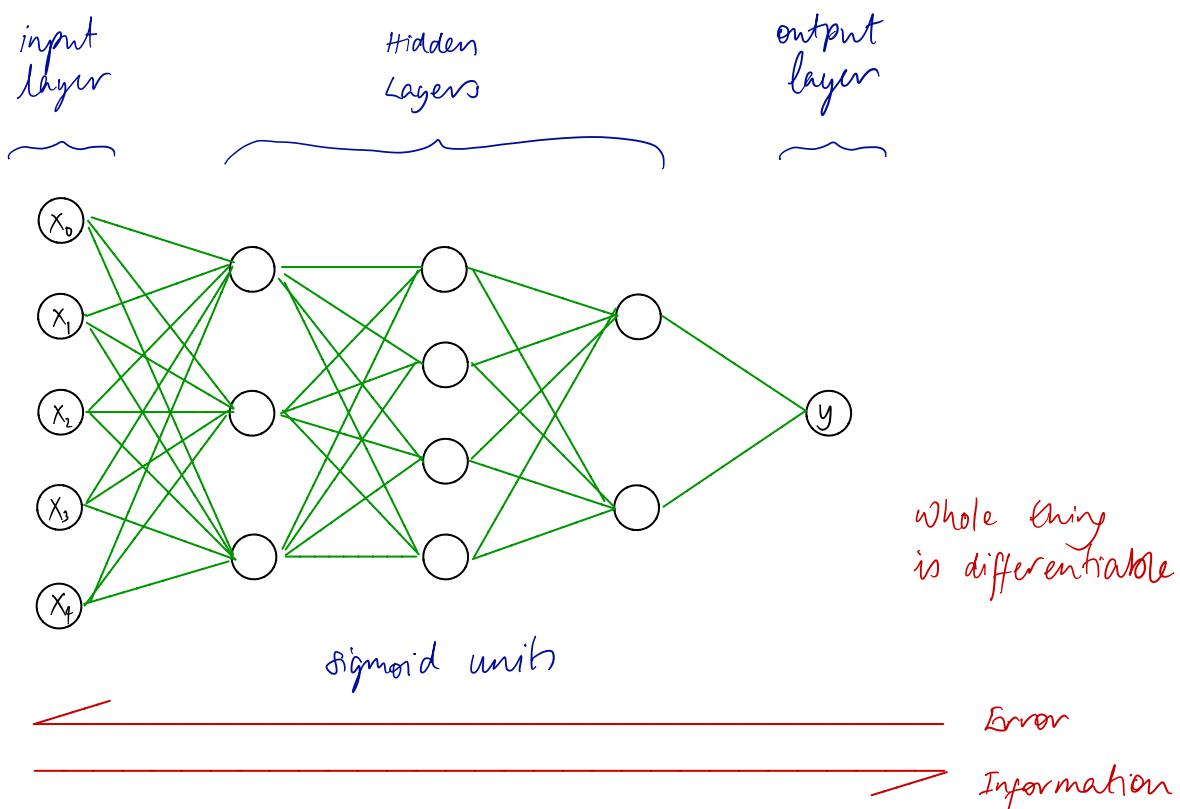
as $a \rightarrow \infty$, $\sigma(a) \rightarrow 1$

as $a \rightarrow -\infty$, $\sigma(a) \rightarrow 0$

$$\frac{d\sigma}{da} = \sigma(a)[1 - \sigma(a)]$$



* Neural Network Sketch



(Error) Back Propagation : Computationally beneficial organisation of the chain rule

While for a single unit there may well be a global optima for a network there are typically many local optima

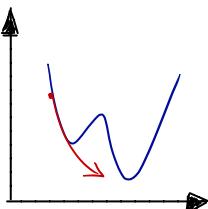
* Optimising Weights

It can been argued that Learning \equiv Optimisation...

Gradient descent can get stuck in a local optima...

Advanced optimisation methods:

- momentum:



- higher order derivatives:

combinations of derivatives rather than just first order derivatives with respect to the weights

- randomised optimisation: to make the optimisation more robust

- penalty for complexity: for a neural network reducing complexity means
 - fewer nodes
 - fewer layers
 - smaller weights

* Restriction Bias

Representational power

Set of hypotheses we will consider

- perceptron - half spaces divided by a hyperplane
- sigmoids
- more nodes, more layers and larger weights

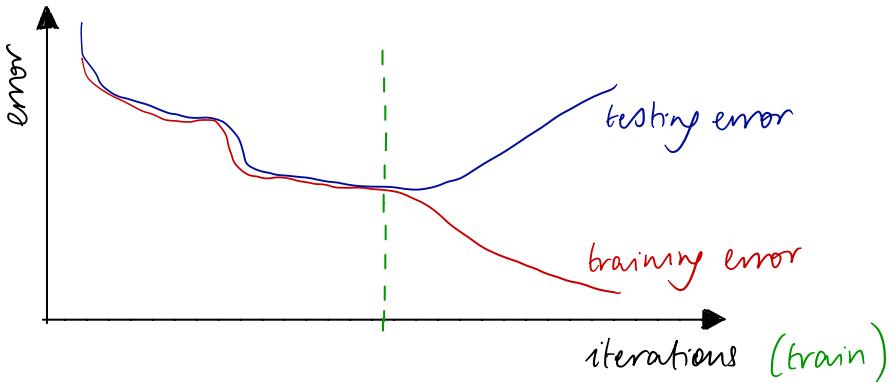
→ Boolean functions: network of threshold like units

→ Continuous functions: need one hidden layer with enough nodes

→ Arbitrary functions: (including discontinuous functions) need two hidden layers

DANGER OF OVERFITTING!

- A given network architecture is restricted in what it can capture
- We can, as before, use cross validation to decide
 - how many hidden layers
 - how many nodes in each layer
 - when to stop training because the weights have gotten too large



* Preference Bias

Algorithms selection of one representation over another

What algorithm?

Gradient descent - Initial weights?

→ small random values

- local minima variability
i.e. to avoid getting stuck in the same local minima when we rerun the training
- low complexity → smaller weights
i.e. to avoid overfitting

Occam's Razor: Entities should not be multiplied unnecessarily

→ All things being equal, choose the less complex solution.

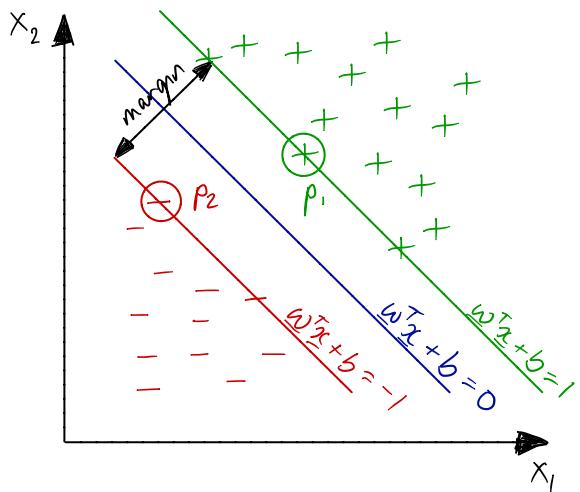
→ Shown to reduce the generalisation error

Kernel Methods & Support Vector Machines

* Support Vector Machines

The 'best' hyperplane correctly classifies the training data without over-committing to it i.e. the plane which fits the data but does not overfit to it.

N.B. This is not overfitting in the complexity sense



Suppose our data is linearly separable and the equation of our discriminating plane is of the form

$$\underline{w}^T \underline{x} + b = 0$$

This time $y \in \{+1, -1\}$ i.e.

$$y = 2H(\underline{w}^T \underline{x} + b) - 1$$

We take the equation of the plane parallel to the decision boundary and above it which touches the closest data point to be $\underline{w}^T \underline{x} + b = 1$ and the corresponding plane below the decision boundary has the equation $\underline{w}^T \underline{x} + b = -1$

[Note at this point \underline{w} and b are unknown - we want to find them.]

We want to maximise the distance between the red and green planes i.e. maximise the 'margin'

Recall \underline{w} is a vector normal to the planes. Given the two data points p_1 and p_2 on the green and red planes respectively, the distance between the planes can be calculated as the projection of the vector $p_1 - p_2$ onto the unit vector $\underline{w}/|\underline{w}|$, i.e. the dot product:

$$\text{margin} = \frac{\underline{w}}{|\underline{w}|} \cdot (\underline{p}_1 - \underline{p}_2) = \frac{\underline{w}^T}{|\underline{w}|} (\underline{p}_1 - \underline{p}_2)$$

We know $\underline{w}^T p_1 + b = 1$ and $\underline{w}^T p_2 + b = -1$

Subtracting : $\underline{w}^T (p_1 - p_2) = 2$

$$\Rightarrow \text{margin} = \frac{\underline{w}^T (p_1 - p_2)}{\|\underline{w}\|} = \frac{2}{\|\underline{w}\|}$$

So we want to maximise $\frac{2}{\|\underline{w}\|}$ while correctly classifying everything i.e.

$$\underbrace{\max \left\{ \frac{2}{\|\underline{w}\|} \right\}}_{\text{maximise margin}} \text{ and } y_i (\underline{w}^T \underline{x}_i + b) \geq 1 \quad \forall i \quad \begin{matrix} \text{the i-th sample} \\ \text{is correctly classified} \end{matrix}$$

equivalently $\min \left\{ \frac{1}{2} \|\underline{w}\|^2 \right\}$ ← always has a unique solution

which is a quadratic programming problem which can then be transformed into the standard form of quadratic programming problem using

$$\underline{w} = \sum_{i=1}^m \alpha_i y_i \underline{x}_i \quad \begin{matrix} \text{sum over} \\ m \text{ samples} \end{matrix}$$

In our new problem we want to find the α s which maximise

$$W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \underline{x}_i^T \underline{x}_j$$

such that

$$\alpha_i \geq 0, \quad \sum_{i=1}^m \alpha_i y_i \neq 0$$

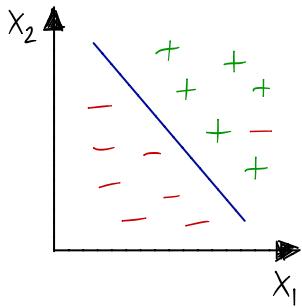
← this term is larger for data points with more similar features

Actually, it turns out that most of the α s are zero because data points which are far away from the decision boundary don't matter as they are not involved in defining it.

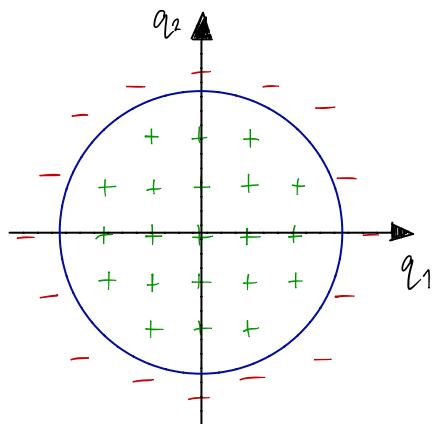
This is a bit like KNN, but solving our quadratic programming problem determines which data points are the important ones i.e. the ones with corresponding non-zero α .

The data points for which the corresponding α s are non-zero are the "SUPPORT VECTORS"

* Linearly Separable



Hence we want to maximise the margin while minimising the classification error



For this problem we can use a trick.
Define the transformation,

$$\Phi(\underline{x}) = (q_1^2, q_2^2, \sqrt{2}q_1q_2),$$

and apply it to all our data points. Then our decision boundary

$$\underline{w}^\top \underline{x} + b = 0$$

get mapped from a plane to an ellipse:

$$\underline{w}^\top \Phi(\underline{x}) + b = 0 \Rightarrow w_1 x_1^2 + w_2 x_2^2 + \sqrt{2} w_3 x_1 x_2 + b = 0$$

Let look at the effect of this transform on $W(x)$

$$W(x) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \underline{x}_i^\top \underline{x}_j$$

our notion of similarity or two points

$$\begin{aligned} \Phi(\underline{x}) \cdot \Phi(\underline{y}) &= (x_1^2, x_2^2, \sqrt{2} x_1 x_2) \cdot (y_1^2, y_2^2, \sqrt{2} y_1 y_2) \\ &= x_1^2 y_1^2 + x_2^2 y_2^2 + 2 x_1 x_2 y_1 y_2 \\ &= (x_1 y_1 + x_2 y_2)^2 = (\underline{x}^\top \underline{y})^2 \end{aligned}$$

So actually we need not calculate $\Phi(\underline{q})$ at all. We can just replace $\underline{x}_i^\top \underline{x}_j$ in $W(x)$ with $(\underline{x}_i^\top \underline{x}_j)^2$

* Kernel

$$W(x) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(\underline{x}_i, \underline{x}_j)$$

Kernel
Domain in Knowledge
Similarity

Example Kernels :

$$K(\underline{x}, \underline{y}) = (\underline{x}^T \underline{y} + c)^P \quad \text{polynomial kernel}$$

$$K(\underline{x}, \underline{y}) = \exp(-\|\underline{x} - \underline{y}\|^2 / 2\sigma^2) \quad \text{gaussian kernel}$$

$$K(\underline{x}, \underline{y}) = \tanh(\alpha \underline{x}^T \underline{y} + \theta) \quad \text{sigmoid type kernel}$$

$K(\underline{x}, \underline{y})$ has to be a valid 'distance' function

Mercer Condition :

$K(\underline{x}, \underline{y})$ must be positive semi definite i.e.

A symmetric function $K(\underline{x}, \underline{y})$ can be expressed as an inner product

$$K(\underline{x}, \underline{y}) = \langle \phi(\underline{x}), \phi(\underline{y}) \rangle$$

for some $\phi \Leftrightarrow K$ is +ve semi definite i.e.

$$\int K(\underline{x}, \underline{y}) g(\underline{x}) g(\underline{y}) d\underline{x} d\underline{y} \geq 0 \quad \forall g$$

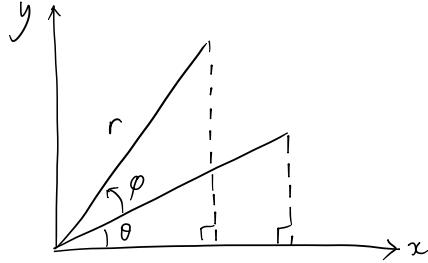
or equivalently

$$\begin{pmatrix} K(\underline{x}_1, \underline{x}_1) & K(\underline{x}_1, \underline{x}_2) & \cdots & K(\underline{x}_1, \underline{x}_n) \\ K(\underline{x}_2, \underline{x}_1) & K(\underline{x}_2, \underline{x}_2) & \cdots & K(\underline{x}_2, \underline{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ K(\underline{x}_n, \underline{x}_1) & K(\underline{x}_n, \underline{x}_2) & \cdots & K(\underline{x}_n, \underline{x}_n) \end{pmatrix}$$

is +ve semi definite.

* Boosting ???

* Article: Equation of an ellipse ...



rotation by φ

$$\begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} r \cos \theta \\ r \sin \theta \end{pmatrix} \rightarrow \begin{pmatrix} r \cos(\theta + \varphi) \\ r \sin(\theta + \varphi) \end{pmatrix}$$

$$\begin{pmatrix} r(\cos \theta \cos \varphi - \sin \theta \sin \varphi) \\ r(\sin \theta \cos \varphi + \cos \theta \sin \varphi) \end{pmatrix} = \begin{pmatrix} \cos \varphi & -\sin \varphi \\ \sin \varphi & \cos \varphi \end{pmatrix} \begin{pmatrix} r \cos \theta \\ r \sin \theta \end{pmatrix}$$

$$\begin{pmatrix} \cos \varphi & -\sin \varphi \\ \sin \varphi & \cos \varphi \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} x \cos \varphi - y \sin \varphi \\ x \sin \varphi + y \cos \varphi \end{pmatrix}$$

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} = 1 \rightarrow \frac{1}{a^2} (x \cos \varphi - y \sin \varphi)^2 + \frac{1}{b^2} (x \sin \varphi + y \cos \varphi)^2 = 1$$

$$\Rightarrow \frac{1}{a^2} (x^2 \cos^2 \varphi - 2xy \cos \varphi \sin \varphi + y^2 \sin^2 \varphi) + \frac{1}{b^2} (x^2 \sin^2 \varphi + 2xy \cos \varphi \sin \varphi + y^2 \cos^2 \varphi) = 1$$

$$\Rightarrow \left(\frac{\cos^2 \varphi}{a^2} + \frac{\sin^2 \varphi}{b^2} \right) x^2 + 2 \cos \varphi \sin \varphi \left(\frac{1}{b^2} - \frac{1}{a^2} \right) xy + \left(\frac{\sin^2 \varphi}{a^2} + \frac{\cos^2 \varphi}{b^2} \right) y^2 = 1$$

We had

$$w_1 x^2 + w_2 y^2 + \sqrt{2} w_3 xy + c = 0$$

$$w_1 = \left(\frac{\cos^2 \varphi}{a^2} + \frac{\sin^2 \varphi}{b^2} \right) \quad w_2 = \left(\frac{\sin^2 \varphi}{a^2} + \frac{\cos^2 \varphi}{b^2} \right)$$

$$w_3 = \sqrt{2} \cos \varphi \sin \varphi \left(\frac{1}{b^2} - \frac{1}{a^2} \right) \quad \text{and} \quad c = 1$$

$$w^T \underline{x} + b = 0 \rightarrow w_1 x_1^2 + w_2 x_2^2 + w_3 \sqrt{2} x_1 x_2 + b = 0$$

$$\Rightarrow x_1^2 + \frac{w_3 \sqrt{2} x_1 x_2}{w_1} + \frac{w_2 x_2^2}{w_1} + b = 0$$

$$\Rightarrow \left(x_1 + \frac{w_3 \sqrt{2} x_2}{2w_1} \right)^2 + \left(\frac{w_2}{w_1} - \frac{w_3^2}{2w_1^2} \right) x_2^2 + b = 0$$

$$\Rightarrow x_1 = \pm \sqrt{\left(\frac{w_3^2}{2w_1^2} - \frac{w_2}{w_1} \right) x_2^2 - b} - \frac{w_3 \sqrt{2} x_2}{\sqrt{2} w_1}$$

Can we plot it if we want?

$$q_1, q_2 : \quad x^2 + y^2 = 1$$

$$\Phi(y) : \quad \begin{aligned} x &\rightarrow x^2 \\ y &\rightarrow y^2 \\ z &\rightarrow \sqrt{2}xy \end{aligned} \quad \left. \begin{array}{l} x^4 + y^4 = 1 \\ z^2 = 2x^2y^2 \end{array} \right\}$$

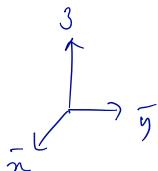
What about the blue circle? what does that get mapped to?

$$(x, y) \rightarrow (\bar{x}, \bar{y}, \bar{z}) = (x^2, y^2, \sqrt{2}xy)$$

$$(x, \sqrt{1-x^2}) \rightarrow (x^2, 1-x^2, \sqrt{2(1-x^2)}x)$$

$$\bar{y} = 1 - \bar{x} \quad \bar{z} = \sqrt{2xy} = \sqrt{2\bar{x}(1-\bar{x})}$$

$$\bar{x} = 1 - \bar{y} \Rightarrow \bar{z}^2 = 2\bar{x}(1-\bar{x}) = 2\bar{x} - 2\bar{x}^2$$



$$\Rightarrow \bar{z}^2 + 2\bar{x}^2 - 2\bar{x} = 0$$

$$\Rightarrow \bar{z}^2 + 2(\bar{x} - \frac{1}{2})^2 = \frac{1}{4}$$

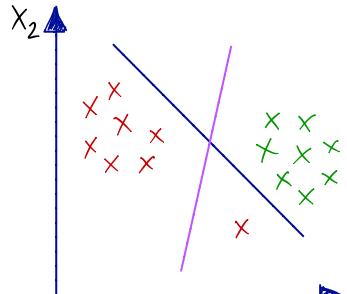
← *an ellipse*

If all these points lie on a plane.

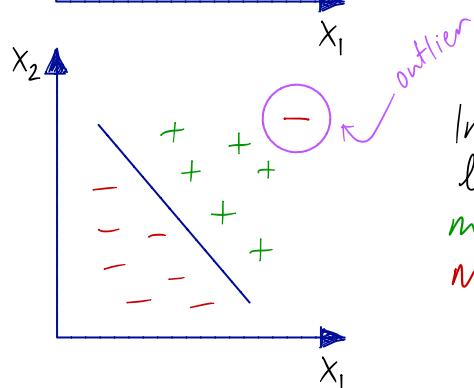
Support Vector Machines

* Margin vs Outliers

Support Vector Machines prioritise classifying data correctly over maximising the margin

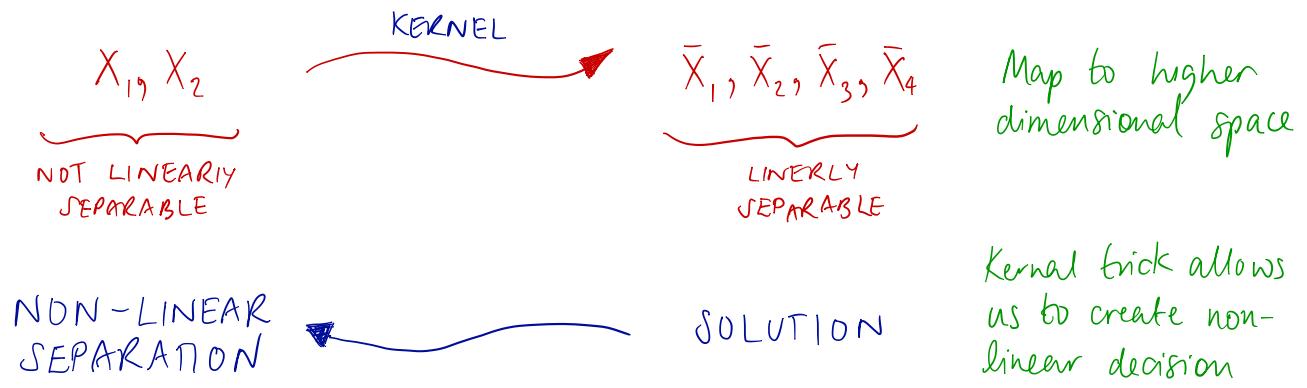


So in this example the SVM would choose the blue decision boundary rather than the purple one



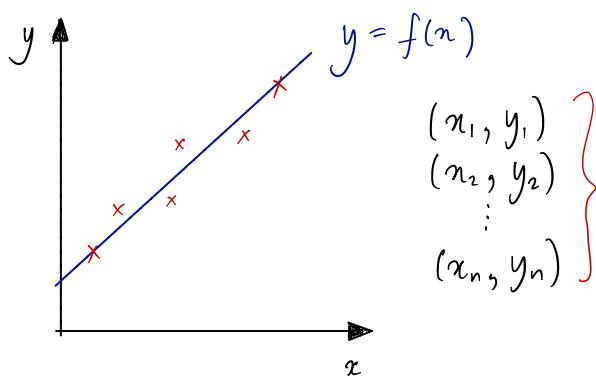
In the case where the data is not linearly separable, SVM will maximise the margin while minimising the classification error

* Kernel Trick



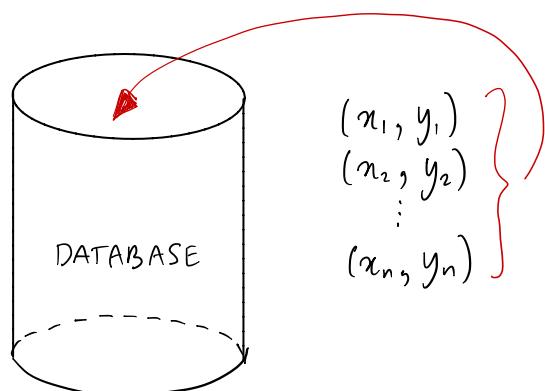
Instance Based Learning

* Instance Based Learning Versus Regression



We use our data to learn a function f which maps our input to our output

Once we have trained our 'machine' we no longer need our original data



In contrast, in Instance Based Learning we keep all our original data in a database.

When we want to make a prediction, we simply lookup the value in our database

Advantages

- Reproduces 'training data'
- Fast
- Simple

Disadvantages

- Overfitting
- Does not generalise well

What if the x we want to predict y for is not in the database? → We use k nearest neighbours...

GIVEN: Training Data
Distance Metric
Number of neighbours
Query point

$$\rightarrow \text{NN} = \{i : d(q, x_i) \text{ k smallest}\}$$

$\rightarrow \text{RETURN}$

- CLASSIFICATION
- REGRESSION

$$D = \{(x_i, y_i)\}$$

$d(x_i, x_j)$ How do we define similarity?

k How many neighbours?

q

What do we do with our KNN?

Weighted	Vote?	How to break ties?
Weighted	Mean?	$\frac{1}{d(q, x_i)}$?

Let's make decisions!

* Given n sorted data points

		Running Time	Space
1NN	Learning	1	n
	Query	$\log_2 n$ <small>Binary Search</small>	1
KNN	Learning	1	n
	Query	$\log_2 n + k$	1
linear regression	Learning	n	1
	Query	1	1

* KNN Bias

Preference Bias - Our belief about what makes a good hypothesis

- Locality \rightarrow Near points are similar i.e. distance metric
- Smoothness \rightarrow Averaging over K points
- All features matter equally \rightarrow

* Some other stuff

\rightarrow distance metric need not be continuous, could be discrete

$\rightarrow k = n$ e.g.

- weighted average
- locally weighted regression

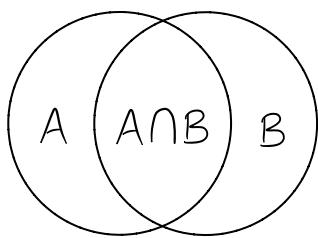
\hookrightarrow Suppose we do locally weighted linear regression. We are able to build more complicated curves using just straight lines

could replace these with:

- Decision Tree
- Neural Network
- Regression

Naive Bayes

* Terminology



Bayes Theorem :

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{\underbrace{P(A)}_{\text{Prior Probability}} \underbrace{P(B|A)}_{\text{Likelihood function: Probability of the evidence given } B}}{\underbrace{P(B)}_{P(\text{cancer} | -ve)}} \quad \begin{array}{l} \text{e.g. } P(\text{cancer}) \\ \text{e.g. } P(\text{cancer test} | \text{Cancer}) \end{array}$$

AKA
 $P(B|A)$ is
Sensitivity

AKA
 $P(\text{cancer} | \text{cancer test})$ is
Specificity

Naive Bayes methods are a set of supervised learning algorithms based on applying Bayes' theorem with the 'naive' assumption of independence between every pair of features.

Bayes Theorem :

$$P(y|x_1, x_2, \dots, x_n) = \frac{P(y) P(x_1, x_2, \dots, x_n | y)}{P(x_1, x_2, \dots, x_n)}$$

Naive Assumption :

$$P(x_i | y, x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n) = P(x_i | y)$$

$$\Rightarrow P(y|x_1, x_2, \dots, x_n) = \frac{P(y) \prod_{i=1}^n P(x_i | y)}{P(x_1, x_2, \dots, x_n)} \quad \begin{array}{l} \text{constant for a given} \\ \text{set of input features} \end{array}$$

$$\hat{y} = \operatorname{argmax}_y \left\{ P(y) \prod_{i=1}^n P(x_i | y) \right\}$$

We then use maximum a posteriori (MAP) i.e. the largest posterior given all the priors, to estimate $\hat{P}(y)$ and our training data to calculate $P(x_i | y)$.

Different Naive Bayes classifiers differ in assumptions they make regarding the distribution of $P(x_i | y)$.

Bayesian Learning

* Introduction

Learn the best hypothesis given data i.e. $\operatorname{argmax}_{h \in H} P(h|D)$

most probable
 h
 D

For each $h \in H$ and a given dataset D we have

$$\text{Bayes Rule : } P(h|D) = \frac{P(D|h) P(h)}{P(D)}$$

affects all computations equally

Maximum A Posteriori (MAP) : Maximum posterior given all our priors

$$h_{\text{MAP}} = \operatorname{argmax}_h \{P(h|D)\}$$

maximum probability hypothesis given our data across all hypotheses

Maximum Likelihood or Maximum A Priori : Maximum prior given our

$$h_{\text{ML}} = \operatorname{argmax}_h \{P(D|h)\}$$

If we assume that all hypotheses are equally probable i.e. uniformly distributed then we have that $h_{\text{MAP}} = h_{\text{ML}}$

* Bayesian Learning in action

1. Given $\{(x_i, d_i)\}$ as noise free examples of c

$$2. c \in H \quad c(x_i) = d_i$$

$$3. \text{Uniform Prior i.e. } P(h) = \frac{1}{|H|} \quad P(h|D) = \frac{P(D|h) P(h)}{P(D)}$$

$$P(D|h) = \begin{cases} 1 & \text{if } d_i = h(x_i) \quad \forall x_i, d_i \in D \\ 0 & \text{otherwise} \end{cases} \quad = \underbrace{\frac{1/|H|}{|VS|/|H|}}_{h \in |VS|} = \frac{1}{|VS|}$$

$$P(D) = \sum_{h \in H} P(D|h_i) P(h_i) = \sum_{h \in VS_{h,D}} 1/|H| = \frac{|VS|}{|H|}$$

Space of hypotheses consistent with D

* Bayesian Learning

- Given $\{(x_i, d_i)\}$

$$h_{ML} = \operatorname{argmax}_h \{P(h|D)\}$$

$$d_i = f(x_i) + \epsilon_i$$

$$= \operatorname{argmax}_h \{P(D|h)\}$$

$$\epsilon_i \sim N(0, \sigma^2) \text{ i.i.d}$$

$$= \operatorname{argmax}_h \prod_i P(d_i|h)$$

$$P(\epsilon_i < x) = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^x e^{-\frac{1}{2}\frac{x^2}{\sigma^2}} dx$$



$$\begin{aligned} h_{ML} &= \operatorname{argmax}_h \prod_i \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(d_i - h(x_i))^2} \\ &= \operatorname{argmin}_h \sum_i (d_i - h(x_i))^2 \quad \text{ssse} \end{aligned}$$

So to find the best hypothesis assuming normally distributed errors we just minimise the sum of squared errors.

* Minimum Description Length

$$\begin{aligned} h_{MAP} &= \operatorname{argmax} \{P(D|h)P(h)\} \\ &= \operatorname{argmax} \{\log_2 P(D|h) + \log_2 P(h)\} \\ &= \operatorname{argmin} \underbrace{-\log_2 P(D|h)}_{\text{like: error}} \underbrace{-\log_2 P(h)}_{\text{like: 'size' (complexity) of } h.} \end{aligned}$$

event A with probability p has length $-\log_2 p$
(Information Theory)

- ← e.g.
• depth of decision tree
• neural network weights

Occlusion Razor

Trade off between 'size' and training error

* Bayesian Classification

Given a new feature vector x what's the most probable classification?

Unfortunately $h_{MAP}(x)$ doesn't always give the most probable classification

Bayes optimal classifier :

$$\underset{v_i \in V}{\operatorname{argmax}} \sum_{h_i \in H} P(v_i | h_i) P(h_i | D)$$

Set of all possible classifications assigned by H

Example :

$$P(h_1 | D) = 0.4 \quad h_1(x) = + \quad \left. \begin{array}{l} P(-|h_1) = 0 \\ P(+|h_1) = 1 \end{array} \right\}$$

$$P(h_2 | D) = 0.3 \quad h_2(x) = - \quad \left. \begin{array}{l} P(-|h_2) = 1 \\ P(+|h_2) = 0 \end{array} \right\}$$

$$P(h_3 | D) = 0.3 \quad h_3(x) = - \quad \left. \begin{array}{l} P(-|h_3) = 1 \\ P(+|h_3) = 0 \end{array} \right\}$$

$$V = \{-, +\} \quad \sum_{h_i \in H} P(-|h_i) P(h_i | D) = 0.6$$

$$\sum_{h_i \in H} P(+|h_i) P(h_i | D) = 0.4$$

Bayes optimal classifier is very costly to compute because the posterior, $P(h_i | D)$, must be computed for each hypothesis $h_i \in H$.

Bayesian Inference

* Conditional Independence

Recall,

$$\text{Independence: } P(X \cap Y) = P(X)P(Y)$$

$$\text{Chain Rule: } P(X \cap Y) = P(X|Y)P(Y) \quad \therefore \quad P(X|Y) = P(X)$$

X is conditionally independent of Y given Z if

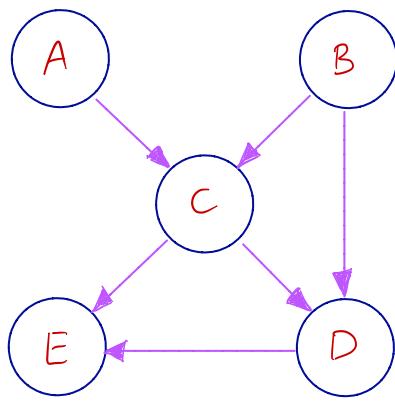
$$\forall x, y, z \quad P(x=x | Y=y \wedge Z=z) = P(x=x | Z=z)$$

or more compactly we write

$$P(x|Y \cap Z) = P(x|Z)$$

* Belief Networks / Bayesian Networks / Graphical Models

* Sampling from and recovering the joint distribution



$$\begin{aligned} A &\sim P(A) \\ B &\sim P(B) \\ C &\sim P(C | A \wedge B) \\ D &\sim P(D | B \wedge C) \\ E &\sim P(E | C \wedge D) \end{aligned}$$

Sampling order should be topological!

Must be acyclic

$$\begin{aligned} P(A \wedge B \wedge C \wedge D \wedge E) \\ = P(A)P(B)P(C | A \wedge B)P(D | B \wedge C)P(E | C \wedge D) \end{aligned}$$

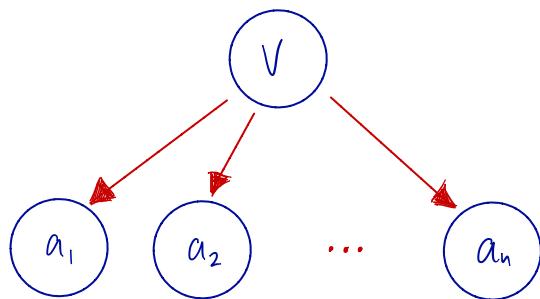
* Inferencing Rules

- Marginalisation: $P(x) = \sum_y P(x \cap y)$

- Chain Rule: $P(x \cap y) = P(x|y)P(y)$

- Bayes Rule: $P(x|y) = \frac{P(x)P(y|x)}{P(y)}$

* Naive Bayes: Special Case - Classify V



e.g. $V = \text{SPAM?}$

$a_1, a_2, \dots, a_n = \text{words}$

value not important

we want to classify V
so only need to know
the most probable value

$$\begin{aligned} P(V | a_1 \cap \dots \cap a_n) &= P(a_1 \cap \dots \cap a_n | V) \underbrace{P(V)}_{\text{P}(a_1 \cap \dots \cap a_n)} \\ &= P(V) \prod_{i=1}^n P(a_i | V) / P(a_1 \cap \dots \cap a_n) \\ \Rightarrow V &= \underset{V}{\operatorname{argmax}} \left\{ P(V) \prod_{i=1}^n P(a_i | V) \right\} \end{aligned}$$

* Why Naive Bayes is Cool

- Inference is cheap
- Few parameters
- Estimate parameters with labelled data
- Connects inference and classification
- Empirically successful

Why does it work?

→ Actual probabilities are not important in classification,
we only want to know the most probable class

→ We calculate our conditional probabilities by counting

$$P(a_i | V) = \frac{\# a_i \cap V}{\# V}$$

Even if the a_i are not independent, the order is preserved.

Ensemble Learning (Bagging & Boosting)

* General Idea

- ① Learn over subset of data \rightarrow rules ← How do you make subsets?
- ② Combine rules ← How do you combine?
Can use any learning algorithm

* Bagging

- ① Uniformly randomly pick data and apply a learner
- ② Take the mean result

A bit like cross validation - reduces variance / overfitting

* Boosting

- ① Subset = 'hardest' examples ?! ← i.e. the ones that get misclassified
- ② Take the weighted mean

$$\text{Error} = \mathbb{P}_{\mathcal{D}}[h(x) \neq c(x)] \quad \text{accounts for distribution of } x, \mathcal{D}$$

Weak Learner: always does better than chance

$$\text{i.e. } \forall \mathcal{D} \quad \mathbb{P}_{\mathcal{D}}[h(x) \neq c(x)] < \frac{1}{2}$$

- Given training examples $\{(x_i, y_i)\} \quad y_i \in \{-1, +1\}$

- For $t=1$ to T

→ Construct D_t :

$$\begin{aligned} D_1 &= \frac{1}{n}, \\ D_{t+1} &= \sum_{i=1}^n D_t(i) \exp\left\{\underbrace{-\alpha_t y_i h_t(x_i)}_{\substack{+1 \text{ when } h \text{ is correct} \\ -1 \text{ else}}}\right\} / Z_t \\ \alpha_t &= \ln(1 - \epsilon_t) / (2\epsilon_t) > 0 \\ Z_t &= \sum_i D_t(i) \exp\left\{-\alpha_t y_i h_t(x_i)\right\} \end{aligned}$$

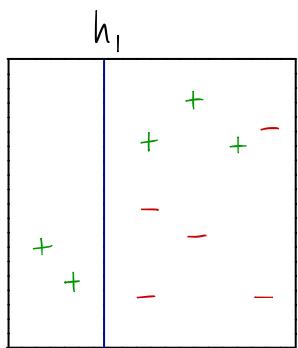
normalisation factor

Harder examples

- Increases the weight on incorrectly classified examples and decreases the weight on correctly classified examples.
- Find weak classifier $h_t(x)$ with small error $\epsilon_t = \mathbb{P}_{D_t}[h_t(x) \neq c(x)]$

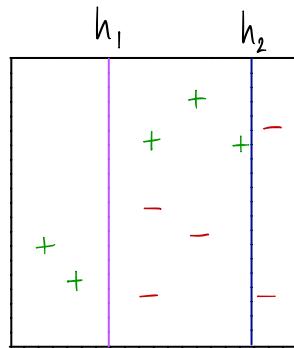
- Output $H_{\text{final}}(x) = \operatorname{sgn} \left(\sum_t \alpha_t h_t(x) \right)$ Weighted mean

* Example $H = \{ \text{axis aligned semi planes} \}$



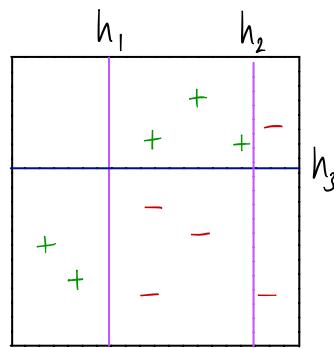
$$\epsilon_t = 0.3$$

$$\alpha_t = 0.42$$



$$\epsilon_t = 0.21$$

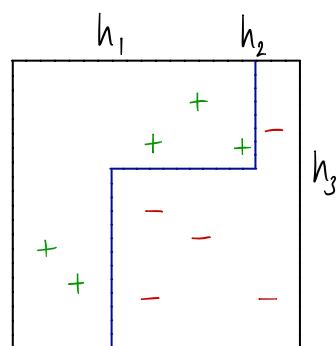
$$\alpha_t = 0.65$$



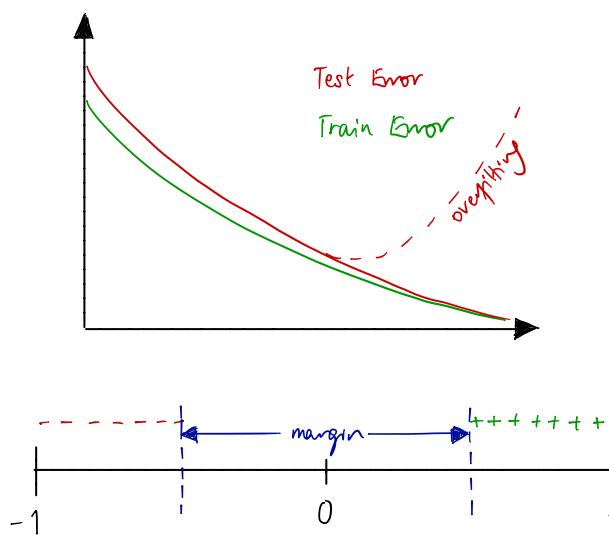
$$\epsilon_t = 0.14$$

$$\alpha_t = 0.92$$

$H_{\text{final}}(x) \rightarrow$



* Overfitting



Boosting does not typically result in the graph we see when overfitting.

$$H_{\text{final}}(x) = \text{sgn} \left(\sum_t \alpha_t h_t(x) \right)$$

$$-1 \leq \sum_t \alpha_t h_t(x) / \sum_t \alpha_t \leq 1$$

Where on the line this lands for any given x expresses how much confidence we have in the classification

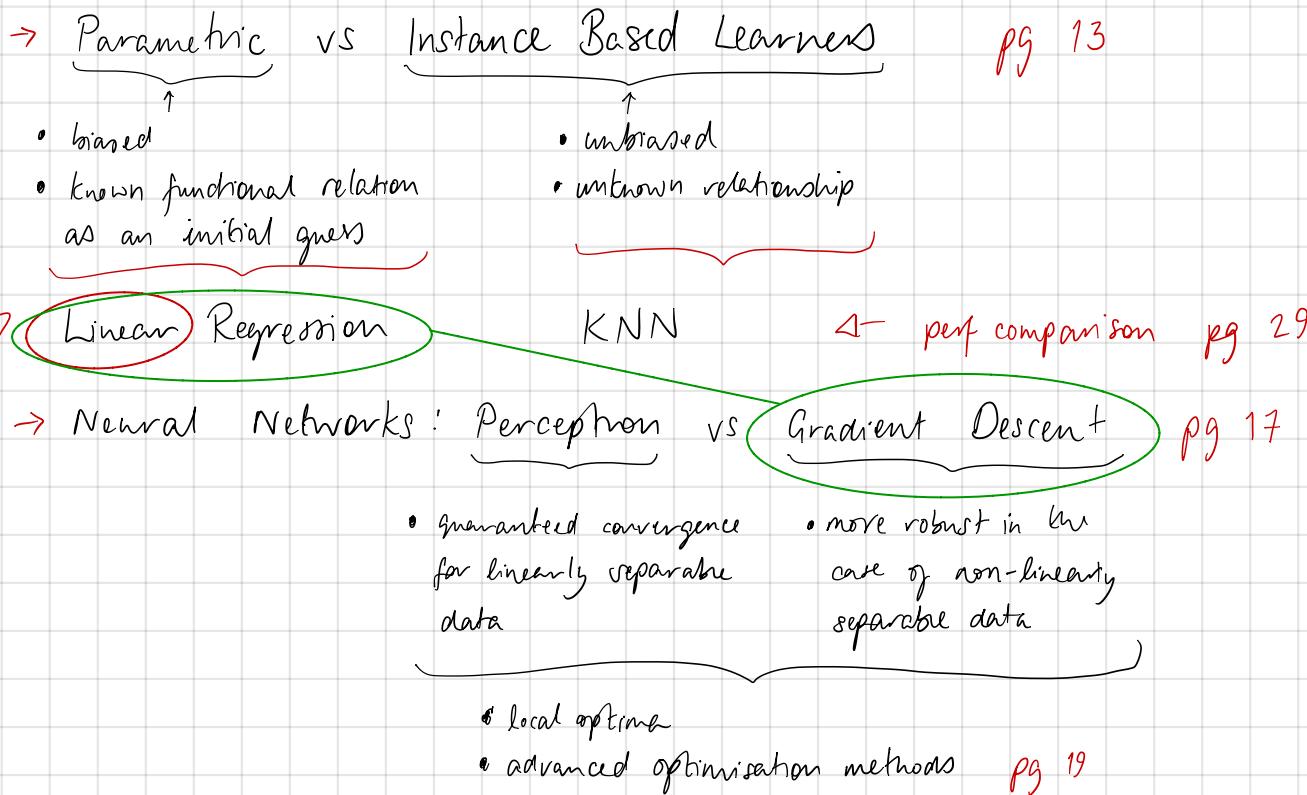
As we increase T we increase the margin

Boosting can overfit e.g. if the learner is a Neural Network with too many layers and nodes

P2: Build a Student Intervention System

* Classifier choices & comparisons from lectures

→ Decision trees



→ Support Vector Machine

→ Naive Bayes

* Discussion points re. performance

- METRICS
- Train vs Query performance KNN vs linear regression
 - F1 score - want to be conservative - intervention is preferable rather than to be avoided
 - baselines?
 - Soft computing - approximate solutions
 - quality improves indefinitely with learning

Comparing performance with complexity for individual algorithms
e.g.

→ Data size

→ Decision Tree depth query time?

* "find the most effective model with the least amount of computation cost (you pay the company by the memory and CPU time you use on their servers)"

Model Evaluation:

1. F1 score
2. Training set size
3. Prediction computational cost }

which do we do more of?

how often do we get more data to train

* Naive Bayes in Laymens terms:

Bayes Theorem:

* Confusion matrix

Actual Class			
		+	-
Predicted Class	+	a True +ve	b False +ve
	-	c False -ve	d True -ve

precision = $\frac{\# \text{ true +ve}}{\# \text{ true +ve} + \# \text{ false +ve}} \leq 1$
 $a/(a+b)$

recall = $\frac{\# \text{ true +ve}}{\# \text{ true +ve} + \# \text{ false -ve}} \leq 1$
 $a/(a+c)$

$F_1 = \frac{2 \text{ precision recall}}{\text{precision} + \text{recall}}$

$F_1 = \frac{2 \frac{\# \text{ true +ve}}{\# \text{ true +ve} + \# \text{ false +ve}} \frac{\# \text{ true +ve}}{\# \text{ true +ve} + \# \text{ false -ve}}}{\frac{\# \text{ true +ve}}{\# \text{ true +ve} + \# \text{ false +ve}} + \frac{\# \text{ true +ve}}{\# \text{ true +ve} + \# \text{ false -ve}}}$

$= \frac{2(\# \text{ true +ve})^2}{\# \text{ true +ve}(2 \# \text{ true +ve} + \# \text{ false +ve} + \# \text{ false -ve})}$

$= \frac{2 \# \text{ true +ve}}{(2 \# \text{ true +ve} + \# \text{ false +ve} + \# \text{ false -ve})}$

$$\text{bad} = 0 \leq F_1 \leq 1 = \text{good}$$

■ Baselines for performance metrics

→ We are interested in identifying students which might fail and thus require intervention so choosing 'fail' or 'intervene' as our **true** class gives the more meaningful F_1 score

→ For a given accuracy score we prefer false +ves to false -ves so in the best case $\text{recall} = 1$ but in general we want $\text{precision} < \text{recall}$

→ Our data is imbalanced i.e. pass rate $\approx 67\%$. So, $\approx 1/3$ of the students require intervention

$$\text{Let pass rate} = p$$

■ suppose we always predict pass

$$\left. \begin{array}{l} \text{accuracy} = p \left(= \frac{2}{3}\right) \\ \text{recall} = 0 \\ \text{precision} = \frac{1}{3} \end{array} \right\} F1 = 0$$

0	0
1-p	p

■ suppose we always predict fail

$$\left. \begin{array}{l} \text{accuracy} = \text{precision} = 1-p \left(= \frac{1}{3}\right) \\ \text{recall} = 1 \\ \Rightarrow F1 = \frac{2(1-p)}{2-p} \quad \left(= \frac{2 \cdot \frac{1}{3}}{\frac{4}{3}} = \frac{1}{2}\right) \end{array} \right.$$

1-p	p
0	0

■ randomly choose pass/fail - choose pass q of the time

$$\begin{aligned} \Rightarrow \text{accuracy} &= (1-q)(1-p) + qp \\ &= 1-q-p + 2qp \\ &= \frac{1}{3} - q \left(1 - \frac{4}{3}\right) = \frac{1}{3} (1+q) \\ \Rightarrow \text{precision} &= \frac{(1-q)(1-p)}{(1-q)(1-p) + p} = 1-p \end{aligned}$$

Actual Class		1	0
Predicted Class	1	$(1-q)(1-p)$	$(1-q)p$
	0	$q(1-p)$	qp
	True -ve	Falre -ve	True -ve

$$\Rightarrow \text{recall} = \frac{(1-q)(1-p)}{(1-p)(1-q+p)} = 1-q$$

$$\Rightarrow F1 = \frac{2(1-p)(1-q)}{2-p-q} = \frac{\frac{2}{3}(1-q)}{\frac{4}{3}-q} = \frac{2(1-q)}{4-3q}$$

Now for a given accuracy we prefer false -ve to false +ve. We note $p=q$ gives

$$\text{precision} = \text{recall} = \frac{1}{3}$$

$$\text{accuracy} = \frac{5}{9} > \frac{1}{2}$$

$$F1 = \frac{\frac{2}{3}}{\frac{4}{3}-2} = \frac{1}{3}$$

* Relationship between accuracy and f1 score

Actual Class		
+	-	
Predicted Class	+	a b True +ve False +ve
	-	c d False -ve True -ve

$$\text{accuracy} = \frac{a+d}{a+b+c+d}$$

$$\text{precision} = \frac{a}{a+b} \quad \text{recall} = \frac{a}{a+c}$$

$$\text{f1 score} = \frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

$$f = \frac{2pq}{p+q} \Rightarrow \frac{df}{dp} = \frac{(p+q)2q - 2pq}{(p+q)^2} \\ = \frac{2q^2}{(p+q)^2} > 0$$

- Increase accuracy

Case (i) increase a

- ⇒ decrease either b (or c)
- ⇒ increase precision (or recall)
- ⇒ increase f1 score

← numerator increases
denominator unchanged

Case (ii) increase d

- ⇒ decrease either b (or c)
- ⇒ increase precision (or recall)
- ⇒ increase f1 score

← numerator unchanged
denominator decreases

- Fixed accuracy

Case (i) increase a decrease d

- ⇒ increase precision and recall
- ⇒ increase f1 score

← numerator and denominator increase by same amount

In conclusion an increase in accuracy ⇒ increase in f1 score

but

an increase in f1 score $\not\Rightarrow$ increase in accuracy

Topic one

Topic two

Topic three

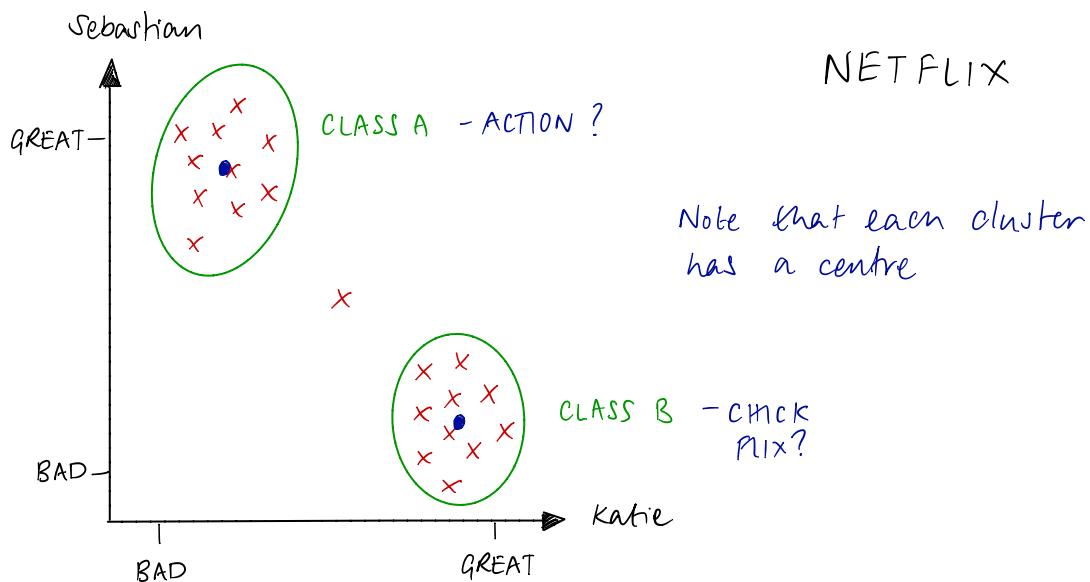
Topic four

Topic five

C3: Unsupervised Learning

Clustering

* Example - Clustering Movies



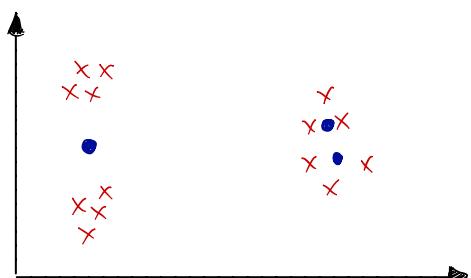
* K means

General Idea

- Decide how many clusters there are
- Assign initial cluster centres at random
- Assign each point to the cluster centre it is closest to
- Move the cluster centres to the location which minimises the sum of square distances to their centres
- **ITERATE**

Possible Problems

- The final solution can be dependent on the initial conditions
 - This can happen if there are no clear clusters e.g. if the data is uniformly spread in space
 - Local minima:



On the left is an example of a stable solution

"Local hill climbing algorithm"

sklearn.cluster.KMeans

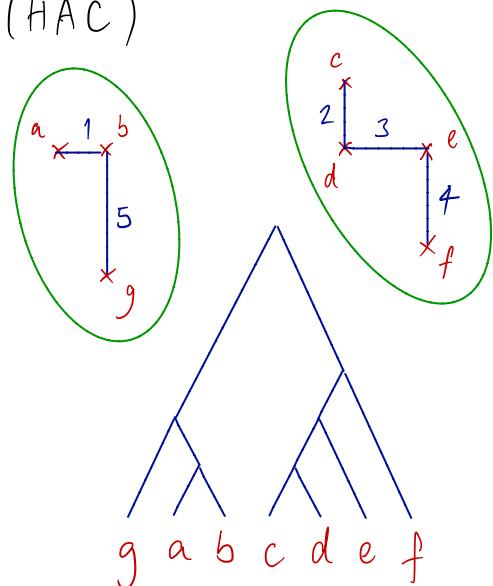
Parameters :
n_clusters -
n_init - no. initialisations
max_iter -

* Single Linkage Clustering (SLC)

Hierarchical Agglomerative Clustering (HAC)

Algorithm :

- consider each object a cluster (n objects)
- define 'intercluster distance' as the distance between the closest two points in the two clusters mean / median
- Merge two closest clusters
- Repeat n-k times to make k clusters

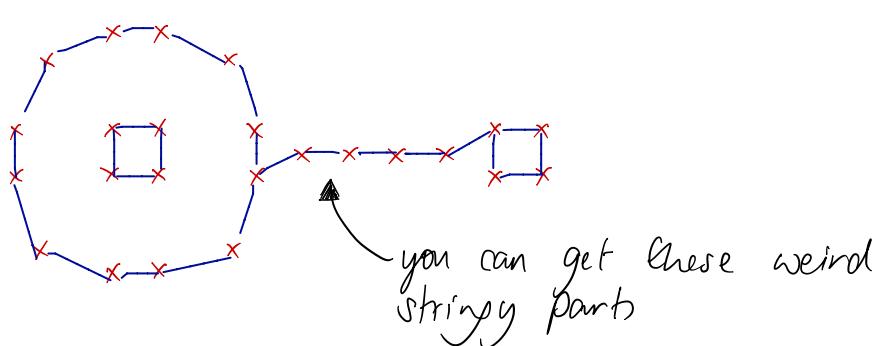


Properties

- Deterministic - you always get the same answer
- If distance = edge length in a graph then this is the same as minimum spanning tree algorithm
- Running time = n^3
 - compute all distances : $\sum_{i=1}^n i = n(n+1) = O(n^2)$
 - compute min at each step : $O(n)$

Issues with single link clustering

Example : $k = 2$



* Soft Clustering

$\times \quad \times \quad \times \quad \quad \quad \times \quad \quad \quad \times \quad \times \quad \times$

Here points can be 'shared' between clusters...

Algorithm:

Assume the data was generated by

1. Select one of k Gaussians [fixed known variance] uniformly
2. Sample x_i from that Gaussian
3. Repeat n times

Task:

Find a hypothesis $h = (\mu_1, \mu_2, \dots, \mu_k)$ that maximizes the probability of the data (Maximum Likelihood)

Maximum Likelihood Gaussian

The maximum likelihood mean of the Gaussian is just the mean of the data.

What if we want to find k means

\rightarrow we use hidden variables

\rightarrow each data point x is really $(x, z_1, z_2, \dots, z_k)$

indicates which cluster x belongs to

Expectation Maximisation

constant for a given i

$$\text{Bayes Theorem: } P(\mu = \mu_j | x = x_i) = \frac{\overbrace{P(\mu = \mu_j)}^z P(x = x_i | \mu = \mu_j)}{P(x = x_i)}$$

$$\mathbb{E}[z_{ij}] = \frac{P(x = x_i | \mu = \mu_j)}{\sum_{j=1}^k P(x = x_i | \mu = \mu_j)}$$

Probability the
data point
belongs to the
cluster

Expectation
(define z from k)

$$\mu_j = \frac{\sum_i \mathbb{E}[z_{ij}] x_i}{\sum_i \mathbb{E}[z_{ij}]}$$

Maximisation
(define μ from z)

$$P(x = x_i | \mu = \mu_j) = e^{-\frac{1}{2}\sigma^2 (x_i - \mu_j)^2}$$

* Properties of Expectation Maximize

- Monotonically non-decreasing likelihood
- Doesn't converge (practically does)
- Will not diverge
- Can get stuck - local minima
- Works with any distributions (provided Exp and Max are solvable)

* Clustering Properties

Clustering algorithms essentially partition data

D = distance matrix

P_D = clustering scheme under distance matrix D

C = clustering or partition

Richness: For any assignment of objects to clusters, there is some distance matrix D such that P_D returns that clustering i.e.

$$\forall C \exists D \mid P_D = C$$

Scale-Invariance: Scaling distances by a positive value does not change the clustering
i.e.

$$\forall D, k \quad P_{kD} = P_D$$

Consistency: Shrinking intracluster and expanding intercluster distances does not change the clustering
i.e.

$$P_D = P_{D'}$$

Examples :

Single Link Clustering which Proprieties when ...

... $\frac{1}{2}$ clusters reached



... clusters are Θ units apart



... clusters are Θ/w units apart
where $w = \max_{i,j} D(i,j)$



* Impossibility Theorem - Kleinberg

No clustering scheme can have all three properties

1. Richness
2. Scale-invariance
3. Consistency

* Feature Scaling

$$X' = \frac{X - X_{\min}}{X_{\max} - X_{\min}} \quad 0 \leq X' \leq 1$$

Algorithms requiring feature scaling quiz

Decision Trees

SVM

Linear Regression

K-Means Clustering

Both these
algorithms
involve calculation
of distances,

Feature Selection

* Motivation

→ Insight into our data

→ Curse of dimensionality

Want an algorithm which takes our full feature set, and returns a subset of features which are 'important'

* Algorithms

Actually we want some function which takes a feature set and returns a score for each feature

Suppose we have a dataset with n samples and m features our function f has the type

$$f: \mathbb{R}^{(n \times m)} \rightarrow \mathbb{R}^m$$

How hard is this problem?

We want a subset of k features from our full set of m features. There are $C(k, m)$ combinations of features we could choose. Recall

$$C(k, m) = \frac{m!}{k!(m-k)!}$$

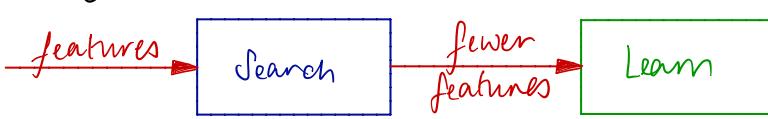
Suppose we don't know how many features to choose, then we also have to optimise over all possible subsets of which there are $\sum_{k=1}^m C(k, m) = 2^m$

possible subsets.

So this problem is exponential in the number of features.

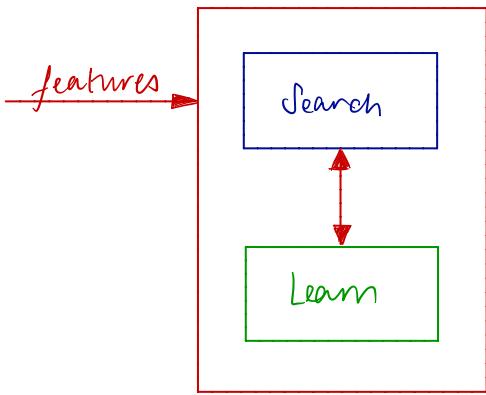
* Filtering & Wrapping

Filtering:



- + Speed
- Isolated features
- Ignores the learning problem / bias

Wrapping:



- + Takes into account model bias
- + Learns
- Slow

Feature selection is done by decision trees - it looks for features which result in the most INFORMATION GAIN.
We can use a decision tree in filtering

For filtering we could look at

- Information gain
 - Variance or entropy
 - "Useful" features
 - Linearly independent features
- } Domain knowledge

For wrapping we could look at

- Hill climbing algorithms
- Randomized optimisation
- Forward searching
- Backward searching

* Relevance

- x_i is strongly relevant if removing it degrades the Bayes Optimal Classifier
- x_i is weakly relevant if
 - it is not strongly relevant
 - ∃ a subset of features S | adding x_i to S improves the Bayes Optimal Classifier
- x_i is otherwise irrelevant

Bayes optimal classifier: $\operatorname{argmax}_{v_i \in V} \sum_{h_i \in H} P(v_i | h_i) P(h_i | D)$

set of all possible classifications assigned by H

Example : Smallest feature set sufficient to get zero training error

a	b	c	d	e	label
0	0	1	1	1	-
0	1	1	1	1	-
1	0	1	0	0	-
1	1	1	0	0	+

Strongly relevant irrelevant Weakly relevant

Decision Tree :

$$a \text{ AND } b \Rightarrow +$$

Perceptron :

$$a + b - c > 0 \Rightarrow +$$

$$e = \text{NOT } a$$

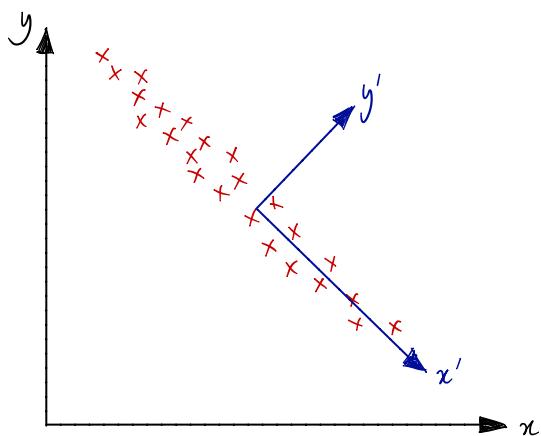
* Relevance vs Usefulness

- Relevance measures the effect on the Bayes Optimal Classifier
~ Information
- Usefulness measures the effect on a particular predictor
~ Error / model

So in the above example feature c is useful for perceptron but not for the decision tree

Principal Component Analysis - PCA

* What is Principal Component Analysis



Given data of any shape, PCA finds a new co-ordinate system, obtained from the original through translation and rotation only. The centre of the new co-ordinate system is at the centre of the data and the x -axis is aligned with the principal axis of variation.

* Measurable vs. Latent Features

Example: Given the features of a house, what is its price?

Measurable Features

- Square footage
- No. of rooms
- School ranking
- Neighbourhood safety

Latent Features

- { - Size
- { - Neighbourhood

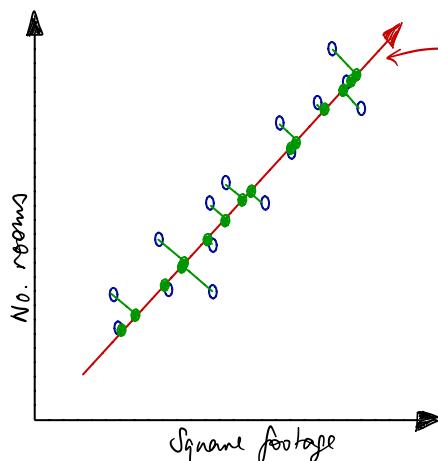
* PCA motivation

- Many features but suspect a smaller number of features may be driving the patterns
- Try making a composite feature ^{principal component} that more directly probes the underlying phenomenon

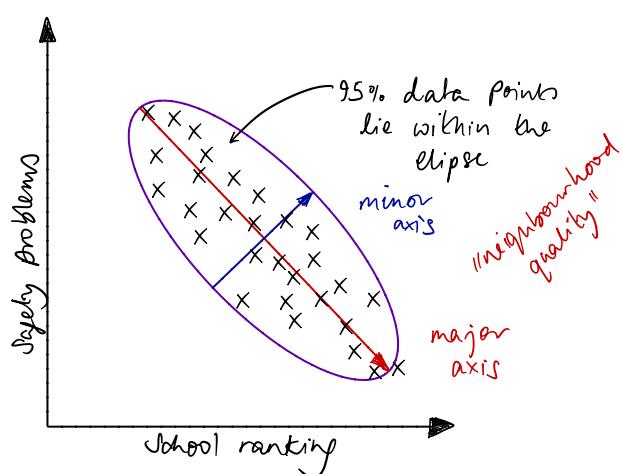
Context

- dimensionality reduction
- unsupervised learning

* Example: Square footage + No. rooms \rightarrow size



We project our data onto the principal axis of variation to reduce the number of dimensions from two dimensions to one.



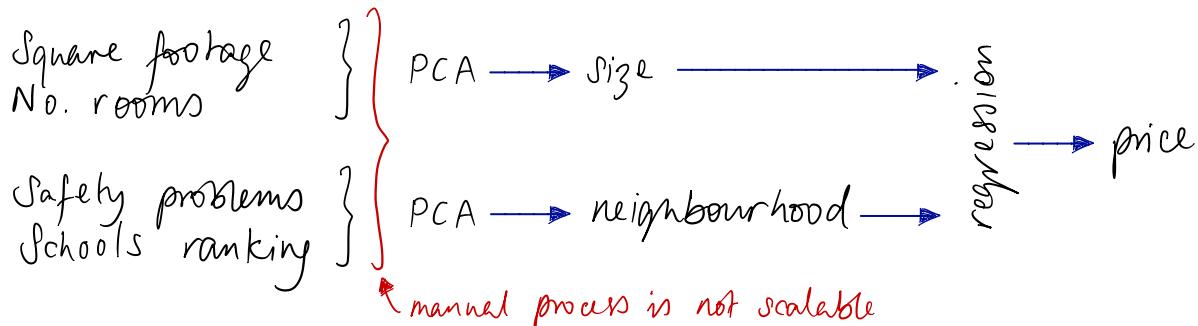
Principal component of a dataset is the direction that has the largest variance.

This retains the maximum amount of information in the original data

Note that the amount of information lost in projecting onto a given axis is proportional to the sum of the distances of the data points from the axis.

Projection onto the direction of maximal variance minimises the information loss

* PCA for Feature Transformation



Actually we could perform PCA on all four features and reduce our features from four to two

* PCA summary

- Systemized way to transform input features into principal components
- Use principal components as new features
- Principal components are directions in the data that maximise variance (minimise information loss) when you project/compress down the data onto them
- The more variance of data along a principal component, the higher that component is ranked
- Principal components are all orthogonal so don't overlap
- Maximum number of principal components
= Number of input features

* When to use PCA

- Latent features driving the patterns in the data
- Dimensionality reduction
 - Visualise high dimensional data
 - Reduce noise
- Make other algorithms (regression, classification) work better by reducing dimensionality
 - e.g. eigenfaces - PCA for facial recognition
 - dimensionality reduction before SVM

* PCA for Facial Recognition

- Pictures of faces generally have high input dimensionality (many pixels)
- Faces have general patterns that could be captured in a smaller number of dimensions (two eyes on top, a mouth on the bottom for example)

Feature Transformation

* Introduction

Feature Transformation: The problem of pre-processing a set of features to create a new (smaller? / more compact?) feature set, while retaining as much (relevant? / useful?) information as possible i.e.

$$f(\underline{x}) : \mathbb{F}^n \rightarrow \mathbb{F}^m$$

Usually,

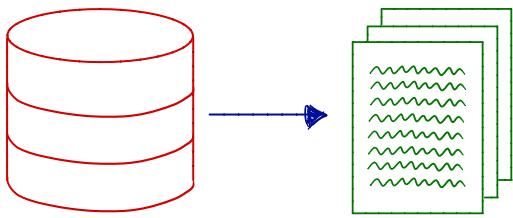
$$m < n \quad \text{and} \quad f(\underline{x}) = P^T \underline{x}$$

$P = n \times m$ matrix
 \Rightarrow new features
are a linear
combination of
original features

N.B. Feature selection is a subset of feature transformation

* Motivation

Example: Information Retrieval - ad hoc query
(Google document retrieval)



Features \rightarrow words & counts
 \uparrow good indicators
problem: \uparrow large feature set

More problems: \rightarrow POLYSEMY: Words can have multiple meanings
 \rightarrow FALSE +VES

\rightarrow Can combine words which eliminate ambiguity

\rightarrow SYNONYMY: Multiple words can have the same meaning
 \rightarrow FALSE -VES

\rightarrow Can combine words with similar meanings into a single feature

[Aside // LISP is a superset that subsumes all languages including natural languages?!]

* Independent Components Analysis (ICA)

Recall PCA looks for correlation to maximise variance and transform features to an ordered set

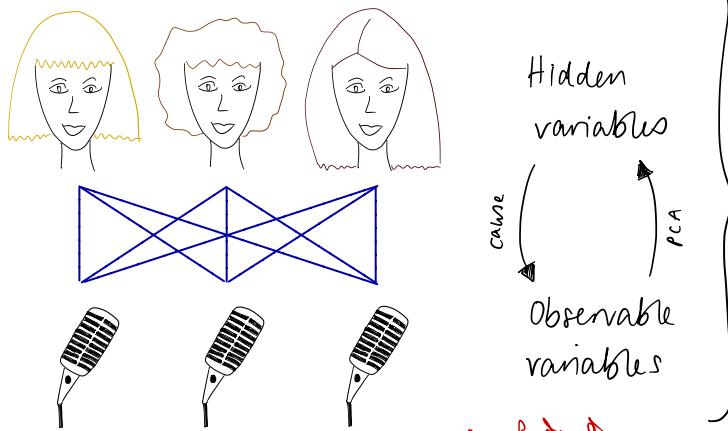
ICA instead maximises independence

ICA transforms features to a statistically independent set of new features

$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \rightarrow \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix} \quad m \leq n \quad y_i \perp y_j \quad \forall i \neq j$$

1) $I(y_i, y_j) = 0$
 2) Minimise $I(\underline{x}, \underline{y})$

Example: Blind Source Separation (Cocktail Party)



There are lots of people talking at a cocktail party and as many microphones dotted around the room. Each mic receives a linear combination of the voices. ICA would allow us to separate the sounds from each person.

* PCA vs. ICA

PCA

- Mutually orthogonal
- Maximal variance
- Ordered features

In the case where the hidden variables are independent and normally distributed

ICA

- Mutually independent
- Minimal mutual information
- Bag of features

Between the original data and the transformed data

Notes:

- Uncorrelated is not the same as statistically independent except for under some very specific distributions, for example the normal distribution
- Distribution which maximises variance is normal so in the case where our hidden variables are normally distributed, ICA = PCA
- Suppose our hidden variables are not normally distributed but are independent, as the number of variables increases the distribution of the linear sum tends to normal (by central limit theorem). Now PCA in this case just looks for normally distributed components so will just find the linear sum. ICA does not assume hidden variables to be normally distributed. For this reason PCA performs poorly for the Blind source separation problem.

* PCA vs. ICA continued ...

PCA

- Performs poorly on BSS
- Non-directional:
 - Gives the same result if the transpose of the transform is used
- Eigenfaces:
 - Brightness
 - Average face
- Linear Algebraic approach

ICA

- Performs well on BSS
 - Directional:
 - Result for the transpose transform differs
 - Eigenfaces
 - Nose } facial
 - Eyes } features
 - Natural Scenes
 - Edges
 - Documents
 - Topics
 - Probabilistic approach
- } finds structure

* Alternatives

- Random Component Analysis (RCA):
 - Generates random directions i.e. P^T is random
 - $n \rightarrow m$ features $m < n$ (deals with the curse of dimensionality problem)
↑
 m is not as small as it would be for PCA
 - RCA is very fast compared to PCA and PCA
- Linear Discriminant Analysis (LDA)
 - Finds a projection that discriminates based on the label
 - Not Latent Dirichlet Allocation which is another approach

New Developments in Machine Learning

* Big Data

Algorithmic challenges from gigantic datasets where even linear algorithms can be slow

* Deep Learning

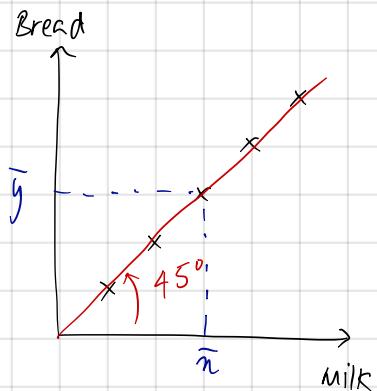
For a long time it's been unclear how adding layers to neural networks can improve performance but recent developments include new techniques for getting signals through multiple layers

* Semi-Supervised Learning

Extract (using unsupervised methods) enough structure that when combined with the labelled data can mimic the situation here you have much more labelled data

P3: Customer Segments

* PCA example :



General rotation by φ :

$$\begin{pmatrix} \cos \varphi & -\sin \varphi \\ \sin \varphi & \cos \varphi \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} x \cos \varphi - y \sin \varphi \\ x \sin \varphi + y \cos \varphi \end{pmatrix}$$

$$\cos -\frac{\pi}{4} = -\sin -\frac{\pi}{4} = \frac{1}{\sqrt{2}}$$

$$\Rightarrow \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} x+y \\ -x+y \end{pmatrix} \quad \rightarrow \uparrow \quad \leftarrow \uparrow$$

point $(1,1)$ goes to $(1,0)$

Topic one

Topic two

Topic three

Topic four

Topic five

C4: Markov Decision Processes

Markov Decision Processes

* Decision Making & Reinforcement Learning

So far we've looked at supervised and unsupervised learning. Reinforcement learning is inspired in part by the behaviourist school of thought in psychology. It is a learning system where rewards and punishments are assigned to different scenarios with different weights to dictate prioritisation.

→ Supervised Learning : $y = f(x)$

Function approximation - Given a set of x and y and a function f that will map some new x to a proper y .

→ Unsupervised Learning : $f(x)$

Clustering description - Given a set of n find an f that gives a compact description of that set of n

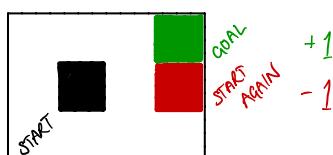
→ Reinforcement Learning : $y = f(x)$

Given x and z find y and f
 $\uparrow s \quad \uparrow r \quad \uparrow a \quad \uparrow \pi^*$

* Markov Decision Processes

PROBLEM

STATES : s



ACTIONS : $A(s), A$ e.g. U/D/L/R

MODEL : $T(s, a, s') \sim P(s' | s, a)$

REWARD : $R(s), R(s, a), R(s, a, s')$
Temporal Credit Assignment

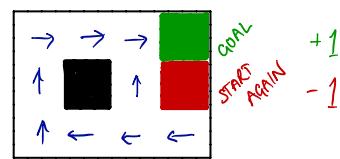
{ Probability of moving in the direction you choose is 80%. There is a 10% chance of moving in either direction parallel to that

e.g. +1, -1 else -0.01

SOLUTION

POLICY : $\pi(s) \rightarrow a$

π^* optimal policy maximises the reward you receive or expect to receive over your lifetime



We consider the following types of problems:

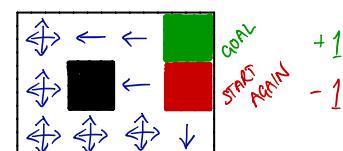
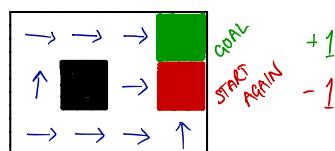
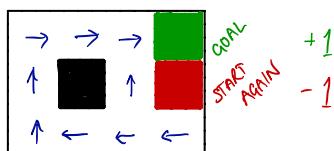
- Markovian Property - Only the present matters
- Stationary models - Time homogeneous

The reward function matters a lot :

reward = -0.04

reward = -2 < 1.6284

reward = +2



* Sequences of rewards

- Stationary Assumption \Rightarrow Infinite time horizon
i.e. $\pi(s) \rightarrow a$ not $\pi(s, t) \rightarrow a$
- Stationary Preferences \Rightarrow Utility of sequences (of states)

A utility is a function which maps a sequence of states to a real number. It allows us to express the preference of one sequence over another.

So, preferences being stationary means

if $u(s_0, s_1, \dots) > u(s_0, s'_1, \dots)$ then $u(s_1, s_2, \dots) > u(s'_1, s'_2, \dots)$

Suppose rewards are all positive, $R(s_t) \geq 0 \forall t$, we can't just sum rewards for our utility function because the sum could be infinite:

$$u(s_0, s_1, \dots) = \sum_{t=0}^{\infty} R(s_t)$$

We have to discount them:

$$U(s_0, s_1, \dots) = \sum_{t=0}^{\infty} \gamma^t R(s_t), \quad 0 \leq \gamma < 1$$

$$U(s_0, s_1, \dots) \leq \sum_{t=0}^{\infty} \gamma^t R_{\max} = \frac{R_{\max}}{1-\gamma}$$

* The optimal policy and Bellman's Equation

The optimal policy is given by

$$\pi^* = \operatorname{argmax}_{\pi} E \left(\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi \right)$$

$$U^\pi(s) = E \left(\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi, s_0 = s \right)$$

$$\pi^*(s) = \operatorname{argmax}_a \left(\sum_{s'} T(s, a, s') U^{\pi^*}(s') \right)$$

$$U^{\pi^*}(s) = R(s) + \gamma \max_a \left(\sum_{s'} T(s, a, s') U^{\pi^*}(s') \right) \quad \text{Bellman's Equation}$$

$\hookrightarrow n$ non-linear equations in n unknowns ($n = \text{number of states } s$)

* Value Iteration

The solution of Bellman's equation can be found using the following algorithm:

- Start with some arbitrary utilities
- Update utilities based on neighbours:

$$\hat{U}_{t+1}^{\pi^*}(s) = R(s) + \gamma \max_a \left(\sum_{s'} T(s, a, s') \hat{U}_t^{\pi^*}(s') \right)$$

- Repeat until converged

Example: Find $U_1^{\pi^*}(x)$ and $U_2^{\pi^*}(x)$

$t=0$	$s=l$ $R=-0.04$ $u=0$	$s=x$ $R=-0.04$ $u=0$	$s=r$ $R=+1$ $u=+1$
		$s=d$ $R=-0.04$ $u=0$	$R=-1$ $u=-1$
	$R=-0.04$ $u=0$	$R=-0.04$ $u=0$	$R=-0.04$ $u=0$

$$\begin{cases} \gamma = 1/2 \\ R(s) = -0.04 \\ U_0(s) = 0 \end{cases}$$

unless marked otherwise

Transition matrix for $s=x$:

$$T(s, a, s') = \begin{pmatrix} 0.8 & 0.1 & 0 & 0.1 \\ 0.1 & 0.8 & 0.1 & 0 \\ 0 & 0.1 & 0.8 & 0.1 \\ 0.1 & 0 & 0.1 & 0.8 \end{pmatrix}$$

$\xrightarrow{e} \xleftarrow{d} \xleftarrow{a} \xrightarrow{w}$

$t=0$	$S=l$ $R=-0.04$ $U=0$	$S=x$ $R=-0.04$ $U=0$	$S=r$ $R=+1$ $U=+1$
	$R=-0.04$ $U=0$	$S=d$ $R=-0.04$ $U=0$	$R=-1$ $U=-1$

$$U_1(x) = -0.04 + \frac{1}{2} \max_a \left\{ \begin{array}{l} \uparrow : (0.8 \times 0) + (0.1 \times 1) + (0.1 \times 0) \\ \rightarrow : (0.8 \times 1) + (0.1 \times 0) + (0.1 \times 0) \\ \downarrow : (0.8 \times 0) + (0.1 \times 1) + (0.1 \times 0) \\ \leftarrow : (0.8 \times 0) + (0.1 \times 0) + (0.1 \times 0) \end{array} \right\}^* = 0.36$$

we note the maximal value is in the direction where the reward is greater i.e. $s=r$

$$U_1(d) = -0.04 + \frac{1}{2} \max_a \left\{ \begin{array}{l} \uparrow : (0.8 \times 0) + (0.1 \times -1) + (0.1 \times 0) \\ \rightarrow : (0.8 \times -1) + (0.1 \times 0) + (0.1 \times 0) \\ \downarrow : (0.8 \times 0) + (0.1 \times -1) + (0.1 \times 0) \\ \leftarrow : (0.8 \times 0) + (0.1 \times 0) + (0.1 \times 0) \end{array} \right\}^* = -0.04$$

with the current second abilities, the best policy at this point is π^* in future iterations

$$U_1(l) = -0.04 + \frac{1}{2} \max_a \left\{ \begin{array}{l} \uparrow : (0.8 \times 0) + (0.1 \times 0) + (0.1 \times 0) \\ \rightarrow : (0.8 \times 0) + (0.1 \times 0) + (0.1 \times 0) \\ \downarrow : (0.8 \times 0) + (0.1 \times 0) + (0.1 \times 0) \\ \leftarrow : (0.8 \times 0) + (0.1 \times 0) + (0.1 \times 0) \end{array} \right\}^* = -0.04$$

$t=1$	$S=l$ $R=-0.04$ $U=-0.04$	$S=x$ $R=-0.04$ $U=0.36$	$S=r$ $R=+1$ $U=+1$
	$R=-0.04$ $U=-0.04$	$S=d$ $R=-0.04$ $U=-0.04$	$R=-1$ $U=-1$

$$U_2(x) = -0.04 + \frac{1}{2} \max_a \left\{ \begin{array}{l} \uparrow : (0.8 \times 0.36) + (0.1 \times 1) + (0.1 \times -0.04) \\ \rightarrow : (0.8 \times 1) + (0.1 \times -0.04) + (0.1 \times 0.36) \\ \downarrow : (0.8 \times -0.04) + (0.1 \times -0.04) + (0.1 \times 1) \\ \leftarrow : (0.8 \times -0.04) + (0.1 \times -0.04) + (0.1 \times 0.36) \end{array} \right\}^* = 0.376$$

once again, the maximal value in the direction where the reward is greatest

Notes:

- Utility values propagate out from their neighbours
- We don't need to know the exact values of $\sum_{s'} T(s, a, s') U_t^{\pi^*}(s')$

for each action only how they compare to each other for the neighbouring states.

Actually what we really want is the optimal policy, not the utility function. To find this, we don't need to know the exact values for the utility function, only the relative values. This leads us to our next algorithm.

* Policy Iteration

The solution of Bellman's equation can be found using the following algorithm

- Start with $\pi_0 \leftarrow \text{guess}$

- Evaluate: given π_t calculate $U_t = U^{\pi_t}$

$$U_t(s) = R(s) + \gamma \left(\sum_{s'} T(s, \pi_t(s), s') U_t(s') \right)$$

$\hookrightarrow n$ linear equations in n unknowns ($n = \text{number of states } s$)

- Improve:

$$\pi_{t+1}(s) = \underset{a}{\operatorname{argmax}} \left(\sum_{s'} T(s, a, s') U_t(s') \right)$$

Notes:

- No max so our new equation is LINEAR
- More expensive to evaluate (matrix inversion - $O(n^3)$)
- Converges faster
- Makes bigger jumps because we are working in policy space rather than utility space.
- Guaranteed to converge - proof is like k-Means

* What does all this have to do with reinforcement learning

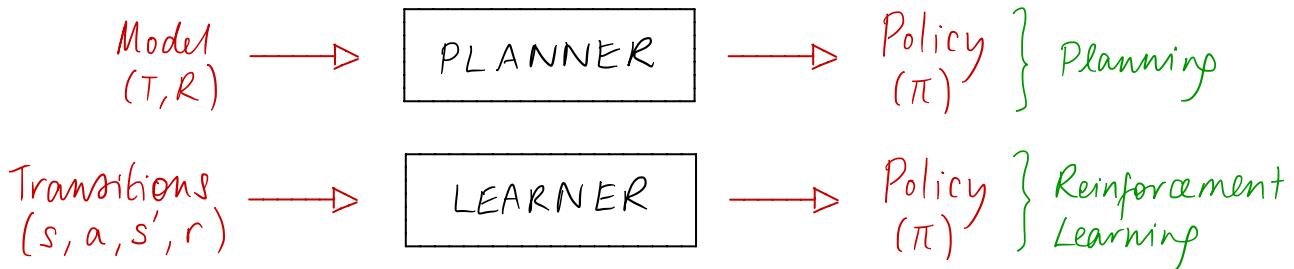
This should make it easier to think about / understand how reinforcement learning works.

In reinforcement learning you don't necessarily know what the rewards or transitions are or even the states or actions for that matter.

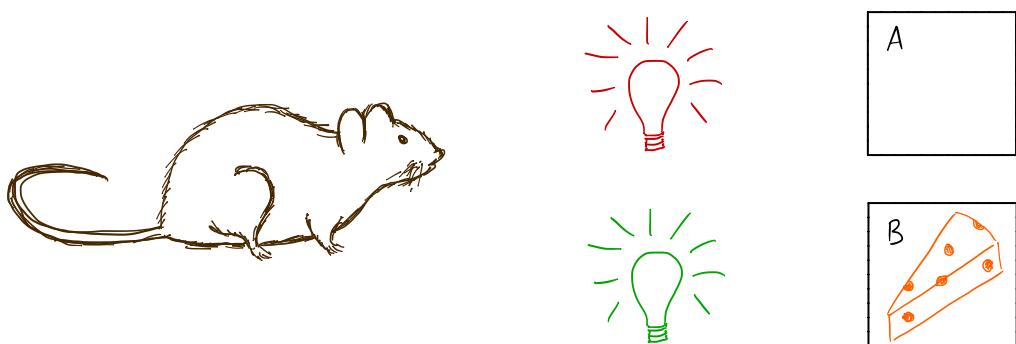
Reinforcement Learning

...In Markov Decision Processes

* Reinforcement Learning API

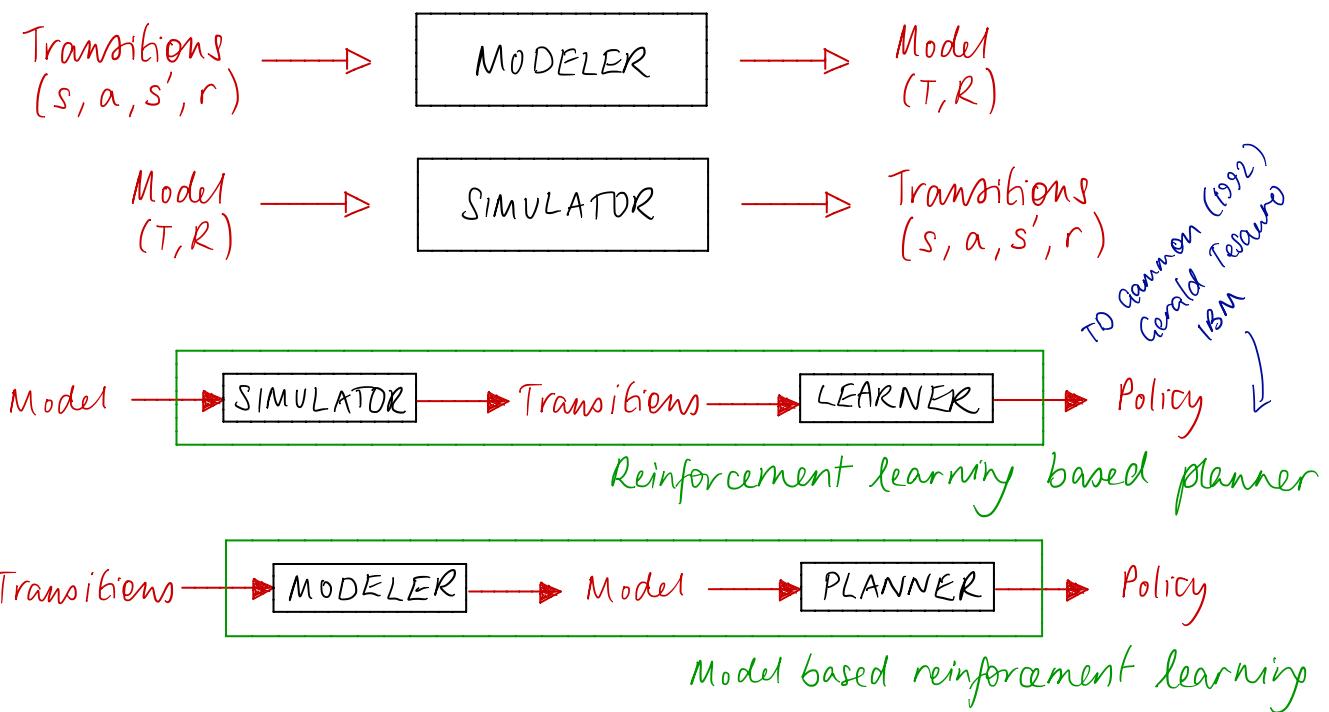


Brief aside: Reinforcement Learning 1930's

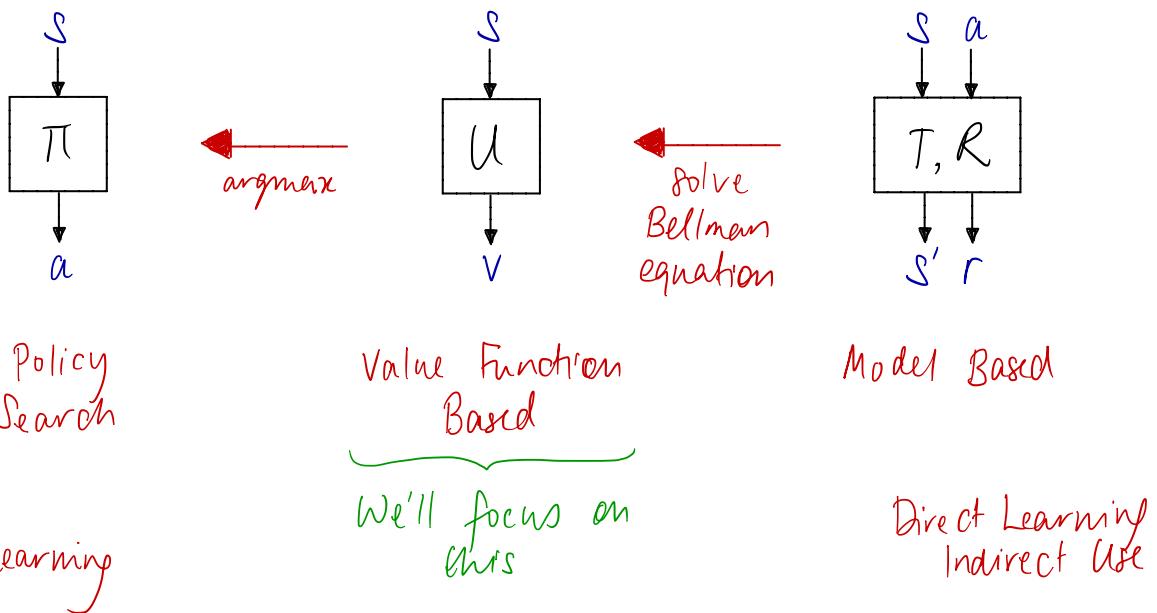


Animal sees stimulus (red/green light) s
 takes an action (open a door) a
 receives a reward (cheese?) r

Strengthens response to the stimulus
 Reward maximisation



* Three Approaches to Reinforcement Learning



* A new kind of value function

Before we had

Value function :

$$U(s) = R(s) + \gamma \max_a \left(\sum_{s'} T(s, a, s') U(s') \right)$$

Policy:

$$\pi(s) = \text{argmax}_a \left(\sum_{s'} T(s, a, s') U(s') \right)$$

Our new value function - Q function :

$$Q(s, a) = R(s) + \gamma \sum_{s'} T(s, a, s') \max_a Q(s', a')$$

Value for arriving in s , leaving via a , proceeding optimally thereafter

Note,

$$U(s) = \max_a Q(s, a)$$

and

$$\pi(s) = \text{argmax}_a Q(s, a)$$

In Q -learning we are evaluating the Bellman equations from data

* Estimating Q from transitions

$$Q(s, a) = R(s) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q(s', a')$$

The problem here is we don't know R or T , we only have transitions:

$$\langle s, a, r, s' \rangle: \quad \hat{Q}(s, a) \xleftarrow{\alpha_t} r + \gamma \max_{a'} \hat{Q}(s', a')$$

Notation:
 $v \leftarrow x$
 $v \leftarrow (1-\alpha)x + \alpha v$

Annotations:
 learning rate
 again we estimate not to denote
 utility of next state
 utility of state

Aside: $v_t \xleftarrow{\alpha_t} x_t \quad x \sim X, \quad \sum_{t=1}^{\infty} \alpha_t = \infty, \quad \sum_{t=1}^{\infty} \alpha_t^2 < \infty$

v_t converges to $\mathbb{E}(x)$

e.g. $x_t = 1/t$
 $\sum_{t=1}^{\infty} 1/t = \ln t, \quad \sum_{t=1}^{\infty} 1/t^2 = \pi^2/6$

Back to estimating transitions...

$$\langle s, a, r, s' \rangle: \quad \hat{Q}(s, a) \xleftarrow{\alpha_t} r + \gamma \max_{a'} \hat{Q}(s', a')$$

↓

doesn't really work

dropping over time

$$\begin{aligned}
 &= \mathbb{E} \left[r + \gamma \max_{a'} \hat{Q}(s', a') \right] \\
 &= R(s) + \gamma \mathbb{E} \left[\max_{a'} \hat{Q}(s', a') \right] \\
 &= R(s) + \gamma \sum_{s'} T(s, a, s') Q(s', a')
 \end{aligned}$$

* Q-Learning Convergence

\hat{Q} starts anywhere $\hat{Q}(s, a) \xleftarrow{\alpha_t} r + \gamma \max_{a'} \hat{Q}(s', a')$

then $\hat{Q}(s, a) \rightarrow Q(s, a)$ soln of the Bellman eqn

if // s, a visited infinitely often, $\sum_{t=1}^{\infty} \alpha_t = \infty, \quad \sum_{t=1}^{\infty} \alpha_t^2 < \infty$

and $s' \sim T(s, a, s'), \quad r \sim R(s)$

* Choosing Actions

Q-learning is a family of algorithms which typically vary in the following ways

1. how actions are chosen
2. how \hat{Q} is initialised
3. how α_t decays

How could we choose actions?

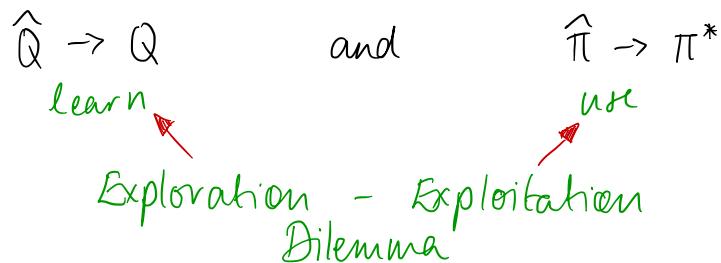
- always choose a_0 → won't explore / learn
- choose randomly → won't use what's learned
- use \hat{Q} (always choose optimal action) → greedy learner, local minima
- random restarting → too slow, we already need to visit s, a infinitely often

We want to combine learning and some randomness in order to optimise over the whole state space. To do this we use a "simulated annealing" like approach where we sometimes take a random action,

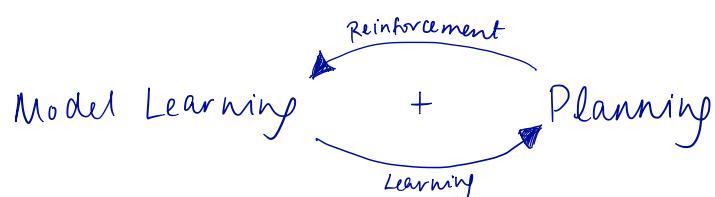
$$\hat{\pi}(s) = \begin{cases} \underset{a}{\operatorname{argmax}} \hat{Q}(s, a) & \text{with probability } 1 - \varepsilon \\ \text{random action} & \text{with probability } \varepsilon \end{cases}$$

* ε Greedy Exploration

In Greedy Limit and Infinite Exploration (GLIE) ε decays,



Fundamental trade-off in Reinforcement Learning



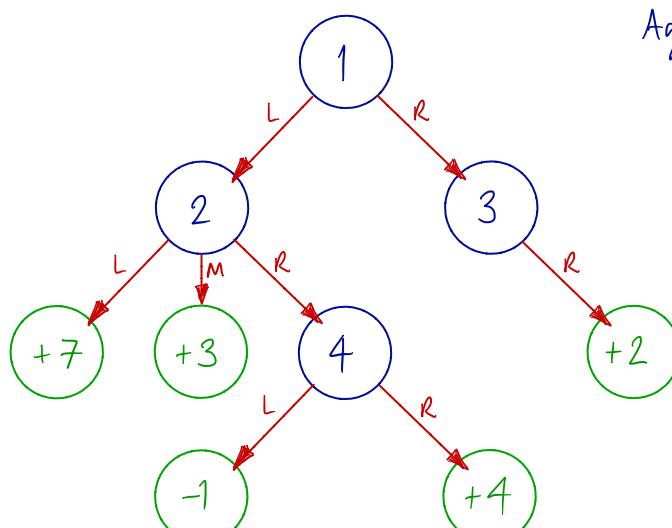
Game Theory

* What is game theory?

- Mathematics of conflict
- Single agent \rightarrow multiple agents
- Economics, politics, sociology, biology, ...
- Increasingly a part of machine learning and artificial intelligence

* A simple game

Two player (agents A and B), zero sum, finite, deterministic game of perfect information



Agent Choice
A

Four states
in blue

Leaves show
the score for
A in green.
The score for
B is minus
the score
for A.

Matrix form of the game

Tells us everything we need
to know about the game.

What's the optimal strategy
for each agent?

		B	②	L	M	R
		A	③	R	R	R
A	①	④				
	L	L	7	3	-1	
	L	R	7	3	4	
	R	L	2	2	2	
	R	R	2	2	2	

For agent A, this is given by the row with the best worst case scenario for A

For agent B, this is given by the column with the best worst case scenario for B or equivalently the worst best case scenario for A.

It turns out that for these types of games, the score or reward from the optimal strategy for both players is the same.

Fundamental Result: Von Neumann

In a

Two player, zero sum, finite, ~~deterministic~~ game of perfect information,

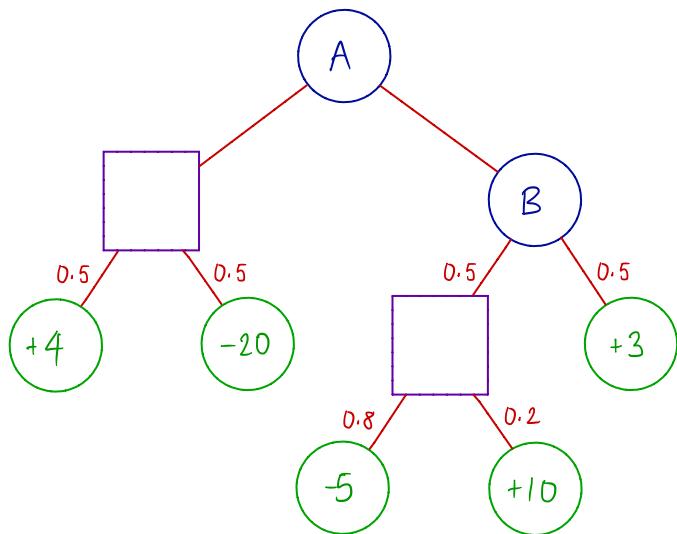
$$\text{Minimax} = \text{Maximin}$$

and

There always exists an optimal pure strategy for each player.

*don't need
the game
to be
deterministic*

* A simple non-deterministic game



Our fundamental theorem still holds (Von Neumann)

Again, all we need is the matrix form of the game

A	B	L	R
	L	-8	-8
	R	-2	3

* A simple non-deterministic game of hidden information

Two player, zero sum, finite, non-deterministic game of hidden information.

Mini Poker:

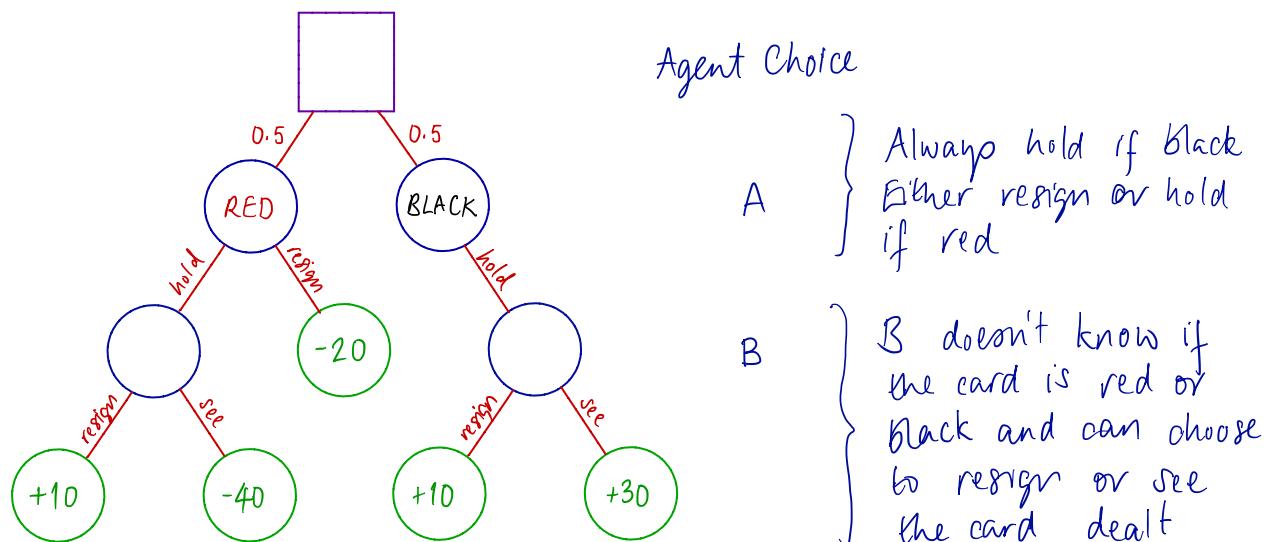
- A is dealt a card, red or black 50%
- A may resign if red: -20\$ for A
else A holds

B resigns: +10 \$

B sees:

if red: -40 \$
if black: +30 \$

Mini Poker Game Tree



Matrix form of the game:

		B	
		resign	see
A		resign	
resign		-5	+5
hold		+10	-5

refers to the strategy in the case A is dealt a red card only

Von Neumann does not hold when there is hidden information.
i.e. $\text{Minimax} \neq \text{Maximin}$

← Pure strategy
we need a mixed strategy

* Mixed Strategy

Here we have a distribution over the strategies.

Example: suppose that A takes the following mixed strategy

If the card is red then p of the time A holds and $1-p$ of the time A resigns.

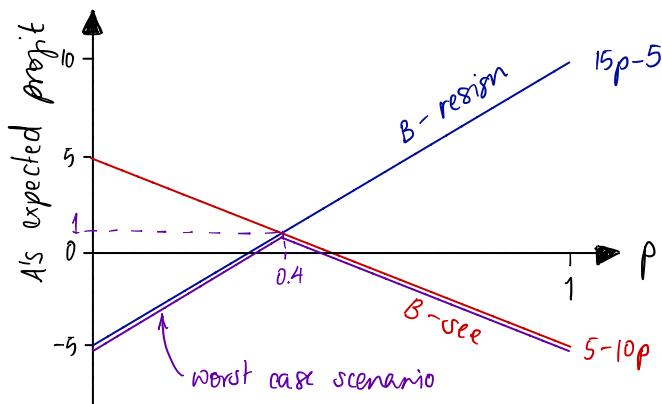
If B always resigns, A's expected profit is

$$10p - 5(1-p) = 15p - 5$$

If B always holds, A's expected profit is

$$-5p + 5(1-p) = 5 - 10p$$

Graphically ...

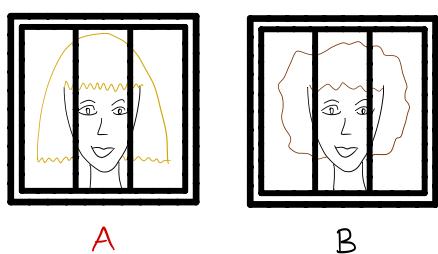


The best, worst case scenario is in this case where the lines intersect, at $p=0.4$. A's expected profit with this strategy is 1.

Note that if B employs a mixed strategy the result is the same.

- Two player non-zero sum non-deterministic game of hidden information

Prisoner Dilemma:



Two prisoners, A and B, are arrested for their involvement in a crime. Each prisoner has the option of either snitching or not.

(A, B)	don't snitch	snitch
don't snitch	$(-1, -1)$	$(-9, 0)$
snitch	$(0, -9)$	$(-6, -6)$

The matrix shows the jail time they get in each case.

Note, the matrix is symmetric.

If we think collectively and combine the jail times it's clear that if neither prisoner snitches the combined jail time is minimised.

If each player simply aims to minimise their own jail time, it's always better to snitch. The problem is this 'selfish' strategy, when employed by both prisoners leads to the worst possible combined jail time. The snitching strategy dominates.

* The Nash Equilibrium (as in John Nash of A Beautiful Mind)

Suppose we have n players which can each choose a strategy from their set of all possible strategies,

$$S_1, S_2, \dots, S_n$$

The collective strategy

$$s^*, s_2^*, \dots, s_n^* \text{ where } s_i^* \in S_1, s_2^* \in S_2, \dots, s_n^* \in S_n$$

is a Nash Equilibrium if and only if

$$\forall i \quad s_i^* = \underset{s_i}{\operatorname{argmax}} \ U(s_1^*, \dots, s_i, \dots, s_n^*)$$

That is to say, assuming all other players stick with the strategy they chose, each player has chosen the best possible strategy and has no motivation to switch.

Example:

		(A, B)	don't snitch	snitch
		don't snitch	(-1, -1)	(-9, 0)
		snitch	(0, -9)	(-6, -6)
				Nash equilibrium

* Some results

In the n -player pure strategy game, if elimination of strictly dominated strategies eliminates all but one combination, that combination is the unique Nash Equilibrium

Any Nash Equilibrium will survive elimination of strictly dominated strategies

If n is finite & $\forall i \ S_i$ is finite \exists at least one (possibly mixed) Nash Equilibrium

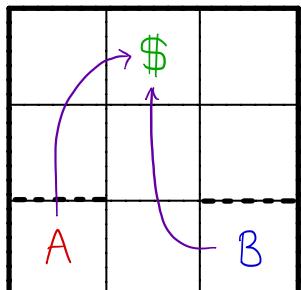
* Final notes

Utility isn't just the reward - Mechanism Design

* Stochastic games and multiagent reinforcement learning

MDP : RL :: STOCHASTIC GAME : MULT-AGENT RL

Example :



- Actions: N, E, S, W, X
- First to reach goal gets \$100
- Both arrive both win
- Semi-wall (50% go through)
- Coin flip if collide

Two Nash Equilibria in this case (one is shown)

Similarly to MDPs we have ... (Shapley)

S: states

A_i : actions for player i

T: transitions

R_i : rewards for player i

γ : discount

S

a, b $a \in A_1, b \in A_2$

T(s, (a, b), s')

$R_1(s, (a, b)), R_2(s, (a, b))$

γ

* Zero sum stochastic games

Bellman Equation for stochastic games:

$$Q^*(s, (a, b)) = R_i(s, (a, b)) + \gamma \sum_{s'} T(s, (a, b), s') \text{ minimax } Q^*(s', (a', b'))$$

$\langle s, (a, b), (r_1, r_2), s' \rangle$: minimax-Q

$$Q_i(s, (a, b)) \leftarrow r_i + \gamma \underset{a', b'}{\text{minimax}} Q_i(s', (a', b'))$$

- Value iteration works
- Minimax Q converges
- Unique solution to Q*
- Policies for the two players can be computed independently
- Update efficient (polynomial time)
- Q function sufficient to specify policy

* General sum stochastic games

$$Q^*(s, (a, b)) = R_i(s, (a, b)) + \gamma \sum_{s'} T(s, (a, b), s') \underset{\text{minimax}}{\underset{\text{Nash}}{\text{minimax}}} Q^*(s', (a', b'))$$

$$\langle s, (a, b), (r_1, r_2), s' \rangle : \underset{\text{minimax}}{\underset{\text{Nash}}{\text{minimax}}} - Q$$

$$Q_i(s, (a, b)) \leftarrow r_i + \gamma \underset{a', b'}{\underset{\text{minimax}}{\underset{\text{Nash}}{\text{minimax}}}} Q_i(s', (a', b'))$$

- Value iteration doesn't work
- ~~minimax~~ Nash Q doesn't converge
- Unique solution to Q^* does not exist
- Policies for the two players can't be computed independently
- Update inefficient (polynomial time) $P = \text{ppad}$
- Q function insufficient to specify policy

Oh dear !!!

* Some ideas for general sum stochastic games

- repeated stochastic games (folk theorem)
- cheap talk \rightarrow correlated equilibria
- cognitive hierarchy \rightarrow best responses
- side payments \rightarrow coco values

P4: Train a Smartcab to Drive

Update Rule for \hat{Q}

$$\hat{Q}(s, a) = (1 - \alpha_t) \hat{Q}(s, a) + \alpha_t (r + \gamma \max_{a'} \hat{Q}(s', a'))$$

Using $\alpha_t = \frac{1}{t}$

$$\Rightarrow \hat{Q}(s, a) = \left(1 - \frac{1}{t}\right) \hat{Q}(s, a) + \frac{1}{t} (r + \gamma \max_{a'} \hat{Q}(s', a'))$$