

# Supervised Learning

## Decision Trees

### \* Terminology :

Instances - input

Concept - function : input  $\rightarrow$  output

Target concept - function which classifies correctly

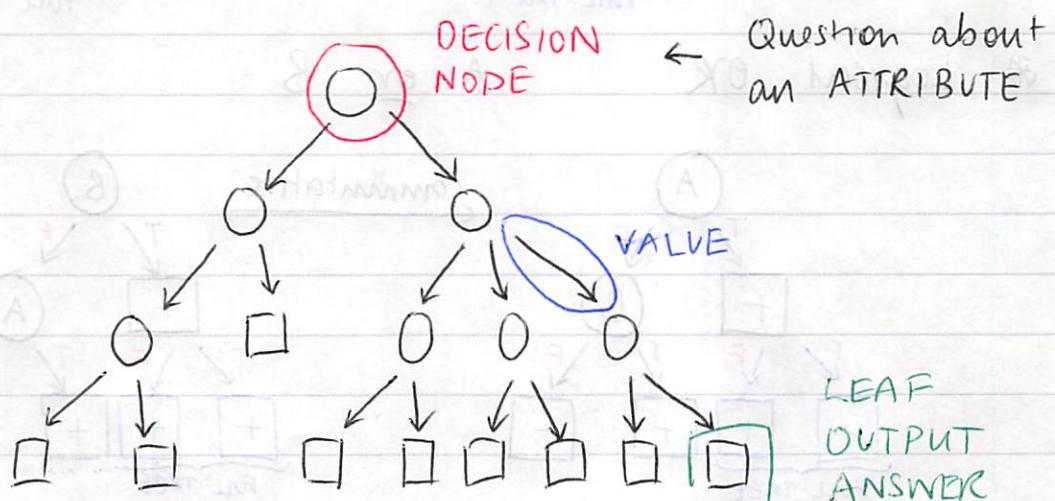
Hypothesis class - class of concepts under consideration

Sample - training set

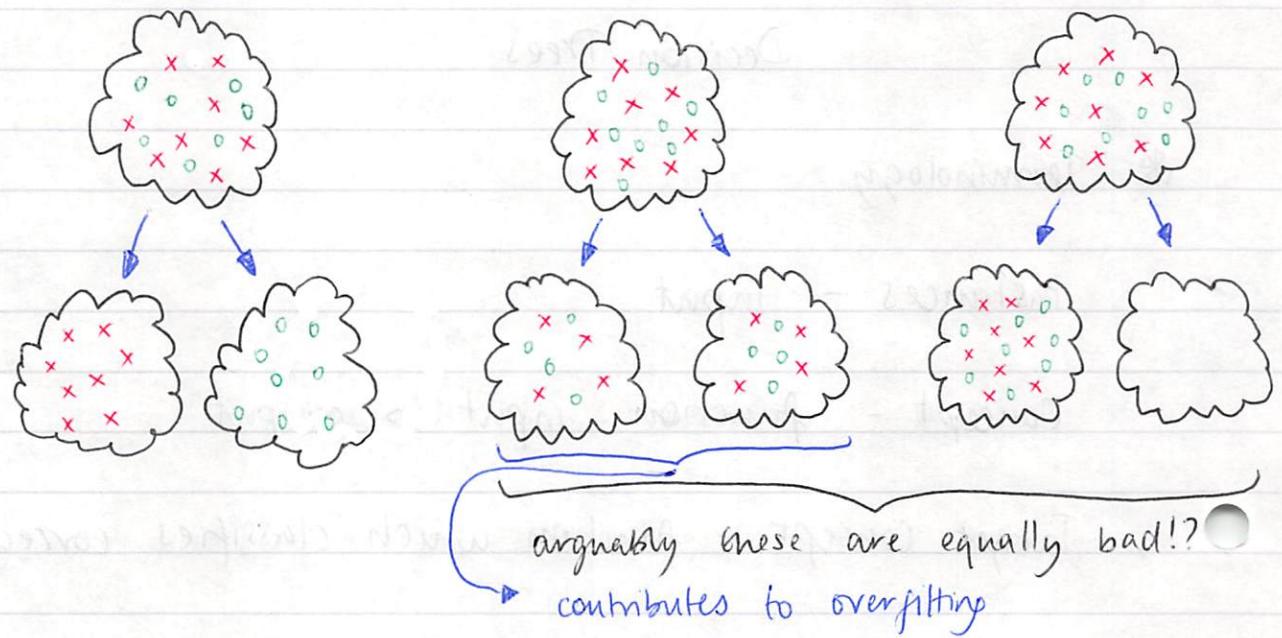
Candidate concept

Testing set - independent of training set

### \* Representation

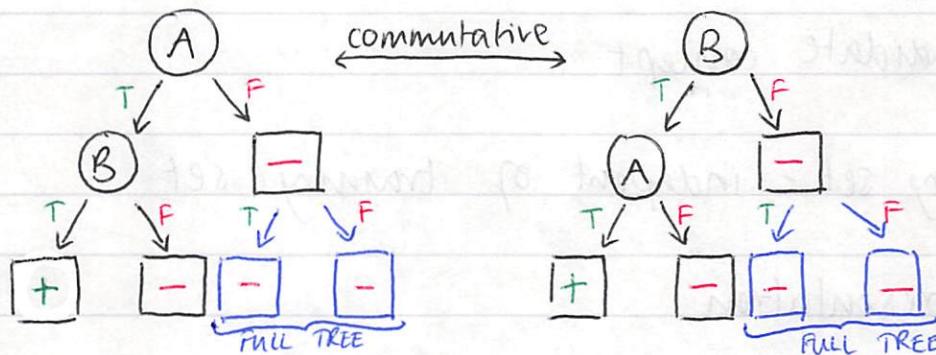


\* Best Attribute : In order of preference :

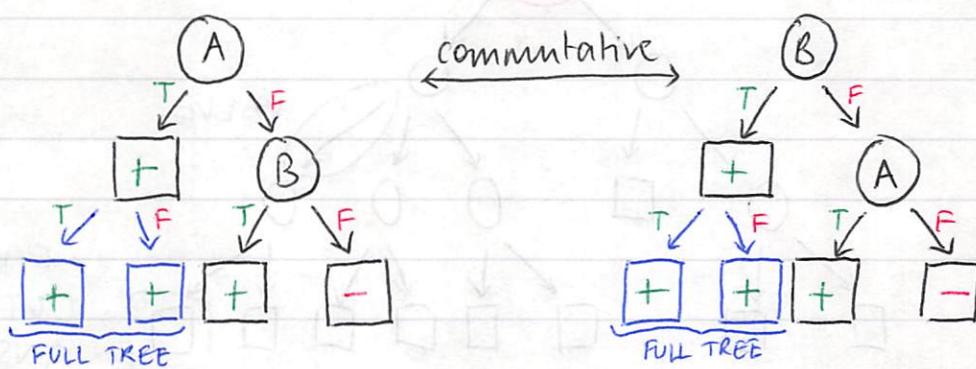


\* Decision Trees : Expressiveness

\* Logical AND      A and B

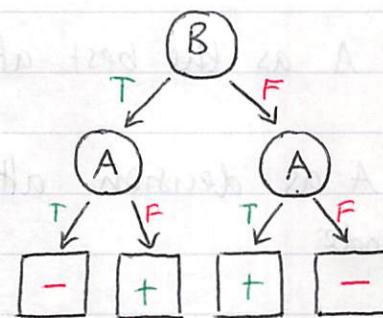
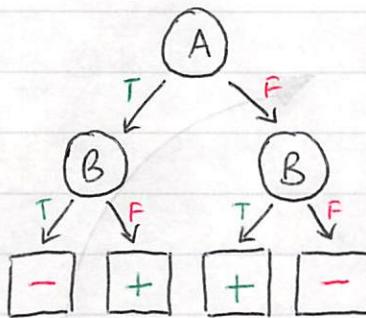


\* Logical OR      A or B

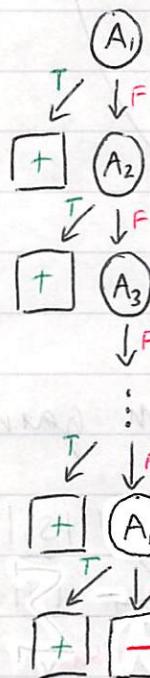


\* Logical XOR

$A \text{ xor } B$



\* n-OR : ANY

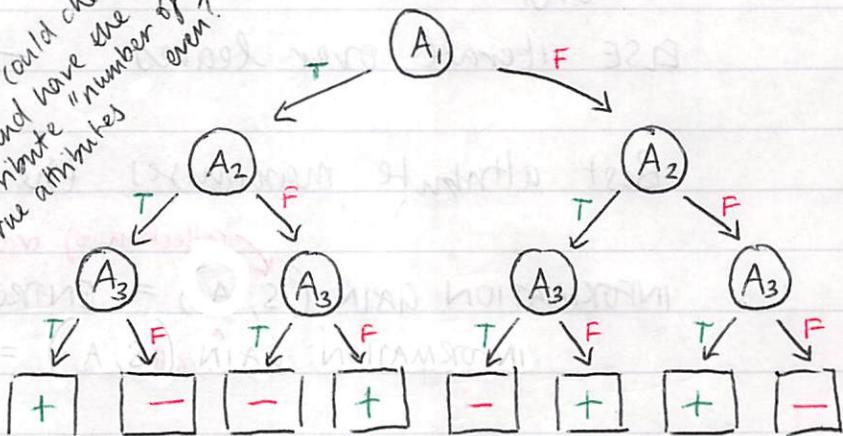


\* n-XOR : ODD PARITY

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 <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	...	A <sub>n</sub>	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

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

$2^n$  rows

→ get to one known  
terminal node found  
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

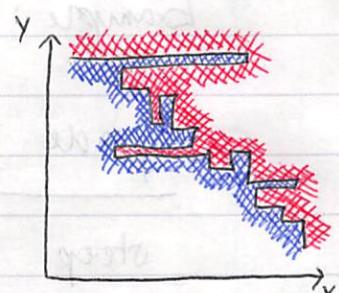
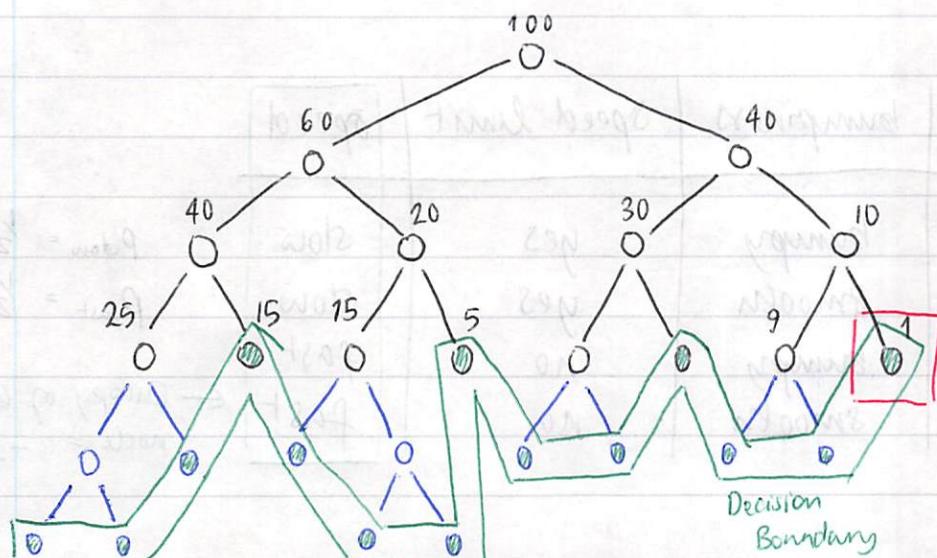
→ Picking the best attribute entropy → variance?

→ Output - average, local linear fit

## More Decision Trees

### \* Decision Tree Parameters

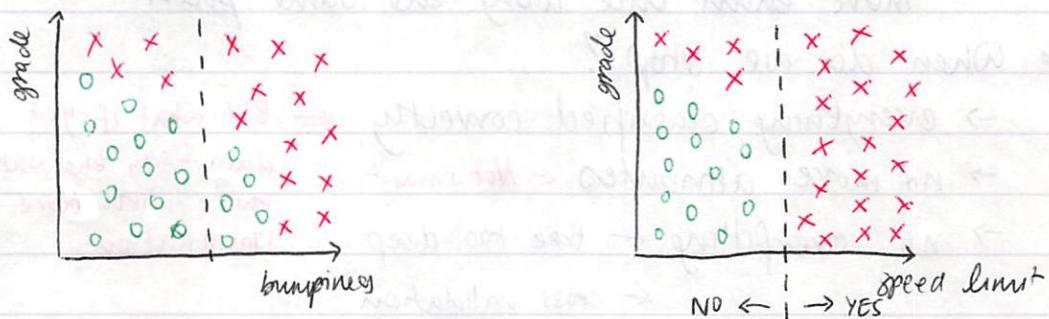
min\_samples\_split = 2 by default



## \* 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

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

$\uparrow$  fraction of examples  
in class  $i$

$\uparrow$  sum over all  
classes  $i$

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

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

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

$\leftarrow$  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

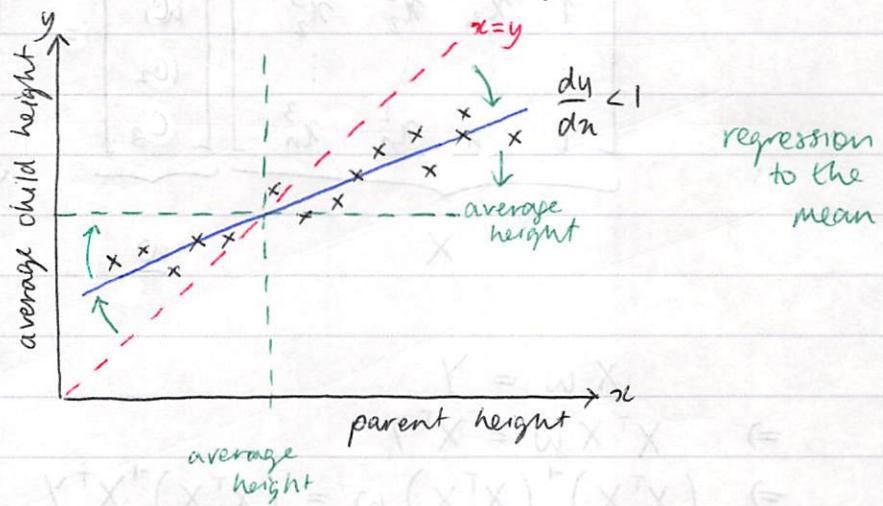
Find the sweet spot.

# 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,

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

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

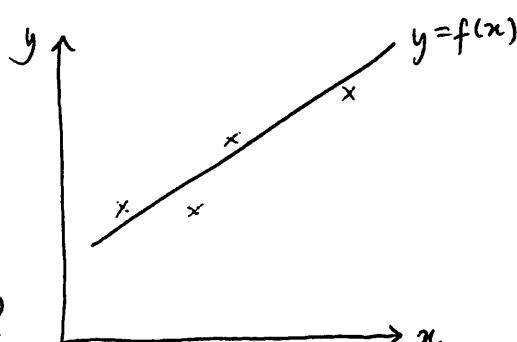
$X \quad w \quad Y$

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

## \* Errors

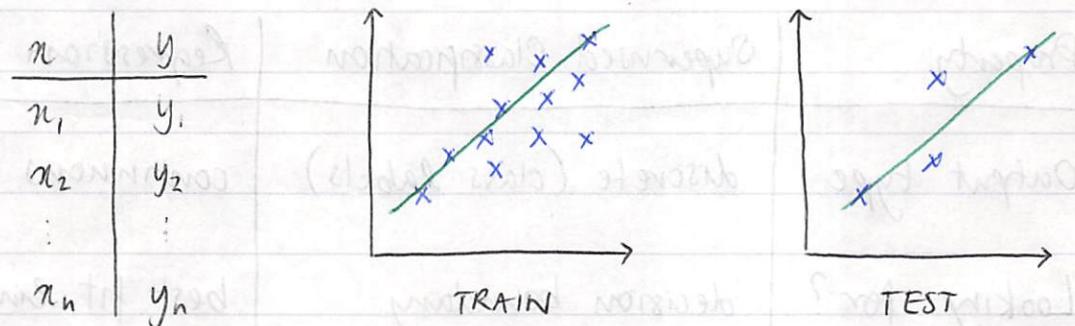
Training data has errors  
so we end up not modelling  
 $f$  but  $f + \epsilon$ .

Where do errors come from?

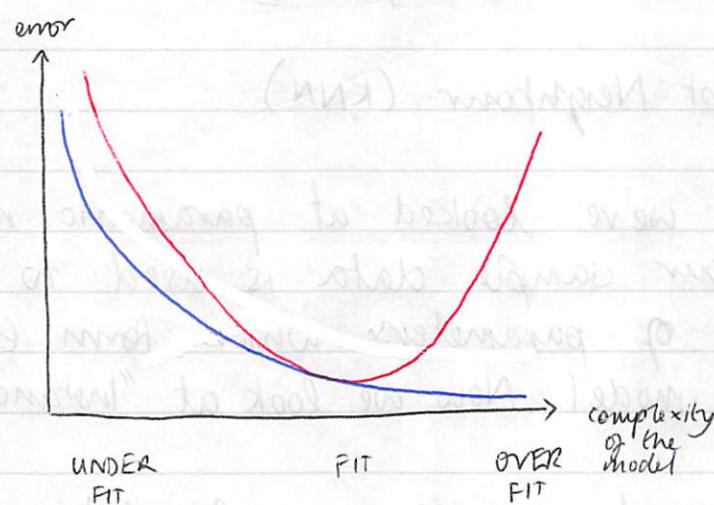


- sensor error
- maliciously - being given bad data
- transcription error
- unmodelled influences

## \* 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 $r^2$

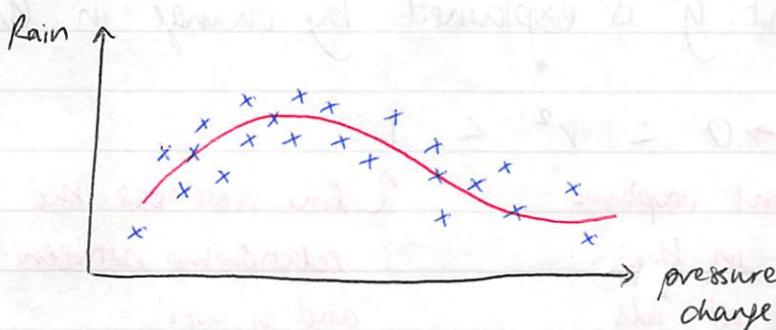
## More Regression

### \* 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.

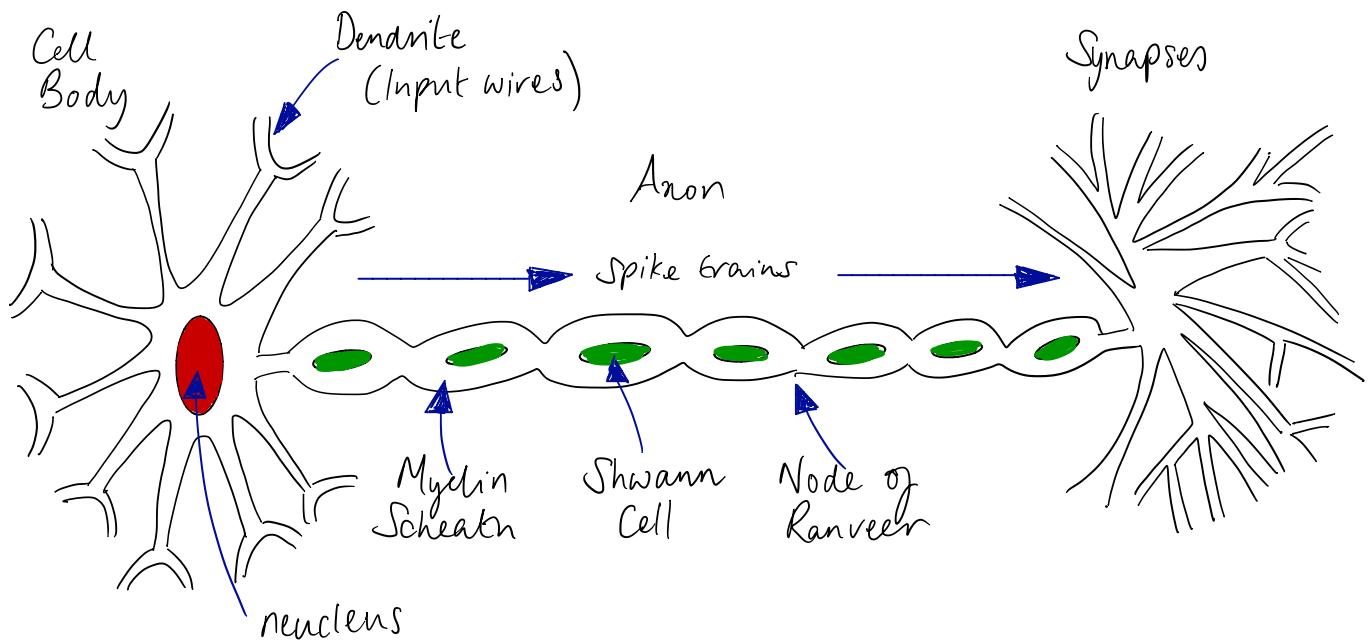
Instance based methods are good when we don't know of what <sup>the</sup> relationship between input and output might be.

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.

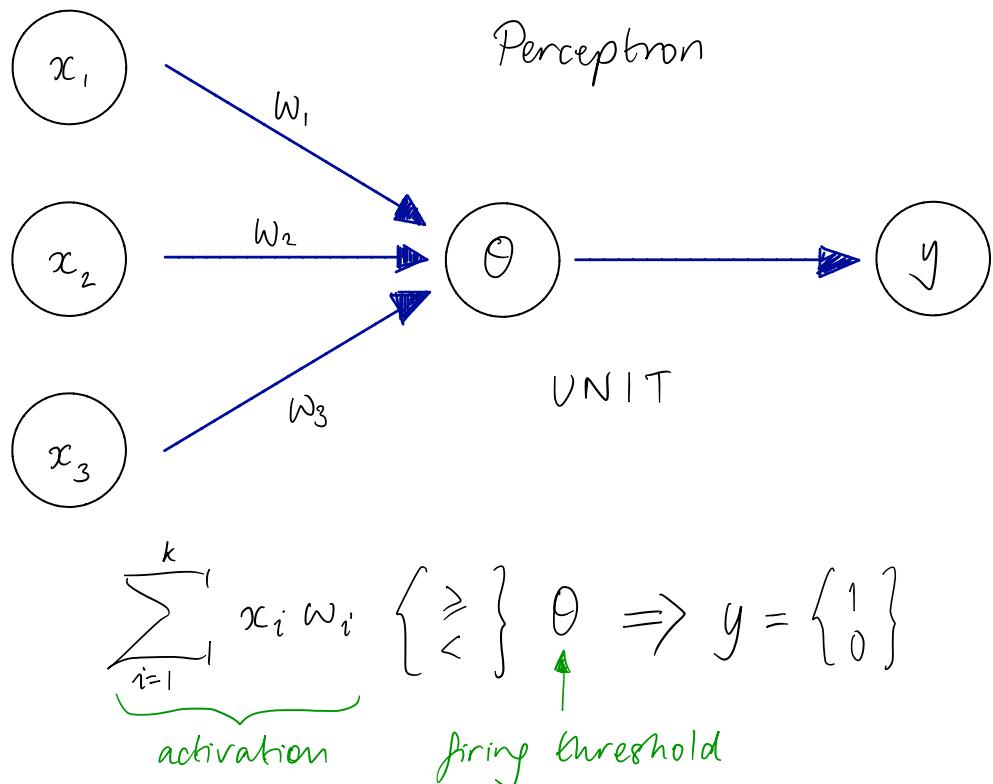
For instance based models training is fast and querying is slow. New data is easily added to update the model.

# Neural Networks

\* Neuron: Component Unit



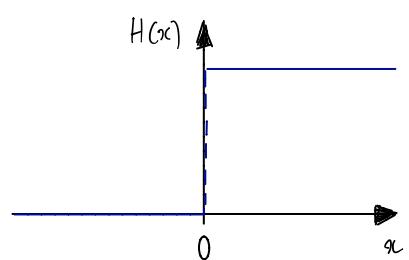
\* Artificial Neural Network



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,

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

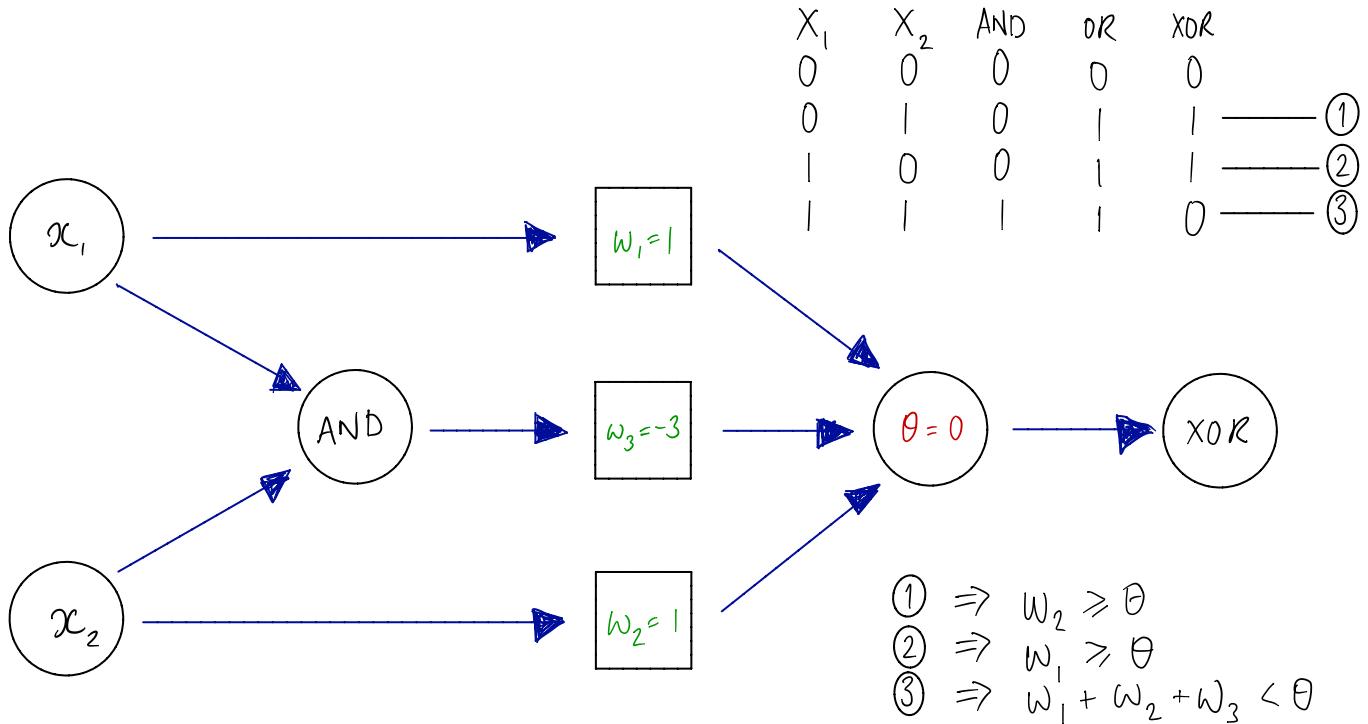
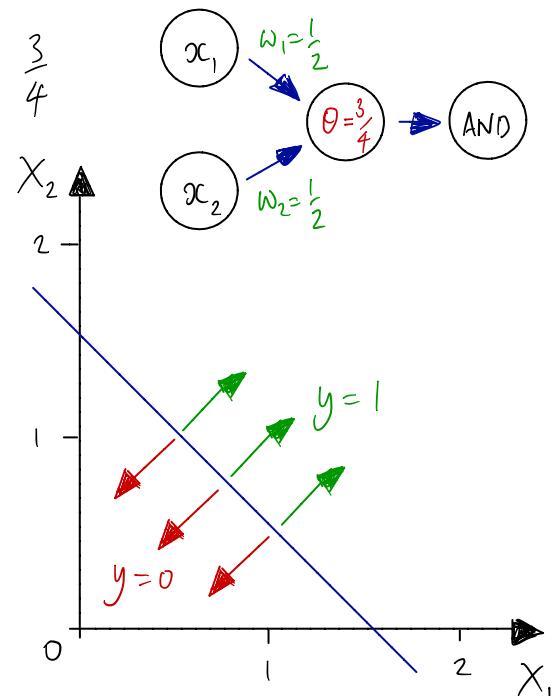
$$\text{activation} = \frac{1}{2} (x_1 + x_2) = \frac{3}{4}$$

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

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  
 $\eta$  = learning rate

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

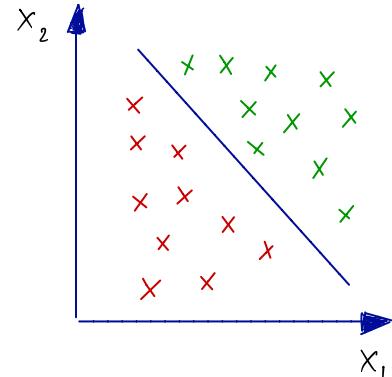
Update Rule :  $w_i = w_i + \Delta w_i$

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

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

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

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

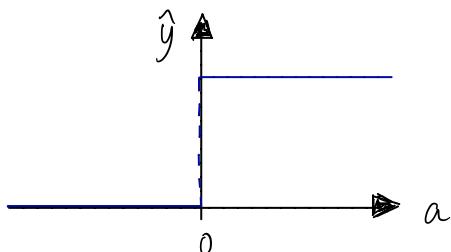
perception :  $\Delta w_i = \eta (y - \hat{y}) x_i, \quad \hat{y} = H(a)$

- 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.

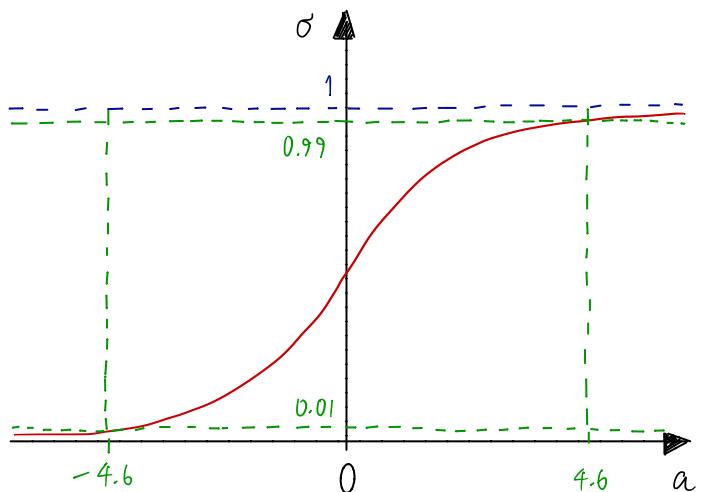
## \* Sigmoid - Differentiable threshold

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

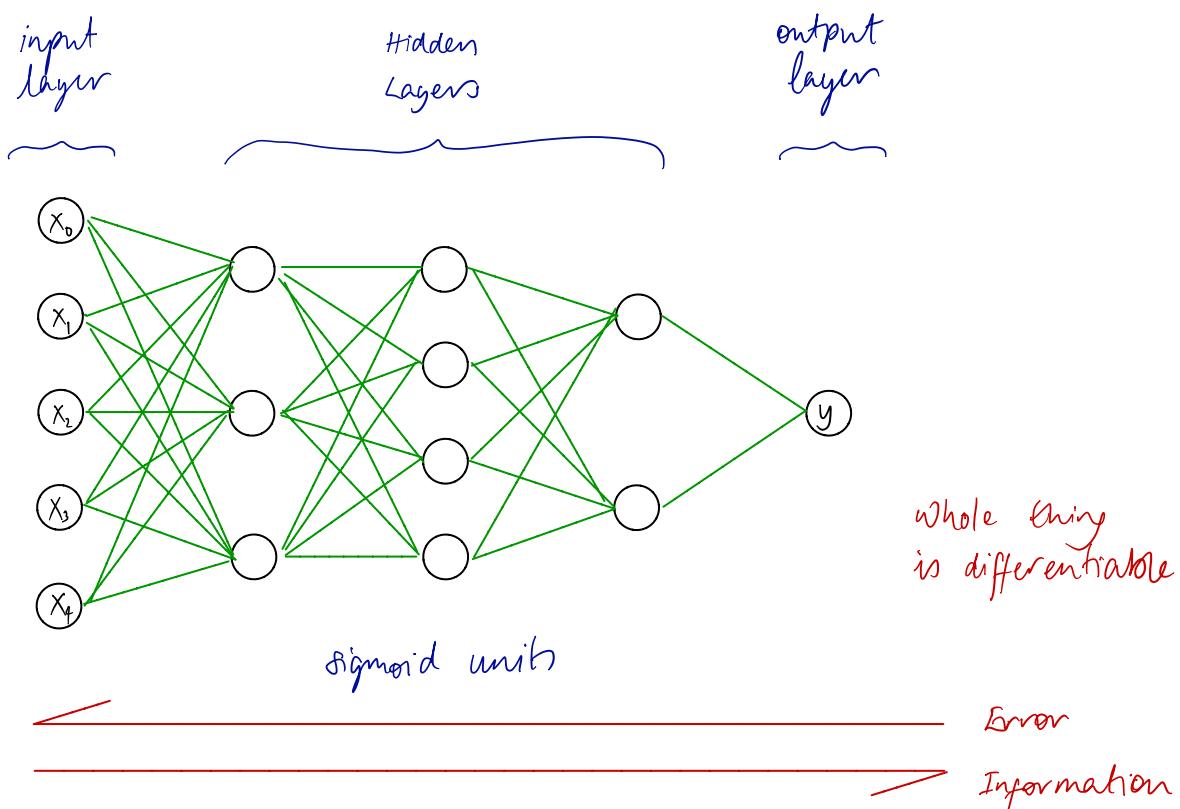
as  $a \rightarrow \infty$ ,  $\sigma(a) \rightarrow 1$

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

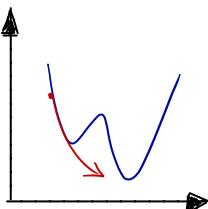
## \* Optimising Weights

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

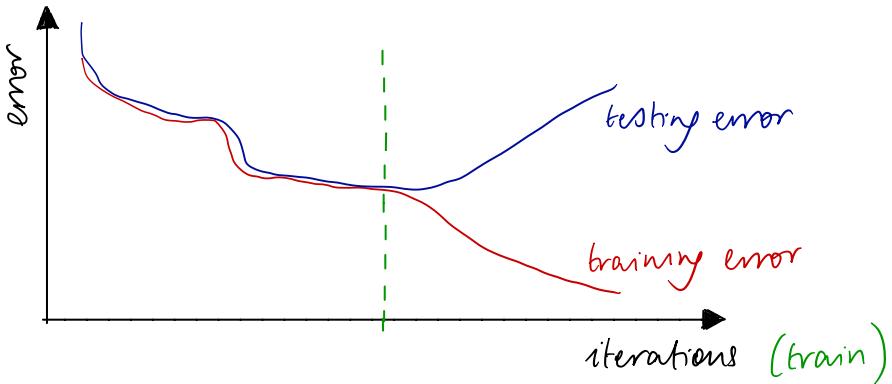
→ Boolean functions: network of threshold like units

→ 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.

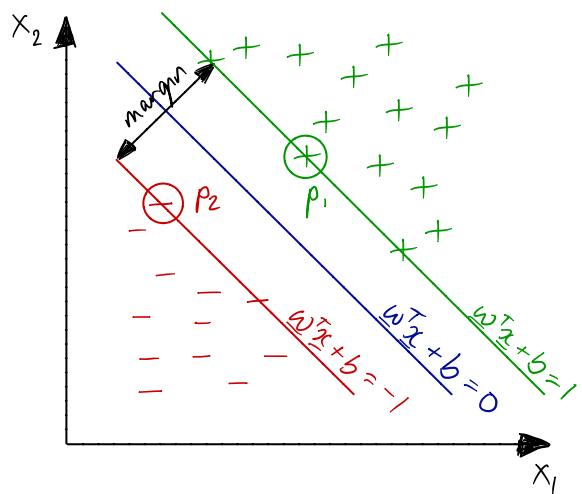
→ Shown to reduce the generalisation error

# 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:

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

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  $\alpha$ s which maximise

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

*This term is larger for data points with more similar features*

such that

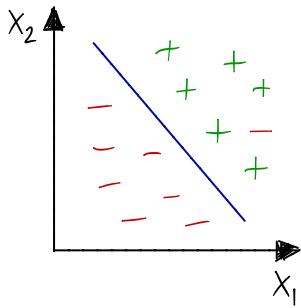
$$\alpha_i \geq 0, \quad \sum_{i=1}^m \alpha_i y_i \neq 0$$

Actually, it turns out that most of the  $\alpha$ 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.

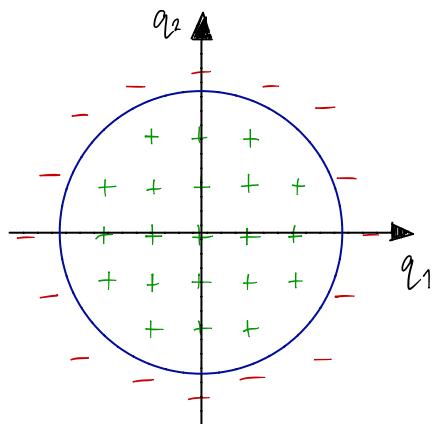
This is a bit like KNN, but solving our quadratic programming problem determines which data points are the important ones i.e. the ones with corresponding non-zero  $\alpha$ .

The data points for which the corresponding  $\alpha$ 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,

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

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

$$\begin{aligned} \Phi(\underline{x}) \cdot \Phi(\underline{y}) &= (x_1^2, x_2^2, \sqrt{2}x_1x_2) \cdot (y_1^2, y_2^2, \sqrt{2}y_1y_2) \\ &= x_1^2 y_1^2 + x_2^2 y_2^2 + 2x_1 x_2 y_1 y_2 \\ &= (x_1 y_1 + x_2 y_2)^2 = (\underline{x}^\top \underline{y})^2 \end{aligned}$$

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.

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

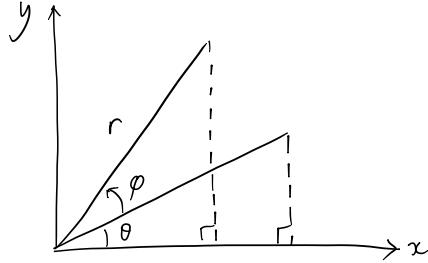
or equivalently

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

is +ve semi definite.

\* Boosting ???

\* Article: Equation of an ellipse ...



rotation by  $\varphi$

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

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

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

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

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

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

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?*

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

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

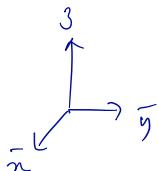
What about the blue circle? what does that get mapped to?

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

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

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

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



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

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

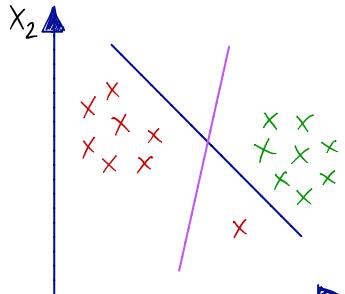
← *an ellipse*

If all these points lie on a plane.

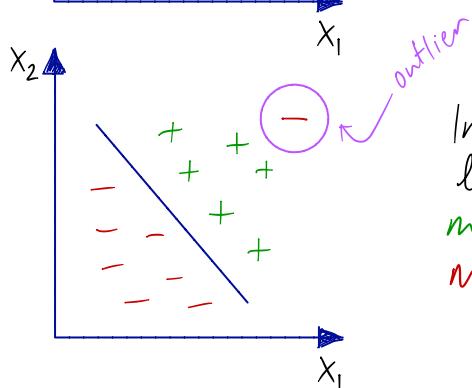
# Support Vector Machines

## \* Margin vs Outliers

Support Vector Machines prioritise classifying data correctly over maximising the margin

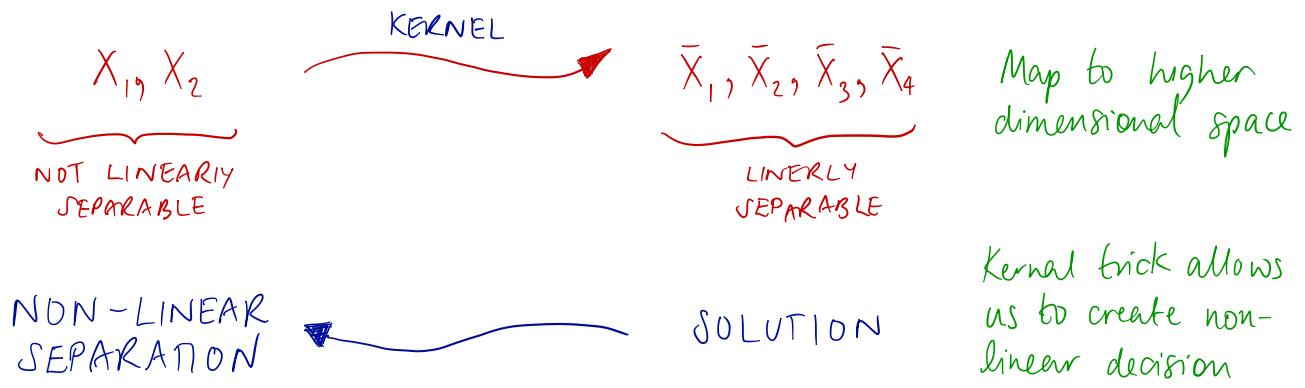


So in this example the SVM would choose the blue decision boundary rather than the purple one



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

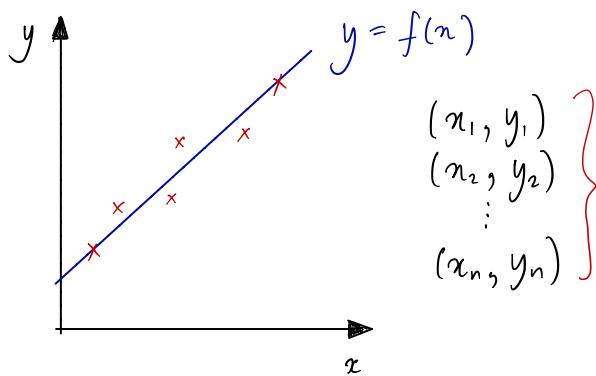
## \* Kernel Trick



Kernel trick allows us to create non-linear decision boundaries

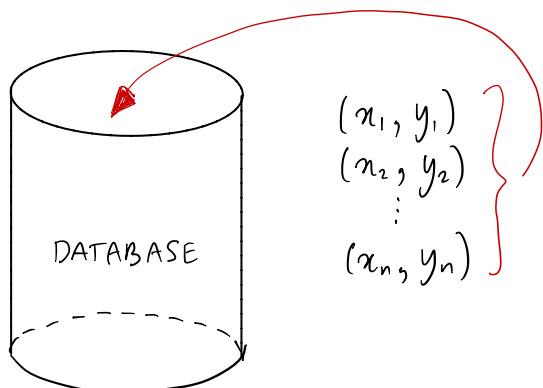
# Instance Based Learning

## \* 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

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

$\rightarrow$  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

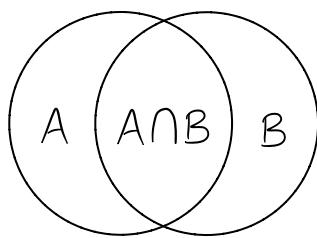
$\hookrightarrow$  Suppose we do locally weighted linear regression. We are able to build more complicated curves using just straight lines

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

$P(\text{cancer} | -ve) = P(B|A)$  is  
Specificity

Posterior probability  
given the evidence  
e.g.  $P(\text{cancer} | +ve \text{ test})$

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

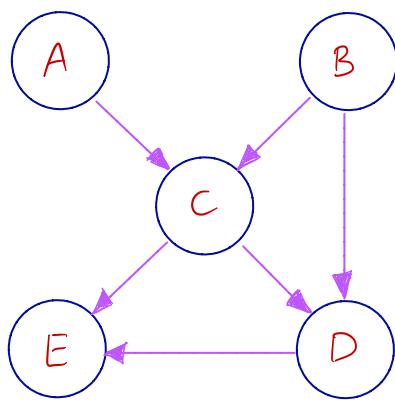
$$\forall x, y, z \quad P(x=x | Y=y \wedge Z=z) = P(x=x | Z=z)$$

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

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

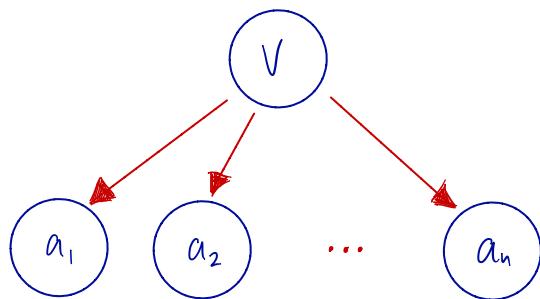
## \* Inferencing Rules

- 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

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

## \* 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 ← 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  $D_t$ :

$$D_t = \frac{1}{n},$$

+1 when  $h$  is correct else -1

$$D_{t+1} = D_t(i) \exp\left\{-\alpha_t y_i h_t(x_i)\right\} / Z_t$$

$$\alpha_t = \ln(1 - \epsilon_t) / (2\epsilon_t) > 0$$

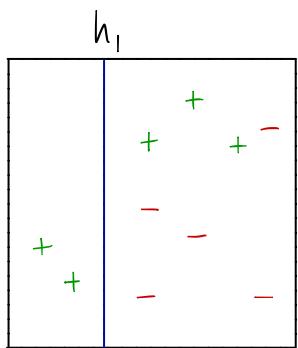
$$Z_t = \sum_i D_t(i) \exp\left\{-\alpha_t y_i h_t(x_i)\right\} \quad \text{normalisation factor}$$

- 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

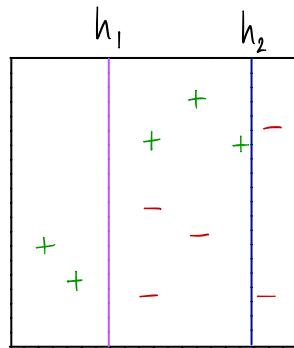
Harder examples

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



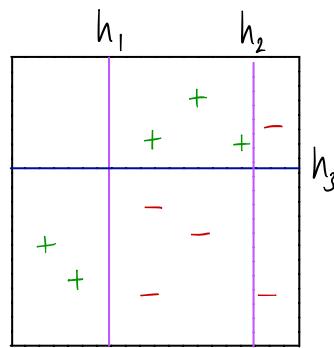
$$\epsilon_t = 0.3$$

$$\alpha_t = 0.42$$



$$\epsilon_t = 0.21$$

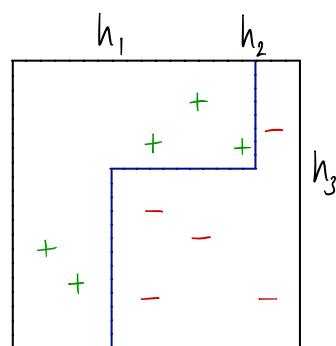
$$\alpha_t = 0.65$$



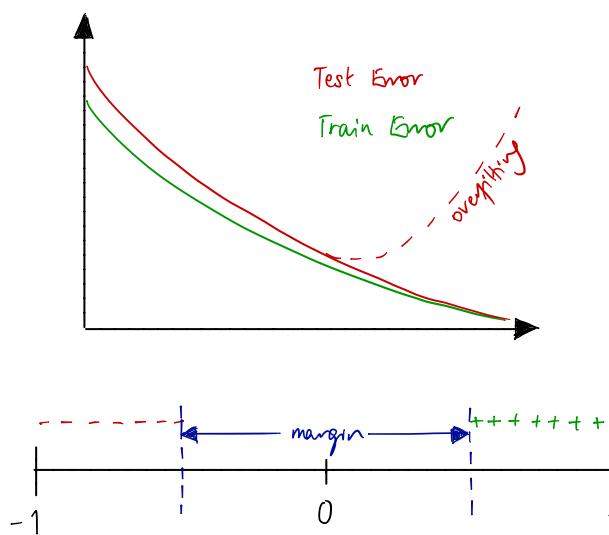
$$\epsilon_t = 0.14$$

$$\alpha_t = 0.92$$

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



\* 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