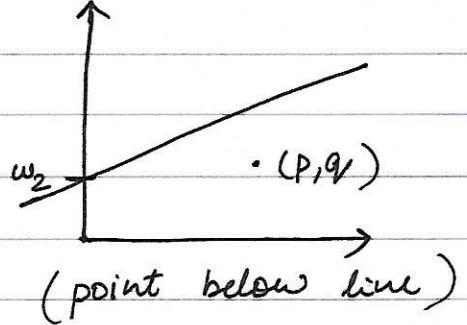
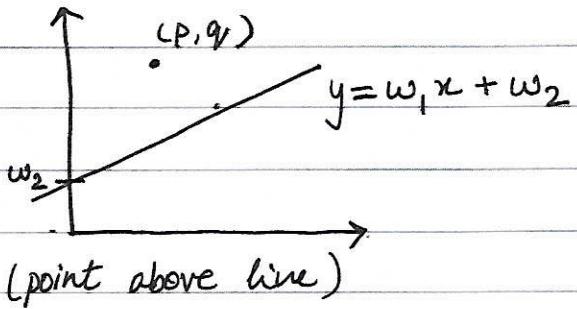


LINEAR REGRESSION

ABSOLUTE TRICK



To bring line closer to the point ; when point above line.

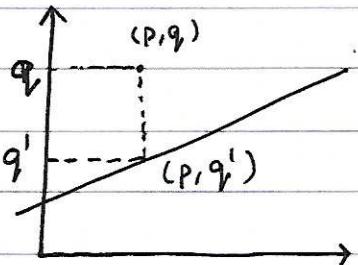
$$y = (w_1 + p\alpha)x + (w_2 + \alpha) \quad \text{where,}$$

α = learning rate

when point below line,

$$y = (w_1 - p\alpha)x + (w_2 - \alpha)$$

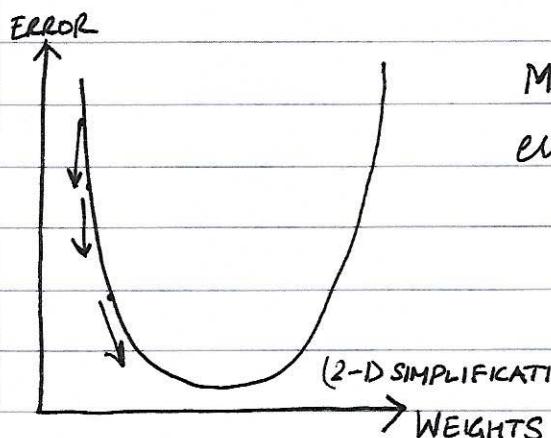
SQUARE TRICK



$$y = (w_1 + p(q-q')\alpha)x + (w_2 + (q-q')\alpha)$$

The equation above also works when point is below the line. $(q-q')$ becomes negative and both slope & y-intercept of line are decreased.

GRADIENT DESCENT

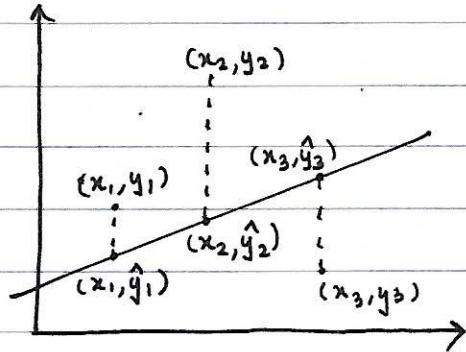


Move in the direction of negative gradient of error function , until the minima is reached.

$$w_i := w_i - \alpha \frac{\partial (\text{Error})}{\partial w_i}$$

(2-D SIMPLIFICATION) (IN REALITY THIS WOULDN'T BE 2-D)

MEAN ABSOLUTE ERROR



$$\text{ERROR} = \frac{1}{m} \sum_{i=1}^m |y_i - \hat{y}_i|, \text{ where,}$$

m: no. of training examples
 (x_i) points

MEAN SQUARED ERROR

$$\text{ERROR} = \frac{1}{2m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

$$\frac{\partial \text{ERROR}}{\partial w_1} = -(y - \hat{y})x$$

$$\frac{\partial \text{ERROR}}{\partial w_2} = -(y - \hat{y})$$

$$\begin{aligned} \frac{\partial}{\partial w_1} \cdot \frac{1}{2m} \sum (y_i - \hat{y}_i)^2 \\ &= \frac{2}{2m} \sum (y_i - \hat{y}_i) \cdot \frac{\partial}{\partial w_1} (y_i - \hat{y}_i) \\ &= \frac{1}{m} \sum (y_i - \hat{y}_i) \cdot \frac{\partial}{\partial w_1} y_i - (w_1 x_i + w_2) \\ &= \frac{1}{m} \sum (y_i - \hat{y}_i) \cdot (-x_i) \end{aligned}$$

Gradient Descent.

$$w_1 = w_1 - \alpha \frac{\partial \text{ERROR}}{\partial w_1}$$

$$= w_1 - \alpha (-y - \hat{y})x$$

$$= w_1 + \alpha (y - \hat{y})x$$

$$w_1 = w_1 + \alpha \cdot (q - q') \cdot p$$

$$w_2 = w_2 - \alpha \frac{\partial \text{ERROR}}{\partial w_2}$$

$$= w_2 - \alpha (-y - \hat{y})$$

$$= w_2 + \alpha (y - \hat{y})$$

$$w_2 = w_2 + \alpha (q - q')$$

The above is same as the square trick to fit a line through a set of points

Similarly, gradient descent with MAE as the error function is equivalent to 'absolute trick' method.

HIGHER DIMENSIONS

	Size	School Quality	...	No. of rooms	Price
House 1	900	6	...	2	\$100K
House 2	560	4	...	1	\$50K
...
House m	2000	7	...	4	\$250K
	(x_1)	(x_2)	\dots	(x_{n-1})	(\hat{y})

$$\hat{y} = \omega_1 x_1 + \omega_2 x_2 + \dots + \omega_{n-1} x_{n-1} + \omega_n$$

MULTIPLE LINEAR REGRESSION

For 'n' predictor variables (or) features :

$$y = m_1 x_1 + m_2 x_2 + \dots + m_n x_n + b$$

CLOSED FORM SOLUTION TO FIND MODEL PARAMETERS (ω_i)

$$E(\omega_1, \omega_2) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 = \frac{1}{2m} \sum_{i=1}^m (\omega_1 x_i + \omega_2 - y_i)^2$$

To minimise error function, we take derivatives and set them to 0

$$\begin{aligned} \frac{\partial E}{\partial \omega_1} &= \sum_{i=1}^m (\omega_1 x_i + \omega_2 - y_i) x_i \\ &= \omega_1 \sum x_i^2 + \omega_2 \sum x_i - \sum x_i y_i \end{aligned}$$

Setting to 0

$$\begin{aligned} \frac{\partial E}{\partial \omega_2} &= \sum (\omega_1 x_i + \omega_2 - y_i) \\ &= \omega_1 \sum x_i + \omega_2 \sum 1 - \sum y_i \end{aligned}$$

$$\omega_1 \sum x_i^2 + \omega_2 \sum x_i = \sum x_i y_i$$

$$\omega_1 \sum x_i + \omega_2 m = \sum y_i$$

2 equations, 2 unknown variables can be solved using any closed form mathematical method. (2 linear equations, 2 variables)

$$\omega_1 a_1 + \omega_2 b_1 = c_1 ; \quad \omega_1 a_2 + \omega_2 b_2 = c_2$$

CLOSED FORM SOLUTION - n DIMENSIONS / FEATURES

$$X = \begin{bmatrix} x_0^{(0)} & x_1^{(0)} & \dots & x_n^{(0)} \\ x_0^{(1)} & x_1^{(1)} & \dots & x_n^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ x_0^{(m)} & x_1^{(m)} & \dots & x_n^{(m)} \end{bmatrix}$$

(rowwise training examples)
 m : training examples
 n : dimensions/features

$$\begin{array}{l} y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \\ (\hat{y}) \end{array} \quad W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

$$E(W) = \frac{1}{2m} (xw - y)^T (xw - y)$$

$$\frac{\partial E}{\partial W} = 2x^T(xw - y)$$

$$W = (x^T x)^{-1} x^T y \quad // \text{Closed form, aka Normal Equation}$$

Closed form solution to find model params/weights would be expensive, if n is large, because finding inverse of matrix $x^T x$ is computationally expensive for high dimension matrices.

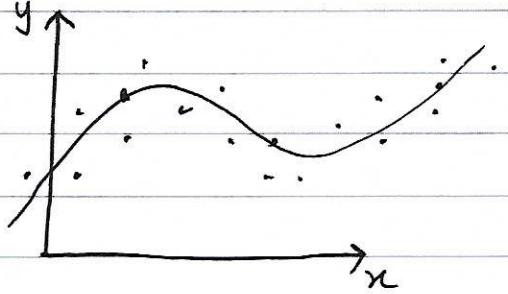
LINEAR REGRESSION WARNINGS

- * LR works best when the data is linear. If not, you'll need to make adjustments (transform your training data), add features, or use another kind of model.

- * LR is sensitive to outliers.

POLYNOMIAL REGRESSION

Transforming features to improve the form/fit of hypothesis function



$$\hat{y} = w_1 x^3 + w_2 x^2 + w_3 x + w_4 \text{ (example)}$$

REGULARISATION

Addreses problem of model 'overfitting'.

L1 Regularisation

Regularisation will remove features from a model (by setting their coeff. to 0) if the penalty for removing them is small. This makes model simpler & generalised

L2 Regularisation

$$-\text{Error}_{\text{new}} = \text{Error}_{\text{without Regularisation}} + \lambda \sum_{i=1}^n |w_i|$$

$$-\text{Error}_{\text{new}} = \text{Error} + \lambda \sum w_i^2$$

- Computationally inefficient (unless data is sparse)

- Computationally efficient

FEATURE SCALING

Transforming training data features into a common range of values.

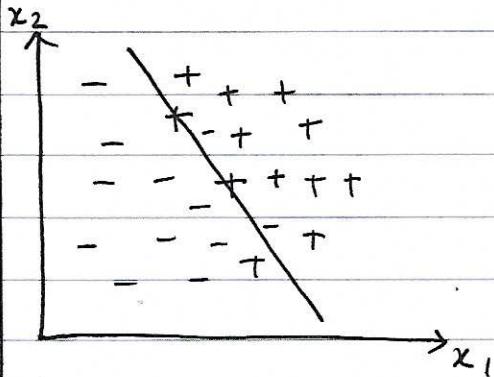
$$x_j := \frac{x_j - \mu_j}{s_j}$$

$$; df['h-new'] = \frac{df['h'] - df['h'].mean()}{df['h'].std()}$$

- can speed-up convergence of gradient descent
- Should always be done with regularised methods

PERCEPTRON ALGORITHM

LINEAR BOUNDARIES



Boundary line :

$$\omega_1 x_1 + \omega_2 x_2 + b = 0$$

$$Wx + b = 0, \text{ where}$$

$$W = (\omega_1, \omega_2) \text{ (weights)}$$

$$X = (x_1, x_2)$$

$$b = \text{bias}$$

Prediction

$$\hat{y} = \begin{cases} 1 & \text{if } Wx + b \geq 0 \\ 0 & \text{if } Wx + b < 0 \end{cases}$$

HIGHER DIMENSIONS

Boundary plane:

$$\omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3 + b = 0$$

$$Wx + b = 0$$

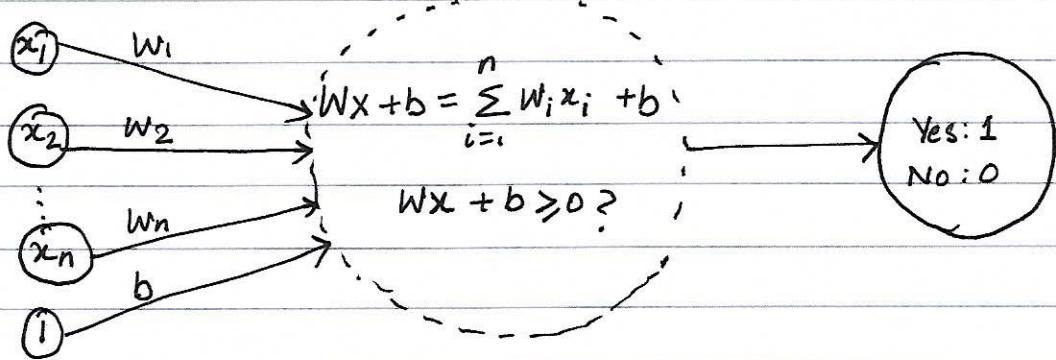
n -dimensional space

$$\omega_1 x_1 + \omega_2 x_2 + \dots + \omega_n x_n + b = 0$$

$$Wx + b = 0$$

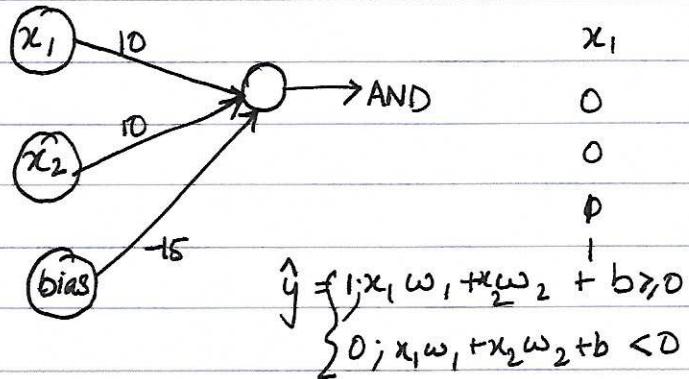
$$\hat{y} = \begin{cases} 1 & \text{if } Wx + b \geq 0 \\ 0 & \text{if } Wx + b < 0 \end{cases}$$

PERCEPTRONS

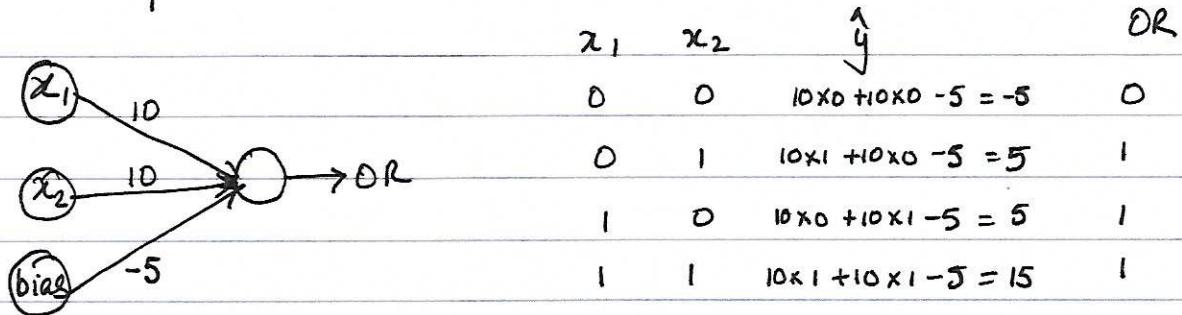


PERCEPTRONS AS LOGICAL OPERATORS

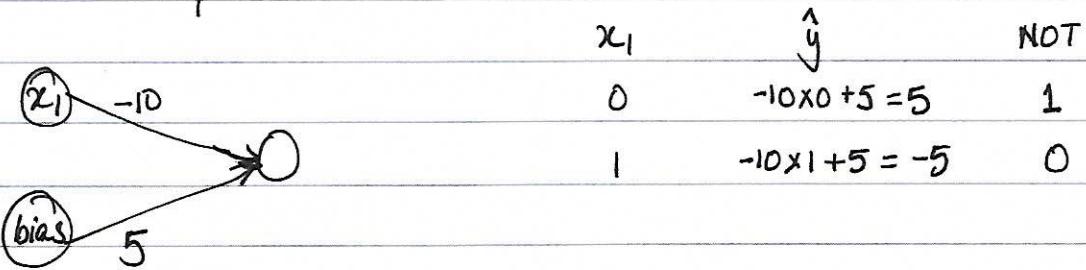
AND Perception



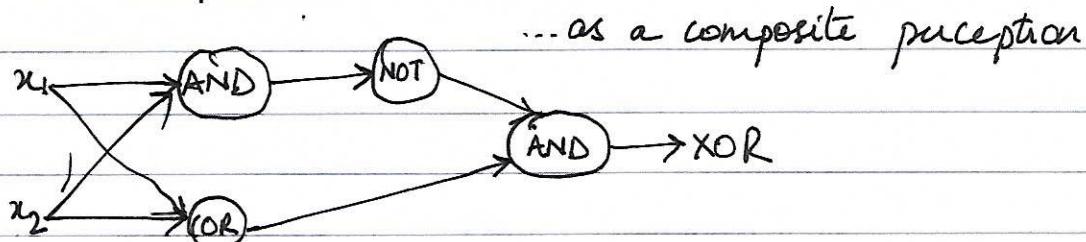
OR Perception



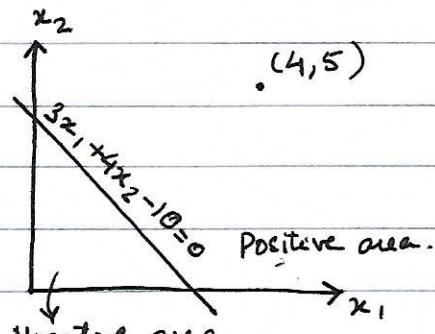
NOT Perception



XOR Perception



PERCEPTRON TRICK



$$\text{Line : } 3x_1 + 4x_2 - 10 = 0$$

Point : $(4, 5)$

Learning rate : $0 \cdot 1$

$$\begin{array}{cccc} & 3 & 4 & -10 \\ -0.4 & \underline{0.5} & 0.1 \\ 2.6 & 3.5 & -10.1 \end{array}$$

$$\text{New line : } 2.6x_1 + 3.5x_2 - 10.1 = 0$$

(if point is below line, add params to coeff. instead of subtracting)

PERCEPTRON ALGORITHM

1. Start with random weights

$$w_1, w_2, \dots, w_n, b$$

2. For every misclassified point (x_1, x_2, \dots, x_n) :

a. If prediction = 0, i.e. in negative area,

For $i = 1..n$

$$\text{Change } w_i = w_i + \alpha x_i$$

$$b = b + \alpha$$

b. If prediction in positive area or prediction = 1,

For $i = 1..n$

$$\text{Change } w_i = w_i - \alpha x_i$$

$$b = b - \alpha$$

DECISION TREES

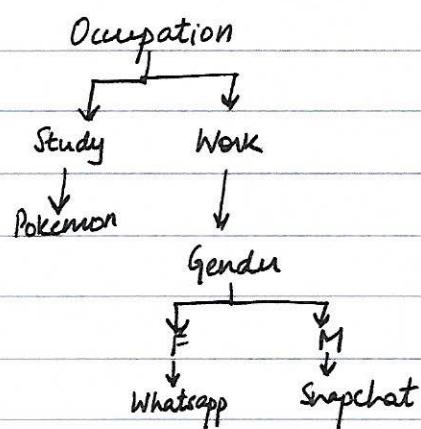
EXAMPLE: RECOMMENDING APPS

(REGRESSION + CLASSIFICATION)

Training Data:

GENDER	OCCUPATION	APP
F	Study	Pokemon
F	Work	Whatsapp
M	Work	Snapchat
F	Work	Whatsapp
M	Study	Pokemon
M	Study	Pokemon

Decision Tree:



ENTROPY

[Decision Trees tend to suffer from high variance or overfitting.]

Ⓐ Ⓑ Ⓒ Ⓓ

High knowledge.

Low Entropy

Ⓐ Ⓑ Ⓒ Ⓓ Ⓔ

Medium knowledge

Medium Entropy

Ⓐ Ⓑ Ⓒ Ⓓ

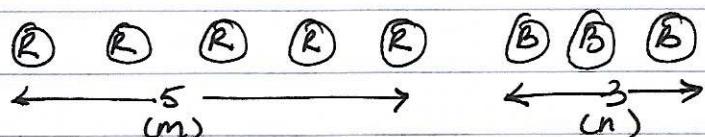
Low Knowledge

High Entropy

Probability of Winning.

	P(red)	P(blue)	P(winning)	$-\log_2 P(\text{win})$
Ⓐ Ⓑ Ⓒ Ⓓ	1	0	$1 \times 1 \times 1 \times 1 = 1$	$0+0+0+0 = 0$
Ⓐ Ⓑ Ⓒ Ⓓ Ⓔ	$\frac{3}{4}$	$\frac{1}{4}$	$\frac{3}{4} \times \frac{3}{4} \times \frac{3}{4} \times \frac{1}{4} = 0.105$	$0.415 + 0.415 + 0.415 + 2 = 3.245$
Ⓐ Ⓑ Ⓒ Ⓓ Ⓔ	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} = 0.0625$	$1+1+1+1 = 4$

Average Entropy



$$(\text{Avg}) \text{ Entropy} = -\frac{5}{8} \log_2 \frac{5}{8} - \frac{3}{8} \log_2 \left(\frac{3}{8} \right) = 0.9544$$

$$\text{General Form} = -\frac{m}{m+n} \log_2 \left(\frac{m}{m+n} \right) - \frac{n}{m+n} \log_2 \left(\frac{n}{m+n} \right)$$

Let $P_1 = \frac{m}{m+n}$; probability that answer is 'm'

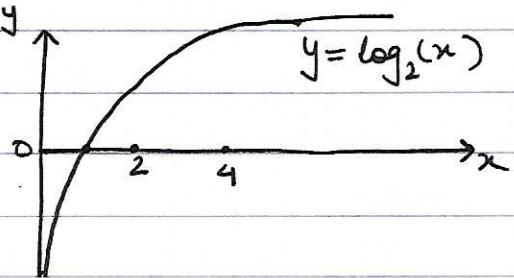
$$\& P_2 = \frac{n}{m+n}$$

$$\text{Thus, Entropy} = -P_1 \log_2(P_1) - P_2 \log_2(P_2)$$

MULTICLASS ENTROPY

$$-P_1 \log_2(P_1) - P_2 \log_2(P_2) - \dots - P_n \log_2(P_n)$$

$$= -\sum_{i=1}^n P_i \log_2(P_i)$$



INFORMATION GAIN

Change in entropy

$$\text{Information gain of a decision step} = \text{Entropy(parent)} - \frac{1}{2} [E(\text{child1}) + E(\text{child2})]$$

Example: Recommending Apps

Gender	Occupation	App			
F	Study	Pokemon			
F	Work	WhatsApp	P	P	P
M	Work	Snapchat	W	W	
F	Work	WhatsApp			S
M	Study	Pokemon			
M	Study	Pokemon			

$$\text{Entropy} = -\frac{3}{6} \log_2\left(\frac{3}{6}\right) - \frac{2}{6} \log_2\left(\frac{2}{6}\right) - \frac{1}{6} \log_2\left(\frac{1}{6}\right) = 1.46$$

Split by gender :	F	M
	P W W	S P P
Entropy	0.92	0.92

∴ Information gain, when split by gender,

$$IG = 1.46 - \frac{1}{2}(0.92 + 0.92) = 0.54$$

Similarly, when split by occupation,

$$IG = 1.46 - \frac{1}{2}(0 + 0.92) = 1$$

∴ It is best to split the original data set by occupation.

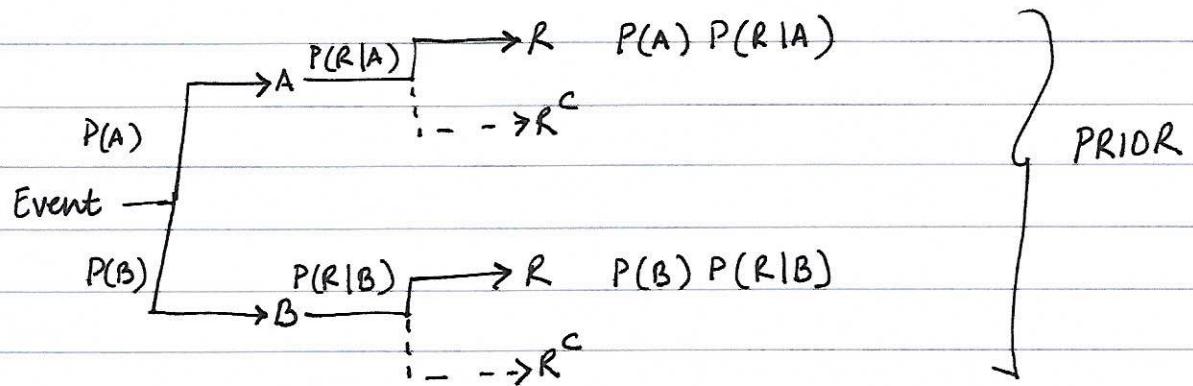
HYPERPARAMETERS

Aspects of the decision tree that are tuned to ensure that it generalises well to new problems.

- Maximum Depth. Large depth often causes overfitting. Too small = underfit.
- Minimum no. of ~~posen~~ samples to split
- Minimum no. of samples per leaf. Small = overfit; Large = underfit.

NAIVE BAYES

BAYES THEOREM



$$P(A|R) = \frac{P(A) P(R|A)}{P(A) P(R|A) + P(B) P(R|B)} \quad \text{] POSTERIOR}$$

$$P(B|R) = \frac{P(B) P(R|B)}{P(A) P(R|A) + P(B) P(R|B)} \quad \text{] POSTERIOR}$$

Alternatively,

Given,

$P(H)$ = Probability of hypothesis

$P(E|H)$ = Probability of evidence given hypothesis

$P(E)$ = Probability of evidence

(For multiple features)

$$P(H|E) = \frac{P(E|H) P(H)}{P(E)} \quad (\text{OR}) \quad P(y|x_1, \dots, x_n) = \frac{P(y) P(x_1, \dots, x_n|y)}{P(x_1, \dots, x_n)}$$

NAIVE ASSUMPTION

$$P(A \cap B) = P(A) P(B) \quad [\text{In practice, this is true if } A \text{ & } B \text{ are independent events}]$$

CONDITIONAL PROBABILITY

$$P(A|B) P(B) = P(B|A) P(A) \Rightarrow P(A|B) \propto P(B|A) P(A)$$

NAIVE BAYES ALGORITHM

$$P(\text{spam} | \text{easy, money}) \propto P(\text{easy, money} | \text{spam}) P(\text{spam})$$

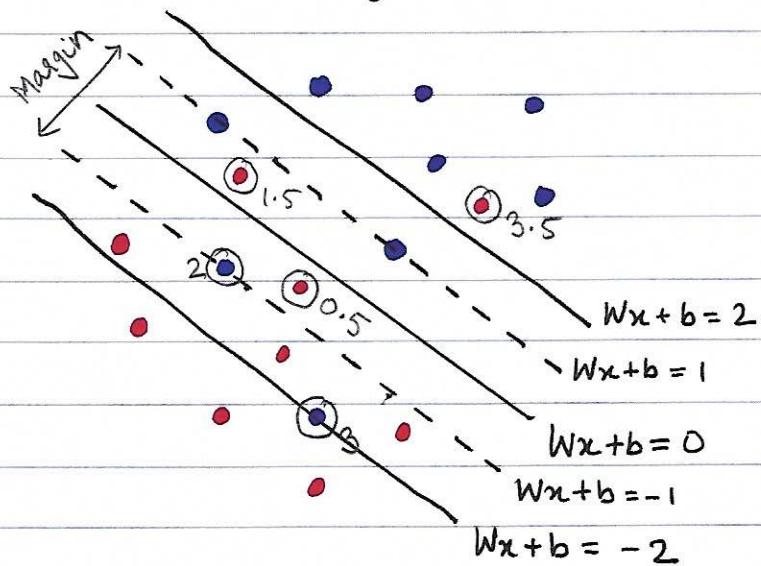
$$\propto P(\text{easy} | \text{spam}) P(\text{money} | \text{spam}) P(\text{spam})$$

(using naive assumption)

SUPPORT VECTOR MACHINES

ERROR FUNCTION

Error = Classification Error + Margin Error.



$$\text{Classification error} = (2+3) + (0.5 + 1.5 + 3.5)$$

i.e. distance of mis-classified points from margins.

Goal is to have as large a margin as possible.

Margin error

- Margin = $\frac{2}{|W|}$ where $|W| = \sqrt{w_1^2 + w_2^2 + \dots + w_n^2}$

- Error = $|W|^2$ (because, we want a large margin to give a small error & vice-versa)

- Formula of margin comes from formula for distance between two parallel lines

$$ax + by + c_1 = 0$$

$$ax + by + c_2 = 0$$

$$\text{distance} = \left| \frac{c_2 - c_1}{\sqrt{a^2 + b^2}} \right|$$

THE C PARAMETER

A hyper parameter used to create a model.

$$\text{Error} = C \cdot \text{Classification Error} + \text{Margin Error}$$

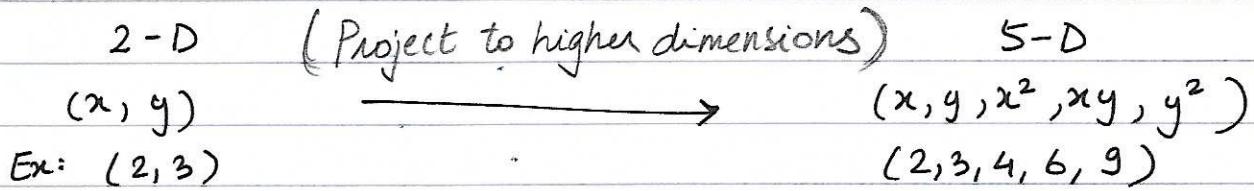
Small C

\Rightarrow Emphasis on margin error \Rightarrow Large Margin \Rightarrow Emphasis on classification error
 May make classification Errors \Rightarrow Classifies points well ; May have small margin

Large C

POLYNOMIAL KERNEL

When a set of points can't be separated linearly, project data points from 2-Dimensions to say, 5-Dimensions.



Now, when data lives in higher dimensional space, we may get lucky and a 4-D hyperplane* may act as a good classification boundary.

(Find Hyperplane)

* Hyperplane : a subspace whose dimensions is one less than that of its ambient space.

If such a hyperplane is found, we project it back to 2-D, obtaining a degree-2 polynomial boundary.

(Project back into higher degree polynomial).

A kernel is simply a set of functions that may act as a classification / decision boundary.

Example :

(1) Linear kernel (x, y)

Candidate functions:

$$2x - y = 1$$

$$5x - 2y = 3, \text{ etc}$$

Degree of polynomial is a HYPERPARAMETER that is tuned to obtain best model.

(2) Polynomial kernel. (x, y, x^2, xy, y^2) (degree 2)

$$x^2 + y^2 = 1$$

$$xy = 1$$

$$y = x^2$$

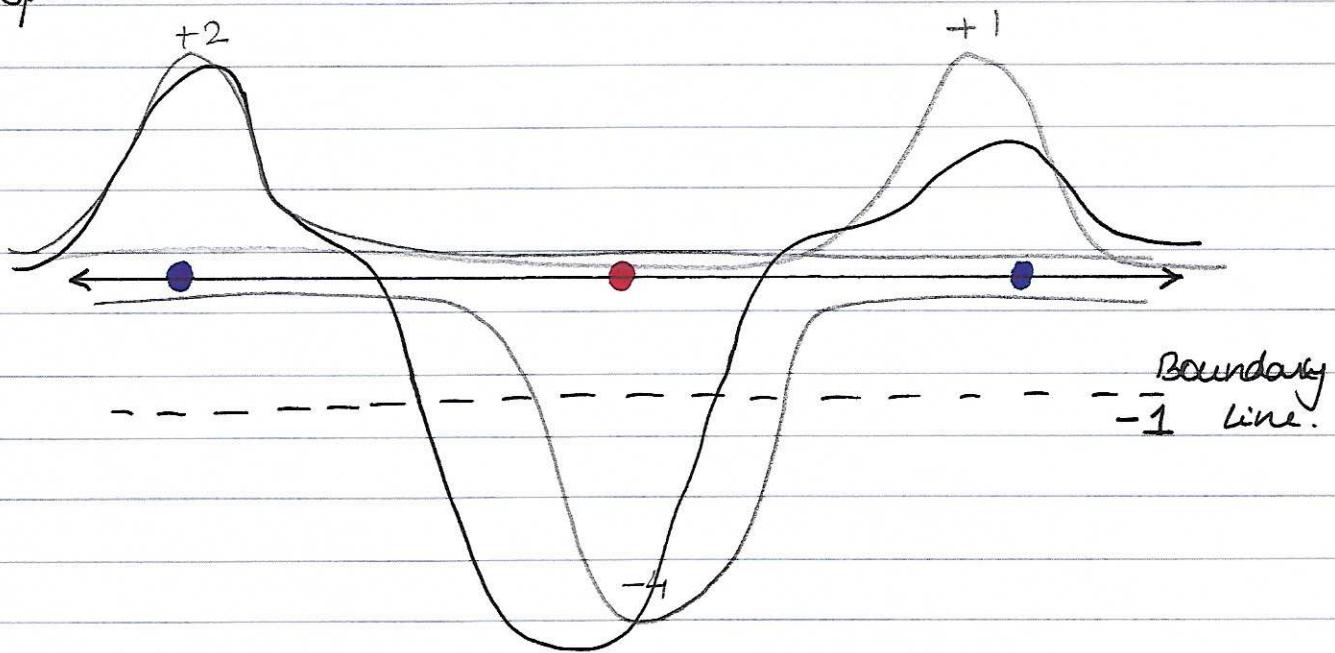
(3) Polynomial kernel $(x, y, \dots, x^3, x^2y, \dots, y^3)$ (degree 3)

$$y = x^3 + 2x^2 - x - 2$$

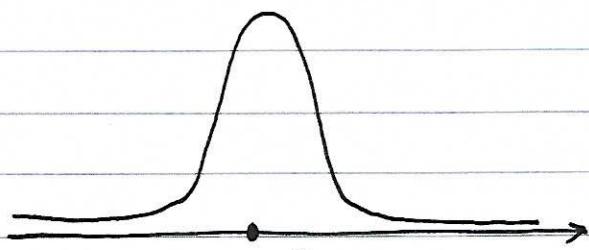
RBF KERNEL

RBF : Radial Basis Function

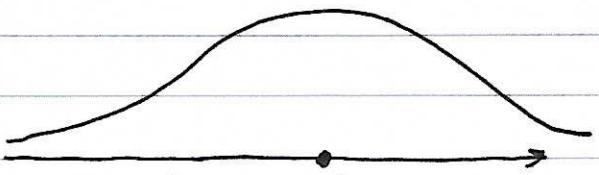
Used to classify points that seem hard to separate in any space.

 γ PARAMETER

γ (Gamma) is a hyperparameter used to tune SPF - RBF kernel models.



Large γ
Tends to overfit



Small γ
Tend to underfit

ENSEMBLE METHODS

(REGRESSION + CLASSIFICATION)

INTRODUCTION

Combining models such that the composite model is better at predicting than the individual models.

Commonly used 'weak' learners in ensemble models are decision trees. (* Decision trees are high variance / overfitting algorithms)

RANDOMNESS

Introducing randomness into high variance algorithms before they are ensembled together. The introduction of randomness combats the tendency of these algorithms to overfit.

Ways to introduce randomness:

(1) Bootstrap the data : sampling the data with replacement and fitting your algorithm to the sampled data

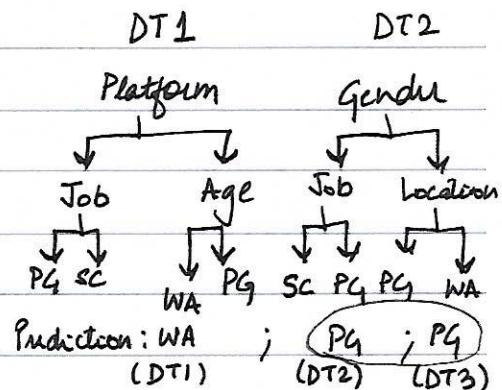
(2) Subset the features : in each split of a decision tree or with each algorithm used in an ensemble, only ~~of~~ a subset of the total possible features are used.

RANDOM FORESTS

Given training data with 'n' features, instead of creating one decision tree using all features, create different trees using a subset of the 'n' features (randomly). Then, get prediction from all trees and the most commonly occurring prediction is the final prediction.

Example -

Gender	Age	Location	Platform	Job	Hobby	App
F	15	US	iOS	School	Tennis	WA
F	25	FR	Andri.	Work	Chess	PG
M	32	UK	Andri.	School	Chess	SC



BAGGING (BOOTSTRAP AGGREGATING)

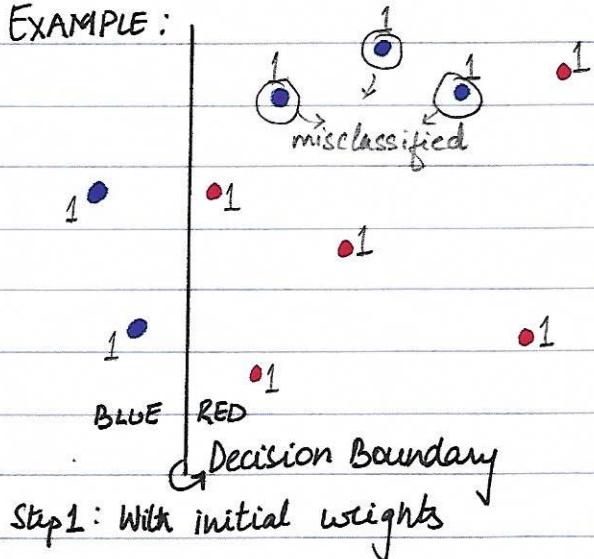
Train multiple 'weak' learners on different subsets of training data.

Then combine the weak learners using an approach such as voting, i.e. prediction made by largest majority of weak learners.

ADABOOST

- (1) Create first weak learner (decision tree) such that it has max accuracy.
- (2) 'Boost' misclassified points by first weak learner, such that the second weak learner will be punished more / generate higher error, if it misclassifies points misclassified by the first learner
- (3) And so on.

EXAMPLE:



Correct : 7
Incorrect : 3

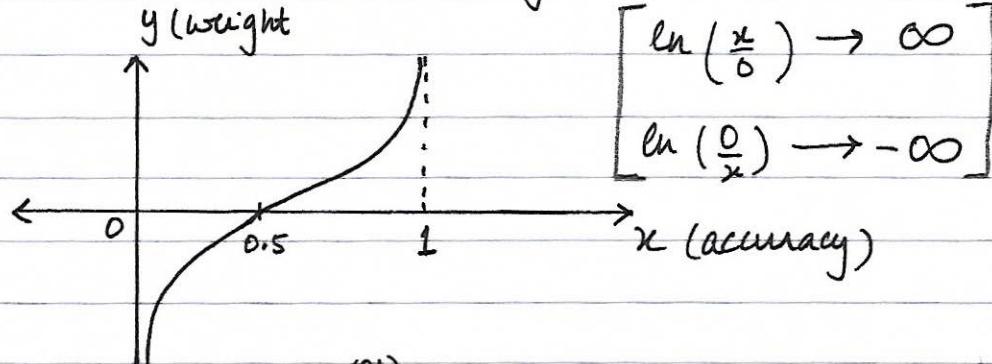
Boost misclassified points by a factor of $7/3$, such that their weight becomes : $1 \times \frac{7}{3} = 2.33$

With new/boosted weights,
Correct : 7 (1×7)
Incorrect : 7 ($\frac{7}{3} \times 3$)

In step 2, create a new weak learner that prioritises points misclassified in step 1. & seq & repeat a few times.
repetitions or steps = Hyperparameters.

To combine the models / weak learners in multiple AdaBoost steps, compute a weight for each weak learner using the following formula.

$$\text{weight} = \ln\left(\frac{\text{accuracy}}{1 - \text{accuracy}}\right) = \ln\left(\frac{\#\text{correct}}{\#\text{incorrect}}\right)$$



Now, prediction = $\begin{cases} 1 & \text{if sum of weights of learners is } > 0 \text{ or positive} \\ 0 & \text{otherwise.} \end{cases}$

Example,

B	Red
-0.84	-0.84

Model 1

$$\text{Weight} = 0.84$$

Blue 1.3
Red -1.3

Model 2

$$\text{Weight} = 1.3$$

Blue	R
1.84	-1.84

Model 3

$$\text{Weight} = 1.84$$

Combining

0.84	-0.84	
1.3	1.3	
1.84	1.84	
+0.84	-0.84	-0.84
-1.3	-1.3	-1.3
1.84	+1.84	-1.84

=

B	B	R
+	-	-

Ensemble methods can also be extended to regression problems. Not just for classification.

MODEL EVALUATION METRICS

CONFUSION MATRIX

	Diagnosed Sick	Diagnosed Healthy
Sick	1000 TRUE POS.	200 FALSE NEG.
Healthy	800 FALSE POS.	8000 TRUE NEG.

False Positive = Type 1 Error | False Negative = Type 2 Error

ACCURACY

$$\text{Accuracy} = \frac{\text{True Pos} + \text{True Neg}}{\text{False Pos} + \text{False Neg} + \text{Total}}$$

PRECISION AND RECALL

Medical Model

False Pos = OK

False Neg = NOT OK

OK if not all are sick

But find all sick people

Spam Detector

False Pos = OK

False Neg = NOT OK

OK if not all spam is found,

But if marked spam, Better be spam!

HIGH RECALL

HIGH PRECISION

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

$$\text{Recall} = \frac{\# \text{True Pos}}{\# \text{True Pos} + \# \text{False Neg}}$$

$$\text{F1 Score (as harmonic mean)} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

F_β SCORE

$$F_\beta = \frac{(1+\beta^2) \cdot \text{Precision} \times \text{Recall}}{\beta^2 \times \text{Precision} + \text{Recall}}$$

If $\beta=0$, then we get precision

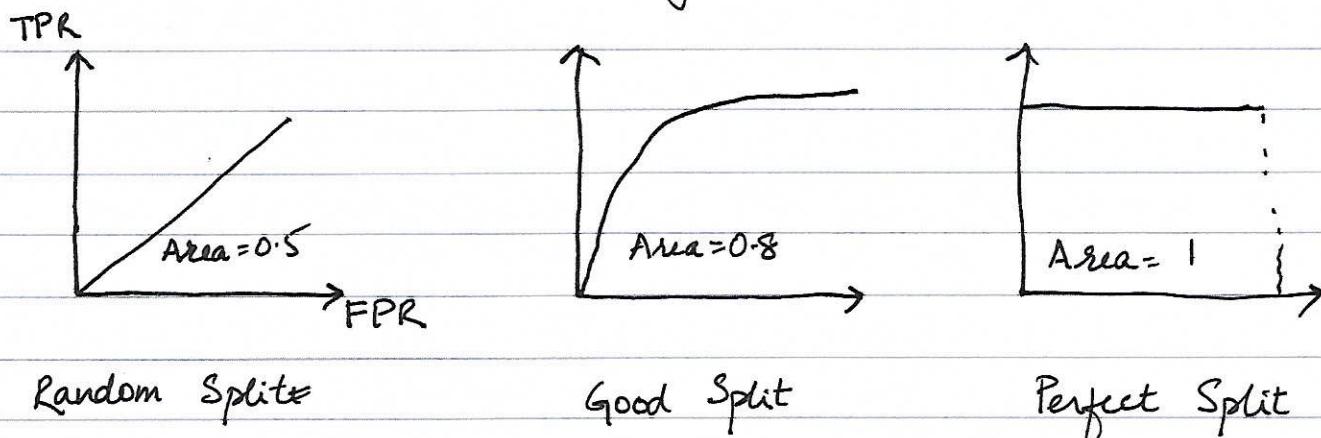
$\beta=\infty$, then we get recall

$\beta=1$, then we get harmonic mean or F1 score

RECEIVER OPERATING CHARACTERISTIC (ROC) CURVE

True Positive Rate = $\frac{\text{True Positives}}{\text{All Positives}}$

False Positive Rate = $\frac{\text{False Positives}}{\text{All Negatives}}$



R² SCORE FOR REGRESSION MODELS

$$R^2 = 1 - \frac{\text{MSE of optimised Linear Reg. Model}}{\text{MSE of simple average Model.}}$$

Bad Model = $R^2 \sim 0$

Good Model = $R^2 \rightarrow 1$