

ML algorithms

1) Linear Regression

Type : supervised learning, Regression

Linear regression models the relationship b/w a dependent variable & one or more independent variables using a linear equation.

The model predicts the output by finding the best fit line that minimizes the diff. b/w predicted & actual value

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \dots \text{ when } \beta \text{ are coefficient}$$

Minimises the mean squared error (MSE) using technique like Ordinary Least Squares.

Use case : Sales forecasting, house price prediction, stock price forecasting.

2) Logistic Regression

Type : Supervised learning, Classification

Logistic regression predicts the probability of a binary outcome.

Instead of a line, it uses a

Sigmoid function to output probability b/w 0 & 1.

Use case : spam detection, disease prediction, credit scoring.

3) Decision Trees :

Type : supervised learning, classification & Regression.

Decision trees split the data into subsets based on features that provides the max. information gain or min. Gini Impurity.

Use case : customer segmentation, risk assessment

4.) Random Forest :

Type : Ensemble Learning, Classification & Regression.

Random Forest is an ensemble of decision trees trained on random samples of data. The model averages the o/p of the individual trees to improve generalisation.

Use cases : Loan approval, image recognition, fraud detection..

5) Support Vector Machines (SVM):

Type : Supervised learning,
classification

SVM finds the hyperplane
that best separates the classes
by maximizing the margin
by support vectors

Use case : Image classification,
text categorisation,
~~3~~

6.) k-nearest Neighbours

Types : supervised learning,
classification,
Regression

Description : k-NN classifies
a new data point by
considering the labels
of the \sqrt{k} closest
points in the training
data

Use cases : Recommender Systems,
anomaly detection,
pattern recognition.

7.) k - Means Clustering :

Type : Unsupervised Learning, Clustering

k - Means divides dataset into k clusters by minimising the variance ~~of~~ within each cluster.

Use cases : Customer segmentation, market research, image compression.

8.) Principal Component Analysis :

Type : Unsupervised Learning, Dimensionality Reduction.

PCA reduces dimensionality by projecting data into a lower - dimensional space while preserving max. variance.

Use cases : Noise reduction, visualization, feature extraction.

ML algorithms

Classical Machine Learning

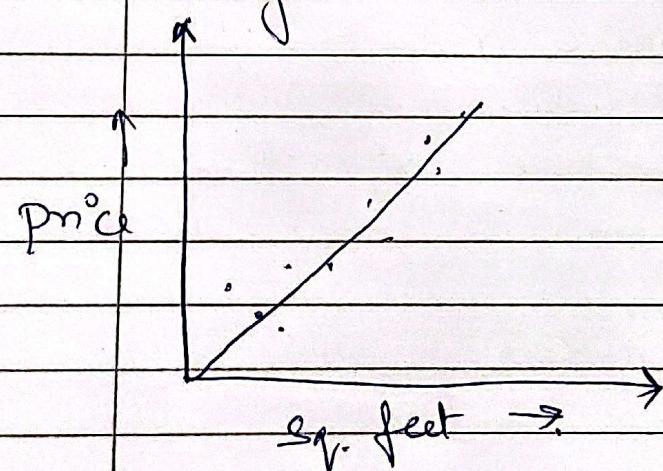
Task driven

Supervised
Learning

Training dataset
where we
know the
true value
for the o/p
dataset (labels)
that we can
train our
algorithm on
to later
predict unknown
data

e.g.

Use linear
regression



Predict the price
of a house.

Data driven

Unsupervised
learning

Any ~~dear~~ learning
problem which
is not
supervised ie.
no truth
about the
data is
known

You can categorise
data based
on some factors
& name them

as you
wish

Supervised Learning

Regression

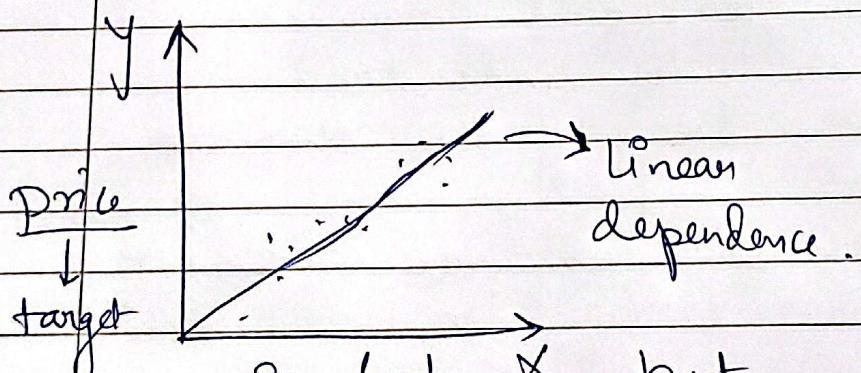
Predict a continuous numeric target variable for a given i/p variable

e.g. Predict the price of a house given a no. of features of the house & determining its relationship to the house

Classification

Predict discrete categorical variable ("label / class")

e.g. we may want to assign a label 'SPAM' or 'NO SPAM' on an email based on its sender, content etc, there could be more than 2 label like junk etc.



Sq. foot. \times but
feature like
feature
age of house
has no effect on Price

Regression

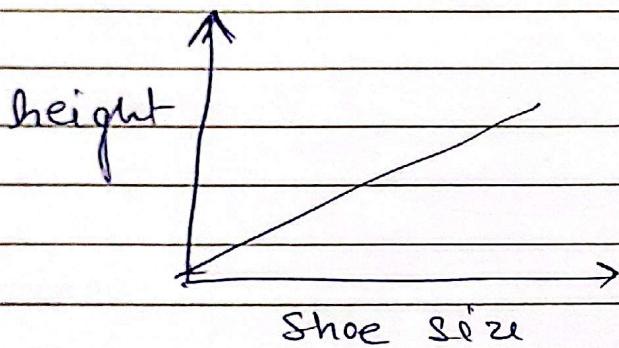
1) Linear Regression

trying to determine linear relationship
b/w 2 variables ie. b/w I/p & O/p

$$Y_i = \beta_0 + \beta_1 X_i$$

↑ ↑ ↑
dependent constant Independent
variable variable

e.g. height & shoe size



for every shoe size increase, person
will be β_1 inches taller
You can include multi-dimensional
data e.g. you could include age,
gender etc. to understand it'
a better idea of the shoe size

Neural networks is an extension
of linear regression.

Classification

Logistical Regression :

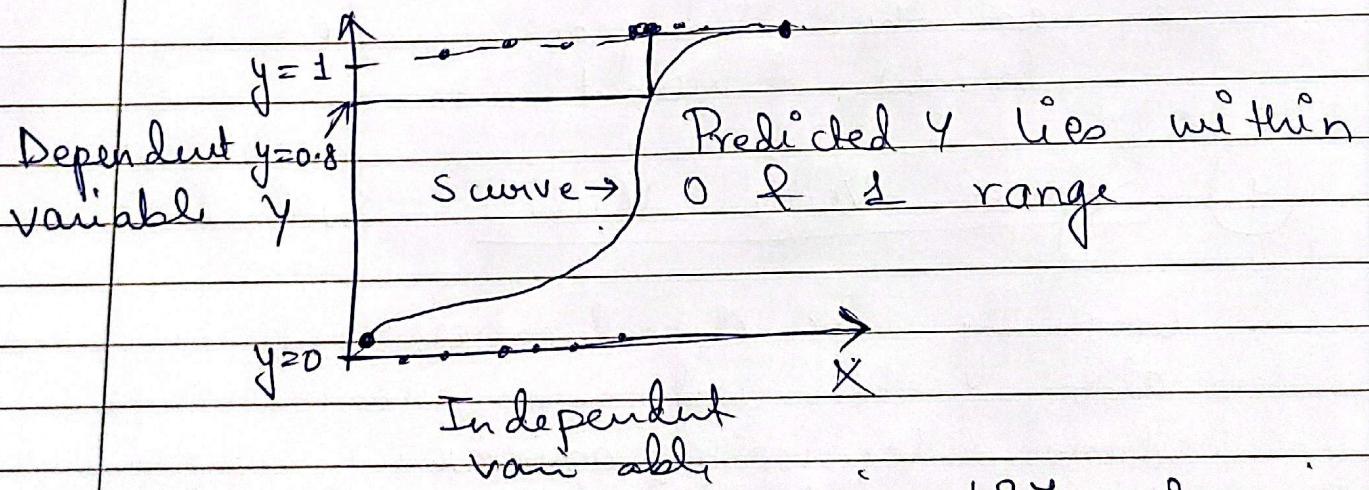
Instead of fitting a line b/w 2 numerical variables with a presumably linear relationship, you now predict a categorical output variable using either categorical or numerical input variables.

Instead of fitting line to data we fit a sigmoid function to the data.

eg. We want to predict ~~1~~^{1 out of 2} classes.

① Gender of a person based on height & weight

So, linear relationship would not work here but sigmoid fn would help in determining the possibility of the data point to fall into one class.

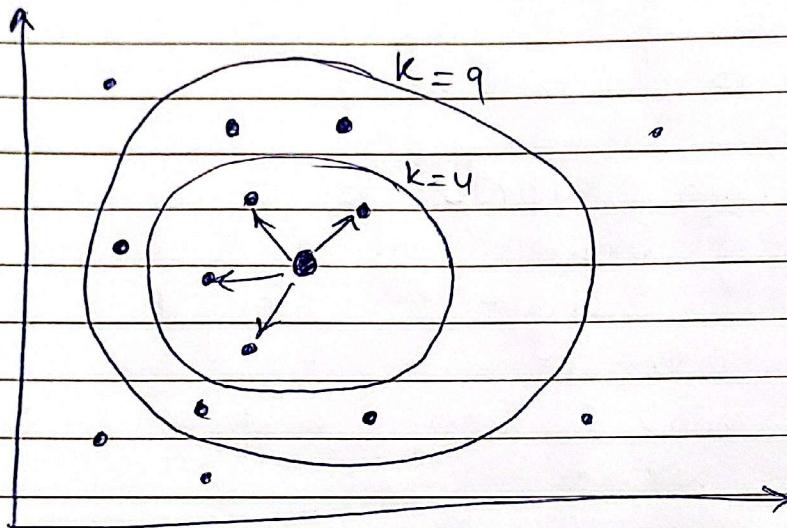


Let's say height is 1.87 meters in a man (\rightarrow likelihood $\rightarrow y = 0.8 \rightarrow 80\%$)

③ K-nearest Algorithm

Can be used for both Classification & regression

No equation
→ for any given data point we will find the target to be the avg. of its k -nearest neighbours.

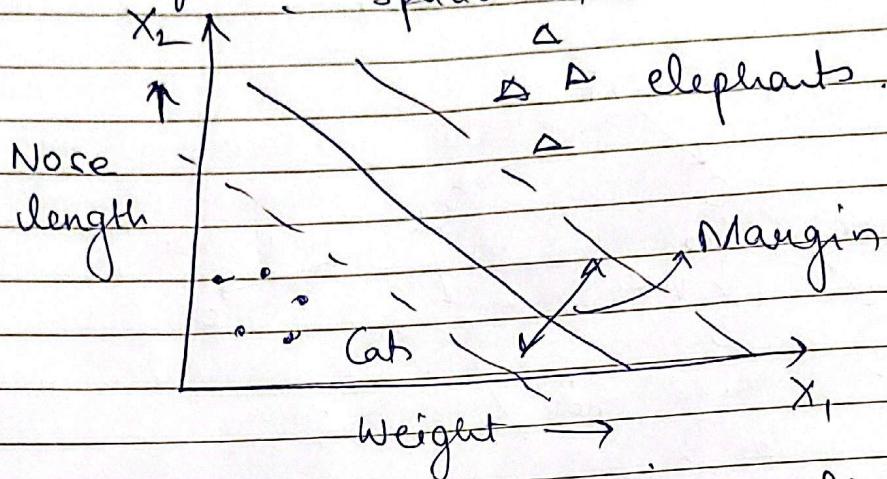


e.g. Gender of a person is the same as the 5 people closest in height & weight to the person.

④ Support Vector Machine (SVM)

Originally for classification, but can also be used for regression, draw decision boundaries that separate data points.

SVM separates classes using the largest margin possible by maximising space b/w each class



Very powerful in high dimensional classes where the decision boundary is called a hyperplane

(S)

Naive Bayes Classifier

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

prob. of B occurring given A has already occurred

prob. of A occurring given B has already occurred

prob. of B occurring

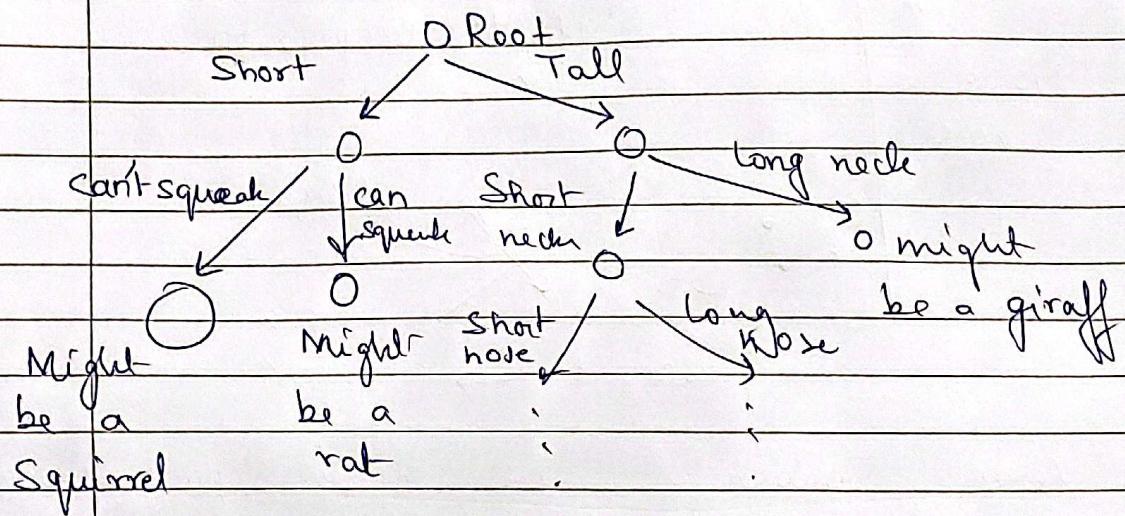
prob. of A occurring

good for Spam classification & other text based classification

6

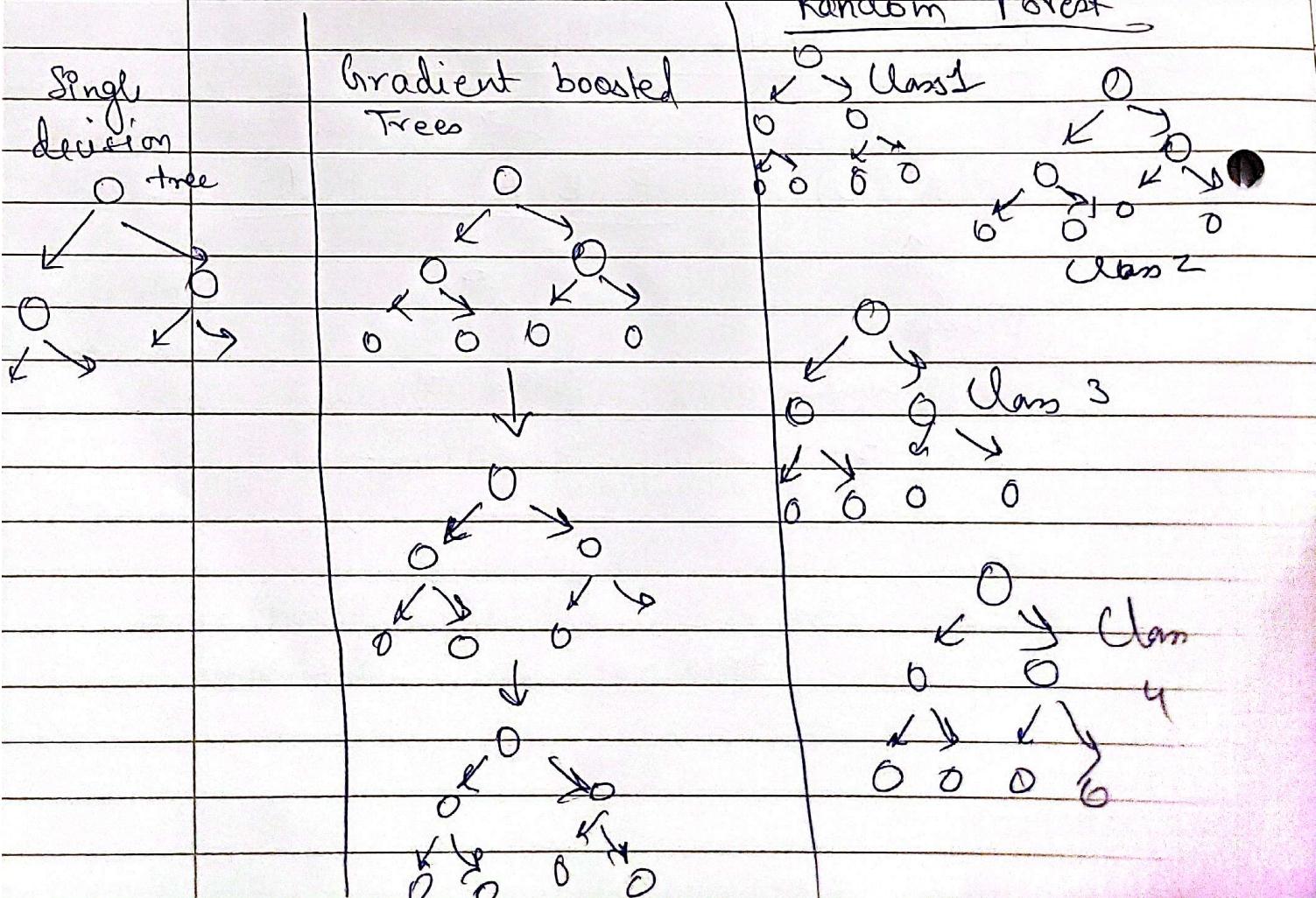
Decision Tree

e.g. decision tree animal classifier



Series of a Yes / No questions, helping us to partition the dataset

Random Forest



Unsupervised Learning.

Clustering
(Divide by similarity)

eg. Targetted marketing

Association
(Identify sequences)

eg. customer recommendation

Dimensionality Reduction
(wider dependencies)

eg. Big data visualization

I) Clustering (diff. from classification as we know the classes we need to predict in classification)

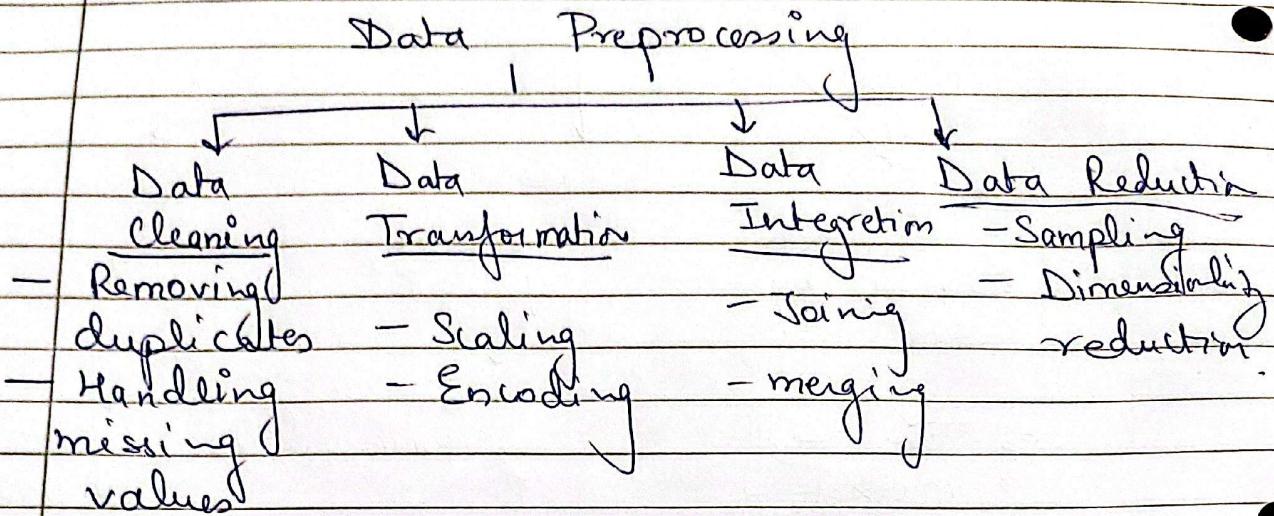
I) k - means clustering

Randomly selecting k centres for your k clusters, & assigning all data points to the cluster center closest to them.

Then recalculate the cluster centers based on the data points now assigned to them & repeat the process until centers of the clusters have stabilised.

II) Dimensionality Reduction

Reduces the dimensions of your dataset, keeping as much info as possible.



III Principle Component Analysis (PCA)

Finds the direction in which most variance in the dataset is retained.