

NAME: RAKSHIT P

1) SVM

According to what I read, for very large datasets SVM or neural networks are preferred but I felt neural networks was a bit out of our scope.

Support vector machines are administered learning models that uses association r learning algorithm which analyze features and identified pattern knowledge, utilized for application classification. SVM can productively perform a regression utilizing the kernel trick, verifiably mapping their inputs into high-dimensional feature spaces.

2) Decision trees

It is used when we have to make a yes/no choice using lot of if conditions based on the features which seems to fit exactly to this type of model and dataset. This is also efficient when there are considerable amount of factors on which final output depends.

3) Naïve Bayes Classifier

Naive Bayes model is easy to build and particularly useful for very large data sets. Along with simplicity, Naive Bayes is known to outperform even highly sophisticated classification methods. In simple terms, a Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature. Even if these features depend on each other or upon the existence of the other features, all of these properties independently contribute to the probability.

4) Logistic regression

This is a classification algorithm and not regression.

Logistic regression performs binary classification, so the label outputs are binary which is exactly what we have to do. Let's define  $P(y=1|x)$  as the conditional probability that the output  $y$  is 1 under the condition that there is given the input feature vector  $x$ . The coefficients  $w$  are the weights that the model wants to learn. Since this algorithm calculates the probability of belonging to each class, you should take into account how much the probability differs from 0 or 1 and average it over all objects as we did with linear regression

5) K-nearest neighbour

It is highly efficient classification algorithm but may not be suitable for large datasets. it takes a bunch of labelled points and uses them to learn how to label other points. To label a new point, it looks at the labelled points closest to that new point (those are its nearest neighbors), and has those neighbors vote, so whichever label the most of the neighbors have is the label for the new point.