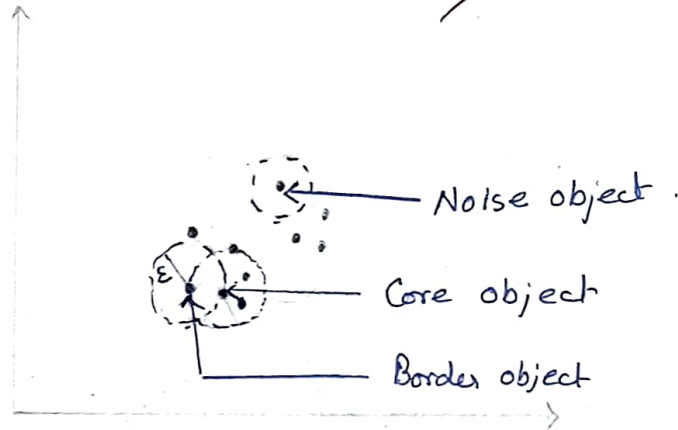1a> i> Core object - An object with atleast MinPt number of (objects) points in
it's ε-neighborhood is called a core object.

Border object - An object which does not have a MinPt number of
objects in it's ε-neighbourhood but is close to a core object. is
called border object.

Noise object - An object which is neither close to a core
object nor has MinPt number of objects in it's ε-neighbourhood
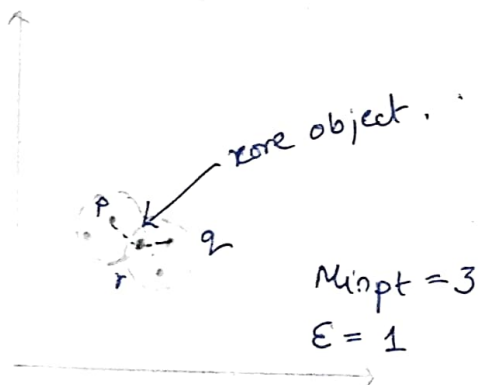is called noise object.

~~ii> Density reac~~



Noise object.

Core object

Border object

MinPt = 4
ε = 1.

(directly)

ii> Density reachability - A point q is said to be ⊥ density
reachable from a point p if p is a core object and q
lies in the ε-neighbourhood of p.

Two
⊙ Density connectivity - ~~A~~ points p and q are said to be
density connected if ~~p is directly density reachable from~~ there
exists a point r such that p is directly density reachable
from r and q is also directly density reachable from r.



core object.

q

r

Minpt = 3
ε = 1

∴ p & q are
density reachable from r.

b> In DBSCAN clustering algorithm,

i) Maximality condition states that if $p \in C$ where $p$ is an object and $C$ is a cluster, and $q$ is directly density reachable from $p$, then $q$ should also be in cluster $C$.

ii) Connectivity condition states that if $p, q \in C$ where $p$ and $q$ are two objects and $C$ is a cluster, then $p$ and $q$ should be density connected, $\forall p, q \in C$.

c> The algorithm for

Let the distance of the $k^{th}$ nearest neighbour be $k$-dist.

If $k$ is smaller than the size of the cluster, the $k$-dist will be small. If cluster size is too small, noise points may get incorrectly labeled with some cluster. However, if noise the clusters become too large then small points clusters will be labeled as noise.

To find the optimal values of the parameters Eps and MinPts, the $k$-dist for all points are computed and points are sorted according to the increasing $k$-dist.

The point where there is a sharp change of $k$-dist, will give the optimal size of cluster. This $k$-dist is takes as Eps and $k$ becomes MinPts.

2a) Given a gold standard data with $c$ class labels, and a clustering with $k$ clusters, denoted by $\{\omega_1, \omega_2 \cdots \omega_k\}$, the purity of the clustering is given by

$$\text{Purity}(\omega_i) = \frac{1}{n_i} \max_j \left(n_j^i\right) \quad \text{where} \quad j \in c$$

The ratio of the size of the dominant class in the $i^{th}$ cluster, $n_j^i$ to the size of the cluster $\omega_i$.

b) We calculate the Rand Index by calculating the number of pairs that are, in the same cluster and same class, that in different class, different clusters and same and different class in the following way

|  | Same clusters. | Different Clusters |
|---|---|---|
| Same labels | A | B |
| Different Labels. | C | D |

$$\text{Rand Index} = \frac{A+D}{A+B+C+D}$$

6) Intercluster distance is the measure of dissimilarity between different clusters in a partitioning and intracluster distance is the measure of dissimilarity within a cluster in the partitioning.

A good partitioning will try to maximize intercluster distance and minimise intracluster distance.

Two intercluster distances are:

i) simple linkage distances : $\delta(s,t) = \min\left\{d(x,y)\right\}_{x\in s, y\in t}$ where $s$ & $t$ are two clusters.

Average linkage distance :
$$\delta(s, t) = \frac{1}{|s||T|} \left\{ \sum_{x \in s, y \in T} d(x, y) \right\}$$

where $|c|$ denotes the size of a cluster $c$.

Two intracluster distances are

i) complete diameter distance given by
$$\Delta s = \max_{x, y \in s} \left\{ d(x, y) \right\} \qquad \checkmark$$

ii) Average diameter distance given by
$$\Delta s = \frac{1}{|s|(|s|-1)} \sum_{\substack{x, y \in s \\ x \neq y}} \left\{ d(x, y) \right\}$$

$\checkmark$

3) Dunn's cluster validation index is a cluster validation technique that combines intercluster and intracluster distances. and for a partitioning $U$ is given by

$$D Index(U) = \min_{1 \leq i \leq c} \left\{ \min_{\substack{1 \leq j \leq c \\ i \neq j}} \left\{ \frac{\delta(x_i, x_j)}{\max_{1 \leq k \leq c} (\Delta X_k)} \right\} \right\}$$

③

$\checkmark$

where $\delta(x_i, x_j)$ is the intercluster distance between $x_i$ and $x_j$ such that $x_i$ and $x_j \in C$ and $\Delta X_k$ is the intracluster distance of a cluster $X_k$.

The larger the value of Dunn's cluster validation index, the better the partitioning.

3a) Support is the ratio of the number of transactions that ~~aa~~ contains ~~itemset~~ the itemset ~~is $\{A\cup B\}$ appears the~~ to the total number of transactions, for an association rule $A \rightarrow B$

$$sup(A \rightarrow B) = Pr(A \cup B)$$

Confidence of the association rule, $A \rightarrow B$ is given by the ratio ~~the~~ of the number of transactions containing ~~the~~ the itemset $\{A\cup B\}$ to the number of transactions containing $A$.

③

$$confidence(A \rightarrow B) = Pr(B/A)$$

An itemset is referred to as ~~aa~~ a frequent item set is ~~to as~~ it's support is more ~~the can~~ than the minimum support and an association rule is referred to as an important rule if it's confidence is more than the minimum confidence value.

b) minsup = 0.3

min conf = 0.8

Items:

$I_1$: Bread $\qquad$ $sup(I_1) = 5/6 = 0.833$

$I_2$: Butter $\qquad$ $sup(I_2) = 3/6 = 0.5$

$I_3$: Milk $\qquad$ $sup(I_3) = 3/6 = 0.5$

$I_4$: Jelly $\qquad$ $sup(I_4) = 1/6 = 0.166$

$I_5$: Coke $\qquad$ $sup(I_5) = 2/6 = 0.33$

The frequent-1-itemset is { Bread, Butter, Milk, Coke } since their support values are greater than minsup.

$F_1 = $ {Bread}, {Butter}, {Milk}, {Coke}

$C_2 = $ { ~~Bread~~ {Bread, Butter}, {Bread, Milk}, {Bread, Coke} ~~a~~ {Butter, Milk}, {Butter, Coke}, {Milk, Coke} }

$Sup(\{Bread, Butter\}) = 3/6 = 0.5$ $\qquad$ $sup(\{Bread, Milk\}) = 2/6 = 0.33$

$\text{sup}(\{ \text{Bread, Coke} \}) = 1/6 = 0.166 < 0.3$ $\quad$ $\text{sup}(\{ \text{Butter, Coke} \}) = 0 < 0.3$

$\text{sup}(\{ \text{Butter, Milk} \}) = 1/6 = 0.166 < 0.3$ $\quad$ $\text{sup}(\{ \text{Milk, Coke} \}) = 1/6 = 0.166 < 0.3$

$F_2 = [ \{ \text{Bread, Butter} \} , \{ \text{Bread, Milk} \} ]$

$C_3 = [ \{ \text{Bread, Butter, Milk} \} ]$

- $\text{sup}(\{ \text{Bread, Butter, Milk} \} ) = 1/6 = 0.166 < 0.3$

$F_3 = \emptyset$

∴ The frequent itemset is $\{ \{ \text{Bread} \} \{ \text{Butter} \} \{ \text{Milk} \} \{ \text{Coke} \}, \{ \text{Bread, Butter} \},$

⑤ $\{ \text{Bread, Milk} \}, \}$.

$\text{conf}(\text{Bread} \rightarrow \text{Butter}) = 3/5 = 0.6 < 0.8$

$\text{conf}(\text{Bread} \rightarrow \text{Milk}) = 2/5 = 0.4 < 0.8$

$\text{conf}(\text{Milk} \rightarrow \text{Bread}) = 2/3 = 0.66 < 0.8$

$\text{conf}(\text{Butter} \rightarrow \text{Bread}) = 3/3 = 1, > 0.8$

∴ The association rule is $\{ \text{Butter} \rightarrow \text{Bread} \}$.

c) The major drawback of the a-priori algorithm is that it does not take importance factor of an itemset into account.

**6a)** **i)**

| Outlook | Humidity | Wind | Play Tennis (class variable) ground truth | Play Tennis (predicted label) |
|---------|----------|------|-------------------------------------------|-------------------------------|
| Sunny | Normal | Strong | Yes | Yes |
| Overcast | Normal | Strong | No | Yes |
| Rain | High | Strong | Yes | No |
| Sunny | High | Weak | No | No |
| Rain | High | Strong | No | No |

Confusion Matrix :.

| Actual/Predicted class class | Play Tennis ⊘ | ¬ Play Tennis |
|------------------------------|----------------|----------------|
| Play Tennis | 1 (True Positive) | 1 (False Negative) |
| ¬ Play Tennis | 1 (False Positive) | 2. (True Negative) |

**④**

**i)** $\text{Precision} = \dfrac{\text{True Positive}}{\text{True Positive} + \text{False Positive}} = \dfrac{1}{2} = 0.5$

**ii)** $\text{Recall} = \dfrac{\text{True Positive}}{\text{True Positive} + \text{False Negative}} = \dfrac{1}{2} = 0.5$

**iii)** $\text{F-score} = \dfrac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}} = \dfrac{2 \times \frac{1}{2} \times \frac{1}{2}}{\frac{1}{2} + \frac{1}{2}} = \frac{1}{2} = 0.5$

**b)** Holdout method — The training and test set is sampled uniformly with the training data being $2/3^{rd}$ of the total number of data samples and testing data being $1/3^{rd}$ of

The model is trained on the training set and evaluated on the test set t times and the accuracy reported is the average of all the t observed accuracies.

ii) Cross-validation - The total number of data samples is divided into k subsets and in each iteration then model is trained on k-1 subsets (leave one out) and tested on the remaining subset.

③ In stratified cross validation ensures that the distribution of the classes in the training set sample is same as that in the original data.

iii) Bootstrap is the process of uniformly sampling the training and dataset training with replacement and from the given data samples. the classifier with it iteratively. It is observed that 63.2 % data goes to the training set and 36.8 % data is unseen by the model. $\frac{99}{6}$

iv) Ensembling classifiers reduce the error rate.
For example, if we ensemble 25 classifiers, each with error rate $E = 0.38$, by ensembling or combining different classifiers by averaging or voting etc reduces the error rate.
The error rate of the ensembled classifier is given by

$$\sum_{i=13}^{25} \binom{25}{i} \varepsilon^i (1 - \varepsilon)^{25-i} = 0.06.$$

The main purpose of ensembles of classifiers is not to find a highly accurate model instead to combine base models which differ in the type of misclassification or errors. It also overcomes the problem of statistical problem of a larger hypothesis space and small number of samples, representational problem etc.

The Adaboost algorithm works as follows :

● Designing :
Boost a
  i) Assign equal weight, $1/N$ to all classes.
     Randomly
     complementary

  ii) Build classifiers and predict the label of the given
      sample.

  iii) Calculate the error rate $\varepsilon$.

  iv) Assign If the error rate $\varepsilon$ is 0 or $\varepsilon$ is greater than or
      equal to 0.5 then then
                terminate the process.

  v) Assign the weight $\varepsilon/(1-\varepsilon)$ to all correctly classified labels.

  vi) Return the weights of the labels.


Classification :
   Initialize weight of all labels with equal weight
For each classifier, the feedback add $-\log(\varepsilon/1-\varepsilon)$ to the
weight of the predicted label misclassified label.

Return the label with highest weight.

5a) $C_1 = $ ~~Play Tennis~~ Yes.

$C_2 = $ ~~Play Tennis~~ No

Naive Bayes classification algorithm makes the assumption that all the attributes are conditionally independent.

Hence $P(X|C_i) = \prod\limits_{x_k \in X} P(x_k | C_i)$

Total number of samples $= 10$

$P(\text{Outlook} = \text{Overcast} | \text{Yes}) = $ ~~$\frac{2}{10} = 0.2$~~

$P(C_1) = 6/10$

$P(\text{Outlook} = \text{Overcast} | \text{No}) = $

$P(C_2) = 4/10$

Since there is no such sample for which outlook is overcast and the class variable is No, we add two samples such that ①

$P(\text{Outlook} = \text{Overcast} | \text{Yes}) = 3/12$.

$P(\text{Outlook} = \text{Overcast} | \text{No}) = 1/12$

$P(\text{Humidity} = \text{High} | \text{Yes}) = 2/10$

$P(\text{Humidity} = \text{High} | \text{No}) = 3/10$

$P(\text{Wind} = \text{Weak} | \text{Yes}) = 4/10$

$P(\text{Wind} = \text{Weak} | \text{No}) = 1/10$

$P(\text{Outlook} = \text{Overcast}, \text{Humidity} = \text{High}, \text{Wind} = \text{Weak} | \text{Yes})$

$$= \frac{3}{12} \times \frac{2}{10} \times \frac{4}{10} \times \frac{6}{10}$$

$P(\text{Outlook} = \text{Overcast}, \text{Humidity} = \text{High}, \text{Wind} = \text{Weak} | \text{No})$

$$= \frac{1}{12} \times \frac{3}{10} \times \frac{1}{10} \times \frac{4}{10}$$

Since, the probability of the class label being Yes is more, the class label will be predicted as Yes.

b) Since, Naive Bayes classification algorithm ~~assumes that~~ the ~~attributes of~~ the conditional independence assumption i.e the ~~occurrence~~ value of an attribute does not depend on the value of ~~another~~ any other attribute. However, this assumption is never true in practise. Hence, there can be ~~so~~ some error in such prediction.

The error can be corrected using Bayesian Belief Networks to predict the label instead.

2) If a feature has continuous values, it is either discretized using different data preprocessing techniques like binning or the Gaussian distribution is predicted using the following formula.

$$\frac{1}{\sqrt{2\pi\sigma^2}} \times e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$

where $\sigma$ is the standard deviation and $\mu$ is the mean of the distribution of the values of the feature.