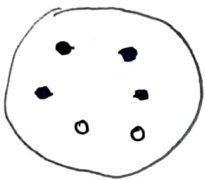2.

Cluster Validation indices.

Purity: ~~of a~~ Purity of cluster is defined as follows.

$$TT_j = \frac{1}{n_i} \max(n_{ij})$$



Purity of the given cluster

$$TT_j = \frac{1}{6} \max(3, 2, 1)$$

$$= \frac{3}{6}$$

$$= \frac{1}{2}$$

$$= 0.5$$

③

Rand Index : Rand Index can be comprehende from the
given analogy.

| | Ground Truth Value with class X | Ground Truth Value not having X |
|---|---|---|
| ~~Cluster~~ Sampling having Class label X | A | C |
| Sample not having Class label X | B | D |

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

b.

Intercluster Distance. Intercluster cluster distance gives the similarity between too different clusters.

Single Linkage intercluster distance. The minimum distance between too data points poi present in the cluster gives single linkage distance.

&. Let us say $S$ & $T$ are too different clusters. So single linkage intercluster distance

$$\Delta_1(S,T) = \min_{\substack{\forall x \in S \\ \forall y \in T}} \left( \delta(x,y) \right)$$

(4)

Complete Linkage Intercluster Distance: The maximum distance between too data points present in too different clusters.

Let us say $S$ & $T$ are too different clusters. So the distance is defined as.

$$\Delta_2(S,T) = \max_{\substack{\forall x \in S \\ \forall y \in T}} \left( \delta(x,y) \right)$$

Intra Cluster Distance. Similarity between too ~~clusters~~ data points present in same cluster.

(diameter cluster distance)

Type 1. Intra cluster distance; for a cluster $S$

$$\Delta_1(S) = \max_{\substack{\forall x \in S \\ \forall y \in S \\ x \neq y}} \left( \delta(x,y) \right)$$

Average Centroid Distance : ① Average centroid
distance is defined. as follows for a cluster $S$.

$$\Delta_1(S) = \frac{1}{\binom{|S|}{(|S|-1)}} \sum_{x \in S} \delta(x, \bar{v})$$

where $\bar{v} = \frac{1}{|S|} \sum_{y \in S} y$

Dunn's Cluster Validation Index. The dunn's cluster
validation index is defined as follows

(3)

$$D\text{index} = \min_{1 \leq i \leq c} \left\{ \min_{\substack{1 \leq j \leq c' \\ i \neq j}} \left\{ \frac{\delta(x_i, x_j)}{\max(\Delta(x_{ic}))} \right\} \right\}$$

To obtain a good cluster the dunn's clustering index
is ~~always~~ ~~maximized~~. minimized. It is always
preferred that the inter cluster distance of any
two cluster is high but the intra cluster dist-
ance is low. ~~Value~~ ~~is related~~ the inter ~~cluster~~ ~~dunn's~~ ~~index the numerator~~ ~~between two~~
~~cluster~~. ~~Dun.~~ ~~Minimum too~~ Minimization of the
~~two~~ term $\frac{\delta(x_i, x_y)}{\max(\Delta(x_{ic}))}$ gives the same.

Q4

$$P = 6, \quad n = 4, \quad I(P,n) = -\frac{6}{10}\log\frac{6}{10} - \frac{4}{10}\log\frac{4}{10}$$

$$= 0.29$$

| Outlook | $P_i$ | $n_i$ | $I(P_i, n_i)$ |
|---------|-------|-------|---------------|
| Sunny | 2 | 2 | 1 |
| Overcast | 2 | 0 | 0 |
| Rain | 2 | 2 | 1 |

$I(2,2) = \frac{4}{10} \times 0.29 = $

$= 0.049 - 0.46 = 1$

$I(2,0) = \frac{2}{10} \times 0.29 = 0$

$= 0.098$

$E(A) = \sum \frac{P_i + n_i}{P+n} I(P_i, n_i)$

$= \frac{4}{10} \times 1 + \frac{2}{10} \times 0 + \frac{4}{10} \times 1$

$= 0.8$

GAIN (OUTLOOK) = 0.29 - 0.8

$= -0.51$

$3\frac{1}{2}$

| Humidity | $P_i$ | $n_i$ | $I(P_i, n_i)$ |
|----------|-------|-------|---------------|
| Normal | 4 | 1 | 0.72 |
| High | 2 | 3 | 0.97 |

$I(4,1) = \frac{5}{10} \times 0.29$

$= 0.145 \quad 0.72$

$I(2,3) = \frac{5}{10} \times 0.29$

$= 0.145 \quad 0.97$

$E(A) = \sum \frac{P_i + n_i}{P+n} I(P_i, n_i)$

$= \frac{5}{10} \times 0.72 + \frac{5}{10} \cdot 0.97$

$= 0.845$

Gain (Humidity) $= 0.855$

| Wind | $P_i$ | $n_i$ | $I(P_i, n_i)$ |
|------|-------|-------|---------------|
| Strong | 2 | 3 | 0.97 |
| Weak | 4 | 1 | 0.72 |

$$I(2,3) = 0.97$$
$$I(4,1) = 0.72$$

$$E(A) = \sum \frac{P_i + n_i}{P+n} I(P_i, n_i)$$

$$= \frac{5}{10} \times 0.97 + \frac{5}{10} \times 0.72$$

$$= 0.845$$

Gain(wind) = −0.555

~~Outlook~~

So ~~both~~ ~~wind~~ and ~~humidity~~ can be the root.

b. Possible terminating Criteria for decision tree:

   i. The generation of decision during the training stops if the obtained gini index impurity at some node is higher than its parent.

(15)

   ii. In ID3 approach, depending on the gain value ~~the~~ generated during the training process. ✗

   iii. During testing when the process reaches the leaf node it returns the possible output and terminates.

c.

A model generally overfits when there is ~~low~~ high variance in the data points but low bias.

The possible causes

    i. There is a high variance present in the dataset.

    ii. There is low bias present in the data-set.

    iii. The model is ~~not~~ trained in such a way, i.e the weights are initialized in such a way., that there model cannot predict unseen data.

Decision tree ~~is~~ classifier is prone to overfit. ~~as~~ Solution to overfitting

    1. Use of Random forest. ~~A~~ In random forest classifier multiple decision trees are used during training. And after the classification is done the majority output is returned as predicted class.

    2. Meta Modelling: Over fitting is also dealt with a technique known as meta modelling. Here. multiple decision tree classifier trained over different instances are taken. And there outputs are in turn fed into a different classifier to predict the values.

5. a.      ~~We need to calculate~~ $P(\text{outlook}/\text{play Tennis} = \text{"yes"})$

$+ P(\text{outlook}/\text{play Tennis} = \text{"no"})$  /outlook $=$

$$P(\text{outlook} = \text{"overcast"}/\text{play Tennis} = \text{"yes"})$$
$$\underline{\phantom{xxx}} \quad P(\text{play Tennis} = \text{"yes"}) /$$
$$=$$

Let us say $X = \langle$ outlook $=$ "overcast", humidity $=$ "high", wind $=$ "weak"$\rangle$

we are to calculate

(4)  ~~$P$~~        $P(\text{play Tennis} = \text{"yes"}/X)$ and $P(\text{play Tennis}$

$= \text{"no"}/X).$

$$P(\text{play Tennis} = \text{"yes"}/X) = \frac{P(X/\text{play Tennis} = \text{"yes"}) \times P\left(\begin{array}{c}\text{play Tennis}\\=\text{"yes"}\end{array}\right)}{~~P(x)~~ \; P(X/PT = \text{"y"}) P(PT = y') + P(X/PT = \text{"n"}) P(PT = \text{"n"})}$$

$$P(\text{play Tennis} = \text{"no"}) = \frac{4}{10} \qquad P(\text{play Tennis} = \text{"yes"})$$
$$= \frac{6}{10}$$

$$P(\text{Outlook} = \text{"overcast"}/\text{play Tennis} = \text{"yes"}) = \frac{2}{2} = 1$$
$$P(\text{Outlook} = \text{"overcast"}/\text{play Tennis} = \text{"no"}) = \frac{0}{2} = 0$$
$$P(\text{~~outlook~~ humidity} = \text{"high"}/\text{play Tennis} = \text{"yes"}) = \frac{1}{4}$$
$$P(\text{humidity} = \text{"high"}/\text{play Tennis} = \text{"no"}) = \frac{3}{4}$$

$$P(winds\ "weak"/playTennis\ "yes") = \frac{4}{5}$$

$$P(wind = "weak" / play\ Tennis = "no") = \frac{1}{5}$$

$$P(X / play\ Tennis = "yes") = \frac{4}{5} \times \frac{1}{4} \times 1 = \frac{1}{5}$$

$$P(X / play\ Tennis = "no") = 0 \times \frac{1}{5} \times \frac{3}{4} = 0$$

$$P(play\ Tennis = "yes" / x) = \frac{\frac{6}{10} \times \frac{1}{5}}{\frac{6}{10} \times \frac{1}{5} + \frac{4}{10} \times 0}$$

$$= \frac{6/50}{6/50}$$

$$= 1$$

$$P(play\ Tennis = "no" / x) = \frac{0 \times \frac{4}{10}}{\frac{6}{10} \times \frac{1}{5} + \frac{4}{10} \times 0}$$

$$= 0$$

∴ The class label of the give sample would be play Tennis = "yes".

**b.** Yes, there is an error in such prediction. Here from the dataset we can see that if outlook = "overcast" then there is no such sample where playTennis = "no". So in that situation $P\left(\text{playTennis} = \text{"no"} \middle/ \text{outlook} = \text{"overcast"}\right)$ is 0. So no matter what the value the other featureset has it will outlook="overcast" will return class label no.

(4)

To deal with this problem there are several methodologies available. One of them is Laplace Correction.

In laplace correct we add a very small value to both numerator and denominator making the probability of the biased class a very small non zero value.

So let us say we have 1000 sample. And in class prediction comes up where the probability distribution become

$$P_1 = \frac{549}{1000} \quad , \quad P_2 = \frac{300}{1000} \quad P_3 = \frac{151}{1000} \quad , \quad P_4 = \frac{0}{1000}$$

After laplace correction the probability distribution becomes.

$$P_1 = \frac{550}{1004}$$

$$P_3 = \frac{152}{1004}$$

$$P_2 = \frac{301}{1004}$$

$$P_4 = \frac{1}{1004}$$

So $P_4$ will not result to zero after Bayes probability prediction, rather it will return a value close to zero.

C.  If any feature in naive bayes classification has a continuous value. Rather using the product we use the gaussian normal distribution

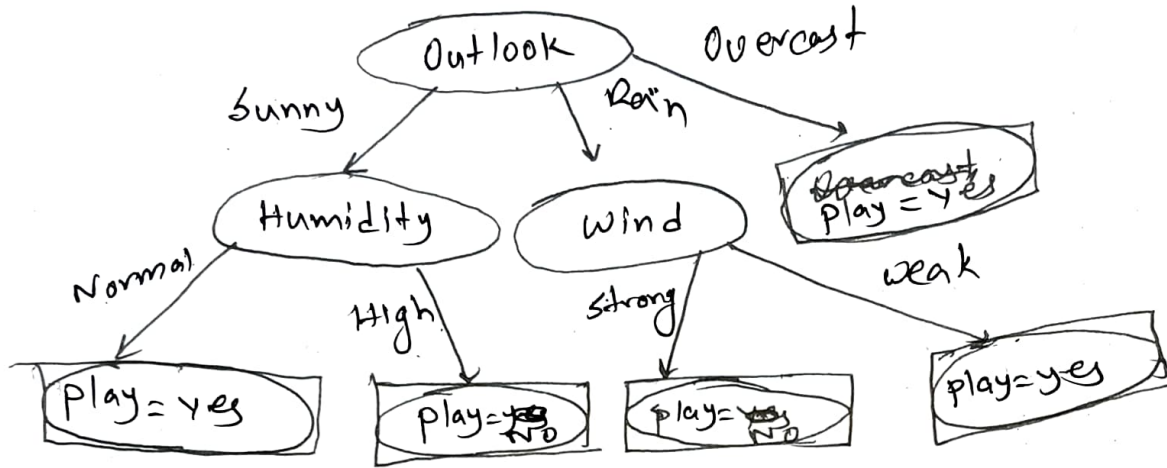$$G(x, \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{(x-\mu)^2}{\sigma}}$$

So for a sample $x_k$, $G(x_k, \mu_k, \sigma_k) = \frac{1 \times e^{\frac{(x-\mu_k)^2}{\sigma_k}}}{\sigma\sqrt{2\pi}}$

③

So $P(X/Y) = \frac{P(Y(x) \times P(x)}{\sum P(x)}$

So in this case first let us consider the $x_k$ has a continuous probal distribution then we use instead of the probability normal value we tot use the gaussian value.

6.a)



Decision tree:
- Outlook
  - Sunny → Humidity
    - Normal → play = Yes
    - High → play = No
  - Rain → Wind
    - Strong → play = No
    - Weak → play = yes
  - Overcast → play = yes

| Outlook | Humidity | wind | Play Tennis CV = yes | play Tennis predisted class |
|---|---|---|---|---|
| Sunny | Normal | Strong | Yes | Yes |
| Rain Overcast | Normal | Strong | No | Yes |
| Rain | High | Strong | Yes | No |
| Sunny | High | Weak | No | No |
| Rain | High | Strong | No | No |

Confusion Matrix.

|  |  | Predicted Class Yes | No |
|---|---|---|---|
| Trues Class | Yes | True positive 1 | False Negative 1 |
|  | No | False positive 1 | False Negative True 2 |

$$\text{Precision} = \frac{TP}{TP + FP}$$

(4)

$$= \frac{1}{1+1} = 0.5 ✓$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{1}{2} = 0.5$$

~~F score =~~

$$F_\beta \text{ score} = \frac{1+\beta^2}{\beta^2} \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

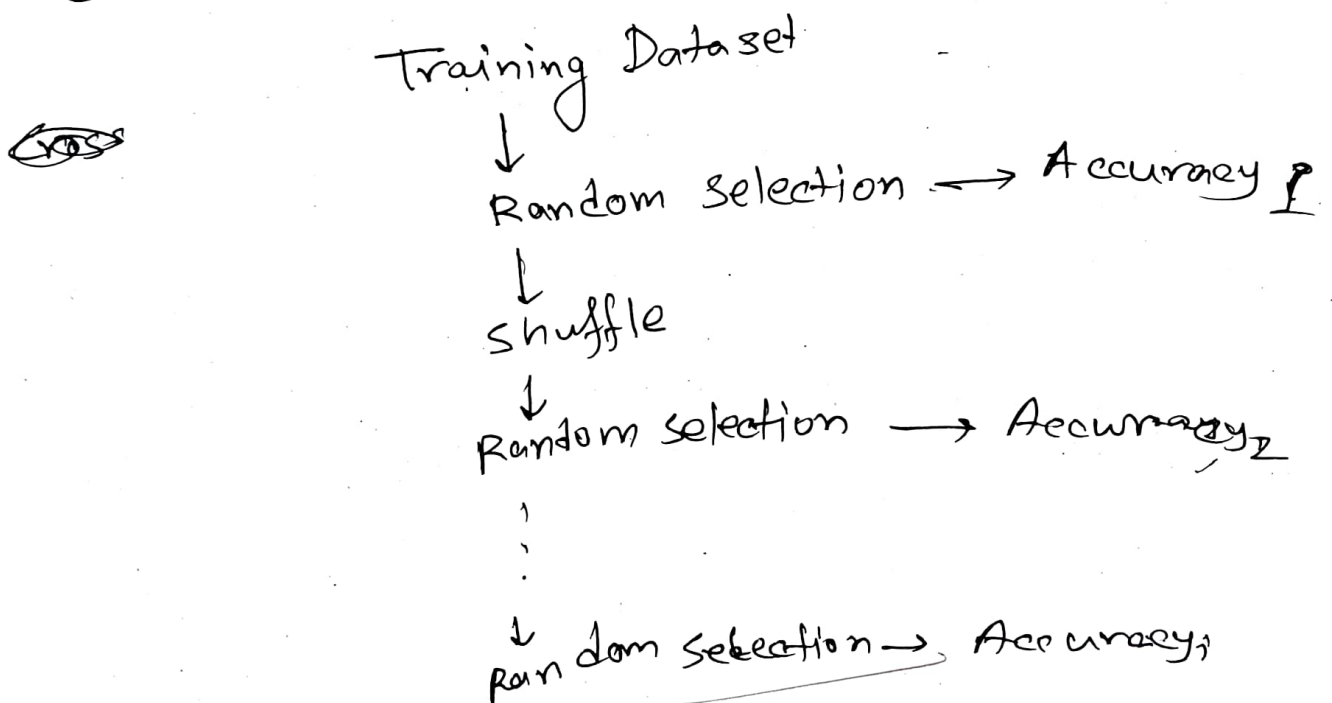$$F_1 \text{ score} = \frac{2 \times 0.5 \times 0.5}{0.5 + 0.5}$$

$$= 0.5$$

b)

9.

Hold out Method: In holdout method from the dataset for a percentage of the dataset is taken for training a the rest of the dataset is used for validation. This select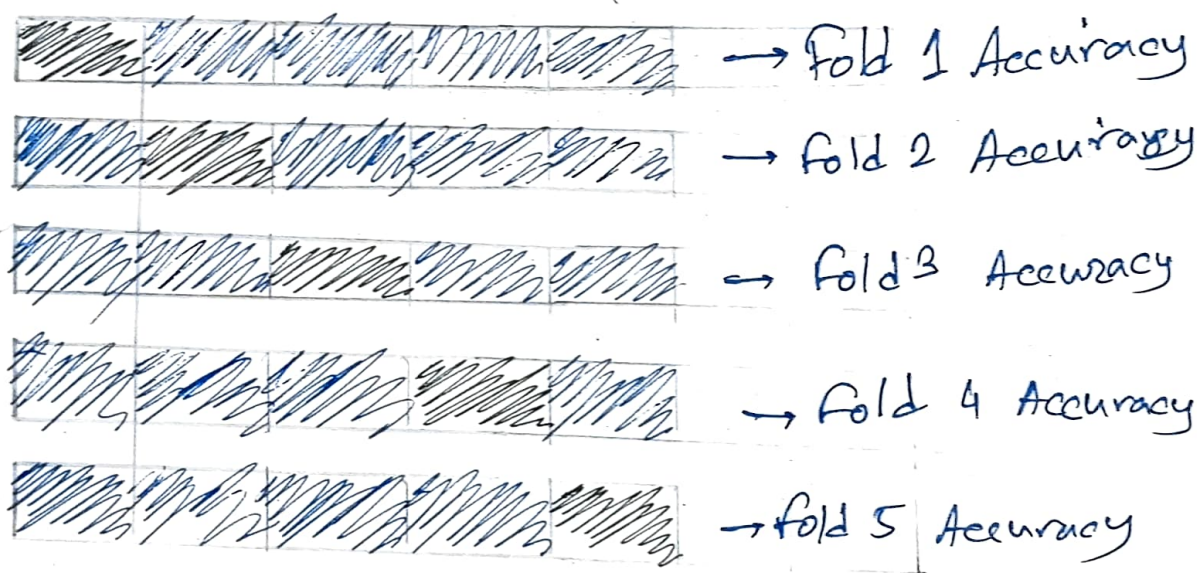ion is completely random and repeated for a number of times. And normally the mean of the all accuracies are reported with an error approximation, i.e the standard deviation.

Training Dataset
↓
Random selection → Accuracy $1$
↓
shuffle
↓
Random selection → Accuracy $2$
|
⋮
↓
Random selection → Accuracy $i$

$$Accuracy = mean (Accuracy_1, Accuracy_2 \ldots Accuracy_i)$$
$$+ std. deviation (Accuracy_1, Accuracy_2$$
$$\ldots Accuracy_i)$$

Cross validation: There are mainly two types of cross validation present. i. K fold cross validation and startified cross validation.

In cross validation the dataset is divided into k folds. And the metrics is calculated K times. First Each of the fold is used as a validator for once, and rest of the folds are used for training

| | |
|---|---|
| ▨▨▨▨▨▨ | → fold 1 Accuracy |
| ▨▨▨▨▨▨ | → fold 2 Accuracy |
| ▨▨▨▨▨▨ | → fold 3 Accuracy |
| ▨▨▨▨▨▨ | → fold 4 Accuracy |
| ▨▨▨▨▨▨ | → fold 5 Accuracy |

$$Accuracy = mean(Folds\ Accuracy) + std\_deviation(Accuracy)$$

Bootstrap Method: In the bootstrap method works by the priniciple "resampling by replacement". Suppose there are $d$ instances present in in a bootstrap. One percentage of the bootstrap is used for training and the rest of the portion which couldnot make it to the boot strap is used for testing. The normal

(3) bootstrap is $0.368$, as $1 - \frac{1}{d} \approx e^{-1} \approx 0.368$

The accuracy is calculated as follows.

$$Acc = (0.368 \, Acc_{train\_set} + 0.632 \times Acc_{test\_set.})$$

---

c.    Main Purpose of Ensemble classifier.

, i. Ensemble classifier Use of ensemble classifier is fruitfull in various ways. (For example a hypothesis is trained over a training set and the accuracy is calculated. But there is a possibillfy that the training set has some unseen instances. So useing a ensemble

classifier outputs to a set of accuracy. and predicted class lables. Depending on that weighted etoo majority the correct class lablel is choosen from a set of class tabels.

Adaboostin: Adaboost is boosting algorithm used in ensemble classifiers. In boosting Adaboost several instances of classtfiers are taken. And then depending on that instan And if a new instance is introduced that instance is give priorily over all. And during traning there exist some particular threshold. If the new value instance works worse compared to that threshold the instance is discarded. the threshold is determined using null hypothesis or t testing

$$t = \frac{err(M_1) - err(M_2)}{\sqrt{var(M_1) - var(M_2)}}$$