**1** Explain how information gain measure is biased towards attributes having large number of values.

Information Gain = $Info(D) - Info_A(D)$  where $A \rightarrow$ attribute chosen

⊗ And it tends toward purity (one class belonging to just yes or no, or maximise towards it)

Now, if our attribute is Unique_ID, it will be diff for every tuple and hence be the best attribute to choose only for decision tree, but in actual it won't give any information on it

Hence, we use information-split $= -\sum \frac{|D_j|}{|D|} \log_2 \frac{|D_j|}{|D|}$

and then Gain Ratio $= \frac{Info Gain}{Info.split}$ to reduce this biaseness towards multivalued attributes. Now, unique_ID will have large no. of splits, hence Info.split huge, $\rightarrow$ Gain Ratio small.

**2** Explain the term "underfitting". Discuss the causes of underfitting and the methods to avoid it.

Underfitting refers to model training when the model is not able to understand and capture the underlying trends & patterns in the dataset. (Low Bias, Perform poor on both training & test Dataset. High Variance)

Causes of Underfitting:
→ Imbalanced Dataset
→ Data leakage ?
→ ~~test~~ complex attributes

To avoid it we can
→ Generalisation, ~~Generalise attributes~~
→ Normalization
→ ~~choose~~ attributes that are actually useful eg. Covi from Patient History dataset instead of it's SSR No.

1

**3** Explain the basic idea of Rainforest (a scalable decision tree classifier).

Decision Tree Classifier ° Creates the decision tree, iterative by greedly with divide-and-Conquer approach.

Rainforest ~~There are certain isu~~ even though decision tree is highly explainable & inherently interpretable, it has certain issues like → Accuracy → Multi-valued attributes, etc

Hence, we scale our ~~①~~ decision tree with diff algo's to reduce this.

Rainforest:                    X    0

Decision tree
↳ MDL
↳ RSS
↳ CHAID

---

**4**

$$P(C \mid A_1 A_2 \ldots A_n) = \frac{P(A_1 A_2 \ldots A_n \mid C) P(C)}{P(A_1 A_2 \ldots A_n)}$$

Interpret the above formula with a real life example.

$$\underbrace{P(C \mid A_1 A_2 \ldots A_n)}_{\text{Posterior Prob}} = \frac{\overbrace{P(A_1 A_2 \ldots A_n \mid C)}^{\text{likelihood}} \; \overbrace{P(C)}^{\text{Prior Prob}}}{\underbrace{P(A_1 A_2 \ldots A_n)}_{\text{evidence}}} \rightarrow \text{Naive Bayes Classifier}$$

Example, fruits can have multiple features like color, taste, cost, etc.

$$P(Apple \mid Color, taste, Cost) = \frac{P(Color, taste, Cost \mid Apple) \; P(Apple)}{P(Color, taste, Cost)}$$

It Calculates the Cond$^n$ Probability using Bayes Classifier.
If it's naive Bayes Classifier and all attributes ~~apple~~ Color, taste, Cost are Independent of each other →

$$P(Apple \mid Color, taste, Cost) = \frac{P(Color \mid Apple) \cdot P(taste \mid Apple) \cdot P(Cost \mid A) \cdot P(A)}{P(Color) \cdot P(taste) \cdot P(Cost)}$$

Hence, it can be read as : Prob of Apple given we know the values of Color, taste & Cost, is equal to Prob. Color, taste, and Cost given we know the fruit is apple, multiplied by prob; of fruit being apple in any case, whole divided by Prob. of any fruit being of this color, taste, and cost.

② Nice!!

**5** Explain the purpose of Laplacian correction in Naïve Bayesian classifier.

In Naïve Bayesian Classifier, as shown in Ques 4 we assume all attributes are Independent of each other.

$$P(c | A_1, A_2 \ldots A_n) = \frac{P(A_1|c_i) \cdot P(A_2|c_i) \ldots P(A_n|c_i) \cdot P(c)}{P(A_1) \cdot P(A_2) \ldots P(A_n)}$$

②

But if any one $P(A_i|c_i) = 0$ then whole Conditional Prob. will become zero.

Hence, in huge datasets, we assume that one row exist for each case (adding one row in huge datasets don't affect that Much).

Hence, ~~Laplace Corr~~ if $P(A_i|c_i) = \frac{N_{ic}}{N_f} = 0$ → No. of elements in $A_i$, $N_f$ → No of total elements

After laplace Corr → $P(A_i|c_i) = \frac{N_i + 1}{N + n}$, assuming n attribute

Now even if $N_i = 0 \Rightarrow P(A_i|c_i) \neq 0$ leading us to prevent the posterior prob. connecting this.

**6** Explain the classification metrics, which will help you to understand the extent of class imbalance aspect of the classifier.

Class Imbalance: ① when a particular class has labels for one side comparatively very huge in number compared to other sides.

Eg. Cancer dataset Containing 97% No and 3% Yes cases. So, even if accuracy = 97% doesn't means its classifying no cases properly.

② Metrics Used: 1) Precision = How many positives are actually true positive

|   | Y | N |
|---|---|---|
| Y | TP | FN |
| N | FP | TN |

$$Precision = \frac{TP}{TP + FN}$$

2) Recall = % of tuples of +ve that are correctly classified as +ve

$$Recall = \frac{TP}{TP + FP}$$

① Specificity ?

For Correction this class Imbalance, we use LIMO, SHAP, etc.

Stream Classifier: Stream of data comes continuously but we have a fixed storage space. It can be TV, running radio, etc.

Issues: • Fixed Storage Space
• Need to process a whole batch at Once
• ~~Context from previous batches is broken~~
  ~~as prev. batches are removed from~~

Ensemble: Create Multiple Models and then ~~dispose~~ predict the Best Model.
either ~~sequentially~~ parallely (Bagging) or ~~Parallely~~ sequential (Boosting)

~~Here~~ Here, what we do is → (Sequential)
1) Let the stream Input come in.
~~2) Then we find the Model Best to this Input.~~
~~3) And remove the model performing worse from the storage~~

2) We update the Model with this Input if it
   • increases the goodness of metrics considered.
   └ else the Input's discarded.

3) The parts in storage performing worse in metric compared to this Input are removed, if this one is added.

4) And hence, the model is trained.

~~This is~~

dividing into multiple chunks
chunks    classifies

2) Then we find the model performing best ~~according to~~ ~~goodness of Metrics to~~ to this Input.

3) ~~And remove~~ If the goodness of Metrics considered for this Model worse than the existing models → Discard this Model.
→ else update this Model in storage and remove the one performing worse from the existing Model.

4) Repeat this until the input comes.

Ensemble 5) Then you will have the best j models at the end, use majority-vote, etc. for choosing the best Model.

intuitively correct

1.5

1.5

Multiclass Classification: Instead of Binary (0|1) classifier, we have more class labels.
eg. One +ve, and the rest all -ve's, etc.

# We perform this classification through (K-NN) (K nearest Neighbours), then it will use euclidean distance. — Hamming distance

Let $110000$ & $000011$ be two ~~values~~ class labels

Dist. from $111111$ same for both

$$\sqrt{\sum (x_i - x_f)^2}$$

But ~~even~~ Numbers completely different.

Hence, we use error correcting codes (for ex. Hamming codes) to correct this.

eg. →

Let 4 attributes be A, b, C, D $<$ $\{1,1,1,1\}$  $\{1,3,2,4\}$

Case I → $P(a) = P(b) = P(c) = P(d) = \frac{1}{2}$  → uni-class

i) Prob $= 4 \times \frac{1}{2} = 2$

ii) Prob $= \frac{1}{2}[1 + 3 + 2 + 4] = 5$

Case 2 → $P(a) = 1$, $P(b) = 3$, $P(c) = 1$, $P(d) = 2$  → multi-class

i) Prob $= $ ~~log~~ $1 + 3 + 1 + 2 = 7$

ii) Prob $= -\log_2 \frac{1}{4} \neq -\log_2 \frac{3}{4} - \log_2 \frac{2}{4} - 2\log_2 \frac{4}{4}$

~~Consistent finalizing~~

So, we use error-correcting codes (for ex → Hamming codes) to correct this.