Name  Siddhant Shah

Subject  DMML

Course & Year  BSc Year 3

Roll Number  BM(202171)

Date  2/3/24

Total No of Pages  4

—·— **Begin here**

※ Question 1 —

No of transaction $\ge 10^{10}$

Size of each transaction $\le 10$

no of tokens: $10^{11}$

Each frequent token occurs in at least $10^{-4} \times 10^{10}$ transaction

$= 10^6$

Maximum such tokens $= 10 \times \dfrac{10^7}{10^6} = 100$

⎰ Each transaction has at most 10 tokens

∴ $n(F_1) \le 100$

Now say all combination $F_2$ is made of union of two elements

Now $n(F_2) \le n(F_1)$  of $F_1$, such that the union is also frequent.

$\Rightarrow \{x,y\} \in F_2 \iff freq(\{x,y\}) \ge 10^6$

and $x, y \in F_1$

Maximum such pairs $= \dfrac{1}{2} \times 10 \times \left(\dfrac{10^7}{10^6}\right) = 50$

↑ length of sequence

$n(F_2) \le 50$

* __Question 3:__ —

The claim is not justified as it doesn't take into account the average value of the attribute.

For example, say ~~$x_1$ pts of education~~ $x_1$ measures distance in cm and $\theta_1 = 3$

and $x_2$ measures length in inches and $\theta_2 = 2$

Taking these values directly will not ~~be~~ give an accurate answer as they have to be normalized.

Another example where this would fail is when we consider $x_1$ that varies between 1 and 10 ~~with~~ with an average value of 5, $\theta_1 = 0.1$

and $x_2$ that varies between 1 and 100 with an average value of 60, $\theta_2 = 0.01$

If we simply compare $\theta_1$ and $\theta_2$, then we get that the first attribute has ~~a~~ a larger contribution. However the average contribution is $\theta_1 \text{Avg}(x_1) = 0.5$ and $\theta_2 \text{Avg}(x_2) = 0.6$ respectively, which contradicts our earlier claim.

(5)

* Question 5 :—

The $i^{th}$ column of $D$, $D^i$ ~~has~~ contains the distances of all the points from $x_i$

$$D^i = \begin{bmatrix} d(x_1, x_i) \\ d(x_2, x_i) \\ \vdots \\ d(x_n, x_i) \end{bmatrix}$$

Running clustering on the columns of $D$ will give us points that are at a similar distance from each of the points in the ~~plan~~ plane.

For example, if $D^i$, $D^j$ are in the same cluster, ~~it with~~ the ~~max do~~ radius $\epsilon$, we get that.

$$\sqrt{(d(x_1, x_j) - d(x_1, x_i))^2 + \ldots + (d(x_n, x_j) - d(x_n, x_i))^2}$$
$$\leq \epsilon$$
$$\implies |d(x_k, x_j) - d(x_k, x_i)| \leq \epsilon \quad \forall \ k = 1, \ldots n$$

Let $k = i$ or $k = j$ to get
$$|d(x_i, x_j)| = d(x_i, x_j) \leq \epsilon$$

This can be done for ~~points~~ points $x_i, x_j$ such that $D^i, D^j$ are in the same cluster.

⑤

* Question 6 :—

Suppose we have $n$ points. We express each point $x_i$ as a linear sum of its $k$ nearest neighbours ($k$ is predecided)

$$x_i = \sum_{j=1}^{n} w_{ij} x_j$$

$w_{ij} = 0$ if $x_j$ is not one of the $k$-nearest neighbours of $x_i$ ~~($w_{ij} = 0$)~~ ($w_{ii} = 0$)

This gives us a matrix $W$.

We have to find $w$ so that the squared distance is minimized

$$\hat{W} = \underset{W}{\text{argmin}} \sum_{i=1}^{n} \left( x_i - \sum_{j=1}^{n} w_{ij} x_j \right)^2$$

This gives us the locally linear embedding.

→ $w_{ij}$ that are such that $x_j$ is a $k$-nearest neighbour are learnable parameters. for all $i = 1, \ldots n$

---

Weights → Embedding

(3)

* <u>Question 2</u> : -

Given a decision tree.

→ For categorical variables, add the values that are permitted for traversal along that path to to a list

Each categorical variable is queried only once

→ For numerical variable, do the same thing but create a range

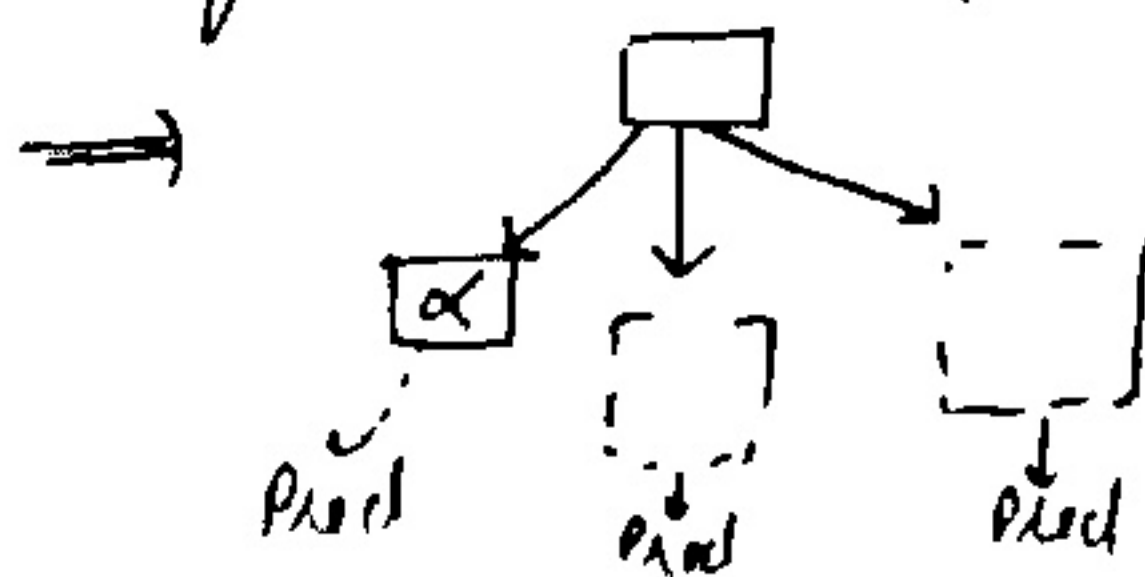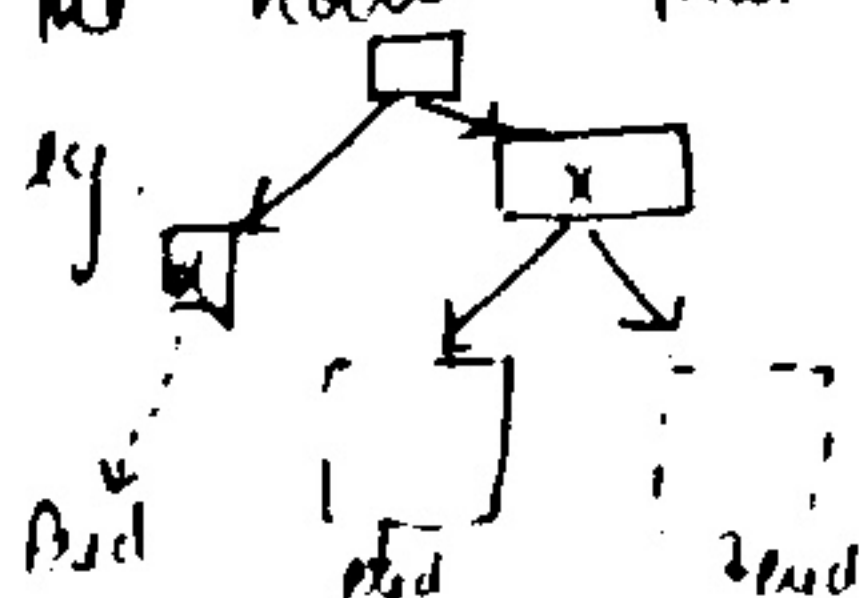If all values $< 10$ are permitted, add $[-\infty, 10)$ to the list

If the ~~value~~ variable is queried again, take the intersection of the previous range for that variable, with the current proposed range.

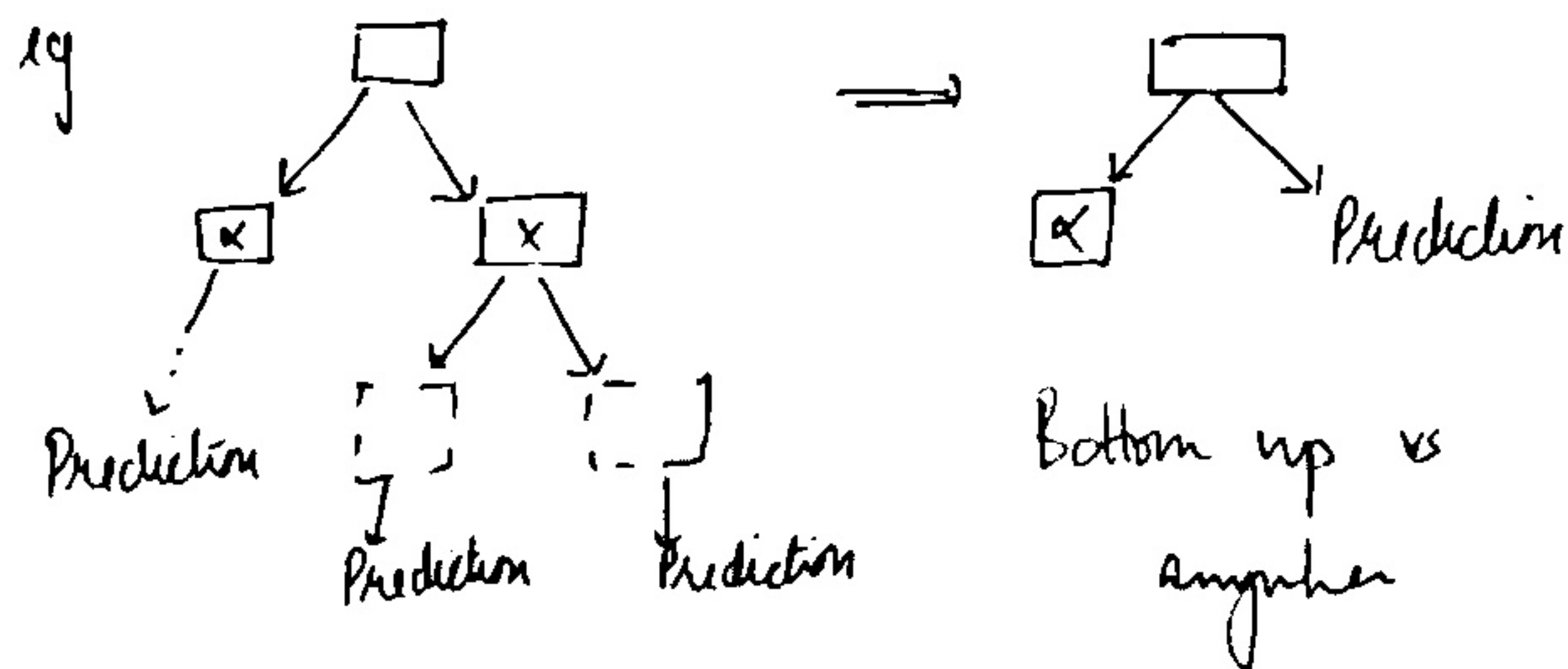→ At the end of the path, we get the prediction of the tree. Add that to the left side of the association rule

Thus, we will be left with rules that look like

$$\{ (\text{categ var } i_1 = \text{value } i_1), (\text{categ var } i_1 = \text{value } i_2, \ldots \ldots, \text{num var } d_1 \in (l_1, m)), \ldots \ldots \} \longrightarrow \text{Prediction}$$

Removing an attribute from the left side would be to not query that attribute at all in the tree but to keep the part of the tree that was under the nodes that ~~was~~ queried the respective attribute.

eg:

The usual methods of pruning either limit the max depth or place a restriction on splitting a node. When we remove a node, we merge all the subnodes into the parent nodes.

eg



Bottom up vs anymore

(4)

---

* **Question 4 :-**

In a ~~non~~ random forest built on the same attributes as those we wish to rank,

For each attribute, we take some metrics to quantify the central tendency of the impurity of any node that queries on that attribute (Mean, max, min)

weighted impurity gain

We rank the attributes by the <u>mean impurity</u>. in decreasing order, as an attribute with lower impurity classifies the data better.

The random forest is better than a single decision tree as the decision tree will only be formed by maximizing the information gain at

each step but won't be able to get a globally min
config for the data due to its high variance.
Averaging over the random classifier increases the
accuracy # and decreases variance.

eg It is possible that a single decision tree won't
create a node that has 50% impurity as there
may be another attribute with lesser impurity.
But it may so happen that after creating ~~that~~ the
earlier node, ~~node~~ adding one more node reduces
the impurity drastically.
This ~~is that~~ tree is likely to be formed in the
random forrest and it should be taken into
account in ranking the attributes.

$$\boxed{3.5}$$

---

* **Question 1 :—**

Each $1^{st}$ item has to occurs in at least $10^{10} \times 10^{-3}$
$$= 10^7 \text{ transactions}$$

Consider that each frequent item occurs exactly $10^7$ times
No. of frequent items $= \left(\dfrac{10^{10}}{10^7}\right) \times 10 = 10^4$

$\therefore n(F_1) \leq 10^4$ ✓

Now for $F_2$, ~~if each item in a transaction~~
if the same transaction is repeated $10^7$ times,
every pair of items in it will be in $F_2$

There can be such $10^3$ ~~paint~~ transactions.

Each transaction gives us such $\binom{10}{2} = 45$ pairs

$\therefore$ Total pairs $= 45 \times 10^3 \geqslant n(F_2)$

No. of elements used $= 10 \times 10^3 = 10^4 \leqslant 10^7$ ~~but this is~~

~~not possible because~~ ~~$10^5 \geqslant 10^4 \geqslant n(F_1)$~~

$\Rightarrow n(F_2) < 4.5 \times 10^4$

$n(F_1) \leqslant 10^4$

$n(F_2) \leqslant 4.5 \times 10^4$ ✓ ⑤