**Write the differences between ID3, CART, C4.5**

Ans:

| Name | Spliting Type | Attribute Type | Missing Values | Pruning Strategy |
|------|---------------|----------------|----------------|------------------|
| ID3 | Information Gain | only Categorical Value | Don't handle missing values | No Pruning |
| CART | Delta Gini | Categorical & Numeric Value. | Handle missing Values | Cost Complexity Pruning is used |
| C4.5 | Gain Ratio | Categorical & Numeric Value | Handle missing values | Error based Pruning. |

☐ Gini ; Binary value নিয়ে কাজ করে,

☐ Gini কম্পিউটেশন দ্রুত করা হয়। জন্য করবে।

**What is tree Pruning.**

Ans: Pruning means removing specific branches to benefit the whole tree. Pruning tree methods address the Problem of overfitting the data.

☐ Training data তে noise থাকলে decision tree তে anomalies দেখা যায়,

☐ Tree Pruning করলে, Performance বাড়ে,

☐ Tree Pruning is being done to increase the efficiency.

☐ দুই ধরনের Pruning -হয়, a. Pre Pruning b. Post Pruning

# Define Pre-Pruning

**Ans:**
A tree is Pruned by halting its Construction early. Upon halting the node becomes a leaf.

# Define Post-Pruning

**Ans:**
Removing sub-trees from a full grown tree.

=> Pre-Pruning আগে tree বানাচ্ছি, তেমন একটা node এর নিচে আর নামছি না, leaf হিসেবে রাখছি।

=> Pre-Pruning এ আগে Construction আরম্ভ করা হয়।

=> Pre-Pruning এ leaf ; most frequent class কে ধরে।

=> Pre-Pruning এ Probability distribution হিসেব করা হয়।

=> Pre-Pruning এর কোনো threshold maintain করা হয়ে থাকে।

# Threshold যদি high হয়, তারপর tree তে এর leaf হয়ে যাবে, তখন tree টি over simplified হয়ে যাবে,

# Threshold যদি low হয়, তারপর tree তে খুব সামান্য পরিমান পরিবর্তন আসবে।

※ Post Pruning -এর বেলায় Pre Pruning থেকও বেশি Computation করতে হয়, কিন্তু Post Pruning most reliable.

## ※ Define Repetition

Ans. It occurs when an attribute is repeatedly tested along a given branch of the tree.

## ※ Define Replication

Ans. Duplicates Subtrees exist within the tree.

※ Pruning -এর নিচু টেকনিক আছে, এরূপ তিন টেকনিক,

    a. Pruning by Cost Complexity

    b. Pruning set

    c. Pessimistic Pruning

---

※ Pruning by Cost Complexity → একটা Post Pruning approach

※ Pruning set →

## Pruning By Cost Complexity

- ✷ CART [ Classification & Regression Trees] / Delta Gini Cost Complexity approach ব্যবহার করে।
- ✷ Cost Complexity একটি ক্রুত্ব সম্পর্ক function.
- ✷ $f(x)$ = Number of decision tree in the leaf node.
- ✷ $f(x)$ = Number of leaf node in decision tree + Error rate.
- ✷ Error rate বলতে ভুলভাবে miss classify করা।
- ✷ Error rate should be equal on less.
- ✷ Cost Complexity এর মধ্যে ক্রুত্ব জিনিস Compare করা হয়, যেগুলো হল
  - a. Present Error Rate
  - b. Leaf node এর সংখ্যা,
- ✷ ছোট tree Pruning By Cost Complexity এর জন্য better
- ✷ Define Error Rate

__Ans.__ The Percentage of the instances misclassified by the tree.

## Pruning Set:

৪ Labled instance এর পরিবর্তে Cost Complexity estimate করা হয়।

৪ এই algorithm set of Pruned trees generate করবে।

৪ $D_L = \{DT_1, DT_2, DT_3, \dots, DT_N\}$ ← Set of Pruned trees.

এই set ৬ minimum error rate, minimum number of leaf এবং Maximum accuracy থাকে।

## Pessimistic Pruning:

৪ এই যুক্ত C4.5 algorithm ৬ use করা হয়।

৪ এখানে বলা যায়

$f(x)$ = Number of leaf node in the decision tree + Error rate চিহ্ন Recall দূরকরণ হয়।

৪ Pessimistic Pruning ; Pruned set ব্যবহার করেনা, এই training set ব্যবহার করে, কারন training set unique.
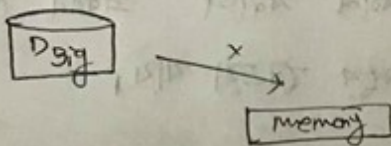
৪ Define Recall

Ans.  Recall means true Positive rate

**⊞ Write the importance of Scalability & Decision Tree Induction**

Ans.

Most often, training data would not fit in the memory. $ID_3$, $C4.5$, CART does well for small data sets.

For this reason tree construction becomes inefficient. As swapping the ~~the~~ the training instances in & out of main & cache memory.



Big data একবারে memory তে আসবে না যায় না। তাই decision tree induction করতে হয়, তখন Scalable approach হিসেবে BOAT এবং Rainforest ব্যবহার করা হয়।

**Rain Forest:**
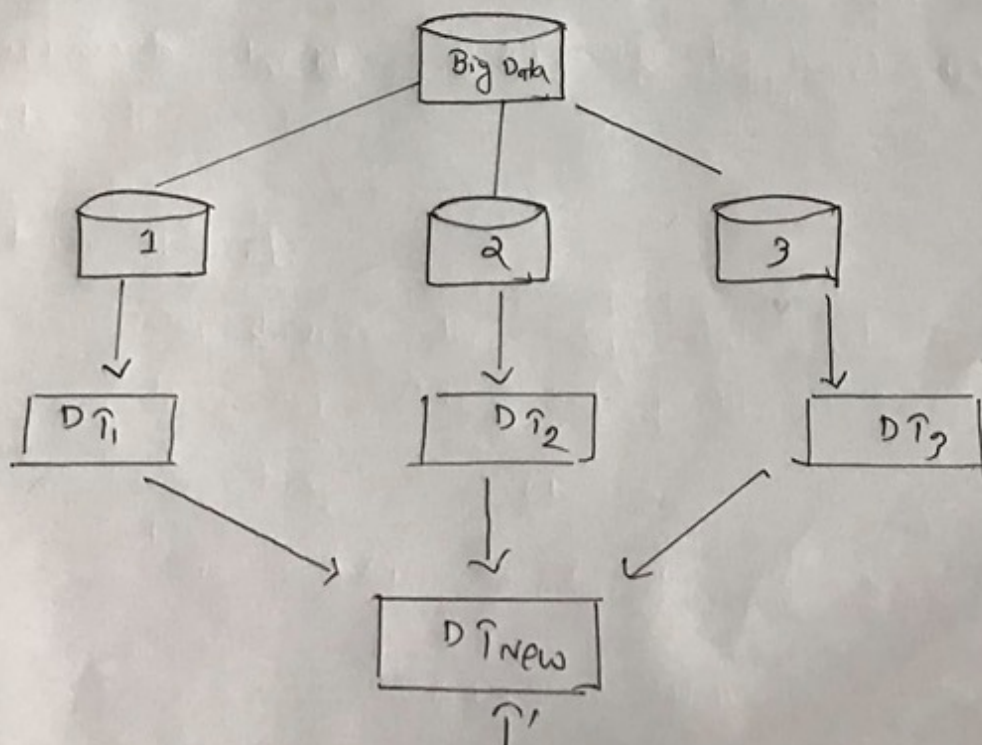
⊞ Divide the big data into small chunks.
⊞ এর AVC-Set [Attribute Value Class label] Maintain করে,

$$D_{Big} = \{x_1, x_2, x_3, x_4, \dots x_n\}$$
$$\downarrow$$
এখানে $x$ গুলো এক chunk/subdata/
Portion of the data.

| Attribute | Class Attribute | | |
|---|---|---|---|
| | Yes | No | |
| Sunny | 2 | 3 | = 5 |
| Overcast | 4 | 0 | = 4 |
| Rain | 3 | 2 | = 5 |
| Total | 9 | 5 | |

$$9 + 5 = 14$$

=> Avc - set এ — মূলত Attribute এর তেলুগুলোর class number - আছে, ফলে সহজে এবার chunk এর Prior Probability বের করে নেওয়া যায়,

=> এই Process এ—

$$D_{Big} = \{ x_1, x_2, x_3, \cdots, x_n \}$$

বারবার update হবতে হয়, তাই Rain-forest অনেক slow Process,

## BOAT: [Bootstrapped Optimistic Algorithm For Tree Construction]

✦ BOAT এ মূলত Big data ছোট Chunk এ বিভক্ত হয়। যা DT₁, DT₂, DT₃, ... ,DTₙ. এগুলো মূলত AVC-set-

✦ এই AVC-set থেকে merge করে decision tree বানানো হয়।

✦ Tree বানানো হয়ে গেলে AVC-set flash করা হয়,

✦ BOAT একটা দুর্দান্ত Incremental Updates এর জন্য ভালো,

✦ BOAT Training data তে insertion এবং update করে, এজন্য এই tree reconstruct করে না,

✦ Gini index ; Boat use করে।

Boal এর figure →

☐ Write the advantage of BOAT over Ram forest.

Ans:

a. New data input.

b. Update the tree

c. No need to start from scratch.

d. Incremental update.

e. It guarantees accuracy.