

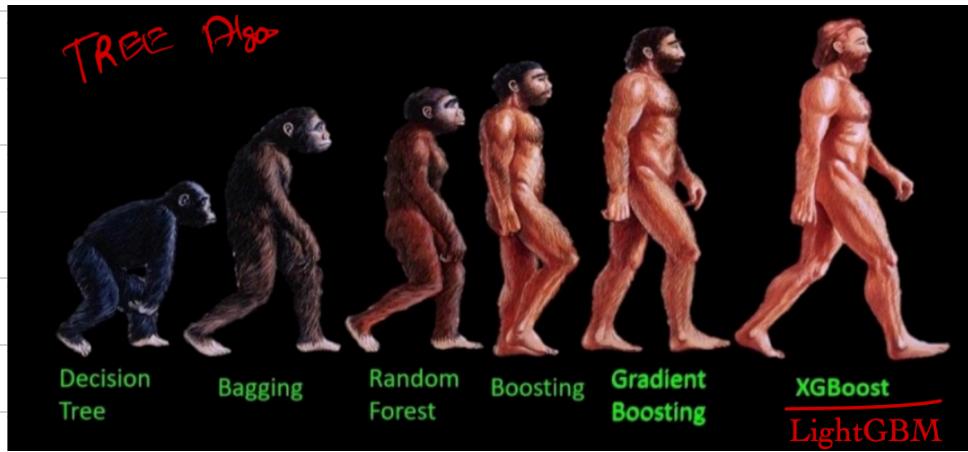
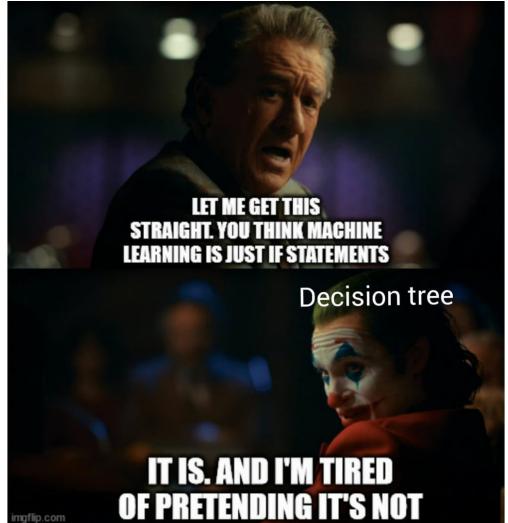
Agenda

- ① Recap for K-Fold cross validate
- ② Algorithm for Decision trees
- ③ Gini Impurity
- ④ How to split on numerical Features
- ⑤ Overfit Vs Underfit
- ⑥ Hyper-Parameter Tuning
- ⑦ Feature Importance
- ⑧ Regression using DT

Algorithm for Decision Trees

What an Amoeba can do ?

- sense many cues (food chemicals, toxins, temperature, stiffness),
- integrate those signals,
- change strategy (e.g., wander randomly until it detects a gradient, then move up it),
- keep short-lived “memories” in its biochemistry (previous exposure changes later responses),
- and navigate messy environments by reorganizing its cytoskeleton on the fly.



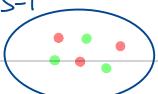
Gini Impurity

$$\text{Entropy} \rightarrow - \sum_{i=1}^K p(y_i) \times \log(p(y_i))$$

$$GI(Y) = 1 - \sum_{i=1}^K p(y_i)^2$$

GI is more sensitive to heterogeneity

S-1

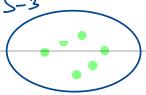


$$P(R) = 3/6$$

$$P(G) = 3/6$$

$$= 1 - ((3/6)^2 + (3/6)^2) = 0.5$$

S-3



$$P(R) = 0$$

$$P(G) = 1$$

$$= 1 - (0^2 + 1^2) = 0$$

A decision tree model uses Gini impurity as the splitting criterion.
If a node has 60 instances of class A and 40 instances of class B,
what is the Gini impurity at that node?

0 users have participated	
A	0.24
B	0.4
C	0.6
D	0.48

End Quiz Now

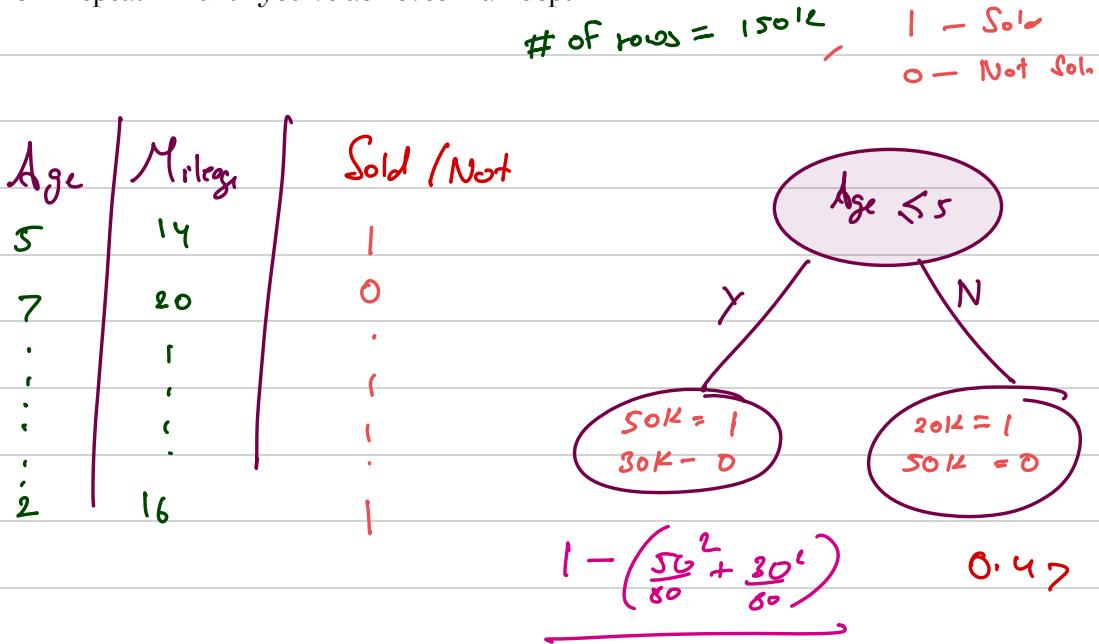
$$(1 - (0.6^2 + 0.4^2)) = (1 - (0.36 + 0.16))$$

Based on all quizzes from the session

 Harshitha Chowd... 1/1 87.43	 Rakesh Karade 1/1 92.20	 Purushottam Ku... 1/1 86.50
4 Deependu Ghosh 1/1 86.07	5 N Nayana 1/1 85.17	6 Aditya Shandilya 1/1 83.67
7 Praveen 1/1 79.23	8 MADEEHA REHMAN 1/1 78.00	9 Snehal Adhikary 1/1 75.57
10 SHASHANK JHA 1/1 73.66		

Algorithm of Decision Trees

1. For all the features in the dataset, and every single value in each of the features, compute information gain.
2. Create a split based on the feature/value pair, which has the highest information gain.
3. Repeat 1-2 until you've achieved max-depth.



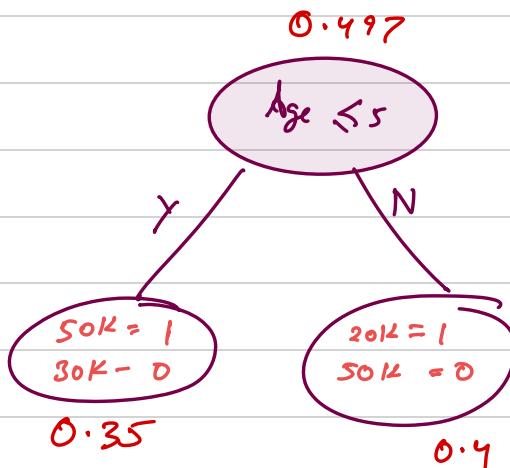
Root node

F_i, C_i — Feature_Value
Path with
highest I.G.

$$\# \text{rows} = 150K$$

$$O - 80K$$

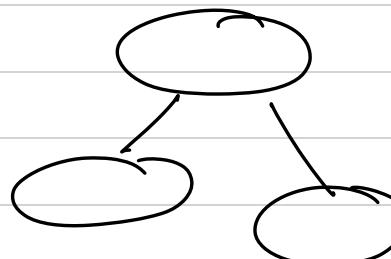
$$I - 70K$$



$$0.35 \times \frac{80}{150} + 0.4 \times \frac{70}{150} = 0.373$$

$$IG = 0.497 - 0.373 = 0.124$$

	f_1	f_2	f_3	f_4	f_5	f_6	y
n_1	Blue	Blue	Blue	Blue	Blue	Blue	Blue
n_2	Blue	Blue	Blue	Blue	Blue	Blue	Red
n_3	Red	Red	Red	Red	Red	Red	Red
n_4	Blue	Blue	Blue	Blue	Blue	Blue	Blue
n_5	Red	Red	Red	Red	Red	Red	Red
n_6	Blue	Blue	Blue	Blue	Blue	Blue	Red
n_7	Blue	Blue	Blue	Blue	Blue	Blue	Red
n_8	Red	Red	Red	Red	Red	Red	Red



For all the blue-rows, I'll use all the features, and values to find out which feature, value pair has highest information gain, and use that for splitting.

Same for the other leaf node, i.e. the red lines.



Brute Force Approach:

1. For each numerical feature:
 - i. Sort numerical feature in ascending ordering.
 - ii. With each unique value in the sorted list, compute IG (Information Gain), store it.
2. Do 1. For all the features
3. Compute Argmax, and find feature, value pair, and do a split, keep doing this until you reach max-depth.

$F_i \rightarrow 0 - 100$ (floating point)

$[0 - 10]$ 0!

$10'$

$[10 - 20]$ 20!

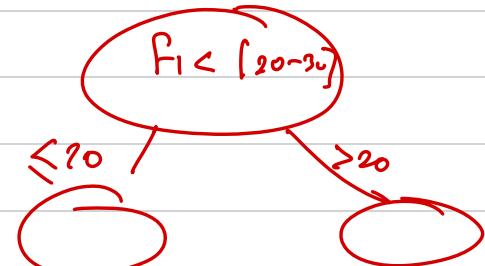
$20'$

$[20 - 30]$ 30!

$30'$

\vdots

$[90 - 100]$



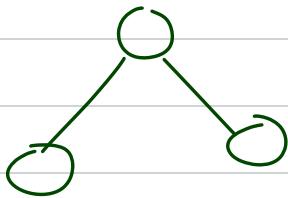
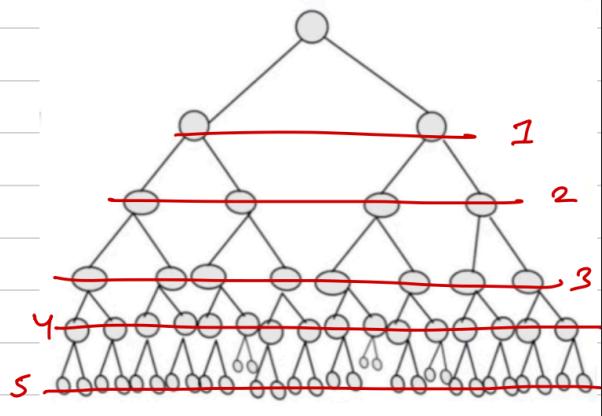
0 - 100

0 - 20

20 - 400

OverFitting and UnderFitting

Max-depth, very high value - overfitting.



underfitting.





prune²

/pruːn/

verb

gerund or present participle: **pruning**

trim (a tree, **shrub**, or bush) by cutting away dead or **overgrown** branches or **stems**, especially to encourage growth.

"now is the time to prune roses"

Similar: [cut back](#) [trim](#) [thin](#) [thin out](#) [pinch back](#) [crop](#) [clip](#) [▼](#)

• cut away (a branch or stem) from a tree, shrub, etc.

"prune back the branches"

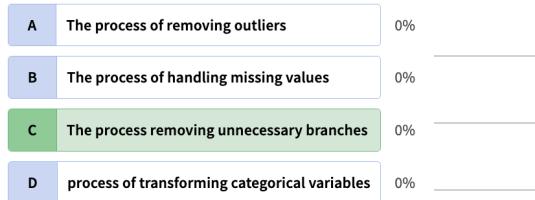
Similar: [cut off](#) [lop \(off\)](#) [chop off](#) [hack off](#) [clip](#) [snip \(off\)](#) [▼](#)

• reduce the extent of (something) by removing **superfluous** or unwanted parts.

"the workforce was pruned"

What is pruning in decision trees?

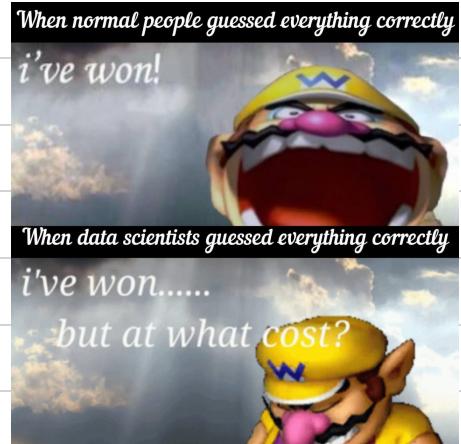
0 users have participated



[End Quiz Now](#)

Based on all quizzes from the session

	Deependu Ghosh	2/2	181.57
	Rakesh Karade	2/2	185.30
	Harshitha Chowd...	2/2	178.20
4	Aditya Shandilya	2/2	175.87
5	Purnishant Kumar	2/2	174.30
6	Snehal Adhikary	2/2	169.23
7	SHASHANK JHA	2/2	168.50
8	Kiran Hebasur	2/2	162.23
9	Praveen	2/2	162.23
10	MADDEHA REHMAN	2/2	161.73



What does each of these hyperparameters do?

- **Criterion:**

- Defines which impurity measure to use.
- By default, it is set to Gini Impurity.

- **max_depth:**

- Defines maximum depth upto which a tree will grow.
- By default, it is **None** i.e.
- it'll grow until all leaves are pure.

- **min_samples_split:**

- Defines the minimum no. of datapoints required to split further.
- Helps control depth which therefore prevents **overfit**.
- By default, it is 2.

- **min_samples_leaf:**

- Defines the minimum no. of samples a leaf node can have.

- **max_leaf_nodes:**

- Defines the maximum no. of leaf nodes a tree can have.
- **max_features:**

 - Selects features to be used while deciding the best split.
 - By default, it considers all the features for split.

- **class_weight:**

- Assign weights to different classes during training.
- Helps in dealing with imbalanced data.

How the performance of a DT will be affected if we tend to increase maximum number of pure leaf nodes?



Which of the given DTs will be affected by an outlier?

4 options

Active Duration (Most preferred: 30 seconds)

Appears for	60 Secs
-------------	---------

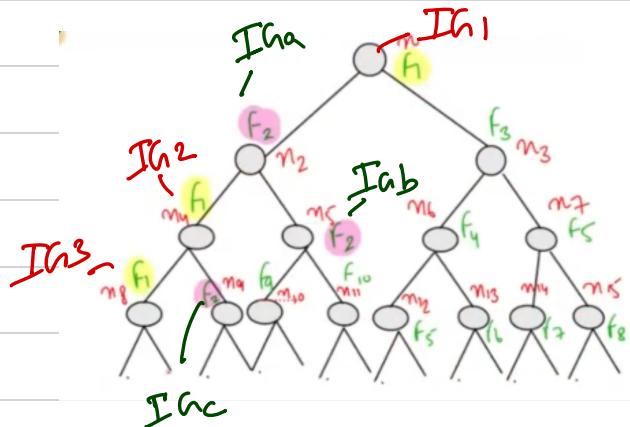
- | | |
|---|------------------|
| A | DT with depth=50 |
| B | DT with depth=5 |
| C | Both of them |
| D | None of them |



Feature Importance

$\text{Imp}(F_1) =$

$$\frac{n_x}{N} \times I_{G1} + \frac{n_y}{N} \times I_{G2} + \frac{n_z}{N} \times I_{G3}$$



$$\text{Imp}(F_2) = \frac{n_2}{N} \times I_{Gm} + \frac{n_5}{N} \times I_{Gb} + \frac{n_9}{N} \times I_{Gc}$$

