

# Decision Trees

10 March 2024 22:06

## Why decision trees?

### 1. Nonlinearity:

Unlike regression models, which assume a linear relationship between the predictors and the target, decision trees can handle complex interactions and non-linearities.

It can do without needing for any transformation of the features.

Decision trees can capture non-linear relationships between predictors and the target.

### 2. Interpretability:

Decision rules represented by the tree structure are easy to understand and visualize, making them accessible to executives who are non-experts.

$$[VP, Director] \leftarrow \text{logit } \underline{\log(\text{odds})}$$

DTs are like a flow chart and business users love flow charts.

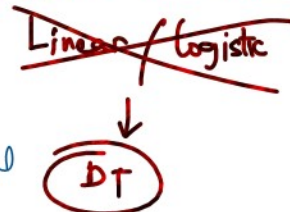
### 3. Scalability:

Decision trees are computationally efficient and scalable to large datasets, making them suitable for real-time applications and large-scale data processing.

Training dataset = 50M rows  $\rightarrow$  on big data

### 4. Handling Mixed Data Types

DTs can handle both numerical and categorical data without the need for one-hot encoding.

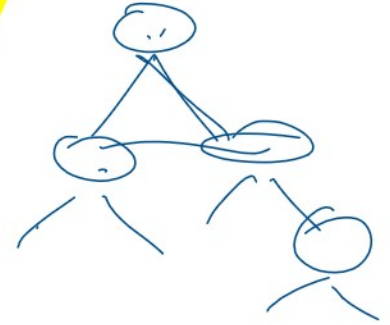


## 5. Robustness to outliers and irrelevant features

DTs can handle outliers & noise in the data along with irrelevant features without significantly impacting the model performance.

[Gini Impurity, Entropy]

↳ to evaluate the important features



## 6. Automatic selection of features

DTs perform automatic feature selection by identifying the most important features at each split.

## Intuition behind Decision Trees

"Decision trees are everywhere."

Sonali → is trying to buy a laptop.

PRICE

BRAND

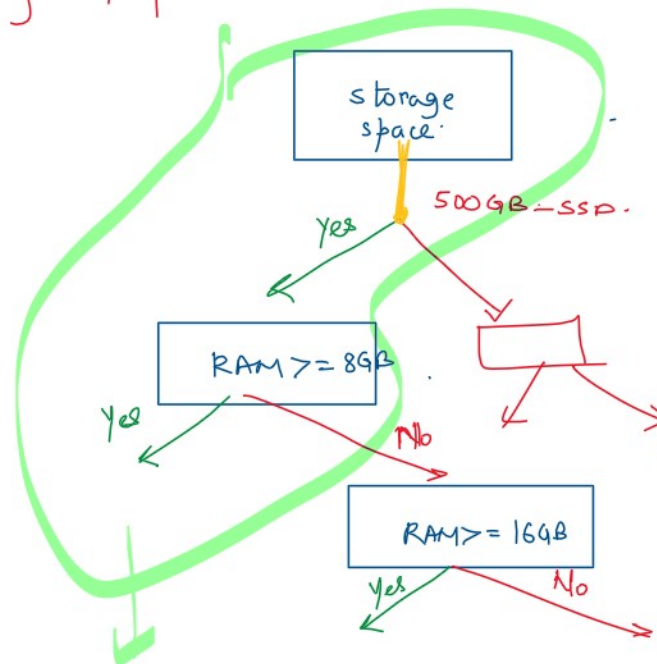
RAM

ROM

CPU

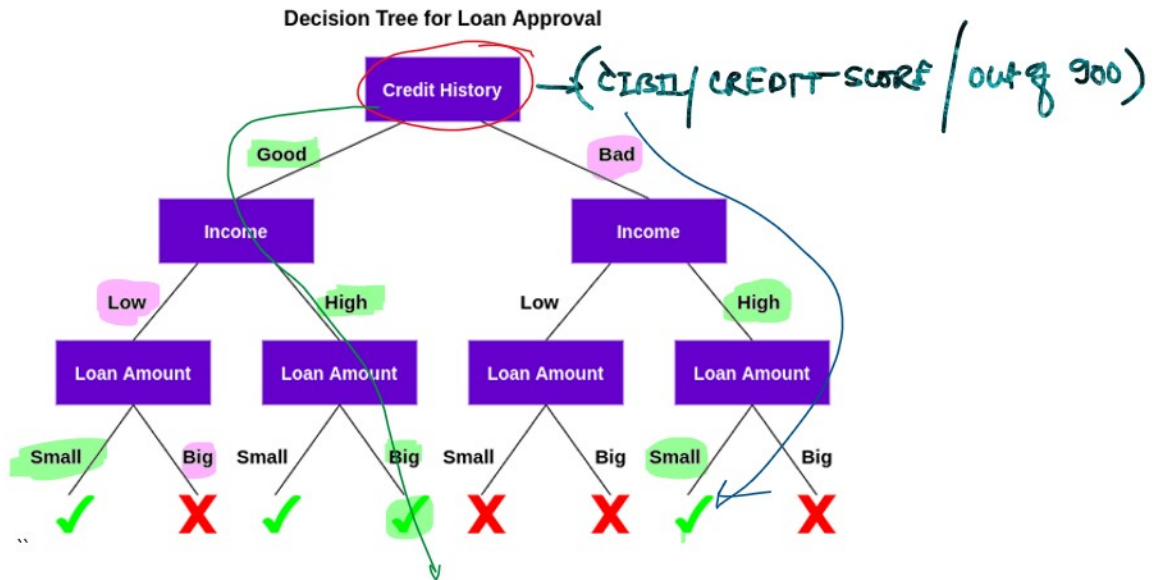
GPU

DISPLAY

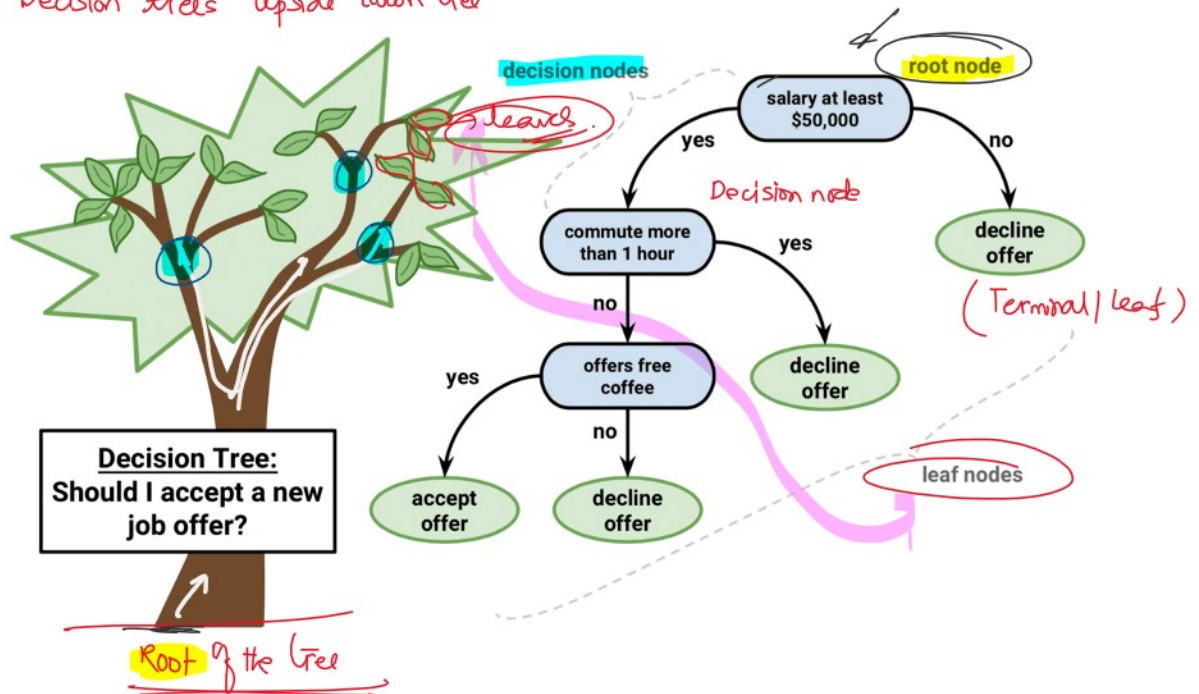


"Buy the laptop"

## # Loan approval Decision Trees



Decision trees upside down tree



## # Root Node:

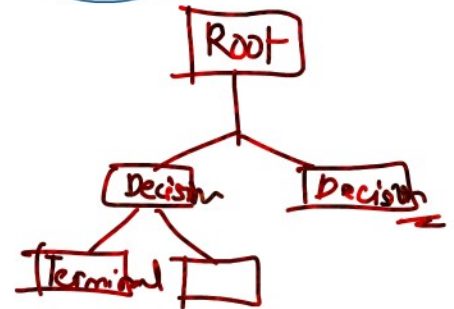
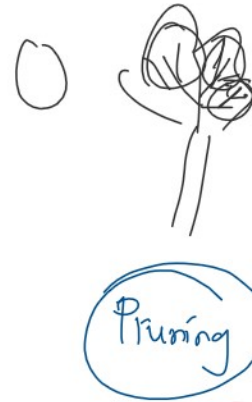
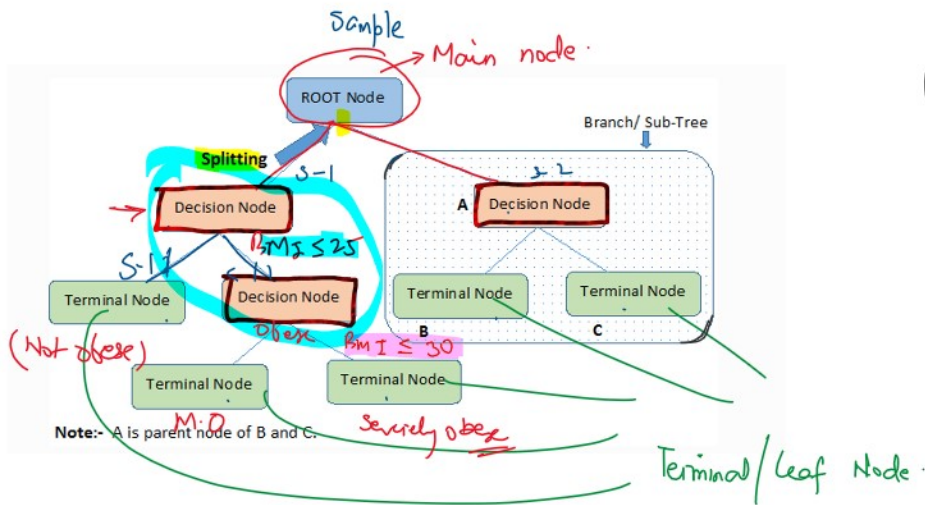
It represents entire population or sample and this further gets divided into two or homogeneous sets.

## # Splitting

It is a process of dividing a node into two or more subnodes.

## # Decision Node:

When a sub-node splits into further sub-nodes, it is called a decision node.



## # Leaf / Terminal Node:

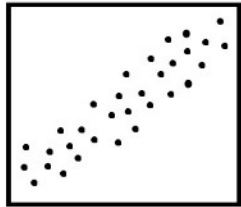
Nodes which do not split is called a leaf / terminal node.

## # Branch / Sub-Tree

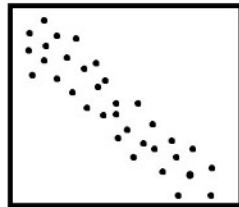
It is a sub-section of entire tree.

Scatter Plot

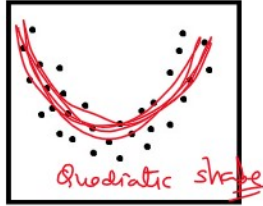




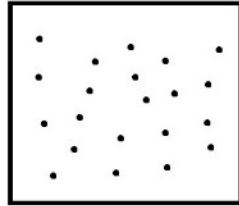
positive linear association



negative linear association



nonlinear association



no association

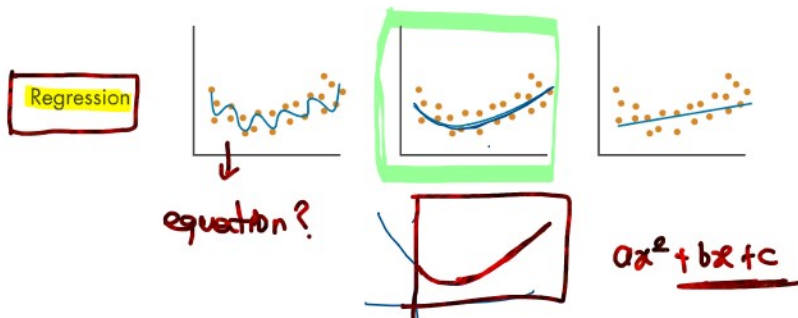
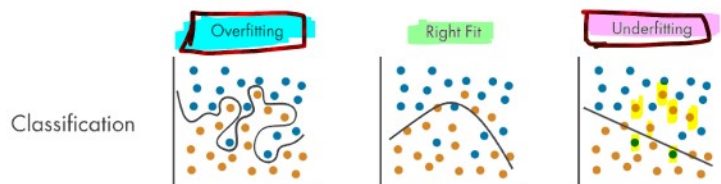
## # Disadvantage / Cons.

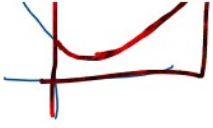
overfitting: one of most practical difficulties for decision tree models.

## # Overfitting and Underfitting

overfitting:

It occurs when a model learns the training data too well (too much), capturing noise or random fluctuations as well that may not be representative of true underlying patterns in the data.



Equation:   $ax^2 + bx + c$

Signs of overfitting:

- Model performs exceptionally well on the training data however it fails to generalize to new/unseen data.  
(performs poorly on the test set)

Reasons for overfitting

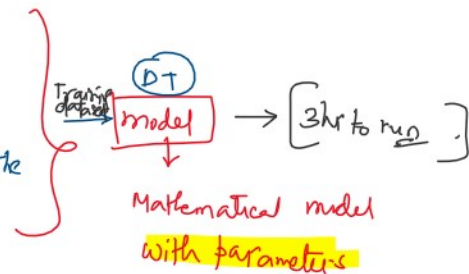
- using a highly complex model with too many parameters
- having insufficient data to support the complexity of the model.

For ex: Need 3 years (36 months) of training data for forecasting use-cases.



- training the model for a lot number of epochs.

one epoch means in ML,  
one complete pass of the  
training dataset through the  
algorithm



Mitigation:

- Use simpler models  $\Rightarrow$  less no. of features

(5-20)



b) do not overkill with the  
no. of training parameters.

- Use/increase the amount of training data.

- Regularization techniques to penalize overly complex models.

Error	Overfitting	Right Fit	Underfitting
Training	Low	Low	High
Test	High	Low	High

all the modelers pursue this

Example: Correlation vs Causality

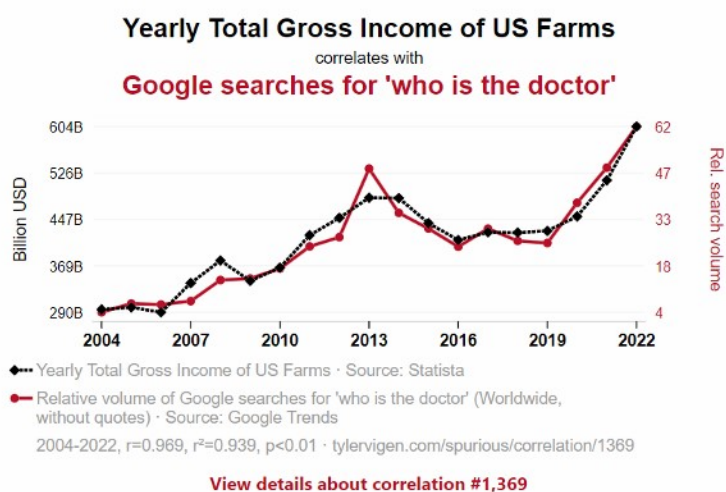
(X): Coding hours  
 (Y): Exam score (out of 100)  
 (Z): Lines of code

Python Exam

[X and Y are highly correlated]

High Z ~~→~~ Better Y (exam score)  
(Not correct)

conclusion: X and Z are highly correlated but they are not related.



✓ Y → X<sub>10</sub> → 80%

[View details about correlation #1,369](#)

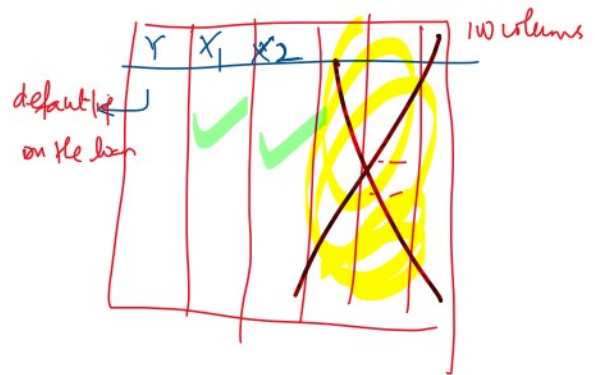
<https://www.tylervigen.com/spurious-correlations>

## # Underfitting

- Underfitting occurs when a model is too simple to capture the underlying patterns in the training data, resulting in poor performance on both training & testing.

## # signs of underfitting

- the model performs poorly on training data
- model is very simple, with almost no features
- it also performs poorly on test set.



## # Reasons for underfitting

- too simple model with too few parameters
- insufficient training or not allowing the model to learn enough during training.

## # Mitigation

- increasing the model complexity by adding more parameters
- ensure sufficient training time and training data
- use a sophisticated model.

✓ cross-validation (cv)

✓ regularization

} to strike the right balance between underfitting & overfitting.



- regularization
  - hyperparameter tuning
- to strike the right balance between underfitting & overfitting.

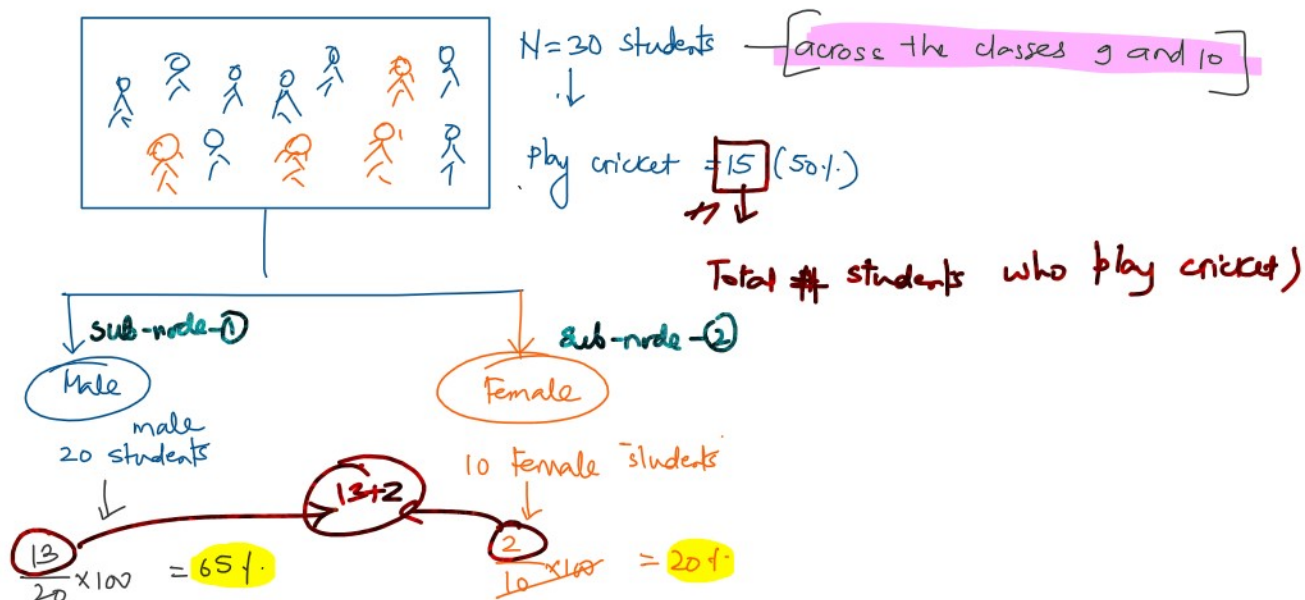
Where to split?

DTs use algorithm(s) to decide to split a node in two or more subnodes basis a criterion - **Gini algorithm**

Gini algorithm

Want to group students based on target variable - playing cricket or not.

Variable #1 Gender: **split on Gender**



# calculate Gini for sub-nodes, using formula.

$$(p^2 + q^2)$$

where  $p$  = probability of success

and  $q$  = probability of failure

and  $q$  = probability of failure

$$\textcircled{1} \text{ Gini index for sub-node male} = (0.65 \times 0.65) + (0.35 \times 0.35)$$

$$p = 0.65 \quad = 0.4225 + 0.1225$$

$$1-p = q = 1-0.65 = 0.35 = 0.5450$$

$$\underline{\underline{0.55}}$$

$$\textcircled{2} \text{ Gini index for sub-node female} = (0.2 \times 0.2) + (0.8 \times 0.8)$$

$$p = 0.2$$

$$= 0.04 + 0.64$$

$$q = 1-p = 0.8$$

$$= 0.68$$

calculated weighted gini index for split at gender

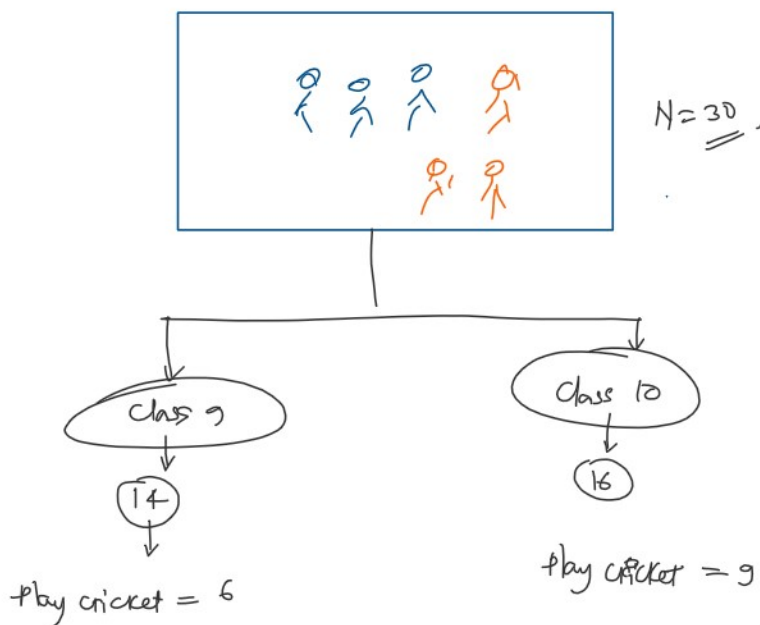
(# students in each M & F)

$$= \left( \frac{20}{30} \times 0.55 + \frac{10}{30} \times 0.68 \right)$$

$$= 20 \times .55 / 30 + 10 \times .68 / 30 = 0.5933$$

split @ gender  $\rightarrow$  Gini index  $\rightarrow \underline{\underline{0.5933}}$ .

Variable #2 class



$$\downarrow$$

$$\text{play cricket} = 6$$

$$\frac{6}{14} \times 100 = 43\%$$

Gini-index for subnode class IX

$$= (0.43 \times 0.43) + (0.57 \times 0.57)$$

$$= 0.51$$

$$\text{play cricket} = 9$$

$$\frac{9}{16} \times 100 = 56.25\%$$

$$9/16 = 0.5625$$

Gini-index for sub-node class X

$$= (0.56 \times 0.56) + (0.44 \times 0.44)$$

$$= 0.51$$

weighted Gini-index is 0.51

Conclusion: Gini index for Gender is higher than split on class, the node (root) split will take place on gender.

$$\text{Gini impurity} = (1 - \text{Gini})$$

$$\text{Gini impurity} \Bigg) \text{ split on Gender} = (1 - 0.59)$$

$$= \underline{\underline{0.41}}$$

Gini  
or  
Entropy