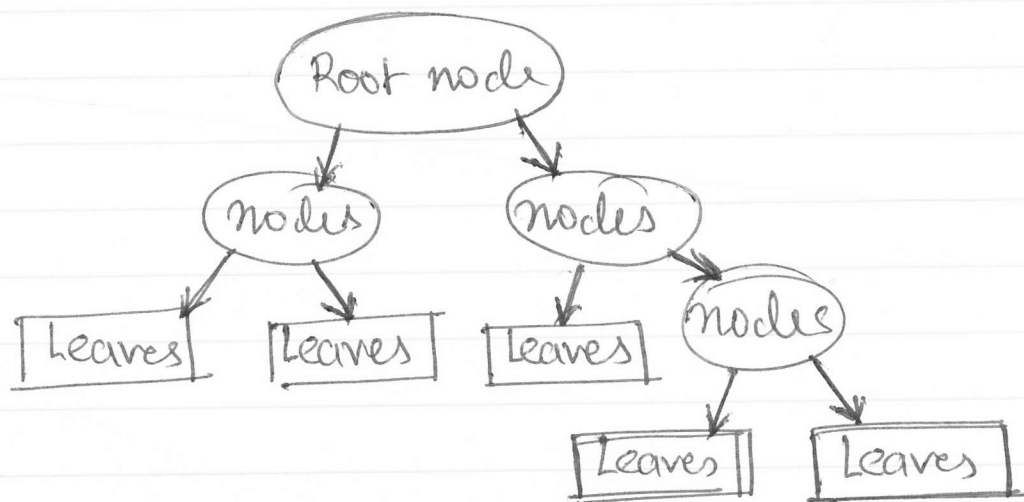


Decision Tree

- most commonly used for classification based on categorical or numerical data that can be discrete or continuous or a combination



— Data Set

Chest Pain (CP)	Good Blood Circulation (BC)	Blocked Arteries (BA)	Heart Disease (HD)
NO	NO	NO	NO
YES	YES	YES	YES
YES	YES	NO	NO
YES	NO	???	YES
⋮	⋮	⋮	⋮

- 3 attributes to predict heart disease
- there are 303 patients (rows) in total
- If the % of missing data is very less ($< 1\%$) we can skip rows with missing data.

— Step 1: Finding Root Node

- Finding how each attribute predict HD.

CP		BC		BA	
True	False	True	False	True	False
HD		HD		HD	
Y	N	Y	N	Y	N
105	39	37	127	92	31
34	125	100	33	45	129

- the nodes (leaf) do not separate the HD with 100% certainty — they are impure nodes

— Step 2: measure node impurity — Gini coefficient

- Calculating Gini coefficient for CP:

$$\text{For 'True' node: } 1 - (\text{Probability of Yes})^2 - (\text{Probability of No})^2$$

$$= 0.395 \text{ — with 144 patients}$$

$$\text{For 'false' node} = 0.336 \text{ — with 159 patients}$$

Total Gini coefficient of CP to predict HD

$$= \left(\frac{144}{144+159} \right) 0.395 + \left(\frac{159}{144+159} \right) 0.336$$

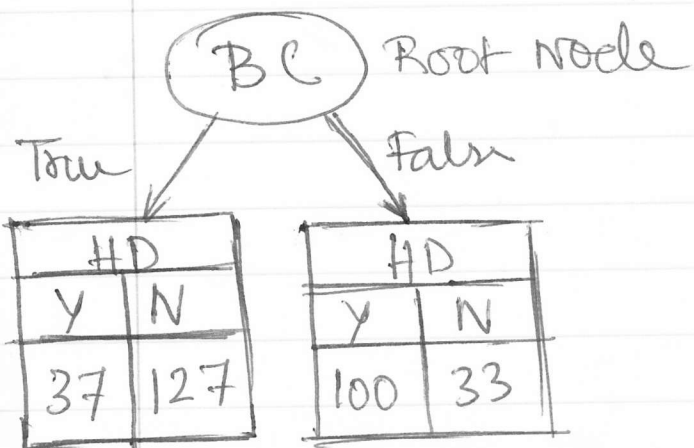
$$\Rightarrow \text{Gini coefficient for CP} = 0.364$$

$$\text{'' '' '' BC} = 0.360 \text{ and}$$

$$\text{'' '' '' BA} = 0.381$$

— Step 3: Finding other nodes + leaf nodes

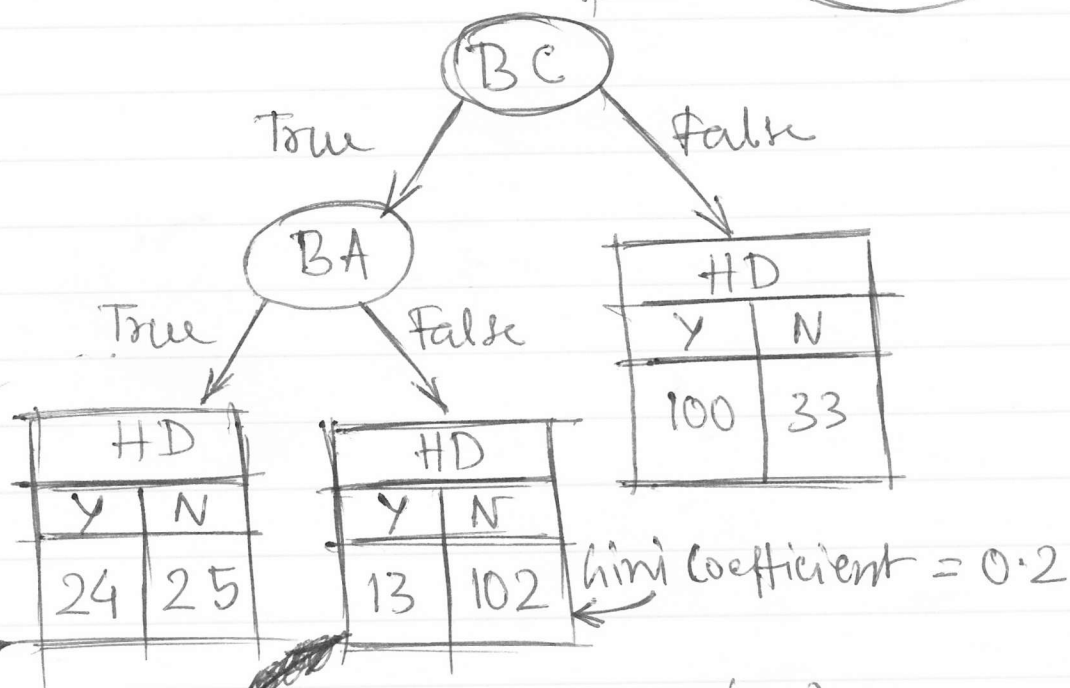
- Pure node has Gini Impurity = 0
- we need to choose the node/attribute with minimum Gini for the Root node — BC



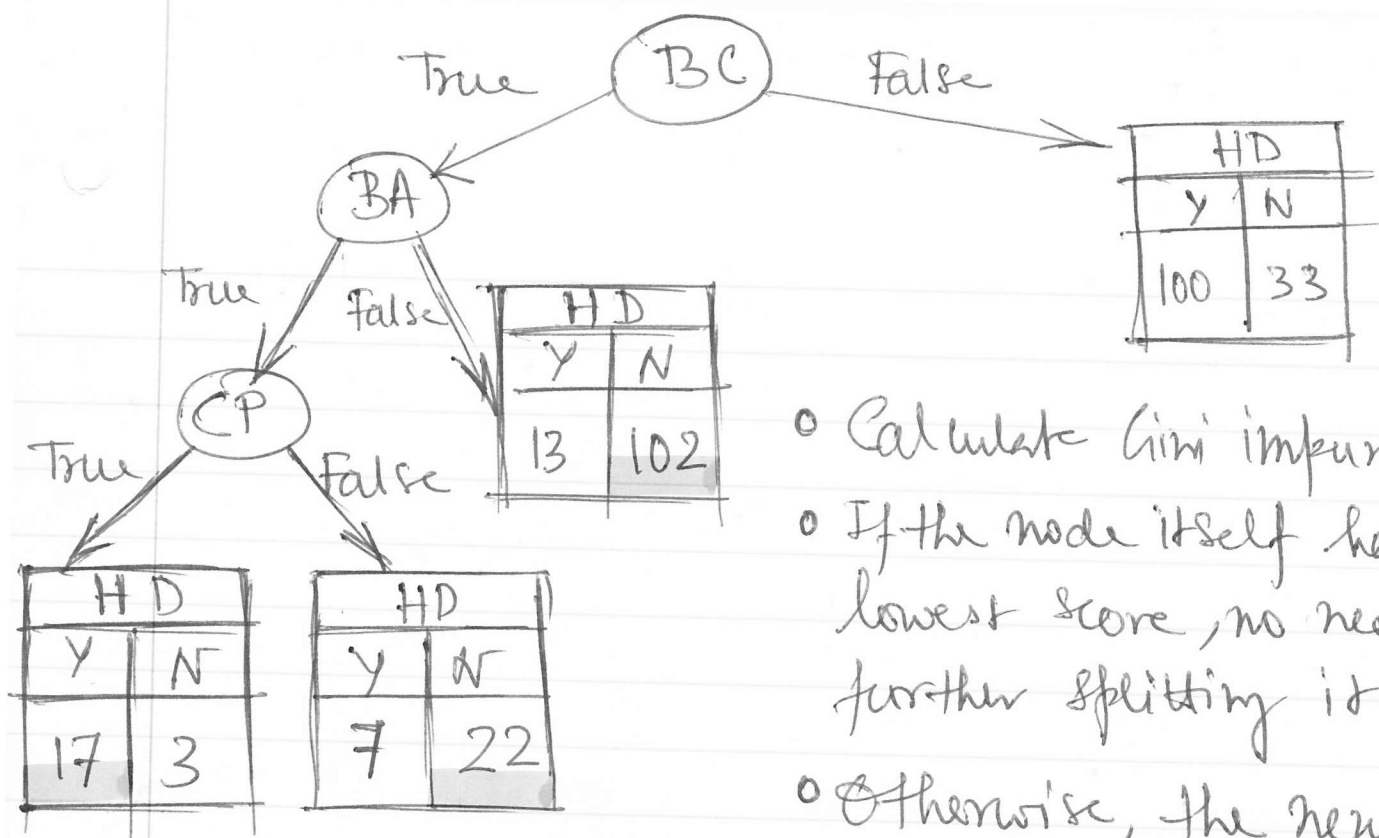
- Find out how well CP and BA separates these 164 (37+127) patients.

Gini coefficient for

CP = 0.3
BA = 0.29

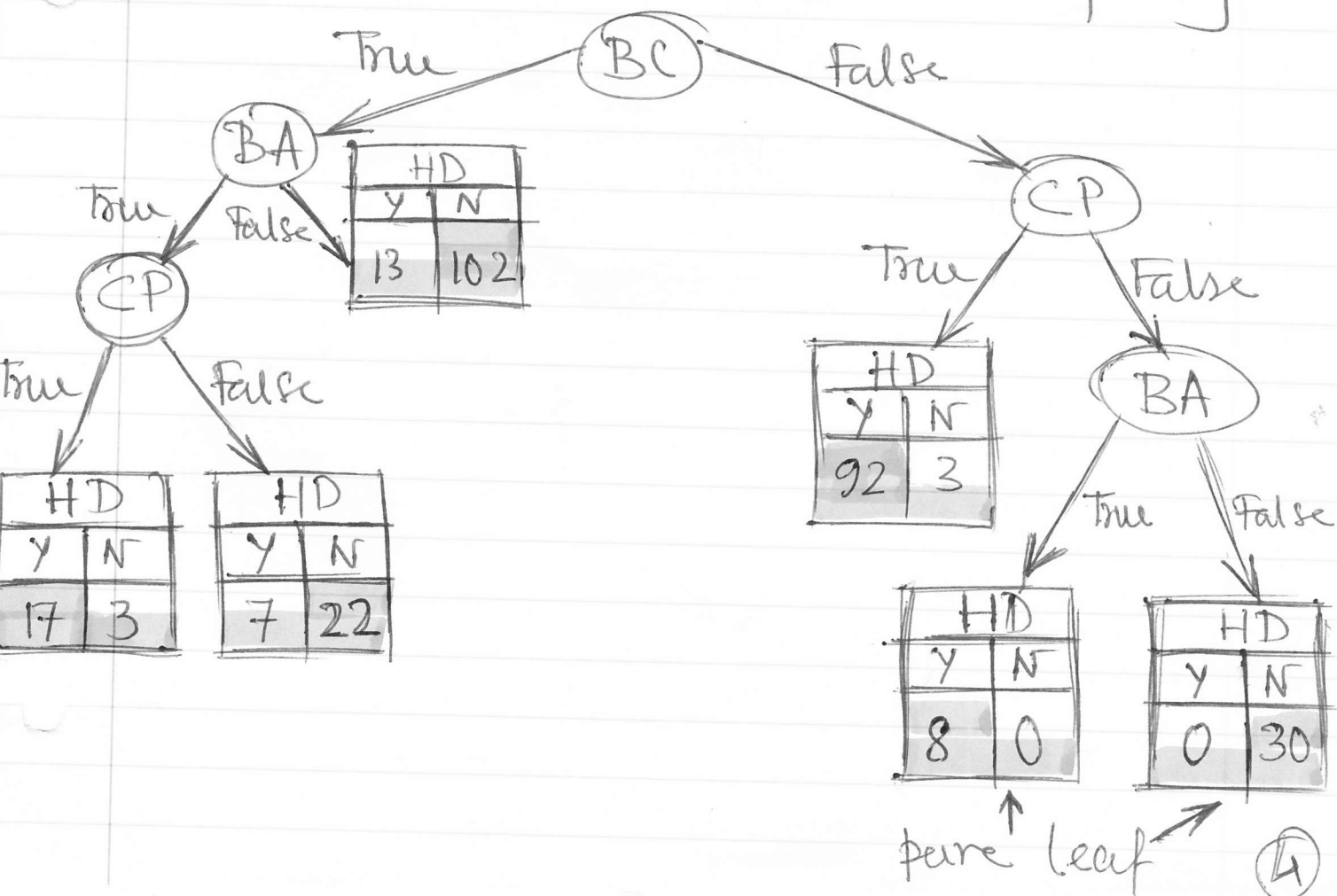


- Apply CP to separate patients (49) with and without HD (3)



- Calculate Gini impurity
- If the node itself has the lowest score, no need to further splitting it
- Otherwise, the next node will be the one with the lowest Gini impurity value.

Final Tree



Random forest

— Decision trees are, in general, prone to overfitting the data, and as a result do not perform well on a diverse set of data than the training data.

— Step 1: Create a "bootstrapped" dataset.

• In bootstrapping, we are allowed to pick the same sample (row) from a dataset multiple times.

Original Dataset

Chest Pain (cp)	Good Blood Circulation	Blocked Arteries (BA)	weight	Heart Disease (HD)
NO	NO	NO	125	NO
YES	YES	YES	180	YES
YES	YES	NO	210	NO
YES	NO	YES	167	YES

Chest Pain	Good Blood Circulation	Blocked Arteries	weight	Heart Disease
YES	YES	YES	180	YES
NO	NO	NO	125	NO
YES	NO	YES	167	YES
YES	NO	YES	167	YES (5)

- Step 2: Create a decision tree using the bootstrap dataset, but using a random subset of column (variable) at each step

- o let's say we are only using 'Good Blood Circulation' and 'Blocked Arteries' (2 random variables) to ~~create~~ ~~can~~ create decision tree to predict heart disease.



- o Say "Good Blood Circulation" is the root node, we again randomly select 2 variables: "Chest Pain" and "Weight" and so on to build a tree

⇒ Make another bootstrap dataset
↳ build another decision tree from random subset of variables.

⇒ At the end, we will have a variety of decision trees (~100) in the random forest.

Predicting using Random Forest:

- o we run the data through each tree and collect their prediction.
- o we make the final prediction on the aggregated data (usually) ~~using~~ based on maximum probability (6)

- Bagging = Bootstrapping dataset

+
Decision making using aggregated
prediction from individual tree.

Evaluating a Random Forest

- Usually, about $\frac{1}{3}$ of the $\textcircled{4}$ original data do not end up in the bootstrapped dataset.
— called Out-of-Bag Dataset
 - Out-of-Bag dataset is used to
 - evaluate the performance of the Random Forest
 - optimize the number of variables used from the bootstrapped dataset to build decision trees
 - Incorrectly classified Out-of-Bag samples are called Out-of-Bag Error.
-