

**Walmart** 

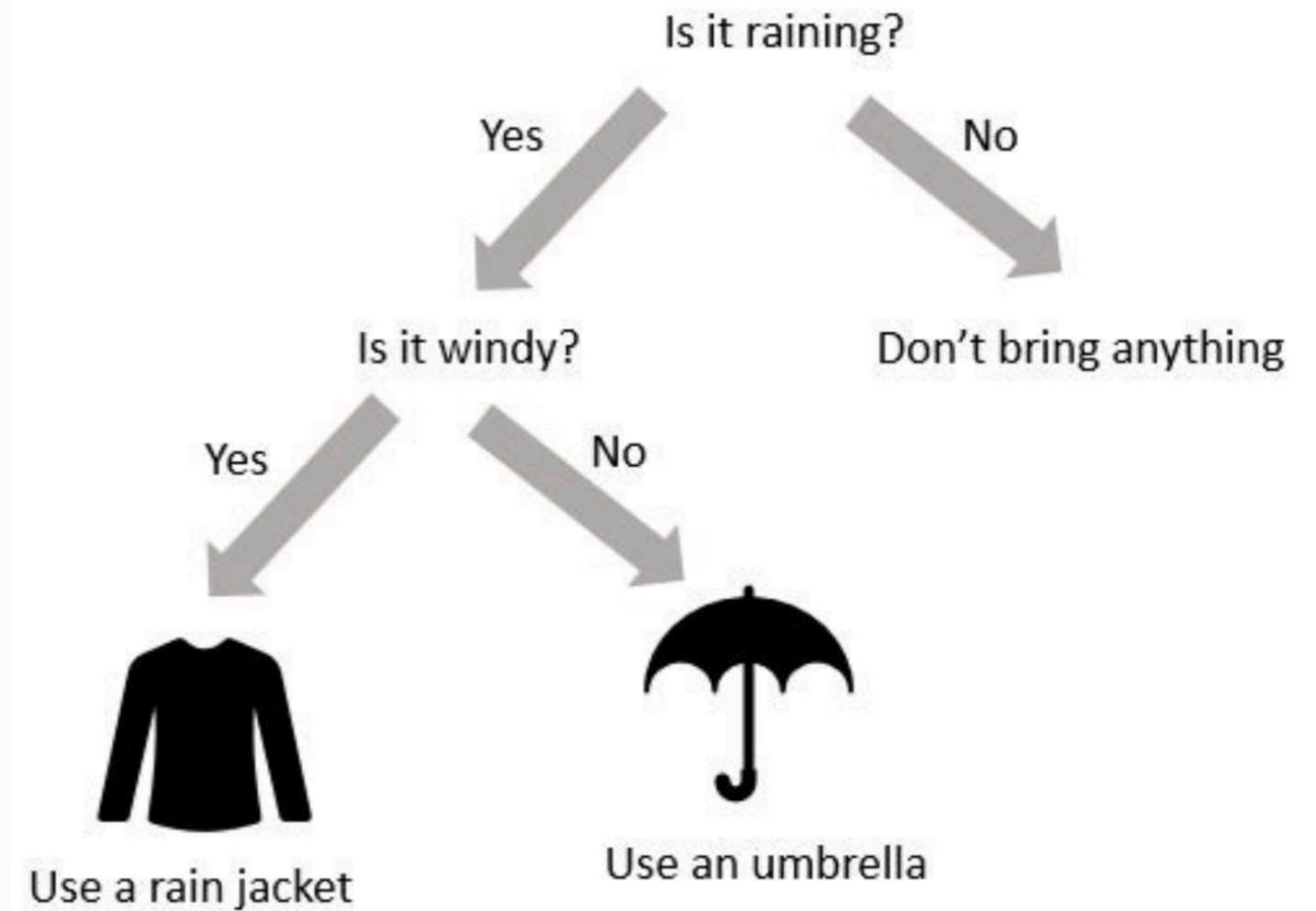
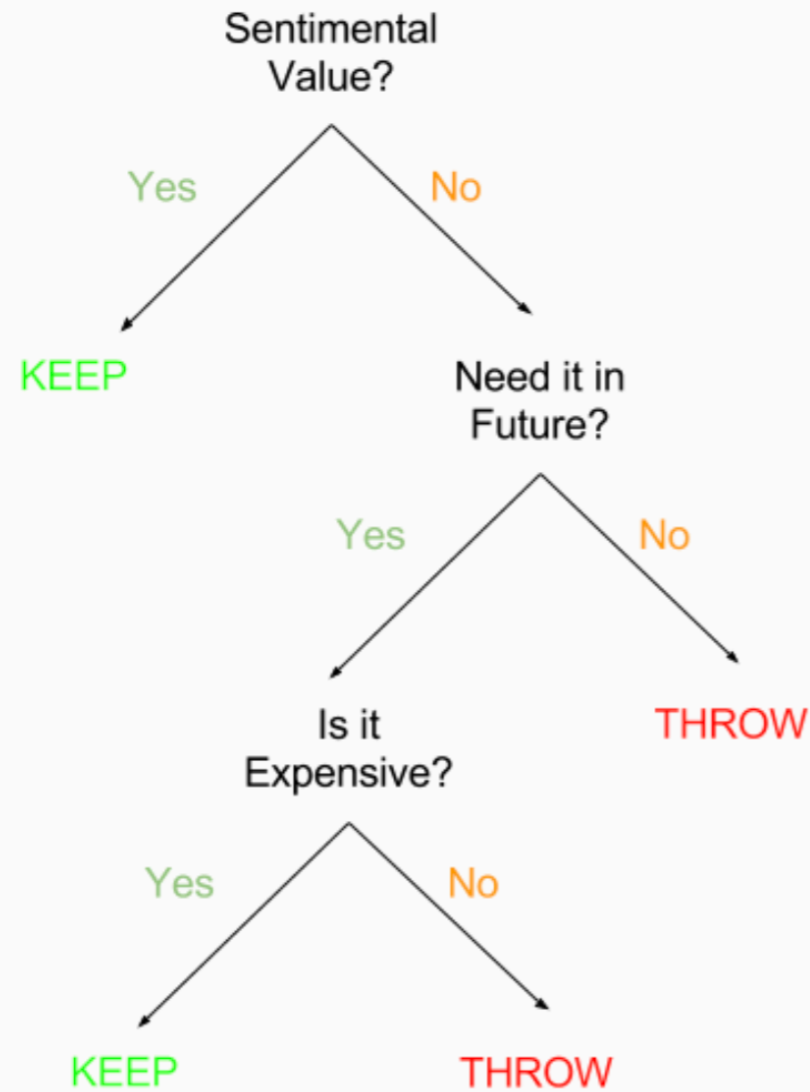
**Decision tree algorithms – 101**

Subhasish Misra & Somedip Karmakar  
**International Data & Analytics**

13/08/20

# DECISION TREE– AN OVERVIEW

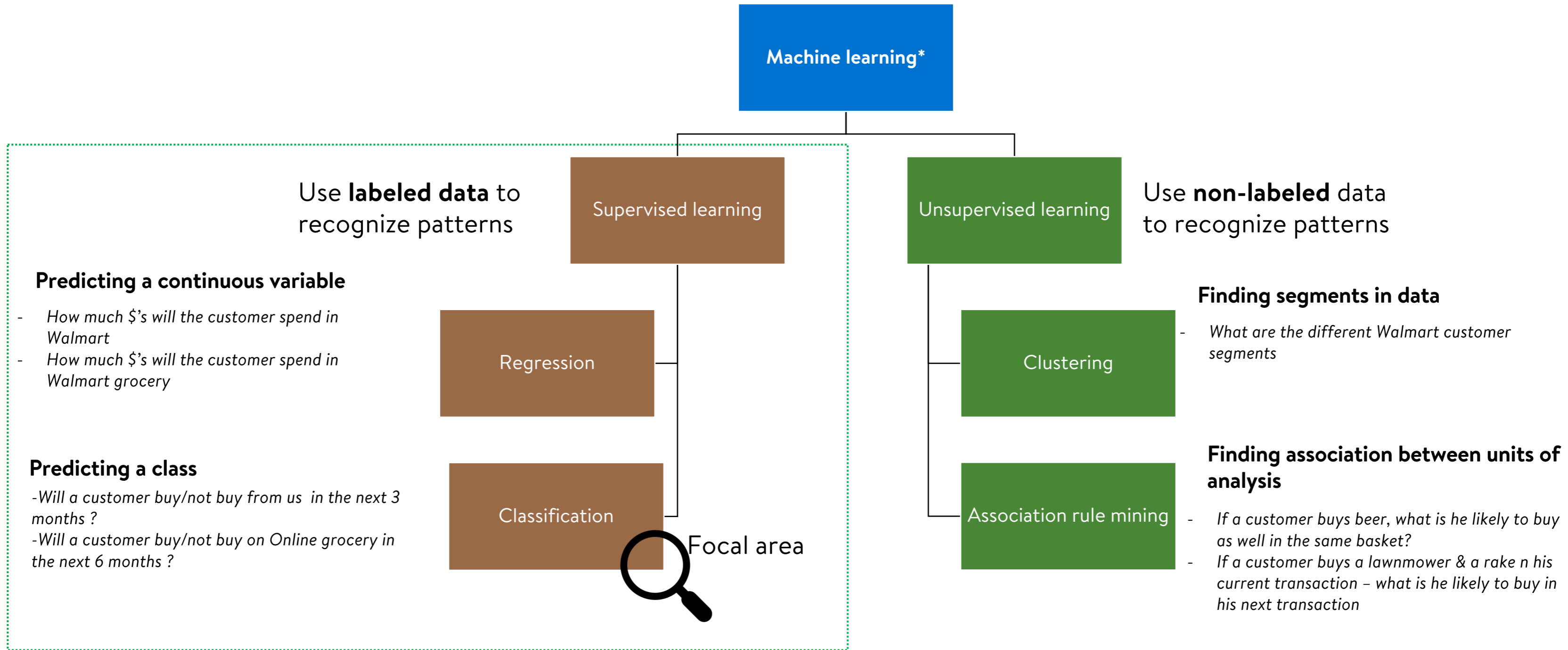
A set of exhaustive if/then rules that helps arrive at a decision !



<https://www.kdnuggets.com/2017/08/machine-learning-abstracts-decision-trees.html>

<https://www.kdnuggets.com/2019/02/decision-trees-introduction.html>

Broadly, the science of helping algorithms learn pattern in data



# Classification: Who wants a riding mower !



Predictors

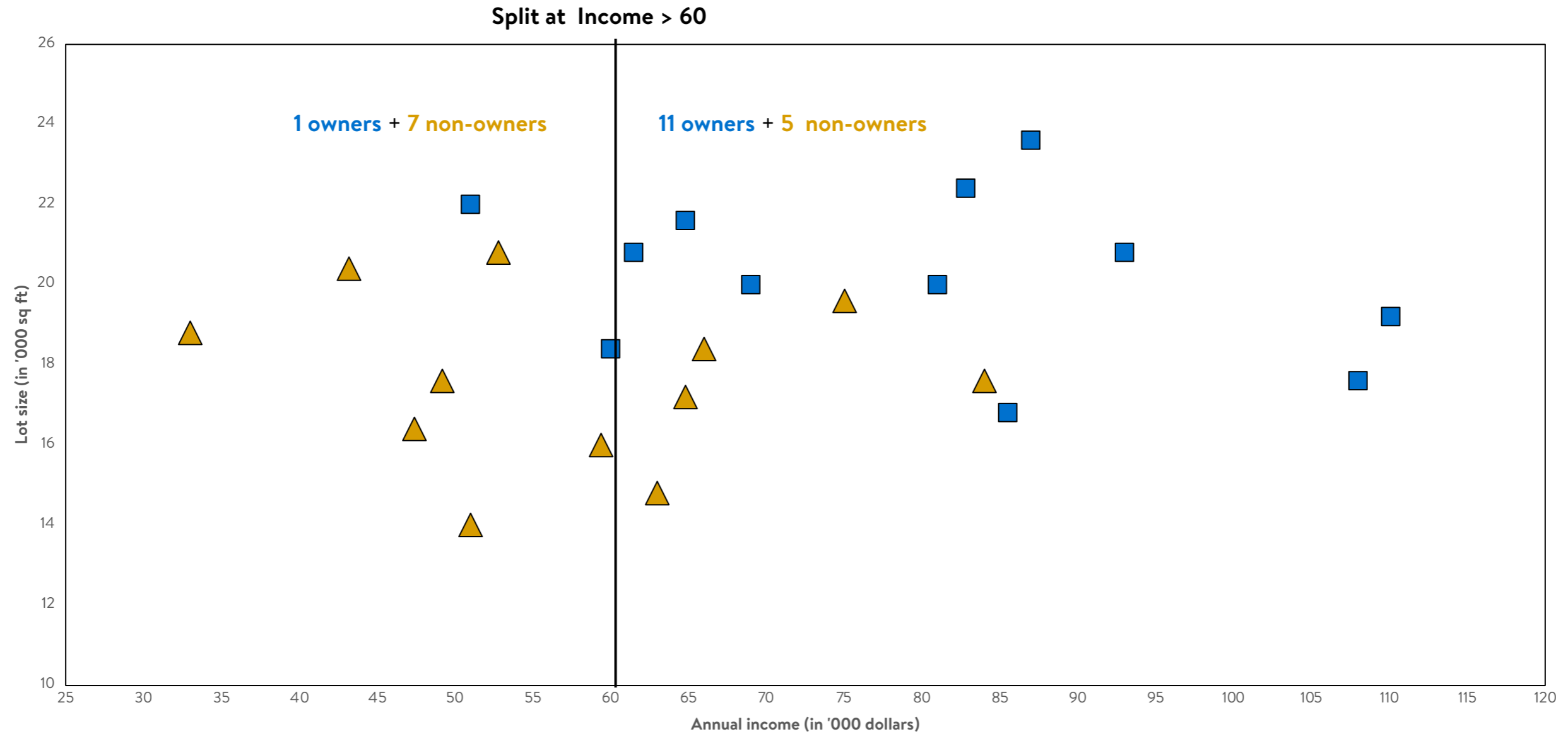
Observation	Income (in '000 \$'s)	Lot_Size (in '000 sqft)
1	60	18.4
2	85.5	16.8
3	64.8	21.6
4	61.5	20.8
5	87	23.6
6	110.1	19.2
7	108	17.6
8	82.8	22.4
9	69	20
10	93	20.8
11	51	22
12	81	20
13	75	19.6
14	52.8	20.8
15	64.8	17.2
16	43.2	20.4
17	84	17.6
18	49.2	17.6
19	59.4	16
20	66	18.4
21	47.4	16.4
22	33	18.8
23	51	14
24	63	14.8

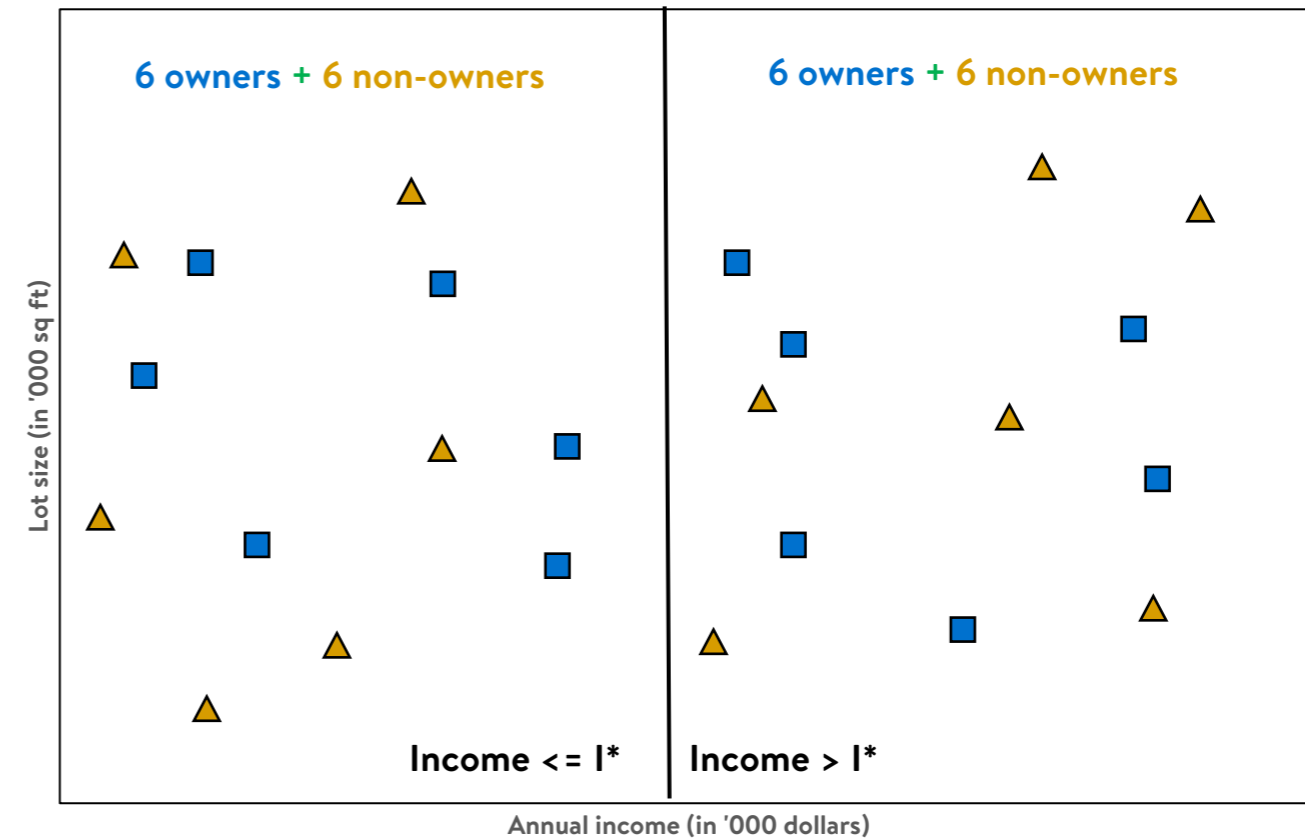
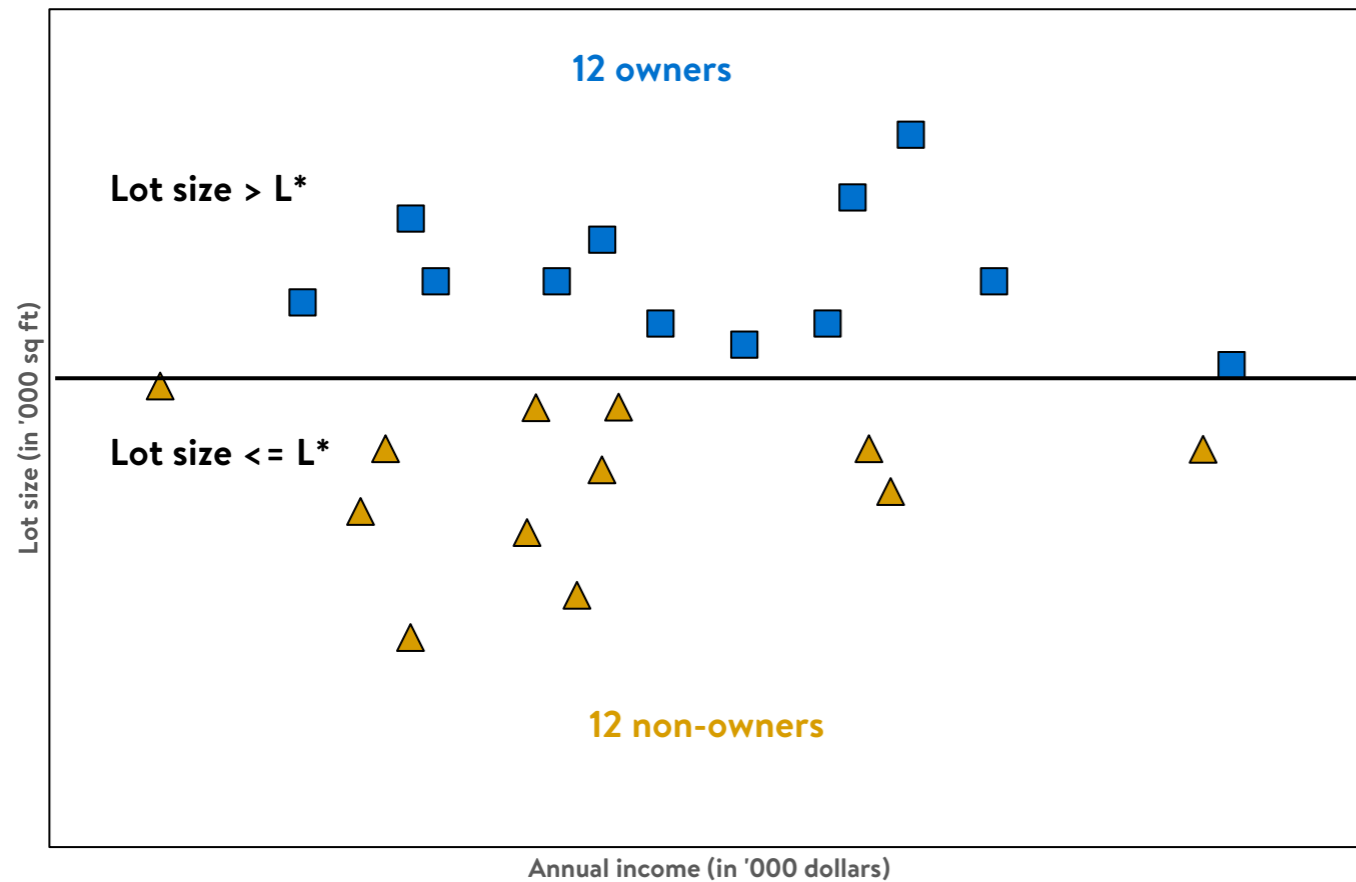
Target variable/or the variable we are trying to predict

Ownership
owner
owner
owner
owner
owner
owner
owner
owner
owner
owner
owner
owner
owner
owner
non-owner
non-owner
non-owner
non-owner
non-owner
non-owner
non-owner
non-owner
non-owner
non-owner
non-owner

Can we learn to 'classify' observations to ownership classes ?

Binary classification (2 classes)  
12 owners  
12 non-owner



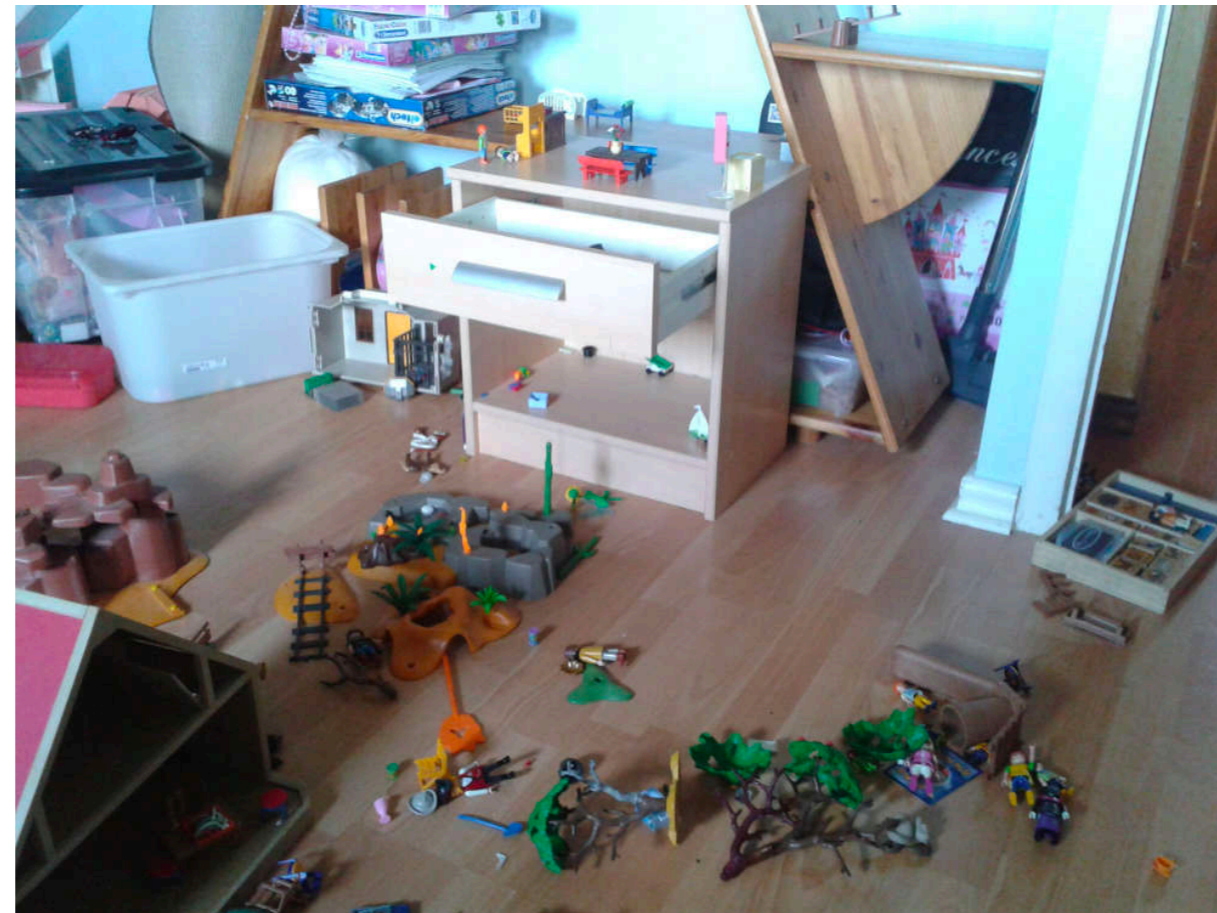


**Ideal scenario:** This split divides the data space to two homogenous ('pure') parts

**Not so Ideal scenario:** The two parts are fairly heterogenous

Entropy is a concept often used in computer science (information theory) and other fields of science too.  
*Entropy is a measure of disorder or heterogeneity.*

***A high on entropy room 😊***



*Translated in the context of data – how do we actually measure it ?*



Entropy:  $-\sum p_k \log_2 p_k$   
 Where  $p_k$  is the proportion of observations in class  $k$  belonging to a rectangle

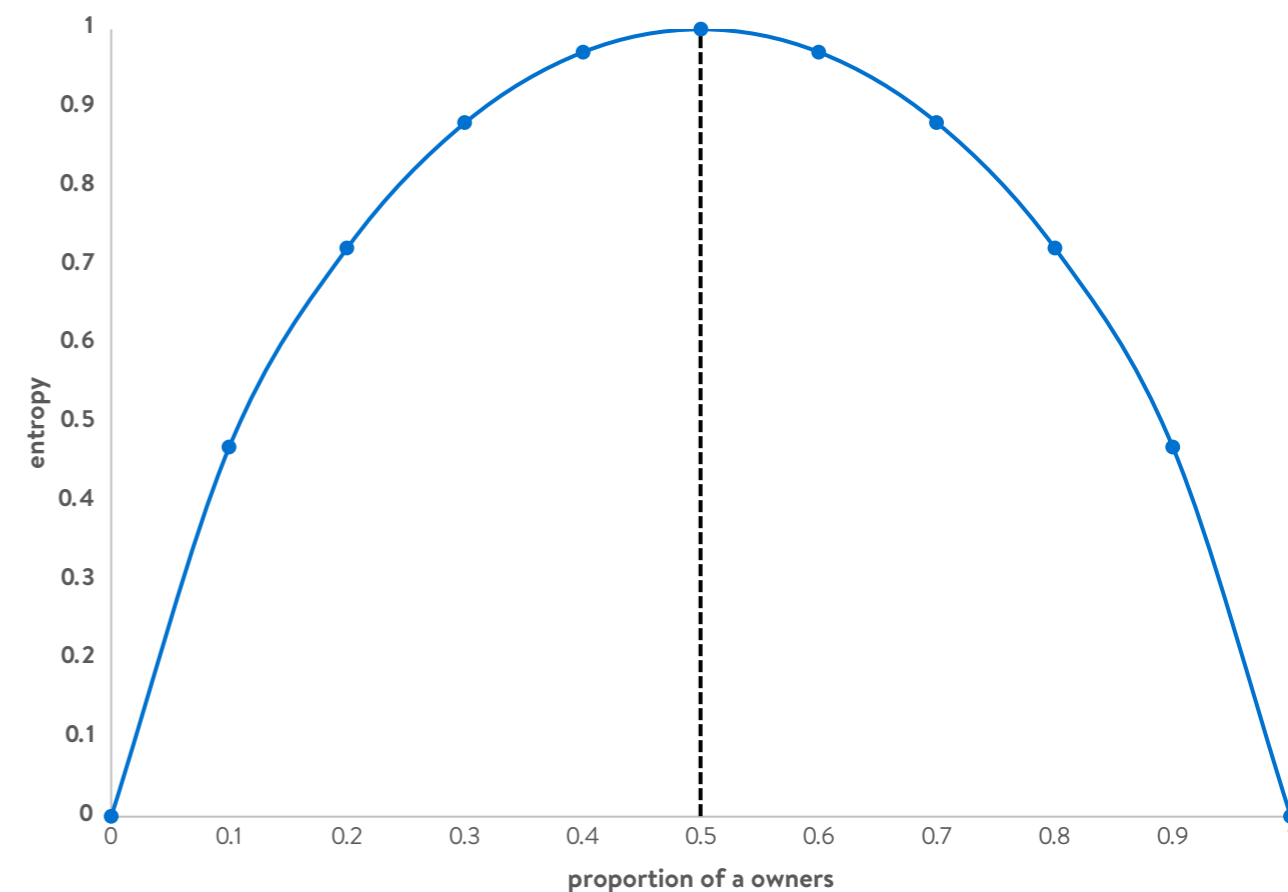


Entropy left block (Ditto for right block):

$$= - \left[ \frac{\text{owners}}{\text{owner} + \text{non-owner}} \cdot \log_2 \left( \frac{\text{owners}}{\text{owner} + \text{non-owner}} \right)^2 + \frac{\text{non-owners}}{\text{owner} + \text{non-owner}} \cdot \log_2 \left( \frac{\text{non-owners}}{\text{owner} + \text{non-owner}} \right)^2 \right]$$

$$= - \left[ \frac{6}{12} \cdot \log_2 \left( \frac{6}{12} \right)^2 + \frac{6}{12} \cdot \log_2 \left( \frac{6}{12} \right)^2 \right]$$

$$= - \left[ .5 \cdot -1 + .5 \cdot -1 \right] = 1$$



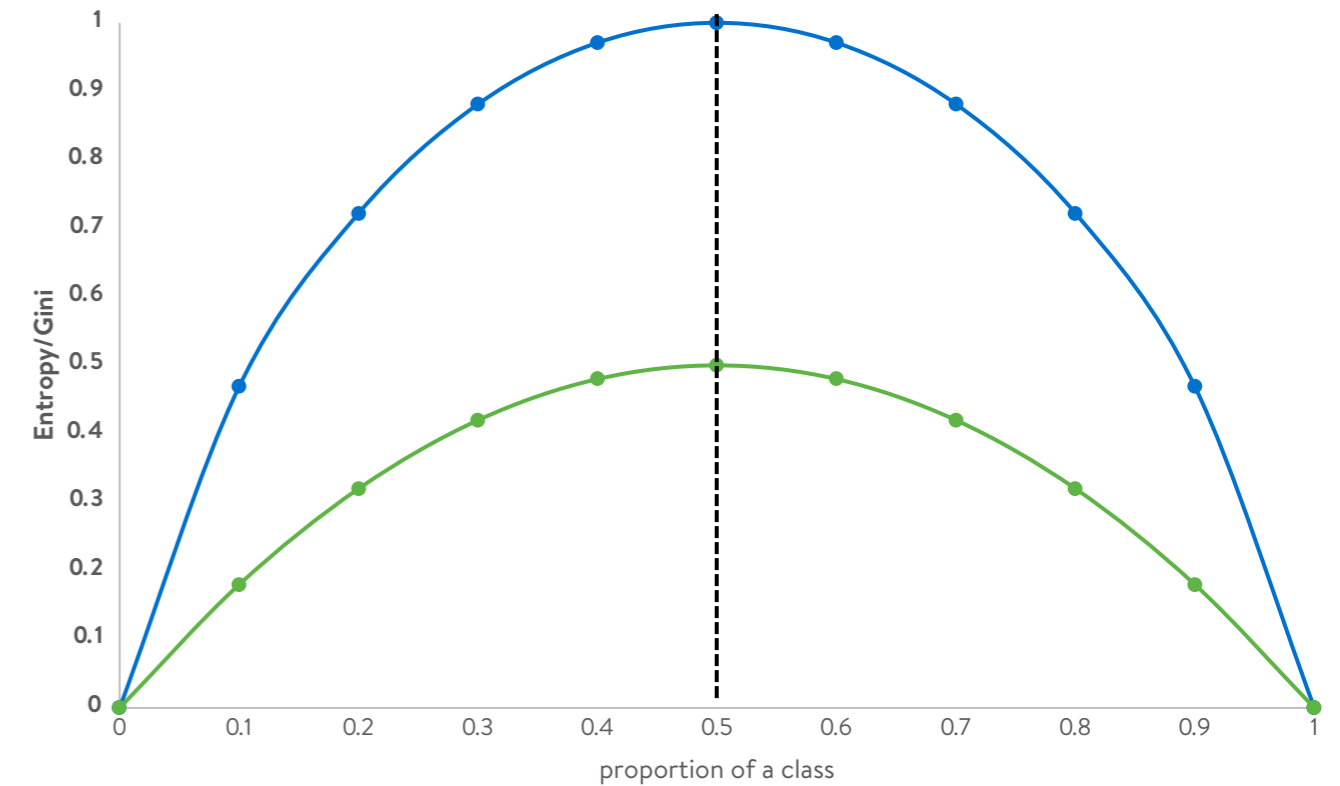
Gets **maximized** when we have equal representation of the classes (for a 2 class problem)

Gini index:  $1 - \sum p_k^2$   
 Where  $p_k$  is the proportion of observations in class k belonging to a rectangle

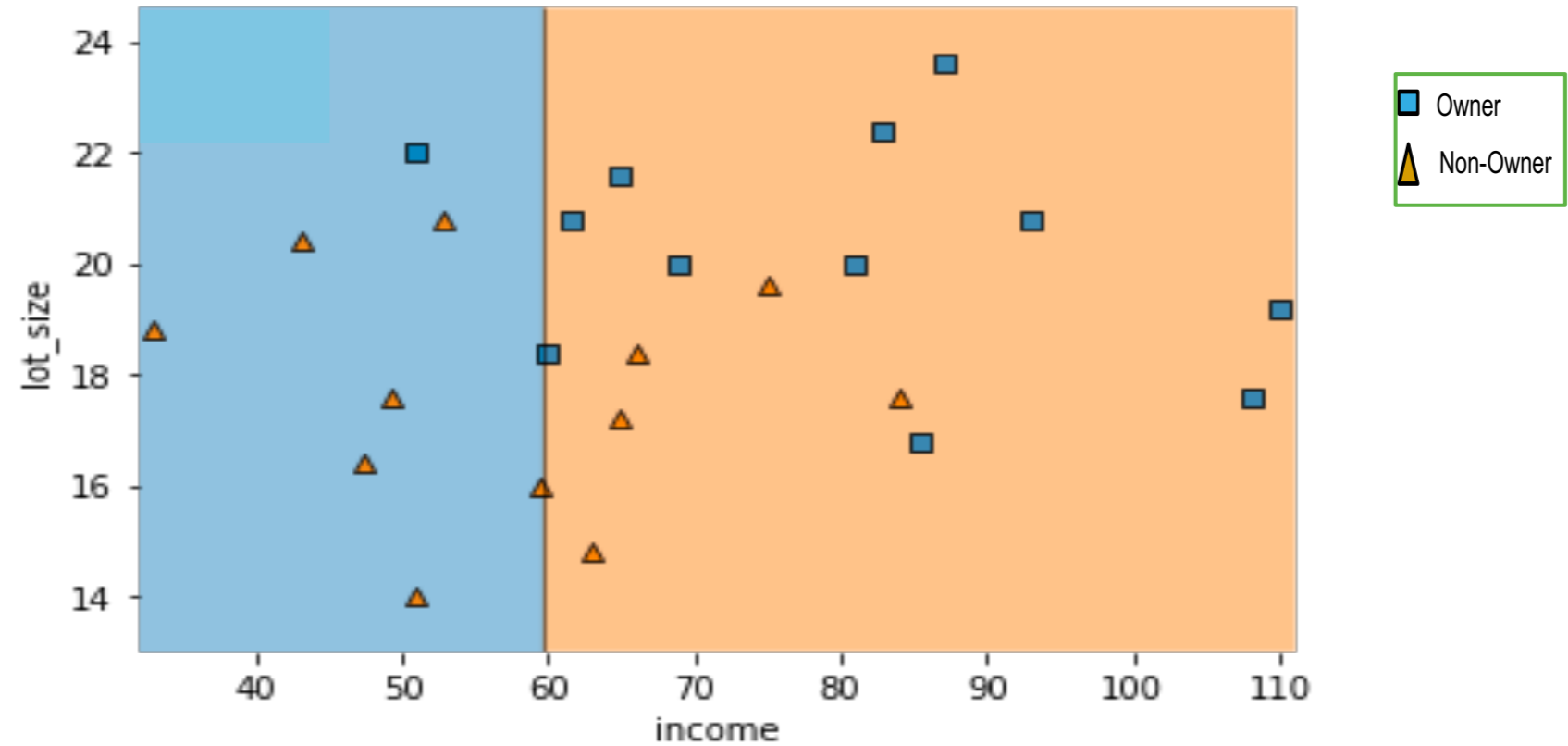


Gini index: Left block (Ditto for right):

$$\begin{aligned}
 &= 1 - (\text{owners}/\text{owner} + \text{non-owner})^2 - (\text{non owners}/\text{owner} + \text{non-owner})^2 \\
 &= 1 - (6/12)^2 - (6/12)^2 \\
 &= 1 - .25 - .25 \\
 &= .5
 \end{aligned}$$



Gets **maximized** when we have equal representation of the classes



$$\begin{aligned} \text{Gini (no split scenario)} &= 1 - (\text{owners/owner} + \text{non-owner})^2 - (\text{non-owners/owner} + \text{non-owner})^2 \\ &= 1 - (12/24)^2 - (12/24)^2 = 0.50 \end{aligned}$$

For a split at Income = 60 \$ scenario:

$$\text{Gini (Left block)} = 1 - (1/8)^2 - (7/8)^2 = \mathbf{0.22}$$

$$\text{Gini (Right block)} = 1 - (11/16)^2 - (5/16)^2 = \mathbf{0.43}$$

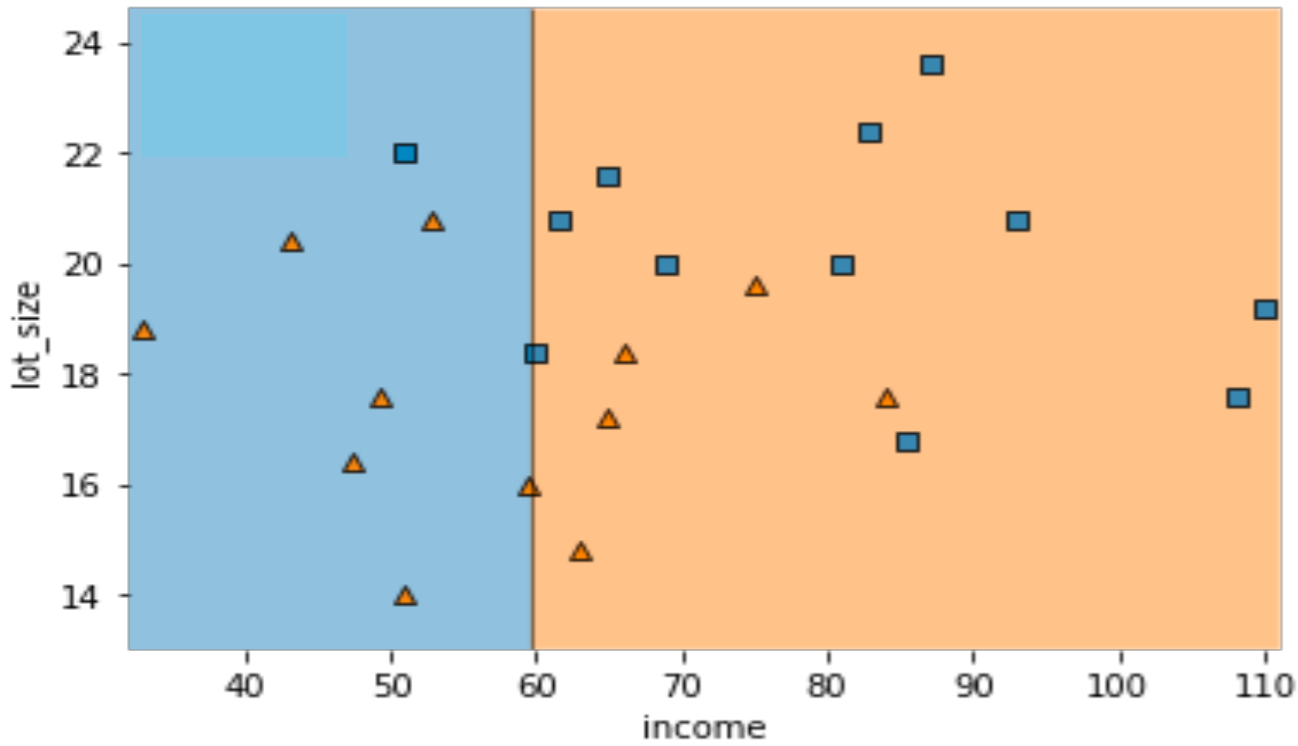
$$\begin{aligned} \text{Overall Gini} &= (\text{obs in left block} / \text{Total obs}) * \mathbf{0.22} + (\text{obs in right block} / \text{Total obs}) * \mathbf{0.43} \\ &= (8 / 24) * \mathbf{0.22} + (16/24) * \mathbf{0.43} = .35 \end{aligned}$$

$$\text{Information gain} = \text{Gini}_{\text{No split}} - \text{Gini}_{\text{after split}} = .50 - .35 = .15$$



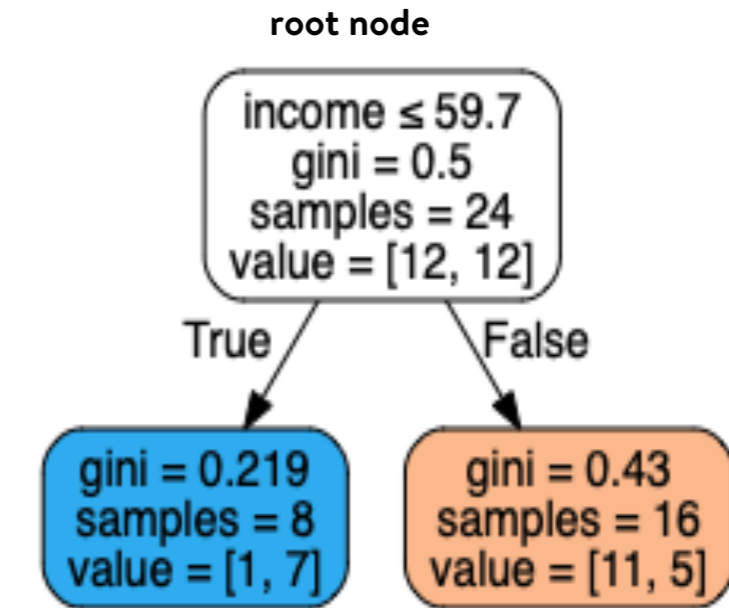
- Do this for all possible split points:
  - ✓ Across the two variables
  - ✓ Compare information gain
  - ✓ Choose split having max gain

Turns out income > 59.7 is best first split



Decision Tree representation

Tree stump – A simple tree with one split



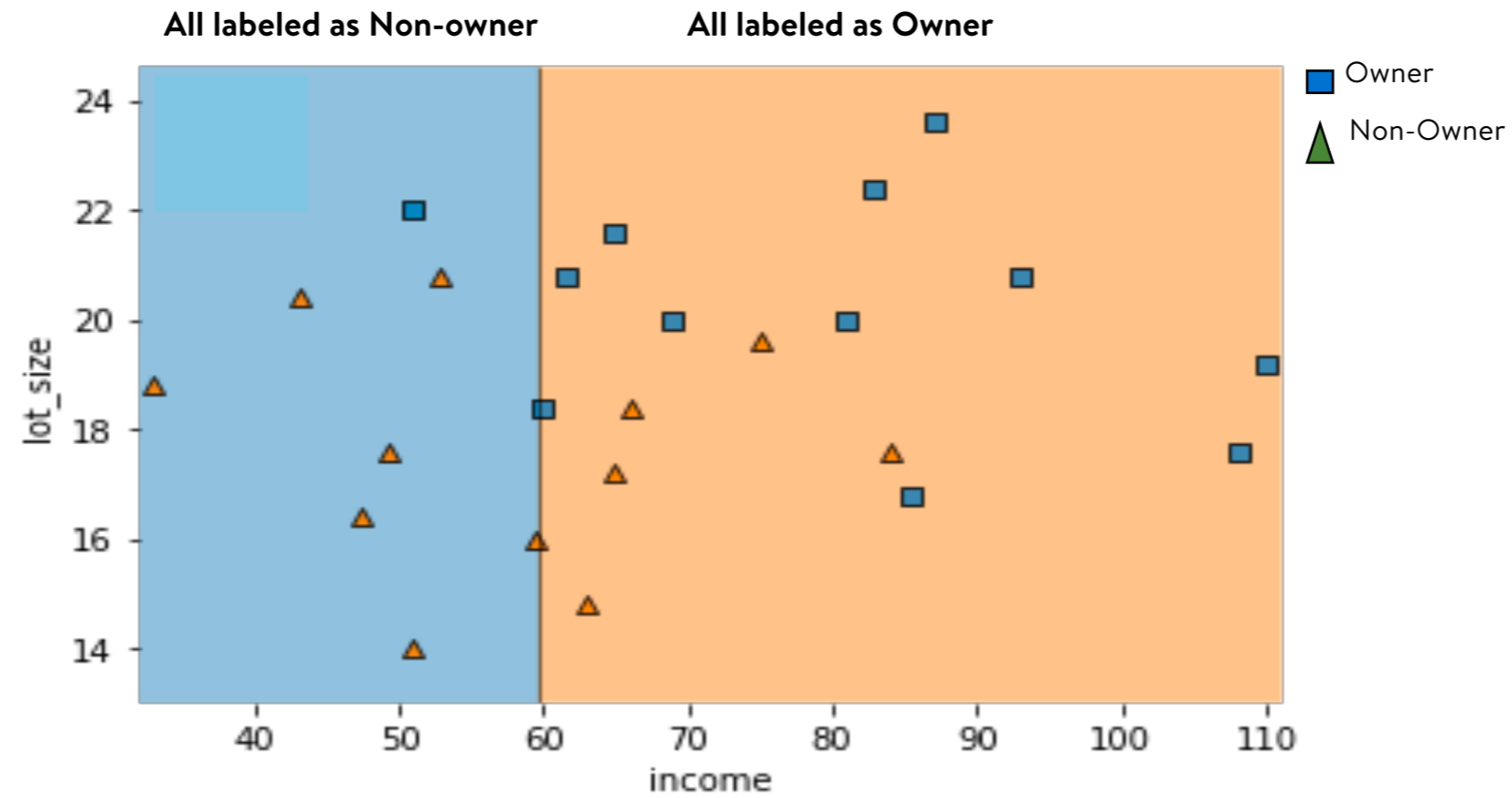
Depth = 1

- Gini index has a value of .21
- There are 8 observations in the node
- For income <= 59.7 bucket there are:
  - 1 Owner
  - 7 Non-Owners

- Gini index has a value of .43
- There are 16 observations in the node
- For income > 59.7 bucket there are:
  - 11 Owners
  - 5 Non- Owners

Proportion of Non-owner is higher than .5. observations in node classified as **non-owners**

Proportion of Owner is higher than .5. Observations classified as **owners**

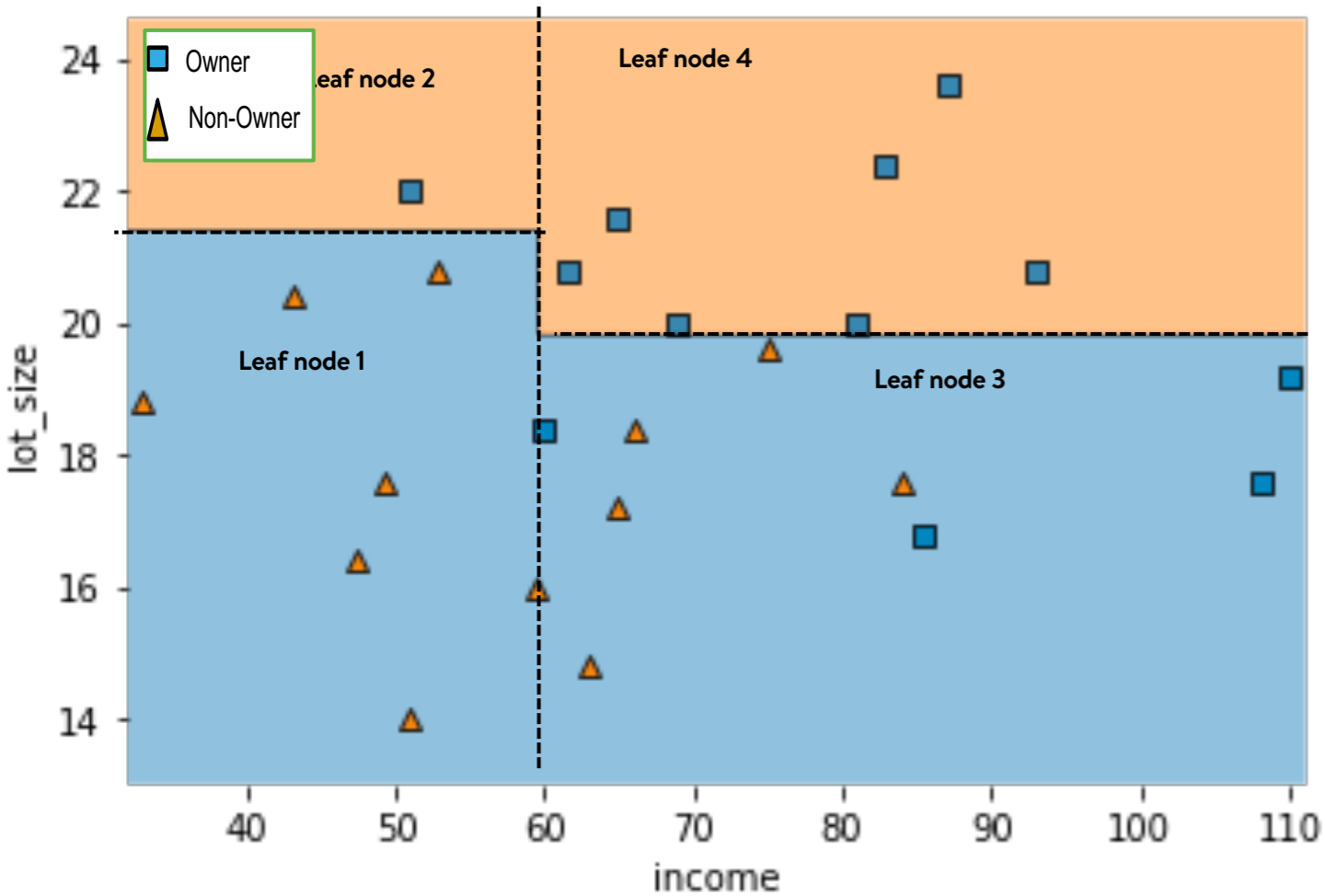


**Accuracy** = Proportion of correctly classified observations =  $\frac{\text{Correctly classified observations}}{\text{Total obs}} = \frac{7 + 11}{24} = 75\%$

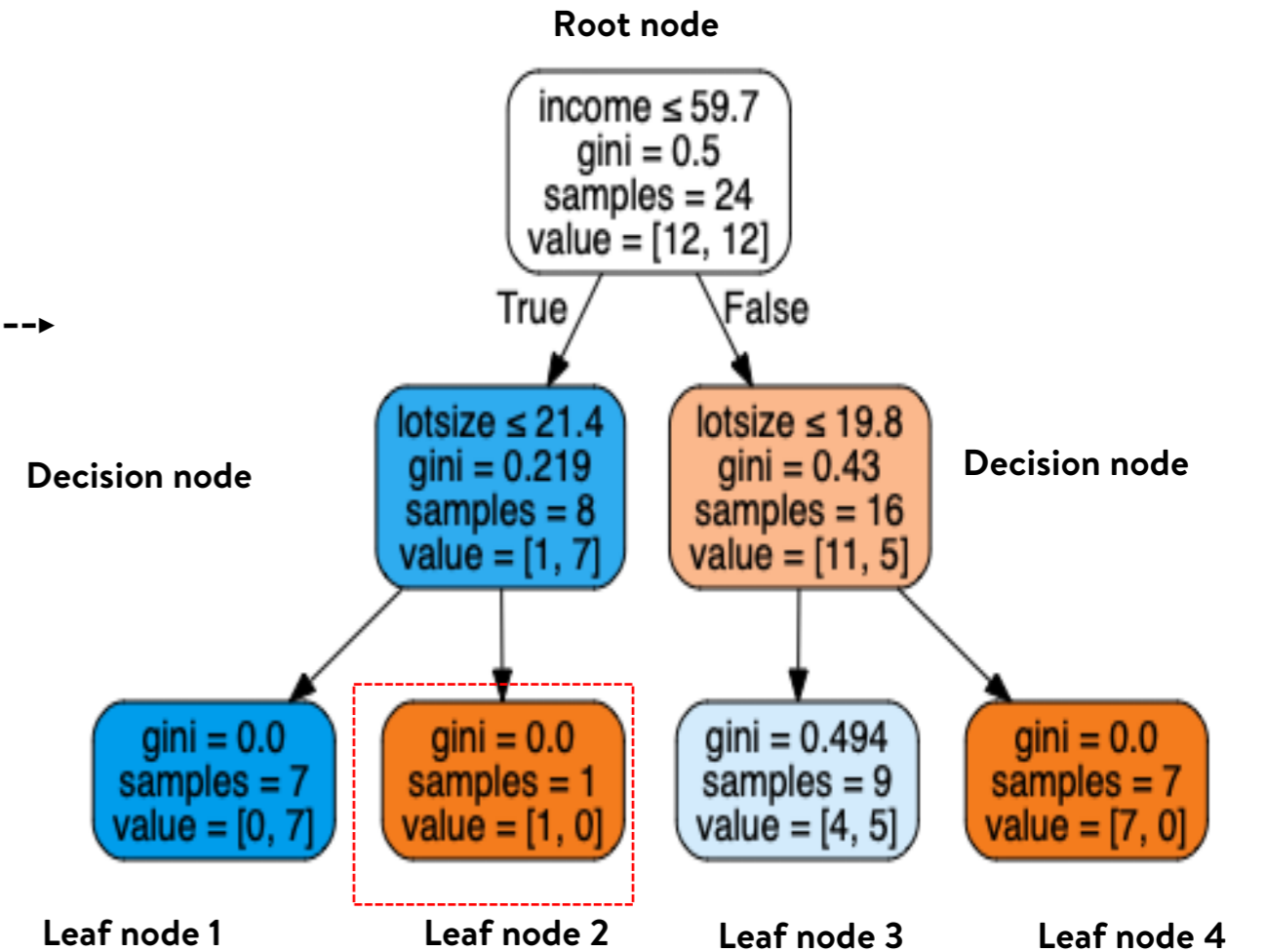
**Sensitivity** = Proportion of the owner's who were correctly classified =  $\frac{\text{Correctly classified owners}}{\text{Total owners}} = \frac{11}{12} = 92\%$

**Specificity** = Proportion of non-owner's who were correctly classified =  $\frac{\text{Correctly classified non-owners}}{\text{Total non-owners}} = \frac{7}{12} = 58\%$

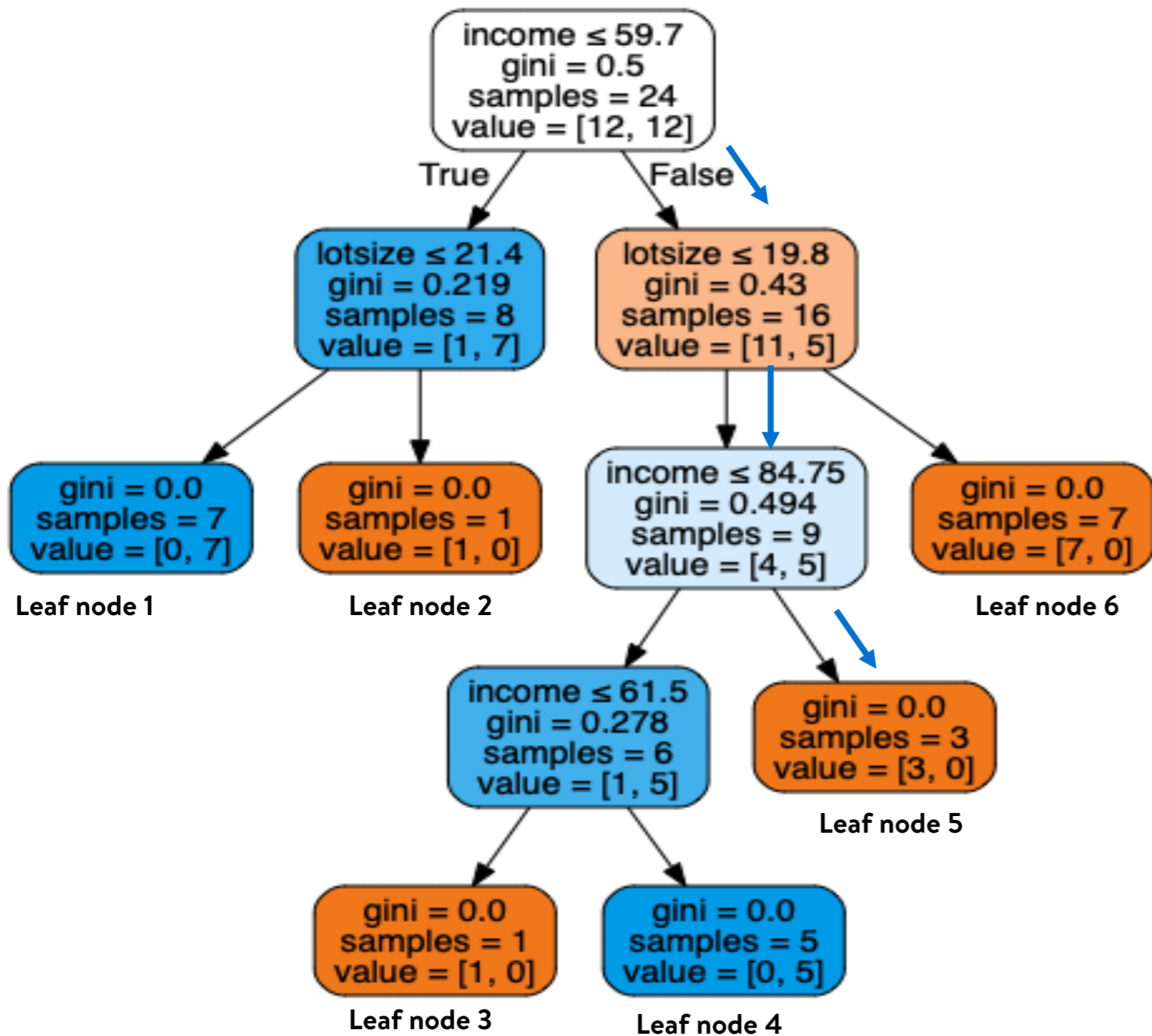
The algorithm 'greedily' searches the space for the next optimal split



Decision Tree representation



# The fully grown tree – interpreting the rules



**Classification rules:**

Leaf node1: If ( income <=59.7) & (Lot size is <=21.4) -> non-owner

Leaf node 6: If ( income > 59.7) & (Lot size is > 19.8) -> owner

Classifying an observation: household 7 ( income = 110.1, Lot size = 19.2)  
Classified as : Owner  
Actual class: Owner

Accuracy (the % of correctly classified observations): 100%

\* Note: There are a gamut of tree algorithms that are available (ID3, C4.5, CHAID, CART etc.)  
The one used here is closest to CART (Leo Breiman et al, 1984)

# Who should we offer a loan ?

- A relatively young American bank is growing rapidly in terms of overall customer acquisition.
- Majority of these are customers with varying sizes of relationship with the bank.
- The customer base of Asset customers is quite small, and the bank WANTS to grow this base rapidly to bring in more loan business.

5000 observations x 11 Predictors

Target variable

Age	Experience	Income	Family	CCAvg	Education	Mortgage	Securities Account	CD Account	Online	CreditCard	Personal Loan
25	1	49	4	1.60	1	0	1	0	0	0	0
45	19	34	3	1.50	1	0	1	0	0	0	0
39	15	11	1	1.00	1	0	0	0	0	0	0
35	9	100	1	2.70	2	0	0	0	0	0	0
35	8	45	4	1.00	2	0	0	0	0	1	0
37	13	29	4	0.40	2	155	0	0	1	0	0
53	27	72	2	1.50	2	0	0	0	1	0	0
50	24	22	1	0.30	3	0	0	0	0	1	0
35	10	81	3	0.60	2	104	0	0	1	0	0
<b>34</b>	<b>9</b>	<b>180</b>	<b>1</b>	<b>8.90</b>	<b>3</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>1</b>
65	39	105	4	2.40	3	0	0	0	0	0	0

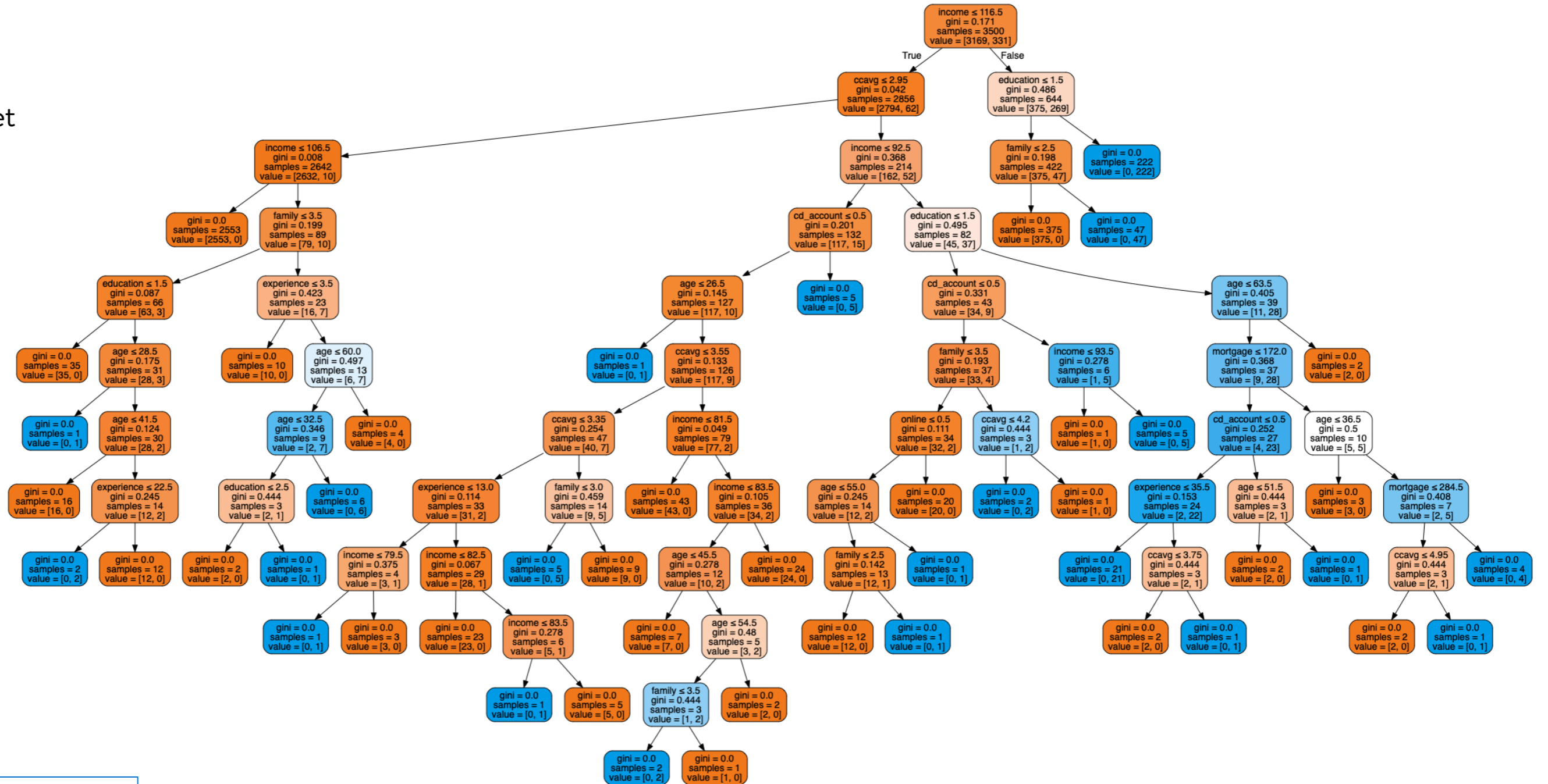
9 % of the customers are loan customers in this dataset

- We'll keep a random set of 70% of this data for calibrating (or training the tree)
- The remaining 30% for testing \* - to how does the model fare in the wild

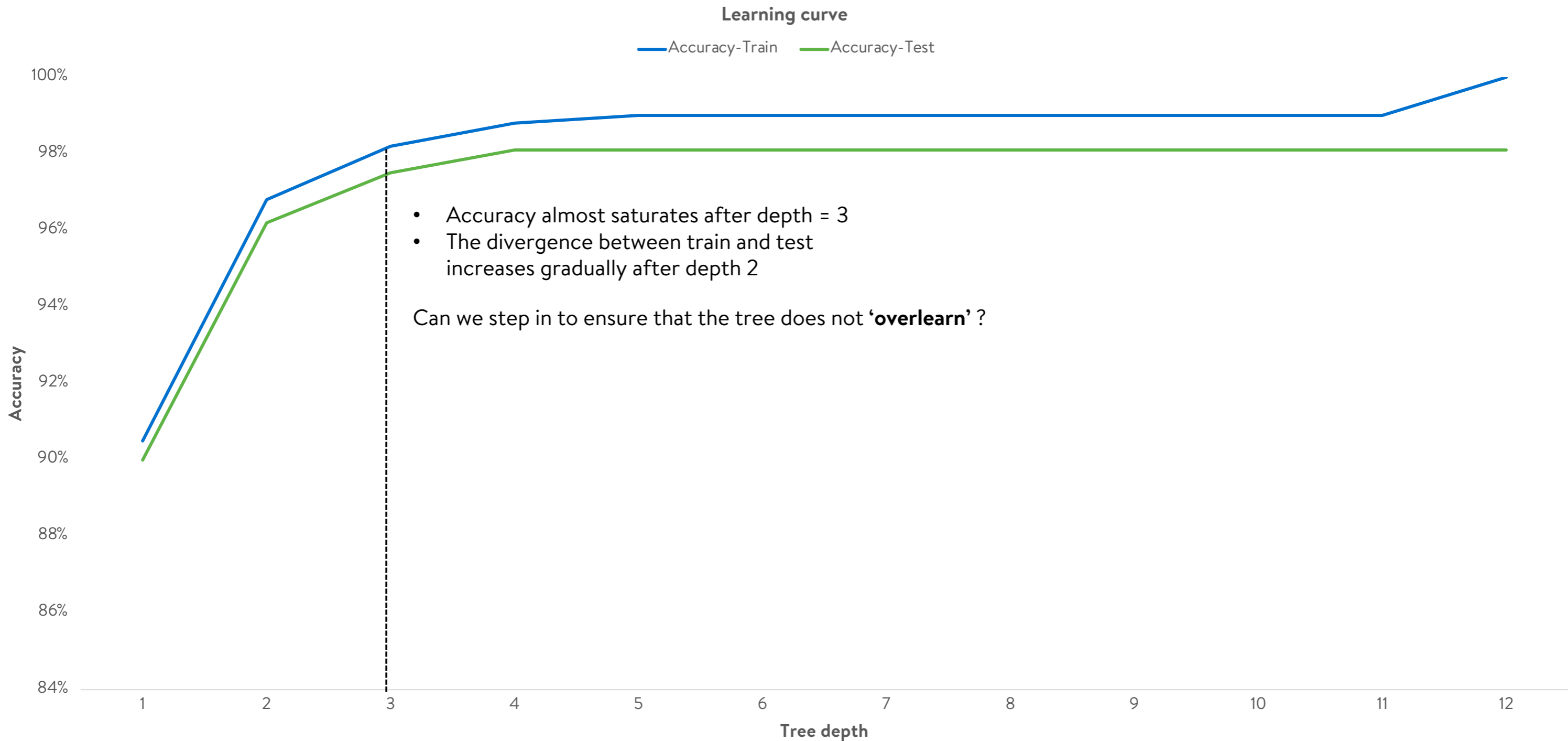


# Our fully grown tree for loan grant is complex !

\* For the training set



**Accuracy (Training): 100%**  
**Accuracy (Test): 98%**



 **prune**<sup>2</sup>  
/pru:n/

*verb*

gerund or present participle: **pruning**

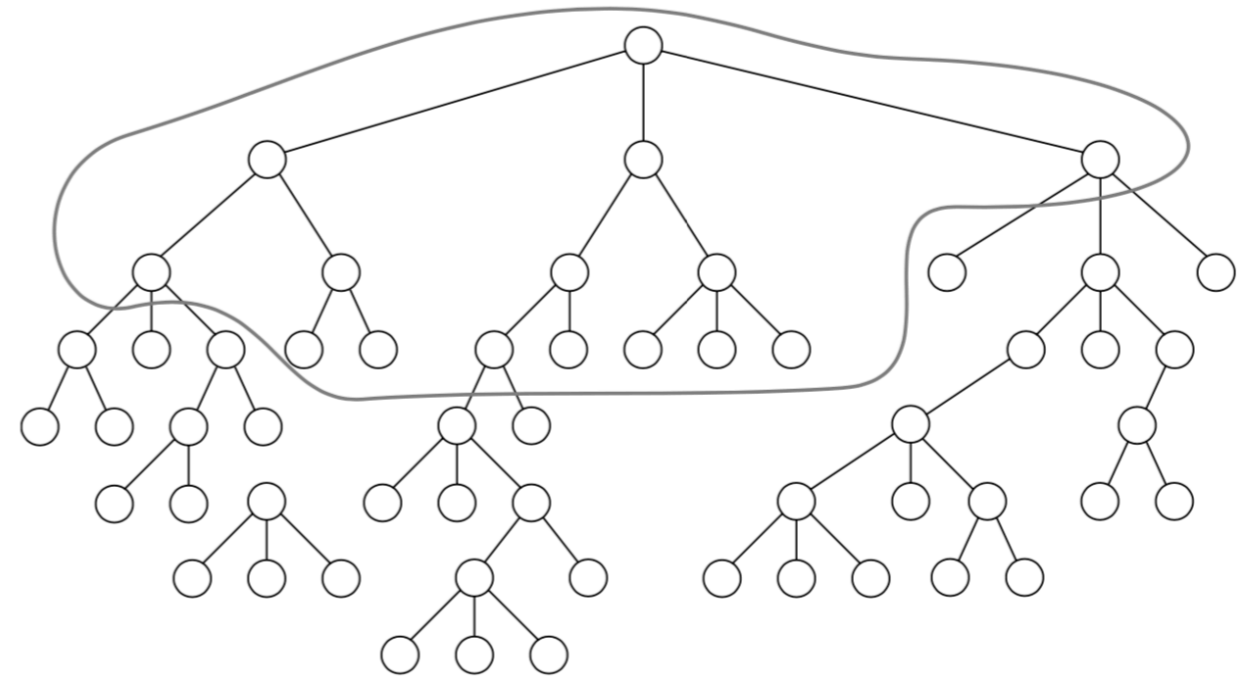
trim (a tree, shrub, or bush) by cutting away dead or overgrown branches or stems, especially to encourage growth.

"now is the time to prune roses"



\*Stihl Shop Greenburg – some rights reserved

**Inside a complex tree, there are simpler, more stable trees.**

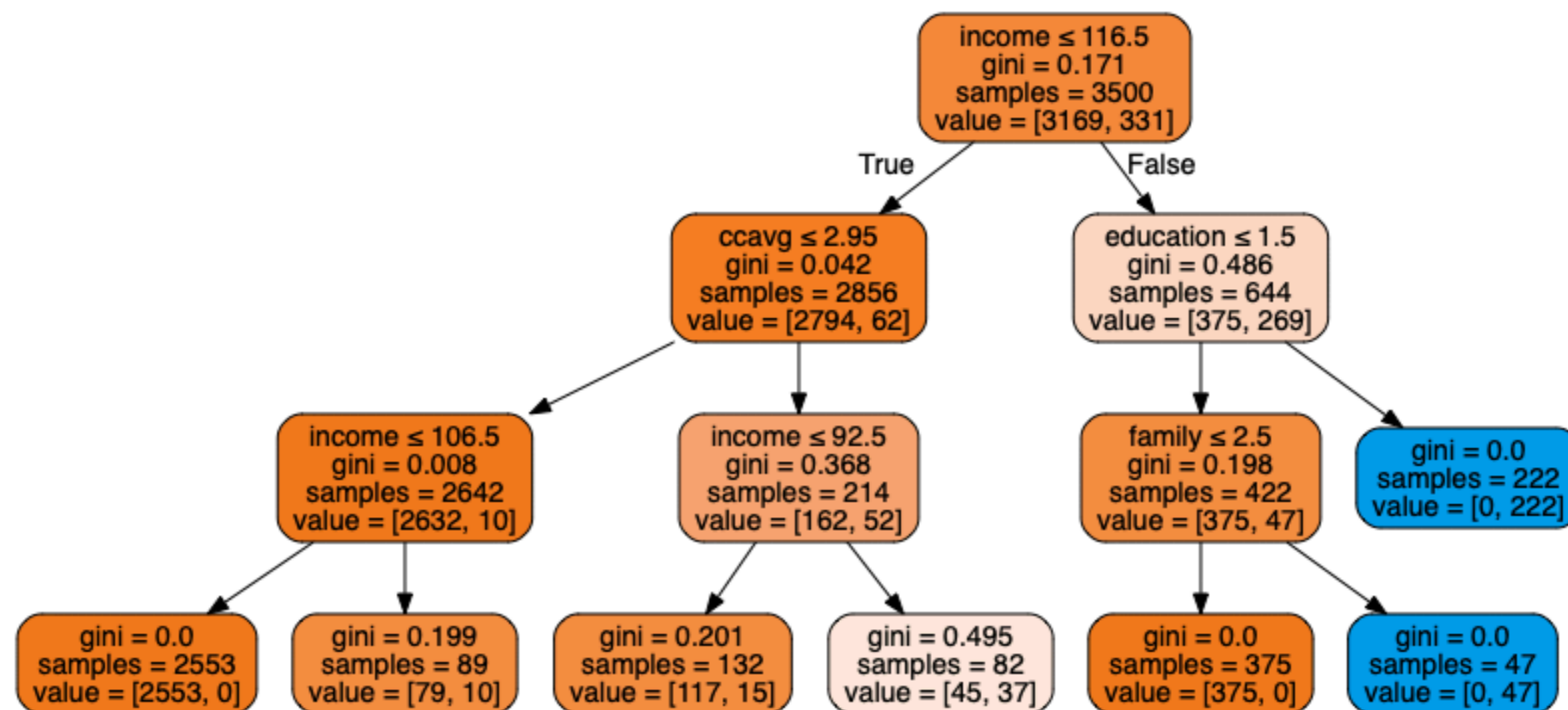


\*Depiction from 'Data mining techniques for marketing, sales, CRM', 3<sup>rd</sup> ed – Berry et al

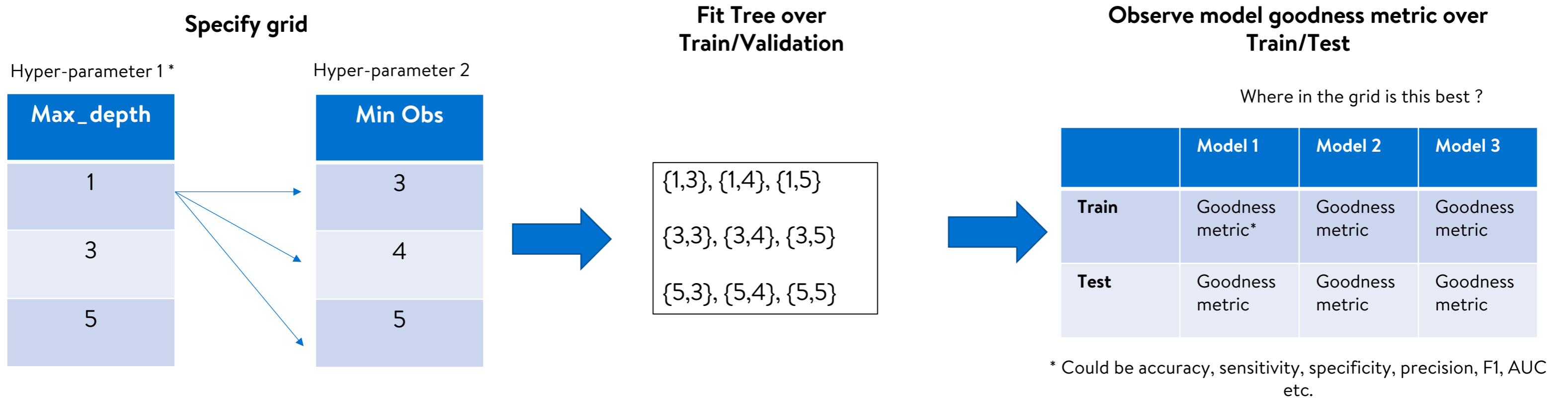
For example, we've set, maximum depth = 3 here. **The train & test accuracy is around 98% here**

There are multiple ways to prune beyond just this – for e.g. :

- **Minimum # of observations in a leaf node**
- **Min # of observations to continue splitting**
- **Min decrease in impurity**



## CARTESIAN GRID SEARCH



## OTHER SEARCH STRATEGIES

- Random Grid search : Searches the space of parameters randomly, not exhaustive. Computational cheaper.
- Bayesian grid search : Keep track of past evaluation results which they use to form a probabilistic model mapping hyperparameters to a probability of a score on the objective function:

\* <https://machinelearningmastery.com/difference-between-a-parameter-and-a-hyperparameter/>

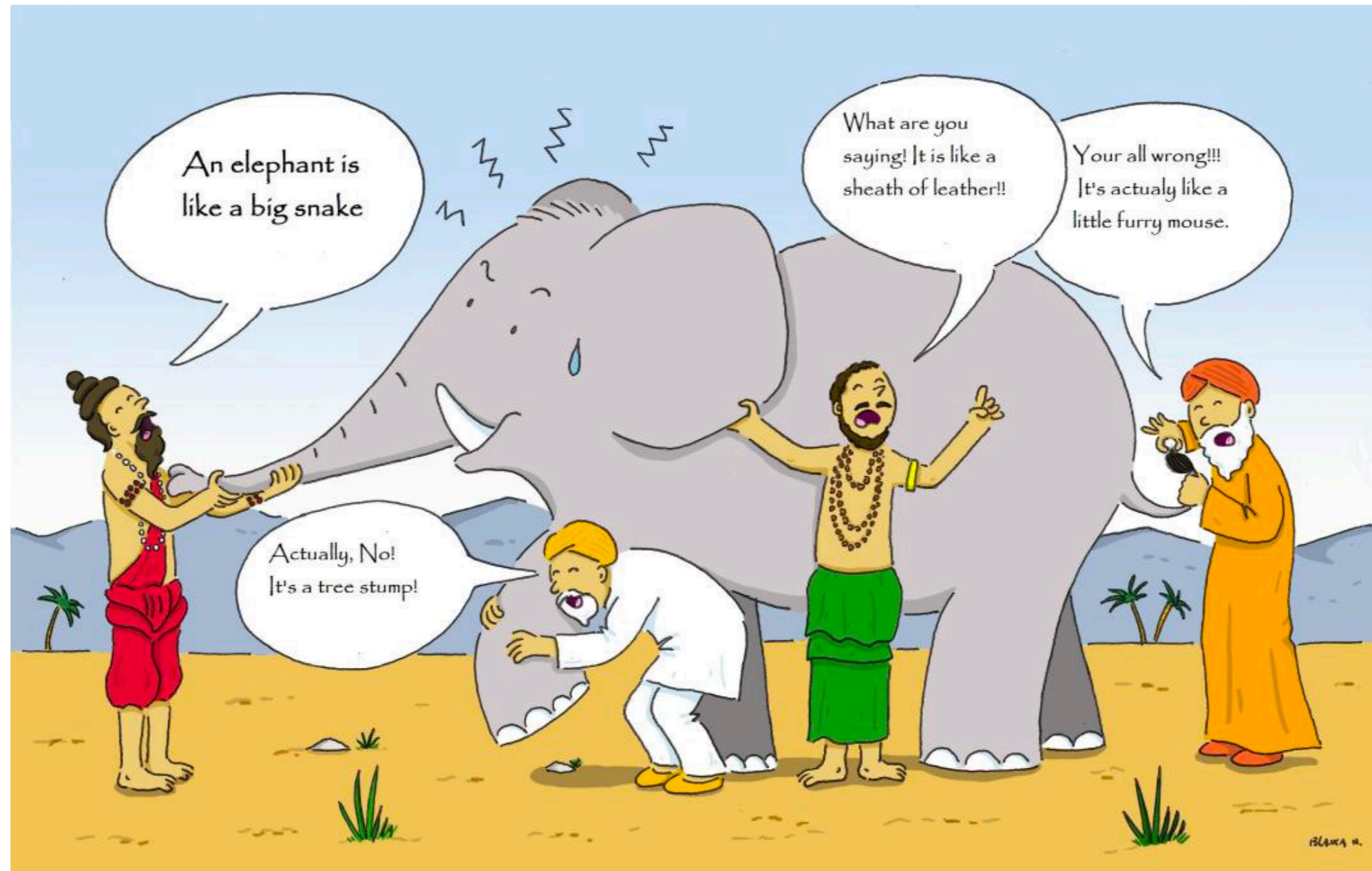
## Advantages :

- Easy to interpret and visualize.
- Can model complex patterns quite well.
- Needs limited assumptions – mainly data driven
- Can be used for classification as well as regression problems

## Disadvantages:

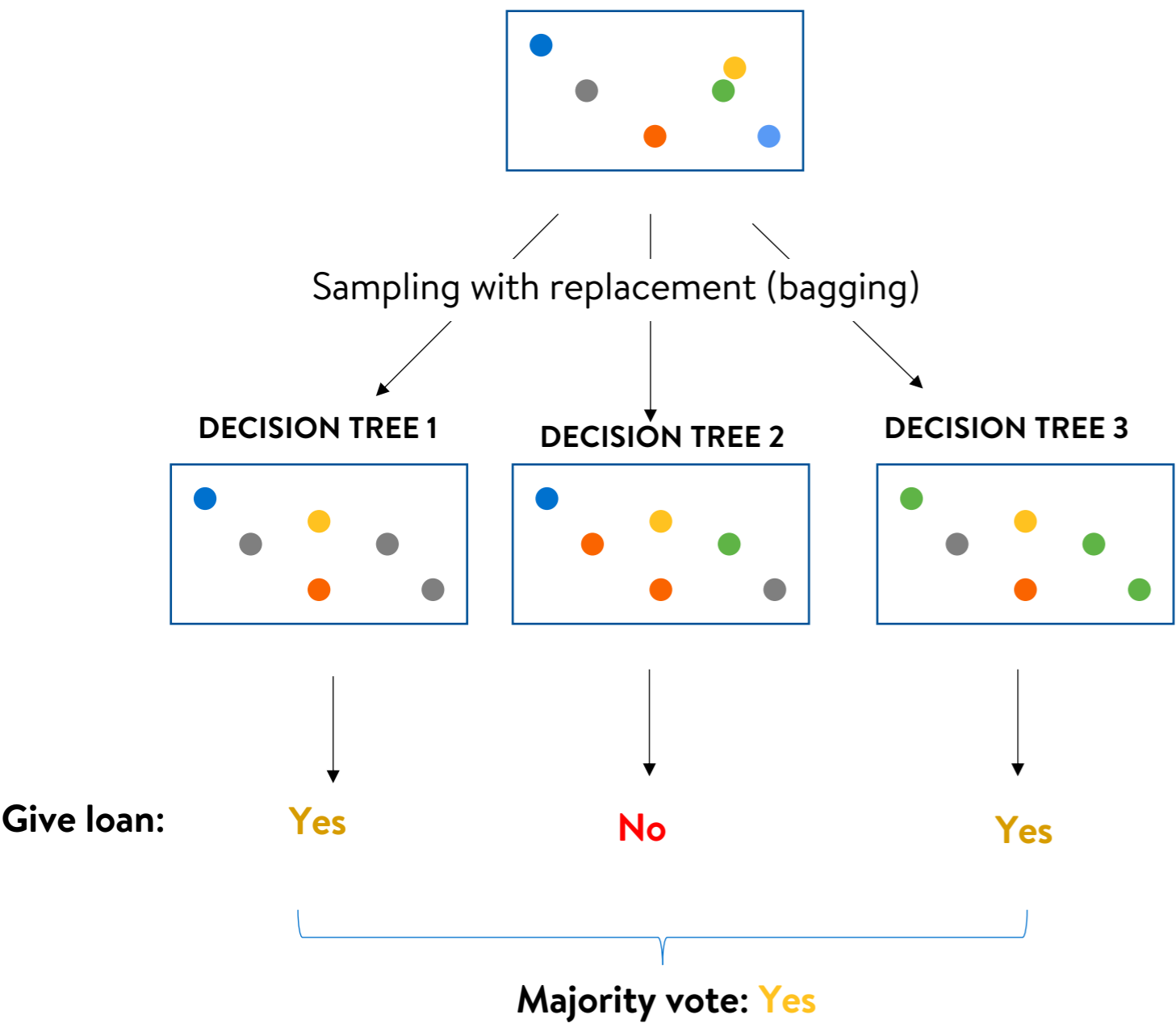
- High tendency to overfit to the data used for training
- Small variation(or variance) in data can result in the different decision tree.

# The blind men and the elephant

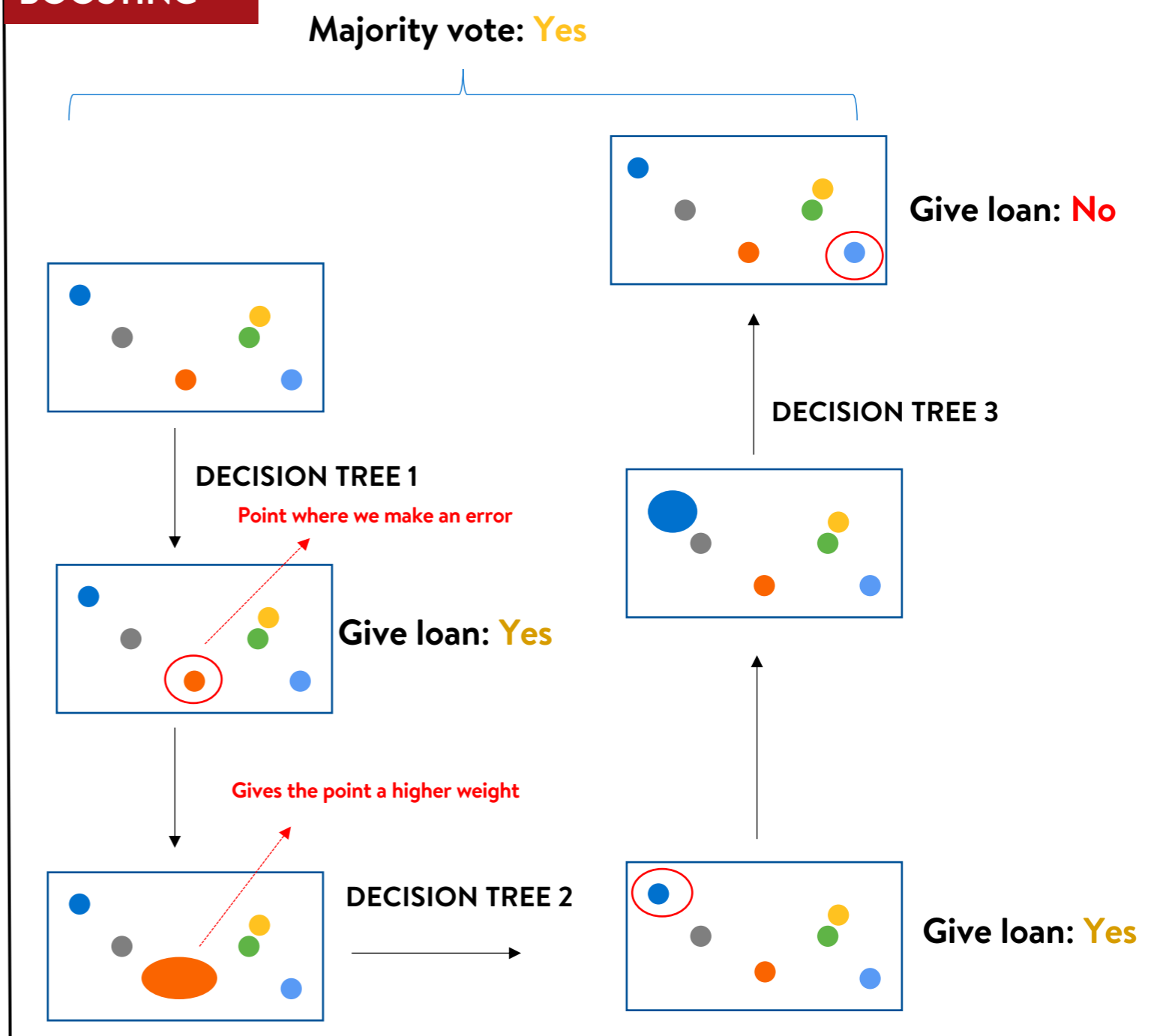


\* <https://medium.com/diogo-menezes-borges/ensemble-learning-when-everybody-takes-a-guess-i-guess-ec35f6cb4600>

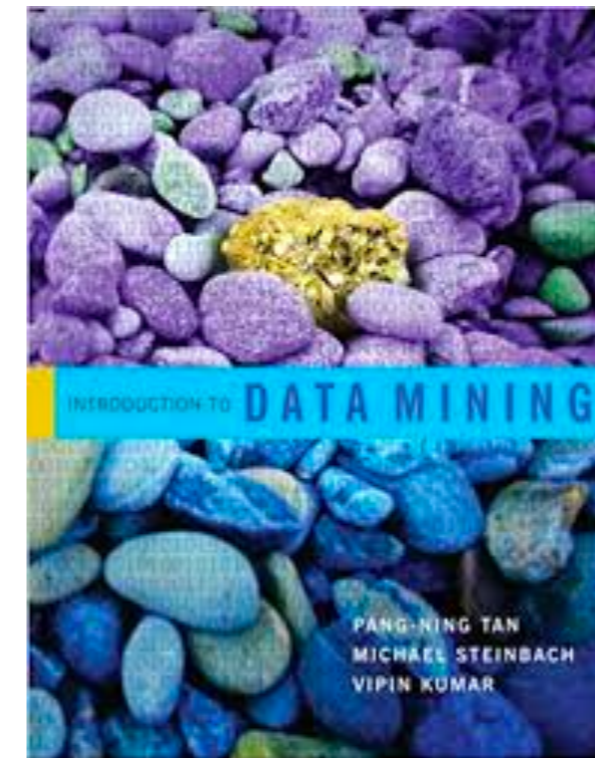
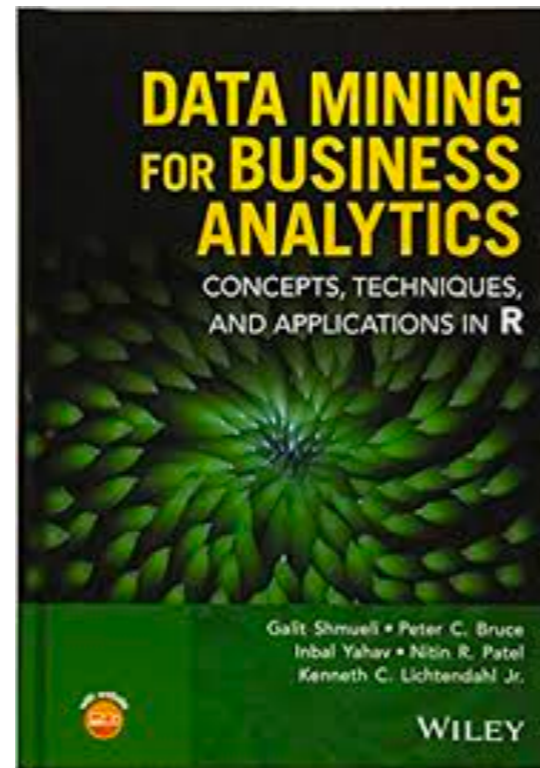
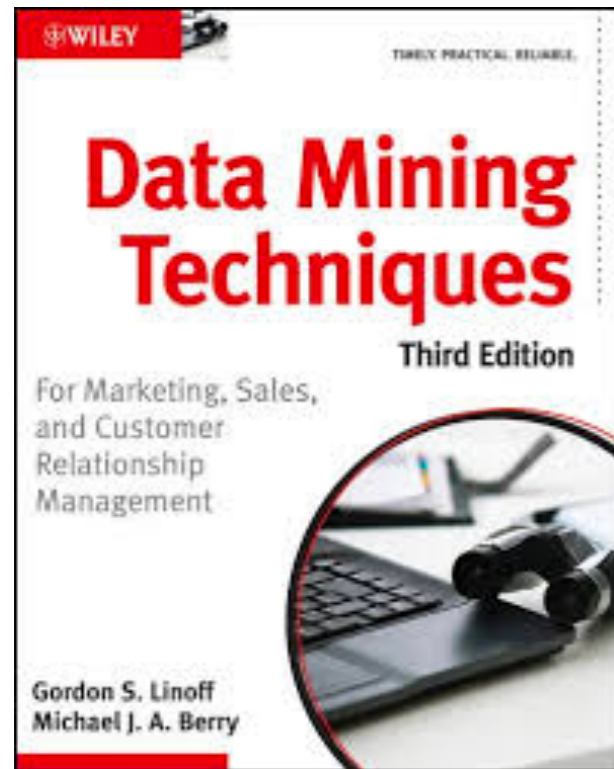
## BAGGING



## BOOSTING







Git repo for session data + code - [https://gecgithub01.walmart.com/smisra/TechByte\\_DTree\\_session](https://gecgithub01.walmart.com/smisra/TechByte_DTree_session)

# WALMART USE CASE

## Omni-Channel Retail Fraud Detection: Refunds, Cancellations, Discounts, Collusions



### Mission

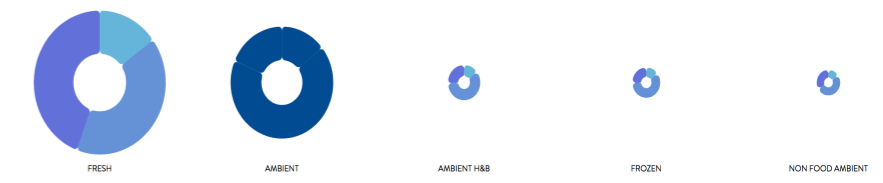
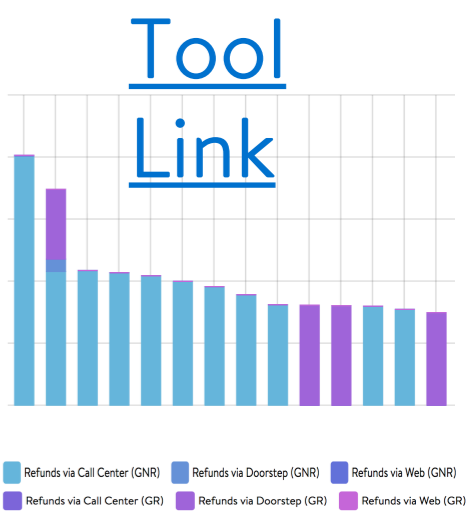
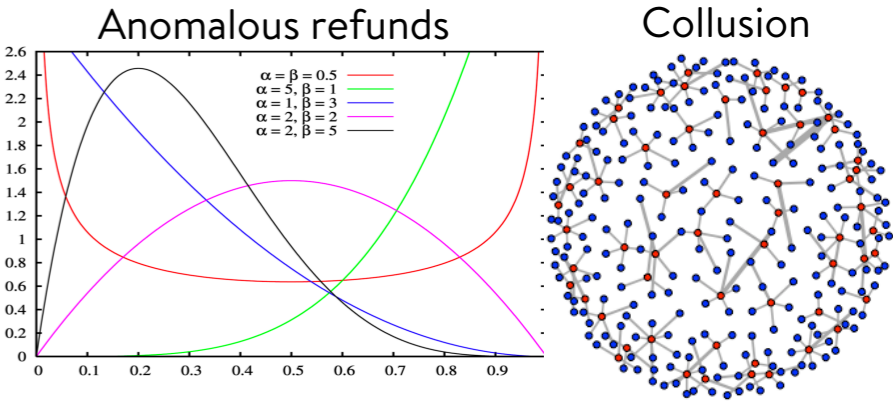
**Detect** and **Prevent** refund and cancellation fraud by customers and collusion with drivers and employees, and identify process improvement opportunities.



### Features

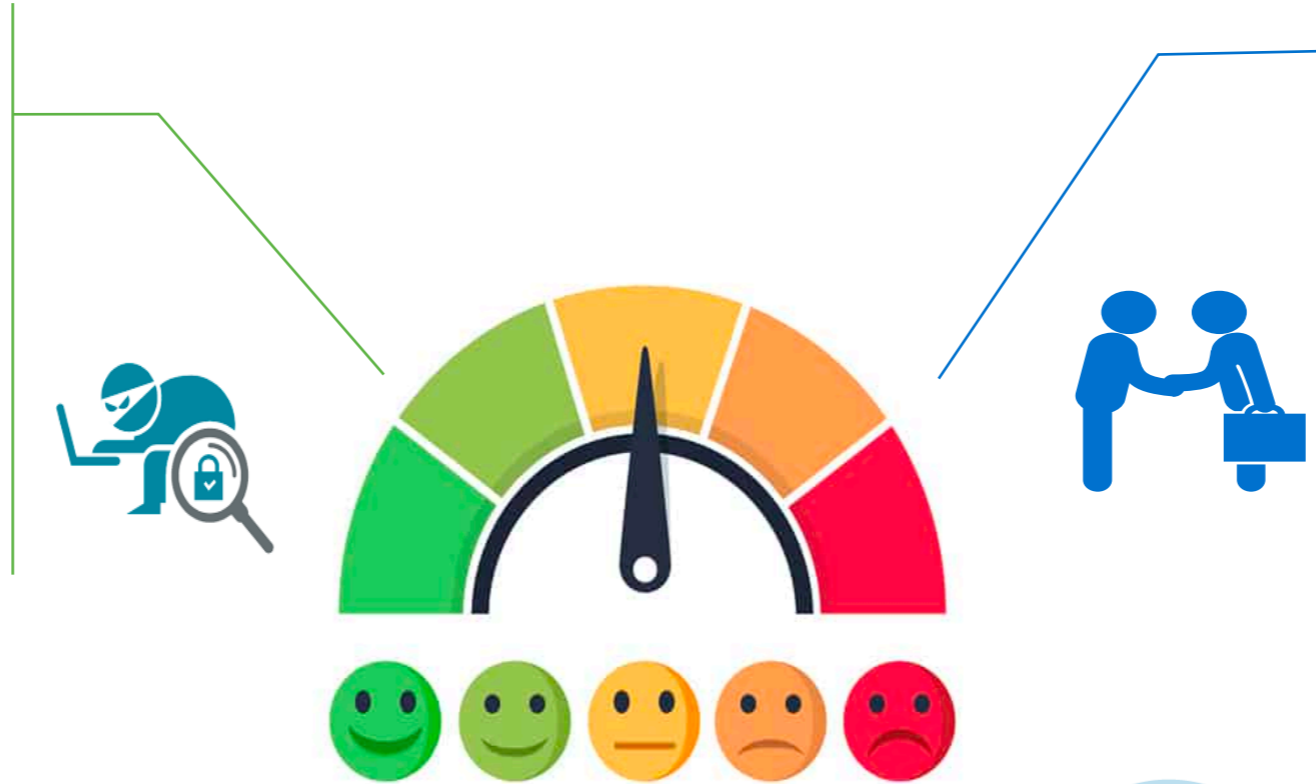
#### KEY BENEFITS

- ❖ Highlights instances of cancellation and refund **abuse** by customers
- ❖ Identifies **collusions** among customers, drivers, store associates
- ❖ Identifies cases of colleague **discount** & reselling abuse
- ❖ Risk assessment by geographic locations & merchandising hierarchy
- ❖ **Prioritizes** cases to take appropriate action by ML generated risk scores
- ❖ Discovers common fraud **modus operandi** to mitigate future risk



## Fraud Risk KPIs

- High Refund Amount/Frequency
- High Risky Cancellations
- High Refund Rate
- GNR Refunds/proportion
- High refunds through Web, Call Center or Doorstep
- Recency

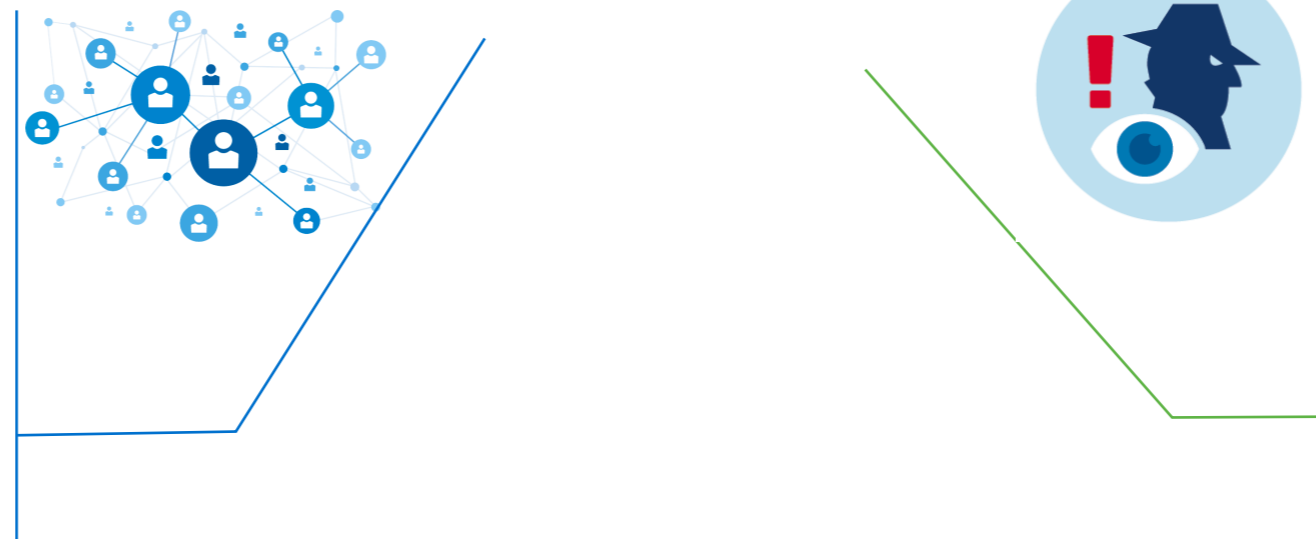


## Multi-party Collusion

- Always refunding with same driver with high refund amount
- Always cancelling orders by same employee after Pick complete

## Risk by Association

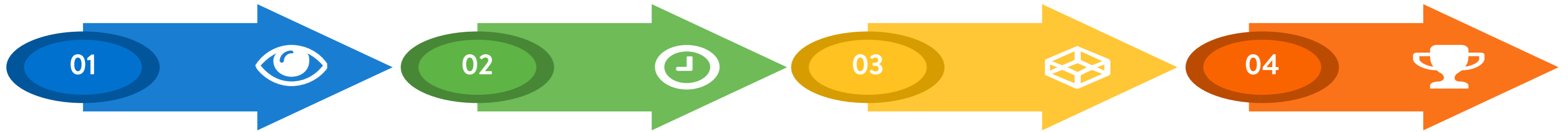
- Share of High-refunding Stores, cities, postcodes
- Refunds made in high value items and in high risk stores



## Suspicious Behavior

- Repeated refunds of same item
- Recent spike in refunds
- Doesn't return Damaged Items or Unwanted Substitutes
- Refunds at a higher price

# How risk scores are computed?



## Metric Aggregation

Customer level aggregated refund, cancellations, collusion, recency

~60 Features

## Feature Selection & Cleaning

Remove highly correlated variables to remove multi-collinearity, normalize the metrics

~40 Features

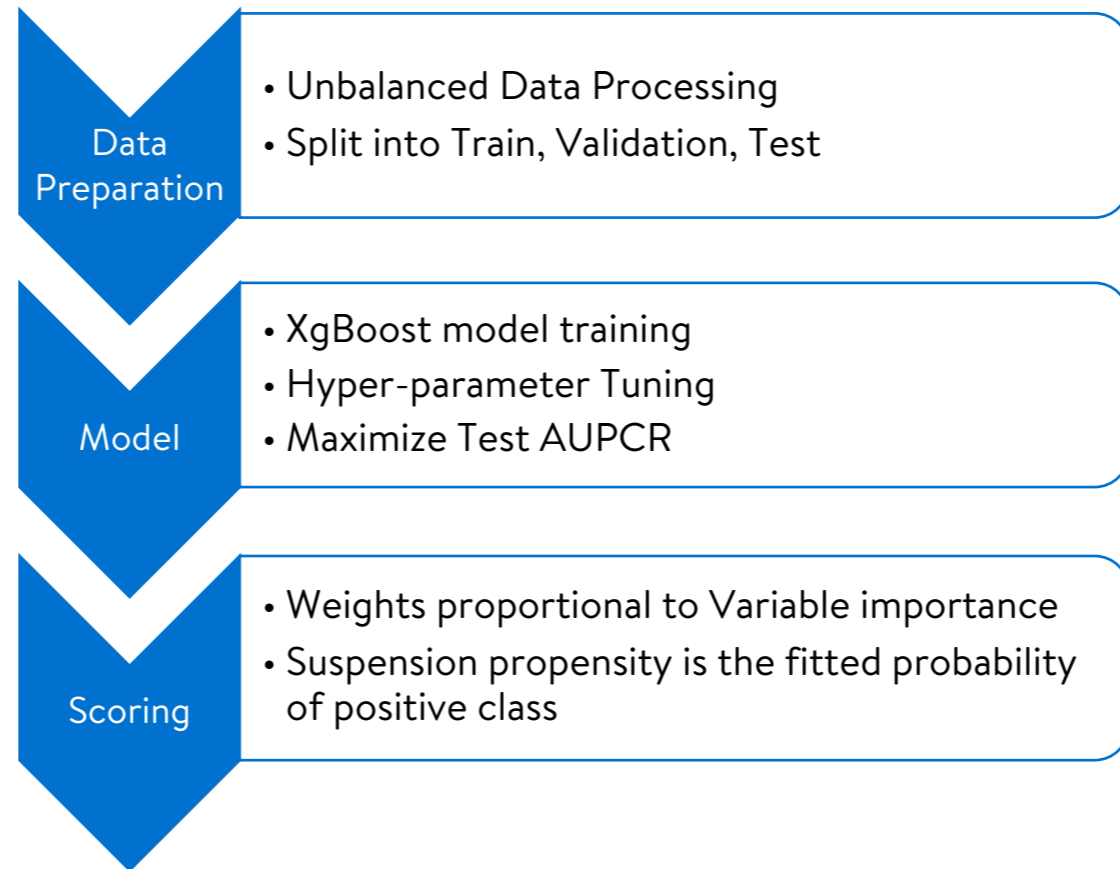
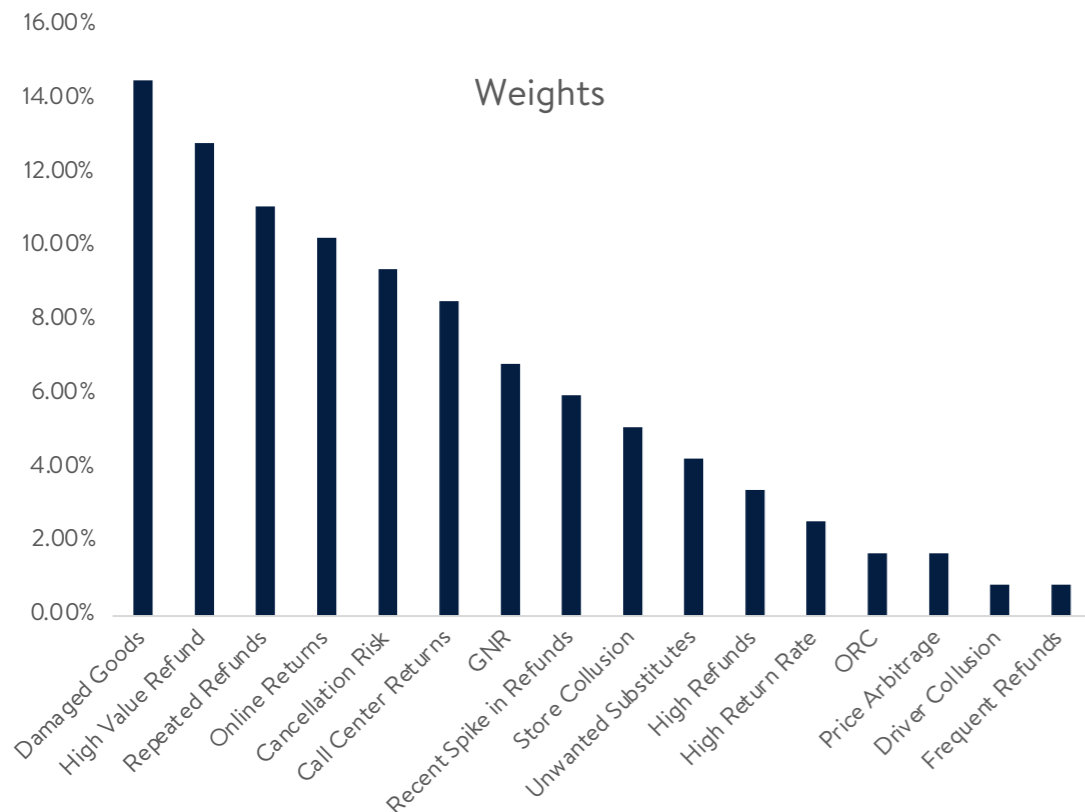
## Choice of Weights

Risk score is weighted aggregate of metrics  
Weights are based on

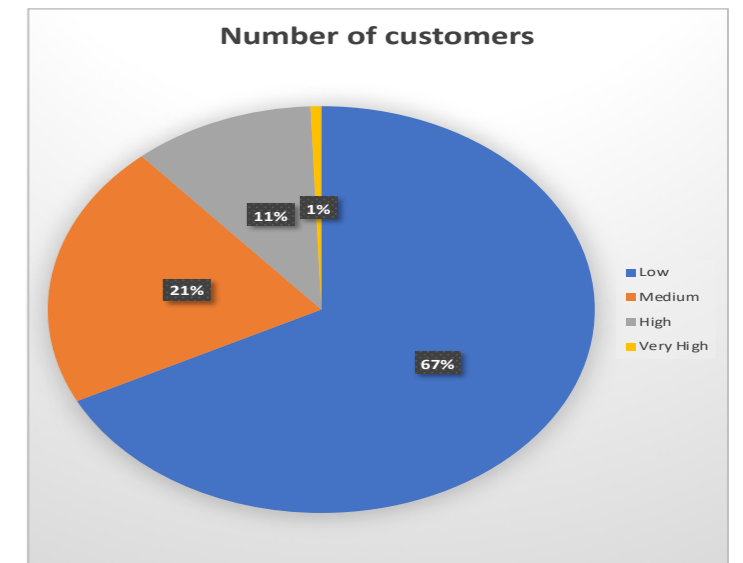
1. Precision
2. Suspension propensity Importance

## Risk Buckets and Reasons

Risk buckets are dynamically chosen from the Risk Scores. Features contributing significantly to the score are the risk reasons.



Metric weights are inversely proportional to their precision

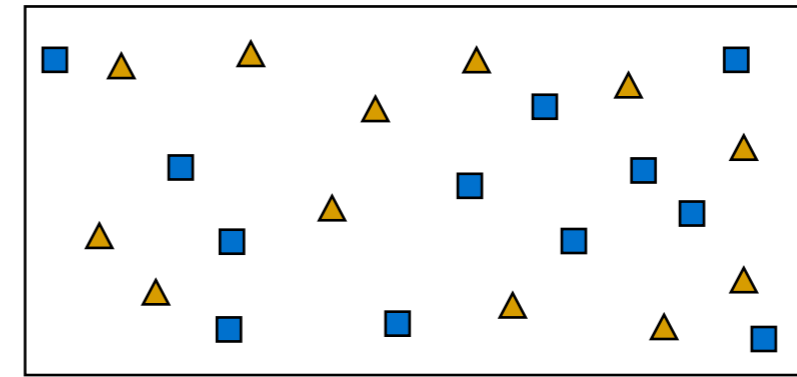
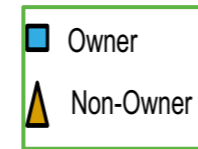
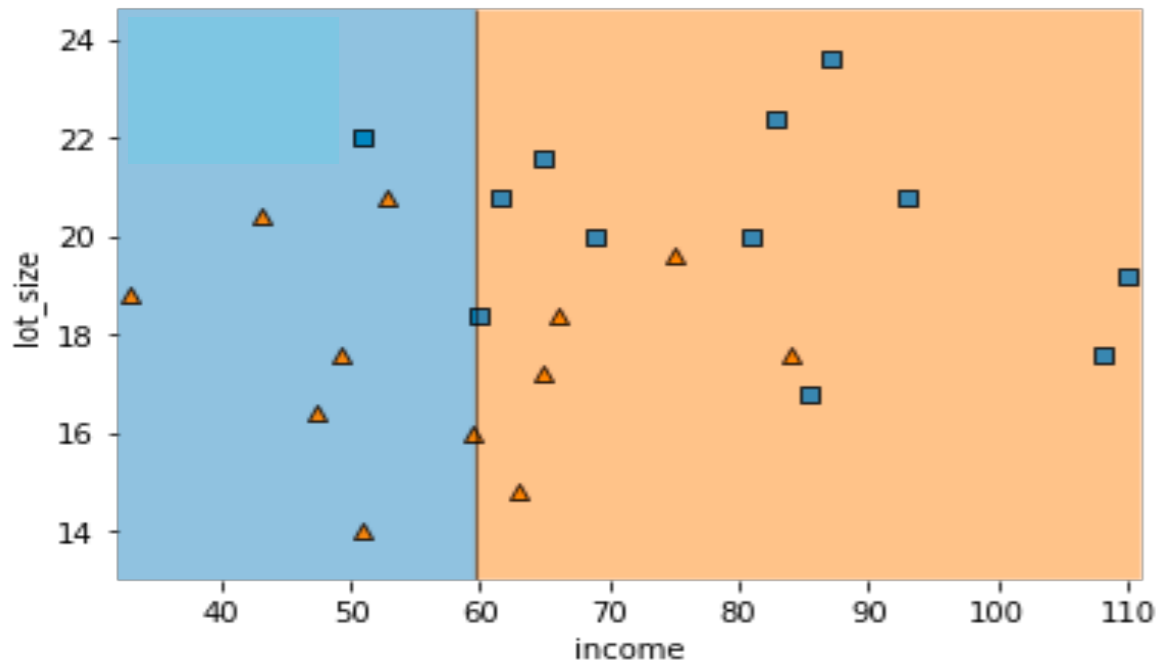


**Thank you !**

**Appendix:**

## Tree stump – A simple tree with one split

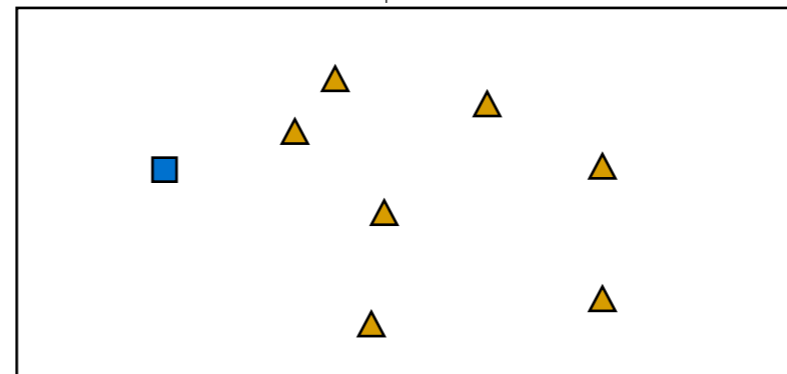
Decision node



Is income  $\leq 59.7$

YES

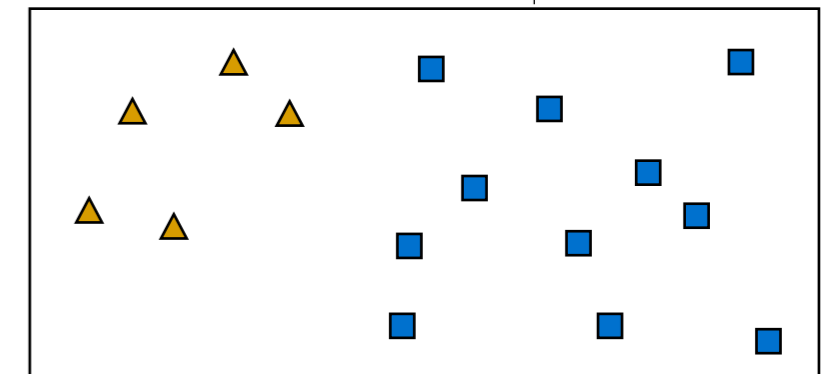
NO



For income  $\leq 59.7$   
bucket there are:

- 7 Non-owners
- 1 Owner

**Proportion of Non-owner is higher than .5.** All observations classified as **non-owners**



For income  $> 59.7$   
bucket there are:

- 5 Non-owners
- 11 Owners

**Proportion of Owner is higher than .5.** All observations classified as **owners**





This split happens randomly – that is, the general characteristics\* of the predictors and the target variable remains the same across both these sets.

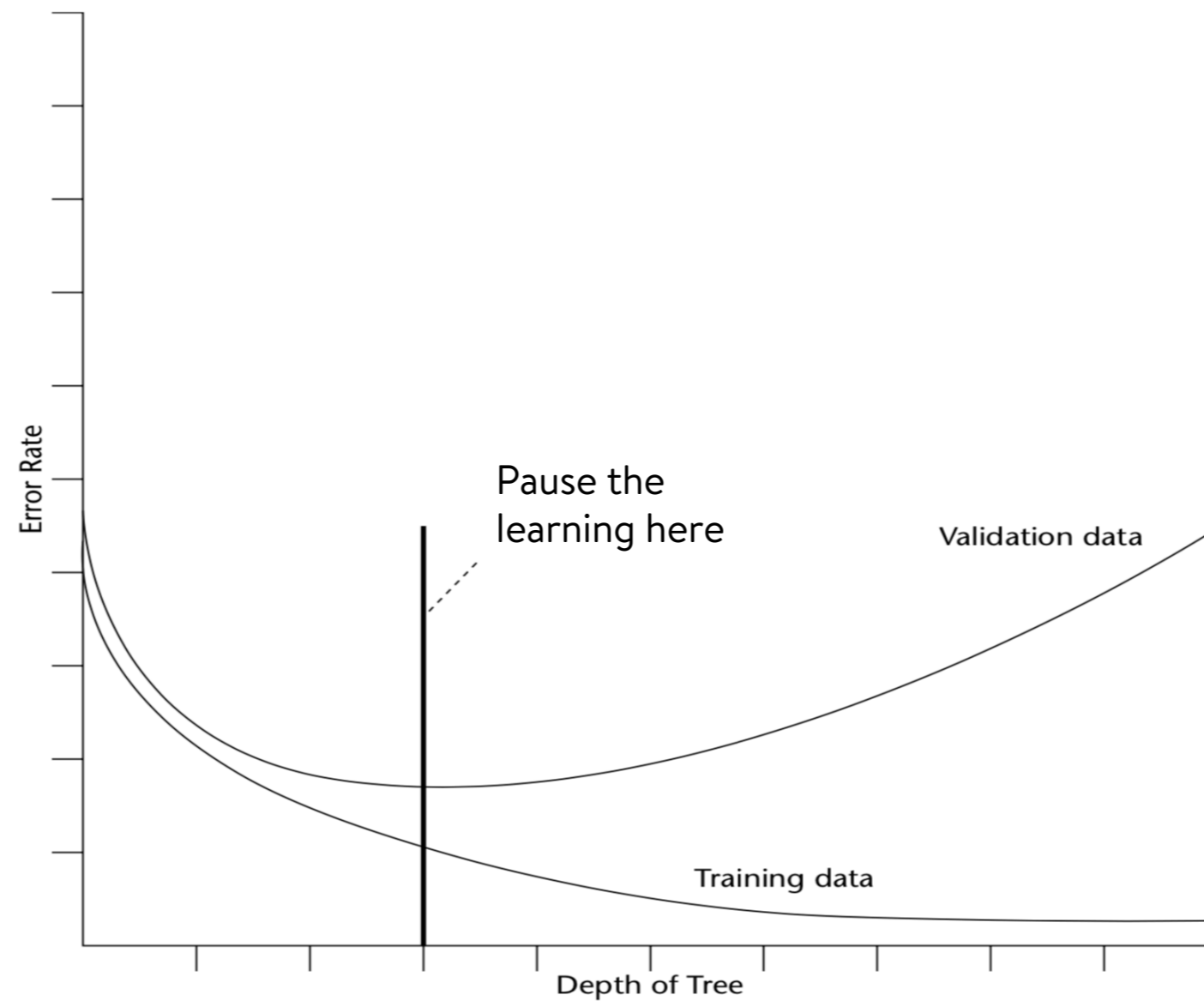


- Calibrate model
- Evaluate model on training data



- How does the model fare in the in the wild
- Does the model generalize

In general a trade-off : Cost = Error + **cost complexity** \* number of leaves in a tree



\*Depiction from 'Data mining techniques for marketing, sales, CRM', 3<sup>rd</sup> ed – Berry et al

# Architecture Diagram for Risk Score Model

