

Classification : Basic Concepts and Techniques (1)

I. Introduction to Classification

↳ What is Classification

Classification is a supervised learning technique that assigns predefined labels to data based on a training dataset.

- Training Data: A collection of records with known class label.
- Attribute (x): Also called predictors, independent variables or features.
- Class Labels (y): Also called responses, dependent variables, or outputs.

↳ General Approach for Classification

1. Model Construction

Learn a classification model from the training dataset.

• Output: Decision rules, decision trees, or mathematical models

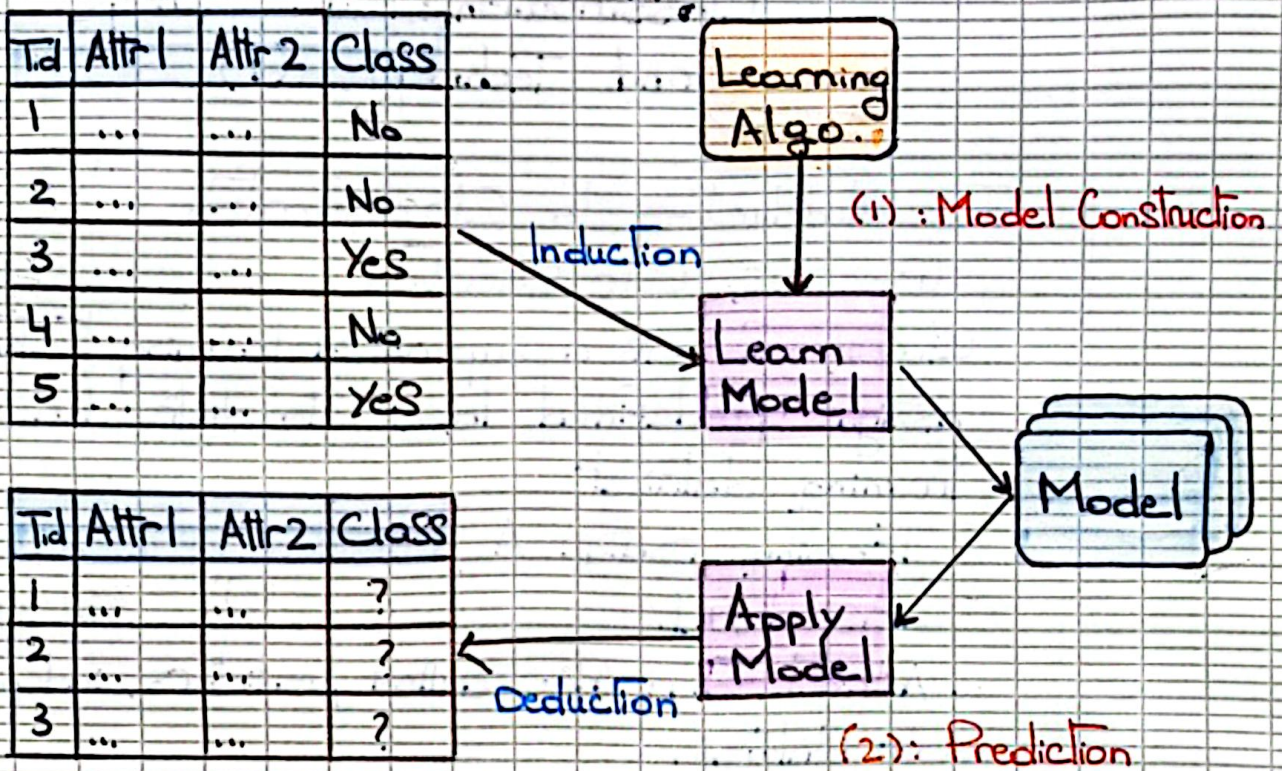
• Example:

Rule: if age > 30 and income > \$50K, then "approved loan"

2. Model Prediction

• Use the model to classify new data points (test)

• Measure accuracy: $Acc = \frac{\text{Correct pred}}{\text{Total pred}}$



II. Classification Techniques

↳ Overview

Classification techniques are broadly divided into two categories:

- 1) **Base Classifiers**: Individual algorithms for classi.
- 2) **Ensemble Classifiers**: Combine multiple base classifiers to improve accuracy

↳ Base Classifiers

1. **Decision Tree-Based Methods**

• **Description**: Models data as a tree structure with nodes representing decisions and branches indicating outcomes

• **Key Algorithm Examples**: ID3, C4.5, CART

Example:

↳ Decision: "Is income > \$30K"

↳ Outcomes: "Loan Approved" or "Loan Rejected"

2. Rule-based Methods

3. Nearest-neighbor

4. Neural Networks, Deep Neural Nets

5. Naive Bayes

6. Support Vector Machines (SVM)

↳ Ensemble Classifiers

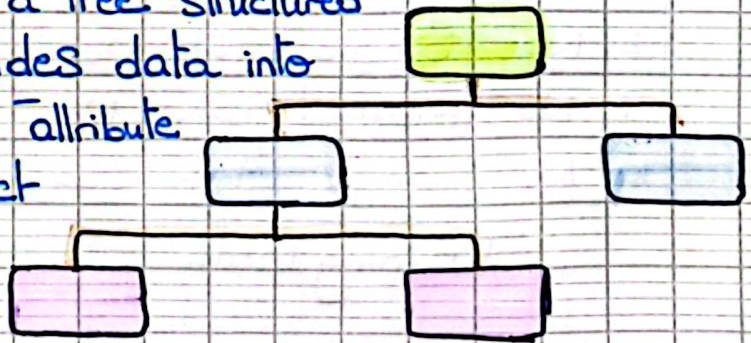
1. Boosting

2. Bagging (Bootstrap Aggregation)

III. Decision Trees and their Splitting Criteria

↳ What are Decision Trees?

A decision tree is a tree structured classifier that divides data into subsets based on attribute values, helping predict outcomes



↳ Nodes:

↳ Internal / split nodes represent attributes

↳ Branches represent decision rules

↳ Decision / Leaf nodes represent class labels

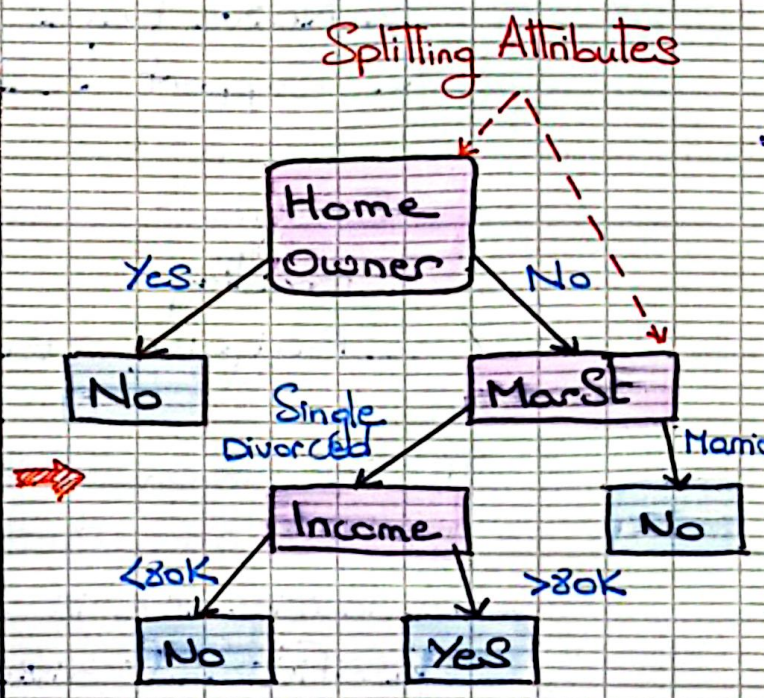
↳ Root Node: Starting point of the tree

Process:

- 1) Start from the root.
- 2) Split data based on attributes.
- 3) Continue until leaf nodes contain homogeneous classes or meet stopping criteria.

Example of a Decision Tree

ID	Home Owner	Marital Status	Income	Class
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	35K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Training Data

Model: Decision Tree

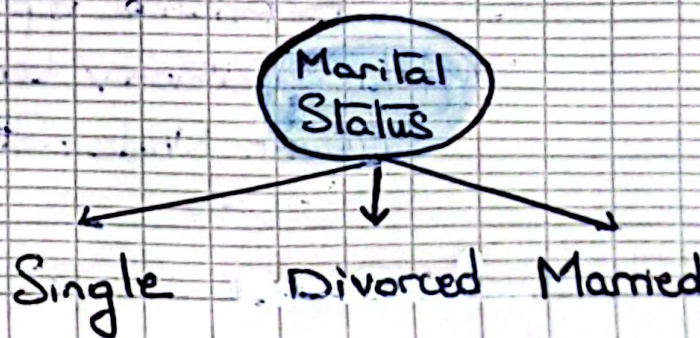
There could be more than one tree that fits the same data.

Splitting Criteria

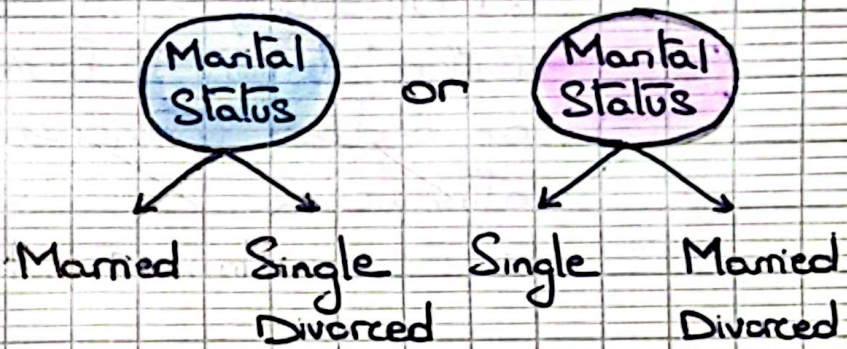
Nominal Attributes

Multi-way Split

Create a branch for each distinct value

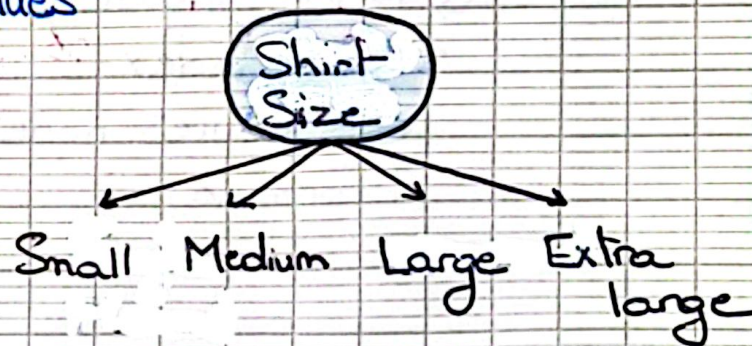


→ **Binary Split**
 Divides values into two subsets

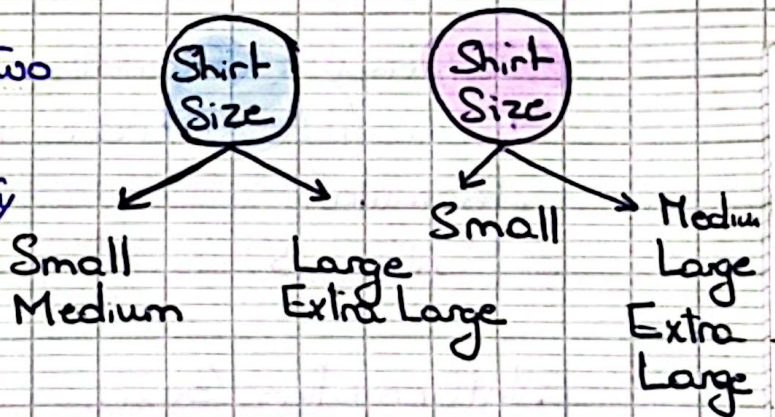


2. **Ordinal Attributes**
 Maintains order among values

→ **Multi-way Split**
 Use as many partition as distinct values

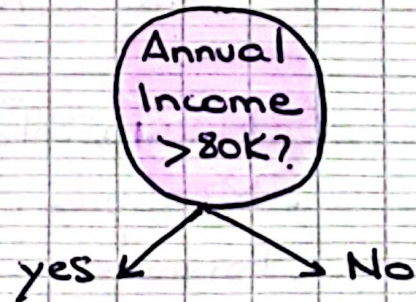


→ **Binary Split**
 Divides values into two subsets
 Preserve order property among attribute values

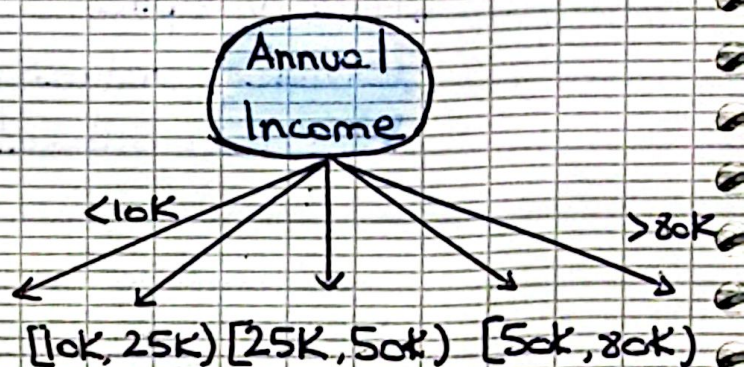


3. **Continuous Attributes**

Typically requires identifying
 → **Binary Split** ($A < v$) or ($A \geq v$)
 Consider all possible splits and finds the best cut
 Can be more compute intensive.



↳ Multi-way Split
 Discretization to form
 an ordinal categorical
 attribute



↳ Node Impurity and Splitting Criteria

1. Gini Index

↳ Measures node impurity, based on class probabilities

• Formula:

$$\text{Gini}(t) = 1 - \sum_{i=1}^n p_i^2$$

Proportion of instances in class i

2. Entropy (Information Gain)

• Measures disorder in a dataset

• Formula:

$$\text{Entropy}(S) = - \sum_{i=1}^n p_i \log_2(p_i)$$

• Information Gain:

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \cdot \text{Entropy}(S_v)$$

3. Misclassification Error

• Measures the proportion of misclassified records

• Formula:

$$\text{Error}(t) = 1 - \max(p_i)$$

→ Finding the Best Split

1. Compute impurity (e.g., Entropy, Gini Index) before splitting (B)
2. Compute impurity after splitting (A)
 - Weighted average of child node impurities
3. Compute Information Gain:

$$\text{Gain} = B - A$$
4. Choose the attribute with the highest gain

→ Example using Entropy

1. Compute impurity before splitting

$$B = \{Y^6, N^4\}$$

$$\begin{aligned} E(B) &= -p(Y) \log_2 p(Y) \\ &\quad - p(N) \log_2 p(N) \\ &= -6/10 \log_2 (6/10) \\ &= 0.44 + 0.53 = 0.97 \end{aligned}$$

2. Compute Gain of each attribute

→ Terrain

$$1) \text{ Trail} = \{Y^2, N^3\}$$

$$\begin{aligned} E(\text{Trail}) &= -2/5 \log_2 (2/5) - 3/5 \log_2 (3/5) \\ &= 0.53 + 0.44 \\ &= \mathbf{0.97} \end{aligned}$$

Id	Terrain	Unicycle Type	Weather	Go-for Ride?
1	Trail	Normal	Rainy	No
2	Road	Normal	Sunny	Yes
3	Trail	Mountain	Sunny	Yes
4	Road	Mountain	Rainy	Yes
5	Trail	Normal	Snowy	No
6	Road	Normal	Rainy	Yes
7	Road	Mountain	Snowy	Yes
8	Trail	Normal	Sunny	No
9	Road	Normal	Snowy	No
10	Trail	Mountain	Snowy	Yes

$$\text{Entropy} = - \sum_{i=0}^n p_i(t) \log_2 p_i(t)$$

$$2) \text{ Road} = \{Y^4, N^1\}$$

$$E(\text{Road}) = -\frac{4}{5} \log_2 \left(\frac{4}{5}\right) - \frac{1}{5} \log_2 \left(\frac{1}{5}\right) \\ = 0.26 + 0.46 = 0.72$$

3) Average Entropy of Terrain

$$E(\text{Terrain}) = \left(\frac{5}{10}\right) E(\text{Trail}) + \left(\frac{5}{10}\right) E(\text{Road}) \\ = 0.485 + 0.36 = 0.845$$

$$\Rightarrow \text{Gain}(S, \text{Terrain}) = E(B) - E(\text{Terrain}) \\ = 0.97 - 0.845 \\ = 0.125$$

→ Unicycle-Type

$$1) \text{ Normal} = \{Y^2, N^4\}$$

$$E(\text{Normal}) = -\frac{2}{6} \log_2 \left(\frac{2}{6}\right) - \frac{4}{6} \log_2 \left(\frac{4}{6}\right) \\ = 0.53 + 0.39 = 0.92$$

$$2) \text{ Mountain} = \{Y^4, N^0\}$$

$E(\text{Mountain}) = 0$ (maximal homogeneity → minimal impurity → minimal entropy)

3) Average Entropy

$$E(\text{U-T}) = \left(\frac{6}{10}\right) E(\text{Normal}) + \left(\frac{4}{10}\right) E(\text{Mountain}) \\ = 0.552$$

$$\Rightarrow \text{Gain}(S, \text{Unicycle-Type}) = E(B) - E(\text{U-T}) \\ = 0.97 - 0.552 \\ = 0.418$$

→ Weather

1) Rainy: $\{Y^2, N^1\}$

$$E(\text{Rainy}) = -\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right)$$

$$= 0.39 + 0.53$$

$$= \mathbf{0.92}$$

2) Sunny: $\{Y^2, N^1\}$

$$E(\text{Sunny}) = 0.92$$

3) Snowy: $\{Y^2, N^2\}$

$$E(\text{Snowy}) = \mathbf{1}$$

(minimal homogeneity → maximal impurity → maximal entropy)

4) Average Entropy

$$E(\text{Weather}) = \left(\frac{3}{10}\right)E(\text{Rainy}) + \left(\frac{3}{10}\right)E(\text{Sunny})$$

$$+ \left(\frac{4}{10}\right)E(\text{Snowy})$$

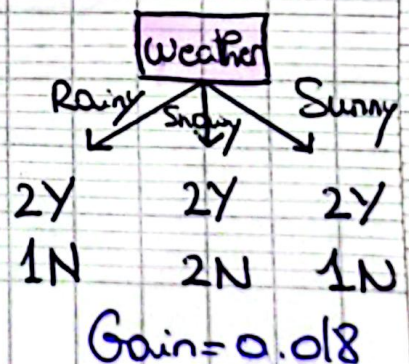
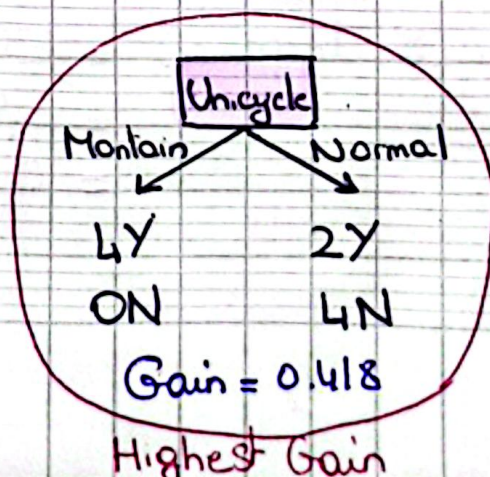
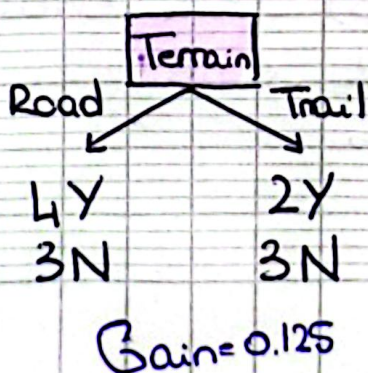
$$= \mathbf{0.952}$$

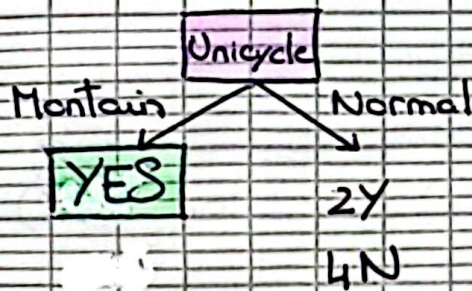
$$\Rightarrow \text{Gain}(S, \text{Weather}) = E(B) - E(\text{Weather})$$

$$= 0.97 - 0.952$$

$$= 0.018$$

3. Choose the attribute with the highest gain





Base case: All data belong to the same class, stop and create a leaf node with that label

→ Before Splitting

$$B = \{Y^2, N^4\}$$

$$E(B) = -\frac{2}{6} \log_2 \left(\frac{2}{6}\right) - \frac{4}{6} \log_2 \left(\frac{4}{6}\right)$$

$$= 0.53 + 0.39$$

$$= 0.92$$

$$= 0.92$$

Id	Terrain	Unicycle	Weather	Go For Ride?
1	Trail	Normal	Rainy	No
2	Road	Normal	Sunny	Yes
5	Trail	Normal	Snowy	No
6	Road	Normal	Rainy	Yes
8	Trail	Normal	Sunny	No
9	Road	Normal	Snowy	No

→ Terrain

$$1) \text{ Trail} = \{Y^0, N^3\}$$

$$E(\text{Trail}) = 0$$

$$2) \text{ Road} = \{Y^2, N^1\}$$

$$E(\text{Road}) = -\left(\frac{2}{3}\right) \log_2 \left(\frac{2}{3}\right) - \left(\frac{1}{3}\right) \log_2 \left(\frac{1}{3}\right)$$

$$= 0.39 + 0.53 = 0.92$$

3) Average Entropy

$$E(\text{Terrain}) = \left(\frac{3}{6}\right) E(\text{Trail}) + \left(\frac{3}{6}\right) E(\text{Road})$$

$$\Rightarrow \text{Gain}(S, \text{Terrain}) = E(B) - E(\text{Terrain})$$

$$= 0.92 - 0.46$$

$$= 0.46$$

Our new dataset S filtered on "Normal"

→ Weather

1) Rainy = $\{Y^1, N^1\}$
 $E(R) = 1$

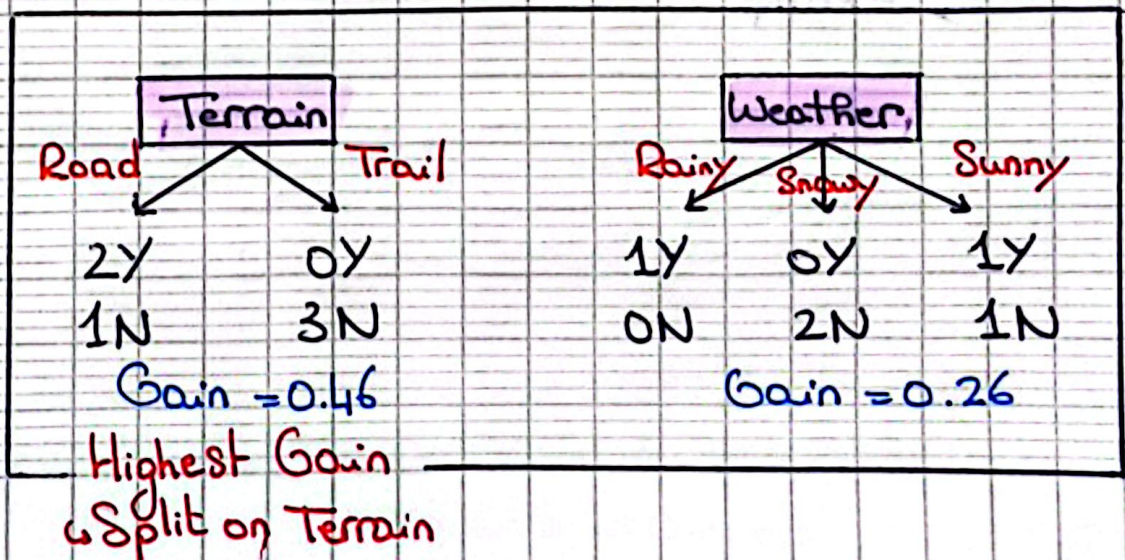
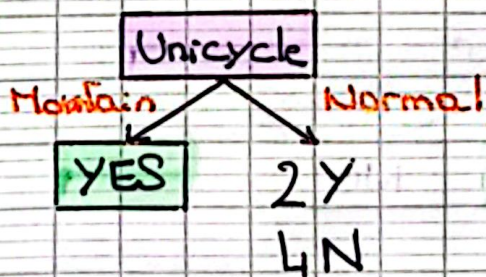
2) Sunny = $\{Y^1, N^1\}$
 $E(\text{Sunny}) = 1$

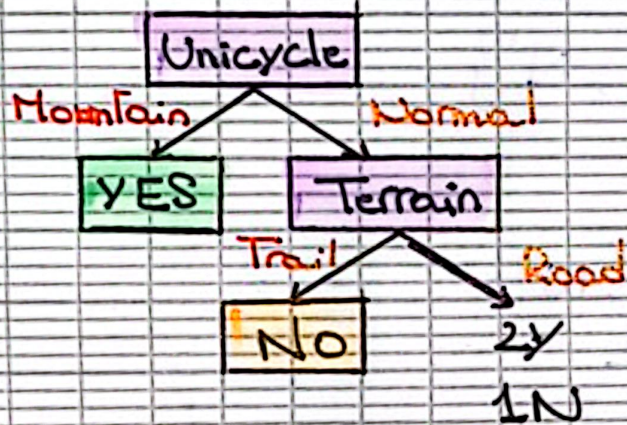
3) Snowy = $\{Y^0, N^2\}$
 $E(\text{Snowy}) = 0$

↳ Average Entropy

$$E(\text{Weather}) = 4/6 = 0.66$$

$$\begin{aligned} \Rightarrow \text{Gain}(S, \text{Weather}) &= E(B) - E(\text{Weather}) \\ &= 0.92 - 0.66 \\ &= 0.26 \end{aligned}$$

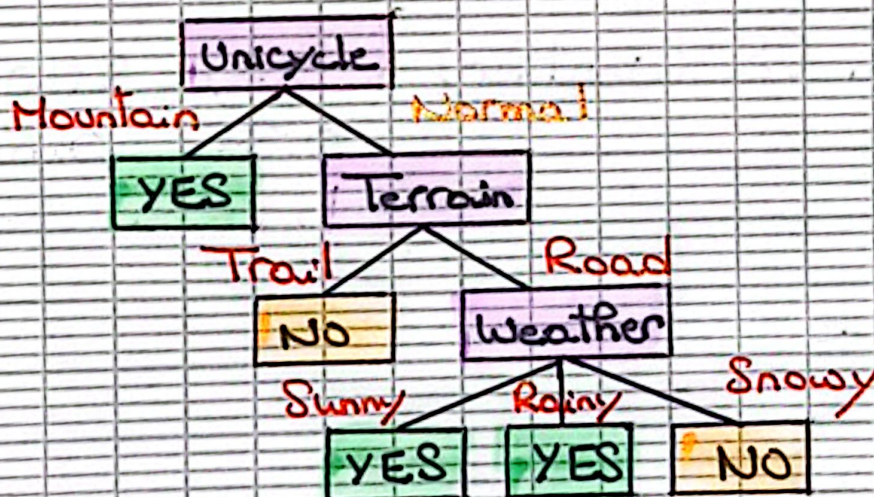




Id	Terrain	Unicycle	Weather	Go For Ride?
2	Road	Normal	Sunny	YES
6	Road	Normal	Rainy	YES
9	Road	Normal	Snowy	No

Our new dataset 'S'
Filtered on "Normal" and
"Road"

⇒ The Final Decision Tree Model



Example Using Gini-Index

Family = F

Sport = S

None = N

Compute Impurity
before Splitting

$$B = \{F^8, S^4, N^4\}$$

$$G.I(B) = 1 - \left[\left(\frac{8}{16}\right)^2 + \left(\frac{4}{16}\right)^2 + \left(\frac{4}{16}\right)^2 \right]$$

$$= 0.625$$

Compute G.I for
each attribute

Income

1) Low: $\{F^4, S^2, N^2\}$

$$G.I(L) = 1 - \left[\left(\frac{4}{6}\right)^2 + \left(\frac{2}{6}\right)^2 + \left(\frac{2}{6}\right)^2 \right]$$

$$= 0.44$$

2) Medium: $\{F^2, S^0, N^2\}$

$$G.I(M) = 1 - \left[\left(\frac{2}{4}\right)^2 + \left(\frac{2}{4}\right)^2 \right]$$

$$= 0.5$$

3) High: $\{F^2, S^4, N^0\}$

$$G.I(H) = 1 - \left[\left(\frac{2}{6}\right)^2 + \left(\frac{4}{6}\right)^2 \right] = 0.44$$

$$\hookrightarrow \text{Average (Income)} = \frac{6}{16} (0.44) + \frac{4}{16} (0.5) + \frac{6}{16} (0.44)$$

$$\Rightarrow G.I(\text{Income}) = G.I(B) - G.I(\text{Income})$$

$$= 0.625 - 0.485$$

$$= 0.175$$

Id	Income	No. of Child	Car
1	Low	2 or more	Family
2	Low	2 or more	Family
3	Low	2 or more	Family
4	Low	2 or more	Family
5	Medium	2 or more	Family
6	Medium	2 or more	Family
7	High	2 or more	Family
8	High	2 or more	Family
9	High	0	Sport
10	High	0	Sport
11	High	1	Sport
12	High	1	Sport
13	Medium	0	None
14	Medium	0	None
15	Low	2 or more	None
16	Low	2 or more	None

→ Nb. of Children

1) 2 or more : $\{F^8, S^0, N^2\}$

$$G.I(2 \text{ or more}) = 1 - \left[\left(\frac{8}{10} \right)^2 + \left(\frac{2}{10} \right)^2 \right] = 0.32$$

2) 0 : $\{F^0, S^2, N^2\}$

$$G.I(0) = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] = 0.5$$

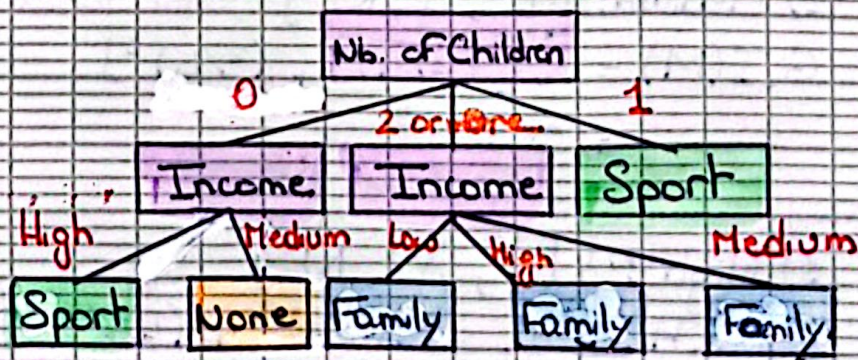
3) 1 : $\{F^0, S^2, N^0\}$

$$G.I(1) = 1 - \left[\left(\frac{2}{2} \right)^2 \right] = 0$$

$$\begin{aligned} \hookrightarrow \text{Avg (Nb. of Children)} &= \frac{10}{16} (0.32) + \frac{4}{16} (0.5) \\ &\quad + \frac{2}{16} (0) \\ &= 0.325 \end{aligned}$$

$$\begin{aligned} \Rightarrow G.I(\text{Nb. of Children}) &= 0.625 - 0.325 \\ &= 0.3 \end{aligned}$$

⇒ Nb. of Children has the higher G.I., so we split on it first



Example: Computing Gini Index for Continuous Attributes

- Sort Income in increasing order
- Find the midpoint between each pair of values
 $(a_i + a_{i+1}) / 2$
- Choose the split position that has the least G.I or E

Id	Home owner	Marital Status	Income	Default
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

		Annual Income																
		60	70	75	85	90	95	100	120	125	220							
		65	72	80	87	92	97	110	122	172								
		<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>			
YES		0	3	0	3	0	3	1	2	2	1	3	0	3	0	3	0	
NO		1	6	2	5	3	4	3	4	3	4	3	4	3	5	2	6	1
Gini Index		0.396	0.374	0.34	0.414	0.4	0.3	0.34	0.37	0.396								

↑ minimum so we split on (97)

$$\begin{aligned}
 G.I(<=65) &= 1 - \left[\left(\frac{1}{1}\right)^2 + \left(\frac{0}{1}\right)^2 \right] = 0 \\
 G.I(>65) &= 1 - \left[\left(\frac{3}{9}\right)^2 + \left(\frac{6}{9}\right)^2 \right] = 0.44 \\
 G.I(<=72) &= 1 - \left[\left(\frac{0}{2}\right)^2 + \left(\frac{2}{2}\right)^2 \right] = 0 \\
 G.I(>72) &= 1 - \left[\left(\frac{3}{8}\right)^2 + \left(\frac{5}{8}\right)^2 \right] = 0.468 \\
 G.I(<=80) &= 1 - \left[\left(\frac{0}{3}\right)^2 + \left(\frac{3}{3}\right)^2 \right] = 0 \\
 G.I(>80) &= 1 - \left[\left(\frac{3}{7}\right)^2 + \left(\frac{4}{7}\right)^2 \right] = 0.489 \\
 \dots
 \end{aligned}$$

} Avg
 $\frac{1}{10} \times 0 + 9/10 \times 0.44 = 0.396$
 } Avg
 0.3744
 } Avg
 0.34