

Problem 1 coding problem

Problem 2 The time complexity for training a decision tree is generally $O(n \cdot m \log n)$ where:

n = the number of training instances

m = the number of features

We are given:

Training time for $n=1000000$ instances is 1 hour.

Assume the number of features m remains the same, and now we estimate for $n=10,000,000$ instances:

The time complexity scales with $n \log n$ (in this case since m is the same):

$$\frac{\text{Time}_{\text{new}}}{\text{Time}_{\text{old}}} = \frac{n_2 \log n_2}{n_1 \log n_1} = \frac{10^7 \log 10^7}{10^6 \log 10^6} = \frac{10 \log 10^7}{\log 10^6} = \frac{10 \cdot 7 \log 10}{6 \log 10} = \frac{70}{6} \approx 11.67 \Rightarrow$$

$$n_2 = 10,000,000$$

$$n_1 = 1,000,000$$

$$\Rightarrow \text{Time}_{\text{new}} = 11.67 \cdot \text{Time}_{\text{old}} = 11.67 \cdot 1 = 11.67 \text{ hours} = 11 \frac{2}{3} \text{ hours}$$

Answer: 11 hours 40 minutes

Problem 3 Gini index, entropy, and classification errors are impurity measures used in decision tree algorithms to decide how to split data at each node. They help quantify how "mixed" or "pure" a node is.

Gini Index: it measures the probability of incorrectly classifying a randomly chosen element if it was randomly labeled according to the distribution of labels in the node. The greater the value of Gini index, the greater the chances of having misclassifications.

$$G_{\text{ini}} = 1 - \sum_{i=1}^c (p_i)^2 \quad \text{where } p_i = \text{the frequency of class } i \text{ at node } t, \text{ and } c \text{ is the total number of classes}$$

(The example is on the next page)

Market Sentiment	Liquidity	Volatility	Return
Bullish	High	Low	Up
Bearish	Low	High	Down
Bullish	High	Low	Up
Bullish	Low	High	Down
Bearish	High	Low	Up
Bullish	Low	Low	Down
Bearish	High	High	Down
Bullish	Low	Low	Down
Bullish	High	Low	Up

The Gini index for a split is

$$Gini = 1 - \sum_{i=1}^k (P_i)^2$$

where P_i = proportion of class i in the subset

• Calculating the Gini index for Market Sentiment (MS)

$$\text{Bullish: } 6 \text{ number of times out of 9} \quad | \Rightarrow P(MS = \text{Bullish}) = 6/9 = 2/3$$

$$\text{Bearish: } 3 \text{ number of times out of 9} \quad | \Rightarrow P(MS = \text{Bearish}) = 3/9 = 1/3$$

$$\bullet P(\text{Return} = \text{Up} | MS = \text{Bullish}) = \frac{P(\text{Return} = \text{Up} \cap MS = \text{Bullish})}{P(MS = \text{Bullish})} = \frac{\#(\text{Return} = \text{Up} \cap MS = \text{Bullish})}{\#(MS = \text{Bullish})} = \frac{3}{6} = \frac{1}{2}$$

$$\bullet P(\text{Return} = \text{Down} | MS = \text{Bullish}) = \frac{\#(\text{Return} = \text{Down} \cap MS = \text{Bullish})}{\#(MS = \text{Bullish})} = \frac{3}{6} = \frac{1}{2}$$

$$\Rightarrow \text{Gini Index} = 1 - \left[\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right] = 1 - \left[\frac{1}{4} + \frac{1}{4} \right] = 1 - \frac{1}{2} = \frac{1}{2}$$

$$\bullet P(\text{Return} = \text{Up} | MS = \text{Bearish}) = \frac{\#(\text{Return} = \text{Up} \cap MS = \text{Bearish})}{\#(MS = \text{Bearish})} = \frac{1}{3}$$

$$\bullet P(\text{Return} = \text{Down} | MS = \text{Bearish}) = \frac{\#(\text{Return} = \text{Down} \cap MS = \text{Bearish})}{\#(MS = \text{Bearish})} = \frac{2}{3}$$

$$\Rightarrow \text{Gini Index} = 1 - \left[\left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right] = 1 - \left(\frac{1}{9} + \frac{4}{9}\right) = 1 - \frac{5}{9} = \frac{4}{9}$$

$$\text{Gini Index} = P(MS = \text{Bullish}) \cdot \frac{1}{2} + P(MS = \text{Bearish}) \cdot \frac{4}{9} = \frac{2}{3} \cdot \frac{1}{2} + \frac{1}{3} \cdot \frac{4}{9} = \frac{1}{3} + \frac{4}{27} = \frac{13}{27} \approx 0.48148 \quad \textcircled{2}$$

$$\textcircled{2} \quad \text{Gini Index for MS} = 0.48148$$

• Calculate the Gini Index for Liquidity (L)

$$\text{High: } 5 \text{ number out of 9} \quad | \Rightarrow P(L = \text{High}) = 5/9$$

$$\text{Low: } 4 \text{ number out of 9} \quad | \Rightarrow P(L = \text{Low}) = 4/9$$

$$\bullet P(\text{Return} = \text{Up} | L = \text{High}) = \frac{\#(\text{Return} = \text{Up} \text{ and } L = \text{High})}{\#(L = \text{High})} = \frac{6}{5}$$

$$\bullet P(\text{Return} = \text{Down} | L = \text{High}) = \frac{\#(\text{Return} = \text{Down} \cap L = \text{High})}{\#(L = \text{High})} = \frac{1}{5}$$

$$\Rightarrow \text{Gini Index} = 1 - \left[\left(\frac{6}{5}\right)^2 + \left(\frac{1}{5}\right)^2 \right] = 1 - \left(\frac{36}{25} + \frac{1}{25}\right) = \frac{8}{25}$$

- $P(\text{Return} = \text{Up} | L = \text{Low}) = \frac{\#(\text{Return} = \text{Up} \wedge L = \text{Low})}{\#(L = \text{Low})} = \frac{0}{5} = 0$
- $P(\text{Return} = \text{Down} | L = \text{Low}) = \frac{\#(\text{Return} = \text{Down} \wedge L = \text{Low})}{\#(L = \text{Low})} = \frac{5}{5} = 1$

$\Rightarrow Gini \text{ Index} = 1 - (1+0) = 0$

$$Gini \text{ Index (for liquidity)} = \frac{5}{9} \cdot \frac{3}{25} + \frac{4}{9} \cdot 0 = \frac{8}{45} \approx 0.1777 \Rightarrow Gini \text{ Index for liquidity} = 0.1777$$

- Calculating the Gini Index for Volatility (V)

High: 3 number out of 9 | $P(V = \text{High}) = \frac{1}{3}$

Low: 6 number out of 9 | $P(V = \text{Low}) = \frac{2}{3}$

- $P(\text{Return} = \text{Up} | V = \text{High}) = \frac{\#(\text{Return} = \text{Up} \wedge V = \text{High})}{\#(V = \text{High})} = \frac{0}{3} = 0$

- $P(\text{Return} = \text{Down} | V = \text{High}) = \frac{\#(\text{Return} = \text{Down} \wedge V = \text{High})}{\#(V = \text{High})} = \frac{3}{3} = 1$

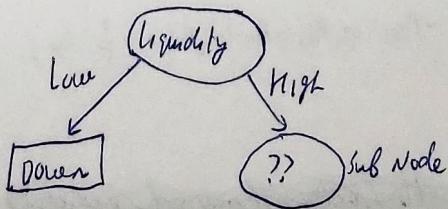
- $P(\text{Return} = \text{Up} | V = \text{Low}) = \frac{\#(\text{Return} = \text{Up} \wedge V = \text{Low})}{\#(V = \text{Low})} = \frac{4}{6} = \frac{2}{3}$

- $P(\text{Return} = \text{Down} | V = \text{Low}) = \frac{\#(\text{Return} = \text{Down} \wedge V = \text{Low})}{\#(V = \text{Low})} = \frac{2}{6} = \frac{1}{3}$

$$Gini \text{ Index (for Volatility)} = \frac{1}{3} \cdot 0 + \frac{2}{3} \cdot \frac{1}{3} = \frac{2}{27} \approx 0.2962 \Rightarrow Gini \text{ Index for Volatility} \approx 0.2962$$

Features	Gini Index
Market sentiment	0.48148
Liquidity	0.1777
Volatility	0.2962

Since Liquidity has the lowest Gini-Index values, we will use it to split the dataset as follows:



We will repeat the same procedure to determine the sub-nodes or branches of the decision tree.
We will calculate the Gini Index for the "High" branch of Liquidity, as follows.

Market Sentiment	Liquidity	Volatility	Return
Bullish	High	Low	Up
Bullish	High	Low	Up
Bearish	High	Low	Up
Bearish	High	High	Down
Bullish	High	Low	Up

And we will do the same process on this new table, and continue building the decision tree.

- Entropy: Entropy is a measure of impurity or randomness in a dataset. It can be defined as a method to measure the impurity or uncertainty. A system or model with lowest entropy is considered is considered better than the other with high entropy.

The mathematical formula for entropy is defined as:

$$E(S) = - \sum_{i=1}^k p_i \log_2 p_i \quad \text{where } S \text{ is a dataset (or node), } k \text{ is the number of classes, } p_i \text{ is the probability of a class label } i \text{ in our data.}$$

For example, if our dataset has 2 classes Yes and No, then the entropy can be calculated as:

$$E(S) = -p_{\text{Yes}} \log_2 p_{\text{Yes}} - p_{\text{No}} \log_2 p_{\text{No}}$$

Example

Income (\$)	Experience (years)	loan approval
30,000	2	YES
15,000	10	NO
70,000	5	YES
20,000	8	NO
25,000	1	NO

$$\text{Calculate Entropy} = -p_{\text{Yes}} \log_2 p_{\text{Yes}} - p_{\text{No}} \log_2 p_{\text{No}} = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.97,$$

Classification Error :

Classification Error is a measure of the proportion of misclassified instances in a dataset. In decision trees, it is used to evaluate the purity of a node by considering how often the majority class fails to represent all the data points.

$$\text{Formula : Classification Error} = 1 - \max_{i=1,2,\dots,k} p_i \quad \text{where } k = \# \text{ of classes}$$

p_i = the proportion of classes belonging to class i .

Consider a dataset:

Class	Count
A	7
B	3

$$P_A = \frac{7}{10} = 0.7$$

$$\max(P_A, P_B) = \max(0.7, 0.3) = 0.7$$

$$P_B = \frac{3}{10} = 0.3$$

$$E = 1 - 0.7 = 0.3$$

classification error

Problem

To extract decision rules, we can use the following methods:

- Direct Method: Direct Method extract rules directly from the training data without first building an intermediate model like a decision tree or NN.

These use separable-and-conquer strategies (also called covering algorithms)

Each rule is generated to cover as many examples of one class as possible without excluding others. The covered examples are removed, and the process is repeated.

Direct Method: Sequential Coverage

- 1) Start from an empty rule
- 2) Grow a rule using the learn-one-rule function
- 3) Remove training records covered by the rule.
- 4) Repeat step 2) and 3) until stopping criteria is met.

Indirect Method:

Indirect methods extract rules from an already trained model, such as decision trees, random forest, neural network, or SVM.

It, first, builds a predictive model, then analyzes the model's internal structure or outputs to generate human-readable rules.

Decision Rule Induction is a method for extracting if-then rules from data, typically used for classification tasks. These rules are designed to be interpretable, actionable, and easy to apply. One example of decision rule induction algorithm is RIPPER. It generates a set of if-then classification rules directly from labeled data and is particularly good for datasets with imbalanced classes.

We should use the decision rule induction method, when we need clear, understandable logic (i.e. medical diagnosis, credit scoring, legal decisions, etc). And when the rules can be quickly applied to new data with little computational cost.

Problem 5

• Calculating the Gini-Index for Age:

$$\begin{array}{l} \text{youth: 5 number out of 14} \\ \text{middle-aged: 4 number out of 14} \\ \text{senior: 5 number out of 14} \end{array} \quad \left| \begin{array}{l} P(\text{age} = \text{youth}) = 5/14 \\ P(\text{age} = \text{middle-aged}) = 4/14 \\ P(\text{age} = \text{senior}) = 5/14 \end{array} \right.$$

$$\bullet P(\text{buys-computer} = \text{yes} | \text{age} = \text{youth}) = \frac{\#(\text{buys-computer} = \text{yes} \cap \text{age} = \text{youth})}{\#(\text{age} = \text{youth})} = \frac{2}{5}$$

$\Rightarrow \text{Gini-Index} = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 1 - \left(\frac{4}{25} + \frac{9}{25}\right) = 1 - \frac{13}{25} = \boxed{\frac{12}{25}}$

$$\bullet P(\text{buys-computer} = \text{no} | \text{age} = \text{youth}) = \frac{\#(\text{buys} = \text{no} \cap \text{age} = \text{youth})}{\#(\text{age} = \text{youth})} = \frac{3}{5}$$

$$\bullet P(\text{buys} = \text{yes} | \text{age} = \text{middle-aged}) = \frac{\#(\text{buys} = \text{yes} \cap \text{age} = \text{middle-aged})}{\#(\text{age} = \text{middle-aged})} = \frac{4}{5} = 1$$

$\Rightarrow \text{Gini-Index} = 1 - 0 - 0 = 0$

$$\bullet P(\text{buys} = \text{no} | \text{age} = \text{middle-aged}) = \frac{\#(\text{buys} = \text{no} \cap \text{age} = \text{middle-aged})}{\#(\text{age} = \text{middle-aged})} = \frac{0}{5} = 0$$

$$\bullet P(\text{buys} = \text{yes} | \text{age} = \text{senior}) = \frac{\#(\text{buys} = \text{yes} \cap \text{age} = \text{senior})}{\#(\text{age} = \text{senior})} = \frac{3}{5}$$

$\Rightarrow \text{Gini-Index} = 1 - \left(\frac{3}{5}\right)^2 + \left(\frac{2}{5}\right)^2 = 1 - \frac{13}{25} = \frac{12}{25}$

$$\bullet P(\text{buys} = \text{no} | \text{age} = \text{senior}) = \frac{\#(\text{buys} = \text{no} \cap \text{age} = \text{senior})}{\#(\text{age} = \text{senior})} = \frac{2}{5}$$

$$\text{Weighted Gini-Index for age} = \frac{5}{14} \cdot \frac{12}{25} + \frac{4}{14} \cdot 0 + \frac{5}{14} \cdot \frac{12}{25} = \frac{5}{7} \cdot \frac{12}{25} = \frac{12}{35} \approx 0,34285$$

• Calculating the Gini-Index for Income:

$$\begin{array}{l} \text{high: 4 number out of 14} \\ \text{medium: 6 out of 14} \\ \text{low: 4 out of 14} \end{array} \quad \left| \begin{array}{l} P(\text{income} = \text{high}) = 4/14 \\ P(\text{income} = \text{medium}) = 6/14 \\ P(\text{income} = \text{low}) = 4/14 \end{array} \right.$$

$$\bullet P(\text{buys} = \text{yes} | \text{income} = \text{high}) = \frac{\#(\text{buys} = \text{yes} \cap \text{income} = \text{high})}{\#(\text{income} = \text{high})} = \frac{2}{4}$$

$\Rightarrow \text{Gini-Index} = 1 - \left(\frac{2}{4}\right)^2 + \left(\frac{2}{4}\right)^2 = 1 - \frac{1}{2} = \underline{\underline{\frac{1}{2}}}$

$$\bullet P(\text{buys} = \text{no} | \text{income} = \text{high}) = \frac{\#(\text{buys} = \text{no} \cap \text{income} = \text{high})}{\#(\text{income} = \text{high})} = \frac{2}{4}$$

$$\bullet P(\text{Buys} = \text{yes} | \text{income} = \text{medium}) = \frac{\#(\text{Buys} = \text{yes} \cap \text{income} = \text{medium})}{\#(\text{income} = \text{medium})} = \frac{4}{6} \quad \Rightarrow G_{\text{m-Index}} = -\left(\frac{1}{8}\right)^2 + \left(\frac{3}{8}\right)^2 = 1 - \left(\frac{4}{9} + \frac{1}{9}\right) = \underline{\frac{4}{9}}$$

$$\bullet P(\text{Buys} = \text{no} | \text{income} = \text{medium}) = \frac{\#(\text{Buys} = \text{no} \cap \text{income} = \text{medium})}{\#(\text{income} = \text{medium})} = \frac{2}{6}$$

$$\bullet P(\text{Buys} = \text{yes} | \text{income} = \text{low}) = \frac{\#(\text{Buys} = \text{yes} \cap \text{income} = \text{low})}{\#(\text{income} = \text{low})} = \frac{3}{4} \quad \Rightarrow G_{\text{m-Index}} = -\left(\frac{1}{4}\right)^2 + \left(\frac{1}{4}\right)^2 = 1 - \frac{10}{16} = \underline{\frac{3}{8}}$$

$$\bullet P(\text{Buys} = \text{no} | \text{income} = \text{low}) = \frac{\#(\text{Buys} = \text{no} \cap \text{income} = \text{low})}{\#(\text{income} = \text{low})} = \frac{1}{4}$$

$$\text{Weighted G}_{\text{m-Index}} \text{ for Income} = \frac{4}{16} \cdot \frac{1}{2} + \frac{6}{16} \cdot \frac{4}{9} + \frac{4}{16} \cdot \frac{3}{8} = \frac{1}{7} + \frac{4}{21} + \frac{3}{28} = \frac{12 + 16 + 9}{84} = \underline{\frac{37}{84} \approx 0.44057}$$

• Calculate the G_m-Index for student

yes : 7 out of 14: $\left| \begin{array}{l} P(\text{student} = \text{yes}) = \frac{1}{2} \\ \Rightarrow P(\text{student} = \text{no}) = \frac{1}{2} \end{array} \right.$

$$\bullet P(\text{Buys} = \text{yes} | \text{student} = \text{yes}) = \frac{\#(\text{Buys} = \text{yes} \cap \text{student} = \text{yes})}{\#(\text{student} = \text{yes})} = \frac{6}{7} \quad \Rightarrow G_{\text{m-Index}} = -\left(\frac{1}{7}\right)^2 + \left(\frac{1}{7}\right)^2 = 1 - \frac{37}{49} = \underline{\frac{12}{49}}$$

$$\bullet P(\text{Buys} = \text{no} | \text{student} = \text{yes}) = \frac{\#(\text{Buys} = \text{no} \cap \text{student} = \text{yes})}{\#(\text{student} = \text{yes})} = \frac{1}{7}$$

$$\bullet P(\text{Buys} = \text{yes} | \text{student} = \text{no}) = \frac{\#(\text{Buys} = \text{yes} \cap \text{student} = \text{no})}{\#(\text{student} = \text{no})} = \frac{3}{7} \quad \Rightarrow G_{\text{m-Index}} = -\left(\frac{3}{7}\right)^2 + \left(\frac{4}{7}\right)^2 = 1 - \frac{25}{49} = \underline{\frac{24}{49}}$$

$$\bullet P(\text{Buys} = \text{no} | \text{student} = \text{no}) = \frac{\#(\text{Buys} = \text{no} \cap \text{student} = \text{no})}{\#(\text{student} = \text{no})} = \frac{4}{7}$$

$$\text{Weighted G}_{\text{m-Index}} \text{ for student} = \frac{1}{2} \cdot \frac{12}{49} + \frac{1}{2} \cdot \frac{24}{49} = \frac{1}{2} \cdot \frac{36}{49} = \underline{\frac{18}{49} \approx 0.3673}$$

• Calculating the Gini-Index for Credit-Rating

fair: 8 out of 14 | $P(\text{credit} = \text{fair}) = 8/14$

excellent: 6 out of 14 | $P(\text{credit} = \text{excellent}) = 6/14$

$$\bullet P(\text{buys} = \text{yes} | \text{credit} = \text{fair}) = \frac{\#\{\text{buys} = \text{yes} \wedge \text{credit} = \text{fair}\}}{\#\{\text{credit} = \text{fair}\}} = \frac{6}{8}$$

$$\bullet P(\text{buys} = \text{no} | \text{credit} = \text{fair}) = \frac{\#\{\text{buys} = \text{no} \wedge \text{credit} = \text{fair}\}}{\#\{\text{credit} = \text{fair}\}} = \frac{2}{8} \quad \Rightarrow \text{Gini-Index} = 1 - \left(\frac{6}{8}\right)^2 + \left(\frac{2}{8}\right)^2 = 1 - \frac{3}{8} = \frac{5}{8}$$

$$\bullet P(\text{buys} = \text{yes} | \text{credit} = \text{excellent}) = \frac{\#\{\text{buys} = \text{yes} \wedge \text{credit} = \text{excellent}\}}{\#\{\text{credit} = \text{excellent}\}} = \frac{3}{6}$$

$$\bullet P(\text{buys} = \text{no} | \text{credit} = \text{excellent}) = \frac{\#\{\text{buys} = \text{no} \wedge \text{credit} = \text{excellent}\}}{\#\{\text{credit} = \text{excellent}\}} = \frac{3}{8} \quad \Rightarrow \text{Gini-Index} = 1 - \left(\frac{3}{8}\right)^2 + \left(\frac{3}{8}\right)^2 = \frac{1}{2}$$

$$\text{Weighted Gini-Index for credit rating} = \frac{8}{14} \cdot \frac{5}{8} + \frac{6}{14} \cdot \frac{1}{2} = \frac{3}{14} + \frac{3}{14} = \frac{6}{14} \approx 0,42857$$

Now, we need to calculate the information gain

Step 1: Calculate Entropy for Parent

$$\text{Entropy}(\text{parent}) = -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}} = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} \approx 0.94028$$

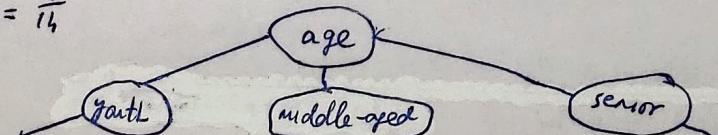
$$P_{\text{yes}} = \frac{9}{14}, P_{\text{no}} = \frac{5}{14}$$

Step 2: Calculate Entropy for children.

$$\text{Age: } P(\text{age} = \text{youth}) = \frac{5}{14}$$

$$P(\text{age} = \text{middle-aged}) = \frac{5}{14}$$

$$P(\text{age} = \text{senior}) = \frac{4}{14}$$



age	income	student	credit-rating	buys
				computer
youth	high	no	fair	no
youth	high	no	excellent	no
youth	medium	no	fair	no
youth	low	yes	fair	yes
youth	medium	yes	excellent	yes

age	income	student	credit-rating	buys
				computer
middle-aged	high	no	fair	yes
middle-aged	low	yes	excellent	yes
middle-aged	medium	no	excellent	yes
middle-aged	high	yes	fair	yes

age	income	student	credit-rating	buys
				computer
senior	medium	no	fair	yes
senior	low	yes	fair	yes
senior	low	yes	excellent	no
senior	medium	yes	fair	yes
senior	medium	no	excellent	no

Calculate Entropy for youth

$$\text{Entropy(youth)} = -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}} = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \approx 0.97095$$

$$P_{\text{yes}} = \frac{2}{5}$$

$$P_{\text{no}} = \frac{3}{5}$$

Calculate Entropy for middle-aged

$$\text{Entropy(middle-aged)} = -P_{\text{yes}} \log_2 P_{\text{yes}} - 0 = 0 \log 1 = 0$$

$$P_{\text{yes}} = \frac{1}{5} = 0$$

$$P_{\text{no}} = \frac{4}{5} = 0$$

Calculate entropy for senior

$$\text{Entropy(senior)} = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \approx 0.97095$$

$$P_{\text{yes}} = \frac{3}{5}$$

$$P_{\text{no}} = \frac{2}{5}$$

Step 3: Calculate average weighted entropy of each children:

$$\text{Weighted Entropy(youth)} = \frac{5}{14} \cdot 0.97095 \approx 0.34678$$

$$\text{Weighted Entropy(middle-aged)} = \frac{1}{14} \cdot 0 = 0$$

$$\text{Weighted Entropy(senior)} = \frac{5}{14} \cdot 0.97095 \approx 0.34678$$

$$\Rightarrow \text{Average Weighted Entropy} = 0.34678 + 0 + 0.34678 = 0.69352$$

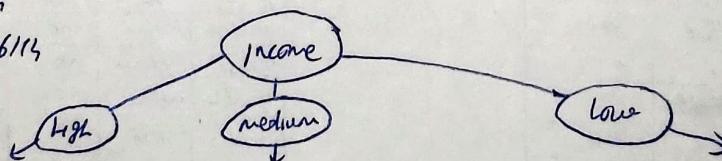
Step 4: Calculate Information gain (IG)

$$\text{Gain}(I, \text{Age}) = \text{Gain}(I) - \frac{1}{3} \text{Average Weighted Entropy} = 0.94028 - 0.69352 = 0.24676$$

Income: $P(\text{income} = \text{high}) = \frac{4}{14}$

Step 2 $P(\text{income} = \text{medium}) = 6/14$

$P(\text{income} = \text{low}) = 4/14$



age	income	student	credit rating	buys computer
youth	high	no	fair	no
youth	high	no	excellent	no
middle-aged	high	no	fair	yes
middle-aged	high	yes	fair	yes

age	income	student	credit rating	buys computer
senior	medium	no	fair	yes
youth	medium	no	fair	no
senior	medium	yes	fair	yes
youth	medium	yes	excellent	yes
middle-aged	medium	no	excellent	yes
senior	medium	no	excellent	no

age	income	student	credit rating	buys computer
senior	low	yes	fair	yes
senior	low	yes	excellent	no
middle-aged	low	yes	excellent	yes
youth	low	yes	fair	yes

Calculate Entropy for high

$$\bullet \text{Entropy}(\text{high}) = -p_{\text{yes}} \log_2 p_{\text{yes}} - p_{\text{no}} \log_2 p_{\text{no}} = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$p_{\text{yes}} = \frac{1}{2}$$

$$p_{\text{no}} = \frac{1}{2}$$

Calculate Entropy for medium

$$\bullet \text{Entropy}(\text{medium}) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0,91829$$

$$p_{\text{yes}} = \frac{2}{3}$$

$$p_{\text{no}} = \frac{1}{3}$$

Calculate Entropy for low

$$\bullet \text{Entropy}(\text{low}) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0,81127$$

$$p_{\text{yes}} = \frac{3}{5}$$

$$p_{\text{no}} = \frac{2}{5}$$

Step 3: Calculate average weighted entropy of each children.

$$\text{Weighted } E(\text{high}) = \frac{1}{7} \cdot 1 = \frac{1}{7}$$

$$\text{Weighted } E(\text{medium}) = \frac{6}{7} \cdot 0,91829 = 0,39355$$

$$\text{Weighted } E(\text{low}) = \frac{1}{7} \cdot 0,81127 = 0,23179$$

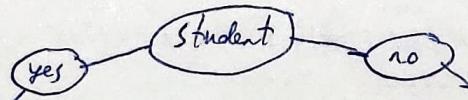
$$\rightarrow \text{Average Weighted } E(\text{income}) = \frac{1}{7} + 0,39355 + 0,23179 = 0,91105$$

Step 4: Calculate Information Gain. (IG)

$$\text{Gain}(R, \text{income}) = \text{Gain}(R) - \text{d Average Weighted } E(\text{income}) \rightarrow 0,94028 - 0,91105 = 0,02923$$

$$\underline{\text{Student}}: \quad p(\text{student} = \text{yes}) = \frac{2}{7} = \frac{1}{2} \quad p(\text{student} = \text{no}) = \frac{5}{7} = \frac{1}{2}$$

Step 2



age	income	student	credit rating	buys computer
senior	low	yes	fair	yes
senior	low	yes	excellent	no
middle	low	yes	excellent	yes
young	low	yes	fair	yes
senior	medium	yes	fair	yes
young	medium	yes	excellent	yes
middle	high	yes	fair	yes

age	income	student	credit rating	buys computer
young	high	no	fair	no
young	high	no	excellent	no
middle	high	no	fair	yes
senior	medium	no	fair	yes
young	medium	no	fair	no
middle	medium	no	excellent	yes
senior	medium	no	excellent	no

Calculate Entropy for yes (student)

$$\bullet \text{Entropy}(\text{yes}) = -p_{\text{yes}} \log_2 p_{\text{yes}} - p_{\text{no}} \log_2 p_{\text{no}} = -\frac{2}{7} \log_2 \frac{2}{7} - \frac{5}{7} \log_2 \frac{5}{7} = 0,59167$$

$$p_{\text{yes}} = \frac{2}{7}$$

$$p_{\text{no}} = \frac{5}{7}$$

Step 3: Calculate average weighted entropy of each children

$$\text{Weighted } E(\text{yes}) = \frac{1}{2} \cdot 0,59167 = 0,295835 \quad \rightarrow \text{Average Weighted } E(\text{student}) = 0,295835 + 0,49261 = 0,788495$$

$$\text{Weighted } E(\text{no}) = \frac{1}{2} \cdot 0,98572 = 0,49261$$

Calculate Entropy for no (student)

$$\bullet \text{Entropy}(\text{no}) = -p_{\text{yes}} \log_2 p_{\text{yes}} - p_{\text{no}} \log_2 p_{\text{no}} = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} = 0,98572$$

$$p_{\text{yes}} = \frac{3}{7}$$

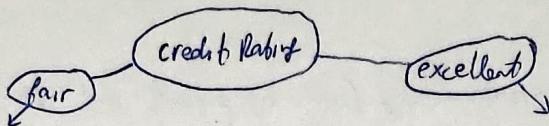
$$p_{\text{no}} = \frac{4}{7}$$

Step 1: Calculate Information Gain (IG)

$$\text{Gain}(R_1, \text{student}) = \text{Gain}(R) - \{\text{Average Weighted Entropy}\} = 0,94028 - 0,789445 = 0,151835$$

Credit-Rating: $P(\text{credit-rating} = \text{fair}) = \frac{3}{14}$ $P(\text{credit-rating} = \text{excellent}) = \frac{6}{14}$

Step 2:



age	income	student	credit rating	buys computer
youth	high	no	fair	no
middle	high	no	fair	yes
senior	medium	no	fair	yes
senior	low	yes	fair	yes
youth	medium	no	fair	no
youth	low	yes	fair	yes
senior	medium	yes	fair	yes
middle	high	yes	fair	yes

age	income	student	credit rating	buys computer
youth	high	no	excellent	no
senior	low	yes	excellent	no
middle	low	yes	excellent	yes
youth	medium	yes	excellent	yes
middle	medium	no	excellent	yes
senior	medium	no	excellent	no

Calculate Entropy for fair

$$\text{Entropy(fair)} = -\frac{3}{7} \log \frac{3}{7} - \frac{4}{7} \log \frac{4}{7} = 0,81127$$

$$P_{yes} = \frac{6}{8}$$

$$P_{no} = \frac{2}{8}$$

$$\text{Entropy(excellent)} = -\frac{1}{7} \log \frac{1}{7} - \frac{6}{7} \log \frac{6}{7} = 1$$

$$P_{yes} = \frac{1}{2}$$

$$P_{no} = \frac{1}{2}$$

Step 3: Calculate average weighted entropy of each division

$$\text{Weighted Entropy(fair)} = \frac{3}{14} \cdot 0,81127 = 0,46358$$

$$\text{Weighted Entropy(excellent)} = \frac{6}{14} \cdot 1 = \frac{6}{14}$$

$$\Rightarrow \text{Average Weighted Entropy(credit-Rating)} = 0,46358 + \frac{6}{14} = 0,89215$$

Step 4: Calculate Information Gain (IG)

$$\text{Gain}(R_1, \text{credit-rating}) = \text{Gain}(R) - \{\text{Average Weighted Entropy(credit-Rating)}\} = 0,94028 - 0,89215 = 0,04813$$

We got the following Gini-Indices and Information-Gain for each features

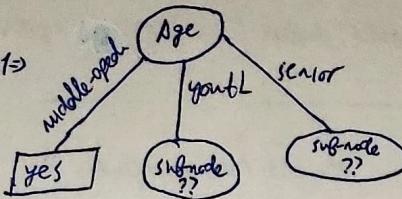
Gini-Index

Feature	Gini-Index	Info-Gain
Age	0,32857	0,23676
Income	0,44067	0,02923
Student	0,3673	0,151835
Credit R.	0,42857	0,04813

Age has the smallest Gini-Index (also the largest Information-Gain), therefore we use age as the splitting attribute for the first split.

When we split by Age feature, we get the following decision tree (so far)

we have $P(\text{buys_computer} = \text{yes} | \text{age} = \text{middle}) = 1 \Rightarrow$
 \Rightarrow we can build the sub-node
 middle-aged immediately



Now, we need to construct the subnodes coming from $\text{age} = \text{youth}$ and $\text{age} = \text{senior}$. For that we need to split the datasets accordingly: (into smaller datasets)

(1)

age	income	student	credit rating	buys computer
youthL	high	no	fair	no
youthL	high	no	excellent	no
youth	medium	no	fair	no
youth	low	yes	fair	yes
youth	medium	yes	excellent	yes

(2)

age	income	student	credit rating	buys computer
senior	medium	no	fair	yes
senior	low	yes	fair	yes
senior	low	yes	excellent	no
senior	medium	yes	fair	yes
senior	medium	no	excellent	no

- Calculating Gini-Indices for (1) ($\text{age} = \text{youthL}$) sub-dataset.

Calculating the Gini-Index for income

$$\begin{aligned} \text{high: } & 2 \text{ out of } 5 & P(\text{income} = \text{high}) &= \frac{2}{5} \\ \text{medium: } & 2 \text{ out of } 5 & P(\text{income} = \text{medium}) &= \frac{2}{5} \\ \text{low: } & 1 \text{ out of } 5 & P(\text{income} = \text{low}) &= \frac{1}{5} \end{aligned}$$

$$\bullet P(\text{buys_comp} = \text{yes} | \text{income} = \text{high}) = \frac{0}{2} = 0$$

$$\bullet P(\text{buys_comp} = \text{no} | \text{income} = \text{high}) = \frac{2}{2} = 1 \quad (2)$$

$$\textcircled{2} \quad G_{\text{Gini-Index}} = 1 - (0^2 + 1^2) = 0$$

$$\begin{aligned} \bullet P(\text{buys_comp} = \text{yes} | \text{income} = \text{medium}) &= \frac{1}{2} & \text{G}_{\text{Gini-Index}} &= 1 - (\frac{1}{2})^2 + (\frac{1}{2})^2 = \frac{1}{2} \\ \bullet P(\text{buys_comp} = \text{no} | \text{income} = \text{medium}) &= \frac{1}{2} & & \\ \bullet P(\text{buys_comp} = \text{yes} | \text{income} = \text{low}) &= 1 & \text{G}_{\text{Gini-Index}} &= 1 - (0^2 + 1^2) = 0 \\ \bullet P(\text{buys_comp} = \text{no} | \text{income} = \text{low}) &= 0 & & \end{aligned}$$

\Rightarrow Weighted Gini-Index for Income (2)

$$\textcircled{2} \quad \frac{2}{5} \cdot 0 + \frac{1}{5} \cdot \frac{1}{2} + \frac{1}{5} \cdot 0 = \frac{1}{5} = 0,2$$

Calculating the Gini-Index for student

$$\begin{aligned} \text{yes: } & 2 \text{ out of } 5 & P(\text{student} = \text{yes}) &= \frac{2}{5} \\ \text{no: } & 3 \text{ out of } 5 & P(\text{student} = \text{no}) &= \frac{3}{5} \end{aligned}$$

$$\begin{aligned} \bullet P(\text{buys_comp} = \text{yes} | \text{student} = \text{yes}) &= \frac{2}{2} = 1 \\ \bullet P(\text{buys_comp} = \text{no} | \text{student} = \text{yes}) &= \frac{0}{2} = 0 \quad \text{G}_{\text{Gini-Index}} = 1 - 1^2 + 0^2 = 0 \end{aligned}$$

$$\begin{aligned} \bullet P(\text{buys_comp} = \text{yes} | \text{student} = \text{no}) &= \frac{0}{3} = 0 \\ \bullet P(\text{buys_comp} = \text{no} | \text{student} = \text{no}) &= \frac{3}{3} = 1 \quad \text{G}_{\text{Gini-Index}} = 1 - (0^2 + 1^2) / 5 = 0 \end{aligned}$$

$\text{G}_{\text{Gini-Index}} = 1 - (0^2 + 1^2) / 5 = 0$

• Calculate the Gini-Index for Credit-Rating

$$\begin{array}{l} \text{fair: 3 out of 5} \\ \text{excellent: 2 out of 5} \end{array} \quad \left| \begin{array}{l} P(\text{credit=fair}) = \frac{3}{5} \\ P(\text{credit=excellent}) = \frac{2}{5} \end{array} \right.$$

$$\begin{array}{l} P(\text{buys=yes} | \text{credit=fair}) = \frac{1}{3} \\ P(\text{buys=no} | \text{credit=fair}) = \frac{2}{3} \end{array} \quad \left| \begin{array}{l} \text{Gini-Index} = 1 - \left(\left(\frac{1}{3} \right)^2 + \left(\frac{2}{3} \right)^2 \right) = \\ = 1 - \frac{5}{9} = \frac{4}{9} \end{array} \right.$$

$$\begin{array}{l} P(\text{buys=yes} | \text{credit=excellent}) = \frac{1}{2} \\ P(\text{buys=no} | \text{credit=excellent}) = \frac{1}{2} \end{array} \quad \left| \begin{array}{l} \text{Gini-Index} = 1 - \left(\frac{1}{2}^2 + \frac{1}{2}^2 \right) = \frac{1}{2} \end{array} \right. \Rightarrow \text{Weighted Gini Index for Credit-Rating} \quad (1) \\ (1) \quad \frac{3}{5} \cdot \frac{1}{3} + \frac{2}{5} \cdot \frac{1}{2} = \frac{12}{45} + \frac{1}{5} = \frac{2}{5} = \frac{12}{15} \approx 0.4666 \end{array}$$

Now, we need to calculate the information - Gain

Step 1: Calculate Entropy for parent

$$E(\text{parent}) = -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}} = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.97095$$

Step 2: Calculate Entropy for children

$$\begin{array}{l} \text{Income: } P(\text{income=high}) = \frac{2}{5} \\ \quad P(\text{income=medium}) = \frac{2}{5} \\ \quad P(\text{income=low}) = \frac{1}{5} \end{array}$$

High				Medium				Low			
Income	Student	Credit Rating	Buys Comp	Income	Student	Credit Rating	Buys Comp	Income	Student	Credit Rating	Buys Comp
high	no	fair	no	medium	no	fair	no	low	yes	fair	yes
high	no	excellent	no	medium	yes	excellent	yes	low	yes	fair	yes

$$\bullet \text{Entropy}(High) = -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}} = 1 \log 1 = 0$$

$$\begin{array}{l} P_{\text{yes}} = 0 \\ P_{\text{no}} = 1 \end{array}$$

$$\bullet \text{Entropy}(Medium) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1.$$

$$\begin{array}{l} P_{\text{yes}} = \frac{1}{2} \\ P_{\text{no}} = \frac{1}{2} \end{array}$$

$$\bullet E(\text{low}) = 1 \log 1 = 0$$

$$\begin{array}{l} P_{\text{yes}} = 1 \\ P_{\text{no}} = 0 \end{array}$$

$$\underline{\text{Step 3: Weighted }} E(High) = \frac{2}{5} \cdot 0 = 0$$

$$\text{Weighted } E(Medium) = \frac{2}{5} \cdot 1 = \frac{2}{5} \quad \left| \begin{array}{l} \text{Average Weighted Entropy}(Income) = 0 + \frac{2}{5} + 0 = \frac{2}{5} = 0.4 \end{array} \right.$$

$$\text{Weighted } E(Low) = \frac{1}{5} \cdot 0 = 0$$

$$\underline{\text{Step 4: }} Gini(\text{High}, \text{Income}) = Gini(Y) - \text{(Average Weighted } E(\text{Income})) = 0.97095 - 0.4 = \underline{0.57095}$$

Student: Step 2: Calculate Entropy for children

$$\begin{array}{l} \text{Student: } P(\text{student=yes}) = \frac{2}{5} \\ \quad P(\text{student=no}) = \frac{3}{5} \end{array}$$

Yes		No					
Income	Student	Credit Rating	Buys Comp	Income	Student	Credit Rating	Buys Comp
high	yes	fair	yes	high	no	fair	no
high	yes	excellent	yes	high	no	excellent	no
medium	yes	fair	no	medium	no	fair	no

Step 3:

$$\bullet E(\text{yes}) = -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}} = 1 \cdot \log 1 = 0$$

$$P_{\text{yes}} = \frac{2}{5} = 0.4$$

$$P_{\text{no}} = 0$$

$$\bullet E(\text{no}) = 1 \cdot \log 1 + 0 = 0$$

$$P_{\text{yes}} = 0$$

$$P_{\text{no}} = 1$$

$$\Rightarrow \text{Average Weighted } E(\text{student}) = \frac{2}{5} \cdot 0 + \frac{3}{5} \cdot 0 = 0$$

$$\underline{\text{Step 4: } I(\text{Gain}(\text{Youth}, \text{student})) = \text{Gain}(Y) - \text{Gain}(\text{student}) = 0.97095 - 0 = 0.97095}$$

Step 2: Calculate Entropy for children.

Credit-Rating

$$P(\text{credit=fair}) = \frac{3}{5}$$

$$P(\text{credit=excellent}) = \frac{2}{5}$$

(credit-rating)

(fair)

(excellent)

Income	student	credit rating	buys comp
high	no	fair	no
medium	no	fair	no
low	yes	fair	yes

Income	student	credit rating	buys comp
high	no	excellent	no
medium	yes	excellent	yes

Step 3:

$$E(\text{fair}) = -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}} = -\frac{1}{3} \log \frac{1}{3} - \frac{2}{3} \log \frac{2}{3} = 0.91829$$

$$P_{\text{yes}} = \frac{1}{3}$$

$$P_{\text{no}} = \frac{2}{3}$$

$$\Rightarrow \text{Average Weighted } E(\text{credit}) = 0.91829 \cdot \frac{3}{5} + \frac{2}{5} \cdot 1 = 0.950974$$

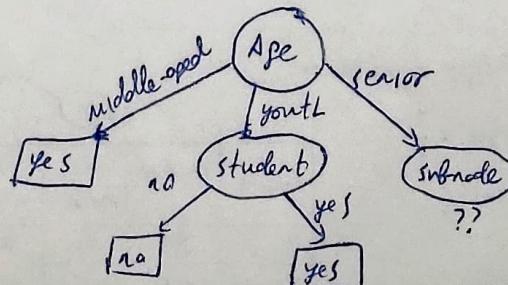
$$E(\text{excellent}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 0$$

$$\underline{\text{Step 4: } I(\text{Gain}(Y, \text{credit-rating})) = \text{Gain}(Y) - \text{Gain}(\text{Credit}) = 0.97095 - 0.950974 = 0.019976}$$

We got the following Gain-index and Infogain:

feature	Gini-index	Infogain
Income	0.2	0.57095
Student	0	0.97095
Credit-rat.	0.46667	0.019976

\Rightarrow student feature has the smallest Gini-index (Also the largest IG=0.97095), hence we take Student as the splitting attribute for the subnode Youth:



Also, we can say $P(\text{buys=yes} | (\text{age=youth} \text{ and } \text{student=yes})) = 1$

$P(\text{buys=no} | (\text{age=youth} \text{ and } \text{student=no})) = 0$ | \therefore we can construct the whole student sub-branch

Calculating Gini-Indices for (2) (age=Senior) sub-dataset

• Calculating the Gini-Index for income

$$\begin{array}{l} \text{high: 3 out of 5} \\ \text{medium: 3 out of 5} \\ \text{low: 2 out of 5} \end{array} \quad \left| \begin{array}{l} P(\text{income} = \text{high}) = \frac{3}{5} = 0.6 \\ P(\text{income} = \text{medium}) = 3/5 = 0.6 \\ P(\text{income} = \text{low}) = 2/5 = 0.4 \end{array} \right.$$

age	income	student	credit rating	buys computer
senior	medium	no	fair	yes
senior	low	yes	fair	yes
senior	low	yes	excellent	no
senior	medium	yes	fair	yes
senior	medium	no	excellent	no

- $P(\text{buys} = \text{yes}, \text{income} = \text{high}) = 0$ (not included in the dataset)
- $P(\text{buys} = \text{no}, \text{income} = \text{high}) = 0$

$$\begin{array}{l} \bullet P(\text{buys} = \text{yes} | \text{income} = \text{medium}) = \frac{2}{3} \\ \bullet P(\text{buys} = \text{no} | \text{income} = \text{medium}) = \frac{1}{3} \end{array} \quad \left| \begin{array}{l} \text{Gini-Index} = 1 - (1\frac{2}{3})^2 + (\frac{1}{3})^2 = 1 - \frac{4}{9} = \frac{5}{9} \\ \Rightarrow \text{Weighted-Gini-Index} = 0.6 \cdot \frac{5}{9} + 0.4 \cdot \frac{1}{9} \end{array} \right. \quad \textcircled{2}$$

$$\begin{array}{l} \bullet P(\text{buys} = \text{yes} | \text{income} = \text{low}) = \frac{1}{2} \\ \bullet P(\text{buys} = \text{no} | \text{income} = \text{low}) = \frac{1}{2} \end{array} \quad \left| \begin{array}{l} \text{Gini-Index} = 1 - (1\frac{1}{2})^2 + (\frac{1}{2})^2 = \frac{1}{2} \\ \Rightarrow \frac{3}{5} \cdot \frac{5}{9} + \frac{2}{5} \cdot \frac{1}{2} = \frac{4}{15} + \frac{1}{5} = \frac{7}{15} \approx 0.46667 \end{array} \right. \quad \textcircled{2}$$

• Calculating the Gini-Index for student

$$\begin{array}{l} \text{yes: 3 out of 5} \\ \text{no: 2 out of 5} \end{array} \quad \left| \begin{array}{l} P(\text{student} = \text{yes}) = \frac{3}{5} \\ P(\text{student} = \text{no}) = \frac{2}{5} \end{array} \right.$$

$$\begin{array}{l} \bullet P(\text{buys} = \text{yes} | \text{student} = \text{yes}) = \frac{2}{3} \\ \bullet P(\text{buys} = \text{no} | \text{student} = \text{yes}) = \frac{1}{3} \end{array} \quad \left| \begin{array}{l} \text{Gini-Index} = 1 - (1\frac{2}{3})^2 + (\frac{1}{3})^2 = \frac{5}{9} \\ \Rightarrow \text{Weighted-Gini-Index} = 0.6 \cdot \frac{5}{9} + 0.4 \cdot \frac{1}{9} \end{array} \right. \quad \textcircled{2}$$

$$\begin{array}{l} \bullet P(\text{buys} = \text{yes} | \text{student} = \text{no}) = \frac{1}{2} \\ \bullet P(\text{buys} = \text{no} | \text{student} = \text{no}) = \frac{1}{2} \end{array} \quad \left| \begin{array}{l} \text{Gini-Index} = 1 - (1\frac{1}{2})^2 + (\frac{1}{2})^2 = \frac{1}{2} \\ \Rightarrow \frac{3}{5} \cdot \frac{5}{9} + \frac{2}{5} \cdot \frac{1}{2} = \frac{4}{15} + \frac{1}{5} = \frac{7}{15} \approx 0.46667 \end{array} \right. \quad \textcircled{2}$$

$$\textcircled{2} \quad \text{Weighted Gini-Index for student} = \frac{3}{5} \cdot \frac{5}{9} + \frac{2}{5} \cdot \frac{1}{2} = \frac{4}{15} + \frac{1}{5} = \frac{7}{15} \approx 0.46667.$$

• Calculating the Gini-Index for Credit-Rating

$$\begin{array}{l} \text{fair: 3 out of 5} \\ \text{excellent: 2 out of 5} \end{array} \quad \left| \begin{array}{l} P(\text{credit} = \text{fair}) = \frac{3}{5} \\ P(\text{credit} = \text{excellent}) = \frac{2}{5} \end{array} \right.$$

$$\begin{array}{l} \bullet P(\text{buys} = \text{yes} | \text{credit} = \text{fair}) = \frac{3}{3} = 1 \\ \bullet P(\text{buys} = \text{no} | \text{credit} = \text{fair}) = \frac{0}{3} = 0 \end{array} \quad \left| \begin{array}{l} \text{Gini-Index} = 1 - 0^2 - 1^2 = 0 \\ \Rightarrow \text{Weighted-Gini-Index} = 0.6 \cdot 0 + 0.4 \cdot 0 = 0 \end{array} \right. \quad \textcircled{2}$$

$$\begin{array}{l} \bullet P(\text{buys} = \text{yes} | \text{credit} = \text{excellent}) = 0 \\ \bullet P(\text{buys} = \text{no} | \text{credit} = \text{excellent}) = 1 \end{array} \quad \left| \begin{array}{l} \text{Gini-Index} = 1 - 0^2 - 1^2 = 0 \\ \Rightarrow \text{Weighted-Gini-Index} = 0.6 \cdot 0 + 0.4 \cdot 0 = 0 \end{array} \right. \quad \textcircled{2}$$

$$\textcircled{2} \quad \text{Weighted Gini-Index for credit-rating} = \frac{3}{5} \cdot 0 + \frac{2}{5} \cdot 0 = 0$$

Now, we calculate Information Gains for Senior-Dataset's features.

Step 1: Calculate entropy for parent:

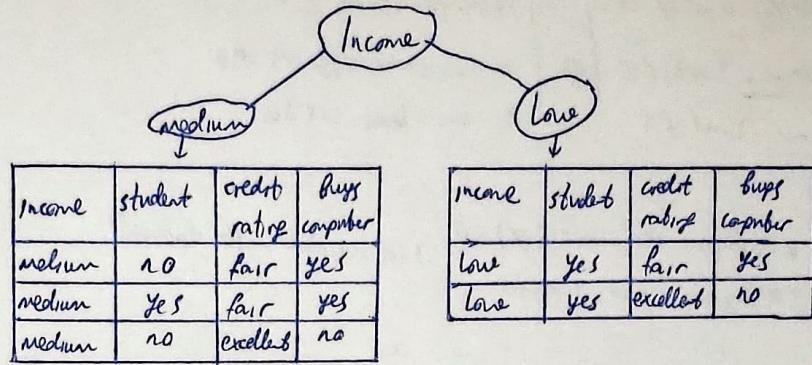
$$\text{Entropy}(\text{parent}) = -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}} = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \approx 0.97095$$

$$P_{\text{yes}} = \frac{3}{5}$$

$$P_{\text{no}} = \frac{2}{5}$$

Step 2: Calculate for children.

$$\begin{aligned} & P(\text{income}=\text{high})=0 \\ \underline{\text{Income:}} \quad & P(\text{income}=\text{medium}) = \frac{3}{5} \\ & P(\text{income}=\text{low}) = \frac{2}{5} \end{aligned}$$



$$\bullet \text{Entropy}(h)=0$$

$$\bullet \text{Entropy}(\text{medium}) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \approx 0.91829$$

$$P_{\text{yes}} = \frac{2}{3}$$

$$P_{\text{no}} = \frac{1}{3}$$

$$\bullet \text{Entropy}(l) = -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}} = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

Step 3:

$$\text{Weighted } E(\text{medium}) = \frac{3}{5} \cdot 0.91829 \approx 0.550974$$

$$\text{Weighted } E(l) = \frac{2}{5} \cdot 1 = \frac{2}{5}$$

$$\Rightarrow \text{Average Weighted Entropy}(\text{income}) = 0.550974 + 0.4 = 0.950974$$

Step 4: Calculate Information Gain

$$I.\text{Gain}(\text{Senior}, \text{income}) = \text{Gain}(S) - \{\text{Average Weighted Entropy}\} = 0.97095 - 0.950974 = 0.019976$$

Step 2: Calculate entropy for children

$$\underline{\text{Student:}} \quad P(\text{student}=\text{yes}) = \frac{3}{5}$$

$$P(\text{student}=\text{no}) = \frac{2}{5}$$

income	student	credit rating	buys computer
low	yes	fair	yes
low	yes	excellent	no
medium	yes	fair	yes

income	student	credit rating	buys computer
medium	no	fair	yes
medium	no	excellent	no

$$\bullet \text{Entropy}(\text{yes}) = -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}} = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \approx 0.91829$$

$$P_{\text{yes}} = \frac{2}{3}$$

$$P_{\text{no}} = \frac{1}{3}$$

$$\bullet \text{Entropy}(\text{no}) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$P_{\text{yes}} = \frac{1}{2}$$

$$P_{\text{no}} = \frac{1}{2}$$

Step 3:

$$\text{Weighted } E(\text{yes}) = \frac{3}{5} \cdot 0.91829 \approx 0.550974$$

$$\text{Weighted } E(\text{no}) = \frac{2}{5} \cdot 1 = \frac{2}{5}$$

$$\Rightarrow \text{Average Weighted Entropy}(\text{student}) = 0.550974 + 0.4 = 0.950974$$

Step 4: Calculate Information Gain.

$$I_{\text{Gain}}(\text{senior, student}) = \text{Gain}(S) - (\text{Average Weighted El}_{\text{student}}) = 0,97095 - 0,950974 = 0,019976$$

Step 5: Calculate Entropy for children

Credit-rating $P(\text{credit}=\text{fair}) = \frac{3}{5}$
 $P(\text{credit}=\text{excellent}) = \frac{2}{5}$

income	student	credit rating	buys computer
medium	no	fair	yes
low	yes	fair	yes
medium	yes	fair	yes

income	student	credit rating	buys computer
low	yes	excellent	no
medium	no	excellent	no

• Entropy(fair) = $1 \log 1 = 0$

$P_{\text{yes}} = 1$

$P_{\text{no}} = 0$

• Entropy(excellent) = 0

$P_{\text{no}} = 1$

$P_{\text{yes}} = 0$

Step 6:

$\Rightarrow \text{Average Weighted Entropy(credit)} = \frac{3}{5} \cdot 0 + \frac{2}{5} \cdot 0 = 0$

$\text{Weighted El(fair)} = 0$

$\text{Weighted El(excellent)} = 0$

Step 7: Calculate Information Gain

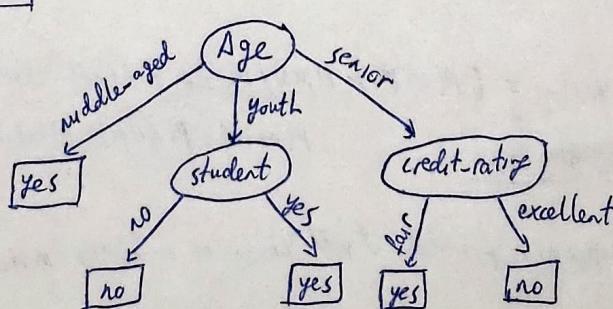
$\text{Information-Gain}(\text{senior, credit-rating}) = \text{Gain}(S) - (\text{Average Weighted El}_{\text{credit}}) \quad (2)$

$\Rightarrow 0,97095 - 0 = 0,97095$

We got the following Gini-Indices and Information Gains

feature	Gini-Index	Info-Gain
income	0,66667	0,019976
student	0,66667	0,019976
credit-rating	0	0,97095

Credit-Rating feature has the smallest Gini-Index = 0 (also the largest Information Gain), therefore we use credit-rating as the splitting attribute for the subnode senior:



Final Decision tree

From the table, we can say: $P(\text{buys}=\text{yes} | \text{age}=\text{senior} \text{ and } \text{credit}=\text{yes}) = 1$ / s) we can build (fully) the subbranch credit-rating.
 $P(\text{buys}=\text{no} | \text{age}=\text{senior} \text{ and } \text{credit}=\text{no}) = 1$

Answer: We got the decision tree!

8)

$$P(\text{buys_computer} = \text{yes}) = \frac{9}{14}$$

$$P(\text{buys_computer} = \text{no}) = \frac{5}{14}$$

The prior probabilities here are the probabilities when $\text{buys_comp} = \text{yes}$ and when $\text{buys_computer} = \text{no}$:

Prior probabilities:

$$P(\text{buys_computer} = \text{yes}) = \frac{9}{14}$$

$$P(\text{buys_computer} = \text{no}) = \frac{5}{14}$$

Now, we need to classify a test instance.

Denote $x = \{\text{age} = A_1, \text{income} = A_2, \text{student} = A_3, \text{credit_rating} = A_4\}$

Since A_1 represents age $\Rightarrow A_1 \in \{\text{youth, middle-aged, senior}\}$

With the same principle: $A_2 \in \{\text{high, medium}\}$

$$A_3 \in \{\text{yes, no}\}$$

$$A_4 \in \{\text{fair, excellent}\}$$

Now, we need to calculate posterior probabilities i.e. calculate $P(\text{yes}|x)$ and $P(\text{no}|x)$.

The posterior probability $P(\text{yes}|x)$ is calculated using Bayes' Theorem, which in Naïve Bayes simplifies

$$\text{to } P(\text{yes}|x) = \frac{P(x|\text{yes}) \cdot P(\text{yes})}{P(x)}$$

$$P(\text{yes}|x) = P(\text{buys_computer} = \text{yes}|x) = P(x|\text{yes}) \cdot P(\text{yes}) / P(x) = P(A_1|\text{yes}) \cdot P(A_2|\text{yes}) \cdot P(A_3|\text{yes}) \cdot P(A_4|\text{yes}) / P(x)$$

Bayes

$$P(\text{no}|x) = P(A_1|\text{no}) \cdot P(A_2|\text{no}) \cdot P(A_3|\text{no}) \cdot P(A_4|\text{no}) / P(x)$$

Since no specific test instance is provided, I will choose an arbitrary instance:

$$\text{let } x = \{\text{age} = \text{youth}, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit_rating} = \text{fair}\}$$

$/P(x)$

$$\text{(Calculate } P(\text{yes}|x) : P(\text{yes}|x) = P(\text{youth}|\text{yes}) \cdot P(\text{medium}|\text{yes}) \cdot P(\text{yes}|\text{yes}) \cdot P(\text{fair}|\text{yes}) \cdot P(\text{yes}) = \frac{2}{9} \cdot \frac{4}{9} \cdot \frac{6}{9} \cdot \frac{6}{9} \cdot \frac{9}{14} \approx 0,0282 / P(x))$$

$$\text{(Calculate } P(\text{no}|x) : P(\text{no}|x) = P(\text{youth}|\text{no}) \cdot P(\text{medium}|\text{no}) \cdot P(\text{yes}|\text{no}) \cdot P(\text{fair}|\text{no}) \cdot P(\text{no}) = \frac{3}{5} \cdot \frac{2}{5} \cdot \frac{1}{5} \cdot \frac{5}{9} \cdot \frac{5}{14} \approx 0,0069 / P(x))$$

$$P(\text{yes}|x) = 0,0282 / P(x) > 0,0069 / P(x) = P(\text{no}|x) \Rightarrow \text{Classify } x \text{ as yes (in this case).}$$

Age	YES	NO
youth	2/9	3/5
middle-aged	4/9	0/5
senior	3/9	2/5

Income	YES	NO
High	2/9	2/5
Medium	4/9	2/5
Low	3/9	1/5

Student	YES	NO
yes	6/9	1/5
No	3/9	4/5

Credit Rating	YES	NO
fair	6/9	2/5
Excellent	3/9	3/5

The above tables are tables of conditional probabilities.

$$\text{For example, for age: } P(\text{age} = \text{youth} | \text{buys_yes}) = \frac{2}{9}$$

$$P(\text{age} = \text{youth} | \text{buys_no}) = \frac{3}{5}$$

$$P(\text{age} = \text{middle} | \text{buys_yes}) = 4/9$$

$$P(\text{age} = \text{middle} | \text{buys_no}) = 0$$

$$P(\text{age} = \text{senior} | \text{buys_yes}) = 3/9$$

$$P(\text{age} = \text{senior} | \text{buys_no}) = 2/5$$

c) RIPPER is a rule-based classification algorithm that generates a set of IF-THEN rules to make prediction. It is particularly effective for datasets with categorical attributes and binary class labels. RIPPER creates α rules for one class by adding conditions that maximize accuracy, then it simplifies each rule by removing conditions that don't contribute to predictive performance. The algorithm repeatedly refines the rules to reduce overfitting and improve generalization.

RIPPER is especially useful when interpretability and concise rule sets are needed. Decision Trees provide a visual and hierarchical understanding of decision-making. Naive Bayes is efficient and performs surprisingly well on many datasets, though it relies on the assumption that attributes are independent.

In our case, we have 14 records in our dataset, that is very small, all attributes are categorical, the target is binary, therefore RIPPER will perform the best because it's designed for rule-based classification on categorical data. It might extract rules like:

If Age = middle-aged THEN buys-computer = yes
or IF student = yes AND credit-rating = fair THEN buys-computer = yes

The second place is Decision Tree. It is still performing well on categorical and small datasets, but without careful pruning, the tree might overfit the small training set.

The third place is Naive Bayes: It assumes independence among features, which doesn't hold for this dataset. Also, the dataset is small, which can make the model unstable.

Problem 6 In the left figure, the model is too simple. It only uses three trees and has a very high learning rate = 1. As a result, it seems to jump too quickly to conclusions and doesn't capture the pattern in the data very well. The predictions are rigid which misses the smooth curvature that the data has. This is underfitting i.e. the model isn't had enough complexity or time to learn the structure in the data. Possible solution is to increase the number of trees, for example use 50 or 100 trees to learn more complex patterns and lowering the learning rate (to 0.1 or 0.05) to allow for more refined learning.

In the right figure, the model does a good job at first place. It uses a low learning rate = 0.1 and a large number of estimators (trees) = 200, which allows it to gradually fit the data. But the model might be overfitting. We see it because it picks up the noise in the data, we can lower the number of trees a bit (making 100-150 trees).

Problem 7 Ensemble classification combines the predictions of multiple base models to improve overall performance. In Ensemble classification we use Hard Voting and Soft Voting.

In Hard Voting, each base classifier predicts a class label, and the final prediction is the one that receives the most votes. Each model makes its own independent decision and vote for a class.

- 1) Most voted class is the winner
- 2) If there is a tie, then soft voting can be used, or weights can be assigned to the classifiers or odd number of classifiers can be used in ensemble.

Soft Voting is an ensemble technique for classification where multiple models don't vote on class labels, instead they provide the class probabilities for each class.

The final prediction is made by averaging the predicted probabilities for each class and selecting the class with the highest average.

Bagging stands for Bootstrap Aggregating, and it's an ensemble method that builds multiple methods on random subsets of the training data with replacement.

First, it randomly samples with replacement from the training set to create several new subsets of the data. Second, train a separate model on each sample of the dataset. Further, it aggregates their predictions. Since sampling is with replacement, some samples may appear multiple times, and some may be left out in each subset. Its advantages are: 1) Reduces overfitting 2) Easy to parallelize 3) Stable and robust.

One of popular examples is Random Forest.

Boosting refers to any ensemble method that can combine several weak learners into a strong learner. The general idea of most boosting methods is to train predictors sequentially, each trying to correct its predecessor.

The most popular boosting methods are AdaBoost, Gradient Boosting and Extreme gradient boosting.

The process

- 1) Train the first model on the training data.
- 2) Evaluate Errors and give higher weights to misclassified samples.
- 3) Train the next model to focus more on difficult examples
- 4) Combine models using weighted voting or gradient optimization.

key feature is that the models are trained sequentially, each one correcting its predecessor.

Its advantages include high predictive power, reduced both bias and variance, often winning all competitors.

Problem 8: Naive Bayes is a probabilistic classification algorithm based on Bayes' Theorem with the assumption that features are conditionally independent given the class. The theorem is:

$$\text{Bayes Theorem: } P(y|x) = \frac{P(x|y) \cdot P(y)}{P(x)}$$

where $P(y|x)$ is the posterior probability of y given x .

$P(y)$ = prior probability of class y

$P(x)$ = marginal probability of feature x .

Prior Probability is the initial probability of an event before seeing any evidence. In our case, it's the probability of "play-footfall=Yes" or "play-footfall=no" without considering any weather conditions.

Posterior Probability is the probability of an event after any evidence. In our case, it's the probability of "play-Football=Yes" or "play-Football=no" with given specific conditions.

To find the posterior probability of x ,

$$x = \{\text{outlook=sunny, Temp=Hot, Humidity=High, Wind=Weak}\}$$

$$\text{We need to calculate } P(\text{yes}|x) = P(x|\text{yes}) \cdot P(\text{yes}) / P(x) = P(\text{sunny}|\text{yes}) \cdot P(\text{Hot}|\text{yes}) \cdot P(\text{High}|\text{yes}) \cdot P(\text{Weak}|\text{yes}) / P(x)$$

By Bayes Theorem

we don't need to find $P(x)$ because it is in the product of both $P(\text{yes}|x)$ and $P(x|\text{yes})$ and we can compare them.

$$P(\text{sunny}|\text{yes}) = \frac{\#(\text{outlook=sunny} \cap \text{play=yes})}{\#(\text{play=yes})} = \frac{1}{3}$$

$$P(\text{Hot}|\text{yes}) = \frac{\#(\text{Temp=Hot} \cap \text{play=yes})}{\#(\text{play=yes})} = \frac{0}{3} = 0$$

$$P(\text{High}|\text{yes}) = \frac{\#(\text{Humidity=High} \cap \text{play=yes})}{\#(\text{play=yes})} = \frac{1}{3}$$

$$P(\text{Weak}|\text{yes}) = \frac{\#(\text{Wind=Weak} \cap \text{play=yes})}{\#(\text{play=yes})} = \frac{2}{3}$$

$$\Rightarrow P(\text{yes}|x) = \frac{1}{3} \cdot 0 \cdot \frac{1}{3} \cdot \frac{2}{3} \cdot P(\text{yes}) \text{ if } P(x) = 0 \Rightarrow P(\text{yes}|x) = 0$$

$$P(\text{No}|x) = P(\text{sunny}|x) \cdot P(\text{Hot}|x) \cdot P(\text{High}|x) \cdot P(\text{Weak}|x) \cdot P(x) = \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{1}{3} \cdot P(\text{No}) / P(x) = \frac{1}{32} \cdot \frac{4}{7} / P(x) = \frac{1}{56} / P(x) \quad (1)$$

$$P(\text{sunny}|x) = \frac{2}{3}$$

$$P(\text{Hot}|x) = \frac{2}{3}$$

$$P(\text{High}|x) = \frac{2}{3}$$

$$P(\text{Weak}|x) = \frac{1}{3}$$

$$(P(\text{No}) = \frac{4}{7}, P(\text{yes}) = \frac{3}{7})$$

$$(2) P(\text{No}|x) = \frac{1}{56} / P(x)$$

In our case $P(\text{No}|x) = \frac{1}{56} / P(x) > P(\text{yes}|x) \Rightarrow$ the Naive Bayes classifier will predict "No" for x .

Problem 9We will use the Hamming Distance

It's the count of mismatches.

$$d_H(x, y) = \sum_{i=1}^n \delta(x_i, y_i)$$

where $\delta(x_i, y_i) = \begin{cases} 0 & \text{if } x_i = y_i \\ 1 & \text{if } x_i \neq y_i \end{cases}$

- $k=2$: We will chose random initial centroids:
- Cluster 1: Row 1 - (Sunny, Hot, High, Weak) = c_1 ,
 - Cluster 2: Row 3 - (Overcast, Cold, Normal, Strong) = c_2

- Now, we need to calculate the distances between the observations and assign each one its cluster.

Row 1: $d_H(c_1, x_1) = 0$ (c_1 := cluster 1)
 $d_H(c_2, x_1) = 4+1+1+1 = 7$

$| 0 < 7 \Rightarrow \text{Assign } x_1 \text{ to } c_1$

Row 2: $d_H(c_1, x_2) = 0+1+1+0 = 2$
 $d_H(c_2, x_2) = 1+1+0+1 = 3$

$| \Rightarrow 2 < 3 \Rightarrow \text{Assign } x_2 \text{ to } c_1$
 $(c_2 := \text{cluster 2})$
 (Row 3)

Row 3: $d_H(c_1, x_3) = 1+1+1+1 = 4$
 $d_H(c_2, x_3) = 0+0+0+0 = 0$

$| \Rightarrow 0 < 4 \Rightarrow \text{Assign } x_3 \text{ to } c_2$

Row 4: $d_H(c_1, x_4) = 0+1+1+1 = 3$
 $d_H(c_2, x_4) = 1+1+0+0 = 2$

$| \Rightarrow 2 < 3 \Rightarrow \text{Assign } x_4 \text{ to } c_2$

Row 5: $d_H(c_1, x_5) = 1+1+0+0 = 2$
 $d_H(c_2, x_5) = 1+1+1+1 = 4$

$| \Rightarrow 2 < 4 \Rightarrow \text{Assign } x_5 \text{ to } c_1$

Row 6: $d_H(c_1, x_6) = 1+0+1+1 = 3$
 $d_H(c_2, x_6) = 0+1+0+0 = 1$

$| \Rightarrow 1 < 3 \Rightarrow \text{Assign } x_6 \text{ to } c_2$

Row 7: $d_H(c_1, x_7) = 1+1+0+1 = 3$
 $d_H(c_2, x_7) = 1+1+1+0 = 3$

$| \Rightarrow 3 = 3 \Rightarrow \text{let's assign to } c_1$

As a result, we got

- Cluster 1: Rows: 1, 2, 5, 7
- Cluster 2: Rows: 3, 4, 6

- Now, we need to update centroids: For cluster 1 and 2 we got the following tables

rows	outlook	Temperature	Humidity	Wind
1	Sunny	Hot	High	Weak
2	Sunny	Normal	Normal	Weak
5	Rainy	Normal	High	Weak
7	Rainy	Normal	High	Strong

rows	outlook	Temperature	Humidity	Wind
3	Overcast	Cold	Normal	Strong
4	Sunny	Normal	Normal	Strong
6	Overcast	Hot	Normal	Strong

From the tables in the previous page, we need to update centroids.

Cluster 1: outlook: {sunny, sunny, rainy, rainy} \Rightarrow Mode = sunny
Temperature: {Hot, warm, warm, warm} \Rightarrow Mode = warm
Humidity: {High, Normal, High, High} \Rightarrow Mode = High
Wind: {Weak, weak, weak, strong} \Rightarrow Mode = weak
 \Rightarrow New C1 = (Sunny, Warm, High, Weak)

Cluster 2: outlook: {overcast, sunny, overcast, } \Rightarrow Mode = overcast
Temperature: {Cold, warm, Hot} \Rightarrow let's pick Mode = Cold
Humidity: {Normal, normal, normal} \Rightarrow Mode = Normal
Wind: {Strong, Strong, Strong} \Rightarrow Mode = Strong
 \Rightarrow New C2 = {Overcast, Cold, Normal, Strong}

• Reassign points

Row 1: $d_n(c_1, x_1) = 0 + 1 + 0 + 0 = 1$ | $\Rightarrow 1 < 3 \Rightarrow$ Assign x_1 to C_1
 $d_n(c_2, x_1) = 1 + 1 + 1 = 3$

Row 2: $d_n(c_1, x_2) = 0 + 0 + 1 + 0 = 1$ | $\Rightarrow 1 < 3 \Rightarrow$ Assign x_2 to C_1
 $d_n(c_2, x_2) = 1 + 1 + 0 + 1 = 3$

Row 3: $d_n(c_1, x_3) = 1 + 1 + 1 + 1 = 4$ | $\Rightarrow 4 > 3 \Rightarrow$ Assign x_3 to C_2
 $d_n(c_2, x_3) = 0 + 0 + 0 + 0 = 0$

Row 4: $d_n(c_1, x_4) = 0 + 0 + 1 + 1 = 2$ | $\Rightarrow 2 < 3 \Rightarrow$ let's assign x_4 to C_2 (original cluster)
 $d_n(c_2, x_4) = 1 + 1 + 0 + 0 = 2$

Row 5: $d_n(c_1, x_5) = 1 + 0 + 0 + 0 = 1$ | $\Rightarrow 1 < 3 \Rightarrow$ Assign x_5 to C_1
 $d_n(c_2, x_5) = 1 + 1 + 1 + 1 = 4$

Row 6: $d_n(c_1, x_6) = 1 + 1 + 1 + 1 = 4$ | $\Rightarrow 4 > 3 \Rightarrow$ Assign x_6 to C_2
 $d_n(c_2, x_6) = 0 + 1 + 0 + 0 = 1$

Row 7: $d_n(c_1, x_7) = 1 + 0 + 0 + 1 = 2$ | $\Rightarrow 2 < 3 \Rightarrow$ Assign x_7 to C_1
 $d_n(c_2, x_7) = 1 + 1 + 1 + 0 = 3$

As a result: we got Cluster 1: 1, 2, 5, 7 | \Rightarrow same as previous iteration \Rightarrow Algorithm has converged
Cluster 2: 3, 4, 6

The final clusters for k=2

Cluster 1

Row#	Outlook	Temperature	Humidity	Wind
1	Sunny	Hot	High	Weak
2	Sunny	Warm	Normal	Weak
5	Rainy	Warm	High	Weak
7	Rainy	Warm	High	Strong

$$\text{last centroid } c_1 = \{\text{Sunny, Warm, High, Weak}\}$$

Cluster 2

Row#	outlook	temperature	humidity	wind
3	overcast	Cold	Normal	Strong
4	Sunny	Warm	Normal	Strong
6	overcast	Hot	Normal	Strong

$$\text{last centroid } c_2 = \{\text{Overcast, Cold, Normal, Strong}\}$$

k=3: We choose random initial centroids :

$$\text{Centroid for Cluster 1 : } c_1 = \text{Row} 1 = \{\text{Sunny, Hot, High, Weak}\}$$

$$\text{Centroid for Cluster 2 : } c_2 = \text{Row} 3 = \{\text{Overcast, Cold, Normal, Strong}\}$$

$$\text{Centroid for Cluster 3 : } c_3 = \text{Row} 5 = \{\text{Rainy, Warm, High, Weak}\}$$

- Calculate the distances and assign clusters.

$$\begin{aligned} \text{Row 1: } d_H(c_1, x_1) &= 0+0+0+0=0 \\ d_H(c_2, x_1) &= 1+1+1+1=4 \quad | \Rightarrow 0 < 1 < 4 \Rightarrow \text{Assign } x_1 \text{ to } c_1 \\ d_H(c_3, x_1) &= 1+1+0+0=2 \end{aligned}$$

$$\begin{aligned} \text{Row 2: } d_H(c_1, x_2) &= 0+1+1+0=2 \\ d_H(c_2, x_2) &= 1+1+0+1=3 \quad | \Rightarrow 2 < 3 \Rightarrow \text{Assign } x_2 \text{ to } c_1 \\ d_H(c_3, x_2) &= 1+0+1+0=2 \end{aligned}$$

$$\begin{aligned} \text{Row 3: } d_H(c_1, x_3) &= 1+1+1+1=4 \\ d_H(c_2, x_3) &= 0+0+0+0=0 \quad | \Rightarrow 0 < 1 < 4 \Rightarrow \text{Assign } x_3 \text{ to } c_2 \\ d_H(c_3, x_3) &= 1+0+1+1=3 \end{aligned}$$

$$\begin{aligned} \text{Row 4: } d_H(c_1, x_4) &= 0+1+1+1=3 \\ d_H(c_2, x_4) &= 1+1+0+0=2 \quad | \Rightarrow 2 < 3 \Rightarrow \text{Assign } x_4 \text{ to } c_2 \\ d_H(c_3, x_4) &= 1+0+1+1=3 \end{aligned}$$

$$\begin{aligned} \text{Row 5: } d_H(c_1, x_5) &= 1+1+0+0=2 \\ d_H(c_2, x_5) &= 1+1+1+1=4 \quad | \Rightarrow 0 < 2 < 4 \Rightarrow \text{Assign } x_5 \text{ to } c_3 \\ d_H(c_3, x_5) &= 0+0+0+0=0 \end{aligned}$$

$$\begin{aligned} \text{Row 6: } d_H(c_1, x_6) &= 1+0+1+1=3 \\ d_H(c_2, x_6) &= 0+1+0+0=1 \quad | \Rightarrow 1 < 3 < 4 \Rightarrow \text{Assign } x_6 \text{ to } c_2 \\ d_H(c_3, x_6) &= 1+1+1+1=4 \end{aligned}$$

$$\begin{aligned} \text{Row 7: } d_H(c_1, x_7) &= 1+1+0+1=3 \\ d_H(c_2, x_7) &= 1+1+1+0=3 \quad | \Rightarrow 1 < 3 < 3 \\ d_H(c_3, x_7) &= 0+0+0+1=1 \end{aligned}$$

② Assign x_7 to c_3

As a result for I iteration, we get:

Cluster 1 Points {2}

Cluster 2: Points : 3,4,6

Cluster 3 Points 5,7

We got the following clusters so far:

Outlook	Temperature	Humidity	Wind
Sunny	Hot	High	Weak
Sunny	Warm	Normal	Weak

cluster 1

Outlook	Temperature	Humidity	Wind
Overcast	Cold	Normal	Strong
Sunny	Warm	Normal	Strong
Overcast	Hot	Normal	Strong

cluster 2

Outlook	Temperature	Humidity	Wind
Rainy	Warm	High	Weak
Rainy	Warm	High	Strong

cluster 3

- Now, we update the centroids.

Cluster 1: Outlook: {Sunny, Sunny} \Rightarrow Mode = Sunny

Temperature: { Hot, warm } \Rightarrow let's choose Mode = Hot

Humidity: { High, Normal } \Rightarrow let's choose Mode = High

Wind: { Weak, weak } \Rightarrow Mode = Weak

\Rightarrow Updated $C_1 = \{ \text{Sunny, Hot, High, Weak} \}$ (3)

(3) updated C_1 is the same as C_1 (unchanged)

Cluster 2: Outlook: { Overcast, sunny, overcast } \Rightarrow Mode = Overcast

Temperature: { Cold, Warm, Not } \Rightarrow let's choose Mode = Cold (as in previous)

Humidity: { Normal, Normal, Normal } \Rightarrow Mode = Normal

Wind: { Strong, Strong, Strong } \Rightarrow Mode = Strong

\Rightarrow Updated $C_2 = \{ \text{Overcast, Cold, Normal, Strong} \}$ \Rightarrow updated C_2 is equal to the previous C_2 (unchanged)

Cluster 3: Outlook: { Rainy, Rainy } \Rightarrow Mode = Rainy

Temperature: { Warm, warm } \Rightarrow Mode = Warm

Humidity: { High, High } \Rightarrow Mode = High

Wind: { Weak, Strong } \Rightarrow Mode = Weak

\Rightarrow Updated $C_3 = \{ \text{Rainy, Warm, High, Weak} \}$ (3)

(3) updated C_3 is the same as C_3 (unchanged)

- We got, that after the first iteration, the cluster centroids remain unchanged from our initial selection, therefore we can say that the algorithm converges!

The final clusters

cluster 2

Outlook	Temperature	Humidity	Wind
Overcast	Cold	Normal	Strong
Sunny	Warm	Normal	Strong

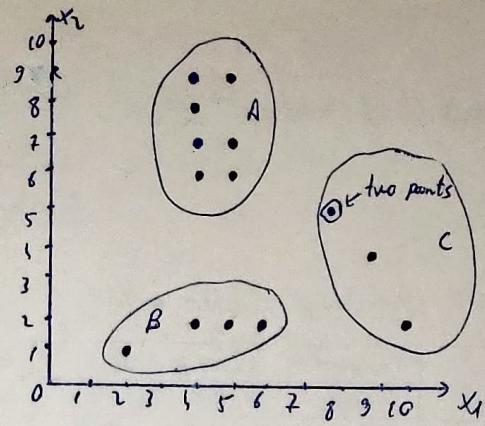
cluster 3

Outlook	Temperature	Humidity	Wind
Rainy	Warm	High	Weak
Rainy	Warm	High	Strong

Problem 10

To choose the appropriate numbers of clusters k , we can first plot the data points and calculate approximate number of groupings. (since the data is small and is from \mathbb{R}^2)

From the graph, we can assume that $k=3$ is a good choice (at first glance)



Step 1: Random Initialization of Centroids.

$$\text{Centroid } A = (4, 8)$$

$$\text{Centroid } B = (4, 2)$$

$$\text{Centroid } C = (9, 4)$$

Now, we need to calculate the distance (Euclidean) of each data points from centroids and do the grouping.

$$1) a_1 = (10, 2) \quad d(A, a_1) = \sqrt{(10-4)^2 + (2-8)^2} = \sqrt{72} = 6\sqrt{2}$$

$$d(B, a_1) = \sqrt{(10-4)^2 + (2-2)^2} = \sqrt{36} = 6 \quad | \Rightarrow 6 < 6 < 6\sqrt{2} \Rightarrow \text{Assign } (10, 2) \text{ to } C$$

$$d(C, a_1) = \sqrt{(10-9)^2 + (2-4)^2} = \sqrt{5}$$

$$2) a_2 = (9, 6) \quad d(A, a_2) = \sqrt{(9-4)^2 + (6-8)^2} = \sqrt{41}$$

$$d(B, a_2) = \sqrt{(9-4)^2 + (6-2)^2} = \sqrt{29} \quad | \Rightarrow 0 < \sqrt{29} < \sqrt{41} \Rightarrow \text{Assign } (9, 6) \text{ to } C$$

$$d(C, a_2) = \sqrt{(9-9)^2 + (6-4)^2} = 0$$

$$3) a_3 = a_4 = (8, 5) \quad d(A, a_3) = \sqrt{(8-4)^2 + (5-8)^2} = 5$$

$$d(B, a_3) = \sqrt{(8-4)^2 + (5-2)^2} = 5 \quad | \Rightarrow \sqrt{2} < 5 = 5 \Rightarrow \text{Assign } (8, 5) \text{ to } C$$

$$d(C, a_3) = \sqrt{(8-9)^2 + (5-4)^2} = \sqrt{2}$$

$$4) a_5 = (6, 2) \quad d(A, a_5) = \sqrt{(6-4)^2 + (2-8)^2} = \sqrt{40}$$

$$d(B, a_5) = \sqrt{(6-4)^2 + (2-2)^2} = 2 \quad | \Rightarrow 2 < \sqrt{40} < \sqrt{40} \Rightarrow \text{Assign } a_5 \text{ to } B$$

$$d(C, a_5) = \sqrt{(6-9)^2 + (2-4)^2} = \sqrt{29}$$

$$5) a_6 = (5, 7) \quad d(A, a_6) = \sqrt{(5-4)^2 + (7-8)^2} = \sqrt{2}$$

$$d(B, a_6) = \sqrt{(5-4)^2 + (7-2)^2} = 1 \quad | \Rightarrow 1 < \sqrt{2} < \sqrt{7} \Rightarrow \text{Assign } a_6 \text{ to } B$$

$$d(C, a_6) = \sqrt{(5-9)^2 + (7-4)^2} = \sqrt{20}$$

$$6) a_7 = (4, 2) \quad d(A, a_7) = \sqrt{(4-4)^2 + (2-8)^2} = 6$$

$$d(B, a_7) = \sqrt{(4-4)^2 + (2-2)^2} = 0 \quad | \Rightarrow 0 < 6 < 6 \Rightarrow \text{Assign } (4, 2) \text{ to } B$$

$$d(C, a_7) = \sqrt{(4-9)^2 + (2-4)^2} = \sqrt{39}$$

$$8) a_8 = (7, 1): \quad d(A, a_8) = \sqrt{(4-7)^2 + (8-1)^2} = \sqrt{53}$$

$$d(B, a_8) = \sqrt{(4-7)^2 + (2-1)^2} = \sqrt{5} \quad \left| \Rightarrow \sqrt{5} < \sqrt{53} < \sqrt{58} \Rightarrow \text{Assign } (7, 1) \text{ to } B \right.$$

$$d(C, a_8) = \sqrt{(2-7)^2 + (3-1)^2} = \sqrt{58}$$

$$9) a_9 = (5, 6): \quad d(A, a_9) = \sqrt{(4-5)^2 + (8-6)^2} = \sqrt{5}$$

$$d(B, a_9) = \sqrt{(4-5)^2 + (2-6)^2} = \sqrt{17} \quad \left| \Rightarrow \sqrt{5} < \sqrt{17} < \sqrt{50} \Rightarrow \text{Assign } (5, 6) \text{ to } A \right.$$

$$d(C, a_9) = \sqrt{(2-5)^2 + (3-6)^2} = \sqrt{50}$$

$$10) a_{10} = (4, 6): \quad d(A, a_{10}) = \sqrt{(4-4)^2 + (8-6)^2} = 2$$

$$d(B, a_{10}) = \sqrt{(4-4)^2 + (2-6)^2} = 4 \quad \left| \Rightarrow 2 < 4 < \sqrt{53} \Rightarrow \text{Assign } (4, 6) \text{ to } A \right.$$

$$d(C, a_{10}) = \sqrt{(2-4)^2 + (3-6)^2} = \sqrt{29}$$

$$11) a_{11} = (5, 7): \quad d(A, a_{11}) = \sqrt{(4-5)^2 + (8-7)^2} = \sqrt{2}$$

$$d(B, a_{11}) = \sqrt{(4-5)^2 + (2-7)^2} = \sqrt{50} \quad \left| \Rightarrow \sqrt{2} < \sqrt{50} \Rightarrow \text{Assign } (5, 7) \text{ to } A \right.$$

$$d(C, a_{11}) = \sqrt{(2-5)^2 + (3-7)^2} = 5$$

$$12) a_{12} = (4, 7): \quad d(A, a_{12}) = \sqrt{(4-4)^2 + (8-7)^2} = 1$$

$$d(B, a_{12}) = \sqrt{(4-4)^2 + (2-7)^2} = 5 \quad \left| \Rightarrow 1 < 5 < \sqrt{53} \Rightarrow \text{Assign } (4, 7) \text{ to } A \right.$$

$$d(C, a_{12}) = \sqrt{(2-4)^2 + (3-7)^2} = \sqrt{53}$$

$$13) a_{13} = (5, 9): \quad d(A, a_{13}) = \sqrt{(4-5)^2 + (8-9)^2} = \sqrt{2}$$

$$d(B, a_{13}) = \sqrt{(4-5)^2 + (2-9)^2} = \sqrt{50} \quad \left| \Rightarrow \sqrt{2} < \sqrt{51} < \sqrt{50} \Rightarrow \text{Assign } (5, 9) \text{ to } A \right.$$

$$d(C, a_{13}) = \sqrt{(2-5)^2 + (3-9)^2} = \sqrt{51}$$

$$14) a_{14} = (4, 8): \quad d(A, a_{14}) = 0$$

$$d(B, a_{14}) = \sqrt{(4-4)^2 + (8-2)^2} = 6 \quad \left| \Rightarrow 0 < 6 < \sqrt{53} \Rightarrow \text{Assign } (4, 8) \text{ to } A \right.$$

$$d(C, a_{14}) = \sqrt{(2-4)^2 + (3-8)^2} = \sqrt{53}$$

$$15) a_{15} = (4, 9): \quad d(A, a_{15}) = \sqrt{(4-4)^2 + (8-9)^2} = 1$$

$$d(B, a_{15}) = \sqrt{(4-4)^2 + (2-9)^2} = 7 \quad \left| \Rightarrow 1 < 7 < 7 \Rightarrow \text{Assign } a_{15} \text{ to } A \right.$$

$$d(C, a_{15}) = \sqrt{(2-4)^2 + (3-9)^2} = 5$$

After the first iteration, we got the clusters:
A : { (5, 6), (4, 6), (5, 7), (4, 7), (5, 9), (4, 8), (4, 9) }
B : { (6, 2), (5, 2), (4, 2), (2, 1) }
C : { (10, 2), (9, 6), (8, 5), (8, 5) }

Now we calculate the means for each cluster to get new centroids.

$$\underline{A}: \begin{aligned} x\text{-mean} &= \frac{5+5+5+5+5+5}{7} = 5,5 \\ y\text{-mean} &= \frac{6+6+7+7+8+9+9}{7} = 7,5 \end{aligned} \quad | \Rightarrow \text{New Centroid - } A = (5,5, 7,5)$$

$$\underline{\beta}: \begin{aligned} x\text{-mean} &= \frac{6+5+4+7}{4} = 5,25 \\ y\text{-mean} &= \frac{7+7+2+4}{4} = 4,75 \end{aligned} \quad | \Rightarrow \text{New Centroid - } \beta = (5,25, 4,75)$$

$$\underline{C}: \begin{aligned} x\text{-mean} &= \frac{10+9+8+8}{4} = 8,75 \\ y\text{-mean} &= \frac{2+4+5+5}{4} = 4 \end{aligned} \quad | \Rightarrow \text{New Centroid - } C = (8,75, 4)$$

• Now we need to again calculate each distance and regroup.

$$1) a_1 = (10,2) \quad d(A, a_1) \approx \sqrt{(10,5-10)^2 + (7,5-2)^2} \approx 7,78 \\ d(B, a_1) \approx 5,76 \\ d(C, a_1) = 2,36 \quad | \Rightarrow 2,36 < 5,76 < 7,78 \xrightarrow{\text{Assign to}} C$$

$$2) a_2 = (9,3) \quad d(A, a_2) = 5,74 \\ d(B, a_2) = 5,77 \quad | \Rightarrow 5,75 < 5,77 < 5,74 \xrightarrow{\text{Assign to}} C \\ d(C, a_2) = 0,75$$

$$3) a_3 = a_4 = (8,5) \quad d(A, a_3) = 4,32 \\ d(B, a_3) = 4,95 \quad | \Rightarrow 4,25 < 4,32 < 4,95 \xrightarrow{\text{Assign}} (8,5) \text{ to } C \\ d(C, a_3) = 1,25$$

$$5) a_5 = (6,2) \quad d(A, a_5) = 5,63 \\ d(B, a_5) = 3,76 \quad | \Rightarrow 1,76 < 3,39 < 5,63 \xrightarrow{\text{Assign}} (6,2) \text{ to } B \\ d(C, a_5) = 3,39$$

$$6) a_6 = (5,2) \quad d(A, a_6) = 5,43 \\ d(B, a_6) = 0,79 \quad | \Rightarrow 0,79 < 5,26 < 5,43 \xrightarrow{\text{Assign}} (5,2) \text{ to } B \\ d(C, a_6) = 5,28$$

$$7) a_7 = (4,2) \quad d(A, a_7) = 5,41 \\ d(B, a_7) = 0,356 \quad | \Rightarrow 0,356 < 5,16 < 5,41 \xrightarrow{\text{Assign}} (4,2) \text{ to } B \\ d(C, a_7) = 5,16$$

$$8) a_8 = (7,1) \quad d(A, a_8) \approx 6.83 \\ d(B, a_8) \approx 2.37 \quad | \quad \xrightarrow{d(B, a_8) < d(A, a_8)} \text{Assign } (7,1) \text{ to } B \\ d(C, a_8) \approx 7.38$$

$$9) a_9 = (5,6) \quad d(A, a_9) \approx 1.52 \\ d(B, a_9) \approx 4.31 \quad | \quad \xrightarrow{\text{Assign } (5,6) \text{ to } A} \\ d(C, a_9) \approx 4.76$$

$$10) a_{10} = (4,6) \quad d(A, a_{10}) \approx 1.388 \\ d(B, a_{10}) \approx 4.26 \quad | \quad \xrightarrow{\text{Assign } (4,6) \text{ to } A} \\ d(C, a_{10}) \approx 5.16$$

$$11) a_{11} = (5,7) \quad d(A, a_{11}) \approx 0.721 \\ d(B, a_{11}) \approx 5.30 \quad | \quad \xrightarrow{\text{Assign } (5,7) \text{ to } A} \\ d(C, a_{11}) \approx 5.77$$

$$12) a_{12} = (4,7) \quad d(A, a_{12}) \approx 0.166 \\ d(B, a_{12}) \approx 5.76 \quad | \quad \xrightarrow{\text{Assign } (4,7) \text{ to } A} \\ d(C, a_{12}) \approx 5.63$$

$$13) a_{13} = (5,9) \quad d(A, a_{13}) = 1.71 \\ d(B, a_{13}) = 7.29 \quad | \quad \xrightarrow{\text{Assign } (5,9) \text{ to } A} \\ d(C, a_{13}) = 6.28$$

$$14) a_{14} = (4,8) \quad d(A, a_{14}) = 0.721 \\ d(B, a_{14}) = 6.25 \quad | \quad \xrightarrow{\text{Assign } (4,8) \text{ to } A} \\ d(C, a_{14}) = 6.27$$

$$15) a_{15} = (4,9) \quad d(A, a_{15}) \approx 1.648 \\ d(B, a_{15}) \approx 7.25 \quad | \quad \xrightarrow{\text{Assign } (4,9) \text{ to } A} \\ d(C, a_{15}) \approx 6.89$$

After second iteration, we got the clusters
A: $\{(5,6), (4,6), (5,7), (4,7), (5,9), (4,8), (4,9)\}$
B: $\{(5,2), (4,2), (6,2), (2,1)\}$
C: $\{(10,2), (9,3), (8,5), (7,5)\}$

Which are the same clusters we got in the previous iteration \Rightarrow Algorithm converges

Final Clusters

$$A: \{(5,6), (4,6), (5,7), (4,7), (5,8), (4,8), (4,9)\}$$

$$B: \{(6,2), (5,2), (4,2), (7,1)\}$$

$$C: \{(10,2), (3,4), (3,5), (8,5)\}.$$

To choose the appropriate number of clusters k , we can use the Elbow Method which involves running k -means for a range of k values (for example $k \in \{1\}$) and plotting the WCSS (Within-Cluster Sum of Squares) for each. The elbow point is where the WCSS curve starts to flatten.