



Data Mining and Machine Learning(SBEN454)



Assignment-6.1

Student Name: Abdulrehman Mahmoud Suliman

Student ID: 1180140

Submitted to: Eng. Peter Salah

Pre-Processing :

- Check that there is no null values.
- Drop id column
- Remove rows with negative values
- Consider any height < 90 and weight < 30 outliers and remove them
- Create new column for body mass index after calculating it from the height and weight
- Drop the height and weight columns
- Calculate the Age with years then categorize them in to three categories
 - category 1 30-50 years
 - category 2 50-60 years
 - category 3 >60 years
- Then remove the outliers from the ap_hi and ap_lo so that
 - ap_hi be from 90 to 180
 - ap_lo be from 60 to 110
 - 1 - Normal ap_hi 90-120 and ap_low 60-80
 - 2 - Elevated ap_hi 120-130 and ap_lo 60-80
 - 3 - At Risk (prehypertension) ap_hi: 130-140 and ap_lo : 80-90
 - 4- High Blood Pressure (hypertension) ap_hi >140 and ap_lo >90
- Tried to combine the ap_hi and ap_lo in one record but found that I will loose to much records that won't fall in any of the categories

- Check that all records are unique and there is no duplication.

Node Class:

This class will be used in the decision tree class as a tree is made up of nodes.

Each node contains information about the:

- attribute or feature
- Threshold upon which the record will be in the left or right tree
- Left data
- Right data
- Information gain of the attribute
- Decision if it is a leaf node

Decision Tree Class:

Functions:

- `__init__`:
 - Initialize the root and defines maximum depth
- Entropy:
 - Calculates the entropy for all classes in a given dataset
 - Returns entropy
- `calculateGain`:
 - calculates the information gain of a given attribute using the entropy function
 - returns gain

- **getPartition:**
 - This function loop over all features and unique categories in the feature column to get the best feature to be represented in this node and the best threshold for each feature for the categories upon which we can divide the data to left and right subtrees
 - Return partition list containing:
 - Feature
 - Threshold of partitioning data
 - Left dataset
 - Right dataset
 - Information Gain
- **buildTree:**
 - This function builds the tree by partitioning the data then checks that the gain of the chosen feature is a positive value and sets the nodes and for the last node or the leaf node it sets one additional thing with is the decision.
 - Returns Node
- **fit:**
 - concatenates the labels column and the data entered by the user and calls the buildTree function.
- **getDecision:**
 - recursively calls itself and get the decisions for the left and right subtree for each node until it reaches the leaf node and gets the final decision.

- **Predict:**
 - This function used to get the test inputs from the user and passes it to the getDecision function to get the predictions.

Why chose this approach?

After searching about how the sklearn library works specially if we have more than two children I found that sklearn only supports binary splitting for many reasons one of them is to support more than one splitting criteria like the gini impurity which only supports the binary splits and it wouldn't make any problem as a series of binary splits can model any number of simultaneous splits.

Comparing the results with sklearn:

By predicting the test samples using the built-in function and the from scratch class I found that the results are nearly the same.

From scratch:

Accuracy = 0.7176456753101664

Built-in :

Accuracy = 0.7171752807667431

$$P = 8 \quad N = 6 \quad \text{Total} = 14$$

$$\begin{aligned} \text{Entropy}(S) &= \frac{P}{P+N} \log\left(\frac{P}{P+N}\right) + \frac{n}{P+N} \log\left(\frac{n}{P+N}\right) \\ &= \frac{8}{14} \log\left(\frac{8}{14}\right) + \frac{6}{14} \log\left(\frac{6}{14}\right) = \underline{\underline{0.989}} \end{aligned}$$

* Entropy for each attribute:-

→ Early registration

	P	n	Entropy
1	4	2	0.918
0	4	4	1

* Calculate Average Information Entropy:-

$$\begin{aligned} I(\text{early registration}) &= \frac{P_1 + n_1}{P+n} \text{Entropy}(1) + \frac{P_2 + n_2}{P+n} \text{Entropy}(0) \\ &= \frac{4+2}{14} \times 0.918 + \frac{8}{14} \times 1 = \underline{\underline{0.965}} \end{aligned}$$

* Calculate Gain

$$\begin{aligned} \text{Gain} &= \text{Entropy}(S) - I(\text{Attribute}) \\ &= 0.989 - 0.965 = \underline{\underline{0.024}} \end{aligned}$$

→ Finished homework II :-

	P	n	Entropy
1	5	2	0.863
0	3	4	0.985

$$I(\text{Finished homework II}) = \frac{5+2}{14} \times 0.863 + \frac{3+4}{14} \times 0.985 = 0.924$$

$$\text{Gain} = 0.989 - 0.924 = \underline{\underline{0.065}}$$

→ Senior :-

	P	n	Entropy
1	5	3	0.954
0	3	3	1

$$I(\text{Senior}) = \frac{5+3}{14} \times 0.954 + \frac{6}{14} \times 1 = 0.974$$

$$\text{Gain} = 0.989 - 0.974 = \boxed{0.015}$$

→ Likes Coffee :-

	P	n	Entropy
1	3	1	0.8113
0	5	5	1

$$I(\text{Likes Coffee}) = \frac{3+1}{14} \times 0.8113 + \frac{10}{14} \times 1 = 0.946$$

$$\text{Gain} = 0.989 - 0.946 = \boxed{0.043}$$

→ Liked the last homework

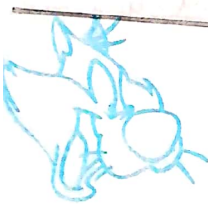
	P	n	Entropy
1	5	4	0.991
0	3	2	0.971

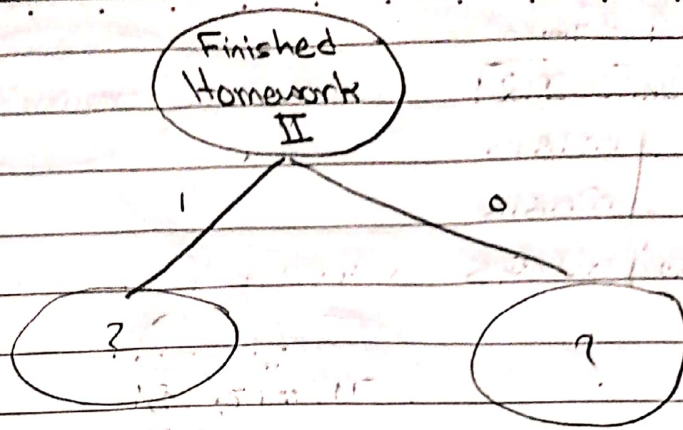
$$I(\text{Liked the last homework}) = \frac{5+4}{14} \times 0.991 + \frac{5}{14} \times 0.971 = 0.984$$

$$\text{Gain} = 0.989 - 0.984 = \boxed{0.005}$$

Attributes	Gain
Early registration	0.024
Finished homework II	0.065
Senior	0.015
Likes coffee	0.043
Liked the last homework	0.005

∴ Finished homework II
will be the root node
as it has the highest
gain.





for (1):

$$P = 5 \quad N = 2 \quad \text{Total} = 7$$

$$\text{Entropy}(S) = -\frac{5}{7} \log\left(\frac{5}{7}\right) - \frac{2}{7} \log\left(\frac{2}{7}\right) = 0.863$$

→ Early registration =

	P	n	Entropy
1	3	0	0
0	2	2	1

$$I(\text{Early registration}) = \frac{3}{7} \times 0 + \frac{2+2}{7} \times 1 = 0.571$$

$$\text{Gain} = 0.863 - 0.571 = \underline{\underline{0.292}}$$

→ Senior =

	P	n	Entropy
1	3	2	0.9709
0	2	0	0

$$I(\text{Senior}) = \frac{5}{7} \times 0.9709 + 0 = 0.6935$$

$$\text{Gain} = 0.863 - 0.6935 = \underline{\underline{0.1695}}$$

→ Likes Coffee :-

	P	n	Entropy
1	1	1	1
0	4	1	0.7219

$$I(\text{Likes Coffee}) = \frac{2}{7} \times 1 + \frac{5}{7} \times 0.7219 = 0.8014$$

$$\text{Gain} = 0.863 - 0.8014 = \underline{\underline{0.0616}}$$

→ Likes the last homeworks

	P	n	Entropy
1	3	2	0.9709
0	2	0	0

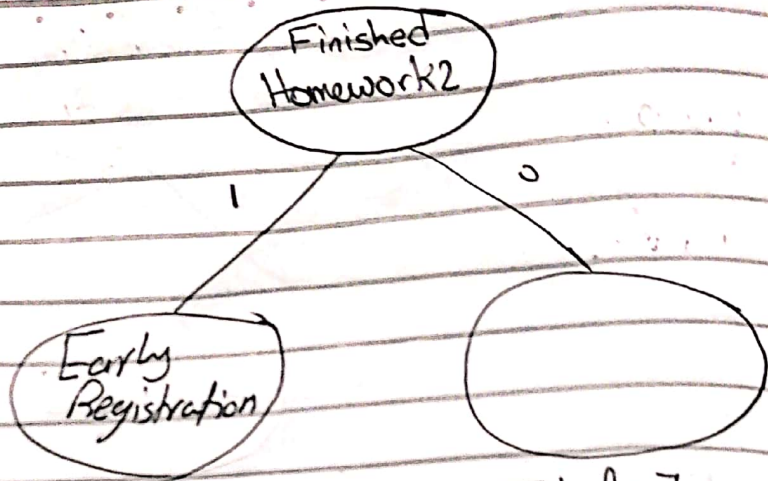
$$I(\text{LLH}) = \frac{5}{7} \times 0.9709 + 0 = 0.6935$$

$$\text{Gain} = 0.863 - 0.6935$$

$$= \underline{\underline{0.1695}}$$

ROX

Attribute	Gain
Early registration	0.202
Senior	0.1695
Likes Coffee	0.0616
Likes last homework	0.1695



for Finished Homeworks II = 0 : $P=3$ $n=4$ Total = 7

→ Early registration:

Entropy(S) = 0.985

	P	n	Entropy
1	1	2	0.918
0	2	2	1

$$I(ER) = \frac{3}{7} \times 0.918 + \frac{4}{7} \times 1 = 0.9668$$

$$\text{Gain} = 0.985 - 0.9668 = \boxed{0.0202}$$

→ Senior

	P	n	Entropy
1	2	1	0.918
0	1	3	0.8113

$$I(\text{Senior}) = \frac{3}{7} \times 0.918 + \frac{4}{7} \times 0.8113 = 0.857$$

$$\text{Gain} = 0.985 - 0.857 = \boxed{0.128}$$

→ Likes Coffee

	P	n	Entropy
1	2	0	0
0	1	4	0.7219

$$I(LC) = \frac{2}{7} \times 0 + \frac{5}{7} \times 0.7219 = 0.515$$

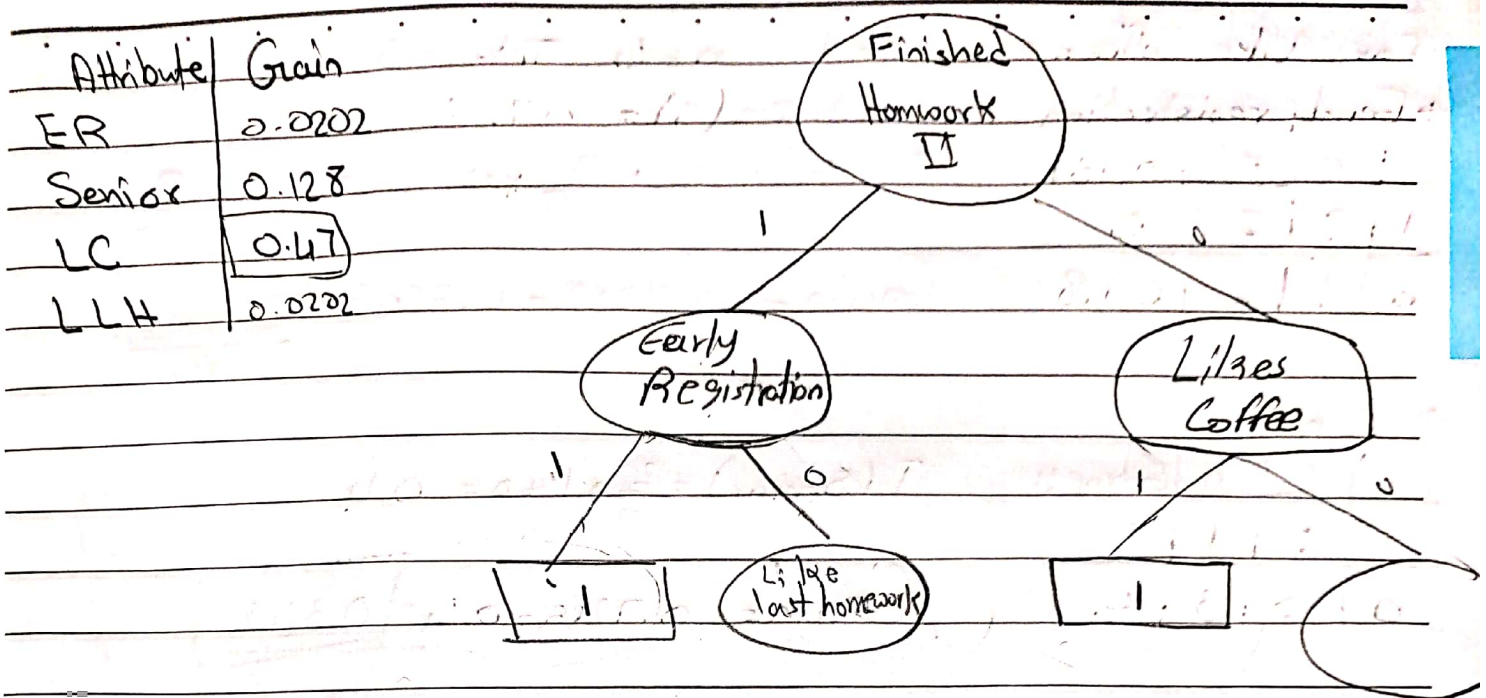
$$\text{Gain} = 0.985 - 0.515 = \boxed{0.47}$$

→ Likes Last Homework

	P	n	Entropy
1	2	2	1
0	1	2	0.918

$$I(LLH) = \frac{4}{7} \times 1 + \frac{3}{7} \times 0.918 = 0.9668$$

$$\text{Gain} = 0.985 - 0.9668 = \boxed{0.0202}$$



for ER = 0 : $P = 2$ $n = 2$ Total = 4 Entropy(S) = 1

→ Senior

P	n	Entropy
1	2	0.918
0	0	0

$$I(\text{Senior}) = \frac{3}{4} \times 0.918 + 0 = 0.6885$$

$$\text{Gain} = 1 - 0.6885 = 0.3115$$

→ Likes Coffee

P	n	Entropy
1	1	1
0	1	1

$$I(\text{LC}) = \frac{2}{4} \times 1 + \frac{2}{4} \times 1 = 1$$

$$\text{Gain} = 1 - 1 = 0$$

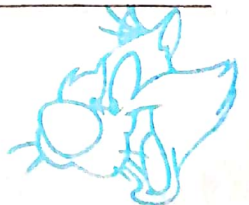
→ Likes last homework

P	n	Entropy
1	2	0.918
0	0	0

$$I(\text{LLH}) = 0.6885$$

$$\text{Gain} = 1 - 0.6885 = 0.3115$$

Attribute	Gain
Senior	0.3115
LC	0
LLH	0.3115



for Like Coffee = 0 $P=1$ $n=4$ Total = 5

→ Early registration

	P	n	Entropy
1	0	2	0
0	1	2	0.918

$$Entropy(S) = 0.7219$$

$$I(ER) = 0 + \frac{3}{5} \times 0.918 = 0.5508$$

$$Gain = 0.7219 - 0.5508 = \boxed{0.1711}$$

→ Senior

	P	n	Entropy
1	1	1	1
0	0	3	0

$$I(Senior) = \frac{2}{5} \times 1 + 0 = 0.4$$

$$Gain = 0.7219 - 0.4 = \boxed{0.3219}$$

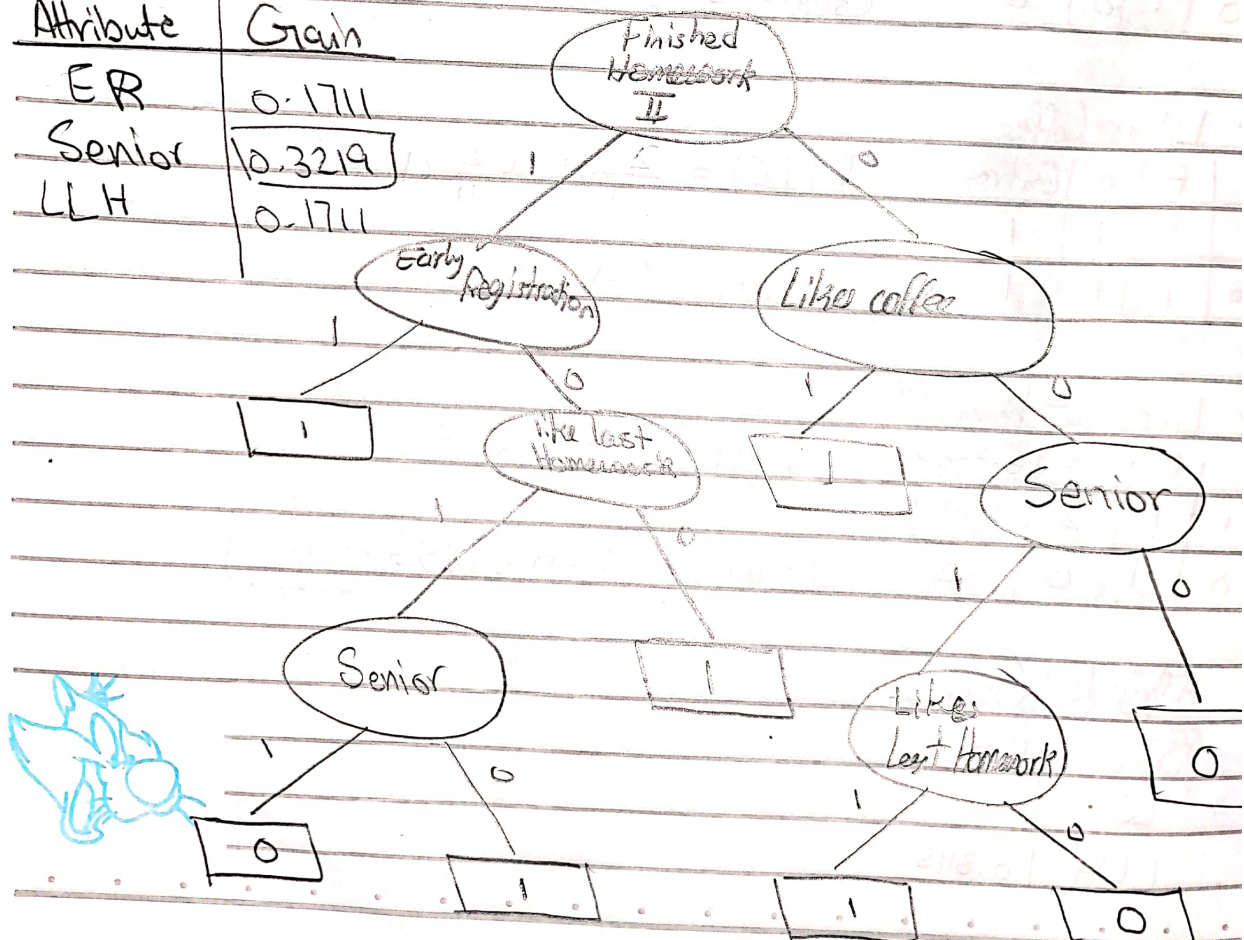
→ Likes Last homework

	P	n	Entropy
1	1	2	0.918
0	0	2	0

$$I(ER) = \frac{3}{5} \times 0.918 = 0.5508$$

$$Gain = \boxed{0.1711}$$

Attribute	Gain
ER	0.1711
Senior	<u>0.3219</u>
LLH	0.1711



for like last homework = 1 $P=1$ $n=2$ Total = 3 Entropy(s) = 0.918

Senior

	P	n	Entropy
1	0	2	0
0	1	0	0

$$I(\text{Senior}) = 0$$

$$\text{Gain} = 0.918 - 0 = 0.918$$

Like coffee

	P	n	Entropy
1	1	1	1
0	0	1	0

$$I(LC) = \frac{2}{3} = 0.67$$

$$\text{Gain} = 0.918 - 0.67 = \underline{\underline{0.248}}$$

for Senior = 1 $P=1$ $n=1$ Total = 2 Entropy(s) = 1

Early registration

	P	n	Entropy
1	0	0	0
0	1	1	1

$$I(ER) = \frac{2}{2} \times 1 = 1$$

$$\text{Gain} = 1 - 1 = 0$$

Like last Homework

	P	n	Entropy
1	1	0	0
0	0	1	0

$$I(LLH) = 0$$

$$\text{Gain} = 1 - 0 = \underline{\underline{1}}$$

