

# ANSWERS

## Exercise 1

- i We start by calculating the entropy of the whole unsplit set and the information gain values that would be achieved for each of the variables if the set were split on that variable.

$$\begin{aligned} \text{entropy}(\text{whole\_set}) &= -P(\text{yes})\log_2 P(\text{yes}) - P(\text{no})\log_2 P(\text{no}) = \\ &= -\frac{7}{12}\log_2 \frac{7}{12} - \frac{5}{12}\log_2 \frac{5}{12} = 0.98 \end{aligned}$$

Alternatively, entropy values for two-class outcomes can be looked up in the table headed "Entropy for subsets with various relative percentage combinations", found in the appendix of the lab sheet. This table gives 'ready-calculated' entropies for different proportions of the two classes in a set. The calculation above corresponds to the following line in the table:

Number of items	Breakdown		Entropy
12	5	7	0.98

From here on we will be looking up entropy values in the table rather than calculating them.

There are three variables on which the top level set could be split: **head shape**, **body shape** and **body colour**. For each of these we calculate the entropies of the variable values and then the information gain that would result if the set were split on that variable. We look up the table to find the entropy for the yes/no split ratio.

- **head shape**

$$\text{entropy}(\text{head shape} : \text{square}) = 0.99 \quad (9/12 \text{ instances, yes/no split: } 5,4)$$

$$\text{entropy}(\text{head shape} : \text{round}) = 0.92 \quad (3/12 \text{ instances, yes/no split: } 2,1)$$

$$\text{IG}(\text{head shape}) = \text{entropy}(\text{whole set}) - 9/12 \times 0.99 - 3/12 \times 0.92 = 0.01$$

- **body shape**

$$\text{entropy}(\text{body shape} : \text{rectangular}) = 0.65 \quad (6/12 \text{ instances, yes/no split: } 5,1)$$

$$\text{entropy}(\text{body shape} : \text{oval}) = 0.92 \quad (6/12 \text{ instances, yes/no split: } 2,4)$$

$$\text{IG}(\text{Temp}) = \text{entropy}(\text{whole set}) - 6/12 \times 0.65 - 6/12 \times 0.92 = 0.2$$

- **body colour**

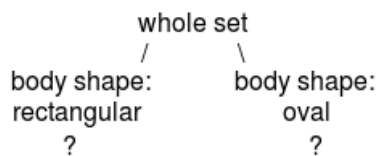
$$\text{entropy}(\text{body colour} : \text{white}) = 0.97 \quad (10/12 \text{ instances, yes/no split: } 6,4)$$

$$\text{entropy}(\text{body colour} : \text{grey}) = 1 \quad (2/12 \text{ instances, yes/no split: } 1,1)$$

$$\text{IG}(\text{Temp}) = \text{entropy}(\text{whole set}) - 10/12 \times 0.97 - 2/12 \times 1 = 0.005$$

The variable that should be used first for splitting the data is **body shape**, as it has the highest information gain, 0.2, with respect to the target variable **Likely to default on mortgage**.

After the first splitting, the tree looks like this:



ii Now we'll continue to split the set until we have a complete tree.

First let's look at the **rectangular** branch, which needs to be split further.

- **head shape**

$$\text{entropy}(\text{head shape} : \text{square}|\text{rectangular}) = 0 \quad (5/6 \text{ instances, yes/no split: } 6,0)$$

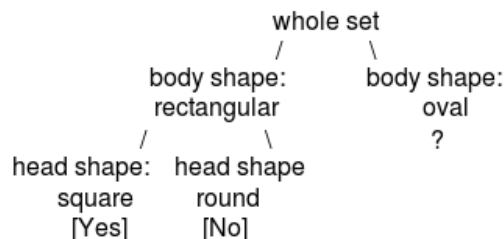
$$\text{entropy}(\text{head shape} : \text{round}|\text{rectangular}) = 0 \quad (1/6 \text{ instances, yes/no split: } 0,1)$$

- **body colour**

$$\text{entropy}(\text{body colour} : \text{white}|\text{rectangular}) = 0 \quad (5/6 \text{ instances, yes/no split: } 6,0)$$

$$\text{entropy}(\text{body colour} : \text{grey}|\text{rectangular}) = 0 \quad (1/6 \text{ instances, yes/no split } 0,1)$$

Note that when the split is between 0 and another number, the entropy is 0. With both attributes above the information gain is maximal, the entropy being 0, so we pick one of the two attributes at random, let's say **head shape**, to split at the **body shape:rectangular** node. Since its entropy is 0, this attribute defines leaf nodes, with values for the target variable.



Now let's look at the **oval** branch, which needs to be split further.

- **head shape**

$$\text{entropy}(\text{head shape} : \text{square}|\text{oval}) = 0 \quad (4/6 \text{ instances, yes/no split: } 0,4)$$

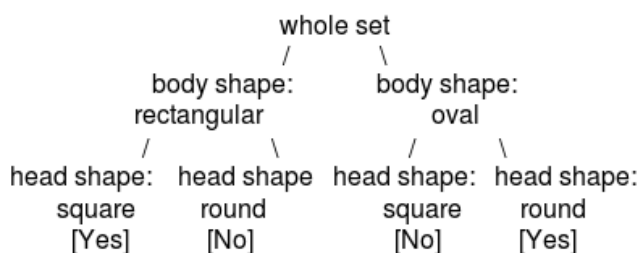
$$\text{entropy}(\text{head shape} : \text{round}|\text{oval}) = 0 \quad (2/6 \text{ instances, yes/no split: } 2,0)$$

- **body colour**

$$\text{entropy}(\text{body colour} : \text{white}|\text{oval}) = 0.72 \quad (5/6 \text{ instances, yes/no split: } 1,4)$$

$$\text{entropy}(\text{body colour} : \text{grey}|\text{oval}) = 0 \quad (1/6 \text{ instances, yes/no split: } 1,0)$$

The attribute **head shape** has an entropy of 0 for **body shape:oval**. We pick it as the next splitting attribute. The decision tree now looks like this:



iii To get the prediction for a customer with **head shape:square**, **body shape:rectangular**, **body colour:gray** we navigate the derived decision tree along the following nodes:

- **whole set**

- `body shape:rectangular`
- `head shape:square`

The last listed node is a leaf node associated with the **Yes** class i.e. the tree predicts that the customer is likely to default on the mortgage.

## Exercise 2

### a) Finding the best splitter attribute for the whole set

For the whole set we read the entropy out of the table (split 5, 9).

$$\text{entropy}(\text{whole\_set}) = 0.94 \quad (\text{split } 5, 9)$$

Now we look at the different attributes to see which one will provide the best Information Gain. We read the entropies from the table, based on the "No"/"Yes" split of the "Play" values in the subset.

#### (a) Outlook

$$\text{entropy}(\text{Outlook} : \text{Sunny}) = 0.97 \quad (\text{subset size: } 5/14, \text{ split: } 2, 3)$$

$$\text{entropy}(\text{Outlook} : \text{Overcast}) = 0 \quad (\text{subset size: } 4/14, \text{ split: } 4, 0)$$

$$\text{entropy}(\text{Outlook} : \text{Rainy}) = 0.97 \quad (\text{subset size: } 5/14, \text{ split: } 3, 2)$$

$$\text{IG}(\text{Outlook}) = \text{entropy}(\text{whole\_set}) - 5/14 * 0.97 - 4/14 * 0 - 5/14 * 0.97 = 0.25$$

#### (b) Temp

$$\text{entropy}(\text{Temp} : \text{Hot}) = 1 \quad (\text{subset size: } 4/14, \text{ split: } 2, 2)$$

$$\text{entropy}(\text{Temp} : \text{Mild}) = 0.92 \quad (\text{subset size: } 6/14, \text{ split: } 4, 2)$$

$$\text{entropy}(\text{Temp} : \text{Cool}) = 0.81 \quad (\text{subset size: } 4/14, \text{ split: } 3, 1)$$

$$\text{IG}(\text{Temp}) = \text{entropy}(\text{whole\_set}) - 4/14 * 1 - 6/14 * 0.92 - 4/14 * 0.81 = 0.03$$

#### (c) Humidity

$$\text{entropy}(\text{Humidity} : \text{High}) = 0.99 \quad (\text{subset size: } 7/14, \text{ split: } 3, 4)$$

$$\text{entropy}(\text{Humidity} : \text{Normal}) = 0.59 \quad (\text{subset size: } 7/14, \text{ split: } 6, 1)$$

$$\text{IG}(\text{Humidity}) = \text{entropy}(\text{whole\_set}) - 7/14 * 0.99 - 7/14 * 0.59 = 0.15$$

#### (d) Windy

$$\text{entropy}(\text{Windy} : \text{True}) = 1 \quad (\text{subset size: } 6/14, \text{ split: } 3, 3)$$

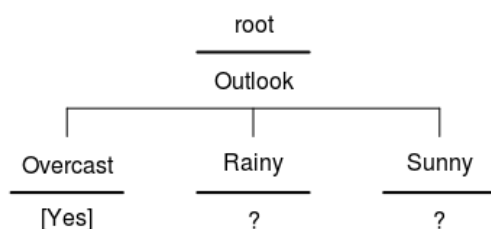
$$\text{entropy}(\text{Windy} : \text{False}) = 0.81 \quad (\text{subset size: } 8/14, \text{ split: } 6, 2)$$

$$\text{IG}(\text{Windy}) = \text{entropy}(\text{whole\_set}) - 6/14 * 1 - 8/14 * 0.81 = 0.05$$

The variable that should be used first for splitting the data is 'Outlook', as it has the highest information gain, 0.25, with respect to the target variable 'Play'.

### b) Building the rest of the tree

After the first split, the tree looks like this:



The Overcast branch has an entropy of 0 and does not need to be split further.

Now we look at the Rainy branch, which needs to be split further.

$$\text{entropy}(\text{Temp} : \text{Mild} | \text{Rainy}) = 0.92 \quad (\text{subset size: } 3/5, \text{ split: } 2, 1)$$

$$\text{entropy}(\text{Temp} : \text{Cool} | \text{Rainy}) = 1 \quad (\text{subset size: } 2/5, \text{ split: } 1, 1)$$

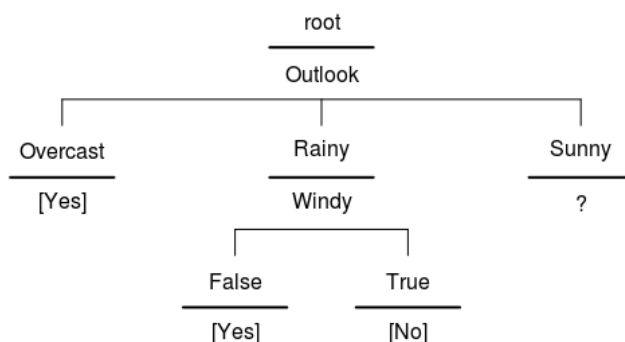
$$\text{entropy}(\text{Humidity} : \text{High} | \text{Rainy}) = 1 \quad (\text{subset size: } 2/5, \text{ split: } 1, 1)$$

$$\text{entropy}(\text{Humidity} : \text{Normal} | \text{Rainy}) = 0.92 \quad (\text{subset size: } 3/5, \text{ split: } 2, 1)$$

$$\text{entropy}(\text{Windy} : \text{True} | \text{Rainy}) = 0 \quad (\text{subset size: } 2/5, \text{ split: } 0, 2)$$

$$\text{entropy}(\text{Windy} : \text{False} | \text{Rainy}) = 0 \quad (\text{subset size: } 3/5, \text{ split: } 3, 0)$$

The IG doesn't need to be calculated as there is only one attribute with entropy of 0 and that is Windy. This means that this attribute provides the greatest information gain and we choose it to split the node 'Rainy'. The decision tree now looks like this:



The False and True branches both have entropy 0 and do not need to be split further.

Finally, we look at the Sunny branch, which does need to be split further.

$$\text{entropy}(\text{Temp} : \text{Hot} | \text{Sunny}) = 0 \quad (\text{subset size: } 2/5, \text{ split } 2, 0)$$

$$\text{entropy}(\text{Temp} : \text{Mild} | \text{Sunny}) = 1 \quad (\text{subset size: } 2/5, \text{ split } 1, 1)$$

$$\text{entropy}(\text{Temp} : \text{Cool} | \text{Sunny}) = 0 \quad (\text{subset size: } 1/5, \text{ split } 1, 0)$$

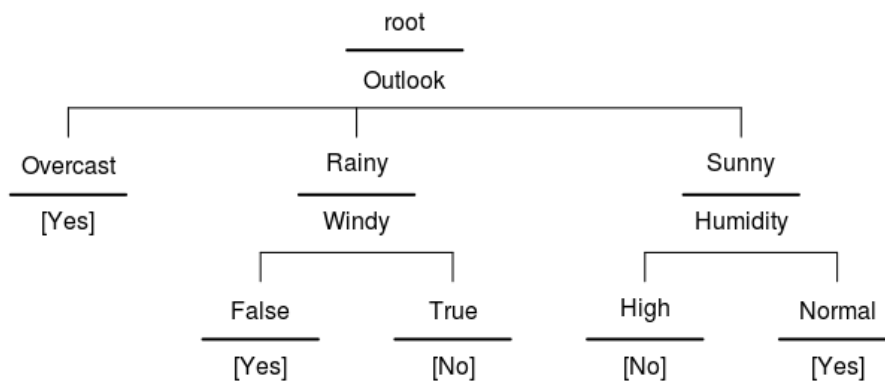
$$\text{entropy}(\text{Humidity} : \text{High} | \text{Sunny}) = 0 \quad (\text{subset size: } 3/5, \text{ split } 3, 0)$$

$$\text{entropy}(\text{Humidity} : \text{Normal} | \text{Sunny}) = 0 \quad (\text{subset size: } 2/5, \text{ split } 2, 0)$$

$$\text{entropy}(\text{Windy} : \text{True} | \text{Sunny}) = 1 \quad (\text{subset size: } 2/5, \text{ split } 1, 1)$$

$$\text{entropy}(\text{Windy} : \text{False} | \text{Sunny}) = 0.92 \quad (\text{subset size: } 3/5, \text{ split } 2, 1)$$

When it's Sunny, the entropy for Humidity is 0 and non zero for the other attributes, so Humidity has the highest IG and should be chosen to split the node 'Sunny'. The decision tree finally looks like this:



c) **Predicting the outcome for a particular instance**

To predict the outcome for a particular instance, we traverse the tree using the attribute values of the instance. The given instance is:

Outlook: Sunny  
Temp: Hot  
Humidity: Normal  
Windy: True

The first split in the tree is by **Outlook** and as the instance has the value **Sunny** for this attribute, at the first branching we take the rightmost branch. The next split is by **Humidity** and the given instance has the value **Normal**, which means that the traversal is along the rightmost branch again, arriving at a leaf node with a probability-1 outcome of "Yes".

**For the given instance, our model predicts that the game will go ahead i.e. Play="Yes".**