

Name: Answer Key

After the first exam in a data mining course, the results of the exam were recorded along with some information about each student. The data is below:

ID	Passed All Assignments	GPA	Language	Passed Exam
1	No	3.1	Python	Yes
2	No	2.0	Python	No
3	Yes	3.5	C++	Yes
4	Yes	2.5	Java	No
5	Yes	3.9	Python	No
6	No	3.3	C++	Yes
7	Yes	3.2	Java	Yes

We want to use the above data to create a decision tree that can predict which students will pass the exam.

What is the class label? passed exam

What are the attributes? passed all assignments, GPA, language

In order to create a decision tree, we need to decide which attribute to split on first. To do this, we must calculate the **gain** of splitting on each of our attributes.

1. Start by calculating the impurity of the parent. Use entropy as the measure of impurity.

$$\begin{aligned}\text{Entropy}(\text{parent}) &= - \left(\frac{4}{7}\right) \log_2 \left(\frac{4}{7}\right) - \left(\frac{3}{7}\right) \log_2 \left(\frac{3}{7}\right) \\ &= 0.985\end{aligned}$$

Name: _____

2. Next, calculate the information gain of splitting on 'Passed All Assignments'. The gain of splitting on 'Passed All Assignments' will be the entropy of the parent minus the entropy of making this split. (We want to know how much the impurity decreases by making this split.)

$$\begin{array}{cc} y:2 & y:2 \\ N:2 & N:1 \end{array} \quad \text{Entropy}(Y) = 1 \quad \text{Entropy}(N) = -\frac{2}{3} \log\left(\frac{2}{3}\right) - \frac{1}{3} \log\left(\frac{1}{3}\right) = 0.918$$

$$\text{Entropy}(\text{split}) = \left(\frac{4}{7}\right)(1) + \left(\frac{3}{7}\right)(0.918) = 0.965$$

$$\text{Gain}_{\text{split}} = 0.985 - 0.965 = 0.02$$

3. Next, calculate the information gain of splitting on 'Language'. The gain of splitting on 'Language' will be the entropy of the parent minus the entropy of making this split. (We want to know how much the impurity decreases by making this split.)

$$\begin{array}{cc} y:1 & N:2 \\ y:2 & N:0 \end{array} \quad \text{Entropy}(\text{Python}) = -\frac{1}{3} \log\left(\frac{1}{3}\right) - \frac{2}{3} \log\left(\frac{2}{3}\right) = 0.918$$

$$\begin{array}{cc} y:2 & N:0 \\ y:1 & N:1 \end{array} \quad \text{Entropy}(\text{C++}) = 0$$

$$\text{Entropy}(\text{Java}) = 1$$

$$\text{Entropy}_{\text{split}} = \left(\frac{3}{7}\right)(0.918) + \left(\frac{2}{7}\right)(0) + \left(\frac{2}{7}\right)(1) = 0.679$$

$$\text{Gain}_{\text{split}} = 0.985 - 0.679 = 0.306$$

Name: _____

4. Next, calculate the information gain of splitting on 'GPA'. Because GPA is a continuous attribute, we need to try different candidate threshold values. Determine all of the candidate thresholds, then calculate the information gain for each of them.

No	No	Yes	Yes	Yes	Yes	No
2.0	2.5	3.1	3.2	3.3	3.5	3.9

↑

↑

① splitting at 2.8

$\leq 2.8: 2$

$> 2.8: 5$

$$\text{Entropy}(\leq 2.8) = 0$$

$$\text{Entropy}(> 2.8) = -\frac{4}{5} \log\left(\frac{4}{5}\right) - \frac{1}{5} \log\left(\frac{1}{5}\right) = 0.722$$

$$\text{Entropy}_{\text{split}} = \frac{2}{7}(0) + \frac{5}{7}(0.722) = 0.516$$

$$\text{Gain}_{\text{split}} = 0.985 - 0.516 = 0.469$$

② splitting at 3.7

$\leq 3.7: 6$

$> 3.7: 1$

$$\text{Entropy}(\leq 3.7) = -\frac{2}{6} \log\left(\frac{2}{6}\right) - \frac{4}{6} \log\left(\frac{4}{6}\right) = 0.918$$

$$\text{Entropy}(> 3.7) = 0$$

$$\text{Entropy}_{\text{split}} = \frac{6}{7}(0.918) + \frac{1}{7}(0) = 0.787$$

$$\text{Gain}_{\text{split}} = 0.985 - 0.787 = 0.198$$

↓
better split
threshold

5. Impurity metrics (like Gini and entropy) favor attributes with more values. Because 'Language' can be split 3 ways, but 'Passed All Assignments' and 'GPA' are only split 2 ways, we should use **gain ratio** to compare the attributes, rather than just gain.

Calculate the **split info** for each of the three attributes. With that, calculate the **gain ratio** for each of the three attributes.

$$\text{Splitinfo}(\text{passed all assignments}) = -\left(\frac{4}{7}\right) \log\left(\frac{4}{7}\right) - \left(\frac{3}{7}\right) \log\left(\frac{3}{7}\right) = 0.985$$

$$\text{Splitinfo}(\text{Language}) = -\left(\frac{3}{7}\right) \log\left(\frac{3}{7}\right) - \left(\frac{2}{7}\right) \log\left(\frac{2}{7}\right) - \left(\frac{2}{7}\right) \log\left(\frac{2}{7}\right) = 1.557$$

$$\text{Splitinfo}(\text{GPA})_{2.8} = -\frac{2}{7} \log\left(\frac{2}{7}\right) - \frac{5}{7} \log\left(\frac{5}{7}\right) = 0.863$$

$$\text{Splitinfo}(\text{GPA})_{3.7} = -\frac{6}{7} \log\left(\frac{6}{7}\right) - \frac{1}{7} \log\left(\frac{1}{7}\right) = 0.58$$

$$\text{Gain ratio}(\text{passed all assignments}) = \frac{0.02}{0.985} = 0.02$$

$$\text{Gain ratio}(\text{Language}) = \frac{0.306}{1.557} = 0.197$$

$$\text{Gain ratio}(\text{GPA})_{2.8} = \frac{0.469}{0.863} = 0.543$$

$$\text{Gain ratio}(\text{GPA})_{3.7} = \frac{0.198}{0.58} = 0.33$$

6. Based on gain ratio, which attribute do you choose to be the first node in the tree?

Based on gain ratio, I would choose GPA with a threshold of 2.8

as it has the highest gain ratio of the three.