

Name: Yola Charara

HW3

1.

Age: 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70

a. Smoothing by bin means using a bin depth of 3:

Equal depth partitioning:

Age: (13, 15, 16), (16, 19, 20), (20, 21, 22), (22, 25, 25), (25, 25, 30), (33, 33, 35), (35, 35, 35), (36, 40, 45), (46, 52, 70).

$$\text{Bin 1: } \frac{13 + 15 + 16}{3} = 14.67$$

$$\text{Bin 2: } \frac{16 + 19 + 20}{3} = 18.33$$

$$\text{Bin 3: } \frac{20 + 21 + 22}{3} = 21$$

$$\text{Bin 4: } \frac{22 + 25 + 25}{3} = 24$$

$$\text{Bin 5: } \frac{25 + 25 + 30}{3} = 26.67$$

$$\text{Bin 6: } \frac{33 + 33 + 35}{3} = 33.67$$

$$\text{Bin 7: } \frac{35 + 35 + 35}{3} = 35$$

$$\text{Bin 8: } \frac{36 + 40 + 45}{3} = 40.33$$

$$\text{Bin 9: } \frac{46 + 52 + 70}{3} = 56$$

Data result after smoothed:

Age: 14.67, 14.67, 14.67, 18.33, 18.33, 18.33, 21, 21, 21, 24, 24, 24, 26.67, 26.67, 26.67, 33.67, 33.67, 33.67, 35, 35, 35, 40.33, 40.33, 40.33, 56, 56, 56.

Smoothing by bin means decreases local variation. It suppresses small fluctuations but at the same time shows the overall trend. Values that are very far from the average are brought nearer to the mean (such as 70 in this case).

b. $\text{IQR} = Q3 - Q1 = 35 - 20 = 15$

$$Q1 - 1.5 * \text{IQR} = 20 - 1.5 * 15 = -2.5$$

$$Q3 + 1.5 * \text{IQR} = 35 + 1.5 * 15 = 57.5$$

Data points outside this range $[-2.5; 57.5]$ are outliers, therefore 70 is an outlier.

$$\text{c. } x' = \frac{x - \min}{\max - \min} = \frac{35 - 13}{70 - 13} = 0.386$$

The normalized value is 0.386

$$\begin{aligned} \text{d. Mean} = \mu &= \frac{\sum age}{n} = \\ &= \frac{13 + 15 + 16 + 16 + 19 + 20 + 20 + 21 + 22 + 22 + 25 + 25 + 25 + 25 + 30 + 33 + 33 + 35 + 35 + 35 + 35 + 36 + 40 + 45 + 46 + 52 + 70}{27} \\ &= \frac{809}{27} = 29.96 \end{aligned}$$

$$\text{Standard deviation} = \sigma = \sqrt{\frac{\sum (x - \mu)^2}{n}} =$$

$$\sqrt{\frac{(13-29.96)^2 + (15-29.96)^2 + (16-29.96)^2 + (16-29.96)^2 + (19-29.96)^2 + (20-29.96)^2 + (20-29.96)^2 + (21-29.96)^2 + (22-29.96)^2 + (22-29.96)^2 + (25-29.96)^2 + (25-29.96)^2 + (25-29.96)^2 + (25-29.96)^2 + (30-29.96)^2 + (33-29.96)^2 + (33-29.96)^2 + (35-29.96)^2 + (35-29.96)^2 + (35-29.96)^2 + (35-29.96)^2 + (36-29.96)^2 + (40-29.96)^2 + (45-29.96)^2 + (46-29.96)^2 + (52-29.96)^2 + (70-29.96)^2}{27}}$$

$$\sigma = 12.70$$

$$\text{Z-score} = z = \frac{x - \mu}{\sigma} = \frac{35 - 29.96}{12.70} = 0.397$$

The z-score for age 35 is 0.397.

$$\text{e. } x = 35; \max \text{ age} = 70$$

$$j \Rightarrow \frac{\max(|age|)}{10^j} < 1$$

$$\frac{70}{10^2} = 0.7 < 1$$

$$\text{So, } j = 2$$

$$x' = \frac{35}{10^2} = 0.35$$

The transformed value of the age 35 is 0.35.

3.

Base info:

Junior count: 113; Senior count: 52; Total count: 165

$$\text{Info} = - \sum p_i \log_2(p_i) = - \frac{113}{165} \log_2\left(\frac{113}{165}\right) - \frac{52}{165} \log_2\left(\frac{52}{165}\right) = 0.90$$

Department info gain:

$$\text{Info sales} = -\frac{80}{110} \log_2 \left(\frac{80}{110} \right) - \frac{30}{110} \log_2 \left(\frac{30}{110} \right) = 0.845$$

$$\text{Info systems} = -\frac{23}{31} \log_2 \left(\frac{23}{31} \right) - \frac{8}{31} \log_2 \left(\frac{8}{31} \right) = 0.824$$

$$\text{Info marketing} = -\frac{4}{14} \log_2 \left(\frac{4}{14} \right) - \frac{10}{14} \log_2 \left(\frac{10}{14} \right) = 0.863$$

$$\text{Info secretary} = -\frac{6}{10} \log_2 \left(\frac{6}{10} \right) - \frac{4}{10} \log_2 \left(\frac{4}{10} \right) = 0.971$$

$$\text{Info department} = \frac{110}{165} (0.845) + \frac{31}{165} (0.824) + \frac{14}{165} (0.863) + \frac{10}{165} (0.971) = 0.850$$

$$\text{Gain department} = 0.90 - 0.850 = 0.05$$

Age info gain:

Only (31_35) has both junior and senior counts. So, all other age datapoints' info are equal to zero.

$$\text{Info (31_35)} = -\frac{44}{79} \log_2 \left(\frac{44}{79} \right) - \frac{35}{79} \log_2 \left(\frac{35}{79} \right) = 0.991$$

$$\text{Info age} = \frac{79}{165} (0.991) = 0.474$$

$$\text{Gain age} = 0.90 - 0.474 = 0.426$$

Salary info gain:

Only (46K_50K) has both junior and senior counts. So, all other salary datapoints' info are equal to zero.

$$\text{Info (46K_50K)} = -\frac{23}{63} \log_2 \left(\frac{23}{63} \right) - \frac{40}{63} \log_2 \left(\frac{40}{63} \right) = 0.947$$

$$\text{Info salary} = \frac{63}{165} (0.947) = 0.362$$

$$\text{Gain salary} = 0.90 - 0.362 = 0.538$$

Root of the decision tree:

The salary gain is the highest gain among the attributes. So, the root node will be the salary.

Subsets of the decision tree:

For this part, the calculations were done in R markdown.

Results:

- Subset for Salary = 46K-50K:
 - Junior = 23, Senior = 40
 - Info = 0.947
 - Info gain department = 0.947
 - Info gain age = 0.947

Both attributes are equal and can work as the next node. If we were to choose, we can pick Age as the first subset after the root for a more organized structure.