## Problem 1

$$\text{Entropy (commute)} = -\frac{3}{13} \times \log_2 \frac{3}{13} - \frac{7}{13} \times \log_2 \frac{7}{13} - \frac{2}{13} \times \log_2 \frac{2}{13}$$
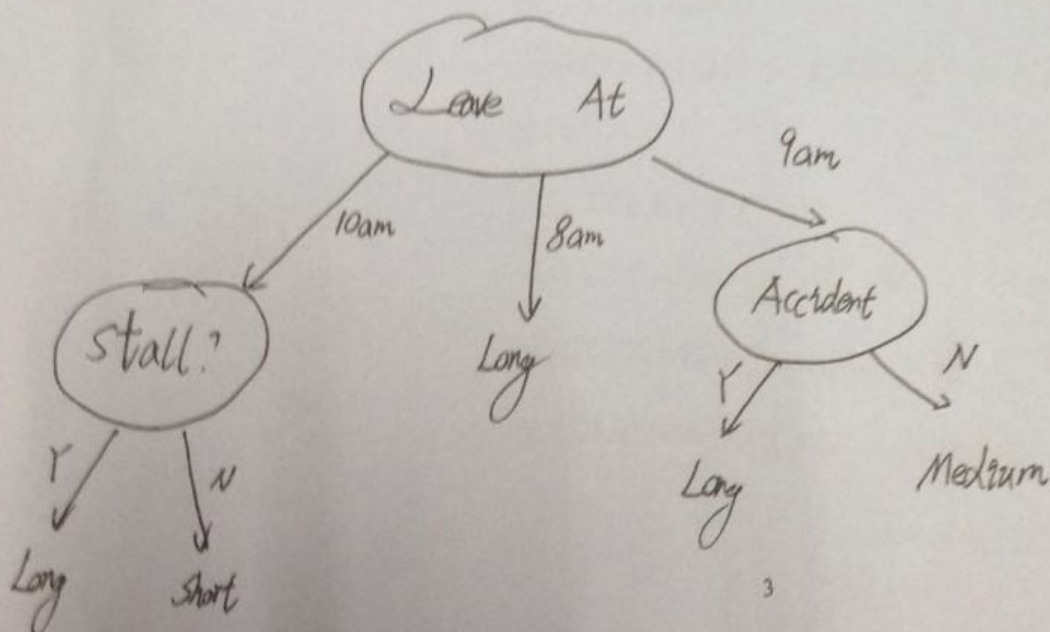
$$= 1.41196$$

$$\text{GAIN (commute/Hour)} = \text{Entropy (commute)} - E \text{(commute/Hour)}$$

$$= \text{Entropy (commute)} - [\frac{4}{13} \times E\text{(commute/Hour=8)} - \frac{3}{13} \times E\text{(commute/Hour=9)}$$

$$- \frac{2}{13} \times E\text{(commute/Hour=10)}]$$

$$= 1.41196 - 0.6 \times 11$$

$$= 0.768999$$

$$\text{GAIN (commute/Accident)} = \text{Entropy (commute)} - \text{Entropy (commute/Accident)}$$

$$= 1.41196 - 0.92357$$

$$= 0.496479$$

$$\text{GAIN (commute/Stall)} = \text{Entropy (commute)} - \text{Entropy (commute/Stall)}$$

$$= 1.41196 - 1.17071$$

$$= 0.248842$$

$$\text{GAIN (commute/Weather)} = 1.41196 - 1.28884 = 0.130719$$

## Decision Tree:

**Problem 2 :**

**For Hour :**

$$\text{Info}(\text{Hour} = 8) = 0$$

$$\text{Info}(\text{Hour} = 9) = -\tfrac{3}{5}\log_2\tfrac{3}{5} - \tfrac{2}{5}\log_2\tfrac{2}{5} = 0.97095$$

$$\text{Info}(\text{Hour} = 10) = -\tfrac{4}{5}\log_2\tfrac{4}{5} - \tfrac{1}{5}\log_2\tfrac{1}{5} = 0.72192$$

$$\text{Info}(\text{Hour}) = \tfrac{3}{13}\times 0 + \tfrac{5}{13}\times 0.97095 + \tfrac{5}{13}\times 0.72192 = 0.6511$$

$$\text{GAIN}(\text{Hour}) = 0.76845$$

**GINI :**

$$\text{GINI}(\text{Hour} = 8) = 0$$

$$\text{GINI}(\text{Hour} = 9) = 1 - (\tfrac{3}{5})^2 - (\tfrac{2}{5})^2 = 0.48$$

$$\text{GINI}(\text{Hour} = 10) = 1 - (\tfrac{4}{5})^2 - (\tfrac{1}{5})^2 = 0.32$$

$$\text{GINI}(\text{Hour}) = 1 - (\tfrac{7}{13})^2 - (\tfrac{4}{13})^2 - (\tfrac{2}{13})^2 = \tfrac{100}{169}$$

$$\text{GINI}(\text{Commute}/\text{Hour}) = \tfrac{3}{13}\times 0 + \tfrac{5}{13}\times 0.48 + \tfrac{5}{13}\times 0.32 = \tfrac{4}{13}$$

$$\triangle \text{GINI}(\text{Hour}) = \cdot\tfrac{100}{169} - \tfrac{4}{13} = 0.28402$$

---

**For Weather :**

$$\text{GAIN}(\text{Weather}) = 0.13071$$

**GINI :**

$$\text{GINI}(\text{Weather} = \text{Sunny}) = \tfrac{4}{9}$$

$$\text{GINI}(W = \text{Cloudy}) = \tfrac{3}{8}$$

$$\text{GINI}(W = \text{Rainy}) = \tfrac{2}{3}$$

~~GINI (W) = 13/18~~

$$\text{GINI}(\text{Commute}/\text{Weather}) = \tfrac{13}{18}$$

$$\triangle \text{GINI}(\text{Weather}) = 0.04043$$

---

**For Accident :**

$$\text{GAIN}(\text{Accident}) = 0.496479$$

**GINI :**

$$\text{GINI}(A = Y) = 0$$

$$\text{GINI}(A = N) = \tfrac{5}{8}$$

$$\text{GINI}(\text{Commute}/\text{Accident}) = \tfrac{5}{13}$$

$$\triangle \text{GINI}(\text{Accident}) = \text{~~0.0994~~}$$
$$0.207\text{ }\%$$

---

**For Stall**

$$\text{GAIN}(\text{Stall}) = 0.24884$$

**GINI**

$$\text{GINI}(S = Y) = 0$$

$$\text{GINI}(S = N) = \tfrac{16}{21}$$

$$\text{GINI}(\text{Commute}/\text{Stall}) = \tfrac{12}{63}$$

$$\triangle \text{GINI}(\text{Stall}) = 0.0994l$$

---

According to "**Theoretical Comparison between the Gini Index and Information Gain Criteria**" (*http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.57.9764&rep=rep1&type=pdf*), the frequency of disagreement of them is only 2%,therefore, it's hard to say which one is better.

But there are still some differences between the two algorithm. For ID3 algorithm and GAIN, log-likelihood might give higher scores to balanced portions when there are too many classes though.However, for CART and GINI, it can be nicer because it doesn't have logarithms and you can find the closed form for its expected value and variance under random split assumption.