

## Lab 05

1. After building the logistic regression model, write the answer for the following questions.

a) The key variables impacting the employee retention are:

- Satisfaction lower: Lower satisfaction increases attrition.
- Time spent in company: Employee with 51 years tend to leave.
- Salary: Low salaries lead to higher turnover.
- No of projects & avg monthly hours: Very high or low values affect retention.



1. The accuracy of logistic regression model is 78.40%. The accuracy overall is good but the model still needs improvements.

## 2. Decision Tree

Consider the following dataset. Calculate entropy & information gain. Use target variable 'classification'. Identify whether splitting node should be on  $a_2$  or  $a_3$  attribute.

→ instance	$a_2$	$a_3$	classification
1	hot	High	No
2	hot	High	No
6	cool	high	No
7	hot	high	No
8	hot	normal	Yes

### Attribute $a_2$

Values ( $a_2$ ) = Hot, Cool

$$S = [1, 4] \text{ Entropy}(S) = -\frac{1}{5} \log_2\left(\frac{1}{5}\right) - \frac{4}{5} \log_2\left(\frac{4}{5}\right) = 0.9209$$

$$S_{\text{hot}} = [1, 3] \text{ Entropy}(S_{\text{hot}}) = -\frac{1}{4} \log_2\left(\frac{1}{4}\right) - \frac{3}{4} \log_2\left(\frac{3}{4}\right) = 0.8112$$

$$S_{\text{cool}} = [0, 1] \text{ Entropy}(S_{\text{cool}}) = 0.0$$

$$\text{Gain}(S, a_2) = \text{Entropy}(S) - \sum_{v \in \{\text{hot}, \text{cool}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$= \text{Entropy}(S) - \frac{4}{5} \text{Entropy}(S_{\text{hot}}) - \frac{1}{5} \text{Entropy}(S_{\text{cool}})$$

$$= 0.9209 - \frac{4}{5} (0.8112) - \frac{1}{5} (0.0)$$

$$= 0.3219$$



Attribute a3 : Value (a3) : High, Normal

$$S_{a3} = [1, 1, 1] \Rightarrow \text{Entropy}(S_{a3}) = -\frac{1}{3} \log_2 \left( \frac{1}{3} \right) - \frac{1}{3} \log_2 \left( \frac{1}{3} \right) - \frac{1}{3} \log_2 \left( \frac{1}{3} \right) = 0.9709$$

$$S_{\text{High}} = [0, 1, 1] \Rightarrow \text{Entropy}(S_{\text{High}}) = 0.0$$

$$S_{\text{Normal}} = [1, 1, 0] \Rightarrow \text{Entropy}(S_{\text{Normal}}) = 0.0$$

$$\text{Gain}(S, a3) = \text{Entropy}(S) - \frac{4}{5} \text{Entropy}(S_{\text{High}}) - \frac{1}{5} \text{Entropy}(S_{\text{Normal}})$$

$$= 0.9709 - \frac{4}{5} (0.0) - \frac{1}{5} (0.0)$$

$$= 0.9709$$

$$\therefore \text{Gain}(S, a2) = 0.3719$$

$$\text{Gain}(S, a3) = 0.9709 \quad (\text{Max gain})$$

⇒ After building decision tree model with the following aims

1. For "iris.csv" dataset

- The accuracy score of model was 1.00 (100%) meaning the model perfectly classified all samples
- The confusion matrix shows that the model made no errors across prediction with the actual class:  $\begin{bmatrix} 10 & 0 & 0 \\ 0 & 9 & 0 \\ 0 & 0 & 11 \end{bmatrix}$
- No misclassification is observed since all diagonal values in confusion matrix are zero.

2. For "petal\_consumption.csv" dataset

- The regression tree structure split data to minimize HSB, with leaf node predicting the avg petal consumption
- ~~most important factors is petal, to v & level.~~
- The regression tree predicts continuous values while a classifier predicts categories.