

CS 484: Introduction to Machine Learning

Autumn 2020 Assignment 3

Question 1 (10 points)

Prove that $E = -\sum_{j=1}^K p_j \log_2(p_j)$ attains its maximum value when $p_j = p_K = 1/K$.

Hint: (1) re-express $E = -\sum_{j=1}^{K-1} p_j \log_2(p_j) - p_K \log_2(p_K)$, (2) use this equality $\sum_{j=1}^K p_j = 1$ in

calculating the partial derivatives $\partial E / \partial p_j, j=1, \dots, (K-1)$, and (3) solve the equations

$$\partial E / \partial p_j = 0, \quad j=1, \dots, (K-1).$$

Question 2 (10 points)

Suppose the predicted target value of a new observation is randomly assigned one of the target categories according to the categories' distribution. Argue analytically that the Gini Impurity is the probability of incorrect classification.

Hint: Probability of incorrect classification is the sum of these products of probabilities $\text{Prob}(\text{Do Not Classify to Category } i \mid \text{given the observation is from the Category } i) \times \text{Prob}(\text{an observation is drawn from the Category } i)$.

Question 3 (10 points)

Argue analytically that a completely impure node yields the highest Gini Impurity.

Question 4, 5, and 6

You will train a decision tree model to predict the usage of a car. The data is the `claim_history.csv` that contains 10,302 observations. The analysis specifications are:

Target Field

- **CAR_USE.** The usage of car. This field has two categories, namely, *Commercial* and *Private*. The *Commercial* category is the Event value.

Nominal Feature

- **CAR_TYPE.** The type of car. This feature has six categories, namely, *Minivan*, *Panel Truck*, *Pickup*, *SUV*, *Sports Car*, and *Van*.
- **OCCUPATION.** The occupation of car owner. This feature has nine categories, namely, *Blue Collar*, *Clerical*, *Doctor*, *Home Maker*, *Lawyer*, *Manager*, *Professional*, *Student*, and *Unknown*.

Ordinal Feature

- **EDUCATION.** The education level of car owner. This feature has five ordered categories which are *Below High School* < *High School* < *Bachelors* < *Masters* < *Doctors*.

Analysis Specifications

- **Partition.** Specify the target field as the stratum variable. Use stratified simple random sampling to assign 70% of the observations to the Training partition, and the remaining 30% of the observations into the Test partition. The random state is 60616.
- **Decision Tree.** The maximum number of branches is two. The maximum depth is two. The split criterion is the Entropy metric.

Question 4 (10 points)

Please answer the following questions about your Data Partition step. You may call the `train_test_split()` function in the `sklearn.model_selection` module in your code.

- (5 points). Please provide the frequency table (i.e., counts and proportions) of the target field in the Training partition?
- (5 points). What is the probability that an observation will be assigned to the Test partition given that `CAR_USE` is *Private*?

Question 5 (30 points)

Please provide information about your decision tree that is trained on the Training partition. You will need to write your own Python program to find the answers.

- (5 points). What is the entropy value of the root node?
- (5 points). What is the split criterion (i.e., feature name and values in the two branches) of the first layer?
- (5 points). What is the entropy of the split of the first layer?
- (5 points). Describe all your leaves (i.e., terminal nodes) in a table. Please include the decision rules and the counts of the target values.

- e) (5 points). What is the Kolmogorov-Smirnov statistic?
- f) (5 points). What is your suggested event probability cutoff value?

Question 6 (30 points)

Please apply your decision tree to the Test partition and then provide the following information. You will choose whether to call sklearn functions or write your own Python program to find the answers.

- a) (5 points). Based on your suggested Kolmogorov-Smirnov event probability cutoff value as the threshold, what is the Misclassification Rate in the Test partition?
- b) (5 points). What is the Root Average Squared Error in the Test partition?
- c) (5 points). What is the Area Under Curve in the Test partition?
- d) (5 points). What is the Gini Coefficient in the Test partition?
- e) (5 points). What is the Goodman-Kruskal Gamma statistic in the Test partition?
- f) (5 points). Generate the Receiver Operating Characteristic curve for the Test partition. The axes must be properly labeled. Also, include the diagonal reference line.