

**CMPT459 Spring 2020**  
**Data Mining**  
**Martin Ester**  
**TAs: Ruijia Mao and Ruchita Rozario**

**Solution of Assignment 1**

Solution posted and discussed in class: January 30, 2020

**Assignment 1.1**

Given the following dataset of bank clients with categorical attributes Age, Salary, City, and Creditworthiness:

Age	Salary	City	Creditworthiness
young	high	Vancouver	bad
medium	medium	Burnaby	bad
young	low	Vancouver	bad
old	medium	Coquitlam	good
old	high	Richmond	good
medium	low	Richmond	bad
young	medium	Vancouver	bad
old	low	Burnaby	bad
young	medium	Coquitlam	good
young	medium	Vancouver	good

We want to use this dataset to train a decision tree classifier to predict the creditworthiness of a bank client, i.e. we are using Creditworthiness as the class label and the other attributes as features. Suppose that the gini index is used to determine the split attributes.

a) Which split attribute does the decision tree algorithm select for the root of the decision tree? Show all the relevant gini indexes that you have to compute.

The algorithm chooses City for the root, since City has the smallest gini index of 0.25.

$\text{gini}(D) = 1 - \sum p_i^2$  for  $i$  from 1 to 2, i.e. for the classes “good” and “bad”

$\text{gini}(D|\text{age}) = 5/10 * \text{gini}(D|\text{age}=\text{young})$   
 $\quad + 2/10 * \text{gini}(D|\text{age}=\text{medium})$   
 $\quad + 3/10 * \text{gini}(D|\text{age}=\text{high})$   
 $= 0.5 * 0.48 + 0.2 * 0.0 + 0.3 * 0.44 = 0.24 + 0.0 + 0.132$   
 $= 0.372$

$\text{gini}(D|\text{age}=\text{young}) = 1 - (3/5)^2 - (2/5)^2 = 1 - 13/25 = 12/25 = 0.48$

$\text{gini}(D|\text{age}=\text{medium}) = 1 - (0/2)^2 - (2/2)^2 = 1 - 1 = 0.0$

$\text{gini}(D|\text{age}=\text{old}) = 1 - (1/3)^2 - (2/3)^2 = 1 - 5/9 = 0.44$

$\text{gini}(D|\text{salary}) = 3/10 * \text{gini}(D|\text{salary}=\text{low})$

$$\begin{aligned}
& + 5/10 * \text{gini}(D|\text{salary}=\text{medium}) \\
& + 2/10 * \text{gini}(D|\text{salary}=\text{high}) \\
& = 0.3 * 0.0 + 0.5 * 0.48 + 0.2 * 0.5 = 0.0 + 0.24 + 0.1 \\
& = 0.34
\end{aligned}$$

$$\text{gini}(D|\text{salary}=\text{low}) = 1 - (3/3)^2 - (0/3)^2 = 1 - 1 = 0.0$$

$$\text{gini}(D|\text{salary}=\text{medium}) = 1 - (2/5)^2 - (3/5)^2 = 1 - 13/25 - 12/25 = 8/25 = 0.48$$

$$\text{gini}(D|\text{salary}=\text{high}) = 1 - (1/2)^2 - (1/2)^2 = 1 - 1/2 = 0.5$$

$$\begin{aligned}
\text{gini}(D|\text{city}) &= 4/10 * \text{gini}(D|\text{city}=\text{Vancouver}) \\
& + 2/10 * \text{gini}(D|\text{city}=\text{Burnaby}) \\
& + 2/10 * \text{gini}(D|\text{city}=\text{Coquitlam}) \\
& + 2/10 * \text{gini}(D|\text{city}=\text{Richmond}) \\
& = 0.4 * 0.375 + 0.2 * 0.0 + 0.2 * 0.0 + 0.2 * 0.5 = 0.15 + 0.0 + 0.0 + 0.1 \\
& = 0.25
\end{aligned}$$

$$\text{gini}(D|\text{city}=\text{Vancouver}) = 1 - (1/4)^2 - (3/4)^2 = 1 - 10/16 = 0.375$$

$$\text{gini}(D|\text{city}=\text{Burnaby}) = 1 - (2/2)^2 - (0/2)^2 = 1 - 1 - 0 = 0.0$$

$$\text{gini}(D|\text{city}=\text{Coquitlam}) = 1 - (0/2)^2 - (2/2)^2 = 1 - 1 = 0.0$$

$$\text{gini}(D|\text{city}=\text{Richmond}) = 1 - (1/2)^2 - (1/2)^2 = 1 - 1/4 - 1/4 = 0.5$$

b) Does the gini index favor attributes with few or with many distinct values, or in other words, is an attribute with few or with many distinct values more likely to have a low gini index? Explain your answer. What rule does this suggest to break ties, i.e. to select one attribute among multiple attributes that all have the same smallest gini index?

The gini index favors attributes with many distinct values. The more distinct values, the smaller the partitions of the training dataset corresponding to the individual attribute values, and the more likely the partitions are pure with respect to the class label. In the extreme case, the partitions contain only one training record and achieve the minimum gini index of 0. This suggests to break ties by selecting, among the attributes that all have the same smallest gini index, the attribute with the smallest number of attribute values.

## Assignment 1.2

Given the above dataset of bank clients with categorical attributes Age, Salary, City, and Creditworthiness. We have trained a Naïve Bayes classifier to predict the creditworthiness of a bank client, i.e. we are using Creditworthiness as the class label. Here are the parameters of the classifier (computed without using Laplacian smoothing, and using C as abbreviation for Creditworthiness):

$$P(C=\text{good})=0.4 \quad P(C=\text{bad})=0.6$$

$$\begin{aligned}
P(\text{Age}=\text{young}|C=\text{good}) &= 0.5 & P(\text{Age}=\text{medium}|C=\text{good}) &= 0.0 & P(\text{Age}=\text{old}|C=\text{good}) &= 0.5 \\
P(\text{Age}=\text{young}|C=\text{bad}) &= 0.5 & P(\text{Age}=\text{medium}|C=\text{bad}) &= 0.34 & P(\text{Age}=\text{old}|C=\text{bad}) &= 0.16
\end{aligned}$$

$$\begin{aligned}
P(\text{Salary}=\text{low}|C=\text{good}) &= 0.0 & P(\text{Salary}=\text{medium}|C=\text{good}) &= 0.75 \\
P(\text{Salary}=\text{high}|C=\text{good}) &= 0.25 \\
P(\text{Salary}=\text{low}|C=\text{bad}) &= 0.5 & P(\text{Salary}=\text{medium}|C=\text{bad}) &= 0.34 \\
P(\text{Salary}=\text{high}|C=\text{bad}) &= 0.16
\end{aligned}$$

$$\begin{aligned}
P(\text{City}=\text{Vancouver}|C=\text{good}) &= 0.25 \\
P(\text{City}=\text{Coquitlam}|C=\text{good}) &= 0.5
\end{aligned}$$

$$\begin{aligned}
P(\text{City}=\text{Burnaby}|C=\text{good}) &= 0.0 \\
P(\text{City}=\text{Richmond}|C=\text{good}) &= 0.25
\end{aligned}$$

$$P(\text{City}=\text{Vancouver}|\text{C}=\text{bad}) = 0.5$$

$$P(\text{City}=\text{Coquitlam}|\text{C}=\text{bad}) = 0.0$$

$$P(\text{City}=\text{Burnaby}|\text{C}=\text{bad}) = 0.34$$

$$P(\text{City}=\text{Richmond}|\text{C}=\text{bad}) = 0.16$$

- a) What class label (Creditworthiness) does the Naïve Bayes classifier predict for an unseen client with Age = young, Salary = low and City = Vancouver? Show the necessary computations.

The classifier predicts the Creditworthiness to be bad.

For good:

$$P(\text{C}=\text{good}) * P(\text{Age}=\text{young}|\text{C}=\text{good}) * P(\text{Salary}=\text{low}|\text{C}=\text{good}) * P(\text{City}=\text{Vancouver}|\text{C}=\text{good})$$

$$= 0.4 * 0.5 * 0.0 * 0.25 = 0.0$$

For bad:

$$P(\text{C}=\text{bad}) * P(\text{Age}=\text{young}|\text{C}=\text{bad}) * P(\text{Salary}=\text{low}|\text{C}=\text{bad}) * P(\text{City}=\text{Vancouver}|\text{C}=\text{bad})$$

$$= 0.6 * 0.5 * 0.5 * 0.5 = 0.075$$

- b) Given the Naïve Bayes assumption and the parameters of the above Naïve Bayes classifier, what is the probability of observing a client with Age = old, Salary = medium and City = Vancouver? Show the necessary computations.

The probability of a client with these attribute values is the marginal probability aggregating over the conditional probabilities of all classes, i.e.

$$P(\text{Age}=\text{old}, \text{Salary}=\text{medium}, \text{City}=\text{Vancouver}) =$$

$$P(\text{Age}=\text{old}, \text{Salary}=\text{medium}, \text{City}=\text{Vancouver}|\text{good}) * P(\text{good})$$

$$+ P(\text{Age}=\text{old}, \text{Salary}=\text{medium}, \text{City}=\text{Vancouver}|\text{bad}) * P(\text{bad})$$

According to the Naïve Bayes assumption,

$$P(\text{Age}=\text{old}, \text{Salary}=\text{medium}, \text{City}=\text{Vancouver}|\text{good}) =$$

$$P(\text{Age}=\text{old}|\text{good}) * P(\text{Salary}=\text{medium}|\text{good}) * P(\text{City}=\text{Vancouver}|\text{good}) =$$

$$0.5 * 0.75 * 0.25 = 0.09375$$

$$P(\text{Age}=\text{old}, \text{Salary}=\text{medium}, \text{City}=\text{Vancouver}|\text{bad}) =$$

$$P(\text{Age}=\text{old}|\text{bad}) * P(\text{Salary}=\text{medium}|\text{bad}) * P(\text{City}=\text{Vancouver}|\text{bad}) =$$

$$0.16 * 0.34 * 0.5 = 0.0272$$

In conclusion,

$$P(\text{Age}=\text{old}, \text{Salary}=\text{medium}, \text{City}=\text{Vancouver})$$

$$= 0.09375 * 0.4 + 0.0272 * 0.6 = 0.0375 + 0.01632 = 0.05382$$

The probability of such a client is 0.05382.