

Problem assignment 8

Due: Wednesday, April 4, 2007

Problem 1. Customer profiling and predictions

In this problem we use the the dataset of consumers preferences for different product brands. The dataset consists of 6 attributes corresponding to individual product categories, the values correspond to different brands in respective categories. For the sake of simplicity, we use numerical labels to distinguish the brands. We assume that each product can appear only in one category. Each customer corresponds to one entry (row) in the dataset. The number of brand products in each product category is listed in the following table.

Product category	1	2	3	4	5	6
Number of brands	5	3	3	4	5	4

The dataset we have available is not complete and some entries reflecting the preferences of the customer towards different brands are missing. The label 0 indicates the missing value. Despite the missingness, our goal is to use the data for prediction. We will use the Naive Bayes model to achieve this goal. We construct the Naive Bayes model in the following way: the root node of the model represents groups of customer (not observable), the leaves are product categories with values corresponding to different brand names. The assumption introduced by the Naive Bayes model is that the preferences of the customer towards brands in different categories are conditionally independent given the group the customer belongs to. The relations between the group and any product category are stochastic and may vary for different individual. The key idea here is that identifying the customers' group helps us to make the prediction about the product preference.

Part a. In the first part we need to learn the parameters of the Naive Bayes model with the hidden class variable and missing values in the dataset. We assume that there are four different groups of customers. Thus the hidden root variable can have 4 different values.

To learn the model we want ML estimate of parameters of the model given data. We will use EM algorithm to obtain them. What follows is a description of EM for this case that will help you to implement the procedure. You will need to **turn in your EM code** in the file *main1a.m*.

EM for the Naive Bayes model

Different customers in the systems are represented through a class variable C . We assume there is a fixed number of classes $M = 4$. Let $X = \{X_1, X_2, \dots, X_n\}$ be a set of relevant attribute variables. Let $D = \{D^1, D^2, \dots, D^N\}$ be a data set such that each D^l is a sample of X with possibly missing values. For example D^l can consist of $\{X_1^l = 3, X_3^l = 1\}$ or $\{X_1^l = 4, X_3^l = 1, X_4^l = 1\}$.

Our goal is find the set of parameters Θ maximizing the likelihood of D , that is, $p(D|\Theta)$. In the case of the Naive Bayes model parameters in Θ correspond to

- a set of priors on classes π_1, \dots, π_M (s.t. $\sum_{j=1}^M \pi_j = 1$).
- a set of parameters θ_{ijk} representing components of conditional distributions $p(X_i = k|C = j)$.

Let H represents the set of hidden variables. These correspond to the hidden class variables C and variables corresponding to missing values in the dataset. For example, if the value for X_2^l is missing in D^l we model it through a variable in H .

The optimization of the likelihood $p(D|\Theta)$ can be carried in a standard way in terms of the loglikelihood optimization. There are various ways to approach this optimization problem. We apply the EM algorithm.

In the EM algorithm we define:

$$Q(\Theta|\Theta') = \sum_{j=1}^M N'_{C=j} \log \pi_j + \sum_{i=1}^d \sum_{k=1}^M \sum_{j=1}^M N'_{C=j, X_i=k} \log \theta_{ijk}$$

where:

$$N'_{C=j} = \sum_{l=1}^N E_{H|D, \Theta'}(\delta_{C^l=j})$$

and

$$N'_{C=j, X_i=k} = \sum_{l=1}^N E_{H|D, \Theta'}(\delta_{C^l=j, X_i^l=k}).$$

The expectation terms are:

$$E_{H|D, \Theta'}(\delta_{C^l=j}) = p(C^l = j|D^l, \Theta')$$

$$E_{H|D, \Theta'}(\delta_{C^l=j, X_i^l=k}) = \begin{cases} p(C^l = j|D^l, \Theta') & \text{if } X_i^l = k \text{ is in } D \\ p(C^l = j|D^l, \Theta')p(X_i^l = k|C^l = j, D^l, \Theta') & \text{if } X_i^l \text{ is missing in } D \\ 0 & \text{otherwise} \end{cases}$$

In EM we maximize $Q(\Theta|\Theta')$. The key advantage here is that this maximization can be carried in the closed form, leading to new parameter estimates:

$$\pi_j = \frac{N'_{C=j}}{\sum_{j=1}^M N'_{C=j}}$$

$$\theta_{ijk} = \frac{N'_{C=j, X_i=k}}{\sum_{k=1}^{|X_i|} N'_{C=j, X_i=k}}.$$

Additional tasks:

- The EM algorithm may take a while to converge as parameters can change very slowly at the end. Thus, one typically stops the EM iterations once only small changes in Q functions are observed. Choose a small constant (0.001) to stop the algorithm.
- Try to run the EM algorithm multiple times on different sets of initial parameters. Analyze and explain your observations in your report.

Part b. Once we learn the model we can use it to do predictions. Use testing dataset to predict the value of attribute 6 for each customer. Note that also here there are missing values. Compute the confusion matrix and mean misclassification error for predicting attribute 6 based on all other 5 attributes. For the prediction use the ML estimates of parameters Θ obtained in part a. To do the prediction simply compute the probability $p(x_6^l = k | D_{mod}^l, \Theta_{ML})$ of attribute 6 having value k for each customer in the test dataset, and do it for all possible values the attribute can take. D_{mod}^l is the data entry for the customer l with attribute 6 removed. Select the value with the best probability. To compute the confusion matrix and the misclassification error compare the predicted and actual value of the attribute. Turn in the program code in file *main1b.m*.

Part c. Assume a predictor in which the value of attribute 6 is selected uniformly at random. What is the worst case mean misclassification error for such a case? Compare the worst case result to the result obtained through the Naive Bayes model.

Extra credit. Run the randomized uniform predictor on the test set. Is the result obtained by the Naive Bayes predictor statistically significantly different from the randomized predictor at significance level 0.05?