

Athens University of Economics and Business

MSc in Business Analytics

Data Mining – Assignment 1

Deadline: 24/5/2021

Group assignment (groups of up to 3 people).

The assignment corresponds to 25% of the total grade of the course.

Discussions between groups are recommended but collaborating on the actual solutions is considered cheating and will be reported.

There will be no extension of the assignment deadline!

Professor: Y.Kotidis (kotidis@aueb.gr)

Assistant responsible for this assignment: I.Filippidou (filippidou@aueb.gr)

Assignment 1

In this assignment you will implement a simple workflow that will assess the similarity between supermarket customers. The workflow will be used to compute, and suggest for any input customer, a list of his/her 10 most similar other customers. Moreover, you will be using these results to predict the rating of a customer. To calculate the (dis)similarity between customers you will first compute the dissimilarity for every given attribute as discussed in lecture “Measuring Data Similarity”. In particular:

1) Import and pre-process the dataset with customers

Download the groceries.csv dataset from moodle. This dataset contains demographic characteristics of supermarket 10000 customers along with a list of groceries they bought. Below is a description of the available attributes:

Customer ID: The unique id of the customer.

Age: The age of the customer.

Sex: Male-Female.

Marital Status: Married, Single, Divorced.

Education: Primary, Secondary, Tertiary.

Annual Income: The annual customer income.

Customer Rating: The rating of the supermarket from the customer (Poor, Fair, Good, Very Good, Excellent).

Persons in Household: Number of persons in the household.

Occupation: The occupation of each customer (retired, housemaid, unemployed, management, entrepreneur, blue-collar, self-employed, services, technician).

Groceries: A list of the customer groceries.

For any numerical missing values, you should replace them with the average value of the attribute in the dataset (keeping the integer part of the average).

2) Compute data (dis-)similarity

To assess the similarity between the customers you could form the dissimilarity matrix for all given attributes. As described in lecture “Measuring Data Similarity”, for every given attribute you first distinguish its type (categorical, ordinal, numerical or set) and then compute the dissimilarity of its values accordingly. For set similarity use the Jaccard similarity between sets. Then, you can calculate the average of the computed dissimilarities in order to derive the dissimilarity over all attributes. Depending on the machine used to implement this assignment you should decide whether it is feasible to compute the dissimilarity matrices, or, have the computations performed on-the-fly for a pair of customers.

3) Nearest Neighbor (NN) search

Using the implementation of the previous step, you will calculate the 10-NN (**most similar**) customers for the customers with ids listed below:

73, 563, 1603, 2200, 3703, 4263, 5300, 6129, 7800, 8555

For this task your script must take as input the customer-id and return the list of her 10 nearest neighbors (**most similar**), along with the corresponding **similarity score**.

An example of the script output for customer id =1 follows:

10 NN for Customer 1	
Customer ID	Similarity Score
7749	
7931	
9514	
628	
6918	
4230	
3148	
4647	
2105	
8050	

4) Customer rating prediction

For this assignment you will implement a classification algorithm which, for a given customer, will predict his rating (poor, fair, good, very good, excellent) for the supermarket. In order to implement the classification for a given customer you need to:

1) Calculate the similarities between the given customer and all other customers and compute his 10-nn (most similar) customers. **IMPORTANT: In the similarity calculations for this step you need to exclude the customer rating attribute.**

2) Based only on the 10 most similar customers computed in the previous step, predict the customer rating using:

- The average rating of the 10 most similar customers (rounded to the nearest integer).
- The weighted average rating of the 10 most similar customers (rounded to the nearest integer).

$$\text{Weighted average rating} = \frac{\sum_{i=1}^{10} \text{similarity}(i) * \text{rank}(\text{rating}(i))}{\sum_{i=1}^{10} \text{similarity}(i)}$$

Where:

- $\text{rating}(i)$ = the rating of the i -th nearest neighbor ($i=1$ for the most similar customer)
- $\text{similarity}(i)$ = the similarity of the i -th nearest neighbor with the given customer

3) For the evaluation of your classification algorithm you will use the 50 first records of the groceries dataset and predict the rating for them. Then, for all $n=50$ records calculate the Mean Prediction Error for both prediction methods.

$$\text{Mean Prediction Error} = \frac{\sum_{i=1}^n |\text{rank}(\text{Predicted rating}(i)) - \text{rank}(\text{True rating}(i))|}{n}$$

Assignment handout:

- 1) A report (pdf) describing in detail any processing and conversion you made to the original data and the reasons it was necessary. The report will also contain examples of how to use your script and its **output to the list of customers provided at step 3 (10-NN and the corresponding similarity scores for every given id)**. Also, in your report you should describe how to use your script and its output **for the classification system at step 4 (for the first 50 records of the dataset) for both prediction methods**. Comment on the mean prediction error of both methods and on any other conclusions you have made. The first page of the report should clearly state the names and student ids of the members of the group. Alternatively you could provide your jupyter notebook.
- 2) Your code. Implementation can be done in any programming language and should be accompanied by the necessary comments and remarks.
- 3) The pdf and the required programs/scripts should be uploaded to moodle until the assignment deadline. You should create a compressed (e.g. zip/tar) file containing the report, your code and any other files required for executing your script (you do not need to include the original dataset). The name of the compressed file should include the student ids of the members of the group.