

Department of Computer Science  
University of Pretoria

Data Mining  
COS 781

Assignment 3

September 26, 2019

## 1 Objectives

This assignment aims to achieve the following general learning objectives:

- To gain experience with a real-world data set, in the context of data mining and exploratory data analysis;
- To experiment with data mining approaches discussed in this course;
- To gain experience in formal, scientific report writing.

## 2 Plagiarism Policy

The Department of Computer Science considers plagiarism to be a serious offence. Disciplinary action will be taken against students who commit plagiarism. Plagiarism includes copying someone else's work without consent, copying a friend's work (even with consent) and copying material from the Internet. Copying will not be tolerated in this module.

For a formal definition of plagiarism, the student is referred to <http://www.ais.up.ac.za/plagiarism/index.htm> (from the main page of the University of Pretoria site, follow the **Library** quick link, and then click the **Plagiarism** link under the **Services** menu). If you have any questions regarding this please consult the lecturer to avoid any misunderstanding.

Note that all assignments submitted for this module implicitly agree to this plagiarism policy, and declare that the submitted work is the student's own work. Assignments may be checked using the Turnitin system. After plagiarism checking, assignments will not be permanently stored on the Turnitin database.

## 3 Submission Instructions

You will have to write a report for this assignment. The report should be in standard PDF format, preferably compiled using L<sup>A</sup>T<sub>E</sub>X. You will have to submit only the report. No additional files of any sort should be submitted. Do not submit files in any other format other than PDF. Failure to follow any of these instructions will result in a zero mark for the assignment.

Upload just your PDF file (named `s99999999.pdf`, where 99999999 is your student number), to the appropriate assignment upload on the course website. Multiple uploads are allowed, but only the last one will be marked. The deadline is **Friday, 18 October 2019, at 23:00**.

## 4 Data Set

The data set that you will work with for this assignment is a real-world data set produced by Statistics South Africa. The data set describes the General Household Survey for 2017, and contains data about South African households and individuals. You will use the data related to individuals.

The course website provides an archive containing an ASCII CSV (comma separated value) file, as well as a PDF file containing important data set meta information. Additional data file formats are available at <http://nesstar.statssa.gov.za:8282/webview/> located in the sub-folder Household Surveys → General Household Survey → General Household Survey (Revised 2017) → 2017 → General Household Survey 2017 (Person file). From this folder, select the Download link on the top right of the page.

## 5 Data Mining

Assignment 1 and Assignment 2 are closely linked. Assignment 2 requires you to perform data preparation on the data set described in Section 4. In Assignment 3 you will then apply a data mining approach to the prepared data set in order to extract useful knowledge. Your data mining should focus on the following points:

- Your knowledge extraction should be based on the data set you prepared for Assignment 2. Ensure that your report for Assignment 2 accurately describes the prepared data that you used as input to the data mining techniques you use for this assignment's analysis. Your data mining analysis must concentrate on the area of knowledge you decided to focus on in Assignment 2.
- You must use at least one of the following four techniques: self-organising maps, data clustering techniques, rule induction methods, or decision trees. You may also choose a second technique from the aforementioned four techniques to supplement your results. You may use two techniques in combination with one another, or use a second technique to verify or add to the results you derive from the application of the first technique.

## 6 Report

You must write a report describing the data preparation you performed in Section 5. The report should be of an academic nature. This means that:

- The report's structure should include all the aspects typically required of an academic paper (these include a title, abstract, introduction, methodology discussion, and conclusions).
- An adequate background discussion should be provided (this means that you must broadly discuss every data mining technique you use, and provide a reference to published sources that describe the techniques. It is **not** acceptable to cite the course slides or Wikipedia).
- It is not necessary to include a discussion on the data set or the data preparation you have done. Simply refer to your Assignment 2 submission where necessary. Your report will primarily focus on the data mining you have performed on the data set. You must also include a background discussion on all the techniques you have used, which covers all the important aspects of the technique. Also provide a conclusion in which you summarise what you have achieved through all the data data mining analysis you have performed.
- Adequate references must be provided (this means that you must cite appropriate sources for all techniques that you discuss, and all the reference details must be correct).
- The tone of language should be formal and scientific, and must use correct spelling and grammar.
- Any figures or graphs included should be clear and of a professional standard. Label all appropriate parts of graphs and other data visualisations, so that they can be easily interpreted (it is usually not sufficient to simply paste a screen shot of a visualisation from a data analysis program or package).

It is recommended that you consult several existing conference and journal papers, as a guide to the type of style you should adopt. There are many such sources freely available online.

Your report should be concise and to the point. Make sure that you leave no gaps in your descriptions of techniques or the procedures you followed. Include all relevant details, even if they seem obvious to you.

It is very strongly recommended that you use the L<sup>A</sup>T<sub>E</sub>X template for IEEE conference papers (available as `LatexSample.zip` in the folder for this assignment, on the course website), to typeset your paper. If you use this template, you should aim for a report (including references) of no more than six pages.

## 7 Marking

The following general breakdown will be used during the assessment of this assignment:

Category	Mark Allocation
Writing style and report structure	10 marks
Background information	10 marks
Data mining discussion	20 marks
References	10 marks
<b>TOTAL</b>	<b>50 marks</b>