

# 31250 Introduction to Data Analytics

## Assignment 2: data exploration and preparation

<b>Due date</b>	<b>11:59pm Friday, 15 May 2020</b>
<b>Marks</b>	Out of 100, weighted to 35% of your final mark.
<b>Submission format</b>	A report in Adobe PDF (preferable) or MS Word Doc. Please also upload the Excel spreadsheet containing your results.
<b>Filename</b>	xxxxxxxx_a2.pdf or xxxxxxxx_a2.doc where xxxxxxxx is your student id. xxxxxxxx_a2.xls for the spreadsheet. You may need to zip files to submit to Turnitin.
<b>Report format</b>	Around 20-25 pages with the information described below. Use 11 or 12 point Times or Arial fonts.
<b>Submit to</b>	Turnitin Assignment Task 2 on UTSOnline/Blackboard.

This assignment is individual work. Each of you will be working with an individual data set that you will be able to download from MS Team. Note the zip package includes data for everybody, named using student ID's. You just need to work on the file belongs to you and can discard the rest ones.

### Scenario

You have just started working as a data miner/analyst in the Analytics Unit of a company. The Head of the Analytics Unit has brought you a data set [a welcome present ;-)). The data set includes two files: description of the attributes and a table with the actual values of these attributes. The Head of the Analytics Unit has mentioned to you that this is some sort of demographic data that a potential client has provided for analysis. The Head of the Analytics Unit would like to have a report with some insights about that data, that he/she could deliver to the client. Your tasks include:

- understanding the specifics of the data set
- extracting information about each of the attributes, possible associations between them and other specifics of the data set.

The tasks in the assignment are specified below.

## Data sets

For this dataset you only have the attribute headings, no descriptions of what they mean. Each student is assigned an individual table with the actual values of these attributes. Please, download the file that is linked to your name.

## Tasks

### 1A. Initial data exploration

1. Identify the type of all the attributes {Age, Employment class, Fnlwgt,....., Salary} (nominal, ordinal, interval or ratio). If it's not clear you may need to justify why you choose the type.
2. Identify the values of the summarising properties for all the attributes including frequency, location and spread (e.g. value ranges of the attributes, frequency of values, distributions, medians, means, variances, percentiles, etc. - the statistics that have been covered in the lectures and materials given). Note that not all of these summary statistics will make sense for all the attribute types, so use your judgement! Where necessary, use proper visualisations for the corresponding statistics.
3. Using KNIME or other tools, explore your data set and identify any outliers, clusters of similar instances, "interesting" attributes and specific values of those attributes. Note that you may need to 'temporarily' recode attributes to numeric or from numeric to nominal. In the report include the corresponding snapshots from the tools and explanation of what has been identified there.

Present your findings in the assignment report.

### 1B. Data preprocessing

Perform each of the following data preparation tasks (each task applies to the original data) using your choice of tool:

- a. Use the following **binning** techniques to smooth the values of the **Age** attribute:
  - equi-width binning
  - equi-depth binning.

In the assignment report for each of these techniques you need to illustrate your steps. In your Excel workbook file place the results in **separate columns** in the corresponding spreadsheet. Use your judgement in choosing the appropriate number of bins - and justify this in the report.

- b. Use the following techniques to **normalise** the attribute **Age**:
  - min-max normalization to transform the values onto the range [0.0-1.0].
  - z-score normalization to transform the values.

In the assignment report provide explanation about each of the applied techniques. In your Excel workbook file place the results in separate columns in the corresponding spreadsheet.

- c. **Discretise** the **Age** attribute into the following categories: Child=0-12; Teenager=13-19; Young adult=20-39; Middle aged=40-59; Old aged=60+. Provide the frequency of each category in your data set.

In the assignment report provide explanation about each of the applied techniques. In your Excel workbook file place the results in a separate column in the corresponding spreadsheet.

- d. Binarise the **Marital status** variable [with values "0" or "1"].

In the assignment report provide explanation about the applied binarisation technique. In your Excel workbook file place the results in separate columns in the corresponding spreadsheet.

## 1C. Summary

At the end of the report include a summary section in which you summarise your findings. The summary **is not** a narrative of what you have done, but a condensed informative section of **what you have found** about the data that you should report to the Head of the Analytics Unit. The summary may include the most important findings (specific characteristics (or values) of some attributes, important information about the distributions, some clusters identified visually that you propose to examine, associations found that should be investigated more rigorously, etc.).

## Deliverables

The deliveries include:

- A report, which structure should follow the tasks of the assignment, and
- An Excel workbook file with individual spreadsheets for each task (spreadsheets should be labeled according to the task names, for example, "1A"). Each of the results of parts (a) through (d) in task 1B should be presented in a separate spreadsheet (and respectively table in the assignment report).

**Report:** In the report include a section (starting with a section title) for each of the tasks in this assignment.

Your report will likely be between 20-25 pages in length using an 11 or 12 point font, including title page and graphs. On average you will require between 15 and 23 hours to complete this assignment.

## **Assessment**

This assignment is assessed as individual work. The assessment criteria are given on MS Team.

## **Relationship to Objectives**

This assessment task addresses the following subject learning objectives (SLOs): 3, 4, 5 and 6

This assessment task contributes to the development of the following course intended learning outcomes (CILOs): A.1, B.1, B.2, B.3 and E.1.

## **Return of Assignments**

We plan to return marked assignments within 3 weeks of submission. Emails will be sent when marking is complete.

## **Academic Standards and Late Penalties**

Please refer to subject outline.