

Assignment 2: Data Analytics

The goal of this assignment is to solve a data analytics problem following the CRISP-DM Process. This being a class assignment rather than a real-life setting, several simplifications will have to be made. Particularly, you will need to make certain assumptions and simplifications in the course of the project, both because a real problem owner and data expert is not available, and deployment of the solution is obviously out of scope. We recommend working on it in groups of two, but you are also allowed to solve it individually

For performing the experiments you can use any Machine Learning platform of your choice, e.g. WEKA, Scikit-Learn, Spark / MLlib, Matlab, R, ... according to your preferences.

- Information on how to obtain and use WEKA, the open Source Machine Learning Platform from Waikato University, is available at <https://wekatutorial.com/>.
- For Scikit-Learn you may consult the tutorial available at <https://scikit-learn.org/stable/tutorial/index.html>
- For **Spark/MLlib** you can consult the manual available at <https://spark.apache.org/docs/latest/ml-guide.html>)
- Structure your **report** for this assignment based on the structure in this assignment paper. Provide detailed documentation of all steps to **ensure reproducibility** of all results based on the information provided.

A) Preparations

- (1) **Select a data set** from the OpenML Machine Learning Repository (<http://www.openml.org>), Kaggle (<https://www.kaggle.com/datasets>), the UCI Machine Learning Repository (<https://archive.ics.uci.edu/ml/index.php>) or a similar benchmark data repository meeting the **following requirements**:
 - posing a classification or regression problem
 - minimum 1.000 instances,
 - minimum 10 attributes,
 - minimum 4 class labels if it is a classification task
 - not an “artificial” dataset, i.e., a dataset consisting of purely synthetic, sampled or interpolated values (e.g. the BNG* datasets on OpenML)
 - where the attributes carry semantics that can be interpreted by you (i.e. not a collection of image files where features still would need to be extracted by you)
 - with a certain variety of attribute semantics and (preferably) also attribute types (i.e. not a data set with just 5000 bag-of-words features or greyscale histogram features of images)
 - where you understand the semantics of the data and the domain so that you can make reasonable assumptions on its use, the goals to be met.
- (2) **Register the dataset** you picked in the TUWEL Wiki. You must make sure that your dataset is not already used by somebody else! (first come, first serve - do it early to get a data set that you also find interesting to work with.)

Assignment 2: Data Analytics

B) Report

Prepare a report (see end of this assignment sheet for formatting and submission information) which documents your analysis containing at least (additional material/sections may be provided if you think they are important in your specific setting) the following sections as a reduced subset of the CRISP-DM process:

(1) Business Understanding

- a. Define and describe a scenario in which a business analytics task based on the data set you identified should be solved
- b. Define and describe the Business Objectives
- c. Define and describe the Business Success Criteria
- d. Define and describe the Data Mining Goals
- e. Define and describe the Data Mining Success Criteria

(2) Data Understanding:

Data Description Report presenting

- a. Data types,
- b. Statistical properties
- c. Data Quality aspects
- d. Visual Exploration of data properties and hypotheses

(3) Data Preparation report

- a. Analyze options and potential for derived attributes (note: if the potential is considered low, these obviously do not necessarily have to be applied for your analysis, but options should be documented)
- b. Analyze options for additional external data sources, attributes that might be useful to better address the business objectives or data mining goals (Note: this description may be hypothetical, i.e. you are not necessarily required to actually obtain and integrate the external data for the analysis)
- c. Describe other pre-processing steps considered, specifying which ones were applied or not applied due to which reason. (e.g. data cleansing, transformations, binning, scaling, outlier removal, attribute removal, transcoding, ...) at a level of detail that ensures reproducibility of changes to the data. (Code may be supplied as supplement to the submission in case you produce your own code)

(4) Modeling

- a. Identify suitable data mining algorithms and select one of these as the most suitable for your experiments, providing a brief justification.
- b. Identify the hyper-parameters available for tuning in your chosen algorithm and select one that you deem most relevant for tuning, providing a brief justification.
- c. Define and document a train / validation / test set split, considering where necessary appropriate stratification, any dependencies between data instances (e.g. time series data) and relative sizes of the respective subsets.
- d. Train the model on the training set and comparing the performance on the validation set to identify the best hyper-parameter setting, explicitly documenting all parameter settings (avoid stating simply to have used “default parameters”, focus on reproducibility of the results you report).

Assignment 2: Data Analytics

- e. Report suitable performance metrics supported, where possible, by figures/graphs showing the tuning process of the hyper parameter.

(5) Evaluation

- a. Apply the final model on the test data and document performance.
- b. Re-train the model with identical hyper-parameters using the full train and validation data and again apply it on the test data, documenting the performance
- c. Identify and document
 - i. state-of-the-art performance from the literature using the same (albeit potentially slightly differently pre-processed) data set from the literature.
 - ii. the expected base-line performance of a trivial acceptor / rejecter or random classifier
- d. Compare the performance achieved with the benchmark and baseline performances (c.f. Section 1e – Data Mining success Criteria) according to different metrics (i.e. overall, but also on per-class level (confusion matrix), micro/macro precision/recall in the case of classification tasks, regression errors in certain parts of the data space, ... (Note your goal is not necessarily to obtain a better result than what has been reported in the state of the art, this is not a grading criterion! On the other hand, if the performance of your classifiers is massively below the stat of the art (or even below a random baseline or trivial acceptor / rejecter) you may want to investigate the reason...))
- e. Compare the performance obtained with the Data Mining success criteria defined in the business understanding phase.

(6) Deployment:

- a. Compare the performance obtained with respect to the needs for addressing the business success criteria and provide recommendations for deployment (fully automatic, hybrid solutions, deploying only for a part of the data space, ...) as well as recommendations for subsequent analysis.
- b. Consider and briefly document potential ethical aspects as well as impact assessment / risks identified in deployment
- c. Document aspects to be monitored during deployment, specifying triggers that should lead to intervention.
- d. Briefly re-visit reproducibility aspects reflecting on aspects well documented and those that might pose a risk in terms of reproducibility based solely on the information provided in this report

(7) Summarize your findings

- a. Briefly summarize your overall findings and lessons learned
- b. **(optional)** Provide **feedback on this exercise** in general: which parts were useful / less useful; which other kind of experiment would have been interesting, ... (this section is, obviously, optional and will not be considered for grading. You may also decide to provide that kind of feedback anonymously via the feedback mechanism in TISS – in any case we would appreciate learning about it to adjust the exercises for next year following a major re-structuring this year based on feedback obtained.)

Assignment 2: Data Analytics

Submission guidelines:

- **Upload ONE [zip/tgz/rar] file** to TUWEL that **contains (1) your report as a PDF file** (no Word files, no TEX sources) **and (2) any auxiliary files needed for reproducing your experiments** (i.e. any scripts, transformation tools, config files etc. that you produced and that represent information not sufficiently documented in the report). You **must follow this naming convention**:
 - o BI2020_gr_<groupno>_<Matnr.1>_<Matnr.2>.zip
 - o Example: A submission of group 99 with 2 students (ids: 00059999, 00039999) looks like this: BI2020_gr_99_00059999_00039999.[zip/tgz/rar]
 - o Example: A submission of a single student (with group no. 1) (id: 00968787) looks like this: BI2020_gr_01_00968787.[zip/tgz/rar]
 - o Apply the same naming convention to the report (but obviously with pdf extension)
- **Follow the ACM formatting guidelines, using the templates provided at** <https://www.acm.org/publications/proceedings-template>. (Conference Proceedings Style File, 2-column Layout) LaTeX recommended, but Word/OpenOffice template is obviously also ok.
- **Put your names, group number and your student IDs in the report!** (as author info)
- **Report page limit: Maximum 12 pages. Focus on the key aspects!**
- **Use graphs** to visualize findings. Do not just print graphs, also describe what they mean.
- **Use tables** to combine findings and other information for maximum overview whenever possible. Describe what you show and explain the data. Clarify, don't mystify.
- Consider issues of **reproducibility**: ensure you provide sufficient information allowing others to re-produce your experiments.
- **Enumerate and label ALL figures, equations and tables** and refer to them in the report --
 - describe, explain and integrate them with the text. It must be clear to the reader what information can be learned from them.

General advice:

- Reserve plenty of **time for “playing” with the data** and start early.
- **Collaboration between groups** is welcome, **but** ensure your group uses a **unique data set**.
- **Collaboration inside the group**: Try to perform at least part of the tasks within the group together. Specifically, discuss the results amongst each other. Subdividing and **solving tasks alone will cost you more time and not meet the goals of the exercise**. Specifically, we discourage completely splitting the assignment into sub-parts distributed across group members. Collaborate, brainstorm and discuss what you find. In an eventual review meeting, **every group member has to demonstrate knowledge of each aspect of the work and the steps taken**.
- Make sure the **structure of the report** follows the **structure of the tasks** provided here.

Submission Deadline: January 24, 2021, 23:55