

## Data Mining Projects. (2016-2017)

There are 3 types of projects:

- **Oriented towards algorithms** (max grade:10)
- **Oriented towards data** (max grade:10)
- **Oriented towards software tools/library** (max grade:8)

### A. Projects oriented towards algorithms

Projects of type A consists of:

- A report describing the particularity of the problem and at least one of the algorithms which can solve that problems (based on the starting bibliography or on other related works) and reporting results obtained by applying the implemented algorithm(s).
- An implementation from scratch of an algorithm (the programming language is at your choice).

Topics for projects of type A:

1. Algorithms for feature selection (e.g. implementation of Relief algorithm or of a greedy-like forward algorithm). Biblio: FeatureSelection folder
2. Algorithms for feature discretization (e.g. implementation of Holte 1R discretizer). Biblio: FeatureDiscretization folder
3. Algorithms for decision trees induction (e.g. implementation of the ID3 algorithm - <http://www.cise.ufl.edu/~ddd/cap6635/Fall-97/Short-papers/2.htm>). Biblio: DecisionTree folder
4. Covering algorithms (e.g. implementation of PRISM algorithm). Biblio: CoveringAlgorithms folder
5. K-Nearest Neighbor (e.g. implementation of the classical kNN based on Euclidean distance). Biblio: kNN folder
6. Naïve Bayes classifier (e.g. implementation of a classifier for data with discrete attributes). Biblio: NaiveBayes folder
7. Multilayer perceptron and backpropagation (e.g. implementation of a one hidden-layer network trained with standard backpropagation tested for XOR). Biblio: MLP+BP folder
8. KMeans (e.g. implementation of the Lloyd variant of the algorithm). Biblio: kMeans
9. Fuzzy c-means (e.g. implementation of the standard version proposed by Bezdek). Biblio: FuzzyCMeans folder
10. Hierarchical agglomerative clustering algorithm (e.g. implementation of single-linkage variant). Biblio: HierarchicalAlgorithms folder
11. DBSCAN (e.g. implementation of a variant of the DBSCAN algorithm). Biblio: DBSCAN folder
12. Apriori algorithm (e.g. implementation of a simple variant of Apriori algorithm). Biblio: Apriori folder

## B. Projects oriented toward data

- datasets from UCI Machine Learning Repository)
- datasets from <https://www.kaggle.com>

Projects of type B consist of:

- A report describing the dataset, the problem to be solved and the used method (mainly based on the papers referred in the dataset description from UCI Machine Learning Repository)
- Description of the processing workflow (the processing steps applied to the dataset), the parameter values which have been used and the results obtained by applying a data mining tool (at your choice – it could be Weka, an R library, a Python library or another platform) to the dataset

### Topics for projects of type B:

13. DBWorld e-mails data set (<http://archive.ics.uci.edu/ml/datasets/DBWorld+e-mails>). **Aim:** classify the e-mails in two categories: conference announcements vs other messages (binary classification task)
14. Microblog PCU data set (<http://archive.ics.uci.edu/ml/datasets/microblogPCU>). **Aim:** identify spammers (binary classification task)
15. SMS Spam Collection (<http://archive.ics.uci.edu/ml/datasets/SMS+Spam+Collection>). **Aim:** classification of SMS messages in spam/ham (binary classification task)
16. Energy efficiency data set (<http://archive.ics.uci.edu/ml/datasets/Energy+efficiency>). **Aim:** predict heating and cooling load in a building (based on a set of other characteristics)
17. GPS trajectories (<http://archive.ics.uci.edu/ml/datasets/GPS+Trajectories>). **Aim:** identify clusters of similar trajectories (clustering)
18. Blog feedback dataset (<http://archive.ics.uci.edu/ml/datasets/BlogFeedback>). **Aim:** prediction of the number of comments in the following 24 h) (regression)
19. Online news popularity (<http://archive.ics.uci.edu/ml/datasets/Online+News+Popularity>). **Aim:** prediction of the number of shares of the news (regression)
20. Student performance dataset (<http://archive.ics.uci.edu/ml/datasets/Student+Performance>). **Aim:** prediction of one of the grades (math, Portuguese or final)
21. AAAI2013 Accepted Papers Dataset (<http://archive.ics.uci.edu/ml/datasets/AAAI+2013+Accepted+Papers>). **Aim:** clustering based on keywords
22. News aggregator dataset (<http://archive.ics.uci.edu/ml/datasets/News+Aggregator>). **Aim:** group news by category
23. House price prediction (<https://www.kaggle.com/c/house-prices-advanced-regression-techniques>). **Aim:** house price estimation starting from some characteristics (regression)
24. Credit card fraud detection (<https://www.kaggle.com/dalpozz/creditcardfraud>). **Aim:** fraud prediction based on user actions (classification)

25. Gender recognition by voice (<https://www.kaggle.com/primaryobjects/voicegender> ). **Aim:** prediction of the gender of a person based on the voice characteristics (classification)
26. Paper clustering (<https://www.kaggle.com/benhamner/nips-2015-papers> ). **Aim:** grouping papers submitted to a conference based on the similarity between their content (clustering)

C. **Projects oriented to software tools**

- **Scikit-learn (Machine Learning in Python)** - <http://scikit-learn.org/stable/index.html#>
- **RATTLE (R Analytical Tool To Learn Easily)** - <http://rattle.togaware.com/>

Projects of type C consist of:

- A report describing the software tool (implemented methods, usage characteristics, facilities)
  - Details on using the software tool for a specific problem and the results obtained by applying it (to the dataset in the example and to another dataset)
27. Classification of text documents using sparse features - [http://scikit-learn.org/stable/auto\\_examples/text/document\\_classification\\_20newsgroups.html#example-text-document-classification-20newsgroups-py](http://scikit-learn.org/stable/auto_examples/text/document_classification_20newsgroups.html#example-text-document-classification-20newsgroups-py)
  28. Clustering text documents using k-means - [http://scikit-learn.org/stable/auto\\_examples/text/document\\_clustering.html](http://scikit-learn.org/stable/auto_examples/text/document_clustering.html)
  29. Face completion with a multi-output estimators - [http://scikit-learn.org/stable/auto\\_examples/plot\\_multioutput\\_face\\_completion.html#example-plot-multioutput-face-completion-py](http://scikit-learn.org/stable/auto_examples/plot_multioutput_face_completion.html#example-plot-multioutput-face-completion-py)
  30. Faces recognition example using eigenfaces and SVMs - [http://scikit-learn.org/stable/auto\\_examples/applications/face\\_recognition.html#example-applications-face-recognition-py](http://scikit-learn.org/stable/auto_examples/applications/face_recognition.html#example-applications-face-recognition-py)
  31. Recognizing hand-written digits - [http://scikit-learn.org/stable/auto\\_examples/classification/plot\\_digits\\_classification.html#example-classification-plot-digits-classification-py](http://scikit-learn.org/stable/auto_examples/classification/plot_digits_classification.html#example-classification-plot-digits-classification-py)
  32. Color Quantization using K-Means - [http://scikit-learn.org/stable/auto\\_examples/cluster/plot\\_color\\_quantization.html#example-cluster-plot-color-quantization-py](http://scikit-learn.org/stable/auto_examples/cluster/plot_color_quantization.html#example-cluster-plot-color-quantization-py)
  33. Vector Quantization Example - [http://scikit-learn.org/stable/auto\\_examples/cluster/plot\\_face\\_compress.html#example-cluster-plot-face-compress-py](http://scikit-learn.org/stable/auto_examples/cluster/plot_face_compress.html#example-cluster-plot-face-compress-py)
  34. Constructing decision trees using R - <http://rattle.togaware.com/rattle-examples.html>

**Remark:** proposals of other problems or datasets which can be solved by using data mining methods are also accepted