

Group 3: Machine Learning Basics Kaggle Contest: Home Credit Default Risk (Ho Chuen Ho, Chun Lok Him Brian)

a. Summary of the report

Home Credit Default Risk is an issue where an institution needs to predict whether a client is likely to default. In Group 3's project, datasets from Kaggle are used and various machine learning algorithms are used on the dataset. For simplicity, they used mostly the variables listed in `application_train.csv`.

Group 3 adopted four models for prediction, which include Naïve Bayes, Adaboost, Stochastic Gradient Descent and Light Gradient Boosting. Among them, Adaboost and Light Gradient Boosting achieves the highest accuracy of around 73 - 74%. Nonetheless, since Adaboost includes much more features, it has the longest running time. Group 3 therefore came up with a conclusion that Light Gradient Boosting is the best model for this project, which is the same as the results our group arrived at.

b. Describe the strengths of the report

Strong and supportive data visualization (bar chart of importance, heatmap of variables correlation) for the data analysis / processing section. Variety of models are used for example Naive Bayes, Adaboost, Stochastic Gradient Descent and Light Gradient Boosting. Analytic skills are shown in model comparison taking the variable importance and correlation as support.

c. Describe the weaknesses of the report

Although the data visualization is supportive, they should not be exclusive for data analysis and model comparison. It would be a good choice if group 3 show the imbalance situation of dataset with a suitable chart

- d. Evaluation on quality of writing (1-5): 5, the flow of the presentation, from data analysis, model selection to future improvement, was well delivered and clearly explained to us. The format of the presentation using Prezi was good as well.
- e. Evaluation on presentation (1-5): 5, the presentation looked refreshing using Prezi. Also, the group illustrated the algorithms and listed the respective advantages and disadvantages for each model clearly. Overall the presentation was well delivered.

- f. Evaluation on creativity (1-5): 5, traditional machine learning models (e.g. Does the work propose any genuinely new ideas? Is this a work that you are eager to read and cite? Does it contain some state-of-the-art results? As a reviewer you should try to assess whether the ideas are truly new and creative. Novel combinations, adaptations or extensions of existing ideas are also valuable.

- g. Confidence on your assessment (1-3): 3