
CSC 522 - Comparison of Encoding and Classification Methodologies for Prediction of Grant Applications

Priyaranjan Behera
Department of Computer Science
North Carolina State University
pbehera@ncsu.edu

Sai Sri Harsha Kunapareddy
Department of Computer Science
North Carolina State University
skunapa@ncsu.edu

1 Background

Around the world, the pool of funds available for research grants is steadily shrinking (in a relative sense). In Australia, success rates have fallen to 20-25 per cent, meaning that most academics are spending valuable time making applications that end up being rejected. This problem was hosted as a competition by University of Melbourne to address their inefficiencies with the grant applications. There is also a hope of discovering the most important criteria that are required to succeed in a grant application.

1.1 Problem

The university has provided a dataset containing 249 features, including variables that represent:

- Type and size of the grant: It includes the basic information of the Grant Application like Application Month, Year, Contract Value, Sponsors, etc.
- General area of study in the grant application: It includes the information on the study that is intended to be done with the grant like research fields, courses and disciplines class and socio economic objective class.
- De-identified Information of the Investigators: It includes personal data of the investigators like number of publications in different levels of journals, past history in grant applications, experience, educational qualifications, etc.

The dataset contains multi-variate data with a few of them being continuous variables and a few categorical. There is a variable number of investigators in an application and thus, we need to aggregate the person data to create an efficient model. We implemented several of the classification techniques covered in the CSC522 course to arrive upon an efficient model.

1.2 Literature Survey

As the dataset contains a variable number of person attributes depending on the number of investigators in a grant application, it would lead to multiple missing data fields in applications where the investigator count is less. According to Gerhard Svolba [2], we need to create a one-row-per-subject data mart for most of the analytical methods that we need to proceed with. Accordingly, we need to aggregate the person data using mean, median, standard deviation, the quartiles, or special quantiles, etc to create the input rows.

We looked at approaches taken to solve similar problems which contains both continuous and categorical data [6] and found that decision trees and Naive Bayes are popular choices. On the contrary, methods like neural networks, SVM, etc. cannot work with categorical data. We implement decision trees as well as bagged and boosted trees for more efficiency.

Daniele et al. [4] focused on the transformation of categorical data to continuous/binary data so that we can use the analytical methods which cannot process categorical data. While traditionally binary encoding is used for conversion of categorical data to binary, for categories with high cardinality a probabilistic approach is suggested.

2 Methods

The analysis of the data required extensive preprocessing steps because of its multi-variate nature and variable number of attributes. We implemented different type of classification techniques to compare the models and find the determining attributes for the classification.

2.1 Preprocessing

To handle the categorical values in the data, we indexed the attributes with numerical values after finding the unique values for each of the attributes. As per the approach suggested by Gerhard Svolba [2], we aggregated the variable number of investigators specific data for each of the application which is not null and found the minimum, maximum for all the attributes. While for continuous attributes we also calculated mean, median, sum of the values to create new features. We then created plots of the attributes with respect to the output to visually inspect any correlation of the attributes with the results.

Since classification techniques like SVMs, Neural Networks, cannot handle categorical data, we implemented binary encoding to convert the categorical data into binary attributes which can be used in the analytical techniques. Further, to handle high cardinality, we implemented a probabilistic feature creation as specified at [4].

2.2 Classification

We implemented most of the classification techniques covered in CSC522 and compared their efficiencies.

2.2.1 Decision Tree and Random Forest

According to our survey of approaches employed for similar problems, decision tree turns out to be a popular choice. This attributes to the fact the it can handle categorical values without any further processing. We also implemented bagged trees using the dataset to get a more efficient classification. We also used the error rate obtained by changing values of attributes in a random forest to determine the most important factors which decide the success of a grant application.

2.2.2 Neural Networks

Since the neural networks can handle a high number of attributes, this will serve as an ideal classifier particularly when we use the binary encoding with a high cardinality of the categorical attributes. However, we cannot estimate the determining factors of an application though this method.

2.2.3 SVMs

SVM based classifiers are known to work well with encoded data and thus, we executed the SVM based classifiers on the binary and probabilistic encoded data. We used the linear, quadratic and cubic kernels to find an effective model.

3 Plan

According to the problem statement of the competition, we aimed to find a model which has a greater accuracy while minimizing the false positives in the output.

3.1 Hypothesis

We worked on implementing models of the classifications techniques we have been familiarized with to find a model with lesser number of false positives and higher accuracy. Thus, we focused on the values of True Positive Rate, False Positive Rate, Precision and Accuracy of the models. Most importantly, we used the ROC (Receiver Operating Curve) to determine the effectiveness of the model by comparing the area under the curve for the models. This will be more helpful compared to other criteria as:

- It is insensitive to unbalanced dataset where there is a disparity in the number of test cases outputs.
- We consider all the cut-offs in a model to find its effectiveness.

We worked on experiments to test the technique of conversion of categorical attributes to probabilistic values for high cardinality attributes[4]. We considered the same evaluation parameters mentioned above to compare with models which use binary encoding.

We also worked on finding the most determining factors which decides the fate of an application. For this, we used the characteristic of the random forest libraries which find the effective factors in an input by measuring the error rate when an attribute value is changed.

3.2 Experimental Design

1. *Data Indexing*: Since there are a large number of categorical attributes in the data, indexing was done to standardize them into integral values.
2. *Dimensionality Reduction and Data Creation*: Since the data has variable number of person data in each of the grant application, there was a need to aggregate the data to create similar rows for comparing applications. We used different aggregation functions like minimum, maximum, mean, median, sum of the attribute values for each of the person within a data row[2] to achieve this.
3. *Data Visualization*: Effectiveness of the attributes to determine the result of an application was studied by creating scatter plots and histograms for each of the attributes in the input matrix.
4. *Handling of Categorical Attributes*: While techniques involving decision trees and naive bayes work well with categorical values, other methods need binary or continuous inputs. Thus, we will be adopting two methods to use the data in other classifiers:
 - Binary encoding: In this technique, an indexed categorical variable with N cardinality is converted to its binarized index value.
 - Probability Substitution: In this technique, the probability of getting the output as 1 is calculated for each of the values of the categorical variable. These values are used instead of the original variable values for the classifications. [4]
5. *Classifications*: We used the different inputs obtained in the above steps to train and test the classifiers.
 - Decision Tree
 - Random Forest
 - Boosted Trees
 - Neural Networks
 - SVMs

Primarily, libraries in Matlab will be used for this purpose where tools like 'classificationLearner', 'nntool' will be used in addition to other basic classifier classes.

6. *Factor Determination*: The random forest library explicitly determines the most important factor which decides the output. We will be using this to find out the factors which can improve the credentials of an application.
7. *Evaluation and Results*: Confusion matrices will be drawn for each of the classifications done in the experiments to find the accuracy, recall, precision and false positive rate. Additionally, ROC will be generated from the GUI interfaces in 'classificationLearner' and 'nntool'.

4 Experiments

After the pre-processing steps which included indexing and generation of new features, we resorted to visualization of the data to find any correlation of attributes with the results (As shown in Fig.1).

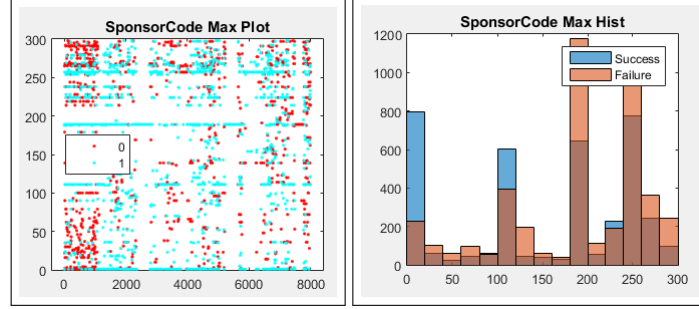


Figure 1: Scatter Plot and Histogram for Sponsor Codes.

We take an example of Sponsor code here, for which we plot the data and see a discerning pattern. We can see that for specific sponsor codes, the chances of getting a grant is more.

But we could not find any conclusive pattern in any of the plots and thus, moved on with the classification. Since decision trees doesn't have much problem with the number of dimensions we moved ahead with building decision tree and random forest classifiers. For the decision tree classifier we opted for a 5-fold cross validation. The confusion matrix along with TPR/FPR rate for the random forest is listed at Fig.2. The accuracy obtained for the decision tree is 82.6% while for the random forest is 85.9%.

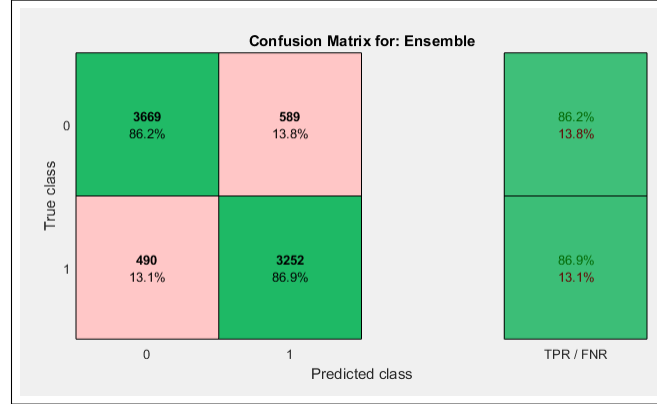


Figure 2: Confusion Matrix and TPR/FPR of Random Forest.

By comparing the models with respect to the ROC curve(Fig.3), we find that the area under the curve for random forest is a little high and thus, random forest is a better model in this scenario.

In addition, we found the important features which determine the outcome of the application. The data is plotted in Fig.4. From the plot we can see that features Contract Value(3), Sponsor Code(1), Application Month(14), Grant Category(2) are the most important factors.

5 Conclusion

The experiments using the categorical values without encoding is done for decision tree, boosted trees and random forest, where we found boosted tree to be a little more accurate. Further, we will work on encoding the categorical attributes and create models based on neural networks, SVM and naive bayes classifiers.

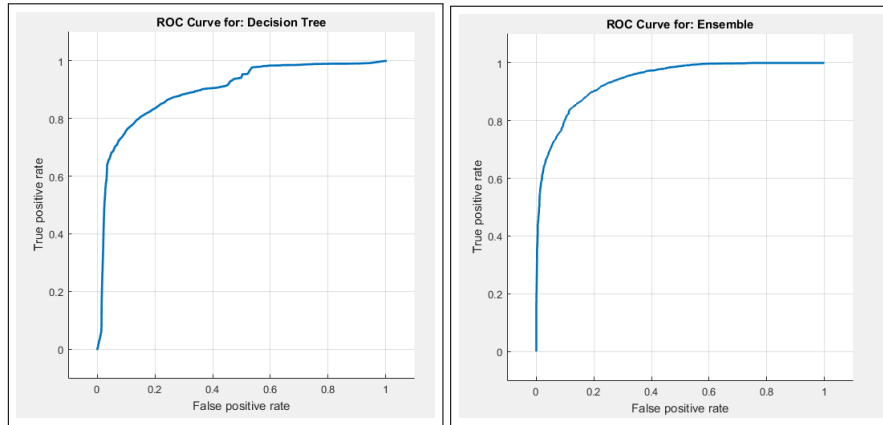


Figure 3: ROC Plot for Decision Tree and Random Forest.

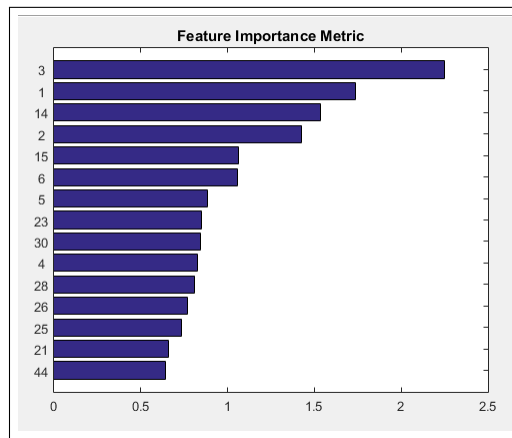


Figure 4: Feature Importance Matrix Generated By Random Forest.

References

- [1] Top 10 algorithms in data mining. Knowl. Inf. Syst. 14, 1 (December 2007), 1-37. DOI=<http://dx.doi.org/10.1007/s10115-007-0114-2>
- [2] Efficient One-Row-per-Subject Data Mart Construction for Data Mining, Gerhard Svolba, PhD, SAS Austria
- [3] Reliable Early Classification on Multivariate Time Series with Numerical and Categorical Attributes, Cao, Tru et al.
- [4] Daniele Micci-Barreca. 2001. A preprocessing scheme for high-cardinality categorical attributes in classification and prediction problems. SIGKDD Explor. Newsl. 3, 1 (July 2001), 27-32. DOI=<http://dx.doi.org/10.1145/507533>.
- [5] Getting Started with Kaggle Data Science Competitions. Loren Shure
- [6] Predict Grant Applications - Kaggle Competition, <https://www.kaggle.com/c/unimelb>

A Appendix

A.1 Deviation from Initial Proposal

In the proposal, we didn't focus more on the pre-processing steps, particularly encoding of the categorical attributes. After further work on the project, this became a necessity to implement the different kinds of classifiers. Thus, we are looking into implementing binary encoding and a probabilistic encoding after the mid-term report submission.

We also found the way to extract the importance of features for the classification with the help of decision trees and random forest. Thus, our results will also include the features which will decide the classifications most.

A.2 Division of Work

Priyaranjan and Harsha worked in the pre-processing steps and implementation of decision trees and random forest collaboratively. Going forward, Priyaranjan will work on the probabilistic encoding of categorical attributes and the classifications. While, Harsha will work on the binary encoding and the related classifications.