



# KDD DATA PROJECT

## Predicting 'Exciting' DonorsChoose.org projects

### Abstract

Data from the KDD Cup 2014 Challenge is used to correlate “exciting” DonorsChoose.org projects to US states and school poverty levels. A score indicating how exciting a project is relative to the poverty level of a school is also calculated.

Ashish Naik

ashishsnaik@yahoo.com

## Table of Contents

Summary .....	2
Introduction .....	2
Steps Used to Arrive at the Results.....	2
Data Cleaning and Transformations .....	3
Source Files .....	3
Handling Missing Values .....	3
Feature extraction and selection .....	3
Source Files .....	3
Essay data.....	3
Resource Data .....	3
Project Data .....	3
School and Teacher Credibility (Project, Donation, and Outcome Data) .....	4
Modeling Process.....	4
Source Files .....	4
Handling Class Imbalance .....	4
GBM Model .....	4
Test Data Predictions .....	6
Source Files .....	6
Correlating and Ranking Exciting Projects .....	6
Source Files .....	6
Performance on Kaggle Leaderboard .....	7
Additional Comments .....	8
Project Environment .....	8
Executing the code.....	8
Executing R code using R Console.....	9
References .....	9
Appendix .....	10
List of Files.....	10
Figure 1: GBM Parameter Tuning .....	5
Figure 2: GBM ROC Curve - Test-set .....	6
Figure 3: Kaggle public and private scores.....	7
Figure 4: Kaggle leaderboard ranking .....	7

# KDD Data Project

## Summary

A total of 239 school projects are identified as among the Top 5 exciting projects in each of 51 US States, across 3 school poverty levels - highest, high, and moderate.

The correlation and ranking of the "exciting" projects is done based on the average number of underprivileged students that could benefit from the project. A project score was calculated as the ratio of 'number of underprivileged students served by an exciting project' to 'total number of underprivileged students served by all exciting projects in the state'. A project with a higher score, got a higher rank.

The ranking is done such, because it focuses on the (possible) number of underprivileged students reached, rather than merely considering the poverty levels of schools or the total number of students reached.

As an example, in the State of Iowa (IA), with only two top ranking projects, a project from a school with high poverty level is ranked higher than one from a school with highest poverty level. Similar can be observed for the State of New Hampshire (NH) where two projects from school(s) with moderate poverty level are ranked higher than one from a school with highest poverty level.

To predict whether a project is exciting or otherwise, an ensemble of 15 Generalized Boosted Regression Models (GBM) was used.

Please see the attached KDDTopProjects.csv file to view the correlation and ranking results.

## Introduction

This data project uses data from the KDD Cup 2014 Challenge and attempts to correlate the predicted "exciting" DonorsChoose.org projects, in the test set, to US states and school poverty levels identifying the top 2-5 projects per state that would benefit poorer schools. It also attempts to compute a score indicating how exciting a project is relative to the poverty level of a school.

The KDD data sets contain information about: the projects itself (projects.csv), donations received by projects (donations.csv), project text posted by the teachers (essays.csv), the resources requested for a project (resources.csv), and outcomes of the projects - whether exciting or otherwise - (outcomes.csv). Projects posted on or after 2014-01-01 are in the test set and the data sets do not contain information about the donations to or outcomes of these projects.

## Steps Used to Arrive at the Results

No exciting projects were reported before 2010-04-14. This could be because DonorsChoose.org tracked the projects differently or the criteria for exciting projects changed around that time. So, projects posted before April 2010 are not considered for analysis.

Projects posted in 2010 (April and later), 2011, and 2012 are used for model training and those posted in 2013 are used for model evaluation (validation and testing).

## Data Cleaning and Transformations

### Source Files

- kdd\_clean\_data.py
- kddData.R
- kddGBM.R (GBM specific transformations)

Performed cleaning of the 'text' features from donations.csv, essays.csv, and resources.csv using a python script (kdd\_clean\_data.py). This ensured that all records were read by the R code. Did a find-replace to remove '\r\n' character-tuple from the text in resources.csv, before running python script.

### Handling Missing Values

- Boolean values were replaced by 'f' (or numeric zero for outcome data).
- Numeric values were replaced by column median value (or numeric zero for outcome data).
- Non-numeric values were replaced by '#NA#' value.

The missing values for School NCESID, Secondary Focus Subject, and Secondary Focus Area were given special treatment.

- School NCESID - missing values replaced by the numeric factor-level of the corresponding *schoolid*, in order to keep the ncesid same for a school and different across schools.
- Secondary Focus Subject - missing values replaced by Primary Focus Subject, as available.
- Secondary Focus Area - missing values replaced by Primary Focus Area, as available.

Specifically for GBM Model:

- Boolean 't' (true) and 'f' (false) values were converted to numeric 1 and 0 respectively.
- Non-numeric features were replaced by their numeric factor-levels.

## Feature extraction and selection

### Source Files

- kddData.R
- kddGBM.R

The following features were extracted from the data.

### Essay data

Word counts of the *title*, *short\_description*, *need\_statement*, and *essay* were calculated. Missing values received a word count of zero.

### Resource Data

For each project the number of resources requested, number of vendors used, total item quantities, total cost of resources, and average cost per item were calculated.

### Project Data

The number of underprivileged students that could possibly benefit from the project and the month and weekday the project was posted.

Assumption: Underprivileged students in a school are from the average poverty-level percentage within the school's *poverty\_level* range. The students benefited (*students\_reached*) from a project contain that much percent of students that are underprivileged.

### School and Teacher Credibility (Project, Donation, and Outcome Data)

The 'exciting projects' criteria can be influenced (at least in part) by teachers' capabilities, for example, to gather at least one teacher-acquired donor. The criteria can also be influenced by schools' abilities to create an environment that promotes externally funded projects, attract project funding, and build a reputation. To accommodate for this, the following credibility features for both, schools and teachers, were extracted:

- number of projects posted
- number of projects that were fully funded
- number of projects that were exciting
- number of projects that had a great chat
- number of projects that received donations from a thoughtful donor
- average number of donations received
- average number of green donations received
- average number of teacher referred donors
- average number of donors that donated \$100 plus
- average number of non-teacher referred donors that donated \$100 plus

Features that showed near-zero variance in the training data were not considered for analysis.

In order to ease the modeling process, the target variable *is\_exciting* was tagged with 'y\_' and the features were tagged with 'x\_'.

### Modeling Process

#### Source Files

- kddGBM.R

#### Handling Class Imbalance

The training data showed about 12:1 class imbalance between the non-exciting and exciting projects. To handle this, an ensemble learning approach with EasyEnsemble Informed Under-Sampling [1] was used.

#### GBM Model

##### Model Training and Parameter Selection

Area under the ROC Curve (AUROC) was used as the performance metric.

A 3-way data split was used for model training, parameter selection, and performance estimation as mentioned below.

1. The data training data were split into training, validation, and test sets.
  - Training Set: Years 2010 (April and later), 2011, and 2012 data.
  - Validation and Test Sets: Year 2013 data. The data was randomly sampled into two equal and class-balanced sets for validation (parameter tuning) and testing (performance estimation).
  - Training set: 289360 samples
  - Validation set: 65665 balanced samples
  - Test set: 65664 balanced samples
2. 5-fold cross validation was used to tune parameters.

- Training set: 5 resamples of balanced training data, each with 22432 minor + 22432 random major samples.
- Validation set: 5 resamples of 60% of the validation set (39399 random samples each).
- Tuning parameters grid:
  - interaction.depth: 7, 8, 9
  - n.trees: 2000 to 4000 , with 500 step size
  - shrinkage: 0.1, 0.01
  - n.minobsinnode: 10, 50

The GBMParameterTuning.txt file in the submission shows the parameter tuning output. Below *Figure 1: GBM Parameter Tuning* shows the parameter tuning plot.

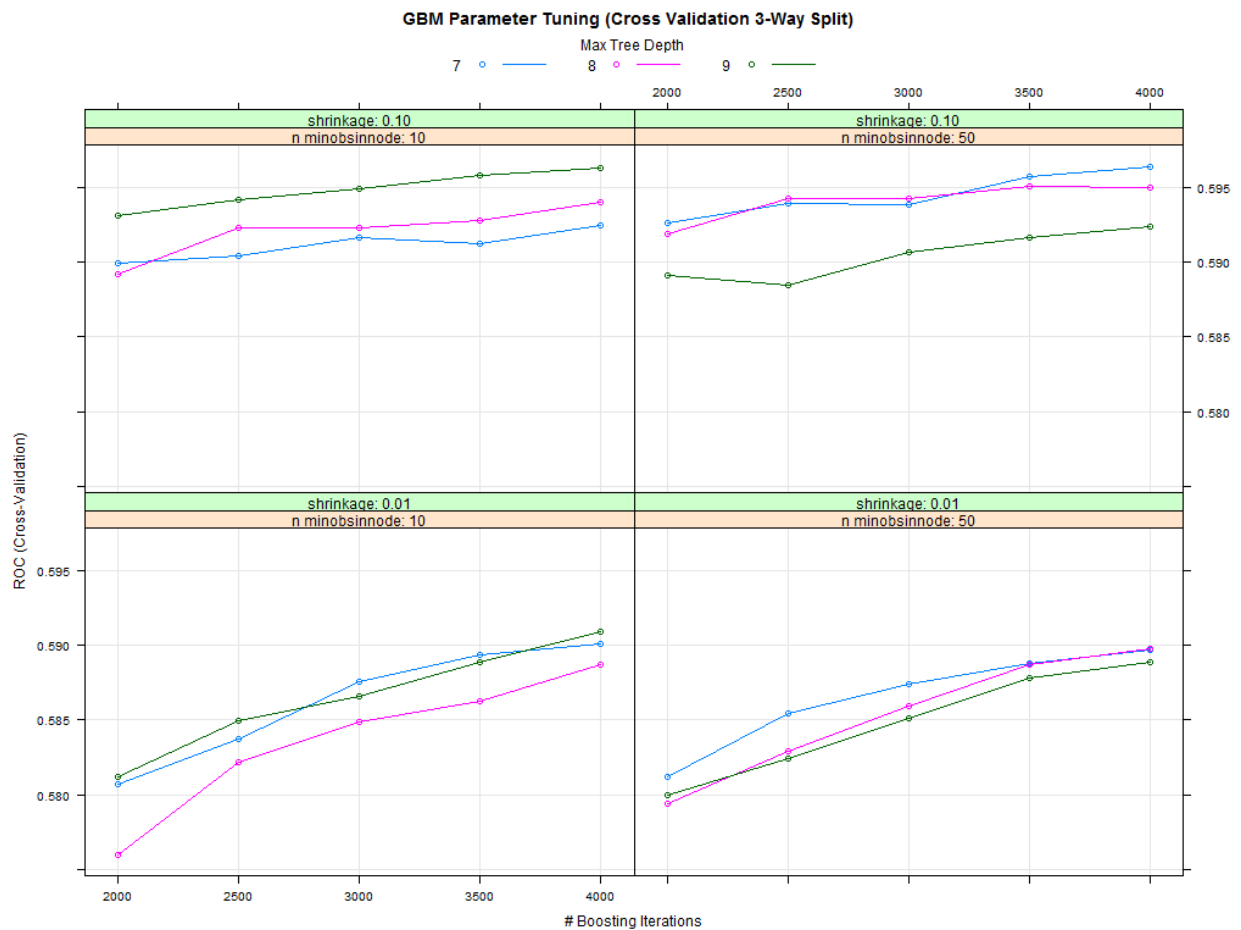
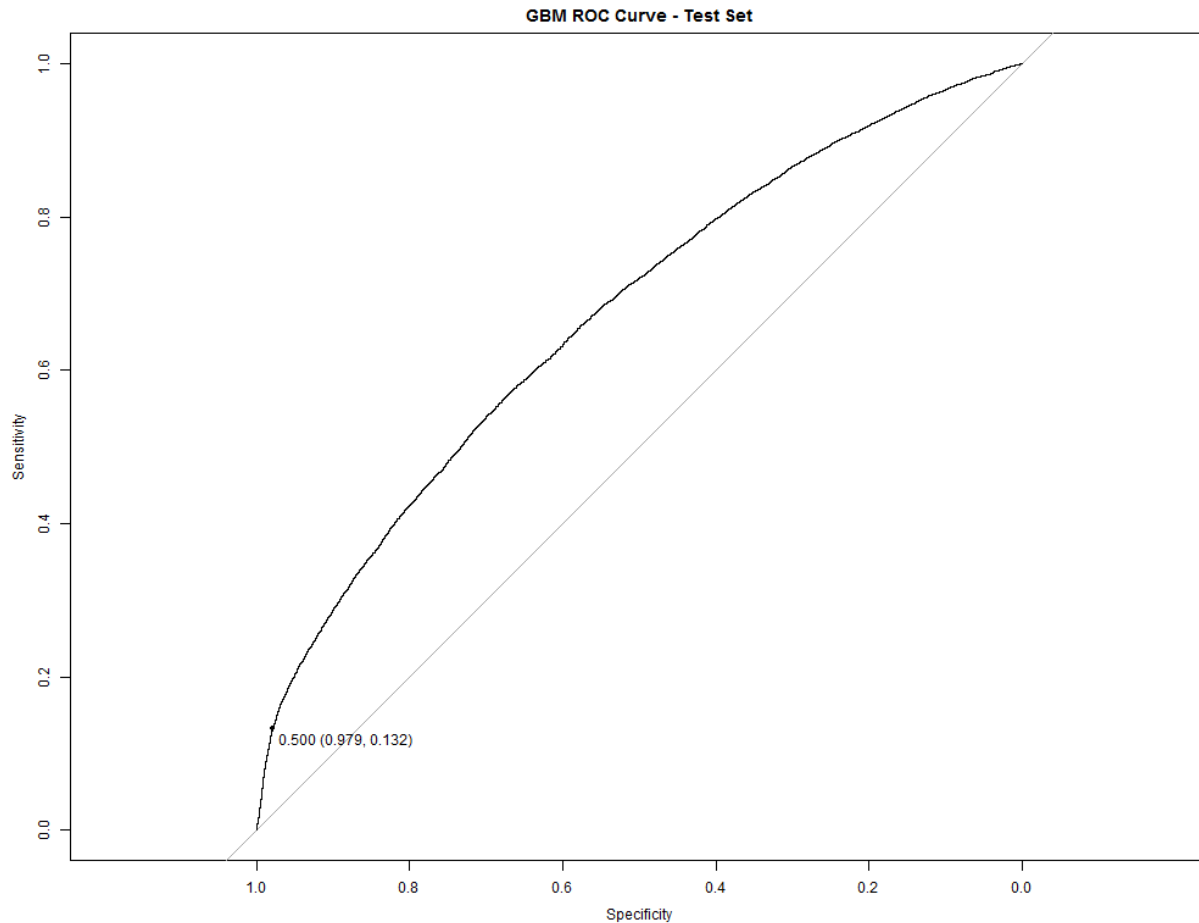


Figure 1: GBM Parameter Tuning

- Performance of the final model was accessed on the test set using the below tuned parameters.
  - Final model parameters:
    - interaction.depth: 7
    - n.trees: 4000
    - shrinkage: 0.1
    - n.minobsinnode: 50

The Test Set AUROC was 0.67.

*Figure 2: GBM ROC Curve - Test-set* below shows the ROC curve.



*Figure 2: GBM ROC Curve - Test-set*

## Test Data Predictions

### Source Files

- kddGBM.R
- kddMySQL.R

Predictions on the test data were made using an ensemble of 15 GBM models trained on balanced training data samples with the tuned parameters. The kddMySQL.R code was used to create a MySQL 'kddcupdb' database and populate the 'predictions' table.

The gbmpredictions.csv.gz file in the submission contains the test set prediction results.

## Correlating and Ranking Exciting Projects

### Source Files

- kddmysql.sql

The above MySQL script was used to create a 'view' of the exciting projects from the 'kddcupdb' .predictions' table and generate the correlation and ranking.

## Performance on Kaggle Leaderboard

'White Noise' is my Kaggle.com screen name.

**KDD Cup 2014 - Predicting Excitement at DonorsChoose.org**  
 Completed • \$2,000 • 472 teams  
 Thu 15 May 2014 – Tue 15 Jul 2014 (10 months ago)

**Your Submissions**  
 You are submitting as part of team [White Noise](#). [Make a submission »](#)

The competition deadline has already passed and you can no longer modify selections. While this competition was active, you could select up to 2 submissions. This information is provided for historical purposes only.

Submission	Files	Public Score	Private Score	Selected?
<b>Post-Deadline:</b> Thu, 21 May 2015 18:03:07 Submission 4 <a href="#">Edit description</a>	<a href="#">gbmPredictions_4.csv.gz</a>	0.56990	0.57114	<input type="checkbox"/>

Figure 3: Kaggle public and private scores

215	↓102	gw0	0.57127	34	Tue, 15 Jul 2014 22:15:31 (-15h)
216	↓61	Zstats	0.57127	7	Tue, 20 May 2014 05:34:10 (-22.9h)
217	↑114	thomaschen639	0.57125	4	Sun, 15 Jun 2014 08:23:17 (-0.3h)
-		<b>White Noise</b>	<b>0.57114</b>	-	<b>Thu, 21 May 2015 18:03:07</b> Post-Deadline
<b>Post-Deadline Entry</b> If you would have submitted this entry during the competition, you would have been around here on the leaderboard.					
218	—	Andrew Carter	0.57100	19	Sun, 13 Jul 2014 02:11:21 (-0.2h)
219	↓33	Dan Mysnyk	0.57084	53	Thu, 10 Jul 2014 19:50:15 (-13.1d)
220	↓15	carpcrap	0.57066	9	Thu, 19 Jun 2014 09:46:53 (-47.5h)
221	↑113	Allenbj	0.57049	3	Tue, 17 Jun 2014 23:03:58 (-26.1d)
222	↓142	C Choo	0.56968	32	Tue, 15 Jul 2014 17:38:30 (-5.7d)
223	↑89	ECNUCS	0.56961	74	Thu, 10 Jul 2014 13:45:04 (-9.1d)

Figure 4: Kaggle leaderboard ranking



## Additional Comments

Due to time constraints, the following, which would help increase the prediction accuracy, was not incorporated in the solution.

- NLP based analysis.
- Additional features and/or weights, especially to accommodate the fact that the test data is only from months January through May, while the training data is from January through December.
- Testing other types of models such as Random Forest and SVM and using higher number of trees.

## Project Environment

- R Version 3.1.3 (Packages: caret, gbm, pROC, qdap, plyr), RStudio Version 0.98.507
- Python 2.7.5 (Packages: pandas, shutil, sys, os, re, gc), PyCharm 2.7.3
- MySQL 5.7, MySQL Workbench 6.0 CE
- Microsoft Windows 7 Professional
- Intel(R) Core(TM) i7-2640M CPU @ 2.80GHz, 2801 MHz, 2 Core(s), SSD, 12GB RAM

## Executing the code

To execute the code, following requirements need to be satisfied:

- Directory structure as mentioned below.

<working\_directory>

- R, Python, and MySql Source Code Files (*mandatory*)
- /data/data\_files (*mandatory, should contain kdd input data files*)
- /data/input (*optional, created by python script to copy cleaned data files*)
- /data/rdata (*optional, created by kddCommon.R*)
- /data/rmodels (*optional, created by kddCommon.R*)
- /output (*optional, created by kddCommon.R*)
- Copy R, Python, MySQL source files to the <working\_directory>
- Change line 3 in kddCommon.R to set the working directory.
- In kddMySQL.R source file, change lines 25 and 30 (dbConnect) to provide the appropriate MySQL username, password, and host.
- The required Python and R packages mentioned in the Project Environment must be installed.

Sequence of execution:

Remove '\r\n' character-tuple from the text in resources.csv. Execute the source files in the following sequence:

- 1) kdd\_clean\_data.py
- 2) kddmain.R
- 3) kddmysql.sql

The kddMySQL.R code (executed through kddmain.R) will create a 'kddcupdb' database, if it does not exist, and also a create kddcupdb.predictions table. Approximate execution time with the aforementioned project environment is about 24 hours. The kddmysql.sql will show the correlation and ranking result.

**NOTE:** If kddcupdb.predictions already exists, IT WILL BE OVERWRITTEN.

### Executing R code using R Console

```
> setwd("<working directory>")  
> con <- file("kddmain_output.log")  
> sink(con, append=TRUE)  
> sink(con, append=TRUE, type="message")  
> source("kddmain.R")
```

## References

1. Xu-Ying Liu, Jianxin Wu, and Zhi-Hua Zhou, "Exploratory Undersampling for Class-Imbalance Learning," IEEE TRANSACTIONS ON SYSTEMS, MAN AND CYBERNETICS – PART B, 2008
2. Kaggle KDD Cup 2014 Competition Forum, <https://www.kaggle.com/c/kdd-cup-2014-predicting-excitement-at-donors-choose/forums>
3. <http://www.stackoverflow.com>
4. <http://www.stackexchange.com>

## Appendix

### List of Files

1. KDDDataProject\_AshishNaik.pdf – This document.
2. KDDTopProjects.csv – Correlation and ranking results.
3. gbmPredictions.csv.gz – Test set prediction results.
4. kdd\_clean\_data.py – Python code used to clean the KDD data.
5. kddCommon.R – R common definitions.
6. kddData.R – R code used to read KDD data, extract and transform features.
7. kddGBM.R – R code used for GBM modeling and predictions.
8. kddmain.R – Main R file used to run the R code.
9. kddMySQL.R – R code used to create the KDD database and predictions table.
10. kddmysql.sql – MySQL script used to generate exciting projects correlation and ranking.
11. GBMParameterTuning.txt – GBM parameter tuning output from R.