DESCRIBING & PREDICTING THE SUCCESS RATE OF

David Forteguerre, Ziwei Pi, Haoxuan Shi

1. PROJECT OVERVIEW AND OBJECTIVES

Online fundraising platforms are becoming more and more popular. However, the majority of projects turn out to be unsuccessful

Goals

- Use *descriptive statistics* to learn more about already completed projects
- Use *machine learning* to find the factors that play a key role in determining the outcome of a project (i.e. success or failure) and make predictions

2. DATA DESCRIPTION

Data source: Kickstarter Projects dataset, available on Kaggle

	name	category	main_categor	currency	country	status	launched_date	deadline	days_allotted	fundraising_goal	fundraising_pledged	completion_percentage	backers	
2	0 The Songs of	f Poetry	Publishing	GBP	GB	failed	8/11/15	10/9/15	59	1533	0	0		0
3	1 Greeting From	r Narrative Film	Film & Video	USD	US	failed	9/2/17	11/1/17	60	30000	2421	8	1	5
1	2 Where is Har	n Narrative Film	Film & Video	USD	US	failed	1/12/13	2/26/13	45	45000	220	0		3
5	3 ToshiCapital	F Music	Music	USD	US	failed	3/17/12	4/16/12	30	5000	1	0		1
3	4 Monarch Esp	Restaurants	Food	USD	US	successful	2/26/16	4/1/16	35	50000	52375	104	. 22	24

2,000.00 \$

787.00 \$

Data cleaning:

- Renamed and reordered columns
- Fixed data types
- Parsed variables with dates
- Only kept projects that failed or were successful (88% of original data)
- Removed country named 'N,0""

3. DATA EXPLORATION

fundraising_goal histogram

- Dropped redundant columns/incorrect data
- Removed rows that had a fundraising goal of \$0
- Created columns (day_allotted, completion_%)
- Checked for NAs

4. MODEL DESCRIPTION

Classification task with a binary outcome (success, failure)

Features used for machine learning

- 1st round of models: main_category, country, launch_month, launch_day, launch_weekday, days_alloted, fundraising_goal
- 2nd round of models: main_category, name, country, launch_month, launch_day, launch_weekday, days_alloted, fundraising_goal

Feature engineering

- main_category, country, launch_month, launch_day, launch_weekday → dummy variables
- days_allotted → untouched
- fundraising_goal → log transformation (common with positive right-skewed data w/ outliers)

(cleaned dataset overview)

after

330,244

10427700

days_allotted histogram

219382

before

377,364

CLEANING

columns

rows

- ✓ → lowercased, punctuation removed with regex, tokenized
- ✓
 →
 stopwords removed (list taken from tm package in R)
- ✓ → vectorized, converted to term frequency—inverse document frequency (tfidf)

Number of features without text: **56** Number of features with text: **56 + 92,680 tokens**

Models used in round #1 (without text features)

- **Logistic regression** (scaled data)
- $\alpha = 0$ (no elastic net regularization) \varnothing
- $\lambda = 0.02$ $\alpha = 0.2$
- $\lambda = 0.1$ $\alpha = 0.4$
- Random forest (unscaled data)
 - maxDepth=1 numTrees=60
 - maxDepth=6 numTrees=100
 - maxDepth=6 numTrees=80

Naïve Bayes

- smoothing=1.0 *on unscaled data ∜*
- smoothing=5.0 on unscaled data
- smoothing=1.0 on scaled data

Models used in round #2 (with text features)

- **Logistic regression**
 - With same parameters as found for the best model in round #1 for maximizing run time

Random forest

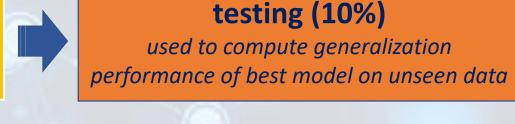
- Eventually dropped as a candidate algorithm for round #2 due to the large number of text features (we conducted many experiments – the server would always crash after a few hours)
- Naïve Bayes
 - With same parameters as found for the best model in round #1 for maximizing run time

5. MODEL COMPARISON METRICS

The data was split into **3 datasets**:







(dataset after feature engineering)

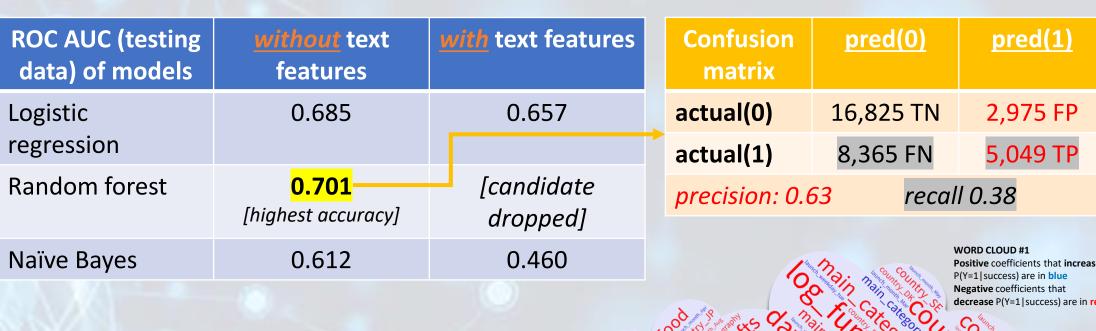
AUC ROC curve was used for evaluating the performance of models

6. RESULTS – PREDICTION PERFORMANCE

PROJECTS

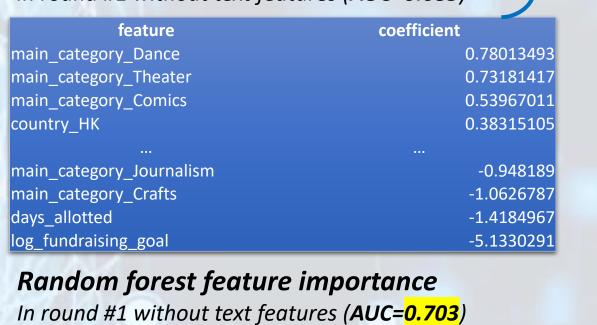


2. testing performance:





Logistic regression feature importance In round #1 without text features (AUC=0.685)



undraising goal 0.195647 0.0477046 0.0462787 0.0396993 0.0324568 0.0251637 _category_Comics 0.0227238 _category_Food



8. CONCLUSION

- Performance was quite high overall with random forest (approx. 70%)
- As a result, success rate may not only depend on the data and could also greatly vary based on factors that are unavailable to us/specific to each project
- Text features did not help improve performance with these models

FUTURE WORK (a data science project is never over! # • 1

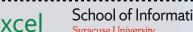
- Rerun logistic regression on text features using elastic net regularization despite the very high computational cost, and utilize/implement more algorithms which allow for more complex learning (e.g. neural networks)
- Attempt to predict fundraising_pledged and number of backers (i.e. numeric variables with a linear relationship instead of categorical ones/classification outcomes)
- Retrieve more data/projects to improve performance











117