

# KICKSTARTER

Mahmoud Belhaj  
Cen Liang  
Jonas Cristens



# The dataset

1. ID = Internal kickstarter ID
2. Name = Name of the project
3. Main\_category = Category of the campaign
4. Category = Sub category of the main\_category
5. Currency = Currency used to support the project
6. Deadline = Deadline for crowdsourcing
7. Goal = fundraising goal
8. Launched = When is the project launched
9. Pledge = How much money was gathered by the crowd
10. State = The state of the project (e.g. success, failure)



# The dataset

1. Backers = Number of people that backed the project
2. Country = The country of origin
3. Usd\_pledged = Pledged amount in USD done by Kickstarter
4. Usd\_pledged\_real = Pledged amount in USD using currency converter
5. Usd\_goal\_real = Goal amount in USD

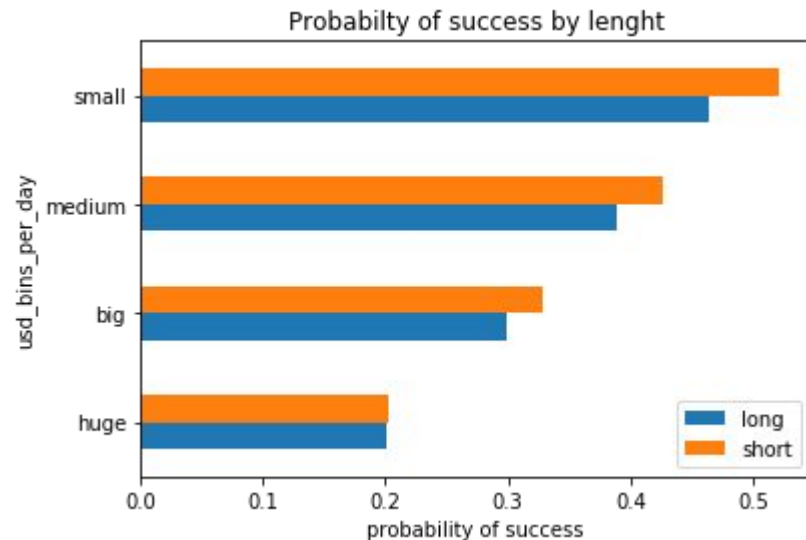
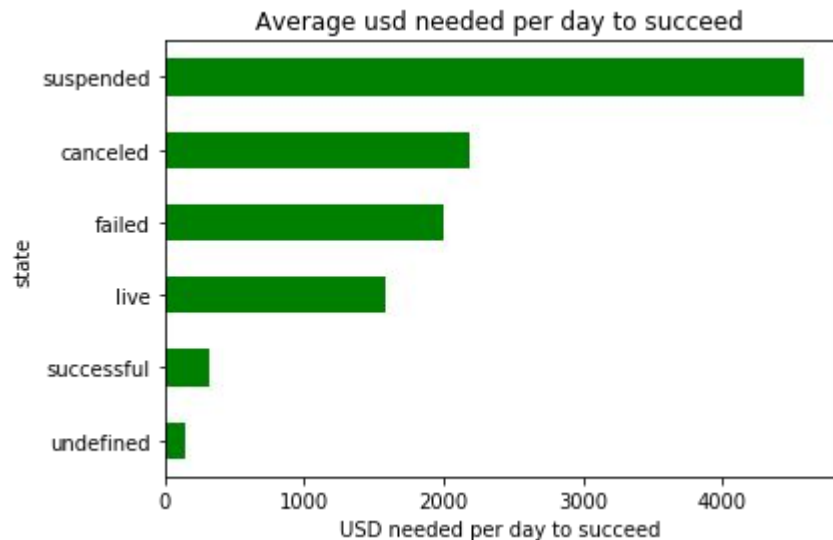


# Feature engineering

- Number of characters in the name of the project
- Number of words
- Number of syllables
- Launched week, month and weekday
- Days to collect the money
- Average amount needed per day to succeed
- Average amount donated per backer
- Difference between mean goal per category and the goal of the project



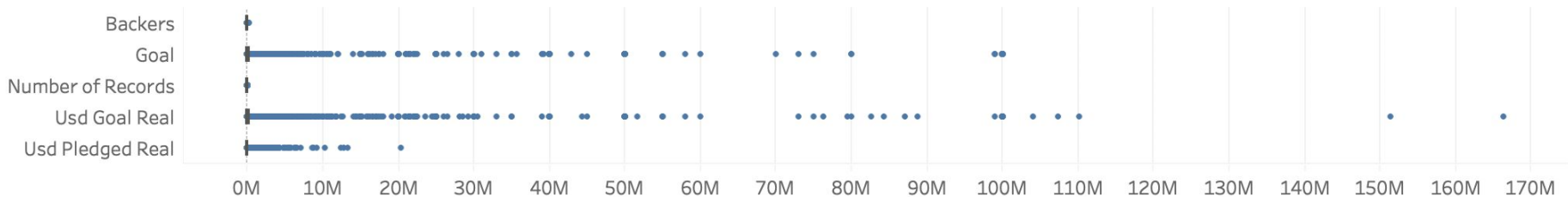
# Daily Goal



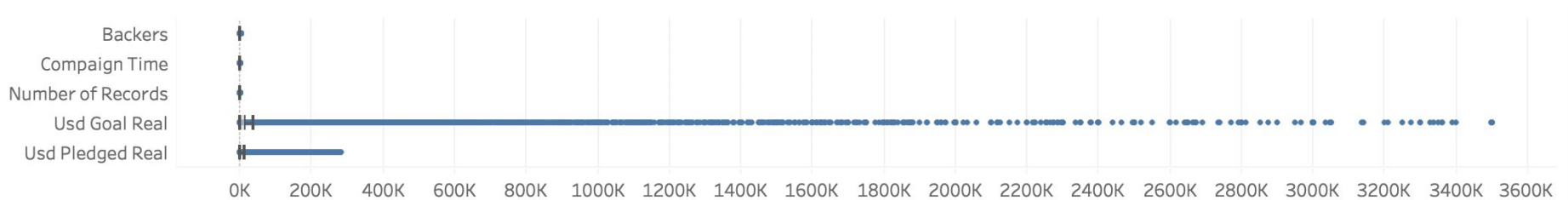


# Data cleaning

Before



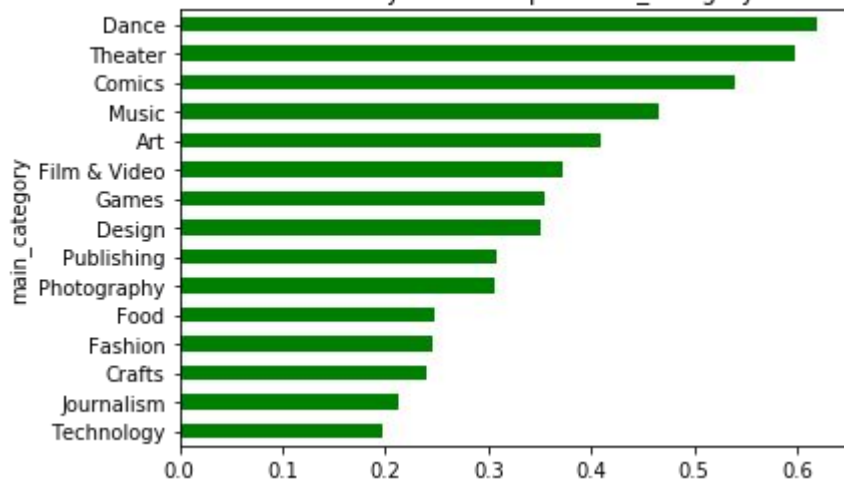
After



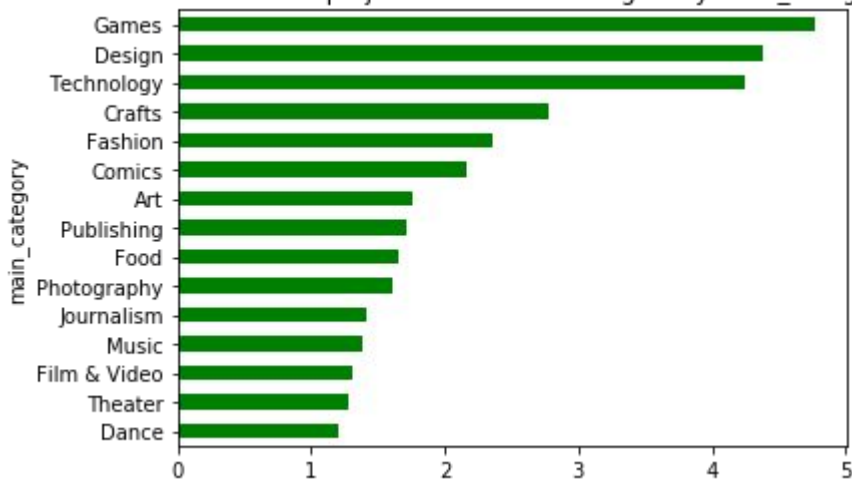
# Categories



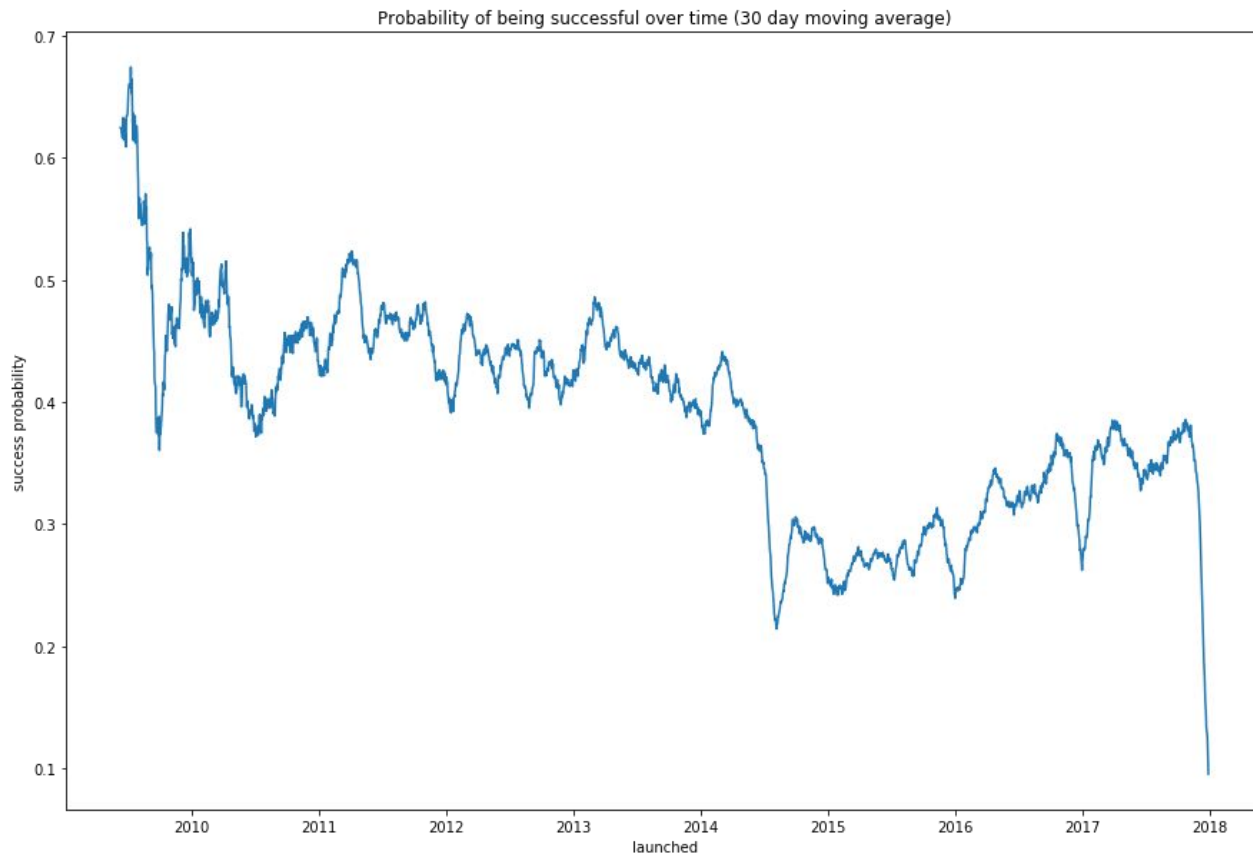
Probability of succes per main\_category



Most successful projects over 100 dollar goal by main\_category



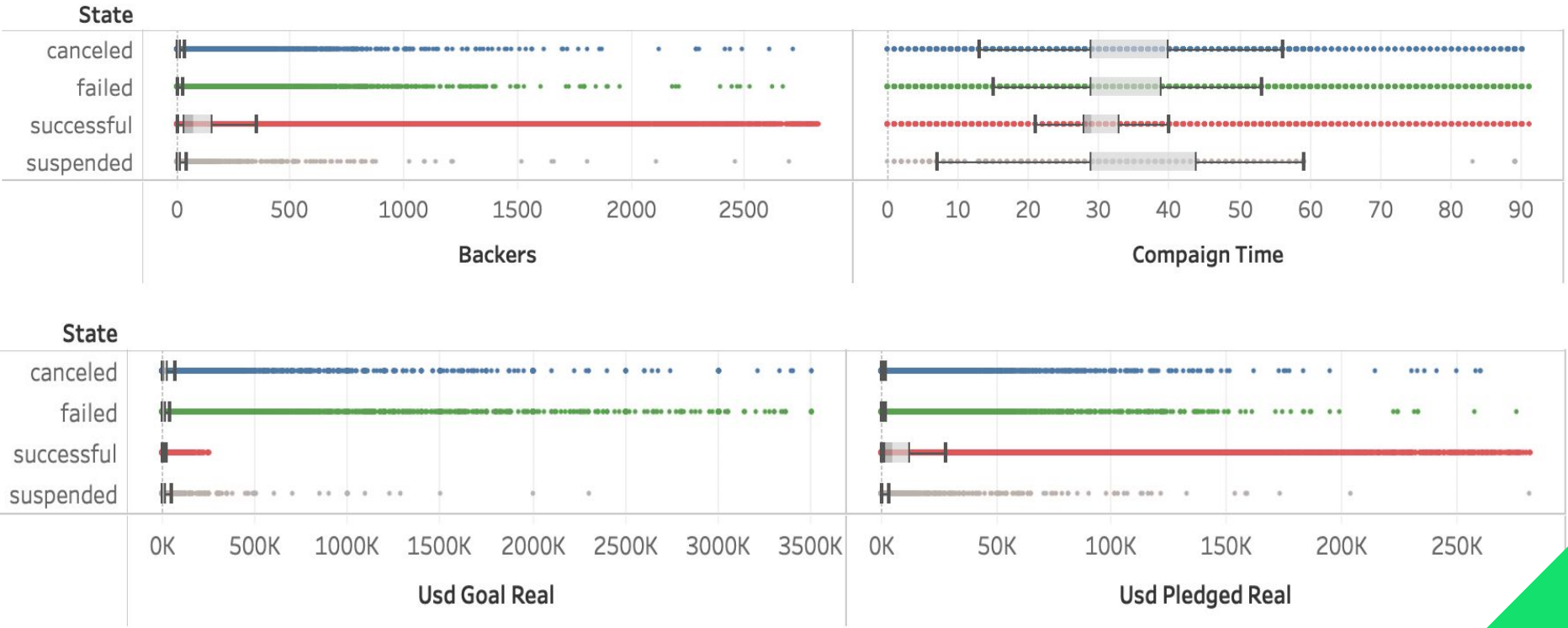
# EDA







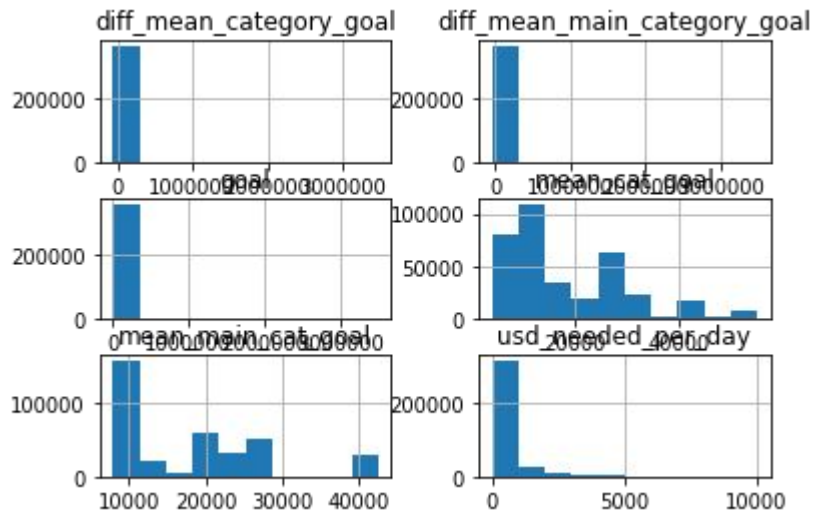
# Successful v.s. Others



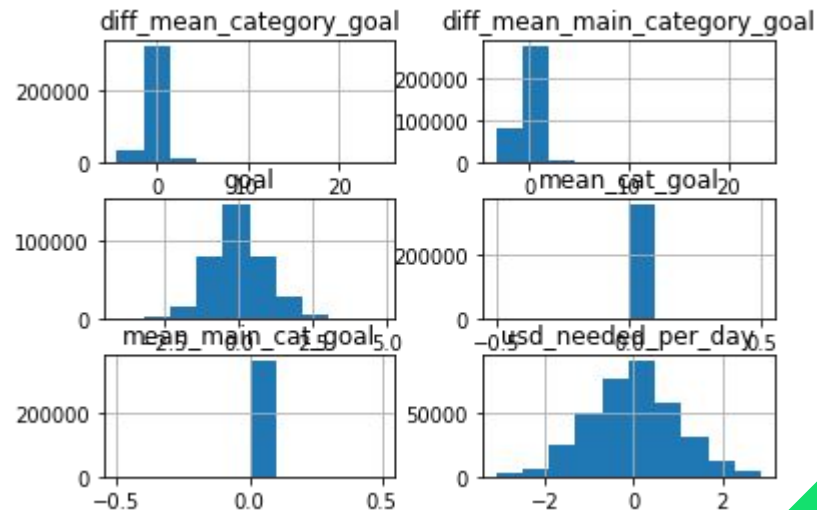
# Standardization of the data

```
1 trans = PowerTransformer().fit_transform(X = df[['goal', 'usd_needed_per_day'\n2 , 'mean_cat_goal', 'mean_main_cat_goal'\n3 , 'diff_mean_category_goal', 'diff_mean_main_category_goal']])
```

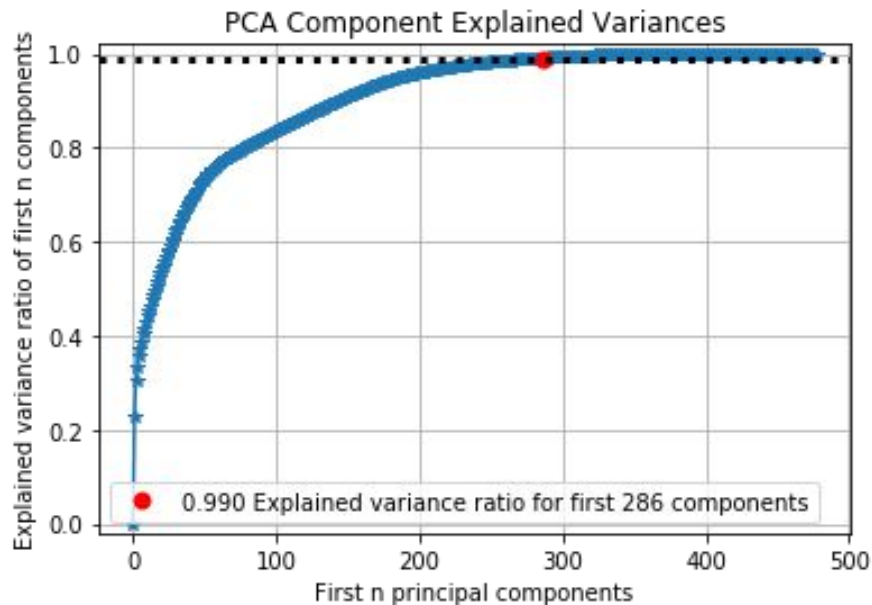
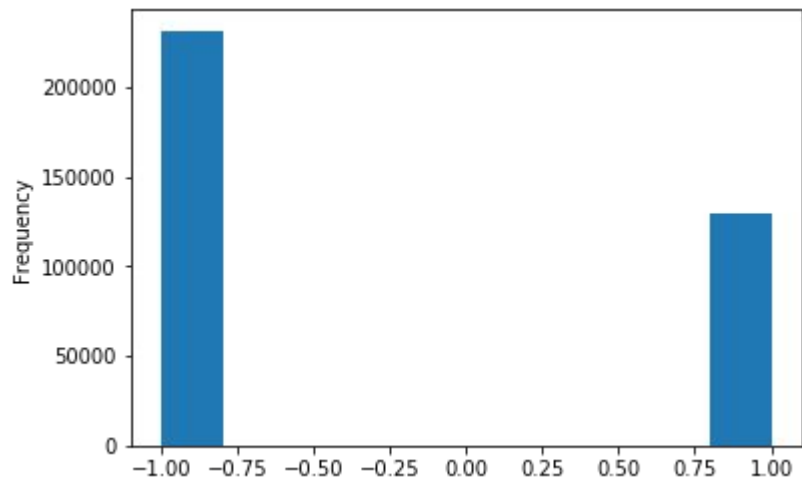
Before



After



# Oversampling the data & PCA





# Experiments

| Model                        | Accuracy       |
|------------------------------|----------------|
| Logistic regression          | 69.2%          |
| Extra Tree Classifier        | 68%            |
| SGDClassifier                | 68.7%          |
| RBFSampler                   | 66.5%          |
| Neural Network               | 69%            |
| Gradient Boosting Classifier | 69%            |
| XGBoost                      | Takes too long |

# Result: LogisticRegression

```
clf = LogisticRegression(C=0.01, n_jobs=-1)
```

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| -1           | 0.76      | 0.75   | 0.76     | 76354   |
| 1            | 0.57      | 0.58   | 0.58     | 42826   |
| micro avg    | 0.69      | 0.69   | 0.69     | 119180  |
| macro avg    | 0.67      | 0.67   | 0.67     | 119180  |
| weighted avg | 0.69      | 0.69   | 0.69     | 119180  |

