# Kickstarter Project: Predicting Success

Final Report

# Kickstarter

Crowdfunding Platform

$4,000,000,000+ towards projects

168,000+ successfully funded projects

16,000,000+ users

55,000,000+ total pledges

# The Problem

Only 37.22% of Kickstarters have been successful

Creators may have to abandon their first go-round with a Kickstarter campaign

# Proposed Solution

Given what we know about projects such as:

     Description Lengths, Goal Amount, Category, Country of Origin,

     Time of Creation, Length of Campaign

We would like to know the success rate of whether or not a project will succeed.

To do this, we'd like a model that can take various inputs to answer whether a project will reach its goal.

# Kickstarter Dataset

Kickstarter does not provide a public API

Web Robots (webrobots.io) scraped Kickstarter to provide this data
- 2019-05-16
- 210,000 rows with 37 columns
- 56 csv files

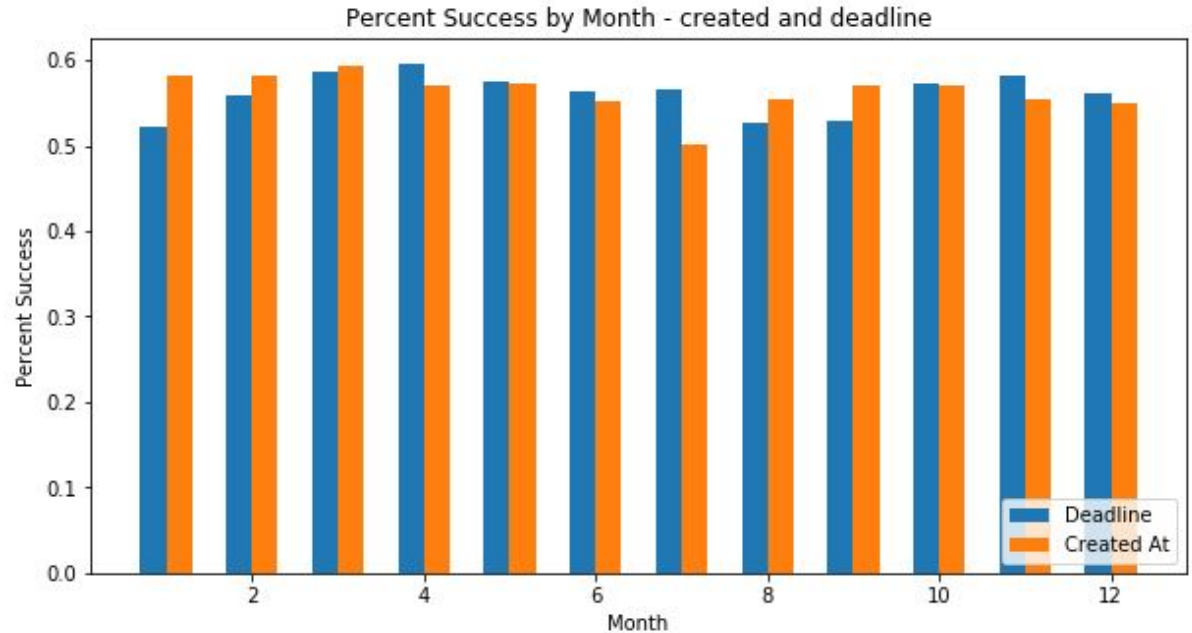Additional web scraping utilizing Google Cloud Platform for more information
- Use URLs from the WebRobots data
- Pull information such as descriptions & rewards

# Exploratory Analysis

# Projects by Month

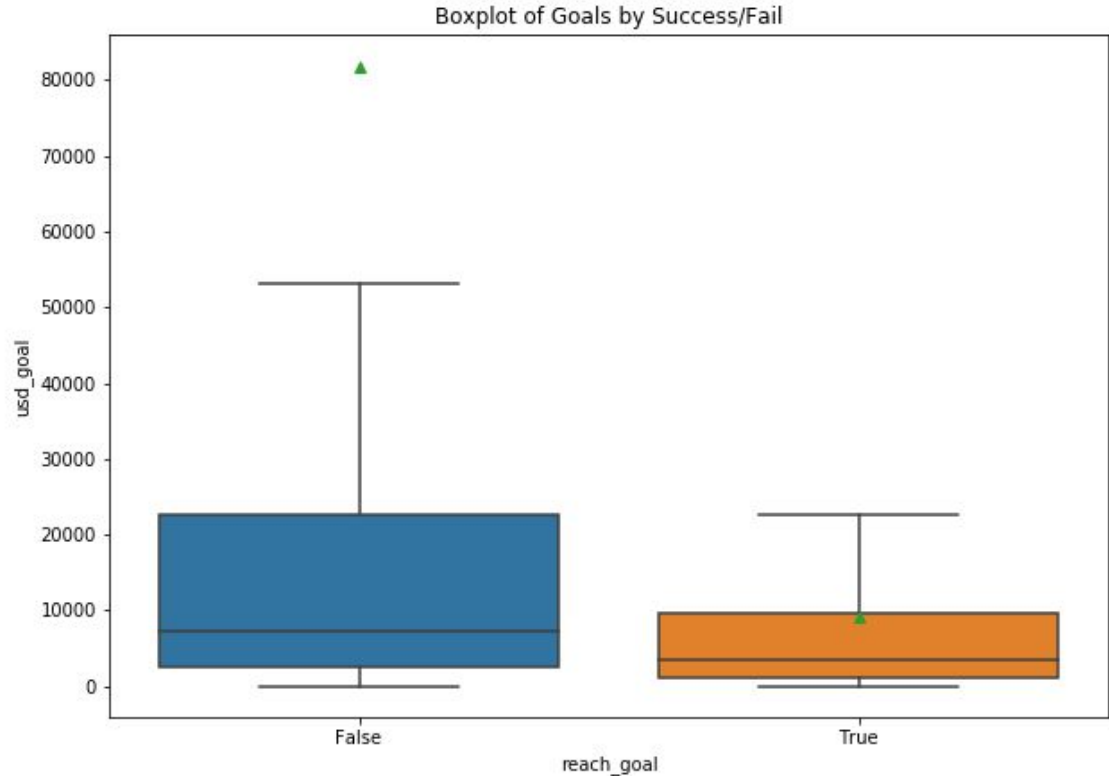Projects started in March are the most successful while projects started in July are the least.

Alternatively, projects finishing in January are the least successful while finishing in April has the highest success.



Percent Success by Month - created and deadline

# Success to Goal (USD)

In general, projects which are successfully funded have much lower mean goals (USD) than projects who have failed.

Additionally, the median of successful projects are also lower than the median of failed projects.
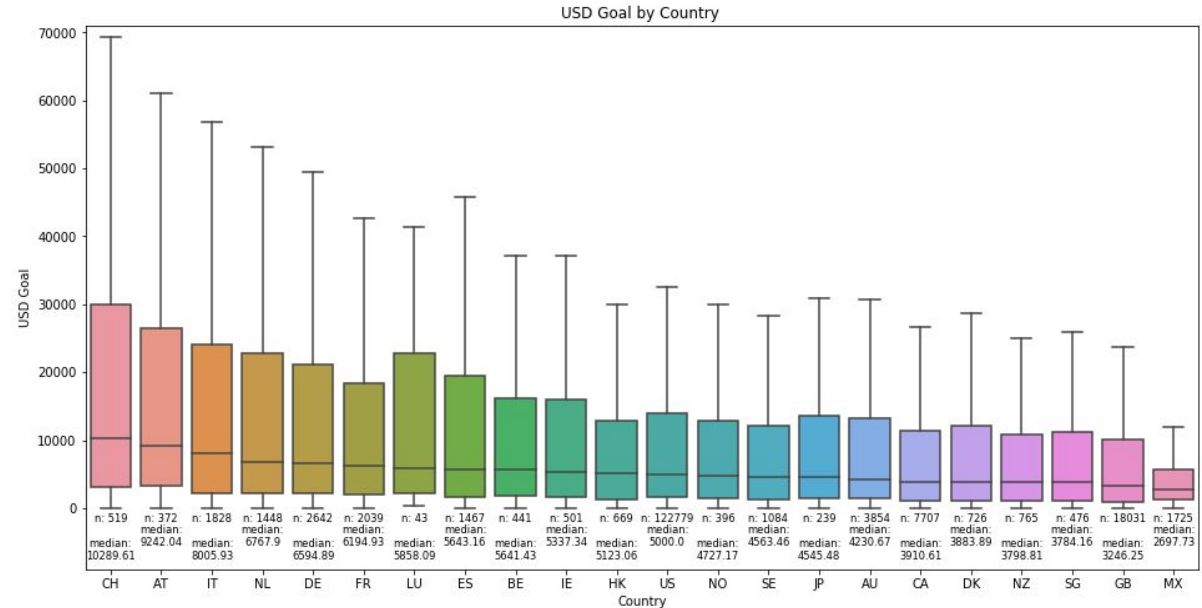
# Country of Origin

In the boxplot of Goals (USD) by Country of Origin, we can see China has the highest median goal at $10,289. The next would be Austria and Italy at $9,242 and $8,005, respectively.
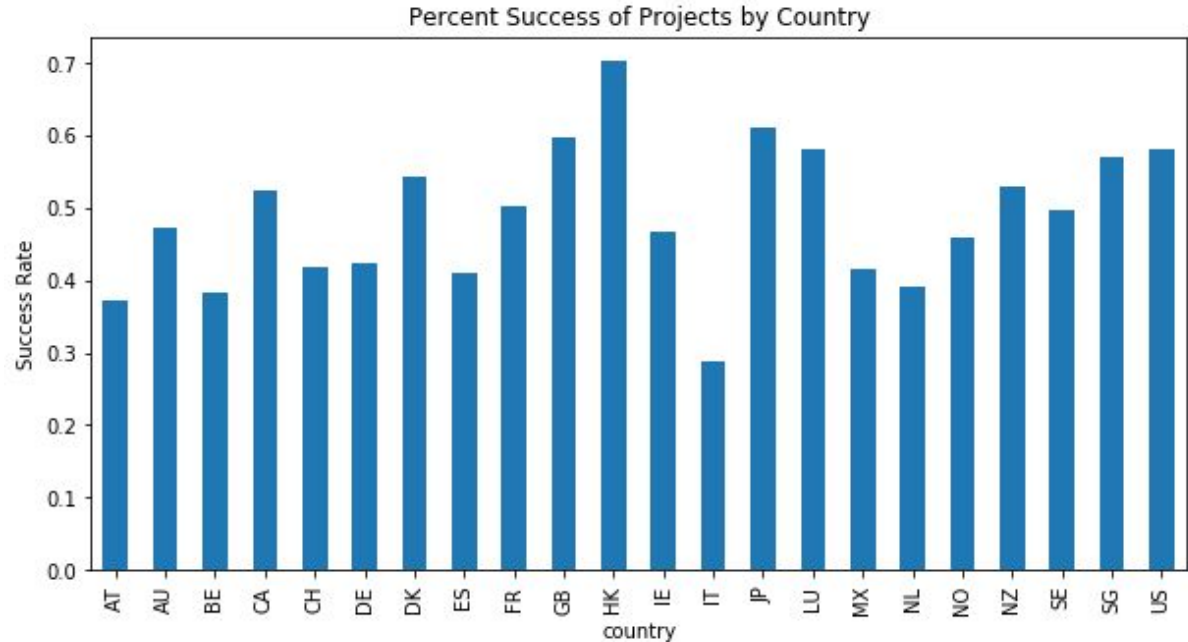
After, median goals seem to balance off around 5,000 to 2,000 USD. The spread also looks to decrease at median goals decrease.



USD Goal by Country

# Country of Origin

However, the three countries who have the highest median goals do not have the highest success rates. Italy actually has the lowest success rate for projects.
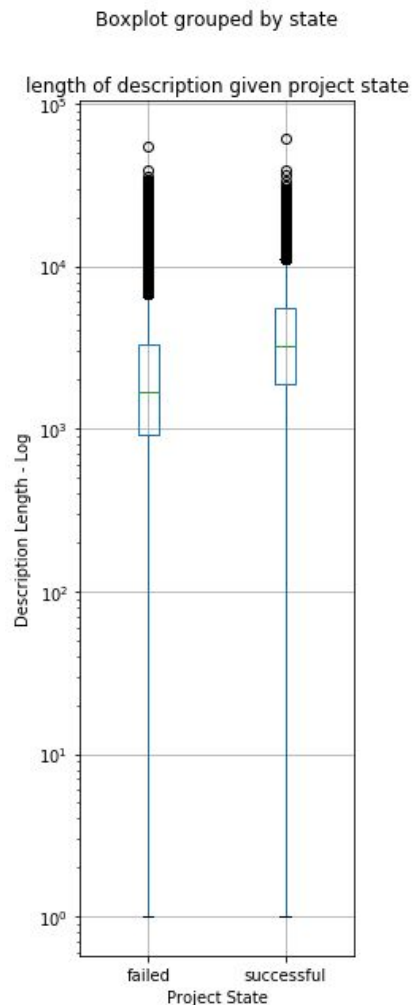
Hong Kong and Japan, countries with a median goal of 5,100 and 4,500 have the two highest success rates but are also on the lower end of the amount of projects from the country.



Percent Success of Projects by Country

# Project Descriptions

In the boxplot for description lengths of successful and failed projects, we can see that the median description length of successful projects is higher than failed projects.
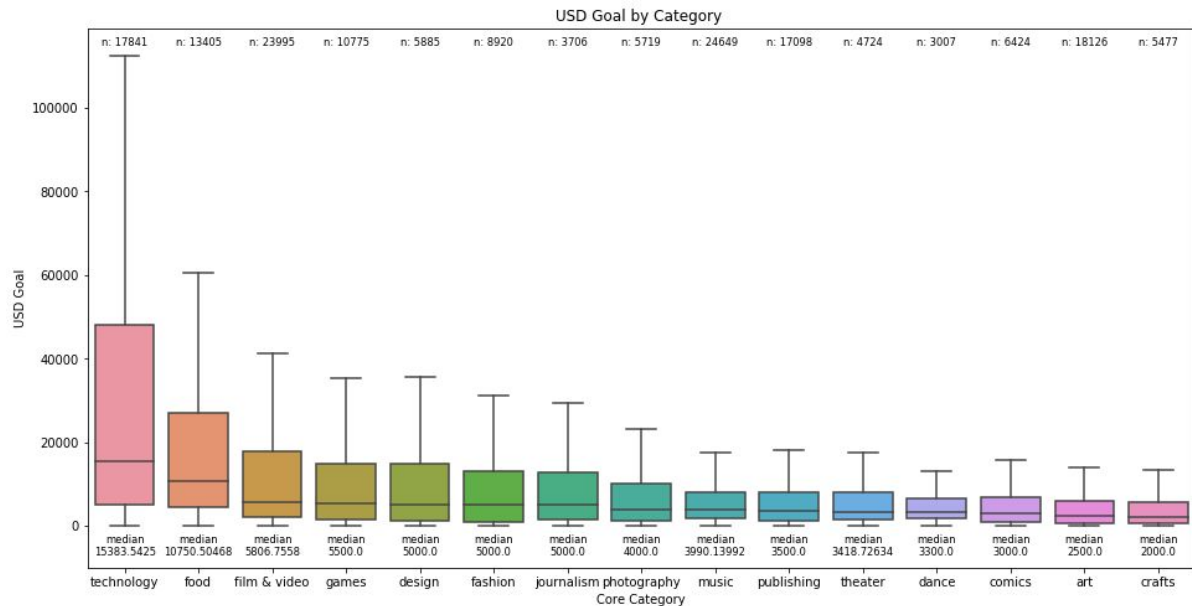
It's worth noting that some projects, both successful and failed, have descriptions of length 0, meaning they have used images and/or video to describe their project.



Boxplot grouped by state

length of description given project state

# Category

Looking at categories to their Goal (USD), we can see that technology has the highest median goal at $15,383. Followed by food, then film & video.

Size of each category varies widely ranging from 24,000 to 3,000 in each category.
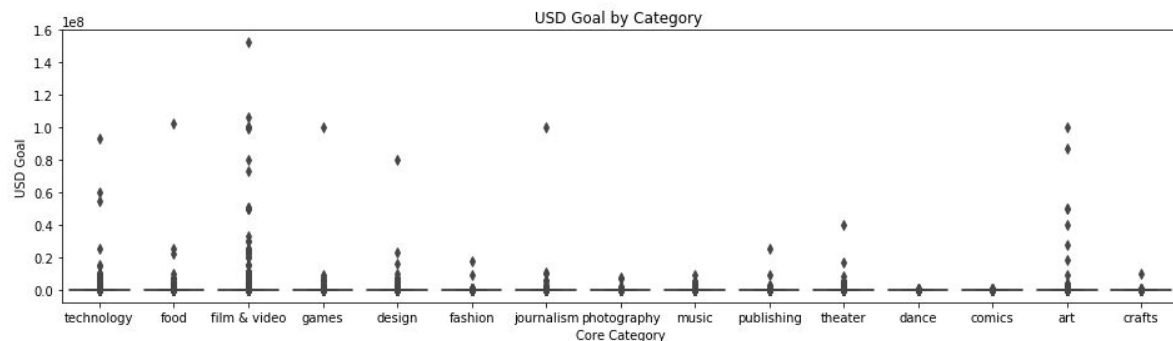


USD Goal by Category

# Category

Looking at the outliers for the previous boxplot, we can see that many categories are afflicted by high outliers such as technology, film & video, and even art.

Film & Video has the highest outlier goal at 1.6 million.

Tech, Film & Video, and Art have outliers across their range but some categories like Journalism are afflicted with gaps of goals.



USD Goal by Category

# Category

Performing a Tukey test, we can see that technology and film & video are the most different from the other categories when it comes to their goals.

Intuitively, this makes sense since tech and videography projects would seem to expend more money than other projects.

| group1 | group2 | meandiff | lower | upper | reject |
|--------|--------|----------|-------|-------|--------|
| art | film & video | 75098.7355 | 38120.565 | 112076.9061 | True |
| comics | film & video | 101276.2274 | 48490.4682 | 154061.9866 | True |
| comics | technology | 64089.8212 | 9415.2691 | 118764.3734 | True |
| crafts | film & video | 99138.275 | 42867.8098 | 155408.7403 | True |
| crafts | technology | 61951.8688 | 3905.8911 | 119997.8466 | True |
| dance | film & video | 101035.8861 | 28345.2743 | 173726.498 | True |
| fashion | film & video | 92658.3509 | 46060.9628 | 139255.7391 | True |
| fashion | technology | 55471.9448 | 6745.2902 | 104198.5993 | True |
| film & video | food | -58887.1959 | -99405.3324 | -18369.0593 | True |
| film & video | games | -72766.196 | -116341.5353 | -29190.8566 | True |
| film & video | music | -97359.664 | -131436.7193 | -63282.6087 | True |
| film & video | photography | -95079.9702 | -150372.6429 | -39787.2974 | True |
| film & video | publishing | -97447.5457 | -135053.6495 | -59841.442 | True |
| film & video | technology | -37186.4062 | -74332.4488 | -40.3636 | True |
| film & video | theater | -72442.0367 | -132252.4474 | -12631.626 | True |
| music | technology | 60173.2578 | 23237.9634 | 97108.5523 | True |
| photography | technology | 57893.564 | 794.9656 | 114992.1623 | True |
| publishing | technology | 60261.1396 | 20046.8471 | 100475.432 | True |

# Category + Description

We can also explore how each category has certain word choices for their project.

**One** for "one of a kind" or "one perfect brew"

**Design** relates to categories like design, fashion, and art.

**Device** is seen in technology

**Recipe** & **Farm** appear in food

# Blurbs

Across all projects, we can see what word choices they use to entice prospective backers to click their campaign.

**One** is used as superlatives i.e. "One of the most"

**World** is used in different manners: the world as a whole and world of their projects
- "world premier", "around the world"
- "world of wine", "art world"

# Blurbs

We can also take a look at predictive features of words using MultinomialNB.

Few have >.90 as their predictors

Unusual ones such as "EDC" for Design and "28mm" for Games.

International keywords appear in Film & Video as "cortometraje" and in Theater as "Teatro"

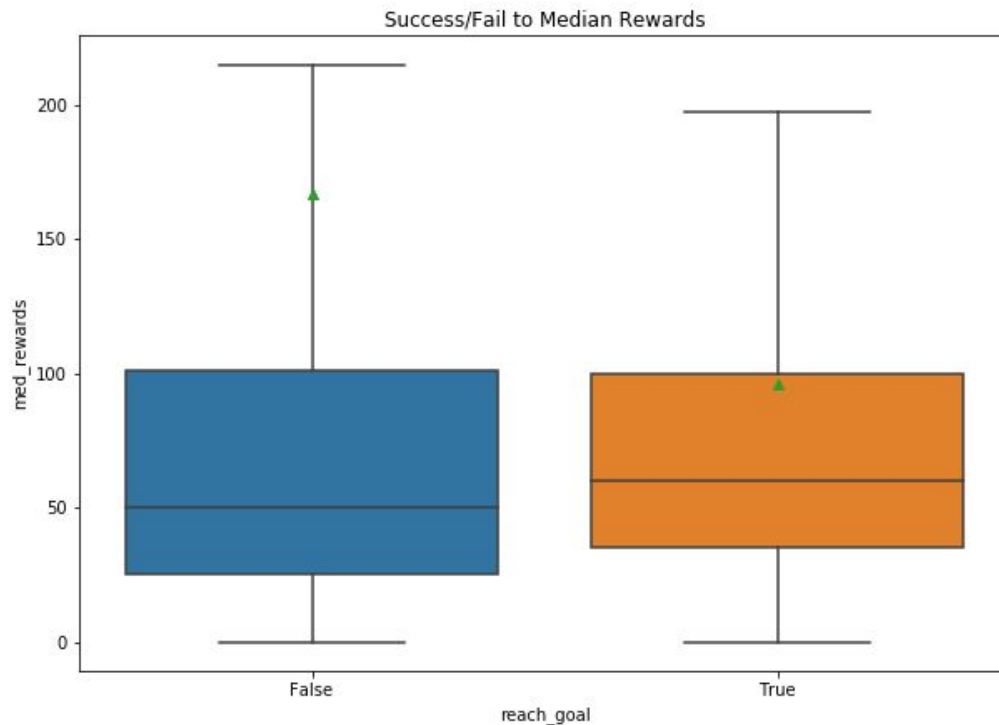| ART: | COMICS: | CRAFTS: |
|---|---|---|
| Sculpture \| 0.93 | Webcomic \| 0.87 | Candles \| 0.92 |
| Mural \| 0.91 | Comic \| 0.78 | Soaps \| 0.79 |
| Installation \| 0.83 | Comics \| 0.72 | Pens \| 0.76 |
| Sculptures \| 0.83 | Graphic \| 0.63 | Candle \| 0.74 |
| Painting \| 0.83 | Manga \| 0.58 | Scented \| 0.74 |
| DANCE: | DESIGN: | FASHION: |
| Choreographers \| 0.77 | Font \| 0.61 | Footwear \| 0.87 |
| Ballet \| 0.71 | Titanium \| 0.60 | Sandals \| 0.79 |
| Dance \| 0.66 | EDC \| 0.58 | Clothing \| 0.78 |
| Choreographer \| 0.65 | Poster \| 0.54 | Shoes \| 0.78 |
| Choreography \| 0.63 | Logos \| 0.53 | Apparel \| 0.77 |
| FILM & VIDEO: | FOOD: | GAMES: |
| Webseries \| 0.93 | gourmet \| 0.93 | platformer \| 0.90 |
| Cortometraje \| 0.91 | bakery \| 0.92 | uspcc \| 0.90 |
| Film \| 0.87 | sauces \| 0.92 | rpg \| 0.88 |
| Mockumentary \| 0.86 | sauce \| 0.89 | 28mm \| 0.88 |
| Animated \| 0.84 | brewery \| 0.89 | strategy \| 0.86 |
| JOURNALISM: | MUSIC: | PHOTOGRAPHY : |
| journalism \| 0.66 | ep \| 0.97 | photobook \| 0.87 |
| news \| 0.50 | album \| 0.97 | nudes \| 0.77 |
| journalists \| 0.45 | cd \| 0.94 | photographing \| 0.73 |
| reporting \| 0.41 | lp \| 0.92 | photographic \| 0.71 |
| coverage \| 0.39 | recording \| 0.92 | photography \| 0.66 |
| PUBLISHING : | TECHNOLOGY : | THEATER : |
| poems \| 0.82 | arduino \| 0.97 | fringe \| 0.74 |
| chapbook \| 0.79 | raspberry \| 0.91 | edinburgh \| 0.69 |
| literary \| 0.76 | wireless \| 0.89 | theatre \| 0.68 |
| rhyming \| 0.74 | bluetooth \| 0.89 | playwrights \| 0.59 |
| essays \| 0.74 | pi \| 0.88 | teatro \| 0.58 |

# Rewards

Looking at a zoomed in view of a scatter plot, we can see that past a certain combination of median rewards and goal do failed projects starter to appear.

At a median reward of $25, projects that have a goal of 10,000 tend to fail. At a median reward of $50, they tend of fail at higher goal levels.

# Rewards

When looking at the boxplot, however, successful and failed projects tend to have the same range in their data. They are more affected by outliers since failed projects tend to have a higher, more unattainable, goal.



Success/Fail to Median Rewards

# In-Depth Analysis of Machine Learning Models

# Scoring Method

Accuracy Score is a good base but does not address our business case

Focused mainly on predicting failure

Precision to improve confidence on determining failure

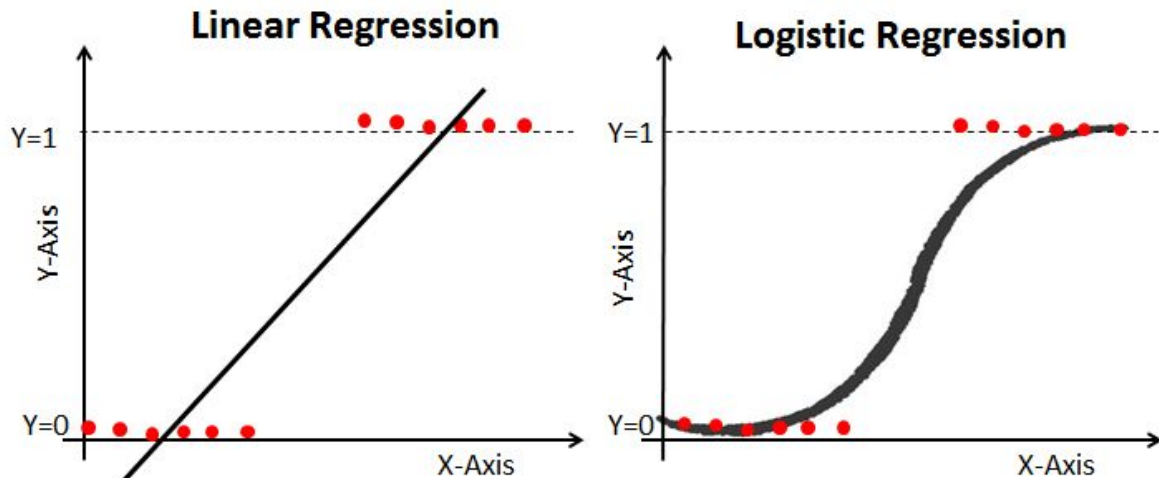Use fbeta_score with a beta of 0.5 to double the weight towards precision

# Logistic Regression

Scale Features using RobustScaler

Get dummy variables using Panda's Get_Dummies

Using our fbeta_score, base LogReg classifier has 73.06% on testing.
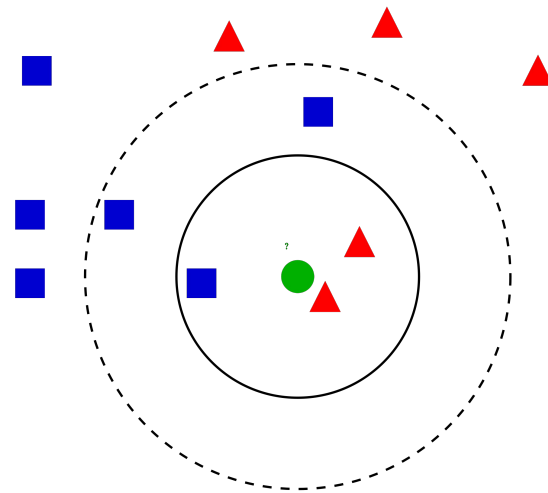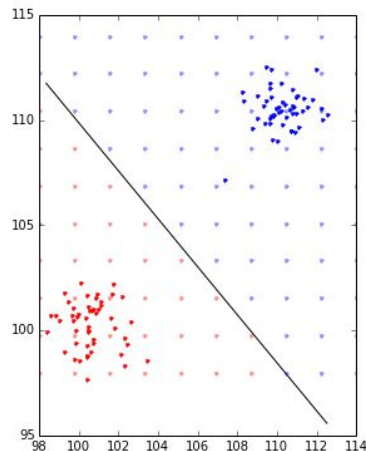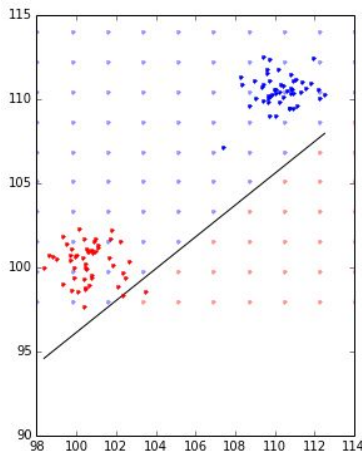
With some parameter tuning, we get a score of 74.40%

# SVM & KNN

Testing an SVM gets a base score of 72.85%.

We get a warning for convergence. Warning does not go away even after 4000 iterations. End up not using this.

KNN stores all data points so it becomes inefficient for our data size. End up not using this.
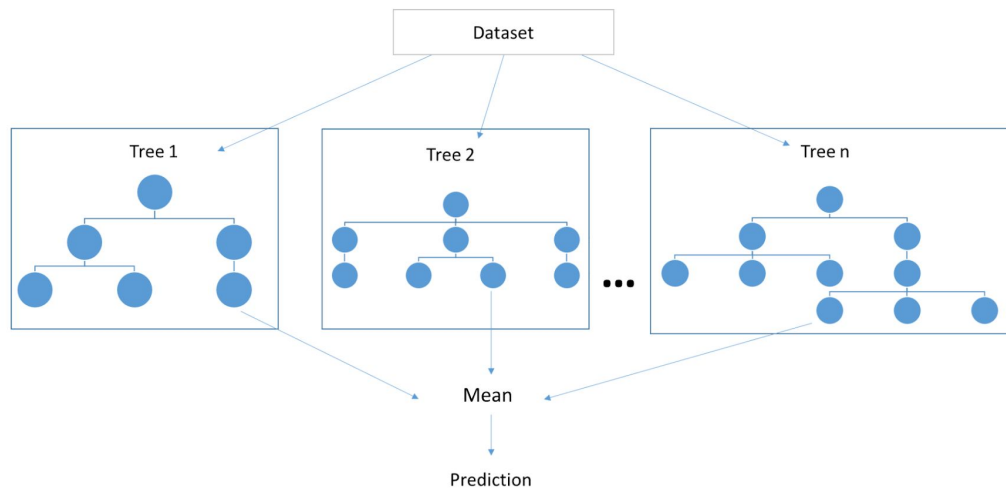
# RandomForestClassifier

Does not need to scale.

Base score: 74.97%

Many parameters- use
RandomizedSearchCV

Best Params Score: 77.72%

Best Params Score through CV:
77.65%

# Final Thoughts & Considerations

With a fbeta score of 77.6%, we can provide future Kickstarter Creators a glimpse into their project by running some of their project details through our algorithm. If it predicts that it will fail, we're confident that there are improvements to be made with the project.

While this is a good first step, there are many more pieces of information we can improve on. We would want to start breaking down the descriptions of projects to get a better sense of keywords that attract success.

We would also want to improve the model itself. While it is nice to know that a project at this point will fail, it will be more helpful to know what specifically about the project causes it to fail.