

# Capstone Project 1: Final Report

## Crunchbase Funding Prediction Model

### Proposal

For my First Capstone Project I decided to use Crunchbase data to build a model to see if it was possible to predict whether or not a company would be funded in the following year based on data at a given time.

### Problem Statement

It is difficult to truly understand what it takes to start a company and successfully raise capital as there are an endless amount of variables many of which cannot accurately be quantified. Using data provided by Crunchbase I hope to gain some level of understanding of the startup funding.

### Why This Topic?

I chose this topic as I have a strong interest in startups. I found the crunchbase data very interesting and wanted to build some sort of model based on what I was learning in the course. In the end I decided to see if it was possible to predict whether a start up would get funding in a specific time frame based on round, company, investor, and other information.

Startups, funding, and venture capital are growing topics in the US and around the world with more and more startups getting funded everyday. I have always loved business and am constantly trying new things and starting new projects but the more I thought about starting or joining a startup the more I realized that I had no clue about the market, risks, and overall viability of getting an idea, raising a capital and being successful. For this reason, I chose to explore this as the topic for my first capstone to try to gain some insight into the startup & funding world.

### Procedure

1. Data Cleaning & Wrangling
2. Exploratory Data Analysis
3. Machine Learning
4. Model Evaluation
5. Analyze Results

## Data Collection & Cleaning Summary

The data I will be using for this project was taken from a GitHub repo. We did have access to the Crunchbase Open Dataset, however, the data only went up to 2013. We were able to find a GitHub repo of someone with access to a more complete version of the data (up to 2015) which we compared to the open data to test for accuracy and decided that the more complete dataset would be better to work with. Although the data is not up-to-date, given access to the full dataset we would be able to run it through the model in the same way to get an updated version of the analysis as the variables are exactly the same.

As we did pull the data from a GitHub repo, no complicated scraping or data extraction was required for the sake of this project so no interesting display of abilities here. However, cleaning and wrangling of the data was extremely time consuming due to the amount of data used and the variability in the data.

## Data Cleaning & Wrangling Summary

The data wrangling process was quite extensive as we needed to accumulate all of the data into features which could be used to train my Random Forest Classifier. This involved transforming all non-numeric data to numeric features as well as filling in missing values and creating new features to help summarize the data more effectively.

### Raw Data

The data was broken up into the following dataset:

	Shape	Numeric Data	Text Data	Datetimes
<b>Investment/Rounds</b>	(212810, 18)	1	16	1
<b>Companies</b>	(51146, 14)	2	9	3
<b>Organizations</b>	(606064, 16)	0	16	0
<b>People</b>	(605630, 15)	0	15	0
<b>Acquisitions</b>	(18968, 18)	1	15	2
<b>IPOs</b>	(1259, 13)	2	8	3

### **Basic Cleaning Steps:**

1. Rename Column names
2. Indicate Null Values
3. Deal (Impute/Drop) with Null Values
4. Strings to Numbers/Datetimes
5. Export Clean Data

### **Basic Wrangling Steps**

1. Create format of dataframe with IDs, Dates, etc
2. Add or summarize numeric values and add to dataframe
3. Text variables as dummies, summarized, or vectorized
4. Export Processed Data

### **Processed Dataset**

Each row represents one fiscal quarter of a company after their first round of funding. Features were chosen to best represent the information provided in the original data with emphasis on some characteristics that seemed more important than others. The target variable will indicate whether or not that company raised a round in the preceding year of that quarter.

There is a large emphasis on Investor data as we used investors to help determine traits of the company that were not quantifiable in the data such as founder character, and confidence. By using investor track record and

As we are planning to use a Random Forest Classifier for this model we will need to represent all non-numeric variables as numeric variables. The random forest is able to ignore noise from less important features but we will do our best to minimize the number of features as we want to be as computationally efficient as possible.

## **Exploratory Data Analysis Summary**

**The EDA for the project was broken up into three main sections.**

1. Investor Analysis
2. Investment Analysis
3. Category Analysis

As there is a wide range of data it took a large amount of exploratory data analysis (EDA) to really get a feel for the entirety of the data and the important features. We split the EDA into 3 separate notebooks to help with organization. The notebook provided includes the analysis of the features finally selected. This being my first big project I did a lot of extra exploration both to be sure that I got everything from the data and to just practice different visualization and analysis tools learned in the course.

## Pt 1. Initial Visualizations

Basic Visualizations were made during the data wrangling of the project to try to find interesting features to add to the model. During this initial exploration I found myself asking different questions about the data which were then grouped into the EDA notebooks that are found in the project folder. These EDA notebooks contain the more in-depth EDA and Inferential Statistics work.

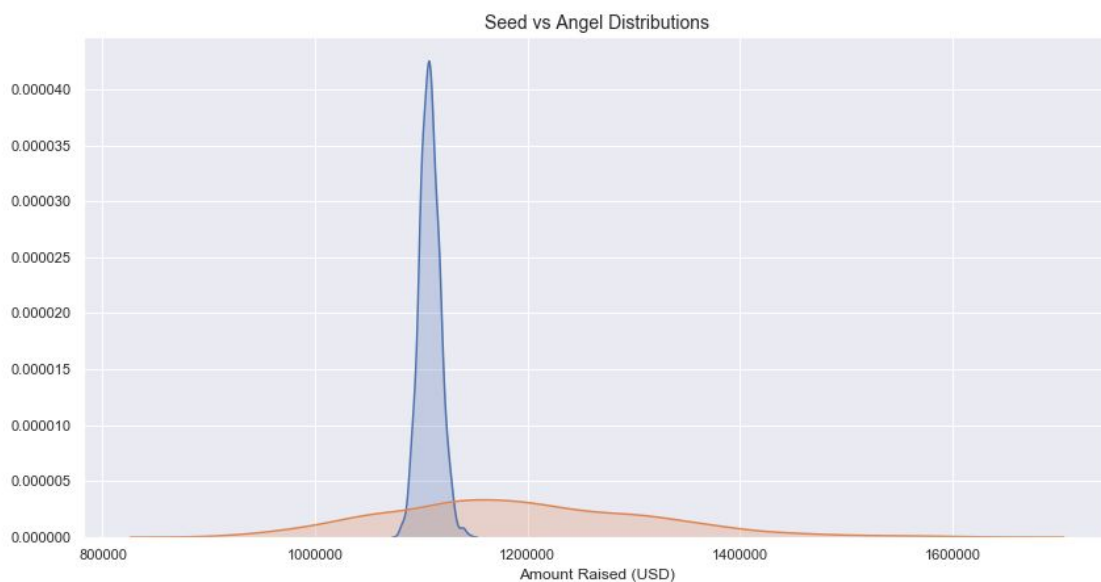
## Pt 2. Inferential Statistics

The EDA Notebooks consist of more detailed visualizations and the inferential statistics. For more clear organization this section was broken up into the three notebooks: Investor, Investment and Category Analysis. The EDA was based off questions that came up during the initial visualization and exploration of the data when choosing features and building the dataset for the machine learning model. The Analysis and Questions are organized as follows:

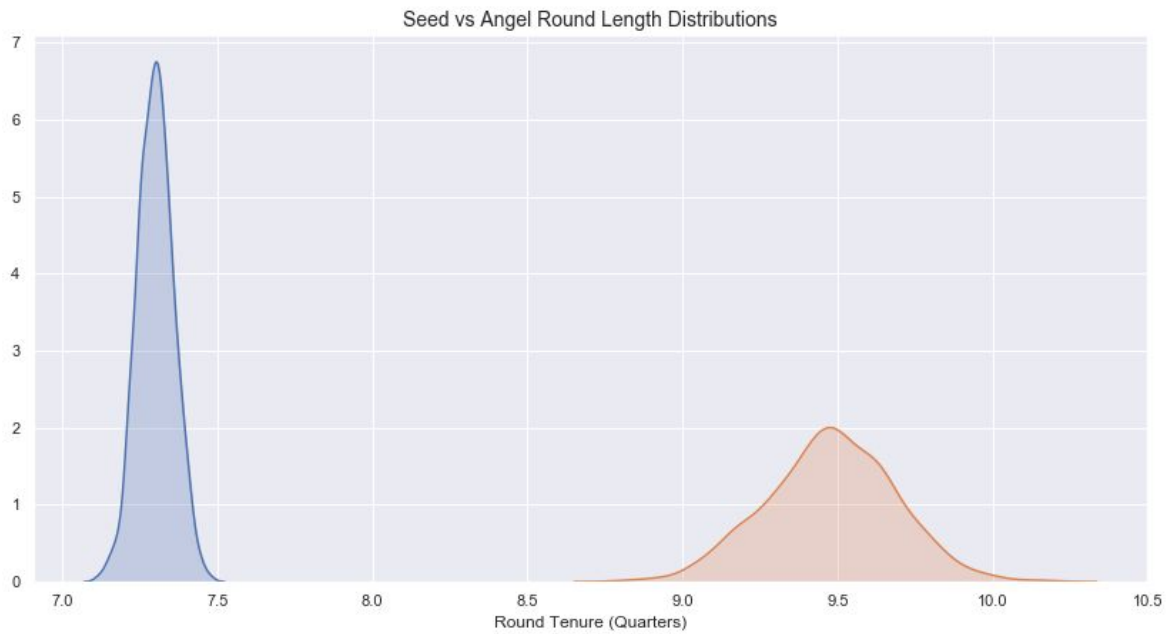
### 1. Investment Analysis

- a. Do Different Rounds Raise Different Amounts?
- b. Do Different Round Types Have Different Tenures?

#### Seed Vs Angel Funding Distribution (Other Rounds In Notebook)



### Seed vs Angel Round Lengths (Other Rounds In Notebook)



### Seed vs Angel Company Tenure at Funding (Other Rounds In Notebook)



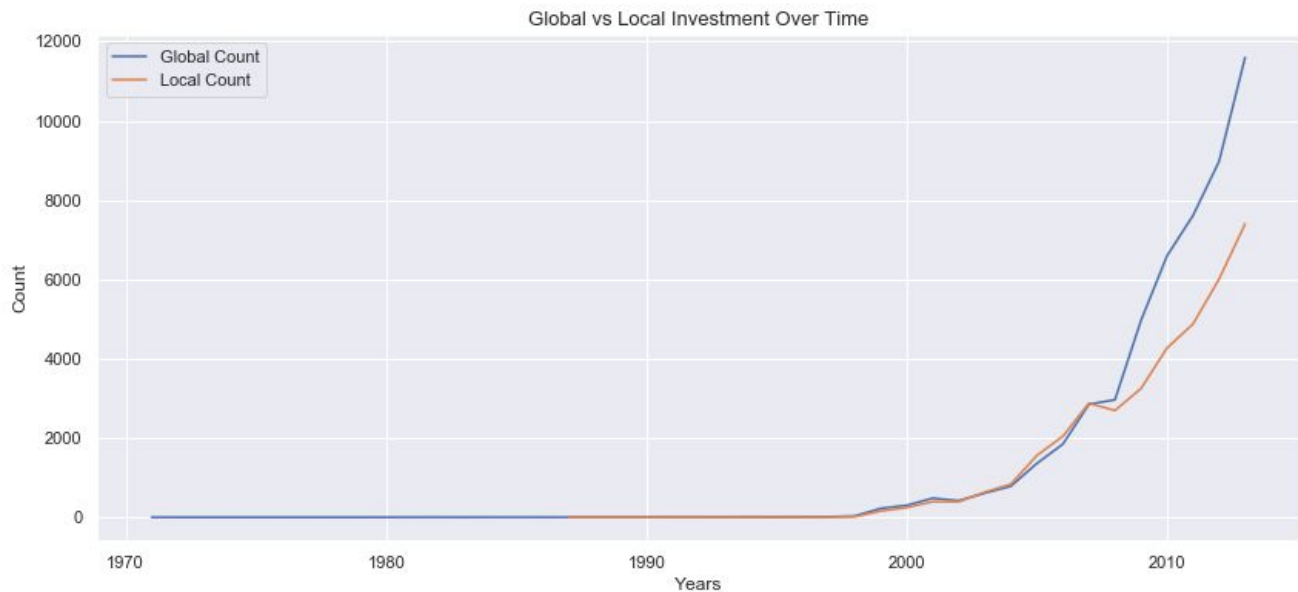
## Investment Analysis: Key Observations

- Amount Funded by Angels Varies More Than Seeds
- Later Round Typically Have Higher Variance in Amount & Tenure. (Not Shown)
- Venture Round Lengths Tends to be between 10 -15 Quarters. (Not Shown)
- Equity Funding is Greater Than Debt Before IPOs but Less After. (Not Shown)

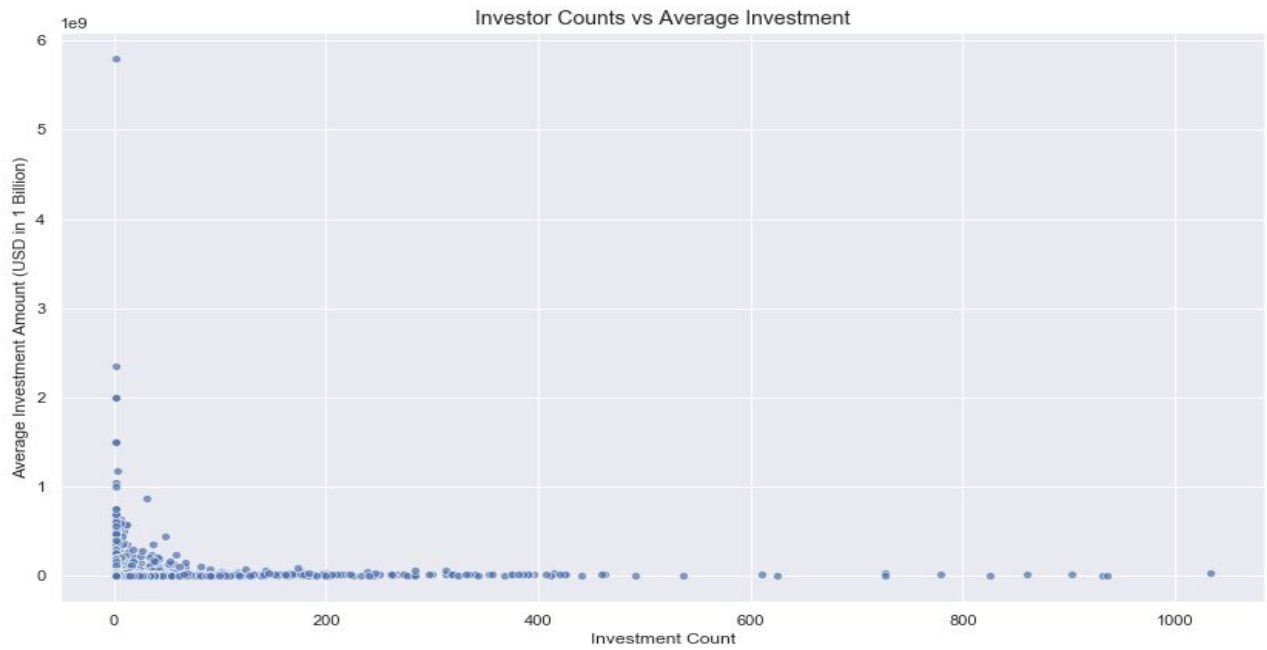
## 2. Investor Analysis

- a. Do Global Investors Fund Differently Than Local Investors?
- b. Do Investors Fund Differently?
  - i. Quantity vs Quality?
  - ii. Focus on Specific Rounds?

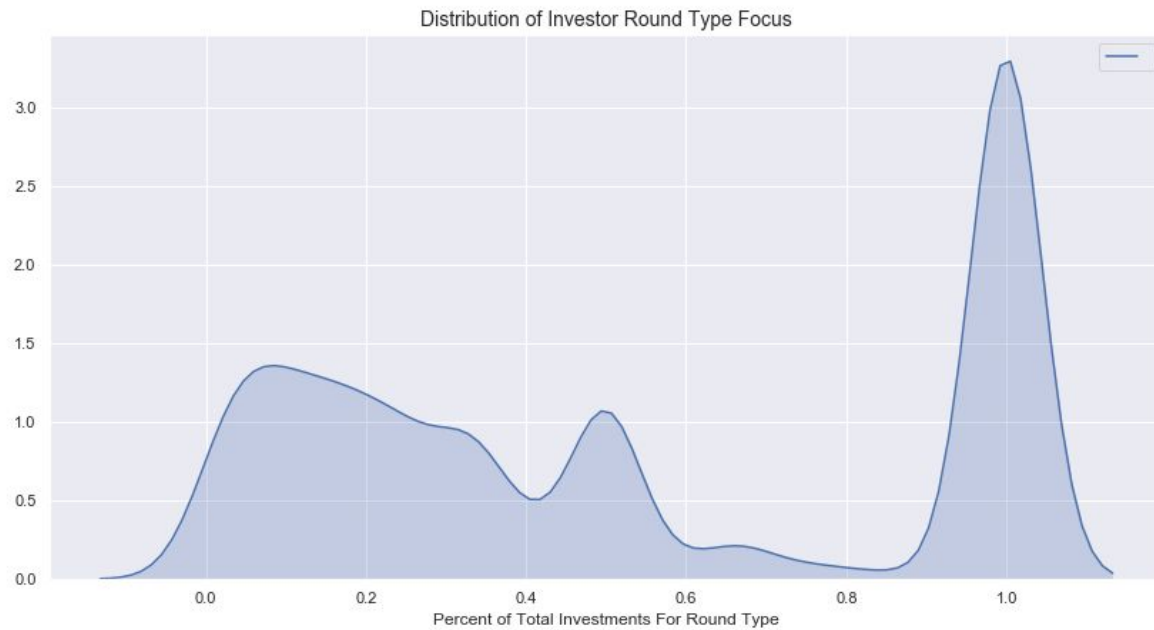
### Analyzing Globalization of Investing (Full Analysis In Notebook)



### Analyzing Difference in Investment Strategies (Full Analysis in Notebook)



### Do Investors Focus On Specific Round Types (Full Analysis In Notebook)



## Investor Analysis: Key Observations

- There are more gaps in data of global investors
- Difference in invested amount of global and local investors is **NOT** statistically significant
- Investors do NOT focus on just **many small** or **few large** investments.
- **Most investors focus on more than one round type.**
- **Seed investors** tend to **focus more** on seed funding.

### 3. Category Analysis

- a. How has funded categories changed in the past 30 years?





## Category Analysis: Key Observations

- **1980s:** Manufacturing, Services, Designers, Automotive, Technology
- **1990s:** Software, Technology, Curated Web, Service, Internet
- **2000s:** Software, BioTech, Enterprise Software, Mobile, Curated Web
- **2010s:** Software, Enterprise Software, Mobile, Curated Web, Commerce

## The Model

**Model Type:** Random Forest Classifier

The predictive model I decided to use in this project is the Random Forest Classifier. This seemed to be the best model for this problem as the model is **flexible, prevents overfitting, is good with high feature and sample counts and does not require scaling.**

### **Data:**

The data used for the model consists of **111 numeric features** and **2 text features** (Company Description and Categories).

### **Pipeline Steps:**

The first step of the pipeline is to prepare the data for the model. Because we have both numeric and text data we deal with these separately and combine them after processing with a FeatureUnion. We then take the processed data and run it through our model.

#### **1. FeatureUnion**

Inside our feature union we fill the missing numeric values with a Simple Imputer. And vectorize the 2 text features using a **HashingVectorizer** looking at **1 and 2 ngrams**. Another option would have been a CountVectorizer but because the amount of text and companies, the computational cost of having the fully vectorized text would have been too large. To reduce features even further we ran the vectorized outputs through **SelectKBest** to get only the **top 3000 features vectors**.

#### **2. Random Forest Classifier**

With the data processed we are left with **3000 hash vector features and 111 numeric features** which we ran through our model. **The target variable in this model is a prediction of whether or not the company in a given row will get funded in the coming year.** After testing the model with different hyperparameters we were not able to find one that performed better than the **default parameters** so in the end the final model uses these parameters.

## Model Evaluation

In the evaluation I looked at different evaluation metrics to show model performance overall as well as looking deeper into the specific errors by round type to see if some round types were more difficult to predict.

**Accuracy:** 0.94

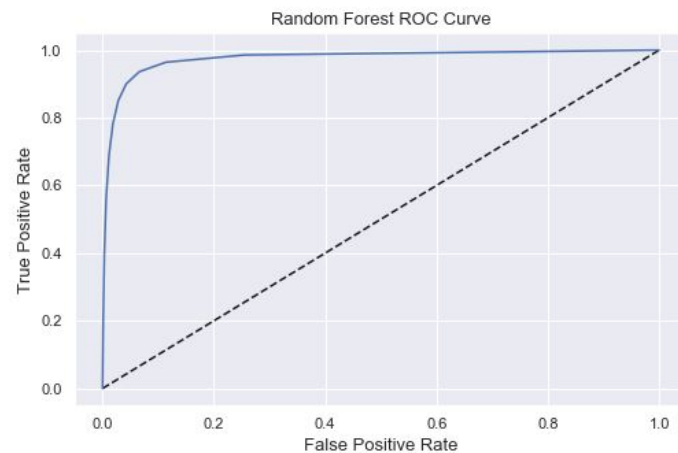
**Avg CV Score:** 0.89

**Avg Precision:** 0.94

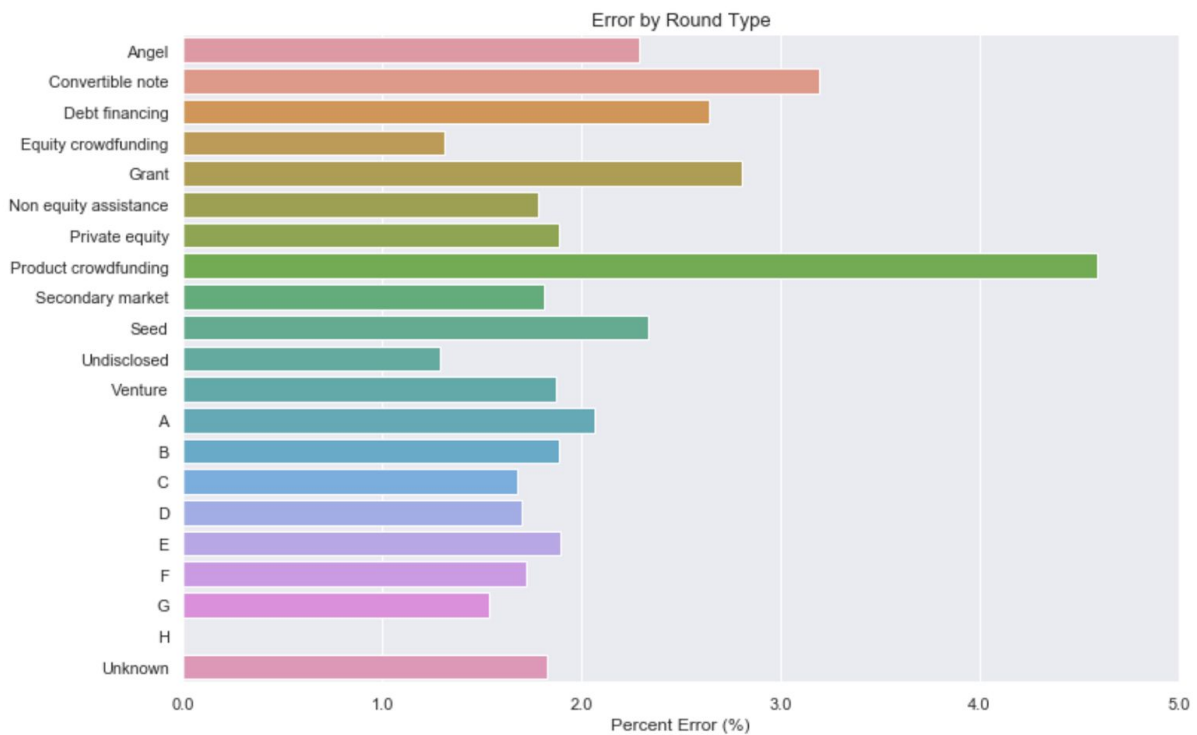
**Avg Recall:** 0.94

**ROC Curve:** *Figure on Left*

**AUC Score:** 0.975



### Error by Round Type



## Precision vs Recall

Based on the classification report below we can see that the model had some trouble with recalling funded targets.

**Precision:** 91% of Predicted Funded Were Actually Funded.

**Recall:** 78% of Actual Funded Rounds Were Identified as Funded.

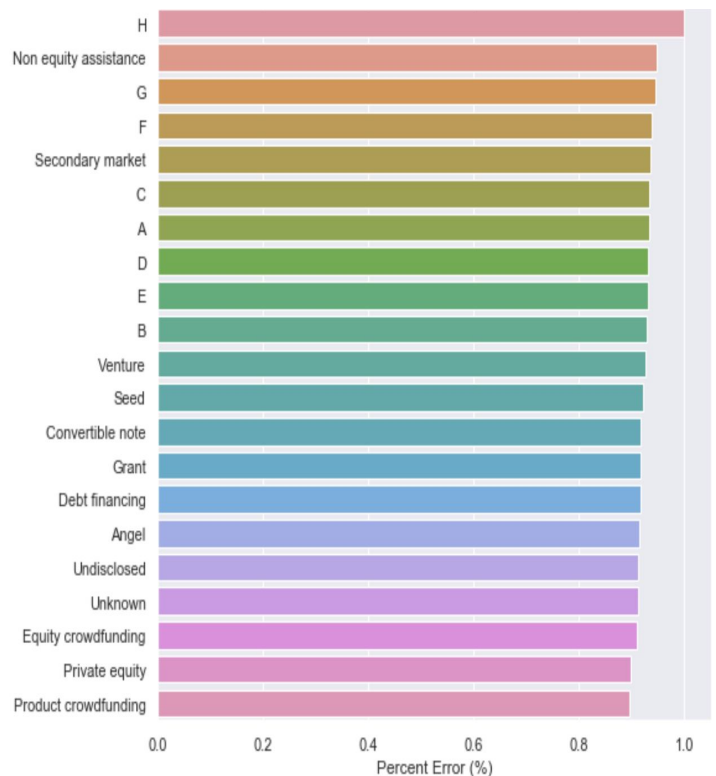
### Are We Happy With This?

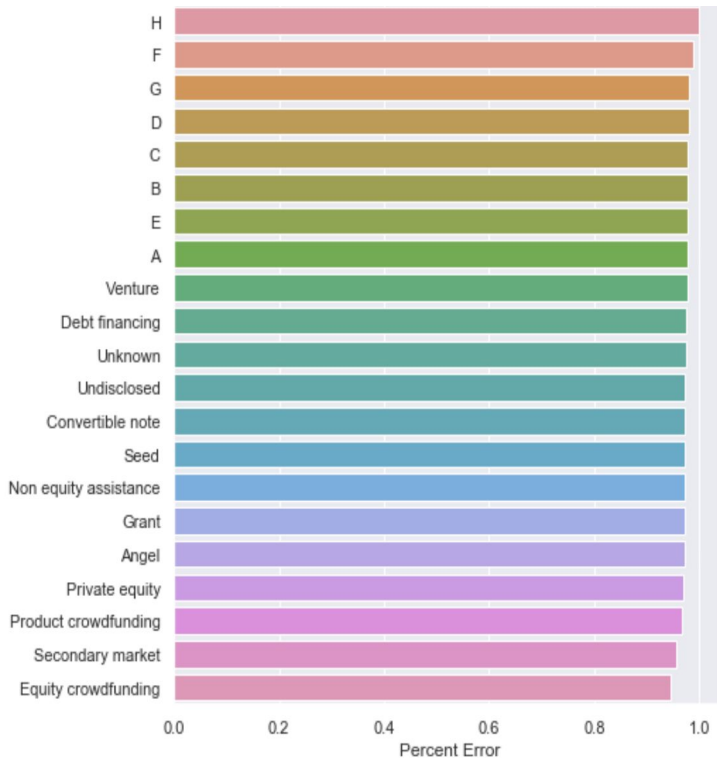
Higher recall would be ideal but in this case precision *is more important* than recall. **We would rather fund a company and be correct than fund all of the correct companies.** (or at least find)

	precision	recall	f1-score	support
0.0	0.95	0.98	0.96	380907
1.0	0.91	0.78	0.84	93510
micro avg	0.94	0.94	0.94	474417
macro avg	0.93	0.88	0.90	474417
weighted avg	0.94	0.94	0.94	474417

### Precision by Round Type

### Recall by Round Type





## Conclusion

Based on the data we had, we were able to build a model that predicts whether a company will get funded in the coming year at an accuracy of about 94%. The model had some difficulty with recalling all fiscal quarters of companies which would be funded in a year but the company quarters it did predict were correct 91% of the time. In this case precision is much more important than recall as for any use case we do not need to know ALL companies that are going to be funded but we do want to be sure that the companies we do think are worth funding are worth funding.

## References and

## Resources

### GitHub Repository

[https://github.com/daikiminaki/Springboard/tree/master/Capstone\\_Project\\_Crunchbase\\_Funding](https://github.com/daikiminaki/Springboard/tree/master/Capstone_Project_Crunchbase_Funding)

### GitHub Presentation

[https://github.com/daikiminaki/Capstone-Project-Crunchbase-Funding-Prediction-Model/blob/master/Capstone\\_1\\_Final\\_Presentation.pdf](https://github.com/daikiminaki/Capstone-Project-Crunchbase-Funding-Prediction-Model/blob/master/Capstone_1_Final_Presentation.pdf)

### Email Address

[dminaki95@gmail.com](mailto:dminaki95@gmail.com)