

IPO Investment via Text Mining

Introduction to NLP @ Open University 2018b, Final Project

Andrew Kreimer

kreimer.andrew@gmail.com

<https://github.com/algonell/IPOMiner>

Abstract

Let us introduce IPOMiner: Python utilities to predict future performance of upcoming IPO (Initial Public Offering). The project is a collection of data-sets and Python code to perform Text Mining on raw SEC S-1 filings. The goal of this project is to apply Text Mining tools and techniques to spot investment opportunities in upcoming IPO. The project is intended for traders and researchers as potential source for alpha generation.

1 Purpose

The project incorporates Text Mining tools and techniques to find potential trade opportunities in upcoming IPO (by buying or selling equity). The project transforms S-1 forms to buy/sell signals via summarization, sentiment analysis, keywords analysis and classification (Jurafsky and Martin, 2014; Barnes et al., 2017; Khan et al., 2015). This project is targeting only English text corpora and US equities markets.

The purpose of this project is to apply Text Mining tools and techniques in order to spot investment opportunities in upcoming IPO (Initial Public Offering). This project is interesting as we try to transform biased texts into actionable trading signals and bridge the gap between accountants and retail traders.

Public companies with good fundamentals encourage investments. Contrary, companies with bad fundamentals or hidden details about failing products and unclear future will try to obfuscate it in their earnings reports to hold investors in. IPO (Initial Public Offerings) are different as we have no earnings report or true intrinsic data, rather hype and biased news. The single most objective

document is the S-1 form created by accountants which is hard to understand.

2 Organization

The system is comprised of three main modules. The first module is responsible for IPO data retrieval. This module uses the EDGAR system provided by SEC (Securities and Exchange Commission) to download and store the raw S-1 documents (EDG; SEC). This module is also responsible for historical price data retrieval via Yahoo Finance (Yah). This module is mainly used for batch processing and should be run monthly.

The second module is responsible for Text Mining and model learning. The module loads, cleans, transforms and prepares our raw data for training. In addition this module learns and stores our predictive model: an IPO investment classifier. This module is also mainly used for batch processing and should be run monthly.

The third module is our production/real-time environment. The module is responsible for real-time classification of upcoming IPO. This module should be run daily or weekly. Our existing datasets are also available for enhancements via Feature Engineering (Kag).

All of our components rely on each other in a sequential order. The raw data retrieval output is an input to our Text Mining module which is an input to our real-time classification module. Data retrieval is our heaviest module and usually takes 2 days to finish on a single machine. A zipped raw data collection for more than 2400 S-1 filings (5.5GB) is included and cuts the running times for assets prior to 2018 (IPO).

3 Example Input

The input of our system is a collection of raw S-1 filings documents indexed by the EDGAR sys-

DROPBOX, INC.
Our Business
Our modern economy runs on knowledge. Today, knowledge lives in the cloud as digital content, and Dropbox is a global collaboration platform where more and more of this content is created, accessed, and shared with the world. We serve more than 500 million registered users across 180 countries.
Dropbox was founded in 2007 with a simple idea: Life would be a lot better if everyone could access their most important information anytime from any device. Over the past decade, we've largely accomplished that mission—but along the way we recognized that for most of our users, sharing and collaborating on Dropbox was even more valuable than storing files.
Our market opportunity has grown as we've expanded from keeping files in sync to keeping teams in sync. Today, Dropbox is well positioned to reimagine the way work gets done. We've focused on reducing the inordinate amount of time and energy the world wastes on "work about work"—tedious tasks like searching for content, switching between applications, and managing workflows.
We want to free up our users to spend more of their time on the work that truly matters. Our mission is to unleash the world's creative energy by designing a more enlightened way of working.
We believe the need for our platform will continue to grow as teams become more fluid and global, and content is increasingly fragmented across incompatible tools and devices. Dropbox breaks down silos by constraining the flow of information between the products and services our users prefer, even if they're not our own.
By solving these universal problems, we've become invaluable to our users. The popularity of our platform drives viral growth, which has allowed us to scale rapidly and efficiently. We've built a thriving global business with over 11 million paying users.
Our revenue was \$603.8 million, \$844.8 million, and \$1,106.8 million in 2015, 2016, and 2017, respectively, representing an annual growth rate of 40% and 31%, respectively. We generated net losses of \$322.9 million, \$210.2 million, and \$111.7 million in 2015, 2016, and 2017, respectively. We also generated positive free cash flow of \$137.4 million and \$305.0 million in 2016 and 2017, respectively, compared to negative free cash flow of \$63.5 million in 2015.

Figure 1: Dropbox S-1/A Filing (SEC)

Company Name	Symbol	Market	Price	Shares	Offer Amount	Date Permitted
AVALARA INC	AVLR	NYSE	\$24	7,500,000	\$180,000,000	6/15/2018
PUXIN LTD	NEW	NYSE	\$17	7,200,000	\$122,400,000	6/15/2018
VERRICA PHARMACEUTICALS INC.	VRCA	NASDAQ Global	\$15	5,000,000	\$75,000,000	6/15/2018
US XPRESS ENTERPRISES INC	USX	NYSE	\$16	18,056,000	\$288,896,000	6/14/2018
CHARAH SOLUTIONS, INC.	CHRA	NYSE	\$12	7,352,941	\$88,235,292	6/14/2018
FAR POINT ACQUISITION CORP	FPACU	NYSE	\$10	55,000,000	\$550,000,000	6/12/2018
GS ACQUISITION HOLDINGS CORP	GSAHU	NYSE	\$10	60,000,000	\$600,000,000	6/8/2018
MEIRAGTX HOLDINGS PLC	MGTX	NASDAQ Global Select	\$15	5,000,000	\$75,000,000	6/8/2018
AMBOW EDUCATION HOLDING LTD.	AMBO	NYSE MKT	\$4.25	1,800,000	\$7,650,000	6/1/2018

Figure 2: Upcoming IPO (NAS)

tem (EDG). S-1 filings are mandatory for private companies willing to go public. S-1 forms combine fundamental and financial data provided by accountants and should be objective and deterministic. Figure 1 shows an introduction of a raw filing.

The system uses raw IPO filings to create a unified data-set which contains multiple features and core meta-data. The system then builds a classifier for real-time prediction of upcoming IPO. Figure 2 shows a sample upcoming IPO at NASDAQ (NAS).

Finally, our system is used for real-time classification of upcoming IPO in order to spot potential investment opportunities. It is possible to buy or sell the underlying stock for various holding periods. Figure 3 shows two IPO of DBX and BTAI and their respective performance over a long-term holding period.

It is clear that DBX could be a profitable buy and hold, while BTAI could be a profitable short sell. The system provides classified directions and the underlying probabilities in order to discover upcoming IPO as potential trades similar to DBX and BTAI (tas; NAS).



(a) DBX



(b) BTAI

Figure 3: IPO Performance (tas)

boxers get a key-role meaning, although dropbox is a company name and dropboxers are essentially users. This issue is hard to spot systematically and must be addressed manually.

4.2 Machine Learning

Our project incorporates various Machine Learning tools and techniques. In terms of data preprocessing we use one-hot encoding to avoid ordered categorical values such as month, market and quarter. In addition, we incorporate standardization to normalize our data points as share amount and offering price vary across listings (Larose, 2006; Kag).

Once a clean and normalized data-set is available, we integrate Ensemble Classification. Our problem is a supervised binary classification task. Our ensemble is comprised of Logistic Regression and Random Forest classifiers (Larose, 2006; Kag). Ensemble Classification is a widely used methodology to improve simple classifiers generalization and boost performance (Kag).

4.3 Supervised Learning

Our problem is a supervised learning binary classification (Larose, 2006). As mentioned before, we collect four target variables (1D, 1W, 1M and 3M) performances based on open prices. We collect percent changes as open prices of respective period since the the first trading day, minus the first trading day opening price. The percentages are then converted to binary variables as positive performance and negative performance (1 and 0 respectively). Note that predicting performance as continuous values, which is known as regression problem in Machine Learning, has not been addressed in this project.

4.4 Randomness

Standard Machine Learning problems provide idealistic world problems in terms of perfect train-test splits and out of sample validation (Larose, 2006; Kag). Real-life applications on the other hand usually have lower generalization and predictive power, particularly in financial markets. The most common issue is the non-repeated data distributions and existence of market momentum (Narang, 2013; Johnson, 2010).

4.5 Hyper-Parameter Tuning

Hyper-parameter tuning is the process of fixing various parameters for our models within the

train-test split in order to increase generalization and predictive power. Hyper-parameter tuning is widely used in Machine Learning tasks (Kag). Contrary, in financial markets it is always advised to keep the models and learning as simple as possible (Kag; Narang, 2013; Johnson, 2010).

5 Discussion

Previous works have shown significant integration of news related to the IPO and their influence on performance. IPO S-1 filings tend to be diversified as various companies make them. The addition of more differences and non-unified data source increases noise. Eventually, the tone and biased jargon tend to influence retail investors (Feuerriegel et al., 2014).

In addition, previous works in the field show various applications for subsets of financial documents. For instance annual reports tend to present wider and more reliable company overview (Kloptchenko et al., 2004). On the other hand, news, tone and biased jargon tend to influence retail investors (Feuerriegel et al., 2014). Another important signal boosting of buying or selling is provided by 8-K documents (Lee et al., 2014). Finally, a brief correlation is described in IPO filings and higher management opinions (Deokar and Tao, 2015).

5.1 Management Data Integration

A research has shown the ability to integrate management opinions and summaries to evaluate upcoming IPO performance (Deokar and Tao, 2015). The research showed how genuine information about the companies can be extracted by combining final IPO statements and management discussion analysis. The combination provides boosted sentiment signal and better predictive modeling. In our project we have been relying only on raw IPO S-1 filings and no external related data.

5.2 Continuous Sentiment

A research of social media impact and performance of equities have shows a significant momentum effect (Makrehchi et al., 2013). The research analyzed aggregated twitter feeds grouped by equities performance (positive and negative) to reveal patterns. They showed that drastic performance of companies tend to continue the following trading days. Regarding IPO investments the moment effect is even more significant and thus

can be another potential data source for our problem.

Another research showed text analysis of 8-K documents and the ability to classify performance. 8-K documents are legally required documents for major corporate changes such as bankruptcy and CEO changes. Text analysis of major financial data shows a high correlation to discrete class variable (going up, going down and staying the same) for the proceeding week (Lee et al., 2014). In our project we have not been incorporating social media signals nor 8-K documents.

5.3 Hype Analysis

Previous works in the fields have been analyzing hype and biased impact of financial news on equity performance (Johnson, 2010; Narang, 2013). A research showed how biased and manipulated language can mislead investors and intensify post IPO performance. In addition, the researched showed how language and jargon can tilt risk to reward ratios as seen by investors (Feuerriegel et al., 2014). Another research showed how opinion mining can help businesses to improve and increase sells. An important aspect of this research is messages and comments analysis regarding the topic explored. The research shows how user opinion is important for boosting our predictive modeling (Chen and Zimbra, 2010).

5.4 Language Specificity

A similar research has been conducted for Turkish companies (Basti et al., 2015). The researched showed how listing data and various over-pricing and under-pricing indicators influence first week performance of IPO. Turkish jargon and format allowed better predictive modeling. Similarly to our research the analysis is based on a single market place. On the other hand, the research is dedicated to a small market with mainly Turkish companies, whereas the US market is a major market place for companies all over the world (Johnson, 2010).

5.5 Cluster Analysis

In contradiction to previous works in the filed, we have tried to incorporated clustering features. Clustering via KNN (K Nearest Neighbours) is widely used strategy for feature improvement (Kag; Larose, 2006). By analyzing our data-set and looking for nearest neighbours we can add new sets of features to existing features (horizontal concatenation). This enhancement adds more

		1D	1W	1M	3M
AUC	LR	0.516569	0.516569	0.516569	0.516569
	RF	0.538337	0.538337	0.538337	0.538337
f1	LR	0.735849	0.735849	0.735849	0.735849
	RF	0.726368	0.726368	0.726368	0.726368
log loss	LR	0.663721	0.663721	0.663721	0.663721
	RF	0.657215	0.657215	0.657215	0.657215

(a) With Clustering Features

		1D	1W	1M	3M
AUC	LR	0.544509	0.544509	0.544509	0.544509
	RF	0.578622	0.578622	0.578622	0.578622
f1	LR	0.732673	0.732673	0.732673	0.732673
	RF	0.735751	0.735751	0.735751	0.735751
log loss	LR	0.655687	0.655687	0.655687	0.655687
	RF	0.640521	0.640521	0.640521	0.640521

(b) Without Clustering Features

Figure 7: Clustering Features (Pedregosa et al., 2011)

meta-data for existing features. Although this methodology usually improves classifiers, in our problem it did not help to improve. Figure 7 shows our experiments with and without clustering features.

5.6 Word2vec Integration

Our best improvement was achieved by incorporating keywords analysis and cross distances between companies. As we showed, simple cluster analysis is too wide to provide some kind of signal. Word2vec models essentially transform words to n-dimensional matrices (Jurafsky and Martin, 2014). The sparse representation better describes words relations. We then apply a series of aggregations on small keywords groups and eventually create better relationship features between companies and their relative performance. Figure 8 shows the combined Word2Vec build from all IPO keywords.

5.7 Multiple Target Variables

The previous works that has been covered usually incorporate single class variable and evaluate the entire performance on a single outcome. In this work we classify performance and provide the underlying probabilities for multiple periods. The combination of multiple products in different periods of expiration can out-perform single direc-

		1D	1W	1M	3M
AUC	LR	0.479237	0.479237	0.479237	0.479237
	RF	0.512065	0.512065	0.512065	0.512065
f1	LR	0.696864	0.696864	0.696864	0.696864
	RF	0.695971	0.695971	0.695971	0.695971
log loss	LR	0.708928	0.708928	0.708928	0.708928
	RF	0.7011	0.7011	0.7011	0.7011

(a) Baseline

		1D	1W	1M	3M
AUC	LR	0.517241	0.517241	0.517241	0.517241
	RF	0.553196	0.553196	0.553196	0.553196
f1	LR	0.745455	0.745455	0.745455	0.745455
	RF	0.746411	0.746411	0.746411	0.746411
log loss	LR	0.670979	0.670979	0.670979	0.670979
	RF	0.663169	0.663169	0.663169	0.663169

(b) Sentiment Analysis

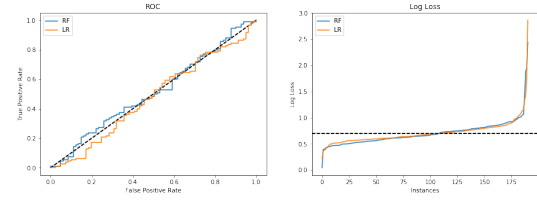
		1D	1W	1M	3M
AUC	LR	0.539319	0.539319	0.539319	0.539319
	RF	0.530908	0.530908	0.530908	0.530908
f1	LR	0.681081	0.681081	0.681081	0.681081
	RF	0.744186	0.744186	0.744186	0.744186
log loss	LR	0.676177	0.676177	0.676177	0.676177
	RF	0.661292	0.661292	0.661292	0.661292

(c) Summarization

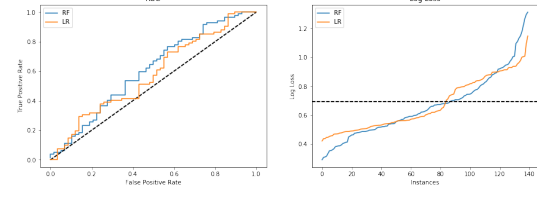
		1D	1W	1M	3M
AUC	LR	0.544509	0.544509	0.544509	0.544509
	RF	0.578622	0.578622	0.578622	0.578622
f1	LR	0.732673	0.732673	0.732673	0.732673
	RF	0.735751	0.735751	0.735751	0.735751
log loss	LR	0.655687	0.655687	0.655687	0.655687
	RF	0.640521	0.640521	0.640521	0.640521

(d) Keywords

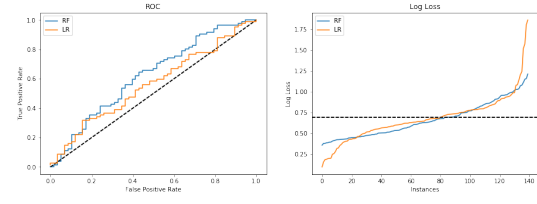
Figure 11: AUC, f1 and log loss (Pedregosa et al., 2011)



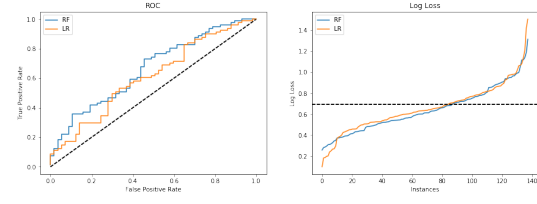
(a) Baseline



(b) Sentiment Analysis



(c) Summarization



(d) Keywords

Figure 12: ROC and logarithmic loss for instances (Pedregosa et al., 2011)

Operating Characteristic) and logarithmic loss for all instances in the test-set. ROC plots show deeper relations between precision and recall by providing various thresholds and respective precision/recall values. Classifiers use 0.5 as a default cut-off point for class labeling, which sometimes make the decision not reliable, particularly in our problem with $AUC < 0.6$. The collection of logarithmic loss values for our test-set shows the general confidence of our predictions and the respective probabilities. Figure 12 shows ROC and multiple log loss values for our test-set. Note the slight improvements over experiments both in terms of ROC and higher AUC and more instances having log loss values less than 0.693 (a coin toss).

6.4 Feature Importance

Finally, we present feature importance evaluation for each experiment. We have been incorporating

Figure 13: Feature Importance (Chen and Guestrin, 2016)

ing historical price data and calculating target variables (1D, 1W, 1M and 3M performance). Another long running batch processes are summarization and keywords extraction which required GCE (Google Compute Engine) assistance due to out of memory issues with particularly large IPO.

8.2 Amendments

Most S-1 filings are followed by S-1/A amendments (and sometimes up to 10 more amendments). This data has not been incorporated in this project. We always chose the most up to date filing. Changes in S-1 filings could improve the long term inspection and opinion of the exchange about the company.

8.3 Repetitive Words

S-1 filings have common keywords that appear over 300 times in most of the documents. The following is a short list of the common words.

share	common
stock	marketing
company	development
month of IPO	previous year to IPO
believe	

Those words can mislead to some importance of particular months. As shown in the baseline model, month or quarter of IPO has no predictive power. Another problem that accountants add many positive words but usually with negative numbers. The right approach would be standardizing all of the S-1 filings to remove the common words, although it is computationally hard.

8.4 Numerical Data

Special care must be taken dealing with numerical data in our S-1 filings. Tables with current cash flow, debt, holdings and company effective risk to reward are hard to integrate. Tables tend to diverse with different colors, styles and shapes. This data must be incorporated manually.

8.5 Domain Knowledge

Many raw data files have implied terminology related to the relevant IPO. In the Dropbox example we have seen dropboxers extracted as a noun or meaningful keyword although essentially it should be replaced with users. It will be interesting to incorporate domain knowledge (both financial and accounting) to our models such as the NLTK Reuters corpus for financial news analysis (NLP).

References

- 5 heroic tools for natural language processing. <https://goo.gl/UwECfQ>.
- Electronic data gathering, analysis, and retrieval system. <https://www.sec.gov/edgar.shtml>.
- Free online word cloud generator and tag cloud creator. <https://www.wordclouds.com/>.
- Ipominer. <https://github.com/algonell/IPOMiner>.
- Kaggle. <https://www.kaggle.com/>.
- National association of securities dealers automated quotations. <https://www.nasdaq.com/>.
- tastytrade. <https://www.tastytrade.com/>.
- U.s. securities and exchange commission. <https://www.sec.gov/>.
- Yahoo finance. <http://finance.yahoo.com/>.
- Mehdi Allahyari, Seyedamin Pouriyeh, Mehdi Assefi, Saeid Safaei, Elizabeth D Trippe, Juan B Gutierrez, and Krys Kochut. 2017. Text summarization techniques: A brief survey. *arXiv preprint arXiv:1707.02268*.
- Jeremy Barnes, Roman Klinger, and Sabine Schulte im Walde. 2017. Assessing state-of-the-art sentiment models on state-of-the-art sentiment datasets. *arXiv preprint arXiv:1709.04219*.
- Eyup Bastı, Cemil Kuzey, and Dursun Delen. 2015. Analyzing initial public offerings' short-term performance using decision trees and svms. *Decision Support Systems*, 73:15–27.
- Hsinchun Chen and David Zimbra. 2010. Ai and opinion mining. *IEEE Intelligent Systems*, 25(3):74–80.
- Tianqi Chen and Carlos Guestrin. 2016. [Xgboost: A scalable tree boosting system](#). *CoRR*, abs/1603.02754.
- Amit Deokar and Jie Tao. 2015. Text mining for studying managements confidence in ipo prospectuses and ipo valuations.
- Stefan Feuerriegel, Jonas Schmitz, and Dirk Neumann. 2014. What matters most? how tone in initial public offering filings and pre-ipo news influences stock market returns.
- Barry Johnson. 2010. *Algorithmic Trading & DMA: An introduction to direct access trading strategies*, volume 200. 4Myeloma Press London.
- Dan Jurafsky and James H Martin. 2014. *Speech and language processing*, volume 3. Pearson London:.

- Aamera ZH Khan, Mohammad Atique, and VM Thakare. 2015. Combining lexicon-based and learning-based methods for twitter sentiment analysis. *International Journal of Electronics, Communication and Soft Computing Science & Engineering (IJECSCE)*, page 89.
- Antonina Kloptchenko, Tomas Eklund, Jonas Karlsson, Barbro Back, Hannu Vanharanta, and Ari Visa. 2004. Combining data and text mining techniques for analysing financial reports. *Intelligent systems in accounting, finance and management*, 12(1):29–41.
- Daniel T Larose. 2006. *Data mining methods & models*. John Wiley & Sons.
- Heeyoung Lee, Mihai Surdeanu, Bill MacCartney, and Dan Jurafsky. 2014. On the importance of text analysis for stock price prediction. In *LREC*, pages 1170–1175.
- Elena Lloret and Manuel Palomar. 2010. Challenging issues of automatic summarization: relevance detection and quality-based evaluation. *Informatica*, 34(1).
- Masoud Makrehchi, Sameena Shah, and Wenhui Liao. 2013. Stock prediction using event-based sentiment analysis. In *Proceedings of the 2013 IEEE/WIC/ACM International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT)-Volume 01*, pages 337–342. IEEE Computer Society.
- Rishi K Narang. 2013. *Inside the Black Box: A Simple Guide to Quantitative and High Frequency Trading*, volume 883. John Wiley & Sons.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Radim Řehůřek and Petr Sojka. 2010. Software Framework for Topic Modelling with Large Corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pages 45–50, Valletta, Malta. ELRA. <http://is.muni.cz/publication/884893/en>.
- Josef Steinberger and Karel Ježek. 2012. Evaluation measures for text summarization. *Computing and Informatics*, 28(2):251–275.