# CS295P: Stock Prediction Project

Group S10
Hanyan Wang, Zhihui Xia, Shuqing Ye

## Introduction

Stock market prediction is the act of trying to determine the future value of a company stock or other financial instrument traded on an exchange. The successful prediction of a stock's future price could yield significant profit. In this project we do **fundamental analysis** on the company profiles and apply machine learning models on the data to find out which features are more important. Then the machine learning will generate a score on each company depending on their data. Our portfolio will pick the top 20 companies with the highest score.

## Data

The training data is collected from Yahoo Finance during the last quarter of 2001, 2006, 2011.
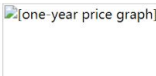
## Approach

Our project is divided into three parts, Parsing, Prediction, and Trade/Adversary. The parsing stage will be responsible for reading the HTML, extracting the data, and generating the database for machine learning. The prediction section will run machine learning models on the database and give a portfolio with 20 best companies. Using the generated trading file, the Adversary will display the result of our portfolio.

### Parsing

#### 1.HTML to CSV

The raw data is a pile of HTML files. We use Jsoup, a Java library, and CSSselector to extract and manipulate data.
Take Apple's profile in 2001-11 for example. It shows much useful information which can be used to do fundamental analysis. We extracted these elements from HTML files, removed unnecessary symbols, converted them into string objects and stored them in csv files.

| Statistics at a Glance -- NasdaqNM:AAPL | | | | | As of 31-Oct-2001 |
| --- | --- | --- | --- | --- | --- |
| **Price and Volume** | | **Per-Share Data** | | **Management Effectiveness** | |
| 52-Week Low on 20-Dec-2000 | $13.625 | Book Value (mrq*) | $11.17 | Return on Assets (ttm) | -0.60% |
| Recent Price | $17.56 | Earnings (ttm) | -$0.14 | Return on Equity (ttm) | -0.96% |
| 52-Week High on 30-Apr-2001 | $27.12 | Earnings (mrq) | $0.18 | **Financial Strength** | |
| Beta | 1.31 | Sales (ttm) | $15.26 | Current Ratio (mrq*) | 3.39 |
| Daily Volume (3-month avg) | 4.99M | Cash (mrq*) | $12.36 | Debt/Equity (mrq*) | 0.08 |
| Daily Volume (10-day avg) | 6.34M | **Valuation Ratios** | | Total Cash (mrq) | $4.34B |
| **Stock Performance** | | Price/Book (mrq*) | 1.57 | **Short Interest** As of 10-Sep-2001 | |
| [one-year price graph] | | Price/Earnings | N/A | Shares Short | 6.64M |
| | | Price/Sales (ttm) | 1.15 | Percent of Float | 2.1% |
| big chart [1d \| 5d \| 3m \| 6m \| 1y \| 2y \| 5y \| max] | | **Income Statements** | | Shares Short (Prior Month) | 6.30M |
| | | Sales (ttm) | $5.36B | Short Ratio | 1.44 |
| | | EBITDA (ttm*) | -$344.0M | Daily Volume | 4.61M |
| 52-Week Change | -14.3% | Income available to common (ttm) | -$37.0M | | |
| 52-Week Change relative to S&P500 | +14.9% | **Profitability** | | | |
| **Share-Related Items** | | Profit Margin (ttm) | -0.7% | | |
| Market Capitalization | $6.16B | Operating Margin (ttm) | -6.4% | | |
| Shares Outstanding | 350.9M | **Fiscal Year** | | | |
| Float | 312.3M | Fiscal Year Ends | Sep 29 | | |
| **Dividends & Splits** | | Most recent quarter (fully updated) | 29-Sep-2001 | | |
| Annual Dividend | none | Most recent quarter (flash earnings) | 30-Sep-2001 | | |
| Last Split: factor 2 on 21-June-2000 | | | | | |

*See Profile Help for a description of each item above.* **M** = millions; **B** = billions; **mrq** = most-recent quarter; **ttm** = trailing twelve months; (as of 30-Sep-2001, except **mrq\*/ttm\*** items as of 29-Sep-2001)

In our implementation, the entry for each company consists of 41 attributes, including 52-Week Low (lowest price within 52 weeks), recent price, beta, market capitalization, DCV (daily cash volume), etc.

Considering the fact that some companies with dots in their names are frequently lacking fundamental data, since they are either foreign stocks or don't trade or work the same way as other companies on the major exchanges, we decided to skip these companies and not to put their data into our database.

## 2. Meric of a Company's performance

We use the formula below to measure a company (stock)'s performance, over an identical period of time:

$$Company\ Return - Market\ Return$$

`Company Return` is a company's monthly rate of return; `Market Return` indicates by how many percent the S&P 500 index changes during that month. By comparing the company with the market, we can evaluate its performance to see if it outperforms its peers, regardless of what its raw return is.

## 3. Output

The figure below shows the fundamental data of every company in 2001-10. Some cleaning and processing will be handled in the next phase.

# Prediction

## 1. Data Processing

We read 9 months data from 2001, 2006 and 2011 and concatenate them. We only select companies with daily cash volume larger than one million.

In the raw data, there are empty cells and special symbols like 'M', 'K', '%', we need to replace empty cells with NaN and convert special symbols to floating-point numbers. After converting, we get something like below:

| | 52-Week Low | Recent Price | 52-Week High | Beta | Daily Volume (3-month avg) | Daily Volume (10-day avg) | 52-Week Change | 52-Week Change relative to S&P500 | C |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 25.875 | 33.17 | 45.71 | 1.17 | 3420000.0 | 2880000.0 | 15.6 | 55.9 | |
| 5 | 13.625 | 16.20 | 27.12 | 1.31 | 5650000.0 | 6490000.0 | −25.5 | −1.7 | |
| 6 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |
| 9 | 9.460 | 12.55 | 18.25 | 0.98 | 183200.0 | 197000.0 | −6.6 | 23.2 | |
| 12 | 7.000 | 9.16 | 14.20 | 1.54 | 391400.0 | 372000.0 | −1.0 | 30.7 | |

Our label is called `delta return`, which is computed as stock return - market return in the corresponding month. We plot the distribution of `delta return` as follows, blue curve

represents our return and black curve represents normal distribution. We want to convert our distribution to the normal one since it helps the training process.



delta return distribution vs Normal distribution

We drop rows and columns with too many missing values. Finally, we fill these missing values with simple imputer to avoid errors in training.

```
#fill missing values

from sklearn.impute import SimpleImputer
my_imputer = SimpleImputer()
train_filled = pd.DataFrame(my_imputer.fit_transform(train))
train_filled.columns=train.columns
train_filled.index=train.index
```

After preprocessing, our data has 23743 entries and 33 columns.

## 2. Model

We use three models to train on our data, namely CatBoost, Decision Tree and Linear Regression. 5 fold cross validation is applied on the data to determine the best model and root mean squared error (RMSE) is set as the loss function.

The following table shows the result for training, the second column is the mean of RMSE and the third column is the stand deviation of RMSE. CatBoost works the best among three and we will use it to train the whole dataset.

| | | |
|---|---|---|
| CatBoost | 0.196775 | 0.019617 |
| DecisionTree | 0.287104 | 0.020307 |
| LinearRG | 0.213732 | 0.033751 |

After training on nine months data in total, we plot the top 20 important features that contribute to the model. All features contribute more or less to the model.

| | Feature Id | Importances |
|---|---|---|
| 0 | 52–Week High | 7.383566 |
| 1 | 52–Week Change | 6.866093 |
| 2 | Short Ratio | 5.515618 |
| 3 | Daily Volume (3–month avg) | 5.287328 |
| 4 | 52–Week Change relative to S&P500 | 5.074821 |
| 5 | Beta | 4.859244 |
| 6 | DCV | 4.667107 |
| 7 | Book Value | 4.391858 |
| 8 | Return Assets | 4.025941 |
| 9 | Shares Outstanding | 3.885802 |
| 10 | Daily Volume (10–day avg) | 3.624967 |
| 11 | Percent of Float | 3.428175 |
| 12 | Price/Sales | 3.411327 |
| 13 | Float | 3.080421 |
| 14 | Shares Short | 2.894444 |
| 15 | Shares Short (Prior Month) | 2.686684 |
| 16 | Current Ratio | 2.552755 |
| 17 | Recent Price | 2.457203 |
| 18 | 52–Week Low | 2.411475 |
| 19 | Cash | 2.369788 |

## 3. Output

We save the trained model under the current directory. When there is a need to predict, we load the model and preprocess the test csv produced by the parser. It is worth noting that this test csv doesn't have label information. We predict the top 20 stocks and export them as "year-month-portfolio.csv" based on the given year and month.

# Trades and Adversary

## 1. Trades

We start from $100,000 and distribute the money evenly over the 20 companies. As a long-term trader, we only trade on a monthly basis -- buy at the beginning of a month, and sell at the end

of that month. To make things easier, each transaction is as close to the market close time, i.e. 16:00, as possible. Here is an example of our trading:

```
2001-10-02 15:59 buy 2647 shares of ARBA
2001-10-02 15:59 buy 356 shares of AVGN
......
2001-10-31 15:59 sell 2647 shares of ARBA
2001-10-31 15:59 sell 356 shares of AVGN
```

## 2. Adversary

We also build an Adversary to mimic our "adversary" in the real market, such as brokers, the exchange and other actors, and evaluate the proposed trades. Adversary is  implemented based on the following rules:

A. All orders are market orders that are filled immediately at the current ask for buys, and at the current bid for sells.
B. You will only be able to buy at most 1% of the DCV at the current ask price. After that, the price goes up by at least X for each 1% of the DCV, where x = bid-ask spread / 2.
C. Similarly, if the trader tries to sell more than 1% of the DCV, then the price the trader gets should decrease by X for each 1% of DCV.
D. If there is no transaction at the time requested for the trade, choose the line with the closest line after. If there is no line after, go for the closest line before.
E. If the bid and ask is N/A, approximate the price as follows. If the whole day daily cash volume is DCV, β is the bid-ask spread, log is the natural logarithm, then the following is a half decent approximation to their relationship:
$$\frac{-5log(\beta)}{log(DCV)} \approx 1$$
Thus the current asking price = current price + β, current bidding price = current price - β.

# Testing Result

For the testing process, we use the data of two years to train the model, and then apply the prediction on the data of the third year. For example, if we wanted to make predictions on 2001-10, we will use data from 2006 to 2011 as training data, then output the portfolio of 2001-10.
Here is our result in each month of Quarter 4 over 3 years, compared with monthly market return.

| Month | Our adversary | Prof's adversary | Portfolio Return | Market Return |
|-------|---------------|------------------|------------------|---------------|
| 2001-10 | $121171.44 | $121583.91 | 21.6% | 0.8% |
| 2001-11 | $100179.65 | $100125.99 | 0.13% | 5.11% |
| 2001-12 | $116653.12 | $117097.12 | 17.1% | 1.61% |
| 2006-10 | $93481.44 | $93807.11 | -6.2% | 3.50% |
| 2006-11 | $96321.79 | $96620.38 | -3.4% | 2.40% |
| 2006-12 | $97162.27 | $97676.89 | -2.3% | 1.55% |
| 2011-10 | $126235.41 | $128915.40 | 28.9% | 14.02% |
| 2011-11 | $105740.69 | $103531.11 | 3.53% | 2.35% |
| 2011-12 | $102122.41 | $102373.86 | 2.37% | 1.05% |

Generally speaking, our portfolio's performance is great. Among 9 months from 2001 to 2011, we obtained positive returns in 6 months. It is worth noting that in October and December 2001, and October 2011 our return is significantly higher than the market, especially in 2001-10, we achieved a 21.6% return, which is 27 times of the market return in that month.

But our portfolio cannot always outperform the market. One major reason is that certain stocks we bought at the beginning of one month, due to delisting, were not trading any more at the end of that month, and tagged $-1 by Adversary, thus reducing our earnings by a large margin. This accounts for the low income in November 2001 and loss in October 2006.

In all, according to the results of Adversary, our portfolio can beat the market and get excess returns.

Please NOTE that for the final prediction application, we will train the data from three years, and generate the final model. Therefore, the result from our final application will be different from the above table.

# Appendix

## Language and Environment

Python 3.6
Java 1.8

Openlab
Google Colab (model training)

## Dependencies

Jsoup >= 1.13.1
OpenCSV >= 5.3
shap >= 0.39.0
scikit-learn >= 0.24.1
Catboost >= 0.24.4
pandas >= 1.1.5
numpy>=1.13.3