# Statistical Arbitrage on S&P 500 Components by Machine Learning Models

Springboard Data Science Career Track Capstone 1, April 27th 2020 Cohort

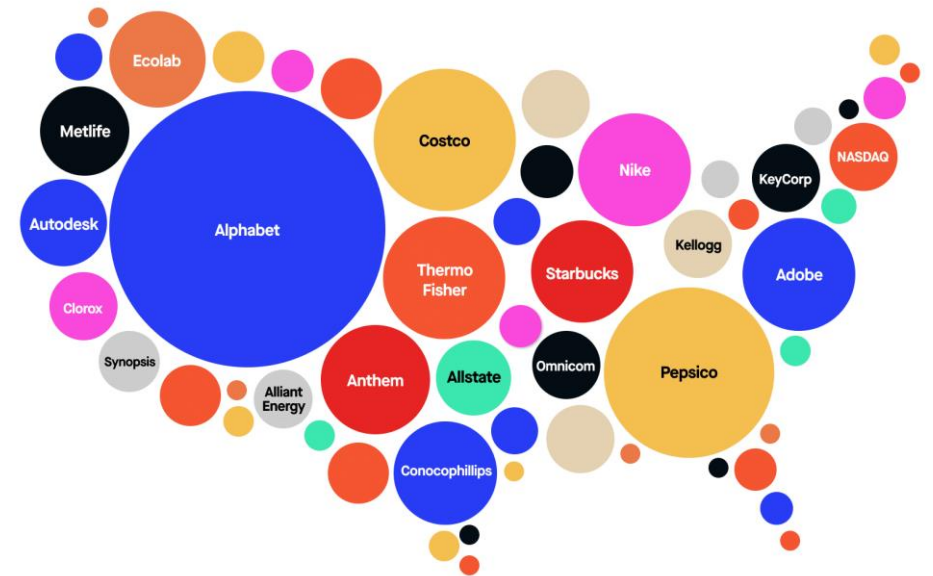Jiaqi Xu

09/23/2020

Thanks for Springboard mentor Kevin Glynn

# The Next 20 minutes is dictated to

- Project Introduction

- Dataset Extract, Transform and Load (ETL)

- Exploration Data Analysis (EDA)

- Machine Learning Models, Parameters Tuning and Prediction Results

- Trading Results and Analysis

- Future Work

# Why S&P 500 Components

- Attractive return - 9.8%/year, including dividends, since inception in 1926 (Wikipedia).

- Mr. Warren Buffet's recommendation in annual meeting 2020.

- A hot topic.

Source: Robinhood Learn

# Previous Related Work

- **Jegadeesh and Titman (1993)** - Early paper for momentum trading. Select stocks: past 6-month returns, holds: 6 months. Excess return: 12.01%/year (1965 - 1989).

- **Krauss et al. (2017)** - A statistical arbitrage strategy: deep neural networks (DNN), gradient-boosted trees (GBT), random forests (RAF) - S&P 500 (1992 – 2015). Return: 0.25%/day, Sharpe ratio: 1.81/year (no transaction cost).

- **Fischer and Krauss (2018)** - Long short-term memory (LSTM) networks: S&P 500 constitutes. Return: 0.46%/day, Sharpe ratio: 5.8/year (no transaction cost).

- **Fischer et al. (2019)** - Logistic regression (LG), RAF: 40 cryptocurrency coins minute (06/18/2018 to 09/17/2018). Return: 0.038%/round-trip trade (after transaction cost: 0.15%/half-turn).

# Goal of This Project

- Optimal machine learning methods: one-step stock movement prediction, updated dataset (2005-01-01 to 2019-12-31) .

- My machine learning methods vs Krauss et al. (2017) and Fischer and Krauss (2018).

- Check: "profits are declining in recent years and there is a severe challenge to the semi-strong form of market efficiency".

# Dataset Extract, Transform and Load (ETL)

- List of historical S&P 500 components - Wharton Research Data Services (WRDS) and Wikipedia.

- Daily data of historical S&P 500 components - Python Yahoo finance module yahoo-finance-1.0.4.

- Missing Components – Delisted from Yahoo finance.

# Exploration Data Analysis (EDA)

- **Dataset**: S&P 500 composites from 2005/01/01 to 2020/01/01 (15 years data).

- **Subsets**: Split dataset into 12 subsets - each subset, 750 days/250 days – formation/ trading. Trading: non-overlapped. Formation: 80%/20% training (5 folder cv)/testing.

| Subset | form_start | form_start | trad_end | Subset | form_start | form_start | trad_end |
|---|---|---|---|---|---|---|---|
| 0 | 1-1-2005 | 1-1-2008 | 1-1-2009 | 6 | 1-1-2011 | 1-1-2014 | 1-1-2015 |
| 1 | 1-1-2006 | 1-1-2009 | 1-1-2010 | 7 | 1-1-2012 | 1-1-2015 | 1-1-2016 |
| 2 | 1-1-2007 | 1-1-2010 | 1-1-2011 | 8 | 1-1-2013 | 1-1-2016 | 1-1-2017 |
| 3 | 1-1-2008 | 1-1-2011 | 1-1-2012 | 9 | 1-1-2014 | 1-1-2017 | 1-1-2018 |
| 4 | 1-1-2009 | 1-1-2012 | 1-1-2013 | 10 | 1-1-2015 | 1-1-2018 | 1-1-2019 |
| 5 | 1-1-2010 | 1-1-2013 | 1-1-2014 | 11 | 1-1-2016 | 1-1-2019 | 1-1-2020 |

Table 1: Formation and Trading periods for 12 non-overlapping trading batches

# Exploration Data Analysis (EDA)

## Feature generation:

- *Input:* $P^s = (P^s_t)_{t \in T}$ - price process of stock s, $s \in \{1, \ldots, n\}$. Simple return $R^s_{t,m}$ over m periods: $R^s_{t,m} = \frac{P^s_t}{P^s_{t-m}} - 1$,

Where $m \in \{\{1, \ldots, 20\} \cup \{40, 60, \ldots, 240\}\}$.

- *Output:* A binary response variable $Y^s_{t+j,j} \in \{0,1\}$ for each stock s.

$$Y^s_{t+j,j} = \begin{cases} 1, & \text{if } R^s_{t+j,j} > \text{median return} \\ 0, & \text{otherwise} \end{cases}$$
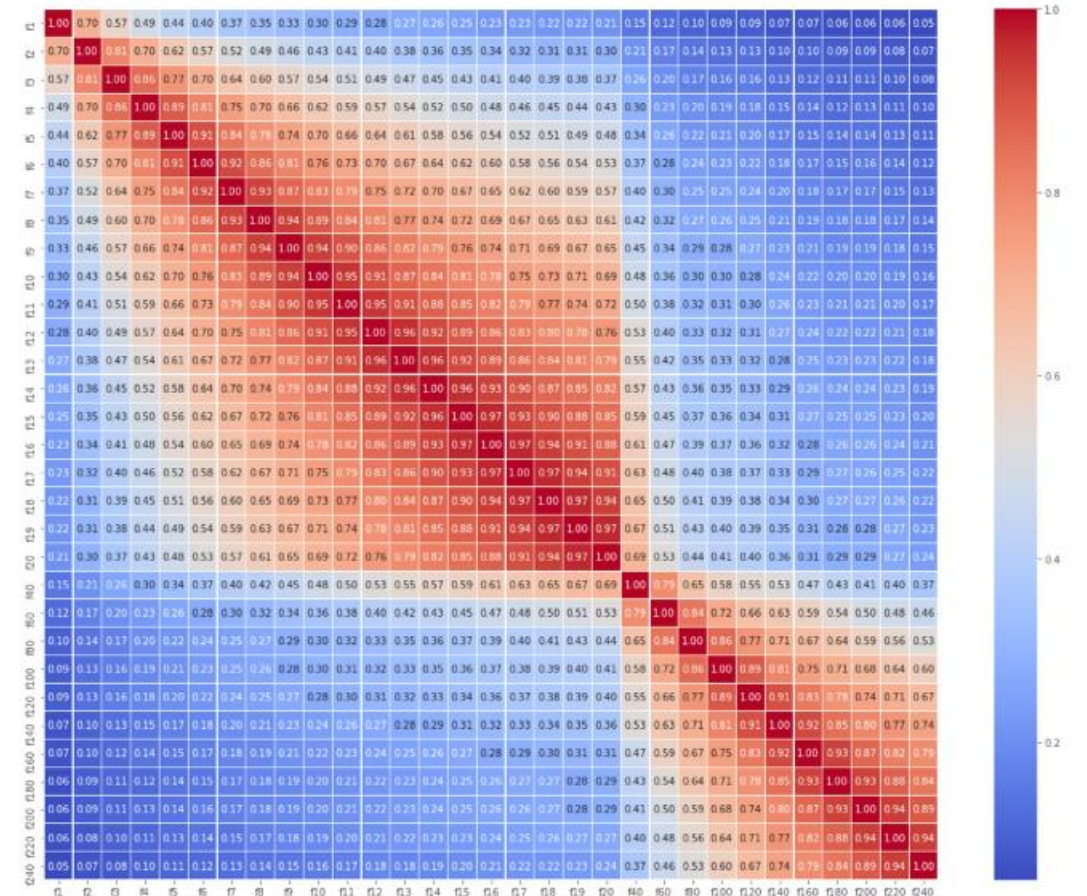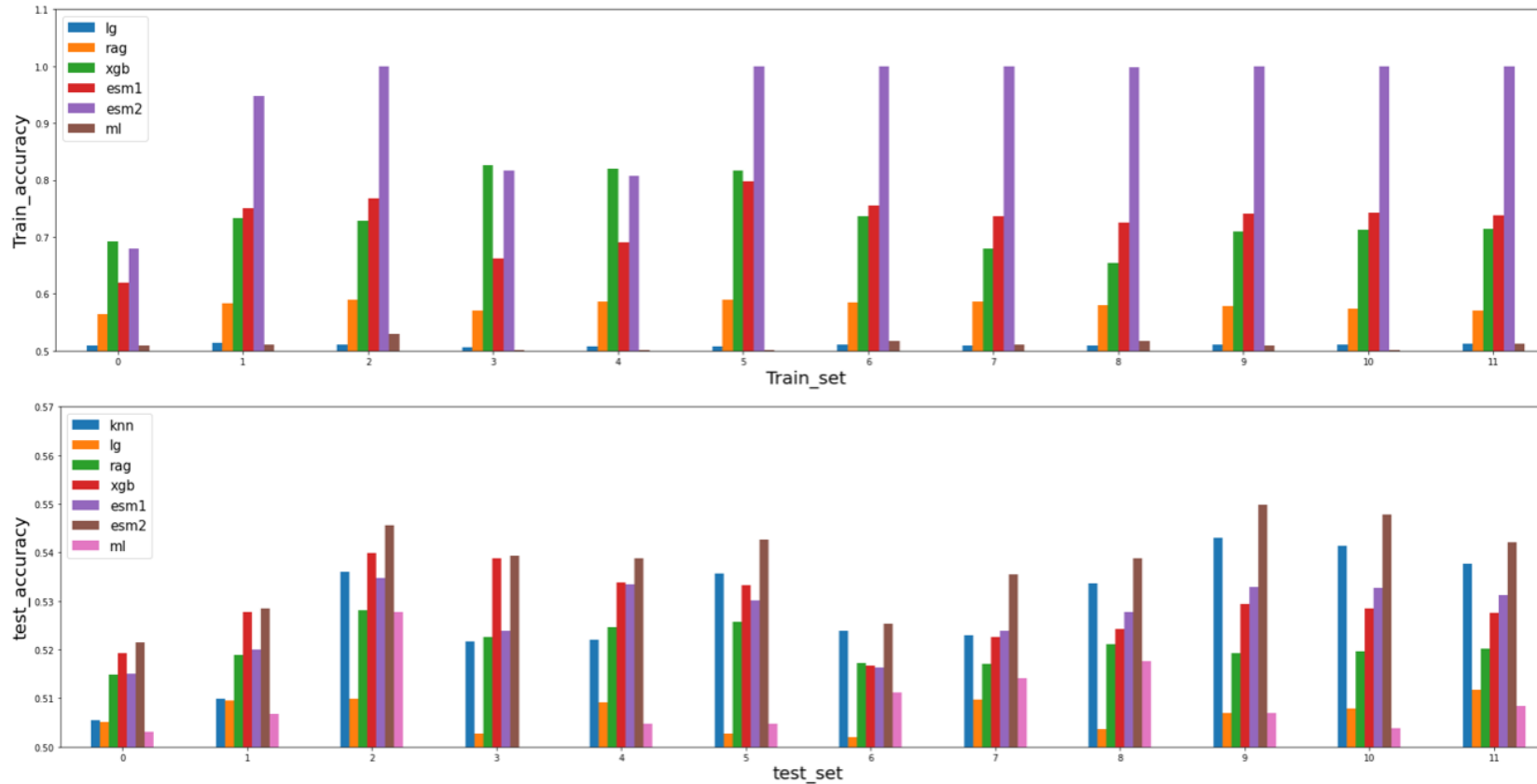
j=1 for 1 period.

## Feature Description:



Figure 1. Feature Correlation Heatmap from Formation Subset 10

# Machine Learning Methods, Grid-search parameters, Results

| Methods | Parameters to be tuned | Be selected values | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | \multicolumn Optimal values for Dataset | | | | | | | | | | | |
| knn | Number of neighbors | {3,5,7,11} | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 5 | 3 | 5 | 5 | 11 |
| | Weight function used in prediction | {'uniform','distance'} | 'uniform' | 'uniform' | 'distance' | 'distance' | 'uniform' | 'uniform' | 'distance' | 'distance' | 'distance' | 'distance' | 'distance' | 'distance' |
| | Metric | {'euclidean', 'manhattan'} | 'manhattan' | 'manhattan' | 'euclidean' | 'euclidean' | 'manhattan' | 'manhattan' | 'euclidean' | 'euclidean' | 'euclidean' | 'euclidean' | 'euclidean' | 'euclidean' |
| raf | The number of trees | {1500,2000,2500} | 2500 | 2500 | 2500 | 2000 | 2500 | 2500 | 2500 | 2500 | 2500 | 2500 | 2500 | 2500 |
| | The maximum depth of the tree | {20,25,30} | 20 | 20 | 30 | 30 | 25 | 25 | 30 | 30 | 25 | 25 | 20 | 20 |
| | The number of features for the best split | {log2(n_features), sqrt(n_features)} | "sqrt" | "sqrt" | "log2" | "log2" | "sqrt" | "sqrt" | "log2" | "log2" | "log2" | "log2" | "log2" | "log2" |
| | The number of samples to draw in bootstrap | {500, 700, 1000} | 700 | 700 | 1000 | 1000 | 700 | 700 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 |
| lg | Inverse of regularization strength | {10^-3, 10^-2,…,10^2,10^3} | 0.01 | 0.1 | 0.001 | 0.001 | 1000 | 0.01 | 0.1 | 0.001 | 0.001 | 0.001 | 10 | 0.001 |
| | Penalization (learning rate) | {'l2'} | "l2" | "l2" | "l2" | "l2" | "l2" | "l2" | "l2" | "l2" | "l2" | "l2" | "l2" | "l2" |
| | Optimization solver | {'newton-cg','lbfgs','sag'} | "sag" | "newton-cg" | "lbfgs" | "newton-cg" | "lgfgs" | "newton-cg" | "sag" | "newton-cg" | "newton-cg" | "sag" | "lbfgs" | "newton-cg" |
| | Maximum iterations | {10000,20000} | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 |
| xgboost | Maximum depth of a tree | {3,5,7} | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 |
| | Minimum sum of instance weight needed in a child | {5,10} | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 10 | 10 | 5 | 5 |
| | Number of trees | {500, 1000} | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 |
| | Penalization (learning rate) | {0.01, 0.02, 0.05, 0.1} | 0.01 | 0.02 | 0.02 | 0.05 | 0.05 | 0.05 | 0.02 | 0.01 | 0.01 | 0.02 | 0.02 | 0.02 |
| ensemble 1 | Average of knn-xgboost | Average | | | | | | | | | | | | |
| ensemble 2 | Stacking of knn-xgboost | Meta-function | | | | | | | | | | | | |
| mlp | structures, early stopping, activate function | 31-10-5-1, 256-128-64-32-1 | 256-128-64-32-1 | | | | | | | | | | | |

Table 2: Grid-search parameters of machine learning methods with the optimal results

# Result from Formation (Train/Test) Window



Figure 2. Train and Test Accuracy for machine learning methods with optimal parameters in all subsets

- ## Train subsets
- Best: stacking ensemble
- Xgboost: subset 0, 1, 3, 4 and 5; average ensemble: others

- ## Test subsets
- Best: stacking ensemble
- Xgboost: subset 0-4;average ensemble: subset 7; knn: others.

# Trading window with Momentum

For each of trading dataset

- Sorting all stocks over the cross-section in descending order with the prediction probabilities (stacking ensemble, table 2).

- Long the highest ("expected winner") k probabilities, short the lowest ("expected loser") k probabilities simultaneously.

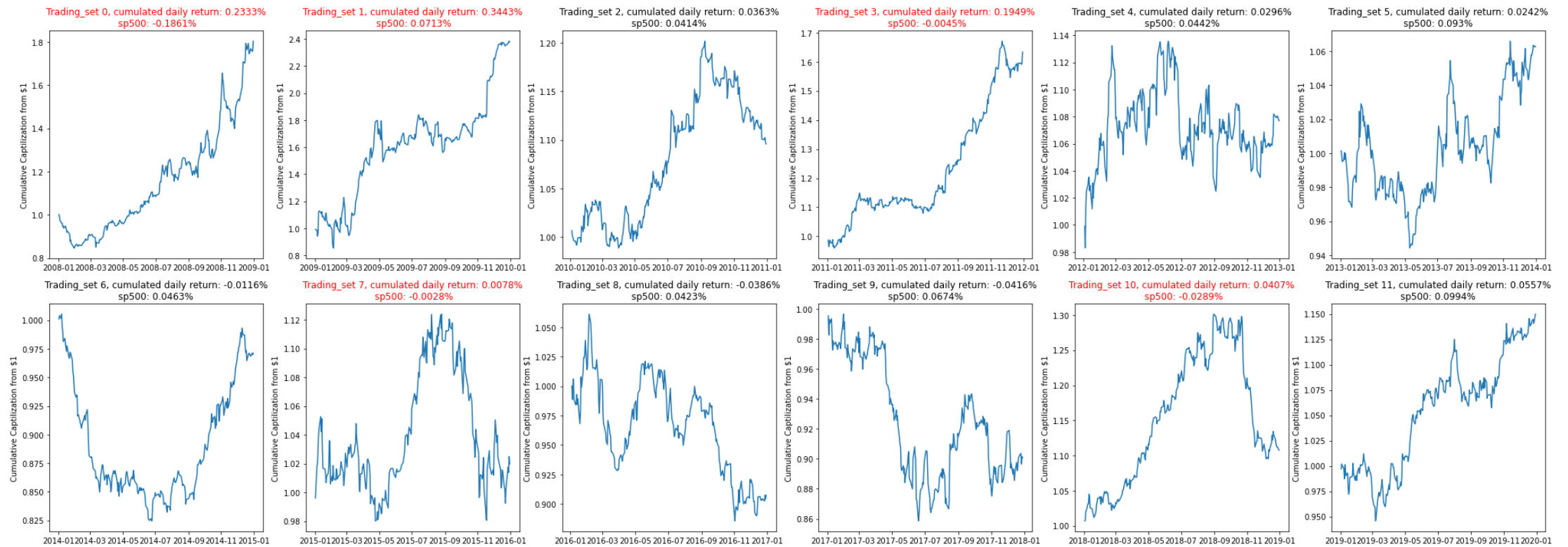# Trading Results of Momentum (3 stocks/pair)



Figure 3: Results of Momentum with 3 stock/pairs
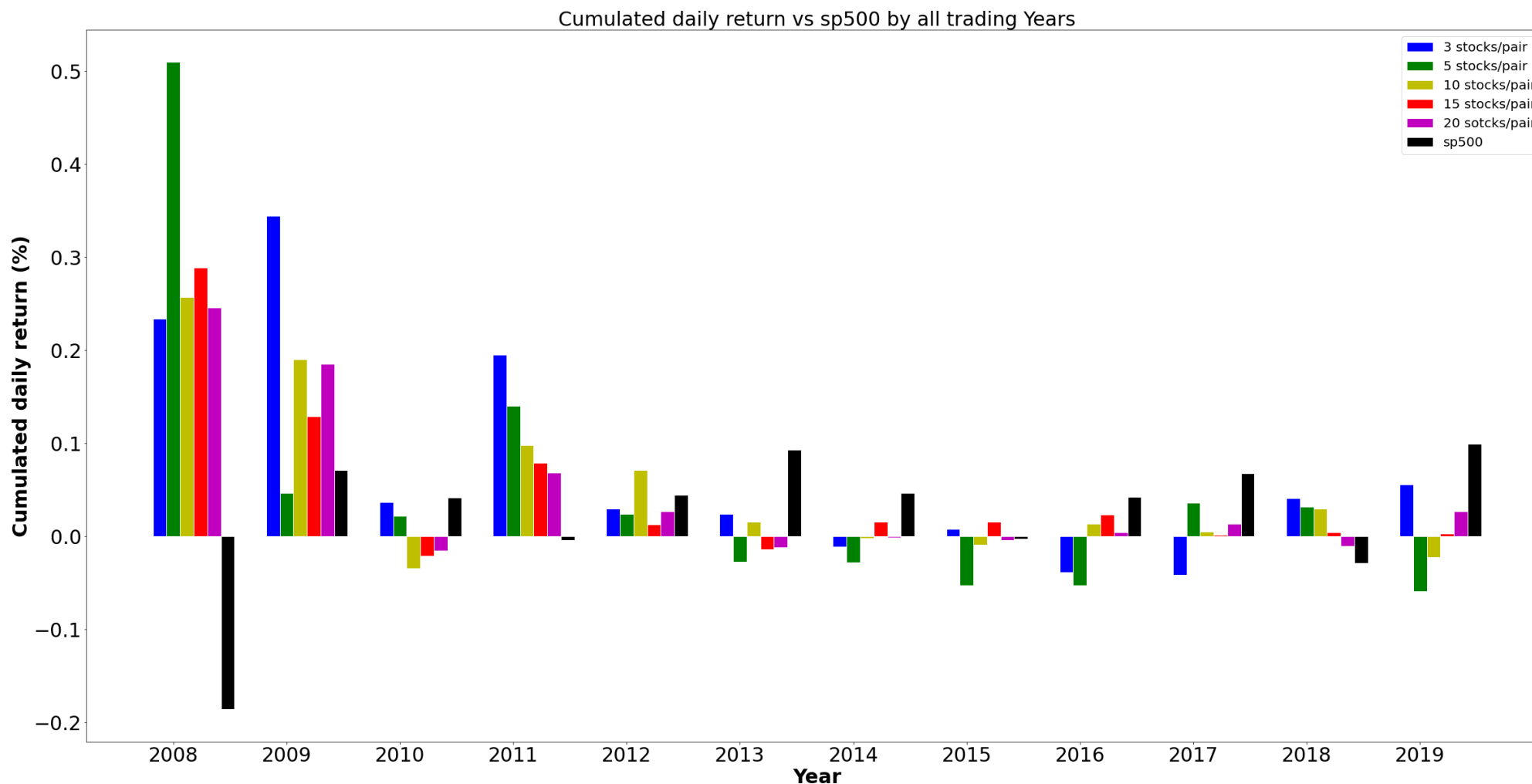
# Stocks K=3, 5, 10, 15, 20/pair vs S&P 500



Figure 4: Results of Momentum with 3,5,10,15,20 stock/pair

# Conclusion

- 2008, 2009 and 2011: beat S&P 500.

- 2013, 2014, 2016, 2017 and 2019: worse than S&P 500.

- Krauss et al. (2017): 0.25% daily return, realized before year 2008 (year 1993 to 2000). In year 2010 to Oct 2015: loss over 50% (with transaction cost). My method: positive daily returns (without transaction cost).

- Confirms hypothesis: "profits are declining in recent years and there is a severe challenge to the semi-strong form of market efficiency."

# Future Work

- Check the data quality: problem with Yahoo finance.

- Other models: recurrent neural network (RNN) and LSTM.

- Try multi-classification

- Use different datasets: Financial Times Stock Exchange (FTSE) 100 or Asian market.

- More criteria (e.g. maximum drawdown), reasons for bad performance years.

- Updated table 3 (tune parameters) and do the momentum trading again.

# Reference

- Jegadeesh, N., & Titman, S. (1993). Returns to buying winners and selling losers: Implications for stock market efficiency. The Journal of finance, 48(1), 65-91.

- Krauss, C., Do, X. A., & Huck, N. (2017). Deep neural networks, gradient-boosted trees, random forests: Statistical arbitrage on the S&P 500. European Journal of Operational Research, 259(2), 689-702.

- Fischer, T., & Krauss, C. (2018). Deep learning with long short-term memory networks for financial market predictions. European Journal of Operational Research, 270(2), 654-669.

- Fischer, T. G., Krauss, C., & Deinert, A. (2019). Statistical arbitrage in cryptocurrency markets. Journal of Risk and Financial Management, 12(1), 31.

# Thank You!

Jiaqi Xu

Email: jiaqiperson@gmail.com

LinkedIn: https://www.linkedin.com/in/jiaqixu1

ResearchGate: https://www.researchgate.net/profile/Jiaqi_Xu10

Project report: https://github.com/jiaqixu/Springboard/blob/master/Capstone/Capstone1/Capstone1_Final_report.pdf