

1.

PJ, you asked me to provide the stock codes and release date time for the 25 Oct 2018 test. I am attaching these information in a spreadsheet which is also the one that I use to analyse the results. Should you need more data such as results in other test cases or all of them, please let me know.

2.

I am afraid I have not been able to find a very comparable paper to my work. I can find either papers that discuss the phenomenon of PEAD from a finance view point, or machine learning papers that try to perform forecast on asset returns under other financial/economical settings. I am not sure if other researchers have put these two things together like I do.

Nonetheless I would like to send you a number of papers that have helped me form ideas. After our meeting I also went back in these paper to specifically look for how they define the dependent variable i.e. the Cumulative Abnormal Return.

- ***Dissecting stock price momentum using financial statement analysis***, published in journal ***Accounting & Finance*** in 2018, provides evidence on the usefulness of financial statement analysis in enhancing the risk-adjusted performance of momentum investing strategies. This is a Finance paper and uses traditional regression methods rather than machine learning models to extensively examine the event driven stock returns and the driving factors. Its goal is not to carry out any forecasting but to understand the weights of each driving factors on the chosen dependent variable, i.e., risk-adjusted return post earnings release. In one of the sections where the paper examines how revision of analysts' consensus forecast on a company's EPS (earnings per share) is driven by numerous factors, the paper states on page 25: "For the forecast period under examination, we then collect the final consensus forecast available **prior to** the earnings announcement for the forecasted fiscal period."
- ***Investor Trading and the Post-Earnings-Announcement Drift***, published in journal ***The Accounting Review*** in 2011, goes through extensive study on PEAD and its possible driving factors in a very similar fashion as the last paper. This paper examines a number of event windows, with some windows being days after an event has happened but the main event window being the three-day window around the earnings release, including days before the events. Please see Table 3 on page 14 and Table 9 on page 28 for this information.
- There are other publications that also use the end of day price prior to an event to evaluate the impact of the event on the subsequent price movements. We can take a look at two Master Theses as example. One is from University of Stavanger, Norway, ***Earnings announcements and stock returns - A study of efficiency in the norwegian capital market***, <https://uis.brage.unit.no/uis-xmlui/handle/11250/183823>, where in section 3.2.4 defines that the starting point of a Cumulative Abnormal Return (CAR) is the last date of the *estimation window*, which is the date when the training period ends, i.e. the date before the event. The other is from UCLA, ***Stock Trend Prediction - Based on Machine Learning Methods***, which also calculates stock returns from one closing price to another.

- Another two academic papers that go through similar discussions with similar definitions on the dependent variable are *Stock Price Reaction to News and No-News: Drift and Reversal After Headlines* from MIT in 2001, and *The stock price reaction to investment news: New evidence from modeling optimal capex and capex guidance* from New York University 2016.

3.

I was wrong when I said there wasn't any credible paper about XGBoost and trading. I had previously looked in our university's online library and Google Scholar but used simple search keywords. I've changed to use more detailed keywords this time and finally found (only) 5 papers that are relevant (having gone through many pages of search results), although none of them are quite similar to mine.

I've provided a quick critique on these papers:

- ***Predicting the Direction of Stock Market Prices using Tree-based Classifiers***, published in journal *North American Journal of Economics and Finance* in 2019, only tests two models, Random Forest and XGBoost, and uses purely historical time series stock data and technical indicators derived from time series stock data as model inputs. It uses data from only 10 companies since Feb 2017. The research goal is to compare the predicted classification rate between these two models over a holding period of 3, 5, 10, 15, 30, 60, and 90 days. However there is little information offered about how the prediction tests have been conducted with no information on the forecast time frame nor precise starting and end points of a stock return calculation period. The important hyperparameters of the models have been chosen by the authors rather than obtained through optimization methods.
- ***Forecasting Stock Market Crisis Events using Deep and Statistical Machine Learning Techniques***, published in journal *Expert Systems With Applications* in 2018, explores a relatively novel and interesting topic of predicting the possibility of a market crash over a 1-day and 20-day horizon across the global markets. Under this guise though it's about forecasting 1-day and 20-day stock market returns and see if they've dropped below a low quantile of historical distribution of stock market returns. By using a vast set of data from global stock markets, bond markets and FX markets, the paper explores a large set of supervised learning models including Logistic Regression, Decision Trees, Random Forest, Support Vector Machines, Deep Neural Network and XGBoost. In the end the paper also draws conclusions by declaring the superiority of certain models including XGBoost over others by examining the forecast results on stock returns through a list of statistical measurement metrics. The paper however does not take care in describing precisely how the dependent variable or the model output is calculated, nor any mentioning of the set up of respective models and systems that encapsulate the models. In the closing stage the paper also talks about using combined models to predict the next financial crisis and shows results. However there is not enough clarity on the set-up under which this particular test is conducted.
- ***Dynamic Weighting Multi Factor Stock Selection Strategy Based on XGboost Machine Learning Algorithm***, included in *2018 IEEE International Conference of Safety Produce Informatization (IICSPI)* conference, is a short paper that uses

XGBoost to dynamically predict the value of a set of seven factors that contribute a stock selection process. Dynamically generated factors are then used to select portfolio of different stocks whose return is measured over a multiple year period. Portfolios of dynamically selected stocks are shown to perform better than benchmark portfolios. However, this paper is vague about many things. It has not clearly described how those factors are actually used to select stocks, nor precisely defined and explained what the inputs are to the XGBoost model. It also hasn't explained very well how the model is set up and how the model hyperparameters are obtained. The paper provides a graph showing how a long-short portfolio performs over a benchmark index over the testing period but there is no reference to long-short portfolios anywhere else in the paper.

- ***Combining Principal Component Analysis, Discrete Wavelet Transform and XGBoost to trade in the financial markets***, published in journal *Expert Systems With Application* in 2019, aims to build a trading system by utilising predicted return of financial instruments one day ahead. The paper gives a very good and clear account of the whole trading system and its constituent components, PCA for input feature dimensionality reduction, Discrete Wavelet Transform for noise reduction, Genetic Algorithm for model optimization and XGBoost as the classification prediction model, and how they work together. For one day forecast the paper also uses the closing price at time t as the starting point and the closing price at $t+1$ as the end point (explained in section 3.2), and measures the model performance over a variety of commodity, stock, and stock index prices.
- The last paper is a UCLA Master's thesis, ***Stock Trend Prediction - Based on Machine Learning Methods***, which has been introduced in the last section.

I have read a lot other reference paper that use different machine learning models and skills in different financial applications. My own impression is that, even a very well written paper which touches upon many different aspects of applying machine learning in Trading can be at its core just about performing prediction on the return of certain financial assets and can be more on the side of investigative and theoretical. A good example is ***Empirical Asset Pricing via Machine Learning*** written by scholars from Yale and Chicago University in 2018, which despite having a number of well discussed topics in it, doesn't even clinically define an asset's excess return (page 8) as used in its context.

In comparison, my own paper is definitely one of those that are much true and relevant to the real market conditions and potential applicability in the market, or at least as a starting point for more advanced research work ahead.

4.

I hope both of you will be ok if I keep the formulae in the paper. Generally all the machine learning papers maintain a certain level of reference to the maths behind the models. To me these formulae are even more relevant because I've listed the model parameters that I seek to optimize in other part of the paper and we are able to locate some of these parameters directly in those formula I've provided which are not simply copied off some text books.

5.

I will compile and propose a list of journals as candidate destinations for future publications. We are going to need them at some point in the future.

6.

PJ, I am still doing some more tests as requested. Specifically I am trying to (a) see if I can find any way to arrive at reasonable prediction results for a prediction whose starting date is after the earnings release such as 1 day after, and (b) see if I can incorporate stock price movement volatility in the results.

I have to be admit experiment (a) is challenging. Just to give an example: A stock's actual return from -1d to 30d could be 10% but the return from -1d to 1d could be 15%, rendering a -5% return from 1d to 30d. In this case it means that although the actual impact of the earnings release event is a positive 10% over 30 days, the stock has actual gone down from 1d to 30d because the majority of the event shock happens on the first day. I'll devote more time and thoughts about how to make this work.

7.

I would like to bring to your attention that, like all machine learning applications in a highly noisy environment, results can not be consistently correct in all situations, at least not at the beginning. For instance, in the recent results document I sent to you (attached to this email) you will find that the result for test 1.3 is the odd one out. There are definitely still patterns that my model is not yet able to capture. Equally important to note that, some of these tests are working with a much reduced set of training points because some of the tests are only working with stock data from a certain industrial sector or stocks of small cap for example. The important message we can say to the audience is that, as supported by all other results, our model definitely can produce the positive results we've been showcasing and we've created a new small line of machine learning research in trading for more future work to add to.