

# Machine Learning Classification Methods and Portfolio Allocation: An Examination of Market Efficiency

Yang Bai  
Kuntara Pukthuanthong  
University of Missouri-Columbia

First Draft: May 1, 2020

This Draft: Oct 2, 2020

## Abstract

We design a novel empirical framework to examine market efficiency through out-of-sample (OOS) predictability. We frame the classic empirical asset pricing problem as a machine learning classification problem. We construct classification models to predict return states. The prediction-based portfolios beat the market in time series and cross-sections with significant economic gains. We directly measure prediction accuracies. For each model, we introduce a novel application of binomial test to test the accuracy of 3.34 million return state predictions. Our models can generate information about the relation between future returns and historical information. The establishment of predictability questions the correctness of prices.

**Key words:** Information theory, portfolio allocation, return state transition, machine learning, classification, artificial neuron network, random forest, dropout additive regression tree, gradient boosting machine, big data

**JEL classification:** C14, C38, C55, G11, G14

*Hal Weizman: “What is the efficient-markets hypothesis and how good a working model is it?”*

*Eugene Fama: “It’s a very simple statement: prices reflect all available information. Testing that turns out to be more difficult, but it’s a simple hypothesis.”*

*Richard Thaler: “I like to distinguish two aspects of it. One is whether you can beat the market. The other is whether prices are correct.”*

— *Are Market Efficient, Chicago Booth Review, Jun 30, 2016*

*“Since 1988, Renaissance’s signature Medallion fund has generated average annual returns of 66 percent. The firm has recorded trading gains of more than one hundred billion dollars. Simons himself is worth twenty-three billion dollars.”*

— *The Man Who Solved the Market: How Jim Simons Launched the Quant Revolution, Gregory Zuckerman, 2019*

## **1 Introduction**

Fama (1970, 1991, and 1998) defines an efficient market as informationally efficient, i.e., prices in an efficient market reflect all available information. In other words, in an efficient market, all information has already been incorporated into prices and there does not exist any pricing error. If there exists no pricing error, prices can only change when new information is available. This implies that no one can benefit from predicting the correction of prices based on any asymmetry of historical information and that there is no such thing as “beating” the market. In short, a necessary condition for an efficient market is that prices are unpredictable. This logic lends us a clear and simple path to examine the market efficiency through the study of predictability.

However, historically, the predictability, especially the out-of-sample (OOS) predictability, is not clear in the literature. Because of the literature’s inability of coming up with a method that delivers robust predictability, Richard Thaler further distinguishes two aspects of market efficiency with his two questions: (1) can we beat the market, and (2) are prices correct? These questions dissect the study of market efficiency to subjects about the predictability and the correctness of prices. He points out that the nonexistence of predictability does not necessarily imply the correctness of prices and vice versa. This way of thinking leaves the room for discussions on the correctness of prices, even if OOS predictability cannot be established. Unfortunately, without clear OOS predictability, discussions on market efficiency are complicated and confusing.

In practice, famous stock investors, such as Jim Simons’ Renaissance Technologies, seem to constantly outperform the stock market through their private trading strategies that predict the future. In other words, there seem existing meaningful predictability that has been explored by these successful

investors. The success of these investors can be due to 3 potential reasons corresponding to 3 forms of market efficiency. First, the market can be completely inefficient. In such a market, prices reflect very limited information and do not even fully reflect the information about past returns. The prices are thus predictable with the past returns. Second, the market can be efficient in the weak form<sup>1</sup>. In a weakly efficient market, prices reflect only the information about the past returns and the prices incorporate the public information including corporate announcements and macroeconomic news with a lag. Investors can take the advantage of the lag and trade on the incorporation of information. Third, the market can be efficient in the semi-strong form. In such a market, prices reflect all public information with no lag, but the investors who have monopolistic access to private information can benefit from trading on the private information. If we can closely examine the way a successful trading strategy beats the market and generate excess profit, we can better understand market efficiency.

For this purpose, we introduce a novel application of machine learning classification methods to the finance literature. We frame the classic asset pricing problem as a machine learning classification problem. Instead of focusing on numeric value predictions, we bucketize the stock returns with the cross-sectional deciles and split the returns into 10 return states. Using the historical information, including the individual stock returns with a lag of at least 1 month, the annual financial information with a lag of at least 6 months, the quarterly financial information with a lag of at least 4 months, the corporate event news with a lag of at least 1 month and the macroeconomic indicators with a lag of at least 1 month, we apply classification methods to predictively classify the future return states of individual stocks and form portfolios based on the predictive classification.

We show that our machine learning classification methods are powerful in portfolio allocation and our portfolios can produce huge OOS profits. Comparing to the empirical findings in the earlier literature, such as Goyal and Welch (2008) and DeMiguel et al. (2009), it is surprising that our classification portfolios, among other machine learning portfolios, can constantly generate OOS gains to a level that the traditional methods cannot achieve. What leads to the success of classification portfolios in OOS comparisons? Among the three potential reasons above that can bring such OOS performance, what is the specific reason that the machine learning portfolios can outperform the market? Given the success of these machine learning methods in the portfolio allocation, are the market prices correct? We choose machine learning classification methods with the specific purpose to look into these questions.

The introduction of the classification methods provides unique benefits for us in both the modeling process and the evaluation process. We take advantage of a clear relation between the classification methods and the information theory. We measure the quality of information extracted from the predictors with cross-entropy and train our models with the optimization goal of reducing information

---

<sup>1</sup>The weak form suggests that today's stock prices reflect all the data of past prices and that no form of technical analysis can be effectively utilized to help investors make trading decisions. Under this form, investors can use fundamental information in the financial statement to identify over- and undervalued stocks to outperform the market. The semi-strong form efficiency theory suggests all information that is public be used to determine a stock's current price and thus investors cannot utilize either technical or fundamental analysis to outperform the market. This form suggests only private information can help investor make profit. The strong form implies all information—both public and private information—is completely accounted for in current stock prices. Investor cannot outperform the market regardless of information we use.

uncertainty. The classification methods also allow us to directly measure model performance through accuracy calculated as the correct proportion of predictive classification. Measuring prediction performance through accuracy is not only easy and explicit but also allows us to conduct formal statistical tests. We further introduce the binomial test to compare our prediction accuracies against the no information accuracy.

The no information accuracy is the highest accuracy that a classifier with no information or limited information can provide. Specifically, the no information accuracy is the accuracy delivered by a naive classifier, which labels the return state of each observation in our sample with the most populated return state in the sample. Based on the rational investor assumption which leads to the efficient market hypothesis, investors have consistent beliefs and enough information about the distribution of macroeconomic variables (Sargent 1994, Barberis and Thaler 2003). If the market is efficient and all information is reflected by prices, investors know about the return distribution but are not able to predict the future beyond the distribution of the returns. The no information accuracy is thus a theory implied benchmark to evaluate whether there exists information about the relation, as captured by a model, between future return states and historical information. Through a set of binomial tests against the no information accuracy, we trace the good performance of our classification portfolios to the statistically significant prediction accuracies. This further implies that our models provide meaningful information about the relation between future return states and historical information. Future return states are thus conditional on historical information and the predictability is established.

Historically, the asset pricing literature was not fortunate enough to come up with strategies that can constantly beat the market in the OOS comparisons in terms of both economic and statistical metrics. For example, Goyal and Welch (2008) examine the market return predictability of the popular predictors and conclude that the popular predictors are not systematically better than the historical mean in the OOS prediction comparisons. DeMiguel, Garlappi and Uppal (2009) examine a range of traditional methods and conclude that the portfolio allocation based on these methods is not systematically better than the naïve portfolio allocation. For a while, the literature could not confirm the OOS performance of any strategy. This inability to replicate the OOS predictability seemed to support market efficiency.

In the recent development of finance machine learning, many methods are identified to be able to constantly produce OOS gains that are of multiple folds of what the market can provide. The examples include Rossi (2018), Gu, Kelly and Xiu (2020) and Chen, Pelgers and Zhu (2020), etc. Following the steps of finance machine learning literature, we propose a new way of creating successful machine learning trading strategies that deliver clear OOS predictability and thus profitability. Built upon our strategies with clear OOS performance, we examine the market efficiency through a thorough analysis and attempt to provide novel insights.

In summary, our findings supply profound insights. First, the significance of accuracies indicates the existence of statistically meaningful predictability. Second, if there exists absolutely no pricing errors, our models should not be able to deliver OOS predictive accuracies that are higher than the

no information accuracy. In other words, the statistically significant OOS predictability of our models suggests that there does exist pricing errors. Third, coupled with the OOS economic metrics that demonstrate the profitability of our classification portfolios, the significance of accuracies indicates that the pricing errors are not minimum and arbitrage opportunities do exist across trading periods. The market corrects the errors in the latter trading periods. Investors applying modeling methods similar to ours can make profits systematically higher than what the market can offer with significantly better risk-return tradeoffs. Therefore, the prices may not be correct. Fourth, this further implies that sophisticated investors can generate information about market prices through their ability to use complex analytical tools. The generated information may not be available to the public and thus may create information asymmetry that sophisticated investors can benefit from. Fifth, the fact that past returns and past corporate announcements contribute to the OOS predictability questions the weak-form and the semi-strong form market efficiency. Sixth, we also document that there exists a substantial imbalance in the return state transition process. The transitions related to extreme return states are with higher certainty indicating lower market efficiency, which brings up the question about the role that the market segments of extreme return states play in market efficiency. Our findings on market efficiency are consistent with the microstructure literature which shows theoretically that the information efficiency is conditional and a full informationally efficient market is impossible (Grossman and Stiglitz 1980). The investors who devote resources to obtain information are thus compensated by the market. At the same time, our findings on market efficiency are also consistent with recent literature indicating that prices are lazy and information may be included in the prices with lags (Cohen, Malloy and Nguyen 2020).

## 1.1 Contribution

Our contribution to the asset pricing literature is six-fold. First, we make methodological contribution to the empirical asset pricing literature. We introduce the machine learning classification methods specialized in single-label multi-class classification. We take a unique angle and reframe the classic asset pricing problem about risk premium explanation and return predictability as a classification problem on the return state transitions. Instead of focusing on numerical value predictions, we focus on the prediction of probabilities associated with future return states. Specifically, we put individual stock returns into 10 cross-sectional return states and study the transitions of return states conditional on historical information. We demonstrate 2 machine learning model architectures, 4 types of algorithms and 22 models. Specifically, we include shallow neuron networks, deep neuron networks, random forests, dropout additive regression trees and stochastic gradient boosted trees.

Second, we answer Thaler’s question about whether we can beat the market through the novel application of classification methods. In the OOS comparisons, the predictions of the classification models can generate average returns, volatility of the returns, skewness of the returns, Sharpe Ratios (SR), certainty equivalent returns (CEQ), and maximum drawdowns (Max DD) that are better than

what the market can provide. For example, in our combined OOS test covering 196301:201912, our best zero investment long-short portfolio can achieve an OOS SR of 0.76 with equal weights and an OOS SR of 0.45 with value weights. The market portfolio delivers an OOS SR of 0.19 and an OOS SR of 0.21 for the two corresponding weighting schemes during the same time period. Note that our SRs are not adjusted by annualization nor R square. Either adjustments can significantly magnify the SR. Despite of using less training data and including only the stocks listed on 3 major exchanges, the performance of our portfolios are competitive and on par with the performance reported in the literature with other methods and the entire CRSP universe including assets that are not stocks. The SR of 0.76 delivered by our best equal-weight model is higher than the SR of 0.707 by the best equal-weight portfolio reported by Gu et al. (2020). The SR of 0.45 delivered by our best value-weight portfolio is higher than the SR of 0.38 delivered by the best value-weight portfolio reported by Gu et al. (2020). The good performance of our portfolios is not from neither taking high leverage nor the concentration of portfolio weights in the microcap stocks. None of our portfolios requires leverage beyond the relax of short selling constraints. When we eliminate the bottom 5% and 10% capitalization stocks, the performance of our portfolios does not disappear.

Third, our introduction of accuracy as a performance metric and the adoption of binomial test contribute to the predictability literature and expands the toolbox for empirical asset pricing studies. We carefully analyze the in-sample (IS) and the OOS prediction accuracies and provide explanation of the good portfolio performance from the angle of information theory. We introduce the accuracy as a metric to evaluate the overall performance of the classification models. The direct measurement of the accuracy as the correct proportion of predictions is only available to classification problems. In numeric value predictions, all metrics are based on prediction errors and it is hard to directly measure the accuracy of predictions. We trace the good performance of our classification portfolios to the accuracy of the return state predictions.

Fourth and most importantly, we further introduce the binomial test to the asset pricing literature, which enables us to conduct meaningful statistical test on the prediction accuracy. We also introduce the no information accuracy, which is the highest accuracy that a classifier using no information can provide. Across multiple setups, we show that our models deliver statistically meaningful predictability and are time-invariantly applicable to generate predictions of future return states. The binomial tests on prediction accuracies against the no information accuracy has profound meaning to the study of market efficiency. The no information accuracy is an accuracy under the assumption of efficient market. In other words, the no information accuracy is the highest accuracy of prediction assuming that no further information can be generated to describe the relation between future prices and historical information. Therefore, a binomial test on the prediction accuracy against the no information accuracy is not only a test on the predictability but also a test of the market efficiency. The statistical significance found in our binomial tests against the no information accuracy implies the generation of information through our models about future return states based on historical observations. This piece of information is not reflected by the current prices and therefore not shared by the majority of

market participates. At the same time, this piece of information does generate OOS profitability. This naturally brings up the question about the correctness of the prices, i.e., the prices may not be correct as people can make profit with public information.

More specifically, combining the predictability of our models, the information generation and the OOS economic gains, our findings show that the prices will move toward the same direction as what the generated information indicates. This means that the market will gradually incorporate the information known privately to the sophisticated investors and move towards the price level that reflect the information produced by complex tools *ex ante*. Our findings indicate that, across the entire CRSP-COMPUSTAT sample, there are systematic trading opportunities based on historical information to generate excess profits on a monthly basis. Considering the size of the economic gains, the profitability that can be generated by trading on the information from complex tools may not be ignorable. This profit related to the generated information should have been eliminated by arbitrage. Therefore, the current market prices may not be correct or may not fully reflect all public information. This also answers Thaler's second question.

The generation of the new information also has an important implication that is directly related to the strong-form market efficiency. The information generation shows that the sophisticated investors, such as Renaissance Technologies, who are able to understand and use complex tools that are similar to machine learning classification methods, can generate information about future return states from their interpretation of historical observations. They can apply the new information to their trading. The generated information, depending on the analytical tools, are likely to remain unique and monopolistic. In other words, the private information may not need to be insider information that is known to the management team of a firm and can be generated based on analyzing historical information. This means that there is a possibility for sophisticated investors to manually introduce information asymmetry to the market.

We are the first to introduce the binomial test and the accuracy metric in the study of return predictability and market efficiency. By far, the binomial test on the prediction accuracy that indicates the new information generation is also a unique contribution to the finance machine learning literature and the market efficiency literature. The generation of new information also helps explain the good performance of successful portfolio allocation strategies from an information point of view.

Fifth, we compare the easiness of predicting different return state transitions and look into our models' strategy that produces the good prediction accuracy. We demonstrate the true return state transition probability matrix. It shows that the return state transitions are not uniformly distributed. The center of the true transition probability matrix is distributed more uniformly, which indicates a higher level of uncertainty. The corners of the true transition probability matrix are with the highest transition probabilities indicating a lower level of uncertainty. The different levels of uncertainty may imply the different levels of the market efficiency among the segments of the market. Higher uncertainty is related to lower predictability, indicating higher level of market efficiency and vice versa. We show that our models benefit the most from the most certain transitions and almost give

up on the more uncertain transitions. We are among the first to supply the asset pricing literature an economic insight of the success of the machine learning portfolio allocation.

Sixth, through our demonstration of the classification methods, we analyze the training process. We contribute to the identification of the important predictors of return states and the understanding of market efficiency in its weak form and semi-strong form. We construct cross-sectional tests in the spirit of Fama and French (2018) by splitting the CRSP sample into odd number months and even number months. The cross-sectional OOS tests show that our models have good CS OOS explanatory power. In addition, we look into the training process of the CS models and the time series (TS) models. In the training process of both the CS models and the TS models, the industry information, the corporate announcements, the macroeconomic indicators and the historical return information all make important contribution. As all our predictors are lagged by at least 1 month and many of the firm characteristics are lagged by at least 6 months, these findings are interesting to the study of market efficiency, especially considering that we do not update our models in the TS OOS testing periods with at least a time length that is close to 30 years<sup>2</sup>. Coupled with the OOS portfolio performance based on the model predictions, we conclude that the lagged public information, including the past returns and the historical corporate announcements, can help predict future return states. The returns are predictable and there may be a lag for the market to reflect all the public information including the past returns and the historical corporate announcements.

In summary, combined with our points made on the predictability and the information generation, our findings show that there is still room for the market efficiency to improve. We compare two time periods of our time series OOS tests and confirm that the predictability of the models decreases in the second half of the CRSP sample (199201:201912). This can be caused by the smaller number of observations in the first half of the CRSP database or the improvement of the market efficiency in the latter half of the CRSP sample.

## 1.2 Literature

We review the recent development in the finance machine learning literature below. Our review is by no means an exhaustive list of the works in the finance machine learning literature and we try to include the works that are the closest to ours in a chronological order (as of the draft date) based on the publication date for the published papers and the latest update date for the working papers. We categorize the papers based on the model implementations.

---

<sup>2</sup>We are looking into dissecting the insights for weak form and strong form market efficiency by separately developing models that use only past trading information and models that use only past corporate news. Results will be included in the next update of the draft.



### **1.2.1 Characteristics**

Brandt, Santa-Clara and Valkanov (2007) is a pioneer in leveraging characteristics in the portfolio allocation problem. They parameterize the portfolio weight of each stock as a function of the stock's characteristics and they estimate the weights of the stocks included in the portfolio with the maximization of the representative investor's utility. The optimal portfolio relative to holding the market provides an in-sample (IS) CEQ gain of 11.1% and 5.4% OOS CEQ gain. Green, Hand and Zhang (2016) is the pioneer of studying the large number of firm characteristics. They construct a sample of more than 100 firm characteristics based on stock performance and financial information. They show that there are 12 characteristics that are reliably independent in contributing to the return predictability during their sample period and the predictability drops after 2003.

### **1.2.2 Tree Models**

Moritz and Zimmermann (2016) introduce the regression trees in the cross-sectional pricing. Using the regression trees they show that the past short-term returns are the most important predictors for the future returns. They sort the portfolios with tree structure. They show that the conditional portfolio sorts through tree structure improve predictions significantly over Fama-MacBeth regression. Rossi (2018), using boosted trees, documents that the non-linearity of the popular Goyal and Welch (2008) predictors can time the market. He emphasizes that the relation between predictors and the best allocation to risky portfolios is non-linear. Brayzgalova, Pelger and Zhu (2020) demonstrate the advantages of applying pruning in the selection of the sorting methods to improve the empirical asset pricing models.

### **1.2.3 Neuron Networks**

Chen, Pelger and Zhu (2020) focus on the neuron network models and asset pricing. They combine 3 neuron networks and essentially generalize the linear pricing kernel under the framework of neuron networks. They introduce the generative adversarial neuron network models to the playground of fine search of the best SDF by identifying the assets that are hardest to model. They also enforce the non-arbitrage constraint to the loss function in the architecture of the networks. Aubry, Kraussl, Manso and Spaenjers (2020) introduces the machine learning methods to the playground of illiquid assets. Specifically, they apply neuron networks to a data with one million painting auctions based on visual and non-visual characteristics of the art pieces. They show that their methods perform drastically better than the traditional pricing methods. Feng et al. (2019) propose the use of neuron networks in the extraction of hidden features and augment the hidden features in the pricing models.

## 1.2.4 Tree Models and Neuron Networks

Gu et al. (2020) demonstrate the powerful pricing capability of the neuron network models and the tree models. Their experiments show the possibility to double the performance of leading regression-based strategies. They also try to form OOS portfolios by predicting the stock returns first and then forming portfolios the stocks based on predicted return. Their best equal-weight strategy coming from a 4 hidden layer neuron network delivers a shockingly 27.1% return on annualized basis. In a concurrent work of ours, Wolff and Echterling (2020) construct 37 stock characteristics and also characterize the portfolio allocation problem as classification problem. They apply neuron networks and tree models to S&P 500 constituents with 21 years of weekly data. They show that their models can also be applied to STOXX Europe 600. Bianchi, Buchner and Tamoni (2020) apply machine learning methods including neuron networks and tree models to the bond market. They demonstrate the superior performance of the machine learning methods in predicting bond returns.

## 1.2.5 Related Works

The two closest related works to our paper are Gu et al. (2020) and Wolff and Echterling (2020). Our paper is distant in many aspects from Gu et al. (2020). First, our introduction of classification is fundamentally different from the implementation of Gu et al. (2020). In fact, we view the portfolio allocation as a selection problem of the stocks, while Gu et al. (2020) view the portfolio allocation as an estimation problem of the stock returns. In other words, we model on the probabilities of return state transitions conditional on historical information and Gu et al. (2020) model on the numeric value of returns. Taking the neuron network models as an example, our neuron network models all include a soft-max output layer of 10 neurons that gives us the probabilities of a stock being in one of the 10 return states in one period ahead, while neuron network models in Gu et al. (2020) all have a linear output layer of 1 neuron and output a return prediction. Similarly, our boosted tree models are based on multi-class probabilities, while the boosted tree models in Gu et al. (2020) are based on linear regressions. Our output gives a better sense of relative performance of stock returns and the probability of occurrence.

Our paper is also very different from Wolff and Echterling (2020). First, we model on monthly returns covering 196301:201912 including all 26302 stocks out of all 33004 securities. We include 332 predictors covering historical returns, firm characteristics and macro indicators, while Wolff and Echterling (2020) include 37 predictors and 21 years of weekly data (199901:201912). In addition, to the specific purpose of our study, we characterize the returns into 10 return states independent of market return, while Wolff and Echterling use binary categorization with reference to market return.

The most important aspect that distinguishes our paper from Gu et al. (2020) and Wolff and Echterling (2020) is the scope of the studies. We choose machine learning classification methods specifically because of their relation with the information theory and the testing metrics that can be adopted. Beyond the modeling aspects and the predictive power, we attempt to examine market efficiency and

supply answers to Thaler’s questions. We also aim at providing new economic intuition about why the machine learning methods can produce portfolios that outperform the market. Through our comprehensive analysis and the demonstration of our 22 models, we show explicitly that the investors can generate new information and the market is unbalanced in terms of the transition probabilities.

The remainder of the paper proceeds as the follows. We specify the models, metrics and the empirical setup in Section 2. In Section 3, we demonstrate the OOS performance of our classification based portfolios with economic and statistical metrics. We analyze the performance through accuracy and discuss binomial tests. We also provide insights about the portfolio allocation strategies generated by the classification portfolios. In Section 4, we document the cross-sectional explanatory power and predictor contribution. We conclude in Section 5.

## **2 Methodology**

We provide the general description of our methods in this section. We explain the basics of our modeling processes and tests, including our model specifications, validation and hyperparameter tuning, testing metrics and sample split. We try to provide details so that people with limited experience in machine learning know about the terminologies and the model specifications.

### **2.1 Model Training and Validation**

#### **2.1.1 A Brief Introduction to Classification and Information Theory**

A classification problem is a choice making problem. For example, given a picture capturing an animal with 2 possible outcomes, cat and dog, a classification problem can be framed as the question: is the animal in the picture a cat? This is a binary choice question. If the answer is yes, we know that the picture captures a cat. If the answer is no, the picture captures a dog. This is a typical binary classification problem with one class being the cat pictures and the other class being the dog pictures. The task in this classification problem is to find a strategy to label a picture to be either a picture of cat or a picture of dog. A strategy is referred as a classifier or a model in machine learning literature. If we have a classifier that always guess that the animal captured by any picture is a cat, then this classifier is a naive binary classifier. The classification outcome of a given picture, i.e. cat or dog, is called the label of the picture. A classification problem is not limited to have 2 candidate outcomes nor a single label. For example, the question that asks “what is the weather tomorrow?” is a multi-label multi-class classification problem. Specifically, for example, an answer to the question can include 2 labels, one about the weather condition and the other about temperature. The candidate outcome weather conditions can include rainy, snowy, sunny, etc. The candidate outcome temperature can include 3 levels: hot, mild and cold. In this paper, we frame our portfolio allocation practice as a single label multi-class classification problem.

In Table 1, we demonstrate how we frame a portfolio allocation problem as a classification problem. We cross-sectionally rank individual stock returns by trading month, put them into their corresponding deciles and use the deciles as the classes of return states. For example, if a stock falls into the lowest decile in a trading month, we define the true label of the stock as the class of return state 1. A stock in return state 1 means that the stock delivers a return that is among the worst performing returns of the trading month. A stock in return state 10 indicates that the stock are among the stocks delivering the best performing returns of the trading month. In later sections, we refer to the return states with the numbers specified in Table 1. In short, small number returns states indicate bad performing return states while large number return states indicate good performing states. Note that we make the lower bound as the inclusive bound. Therefore, we have slightly unbalanced 10 classes.

In the Section 6 of Claude Shannon's (1948) seminal paper, Shannon introduces the concept of bit as the unit for information and the famous Shannon entropy, or information entropy (entropy hereafter), as the measure of the average level of information, or the amount of randomness, in the unit of bits. An arbitrary parent choice question, for example, can be decomposed into a series of binary choice sub-questions and the entropy summarizes the average number of the binary choice sub-questions to answer such that the parent choice question can be answered. Higher entropy means that there is more uncertainty. The entropy can be defined as

$$Entropy(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-,$$

where  $S$  is a Bernoulli trial with 2 possible outcomes  $\{+, -\}$  that are mutually exclusive and  $p_+$  and  $p_- = 1 - p_+$  represent the probabilities of the two possible outcomes respectively. When the entropy is seen as an uncertainty measure of the information, the greater the entropy is the greater the uncertainty is. In a binary case, the entropy is at its largest when  $p_+ = 1 - p_- = 0.5$ . In other words, a Bernoulli distribution that are close to a binary outcome discrete uniform distribution is with the highest uncertainty and is more likely to yield surprising outcome. When a Bernoulli distribution departs from the binary outcome discrete uniform distribution, the uncertainty decreases and a random draw from the Bernoulli distribution is more likely to deliver an unsurprising outcome. Specifically, for example, consider the weather outcomes of snow versus not snow during winter, Florida has a lot lower uncertainty of snow comparing to the uncertainty level of snow in Missouri, since the probability distribution of snow versus not snow in Florida is skewed to towards snow and concentrated in not snow while the probability distribution of snow is closer to uniform distribution in Missouri.

A generalization of entropy to compare the difference between two probability distributions yields the cross entropy. Given a parent choice question and our best strategy to form binary choice sub-questions, the cross entropy measures the average number of binary choice sub-questions that we need

to answer with our best strategy in the environment with the true probability. In the training process of a machine learning model, we have the observed probability distribution and we will select the best strategy, through comparing the candidate estimates of the distributions, to form the binary choice sub-questions and answer the parent choice question. The best strategy is associated with the lowest cross-entropy and thus also with the lowest information uncertainty. The information uncertainty is associated with the possibility of information loss. In other words, the cross-entropy can be thought as a measure to compare the observed probability distribution and the predicted probability distribution.

### 2.1.2 Loss Function and Optimization

In a training process of multi-class classification problem, we want to minimize the overall errors and balance the by-class performance of the classifier. Popular choices of loss functions includes accuracy, information gain based on entropy, mean error across the classes, etc. We adopt cross entropy as the loss function to minimize the information uncertainty between our predicted distribution and the observed distribution. Specifically, we achieve this through comparing our predicted distribution against the true distribution with cross entropy. The formal definition of cross entropy is:

$$\begin{aligned} \text{Cross Entropy} &= \mathbb{E}_{p_i} \log_2 \frac{1}{\hat{p}_i} \\ &= \sum_i p_i \log_2 \frac{1}{\hat{p}_i}, \end{aligned}$$

where  $p_i$  is the observed probability of outcome  $i$  and  $\hat{p}_i$  is the predicted probability. Higher cross entropy represents higher information uncertainty associated with the use of the predicted probability to approximate the observed probability. In other words, higher cross entropy represents lower information decoding quality with the predicted probability comparing to the observed probability. When we train a classification model, for each iteration of weight update, we aim at minimizing the cross entropy through the adjustment in model weights. When we conduct hyperparameter training in the format of grid search, we also select the best model based on cross entropy. By using cross entropy as the criterion to adjust weights and select hyperparameters, we can obtain a model that reduces the uncertainty of the information extraction.

### 2.1.3 Model Specification, Validation and Hyperparameter Tuning

We describe the models we consider in Table 2. We include 2 architectures covering shallow neuron networks, deep neuron networks, dropout additive regression trees, random forest and gradient boosting machine, total of 22 separate models. In our specification and model training, to understand the effect of the major architectural parameters, which control the model complexity, we do not add the architectural parameters, such as the number of hidden layers in a neuron network or the number

of maximum depth in a tree model, into the hyperparameter tuning process. Through training and evaluating models with the different key parameter specification, we can see a clear trend later that the differences in the key structural parameters have a strong influence on the performance of the models.

Specifically, we include neuron networks with 1 to 4 hidden layers. Our tree models are with maximum depths of 2 to 8. The number of hidden layers controls the complexity of interaction across predictors. The number of maximum depth limits the maximum number of leaves that a tree can grow. The complexity of computation can increase exponentially as we increase the depth of the tree models. The last column in Panel A of Table 2 presents the structural capacity of our models. The numbers in curly brackets correspond to the number of neurons in the specific hidden layer. For example, the model ANN4 128 has 128 neurons in the first hidden layer, 64 neurons in the second layer, 32 neurons in the third hidden layer and 16 neurons in fourth hidden layer. Thus, the numbers of neurons in the curly brackets for ANN4 128 is {128,64,32,16}. For our DART models, beyond the important parameters summarized in Panel A, we also specify the dropout rate as 10%. This dropout rate can help generalize the model. To save on computation resources, we also apply early stopping mechanism to all of our models. If a model in training process does not improve the loss function by at least 0.00001 for 3 consecutive rounds, the training stops.

Panel B in Table 2 presents the additional specification information that applies only to neuron network models. Relu is the popular choice of activation function for the hidden layers in the recent finance machine learning literature. However, we did not use Relu as the hidden layer activation function. We want to avoid the dead neurons in the deeper layers of deeper networks. Our selection of Tanh function is famous for its robustness. Because our interest of the study is the return states as classes, we specify the output layer with SoftMax function, which transforms the inputs from the last hidden layer to the probabilities. We set the output layer to include 10 neurons, corresponding to 10 possible return states. For each instance fed to our neuron networks, each neuron in the output layer will produce a probability representing the likelihood that the instance belongs to the associated return state. In the end, we categorize a stock to one of the return states associated with the highest probability.

In any of our neuron networks, we have 3 layers of transformations starting from the input layer. Consider a neuron network with 1 input layer, 1 hidden layer of 1 neuron and 1 output layer with 10 neurons. Let us denote the input layer as  $X$ . We form transformation through activation function taking linear combination of the input layer as its input. The transformation is defined as  $Z = \sigma(\alpha_0 + \alpha^T X)$ , where  $\sigma$  is the Tanh function. Then, we further transform the output of the Tanh function through another linear combination and connect linear combination with the output layer of SoftMax function. Specifically, we first collect the linear combination taking  $Z$  as input and denote the linear combination as  $T_k = \beta_{0k} + \beta_k Z$ , where  $k$  is the number of neurons in the output layer. Then, we connect  $T$ 's to the observations through the SoftMax function  $g(T) = \frac{e^{T_k}}{\sum_{l=1}^{10} e^{T_l}}$ . During training, we adjust the weights in all layers and bring the SoftMax function to produce probabilities for individual observations as close to the real probabilities as possible. In a tree model, the logic is similar but different. In our tree

models, the training process is done with the sub-sample of individual classes. Each tree will select a class of return state and use the subsample of the training set containing the selected return state to learn and adjust the weights. The weights of different trees are not directly interacting with each other until the summarizing step when the model pulls all the information about individual classes together and conducts a majority vote process to decide the prediction. When the majority vote step takes place, the probabilities of individual classes will be summarized and scaled to reflect the overall probabilities with a summation of 1.

Panel C of Table 2 presents the hyperparameters that we want to tune with our validation strategy. As we separate the architectural parameters from the hyperparameter set to help us understand more about the influence of model complexity, we only have limited number of hyperparameters to tune with our model. Specifically, for neuron networks, we tune L1 regularization parameter which decides the penalty put on the weights similar to the regularization in the lasso regression. Our neuron networks prefer the finite L1 regularization. We tune the sampling rate for training data and the sampling rate for the predictors as a control for generalization of the tree models. We take cross validation as the validation strategy for hyperparameter tuning. We choose cross validation, instead of constructing a separate validation sample as being implemented by Chen et al. (2020) and Gu et al (2020), to take the advantage of the data coverage and avoid the loss of OOS testing observations. Sepcifically, we separate the training data set into 5 subsamples in chronological order and conduct 5-fold cross validation.

## 2.2 Performance Evaluation

In order to better communicate our empirical findings, we describe the metrics that we refer to. For model-based portfolio allocations, it is important for us to understand both the economic performance and the statistical performance. Therefore, we list out both the economic metrics and the statistical metrics.

### 2.2.1 Economic Metrics

The purpose to evaluate a model based portfolio economically is to understand whether the portfolio is successful in terms of commonly used traditional measures. Specifically, we refer to Sharpe Ratio (SR) and Certainty Equivalent Return (CEQ) in the evaluation of the risk-return trade-off. Portfolios with better performance in terms of risk-return trade-off have higher SR and CEQ. We define SR as

$$SR = \frac{\mathbb{E}(R - R_f)}{\sigma(R - R_f)},$$

where  $R$  is the return generated from a portfolio of interest and  $R_f$  is the risk free rate of return. For the long-short portfolios, we define the SR as

$$SR_{long-short} = \frac{\mathbb{E}(R_{long} - R_{short})}{\sigma(R_{long} - R_{short})},$$

where  $R_{long}$  is the return generated from holding the long position of the predicted good performing stocks and  $R_{short}$  is the return generated from holding the long position of the predicted bad performing stocks. We define the long-short SR in this way as the long-short portfolio is a theoretically zero investment portfolio. Following DeMiguel, Garlappi and Uppal (2009), we define CEQ as

$$\widehat{CEQ}_k = \hat{\mu}_k - \frac{\gamma}{2} \hat{\sigma}_k^2,$$

where  $\hat{\mu}_k$  is the estimated mean of the return from the asset  $k$  and  $\hat{\sigma}_k^2$  is the variance of the return.  $\gamma$  in the above expression stands for the risk aversion coefficient and we specify  $\gamma = 1$  following DeMiguel et al (2009) and Goyal and Welch (2008).

In addition to these most popular economic metrics, we also provide basic metrics to evaluate the profit and loss. We specify the cumulative return as

$$Y_{t:t+n} = \prod_{i=t}^{t+n} (1 + R_i) - 1$$

, where  $R_i$  is the return from the portfolio of interest in the month  $i$  and  $n$  stands for the number of periods in the investment window. Our cumulative return is therefore defined as the product of gross return net of the initial investment cost. We take the notation of our cumulative return and include maximum drawdown in our evaluation defined as the following:

$$MaxDD_{t:t+n} = \min_{t:t+n} \left\{ \frac{Y_{i+1} - Y_i}{Y_i} \right\},$$

where  $i$  is a trading month during the investment window  $t : t + n$ . Finally, following Gu et al. (2020), we provide turnover defined as

$$Turnover = \frac{1}{n} \sum_{i=t}^{t+n} \left( \sum_j \left| w_{j,i+1} - \frac{w_{j,i}(1 + r_{j,i+1})}{\sum_k w_{k,i}(1 + r_{k,i+1})} \right| \right),$$

where  $w_{j,i}$  represents the weight of stock  $j$  during month  $i$  in a portfolio.



### 2.2.2 Statistical Metrics

For our classification models, we introduce and report a range of metrics that focuses on model accuracy from both the angle of overall accuracy and the angle of balancing the prediction accuracy across different classes. To better present the statistical metrics, suppose that we have a binary classification problem and a classifier making predictions. Consider the following matrix which compares the true values and the predicted values:

	Reference Positive	Reference Negative
Predicted Positive	A	B
Predicted Negative	C	D

, where the rows indicates the predicted value and the columns are the reference of ground truth. The letters A, B, C and D stand for the number of observations. Specifically, for example, A stands for the number of the observations with the true positive label that also are predicted to have positive label. In such case, A is the number of correctly predicted positive observations, or the true positives. Similarly, D is the number of correctly predicted negative observations, or the true negatives. B and C stand for false positives and false negatives. A matrix that compares the number of predicted observations with the ground truth is called a confusion matrix.

With the basics of the confusion matrix being introduced, we can further introduce the popular metrics that evaluate the performance of a classification model. First, referring back to the confusion matrix example above, we define sensitivity and specificity as

$$\text{Sensitivity} = \frac{A}{A+C}$$
$$\text{Specificity} = \frac{D}{B+D}.$$

Sensitivity measures the accuracy of the predicted positives, while the specificity measure the accuracy of the predicted negatives.

In an ideal situation for a binary classification problem, we want to maximize the overall accuracy or the number of A+D and at the same time keep a balance between making positive and negative predictions. A classic example can be the detection of cancer. The proportion of cancer patient over the entire population who take the cancer screening is a relatively small number. Therefore, in such situation, a classifier can gain a very high accuracy if the classifier just simply predictively label all people in the cancer screening as negative. However, the classifier will then fail to detect any potential cancer patient. Another example is about information. For a extremely skewed distribution such as the chance of raining in Sahara desert, telling the information receiver that the rare event will happen, eg., it is going to rain in Sahara desert, carries more information comparing to telling the information receiver that the outcome associated with the largest possibility is likely to happen, eg., it is not

going to rain in Sahara desert. When we train a machine learning model, the overall accuracy is just one aspect that we care. We also care about whether the model generates new information and has meaningful detection rate for each class. Thus, balancing the true positives and the true negatives is important.

Similar to sensitivity and specificity, we can have prevalence and detection prevalence. The prevalence measures the ground truth percentage of the sample being positive and the detection prevalence measures the predicted percentage of the sample being positive. Specifically,

$$\begin{aligned} \text{Prevalence} &= \frac{A + C}{A + B + C + D} \\ \text{Detection Prevalence} &= \frac{A + B}{A + B + C + D}. \end{aligned}$$

Beyond the above metrics that look into individual aspects of the predictions, there are 3 comprehensive metrics: the F1 score, the balanced accuracy and Cohen's Kappa:

$$\begin{aligned} F1 &= \frac{(1 + \beta^2) \times \text{Precision} \times \text{Recall}}{(\beta^2 \times \text{Precision}) + \text{Recall}} \\ \text{Balanced Accuracy} &= \frac{\text{Sensitivity} + \text{Specificity}}{2}, \\ \kappa &= \frac{p_o - p_e}{1 - p_e} \end{aligned}$$

where  $\beta$  in the F1 score is the type II error. The precision and the recall in F1 score are defined as  $\text{Precision} = \frac{A}{A+B}$  and  $\text{Recall} = \frac{A}{A+C}$ .  $p_o = \frac{A+D}{A+B+C+D}$  in  $\kappa$  is the relative agreement observed between the ground truth and the prediction, while  $p_e = p_+ + p_-$  measures the probability that the agreement between the prediction and the ground truth is random, where  $p_+ = \frac{A+B}{A+B+C+D} \cdot \frac{A+C}{A+B+C+D}$  and  $p_- = \frac{C+D}{A+B+C+D} \cdot \frac{B+D}{A+B+C+D}$ . Note that F1 does not take into account the true negatives and thus have limitation on evaluating the results with consideration on the balance between true positives and the true negatives.

### 2.2.3 Accuracy and Binomial Test: A Novel Empirical Framework

In addition to the above metrics, we also introduce accuracy as a direct statistical metric of prediction performance. An accuracy is defined as the associated model's proportion of correct predictions. We also introduce a formal test on the statistical significance of the prediction accuracy. For any classification problem, after the classifier makes the prediction, we have two types of observations, the correctly classified observations and the incorrectly classified observations. Therefore, the prediction for an observation is a Bernoulli trial for a classifier. The number of correctly classified observations can be seen as the number of successes and the number of incorrectly classified observations can be

seen as the number of failures. Following this logic, the prediction accuracy measured as the number of correctly classified observations over the total number of observations is a type of success rate. This enables us to conduct a standard binomial test to compare the success rates. Specifically, we can test the accuracies of our models against some accuracy of a benchmark classifier and we can further understand whether our models can provide more information than what the benchmark classifier provides.

#### **2.2.4 Selection of Benchmark Classifier**

Since the efficient market hypothesis suggests that the market is efficient and there does not exist relation between future returns and historical information, price changes are decided with new information to be released. Consequently, we want a benchmark classifier that makes predictions based on limited information or no information to reflect the implication of the efficient market hypothesis. In general, there are two types of classifiers we can consider. First, we can consider a random classifier. Second, because the investors are rational and the prices in efficient market reflect the expected future prices, investors have rational understanding over the distribution of prices as required by the consistent beliefs which leads to market efficiency. Thus, we can consider a classifier with distributional information of returns. Discussion about the inclusion of classifiers with distributional information of returns is provided in the next subsection.

We want to be conservative. Therefore, from those classifiers using limited or no information, we want to select a classifier that produces the highest possible accuracy. We consider 5 candidate benchmark classifiers. First, we consider a random classifier that takes into account absolutely no information. The random classifier labels each OOS observation with a random label from return state 1 to return state 10 with equal probabilities. The accuracy level of the random classifier reflects the situation where no relation exists between future return states and historical information. Second, we consider a random classifier that randomly labels OOS observations with probabilities based on the observed IS return state probability mass function. Third, we consider a naive classifier that labels OOS observations with the most populated IS return state. The second and the third benchmarks represent the situation where market knows the IS distributional information of return states. The machine learning literature argues that if the distributional information of the response variable is the only information or the only useful information, classifying all observations predictively into the majority class is the best guess. Our multiple comparison test confirms this argument. Fourth, we consider a random classifier that randomly labels OOS observations with the knowledge about OOS return state distribution. Fifth, we consider an accuracy by a naive classifier that labels OOS observations with the most populated OOS return state. The fourth and fifth benchmarks are enhanced versions of their IS counterparts.

### 2.2.5 Discussion on Classifier with Distributional Information of Returns

There are three reasons for us to consider the benchmarks reflecting the distributional information about return states. First, it is a convention in machine learning literature to test predictability against the naive classifier as a benchmark. Since the distributional based classifier accuracy is easily accessible through our empirical observations and does not involve any modeling or predictors, it is a natural choice to examine whether a model and the associated predictors are working, i.e., whether the model and the predictors provide more information than the distributional information. We follow the convention, considering that the distributional information is minimum amount of information. In fact, the accuracy delivered by a naive classifier is often referred as no information accuracy indicating that the minimum amount of information is reflected by the accuracy. Therefore, the benchmarks including the distributional information do not deviate us from our testing purpose about whether useful information is provided by our models and the historical information that our models use.

Second, we would like to demonstrate how we select the most restrictive benchmark among the candidate classifiers for our binomial tests. The selection is conducted through a multiple comparison test with the Monte Carlo samples of accuracies by the random classifiers and the naive classifiers. The test looks into which classifier provides the highest accuracy on average. The OOS prediction accuracies by each of the random classifiers and the naive classifiers provide the simplest setup to test predictability within our novel testing framework. The selection process contributes unique insights to the market efficiency literature. In fact, through Monte Carlo simulation and the multiple comparison test, we show that the naive classifier using the IS distributional information about return states delivers an accuracy that is statistically higher than what is delivered by the random classifier that uses absolutely no information. This means that the future returns can be better predicted by the random classifier using the basic information about IS probability mass function of the return states. However, it is worth to note that predictability does not necessarily equal to the rejection of market efficiency. A meaningful empirical question about market efficiency from predictability has to come with meaningful profitability. We do not see such profitability with the naive classifiers. We will discuss more about this point in the subsection below.

Specifically, for each iteration of the Monte Carlo simulation, the random classifiers and the naive classifiers predict the testing data set return states according to the mechanism mentioned above. In total, each iteration, the Monte Carlo process samples 4886 return states from the real data set mimicking the average number of stocks in each month of our entire sample. We iterate the simulation for 10,000 times. The accuracies for each classifier then allow us to conduct multiple comparison tests. We introduce Tukey's HSD and Table 3 demonstrates the Tukey's HSD test conducted with time period covering 199201:201912 as the testing sample. Tukey's HSD confirms that the naive classifier with OOS distribution knowledge provides the highest accuracy on average. We select the naive classifier as the benchmark to further test whether our models indeed provide information about the relation between future returns and historical information in the binomial tests.

Third and most importantly, a classifier which considers only the distribution of the return states mimics the investor behavior under the assumption of the market efficiency and the associated accuracy is the benchmark as implied by theory. An efficient market means that the future returns are unpredictable based on historical information and that investors are rational with enough information about the distribution of returns. To reflect the investor behavior, the inclusion of the benchmark based on a naive classifier that considers the distribution is necessary and meaningful. In fact, according to the rational investor assumption which leads to efficient market hypothesis, investors have consistent beliefs (Sargent 1993, Barberis and Thaler 2003). Consistent beliefs means that investors have correct information about the distribution they use to forecast unknown state variables. In other words, investors in efficient market must have enough information to infer the distributions of state variables, including the distribution of returns. Therefore, more precisely, we will want to include the accuracy based on the true OOS return state distribution as it is the best distributional information the rational investors can infer. We generally refer the accuracy provided by naive classifiers with distributional information as no information accuracy hereafter.

### **2.2.6 Binomial Test: A Joint Test**

The binomial test is a joint test and provides unique economic insights to finance machine learning literature and market efficiency literature. Specifically, given any one of our models built based on historical public information to predict future return states, if the prediction accuracy is tested as significantly greater than the no information accuracy, the statistical significance is indicative in at least 3 aspects. First, we can conclude that the prediction accuracy of the model is statistically better than the no information accuracy by the naive classifier. In other words, the combination of the model and the predictors delivers good predictive performance.

Second, since the naive classifier using only basic distributional information provides minimum or no information, the statistical significance of better accuracy indicates that the associated model delivers statistically meaningful information and this information is beyond the basic distributional information. Furthermore, because of using historical information and the specific modeling structure, the meaningful accuracy signals that there is a relation that exists between future return states and lagged predictors and that the relation is at least partially decoded by the related model. Third, if a binomial test presents significance and the associated portfolio strategy can generate profits, then the prices may not be correct. The significant predictability suggests that there exists a relation between the predictability and the information uncovered by the model using historical information. The predictability proves that the future prices will move towards predicted level and the profitability suggests that the piece of information uncovered by the model is useful. In other words, the historical information can generate future profits. As mentioned by Fama (1991), the return predictability does not necessarily mean that the market is inefficient. If the predictability does not allow the generation of profit, the market is efficient in terms of allocating resources and the prices correctly reflects in-

formation. Therefore, costs to question market efficiency, the predictability has to be combined with the meaningful profitability that cannot be offset by friction. We show that our models generate large economic gains for investors and the predictability is significant based on our modeling architectures and historical information.

An important implication of the profitability based on historical information is that the market by large does not share the information. If most of the market participants share this information generated by our models with historical information, the arbitrage process will erase the opportunity to benefit from the model predictions. Therefore, the information generated by our models provide new information about future returns and historical information. At the same time, at least to our models, the market price may not be correct historically for the CRSP-COMPUSTAT sample covering 196301:201912, since correct prices should not allow any investor to generate new information using historical information and apply the generated information to make OOS profit. In summary, the binomial test is a joint test of statistical significance of predictability and market efficiency. We discuss more about the use of the binomial test in Section 3.

## **2.3 Data Construction and Sample Split**

### **2.3.1 Data Components**

Our data universe contains 3,342,486 monthly stock information of 26,302 distinct common stocks with current returns listed on 3 major exchanges covering the time period of 196301:201912. We present the basic summary statistics of our data in Table 4. We construct a modeling sample of 332 lagged predictors. The lagged predictors include the return state, 101 firm characteristics, 2-digit SIC industry indicator, 2-digit SIC industry lagged returns, 125 macro indicators. We augment the macro indicators with 9 market specific predictors based on Goyal and Welch’s data set. We also sort the past 94 numeric firm characteristics and further augment the macro indicators with differences between the top decile median returns and the bottom decile median returns.

Specifically, we fully reconstruct the firm characteristics based on Green et al. (2014) with CRSP and COMPUSTAT. We made the data set to be a completely CRSP centric data with no data elimination if possible. In the end, we only eliminate rows with missing current returns and the securities that are not common stocks with SHRCD of 10, 11 or 12 listed on the major 3 exchanges with EXCHCD of 1, 2 or 3. Figure 1 presents our sample coverage of the CRSP universe with the counterpart reference level based on the entire CRSP database. Note that the reference level includes securities with missing returns and securities that are not common stocks.

Following Green et al. (2014) and Gu et al. (2020), we lag the annual firm characteristics by at least 6 months, we lag the quarterly firm characteristics by at least 4 months and we lag the monthly firm characteristics by at least 1 month. Note that the firm characteristics include variables depending only on historical returns, such as return momentum, and the variables depending on financial information and corporate announcement such as earnings and IPO. This data set of firm characteristics

provides a sound experimental environment for the test of market efficiency using machine learning methods. The lagged characteristics ensure that all the information in our models is historically publicly available information and there is no forward looking information leak in terms of the return based information nor the corporate announcements.

We also include macro time series indicators to allow the models learn about the relative intertemporal position of the market. We obtain the macro indicators of McCracken’s fairly new FRED-MD database from the website of Federal Reserve Bank of St. Louis. The database provides most of the mainstream macroeconomic indicators and is updated on a monthly basis (McCracken and Ng 2016). We retain 125 of the macro indicators in the end and exclude the variables ACOGNO, ANDENOX and TWEXMMTH due to limited availability. We obtain Goyal and Welch’s (2008) data set from Amit Goyal’s website. We include  $d/p$ ,  $d/e$ ,  $svar$ ,  $b/m$ ,  $nits$ ,  $corpr$ ,  $dfy$ ,  $dfr$  and  $ltr$  in our final data. We lag the time series macro indicators by at least 1 month.

For factor model tests, we obtain Fama and French’s (1992) MKT-RF, SMB and HML factors augmented with MOM factor of Carhart (1997). We collect these factors from Kenneth French’s data library. We also include the Hou, Xue and Zhang’s (2015) q 4 factors —  $R\_ME$ ,  $R\_IA$ ,  $R\_ROE$  as well as the  $R\_EG$  factor from the update of Hou, Mo, Xue and Zhang (2018). These factors are collected from Lu Zhang’s investment CAPM website. We include these pricing factors in the factor model tests mainly because of their good performance in the empirical asset pricing literature and their popularity.

### **2.3.2 Data Manipulation and Sample Split for Training and Testing**

We manipulate the data to obtain the most appropriate modeling inputs according to conventional machine learning model requirements. First, scaling variables is an important step to minimize the impact on weight adjustment and tree split because of the scale difference. In the machine learning literature, scaling variables is almost always recommended. We follow this standard practice. We scale the numeric firm characteristics cross-sectionally through normalization at in each trading month, if the field is not missing. We scale firm characteristics cross-sectionally for each time point to capture the relative characteristics difference at a time point across stocks.

Second, we also follow the conventional practice in machine learning literature and fill the missing values in categorical firm characteristics by specifying a new category. This practice does not bias the data set. Instead, it keeps the observations and helps the model to split more indicator specific effects from the general effects by retaining more observations. Third, after scaling, we fill the numeric firm characteristics with 0, which is the cross-sectional mean and median in each trading month. We also fill the time series variables, including the macro indicator variable consumer sentiment index (UMCSENTx), with the latest available values. Despite that there is the possibility of introducing noise, we fill the missing values because the possibility of including more useful information due to retaining more observations is also high. In the modeling phase, we scale the time series predictors through normalization with all available IS history such that the time series predictors can capture the

intertemporal relative position of the market trend and the macroeconomic trend.

To enlarge the testing power, our data covers the time window of 196301:201912, almost everything in the CRSP-COMPUSTAT universe. We split the entire data set in two ways corresponding to our time series tests and the purpose of our cross-sectional tests. For our time series tests, specifically, we want to test on intertemporal applicability of the models fitted with our time series in-sample training process in the OOS time periods. Therefore, we split the data in the middle based on the time length. As we have 57 years of data, we choose the end of the year 1991 as the splitting point as the earlier data has less number of stocks on average. Then, we have 2 time series subsamples, the first subsample covers 196301:199112 and the second subsample covers 199201:201912. Later, we refer the subsamples in terms of the coverage of time periods.

To take the advantage of the long time window, we train our models with the subsample covering 196301:199112 and make predictions with the subsample covering 199201:201912. We also train our models with the subsample covering 199201:201912 and make predictions with the subsample covering 196301:199112. We present the results of the OOS predictions with economic metrics both in subsamples and combined. We can use the latter 28 years to make prediction on the first 29 years, because we focus on extracting the relation between the historical public information and the future return states and our models do not depend on time structure. Therefore, we will not violate the setup of OOS tests.

It is obvious that the current working rules apply to the next period while it is not obvious whether the current working rules can apply to the past. Training models in latter dates and testing model in former dates makes a strong case of testing the time invariant applicability of the models. For example, after the appearance of Black Scholes and Merton (BSM) model, we know that BSM will certainly work for the pricing of stock options in the future after its appearance. On the other hand, it is hard to tell whether BSM would work to the dates prior to its existence. If a model captures the real process of return state transition in general, the model should be applicable to any time periods.

Our time series sample splits focus on the evaluation of time series OOS tests, while we position our cross-sectional sample to focus on overall explanatory ability across the entire time series. Unlike traditional econometrics, it is not ideal to test a machine learning model with the sample that the model is trained with. To overcome this issue, we split the data by odd number months and even number months in the spirit of Fama and French (2018). Our cross-sectional data split permits us to model the entire length of the data set across different macroeconomic conditions and test the fitted models in an OOS testing setup with completely new observations that are not seen by the models. We summarize the sample splits for training and testing in Table 4.

### **3 Time Series OOS Performance**

In this section, we present the OOS performance of our classification-based portfolios. We show that our portfolios perform systematically better than what the market portfolio based on buy-and-hold



strategy can deliver, regardless of the weighting schemes and the market capitalization cutoff points. We compare the performance in the aspects of average return, return volatility, risk-return tradeoff and shortfall. Our demonstration of the superior OOS performance of the classification-based portfolios answers Thaler's first question about whether we can beat the market. We further demonstrate the factor model explanation of the returns generated by the classification portfolios and show that the traditional factor models cannot explain the returns generated by our portfolios.

To carry out the OOS economic metrics, we first train our time series models using the first half of our sample covering 196301:199112 and make predictive classifications on return states in the second half of the sample covering 199201:201912. We then train on the second half of the sample covering 199201:201912 and make predictive classifications on return states in the first half of the sample covering 196301:199112.

As discussed in the sample construction section, we are able to use the latter half of the sample to predict the former half of the sample without any violation of the OOS test setup. By always using lagged information to score models, the models trained with the time series subsample covering 199201:201912 do not see any OOS observations, if we test the models with the time series subsample covering 196301:199112. We combine the OOS predictions in the two windows to form portfolios covering the entire sample period of CRSP-COMPUSTAT universe and provide by far the longest time series OOS test coverage.

As a robustness check, we provide performance tables of separated periods in the appendix along with the results from the portfolios enforcing strict market capitalization cutoffs. Neither splitting the full length of the CRSP-COMPUSTAT universe into two halves nor enforcing the market capitalization cutoffs changes our conclusion about the superior OOS performance of our classification portfolios comparing to holding the market.

After the demonstration of the good time series OOS performance that is time invariant, we dig into the statistical metrics of the accuracy that drives the good economic performance of our portfolios. We show that our portfolios can deliver balanced accuracy that is good comparing to the referencing levels of the class prevalence in our data. We further test the overall accuracy of our models against the no information accuracy of the naive classifier through the binomial test. The test against the no information accuracy as the null hypothesis is a conventional test in the machine learning literature. Following the important interpretation of the binomial test discussed in Section 2, we show that the OOS overall accuracy levels provided by the predictions from our models are statistically significantly larger than the no information accuracy, which indicates the generation of new information about OOS return states based on predictions from the models fitted in training process. We discuss the importance of this test to the understanding of market efficiency. We further look into the prediction accuracy by return state transitions and illustrate the opportunities captured by our models from learning with IS historical observations. In the end, we discuss the by-class statistical metrics, the training performance and the model selection during training.

### 3.1 Performance Evaluation with Economic Metrics

Because of the large number of models under consideration, we rely on visualization to summarize the major results and provide the performance metrics in table format in the appendix. Figure 2 and Figure 3 present the summarized economic metrics of the classification portfolios covering 196301:201912. The dashed blue line is the performance provided by the market portfolio, i.e. buy-hold strategy applied to the entire market. The long portfolio indicates a long position in all stocks predicted to be in the return state 10, or the best return state. The short portfolio indicates a short position in all stocks predicted to be in the return state 1, or the worst return state. The long-short portfolio includes a long position in the return of all stocks predicted in the return state 10 and a short position in all stocks predicted in the return state 1.

Figure 2 presents the OOS performance of the equal-weight portfolios and Figure 3 presents the OOS performance of the value-weight portfolios. Note that we change the weights in the long leg and the short leg separately when implementing a long-short portfolio. The allocation to the long leg and the short leg in a long-short portfolio is always equal, i.e., having a 50%:50% allocation ratio, across weighting schemes. All the returns are fully risk-adjusted against risk free rate and the cumulative returns are calculated as the gross returns net of the initial investment.

The first row of Figure 2 (Figure 3) demonstrates the distributional information of the equal-weight (value-weight) portfolio returns. The long portfolios deliver obviously better monthly average returns comparing to the market monthly average return. The short portfolios deliver obviously worse performance by taking the long position in the predicted worst performing stocks. When combined, the long-short portfolios deliver systematically better average monthly returns comparing to market return.

Over the full length OOS performance evaluation period covering 196301:201912, our neuron networks and tree models successfully distinguish the best performing stocks and the worst performing stocks. The similar conclusion holds in the comparison of the standard deviation. Our long-short portfolios deliver surprisingly small monthly return standard deviations. In the comparison in the return skewness against the market, our portfolios seem pushing the skewness towards positive values. This implies that our portfolios are more likely to realize large gains than large losses. However, if one is interested in short only strategies, this means more significant loss. The kurtosis plot demonstrate that our portfolios deliver returns that are more concentrated on a center, which indicates a reduction of the distributional fat tails. We augment the first row in Figure 2 and Figure 3 with Figure 4 and Figure 5 for equal-weight portfolios and value-weight portfolios respectively. Since the return distributions of the long portfolios are to the right of the return distributions of the short portfolios, Figure 4 and Figure 5 further confirm that our models can distinguish the future best return state stocks from the future worst return state stocks. It is also worth to note that the return distributions are closer to normal for our classification based long and short portfolios.

In Figure 2 and Figure 3, the second row demonstrates SR and CEQ of the full length OOS evalu-

ation covering 196301:201912. Because of holding the predicted worst performing stocks, short portfolios deliver negative SRs. The long portfolios based on complex models deliver SRs that are higher than what the market provides. Our long-short portfolios deliver surprisingly high SRs achieving systematically superior performance in term of the risk-return tradeoff. The best performing long-short portfolio based on our gradient boosting machine with the depth of 8 delivers a SR of 0.76, when we implement equal weight portfolio formation. The SR drops to 0.45, when we implement value weight portfolio formation, still more than doubled the SR delivered by the market. This drop is also less than 50 % mentioned by Chen et al (2020) in their findings when shifting from equal-weight portfolios to value-weight portfolios. The drop of performance is understandable as we cannot gain the returns from the small cap stocks with a 1:1 ratio. However, this does not mean that the portfolio performance is completely driven by the microcap stocks. In the appendix, we check the performance with strict capitalization cutoffs across all metrics, the performance of the sample eliminating the bottom 5 % capitalization stocks and the performance of the sample eliminating the bottom 10 % capitalization stocks are both comparable to the results covering the entire sample. Note that our SRs are calculated with monthly returns instead of annual returns and is not adjusted with modeling R squares, while Chen et al. (2020) use annualized returns in the calculation of SRs and Gu et al. (2020) use modeling R squares to adjust their SRs. Annualizing the returns can smooth the time series and shrink the size of the variance, while adjusting the SR with R squares can magnify the size of numerator and shrink the size of denominator. Both of these modifications can inflate the SR.

We also present the CEQ and overall cumulative returns in natural logs during the full length OOS evaluation period. In terms of economic investor gains, both the CEQs and the cumulative returns show the superior overall performance of our long portfolios and long-short portfolios comparing to what the market can generate during the 57-year investment window. The reason of the similar look of the CEQ subgraph and the cumulative return subgraph is that the construction of the log scale cumulative return is similar to the construction of CEQ. In the evaluation of equal weight portfolios, our best performing long-short portfolio based on our ANN model of 1 hidden layer and 128 neurons deliver a full length period cumulative return of 1,074,286,700 % net of the initial cost. In the evaluation of value weight portfolios, the cumulative return of the long-short strategy drops to 5,210,500 %, which is still 76 times of what the market achieves with value weights. We augment the cumulative return evaluation in our Figure 2 and Figure 3 with Figure 6 and Figure 7. As demonstrated, the cumulative returns of our long-short portfolios surge to a level that both the market cumulative returns and the cumulative returns from holding the predicted worst performing stocks become flat lines.

After the surprisingly good performance demonstrated in the return distributions, the risk-return trade-offs and the economic gains, we look at the largest shortfalls associated with our classification portfolios. We evaluate the shortfalls with the maximum drawdown. The short portfolios are associated with the smallest maximum drawdown. Our long portfolios and the market are associated with a similar level of maximum drawdown. Our long-short portfolios deliver a surprisingly low level of maximum drawdown. In the equal-weight implementation, our best performing long-short portfolio in

term of maximum drawdown is again the long-short portfolio based on the gradient boosting machine GBM8 100. It achieves a maximum drawdown of 7.9 %, only 1/4 of the market maximum drawdown, which documents the historic selloff around the Black Monday in 1987. When switching to value weights, the maximum drawdown is still meaningfully lower for the long-short portfolio based on GBM8 100.

### 3.2 Explanatory Analysis with Factor Models

We further provide alpha as an additional metric of the OOS performance evaluation and look into whether the performance can be driven by common market factors. We regress the return time series of our classification portfolios on the popular factors. We include Fama French 3 Factor model (FF3F), Fama French 3 Factor + MOM model, the original q-factor model (q4) of Hou, Xue and Zhang (2015) and the q5 factor model of Hou, Mo, Xue and Zhang (2018) with the R\_EG factor.

Figures 8 and 9 summarize the factor model tests with our equal weight and value weight portfolios respectively. The orange dashed line indicates the position of zero. The first row shows the alphas, the second row shows the t values associated with the alphas and the last row shows the R square values associated with the factor models. In the case of equal-weight portfolios, clearly, FF3F models and FF3F + MOM models cannot explain any of the long-short portfolios based on our classification models, not even the long-short portfolios generated by our underfitted tree models with a maximum depth of 2. Despite that the q factor models perform slightly better, the q factor models can barely explain only the returns generated by the long-short portfolio based on the clearly underfitted DART2 100. The conclusion with the long only portfolios are similar. In fact, the alphas on our portfolio returns are associated with the t values as large as 20 in the equal-weight implementation of the portfolios. The R squares in many cases are around 1 % indicating the associated factor models do not explain the variation. The factor models' ability to explain the long-short portfolio returns does not get any better systematically when we shift from equal weight to value weight. We provide the full testing results in the appendix.

### 3.3 Model Complexity and The Economic Performance

Gu et al. (2020) mention that in their implementation targeting the numeric values, they found the shallow learning outperform deeper learning. They document that the performance of the neuron networks peaks at three hidden layers and their tree models tend to select trees with few leaves. Our performance evaluation with the economic metrics shows different conclusion.

First, the performance of the neuron network is largely influenced by the total capacity of neurons instead of just the number of layers. Controlling the number of layers, higher number of neurons improve the performance. In the factor tests, we do not see significant deterioration of performance when we increase the number of layers, either. Second, for our tree models, there is an obvious improvement trend in all of our subgraphs in Figure 2, Figure 3, Figure 8 and Figure 9. As we

increase the maximum depth of our tree models, the performance of our tree model based portfolios improves substantially. Specifically, with maximum depth of 2, our tree models seem significantly under-fitted and the associated long-short portfolios can deliver performance metrics that are worse than the counterparts of the market, while our tree models with depth of 8 can help form the best performing portfolios. Third, for each of our models, as we increase the model complexity, we can see a widened gap between the long portfolios and the short portfolios across figures. This shows that the models can better distinguish the future best return state stocks from the future worst return state stocks as we increase the model complexity.

### **3.4 Performance Evaluation with Statistical Metrics**

We have demonstrated that the classification portfolios can beat the market and the superior performance of the classification portfolios in a range of economic metrics. In this subsection, we analyze the classification models behind the good economic performance through statistical metrics. We first show the OOS classification accuracy achieved with our models and test the OOS classification accuracy against the true OOS no information accuracy delivered by the naive classifier. The no information accuracy by the naive classifier is under the assumption that the response variable, the future return state, is independent of the lagged predictors. We also dig into the by-class level accuracy, analyze the overall return transitions and provide new insights about market efficiency. We show that our machine learning models take the advantage of the difference in the information uncertainty to deliver superior statistical performance. We base our discussion in this subsection on the combined OOS predictions covering the entire time series. We provide the tables for separate periods in the appendix for robustness check.

#### **3.4.1 Accuracies and Binomial Tests Against the No Information Accuracy**

The accuracy metric is a unique and natural statistical metric that is available to the classification problems. In numeric predictions, one can only construct metrics based the predictions errors as it is impossible to measure accuracy directly. However, in the classification problems, all of the predictions are similar to combinatorial problems as the predictions are by nature similar to choices and selections. This creates a hardline on whether a predicted choice or selection is correct or incorrect. Column 2 in Table 6 presents the OOS accuracies of our models that predictively decide which stock is more likely to be among the best performing stocks and which stock is more likely to be among the worst performing stocks. We also provide Kappa statistics associated for each of the models in Column 2. It is clear that the models delivering good economic performance are systematically better performing in term of the OOS prediction accuracies as well. For example, the underfitted DRF2 200 model delivers an accuracy level of 15.09%, while the GBM8 100 model, whose associated portfolios are among the best performing portfolios, delivers an accuracy level of 15.72%. The performance rank is similar when we compare the models with Kappa. For example, ANN1 16 delivers a Kappa of 4.88%

while a better ANN model, ANN1 128, whose equal-weight long-short portfolio delivers more than one billion percent net cumulative return, has a Kappa of 5.26%.

Beyond the direct interpretation of the statistical performance, we adopt the classification framework to take the advantage of the accuracy metric as a possible direct proxy that can be used to conduct a statistical test on whether we can foresee the future with historical information. We further introduce the binomial test to the efficient market hypothesis literature and the financial machine learning literature. As mentioned in Section 2, since the accuracy metric is in a natural form of proportion, we can compare the statistical difference between 2 accuracies through a binomial test. Specifically, we test the predictive classification accuracies against the accuracy delivered by the naïve classifier as the null hypothesis.

The binomial test on accuracy has profound meaning in our setup. First, it tests the statistical meaningfulness of the accuracy delivered by a classification model. If a model delivers an OOS accuracy that is statistically significantly higher than the no information accuracy, the model captures statistically meaningful essence of the relation between the future return states and the historical information. We show in Table 6 that all of our models deliver statistically significant OOS accuracies in tests against the no information accuracy.

Second, to the finance literature, the accuracy is a natural statistical metric to evaluate the market efficiency and has special meaning to us as discussed in Section 2. Economically, if the market is efficient in the strong form, then the prices fully reflect all available information, regardless of whether the information is private or public. In this situation, we should not be able to benefit from trading stocks using any historical information. If the market is efficient in the semi-strong form, then the prices fully reflect all historical public information. An investor can only benefit from trading stocks using the private information. In other words, even if we can generate profits in the semi-strong form efficient market, the historical public information should not contribute to the prediction of the returns and only the private information can contribute to the prediction of the returns. If the market is in the weak form of efficiency, we can benefit from trading on the information inclusion as the public information, including the corporate announcements and the historical macroeconomic information, will be incorporated to the prices gradually. However, if the market is efficient in the weak form, past returns should not contribute to the prediction of the future returns. If the market is not efficient at all, the returns are predictable and all types of information can make contribution to the prediction of returns.

In summary, the binomial tests on the prediction accuracies, coupled with the analysis of the predictor contribution, can provide direct indications about whether there exists violation of the 3 forms of market efficiency. Statistically, we should not observe any meaningful OOS accuracy in our predictions, if the market is efficient in the strong form. However, if the prediction accuracy of a model is statistically significantly higher than the no information accuracy, we can conclude that there exists some relation between the historical information and the future return states and the relation is captured by the machine learning model. In other words, the statistical significance of the accuracy

can bring up at least the questions about the market efficiency in the strong form. Therefore, what we present with the binomial test is a direct test of the strong form market efficiency. The contribution of historical predictor variables in the predictability can further help us understand whether we should further question the semi-strong form and the weak form of market efficiency.

We present the binomial test results against the OOS no information rate with the last 4 columns in Table 6. As mentioned in Section 2, the no information accuracy measures the best prediction one can make if the predictors do not have any relation with the response variable in the classification problem. We make it stricter by using the real ground truth OOS no information accuracy based on the OOS return state distribution.

The binomial test through accuracy against no information accuracy is a test about whether the response variable and the predictors are independently distributed. Therefore, the no information accuracy provides a natural tool for us to check if our models are providing OOS information predictively based on the historical information including return information, corporate announcements and macroeconomic indicators. In an ideal case, we want to split our returns cross-sectionally by trading month into evenly distributed classes. However, because we have the thresholds of the quantiles when splitting our returns, each return state includes a slightly different number of stocks cross-sectionally within each date and the difference remains there when we look at the entire time coverage. 10.16% is the portion that the return state 7 takes in our entire sample across 196301:201912 and the return state 7 makes the largest portion of our sample among all return states. Therefore, after observing the entire distribution of the return states, an investor should always bet on return state 7, if the return state transition is truly independent of any of the historical information associated to EMH.

It is obvious that the prediction accuracies delivered with the models trained IS based on historical information successfully surpass the ground truth no information accuracy. With the provided 95% confidence intervals and the p values, it is clear that the accuracies delivered by our models are statistically significantly higher than what the ground truth return distribution can deliver. Without looking into the contribution of predictors, which will be discussed in Section 4, our findings in Table 6 have two important implications. First, statistically, we confirm that there is some relation between the OOS return states and the IS lagged predictors, including the historical return information, the historical corporate announcements and the historical macroeconomic indicators. In other words, as the response of the model, the future return states are not independent of the historical information. Second, economically, we confirm that the market efficiency can be improved as the return states are predictable at a significant enough level through the existence of the relation between the OOS future return states and the historical information. The predictability shows that the historical information is not fully reflected by the prices.

Based on the fact that our models deliver significant OOS accuracies, assuming that between two investors, one investor uses our models and the other investor believes that the return states are independent of historical information and thus relies on the OOS true distribution information of return states to make predictions. Our finding means that the investor using our model can then make OOS

return state predictions that is better than the best guess that the other investor can make after she sees the entire ground truth OOS distribution of the return states.

Because of the relation between historical information and the OOS future return states, our models essentially generate new information that the prices do not reflect. In other words, the investors using our models have information advantage and manually introduces de facto information asymmetry to the market. In practice, this implies that the sophisticated investors, such as Renaissance Technologies, can take the information advantage by generating new information predictively about the future return states based on historical information. In general, the generated information, depending on the choice of analytical tools and the sophistication levels of the investor, may not be publicly available even if the other investors can observe the ground truth distribution of the future return states. We discuss more about the semi-strong form market efficiency and the weak form market efficiency in Section 4 based on the predictor contributions.

### **3.5 Return State Transition Uncertainty and the Prediction Strategy**

We next assess the by-class statistical performance of the OOS predictions. We first present the ground truth about return state transitions and discuss the relative uncertainty across the return state transitions. We then evaluate the prediction accuracy level for each individual return state transition achieved by our classification models in general. We demonstrate the strategy that our models take to achieve the superior OOS performance and discuss the easiness for predictions. We further present the by-class OOS performance for each model with popular classification related statistical metrics in the end of the section.

#### **3.5.1 Return State Transition Uncertainty, By-Transition Performance and Implication on Market Efficiency**

Table 7 presents the information of the real return state transitions. From Panel A, we can see that the return state transition probability based on the entire sample covering 196301:201912 is not evenly distributed. In other words, the uncertainty about the return state transition based on different current state is not the same. In general, the corner transitions in the matrix, such as the transition from the current worst performing state to the future highest performing state, are associated with substantially higher probabilities. While the center transitions of the matrix are associated with more evenly distributed probability around 10%. It is also clear that the transitions from the current middle level performing return states to the future extreme performing states are with relatively lower probabilities.

Panel B presents the monthly mean returns of the associated return state transitions of all stocks in our sample covering 196301:201912. We can see that the better future performing states are with better cross-sectional average returns. However, the panel shows that the transitions from the current middle level return states, i.e. return state 2 to return state 9, do not really deliver very different average returns when they transit from current return state to the new return state. For example, the average



return for a stock transiting from return state 3 to return state 9 is not very different from the average return of a stock transiting from return state 4 to return state 9. However, the corners of the mean return table shows different situation. A stock transiting from return state 1 to return state 10 delivers significantly higher average return comparing to a stock transiting from return state 2 to return state 10. Both Panel A and Panel B show that the extreme return transitions are where the opportunity lays for the investors who want to generate excess trading profit. We show that this is exactly the strategy of our models.

Table 8 presents the average accuracy of time series OOS prediction across our models. It is obvious that our models put high stake in the extreme return state transitions. The transitions to the lowest return state obtain the highest overall accuracy levels in OOS prediction. The models also pay more attention to the transitions to the middle return states and the best performing return states. The OOS predictions from our models seem successfully capturing the better certainty of the transitions to the extreme states and the transitions to the middle performing states. This implies that our classification models believe that the extreme return states and the transitions to the middle return states are more predictable. This is consistent with our observation on the true transition probability matrix.

As our models capture this difference in the levels of uncertainty during IS training process and are able to tell what return state is more predictable, the OOS prediction accuracy by return state transition shows the models' inclination about what stocks are priced more efficiently under the current market condition. Our models gain from the inefficiently priced stocks, because those stocks deliver more certainty about return state transition through unevenly distributed transition probabilities. Based on Tables 7 and 8, our models clearly have systematic preference of betting on the extreme transitions and the transitions to the middle return states. Therefore, specifically, the accuracy by return state transition, combined with the ground truth difference in the probability of transition, indicates that the extremely priced stocks and the middle performing stocks are creating a market segment that is less efficient in terms of current pricing.

### **3.5.2 By-Class Statistical Metrics in OOS Predictions**

We discuss the training and the testing by-class statistical metrics in this subsection. In addition to the exploitation of the uneven distribution of the return state transitions, our models attempt to balance the accuracy for each individual classes between the true positive predictions and the true negative predictions. Due to the specific purpose of our classification models, balancing between the true positive predictions and the true negative predictions is challenging. During the training process, our models create one-versus-all multi-class classification structure, similar to one-hot encoding, and compare each individual return state against the other 9 return states. This introduces the unbalancedness. To the classification models, when looking at the data set with any one of the 10 return states, the data set will have 90 % negative rate. By simply labelling all the data points as negative, the model can achieve

more than 90 % accuracy level for each single return state. However, if we denote negativity with 0, we will be predicting all 0's across the 10 classes. In the end, we will have low overall accuracy and this adds no information at all. Therefore, balancing the true positive rate and the true negative rate becomes crucial for the success of our models.

Table 8 summarizes the key statistical metrics that measure the performance of our classification models in the OOS predictions covering 196301:201912 by class. With the detailed by-class metrics, we confirm that the models sacrifice the accuracy and the detection of the true positives for the return state 2, return state 3 and return state 4, regardless of the modeling architecture, potentially due to the unbalancedness of the data and the different level of uncertainties as discussed above. The prevalence column shows the distribution of the ten classes across the entire OOS testing period covering 196301:201912. As shown, the return state 2, 3 and 4 are among the return states with the lowest number of observation in our entire sample. However, it is important to note that the return state 1 is also among the lowest populated states while the models collectively choose not to sacrifice the accuracy in return state 1. As reflected by the better accuracy across the models for return state 1, the models seem spending more resources to improve the accuracy of prediction in return state 1. This again emphasizes that the uncertainty level of the return states are not equal. The extreme return states are more certain than the middle return states and the models are taking the advantage of this higher certainty. If we look at the summarizing measures that balance the accuracy between the true positives and the true negatives, the balanced accuracies for each of the return states are above 50 % indicating that there are gains beyond simply predictively labeling the individual return state as negative.

### **3.5.3 Training and Model Selection**

With the good performance in OOS economic metrics and statistical metrics, we can learn about the applicability of the models in real practice by looking into the training process and see if we can select the correct models to generate the guideline for portfolio allocation. We briefly discuss performance metrics of the training process in this subsection.

Figures 10 and 11 show the economic performance of equal-weight portfolios for the 2 time series training sets covering 196301:199112 and 199201:201912 respectively. We provide the counterparts of the value-weight portfolios in the appendix. The IS economic performance of our classification portfolios are similar to their positions in the OOS economic performance evaluation. Our best performing portfolios are still the models with more complex structures, such as GBM8 100. Table 10 summarizes the statistical metrics of the IS performance of our models for the 2 time series training sets. The accuracies of our models are substantially higher than what we have for the OOS evaluation, which is expected. The best performing models are also the models with more complex structures. The GBM8 100 can deliver accuracy levels of 23.16 % and 22.55 % in the two training sets respectively. The good performance and the consistency between our IS and OOS models ensure the applicability of our models in real practice. In other words, through the training evaluation, we can rank the

performance of our models correctly and choose the selected model to make OOS predictions.

## **4 Cross-Sectional Explanatory Power and Predictor Contribution**

Next, we discuss the overall explanatory power of our models. We construct a separate sample to test the model structures cross-sectionally. We demonstrate the average rank of variable importance across the two types of modeling architectures, i.e., the neuron networks and the tree models. We connect the predictor contribution to the good performance of our models in the cross-sectional OOS evaluation and further discuss the new insights about EMH based on the variable contributions.

### **4.1 Cross-Sectional Explanatory Power of the Models**

We conduct cross-sectional (CS) tests to examine whether our models can capture the overall return state changes. To do this, we split the CRSP-COMPUSTAT universe into two subsamples including odd number months and even number months in the spirit of Fama and French (2018). We use odd number months as our training sample and make CS OOS tests with even number months where the observations are new to the models. This setup for the IS training and the OOS testing enables us to directly look at the overall explanatory performance of our models across the entire coverage of CRSP-COMPUSTAT universe from 1963 to 2019.

Figures 12 and 13 show the CS OOS performance with economic metrics calculated based on the testing set of our CS sample split, i.e., the even number months that the models have never seen during the IS training process with the odd number months. The performance in CS OSS tests are similar to what we observed in the TS OOS tests. First, our classification based long-short portfolios systematically outperform the market portfolio. Second, we observe the similar increasing trend of performance as we adjust the complexity of the models. Third, the best performing models in terms of the CS OOS economic metrics match the best performing models in TS OOS tests. We see that the single layer neuron network with 128 neurons, ANN1 128, perform the best among the neuron networks and GBM8 100 and DART8 100 perform the best among the tree models.

We also include a table of overall accuracy of the classification. We can see in Table 11 that our model prediction accuracy levels are still higher than the no information accuracy in CS OOS predictions. This means that our models overall obtained information about the OOS return transitions based on looking at the IS observations. In other words, our models can explain the relative performance of stock in a traditional cross-sectional setup to the degree of accuracy that are statistically meaningful.

### **4.2 Variable Importance across Models**

With the good IS and OOS performance demonstrated through both the TS setup and the CS setup, we further look into the predictor contributions measured as variable importance across models and discuss what are the driving predictors that leads to superior performance of our models across the

different setups. We summarize the variable importance with the average values, including the rank of the contribution to each individual models. We separately demonstrate the variable importance for the neuron networks and the tree models, since the models have structural differences in dealing with categorical inputs. We discuss the variable importance of the CS training models presented in Table 12 and include the variable importance of the TS training models in the appendix.

Unlike what is documented by Gu et al. (2020) or Chen et al. (2020), the top contributor in both the neuron networks and the tree models is the idiosyncratic volatility. The monthly average of the daily bid-ask spread divided by average of daily spread makes the second important contribution across the modeling architectures. It is worth to mention that the tree models seem highly dependent on idiosyncratic volatility, bid-ask spread and return volatility, while the contribution made by various top contributing predictors in neuron networks are more balanced. Beyond the top contributors, across the modeling architectures, all types of historical information make important contribution. Specifically, we see that other trading related variables and the industry indicators make huge portion of contribution. At the same time, the historical corporate announcements, such as earning price ratio, IPO status, convertible debt obligation, firm R&D expenses, change of number of analysts, etc., all make substantial contribution to our models. Macroeconomic indicators are also among the top 50 contributing predictors.

The fact that the historical information, including return related information, corporate announcements, can contribute to the model predictability is interesting to the understanding of the market efficiency. The semi-strong form of market efficiency permits the generation of excess profit through information asymmetry and lowers the bar to focus on the speed of which the public information is incorporated to the prices. However, the contribution from corporate announcements coupled with the strong predictability across the setups show that there exists systematic relation between the future returns and the lagged corporate announcement variables. Considering that some of the variables, such as R&D are lagged by at least 6 months, the semi-strong form of market efficiency seems questionable. If the market is efficient in the semi-strong form, it is hard to explain why the corporate announcements from 6 months ago can still help predict the future return states.

In addition, the weak form of market efficiency states that the historical prices and trends cannot predict the future returns. Yet, the top ranked contributors to our models are populated with the past return and trading information. In fact, while the contribution from corporate announcements and the contribution from the past return information have a ratio of 50% : 50% in the tree models, the past return information makes more contribution to our neuron network models comparing to the corporate announcements. This shows that the historical return information are able to help predict the future return states for the individual stocks and thus questions the weak form of market efficiency.

## 5 Conclusion

In this paper, we introduce the machine learning classification methods to the asset pricing literature and examine market efficiency. Taking the advantage of the relation between classification and the information theory, we force the models to extract the information about the relation between the historical information corresponding to the different forms of market efficiency.

We analyze the economic performance of our classification portfolios in terms of OOS return distribution, SR, CEQ and maximum drawdown. Our classification based portfolios beat the market systematically in multiple setups. The best models demonstrate surprisingly good performance across the metrics. We also measure the OOS performance with statistical metrics. We see that the OOS prediction accuracies match the OOS economic performance of the associated classification portfolios. We take advantage and utilize the accuracy metric, which is only applicable to classification problems, as a proxy to study whether the future returns of individual stocks are independent of multiple types of historical information. We introduce the binomial test to the asset pricing literature and document the statistical significance of the classification model accuracy against the no information accuracy.

Our findings about the statistical significance has important implications. First, the accuracy indicates that there are meaningful relation between the future return states and the lagged predictors representing historical information. In other words, the prediction accuracy is statistically meaningful and the future return states are predictable. This is the first time in the literature that the market efficiency is examined through the prediction accuracy as the proxy. Second, our classification models successfully capture the relation between the historical information and the future return states in a predictive format, indicating that the new information beyond the distribution of the return states has been generated with our classification models. This demonstrates the possibility that sophisticated investors can apply complex tools, such as the machine learning classification methods, to generate new information that is not reflected by the market prices. Specifically, the generation of new information about future return states by the sophisticated investors is equivalent to manually introduce the information asymmetry to the market. The sophisticated investors can take the information advantage against the other investors and benefit from it in their trading activities.

We also document important findings on the transitions of the return states. The ground truth return state transitions during the time period of 196301:201912 show uneven levels of uncertainty. The extreme state related transitions are with substantially higher certainty comparing to the other transitions. We show that our models learn about this difference of uncertainty and take the advantage of it. The accuracy by return state transitions and the uneven uncertainty we demonstrate collectively imply the different levels of market efficiency related to different return state transitions.

In term of the contribution of the predictors, we show that historical information including return related information, corporate announcements, macroeconomic indicators, etc. all make important contribution to our models. The fact that the historical information including return information and corporate announcements being able to make contribution to the predictability challenges the weak

form of market efficiency and the semi-strong form of market efficiency. Combined with the implication that the sophisticated investors may be able to generate new information based on historical information, we conclude that there is still room for the market efficiency to improve.

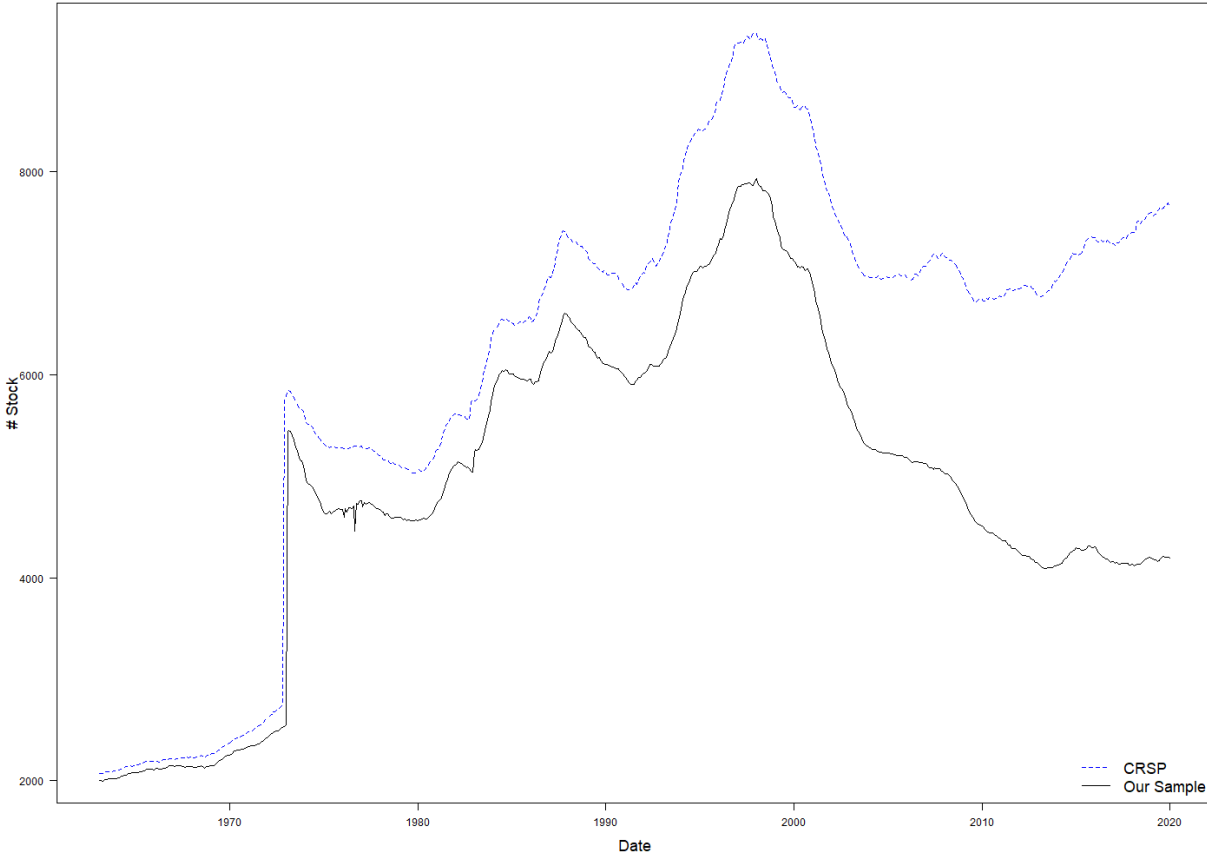
## References

- Barberis and Thaler, 2003, A Survey of Behavioral Finance, ch. 18, p. 1053-1128 in Constantinides, Harris and Stulz eds., *Handbook of the Economics of Finance*, vol. 1, Part 2, Elsevier.
- Bianchi, Daniele, Buchner, Mathhias and Tamoni, Andrea, 2020, Bond Risk Premia with Machine Learning, Queen Mary University of London, University of Warwick and Rutgers working paper.
- Brandt, Michael, Santa-Clara, Pedro, and Valkanov, Rossen, 2007, Parametric Portfolio Policies: Exploiting Characteristics in the Cross-Section of Equity Returns, *Review of Financial Studies* 22(9), 3411-3447.
- Bryzgalova, Svethana, Pelger, Markus, and Zhu, Jason, 2020, Forest Through the Trees: Building Cross-Sections of Stock Returns, London Business School and Stanford University working paper.
- Carhart, Mark M., 1997, On Persistence in Mutual Fund Performance, *Journal of Finance* 52(1), 57-82.
- Chen, Luyang, Pelger, Markus and Zhu, Jason, 2020, Deep Learning in Asset Pricing, Stanford University working paper.
- Cohen, Malloy and Nguyen, 2020, Lazy Prices, *Journal of Finance* 75(3), 1371-1415
- DeMiguel, Victor, Garlappi, Lorenzo, and Uppal, Rama, 2009, Optimal Versus Naïve Diversification: How Inefficient is the 1/N Portfolio Strategy? *Review of Financial Studies* 22(5), 1915-1953.
- Fama, 1969 (May, 1970), Efficient Capital Markets: A Review of Theory and Empirical Work, *Journal of Finance*, 25(2), Papers and Proceedings of the Twenty-Eighth Annual Meeting of the American Finance Association New York, N.Y. December, 28-30, 383-417.
- Fama, 1991, Efficient Capital Markets: II, *Journal of Finance*, 46, 1575-1617.
- Fama, Eugene, and French, Fama, 1992, The Cross-Section of Expected Stock Returns, *Journal of Finance* 47(2), 427-465.
- Fama, Eugene, and French, Fama, 2018, Choosing Factors, *Journal of Financial Economics* 128(2), 234-252.
- Feng, Guan hao, Polson, Nicholas G., Xu, Jianeng, 2019, Deep Learning in Characteristics-Sorted Factor Models, University of Chicago working paper.
- Goyal, Amit, and Welch, Ivo, 2008, A Comprehensive Look at the Empirical Performance of Equity Premium Prediction, *Review of Financial Studies* 21(4), 1455-1508.
- Green, Jeremiah, Hand, John and Zhang, Frank, 2017, The Characteristics that Provide Independent Information about Average U.S. Monthly Stock Returns, *Review of Financial Studies* 30(12), 4389-4436.
- Grossman and Stiglitz, 1980, *American Economic Review*, 70(3), 393-408.
- Gu, Shihao, Kelly, Bryan, and Xiu, Dacheng, 2020, Empirical Asset Pricing via Machine Learning, *Review of Financial Studies* 33(5), 2223-2273.
- Hou, Kewei, Mo, Haitao, Xue, Chen and Zhang, Lu, 2019, Which Factors? *Review of Finance* 23(1), 1-35.

- Hou, Kewei, Xue, Chen, and Zhang, Lu, 2015, Digesting Anomalies: An Investment Approach, *Review of Financial Studies* 28(3), 650-705.
- McCracken, Michael, and Ng, Serena, 2016, FRED-MD: A Monthly Database for Macroeconomic Research, *Journal of Business & Economic Statistics* 34(4), 574-589.
- Moritz, Benjamin and Zimmermann, Tom, 2016, Tree-Based Conditional Portfolio Sorts: The Relation between Past and Future Stock Returns, Ludwig Maximilian University of Munich and University of Cologne working paper.
- Rau, 2011, Market Inefficiency, ch. 18, p. 331-349 in Baker and Nofsinger eds., *Behavioral Finance: Investors, Corporations and Markets*, Wiley.
- Rossi, Alberto, 2018, Predicting stock market returns with machine learning, Georgetown University working paper.
- Sargent, 1994, *Bounded Rationality in Macroeconomics*, A Clarendon Press Publication.
- Shannon, Claude, 1948, A Mathematical Theory of Communication, *Bell System Technical Journal* 27(3), 379-423.
- Wolff, Dominik, and Echterling, Fabian, 2020, Stock Picking with Machine Learning, Darmstadt University of Technology and Deka Investment GmbH working paper.



### Total Number of Stocks in 3 Major Exchanges 196301:201912



**Figure 1 Number of Stocks in CRSP vs Number of Stocks in Our Sample 196301:201912**

Figure 1 presents a comparison of the sample coverage between our data set and the CRSP database. The dashed line represents the number of securities included in the CRSP database and the solid line represents the number of stocks included in our sample. Note that CRSP is a general security database. It includes securities other than stocks of the public firms. In our sample, we include only the stocks listed on NYSE, Amex, and NASDAQ. This figure presents the comparison from January 1963 to December 2019. In total, our sample covers distinct 26302 stocks. On average, our sample covers around 4887 stocks for every trading month. The detailed summary statistics of the sample coverage can be found in Table 4.

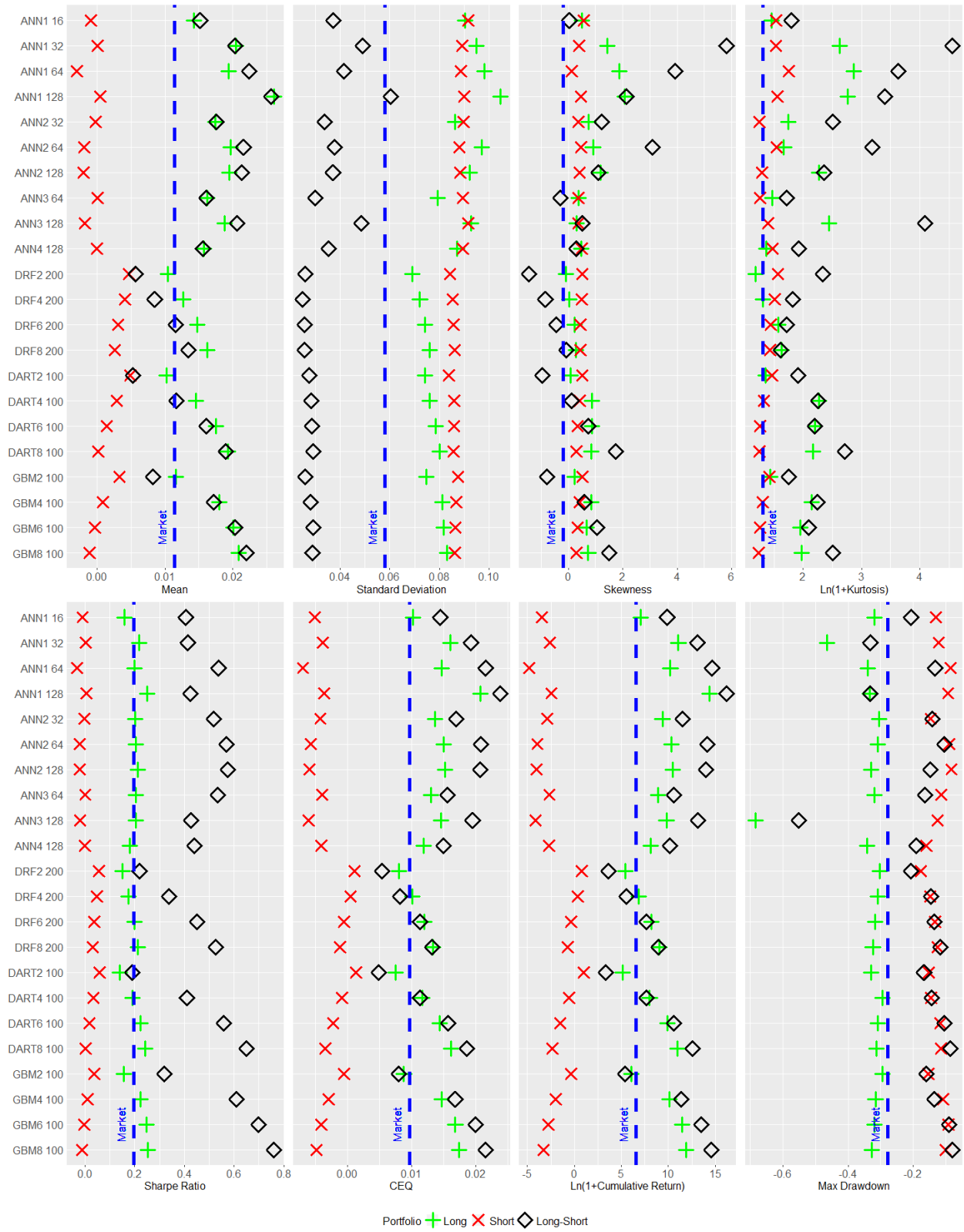


Figure 2 Equal Weight Time Series OOS Portfolio Economic Performance 196301:201912

## Figure 2 (Continues)

Figure 2 summarizes OOS economic metrics of the classification-based portfolios covering 196301:201912 with the equal-weight scheme across different classification models. More specifically, a long portfolio represented by “+” indicates a *long position* in all stocks predicted to be in the return state 10, or the best return state. A short portfolio represented by “X” indicates a *long position* in all stocks predicted to be in the return state 1, or the worst return state. A long-short portfolio represented by “0” indicates a *mixed portfolio including a long position in the return of all stocks predicted in the return state 10 and a short position in all stocks predicted in the return state 1*. In each of the subgraphs, the dashed line indicates the reference performance delivered by the market portfolio, i.e. buy-hold strategy applied to the entire market. The OOS performance of the stocks is equal weighted. The allocation to the long leg and the short leg in a long-short portfolio is always equal, i.e., having a 50%:50% allocation ratio. All the returns are fully risk-adjusted against risk free rate. We measure the plain Sharpe Ratio (SR) as return scaled by standard deviation. We follow the literature and adopt  $\gamma = 1$  in the calculation for CEQ. The cumulative returns are calculated as the gross returns net of the initial investment. Details about the economic metrics are discussed in Section 2.2.1 and the model specifications are presented in Table 2.

The equal weight long-short portfolios based on our classification models systematically outperform the market. As the modeling complexity increases, the performance of our models increases. Our equal weight long-short portfolios deliver surprisingly good performance in term of controlling left-tail risk while delivering the amazing level of economic performance. Specifically, from 196301:201912, the best long-short portfolio across the sample period based on our tree model, GBM8 100, deliver a SR of 0.76 which is comparable to the plain SR of 0.707 achieved by the best portfolio in Gu, Kelly and Xiu (2020) as in their appendix table A.9 (their OOS data, or testing sample, covers from 1987:2016 and they have year-by-year rolling model updating). More surprisingly, the maximum drawdown associated with the long-short portfolio of GBM8 100 is only at 7.88 %. During the same time period, the market portfolio achieves a SR of 0.196 and a maximum drawdown of 27.78 %.

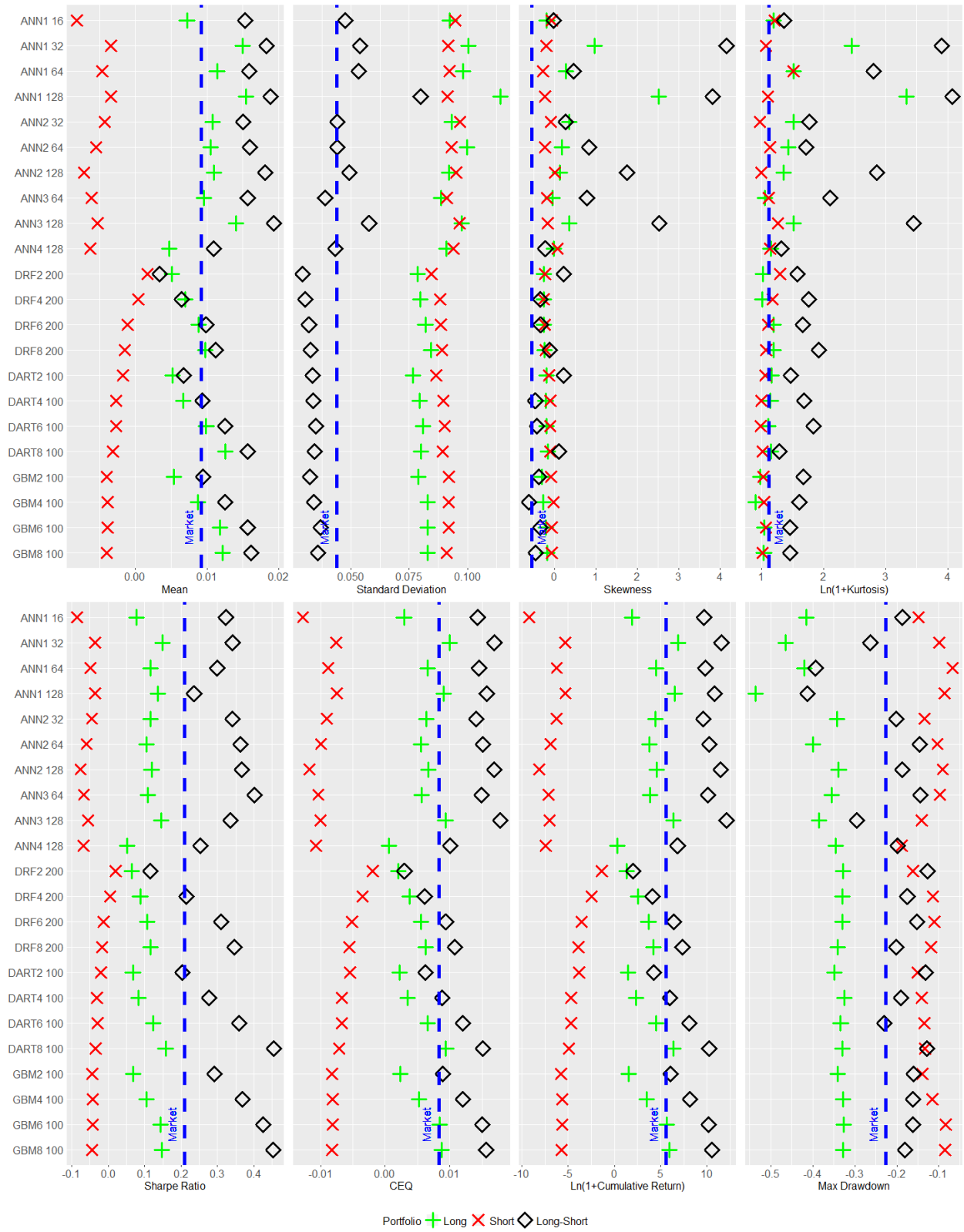


Figure 3 Value Weight Time Series OOS Portfolio Economic Performance 196301:201912

### Figure 3 (Continues)

Similar to Figure 2, Figure 3 summarizes OOS economic metrics of the classification-based portfolios covering 196301:201912 with the value-weight scheme across different classification models. In general, the change of weight scheme does not affect the OOS performance of our classification portfolios.

Specifically, a long portfolio represented by “+” indicates a *long position* in all stocks predicted to be in the return state 10, or the best return state. A short portfolio represented by “X” indicates a *long position* in all stocks predicted to be in the return state 1, or the worst return state. A long-short portfolio represented by “◇” indicates a *mixed portfolio including a long position in the return of all stocks predicted in the return state 10 and a short position in all stocks predicted in the return state 1*. In each of the subgraphs, the dashed line indicates the reference performance delivered by the market portfolio, i.e. buy-hold strategy applied to the entire market. The OOS performance of the stocks is value weighted. The allocation to the long leg and the short leg in a long-short portfolio is always equal, i.e., having a 50%:50% allocation ratio. All the returns are fully risk-adjusted against risk free rate. We measure the plain SR as return scaled by standard deviation. We follow the literature and adopt  $\gamma = 1$  in the calculation for CEQ. The cumulative returns are calculated as the gross returns net of the initial investment. Details about the economic metrics are discussed in Section 2.2.1 and the model specifications are presented in Table 2.

Similar to the equal weight version of the OOS tests in economic metrics, the value weight long-short portfolios based on our classification models systematically outperform the market. As the modeling complexity increases, the performance of our models increases. Our value weight long-short portfolios also deliver good performance in term of controlling left-tail risk while delivering the amazing level of economic performance. Specifically, for the OOS period from 196301:201912, our value weight long-short portfolios based on our tree models, DART8 100, GBM8 100, GBM6 100, deliver OOS SRs of 0.456, 0.453 and 0.425, comparable to the best plain SR of 0.39 achieved by Gu, Kelly and Xiu (2020) as demonstrated in their Table 7. The maximum drawdown of the long-short portfolio based on DART8 100 is 12.7 %. During the same time period, the value weight market portfolio achieves a SR of 0.21 and a maximum drawdown of 22.6 %.

Equal Weight OOS Portfolio Return Density 196301:201912

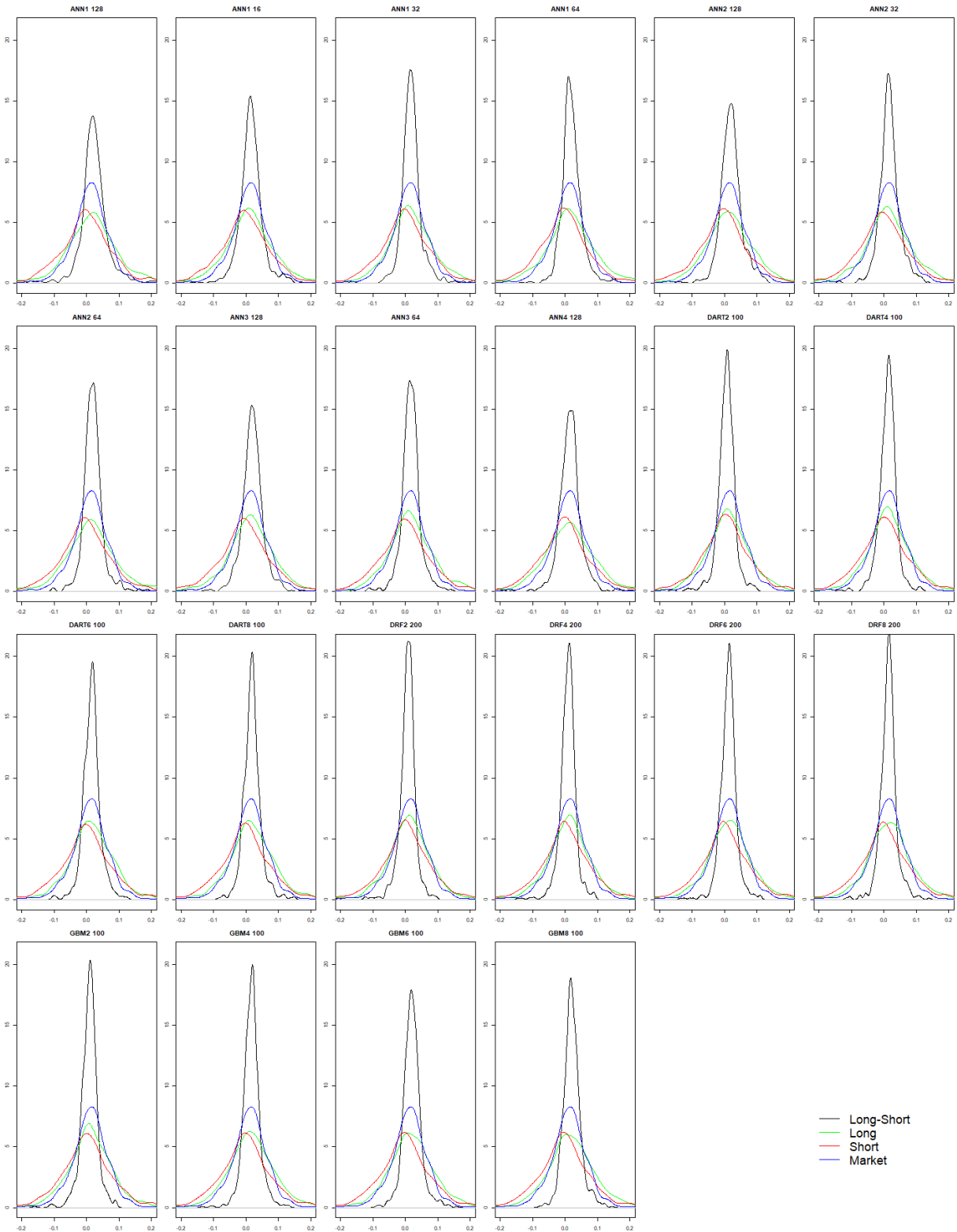


Figure 4 Equal Weight OOS Portfolio Return Distributions 196301:201912

#### Figure 4 (Continues)

Figure 4 presents the overlaid comparison of the OOS portfolio return distributions across the different allocation strategies based on different classification models. The lines in blue, green, red and black represent respectively the OOS return distributions of the market portfolio, the equal weight long position in the predicted lowest return state (labeled as “short” in the figure), the equal weight long position in the predicted highest return state (labeled as “long” in the figure), and the equal weight long-short position. Overall, the gaps between the green lines and the red lines indicate that our models are able to distinguish the worst OOS performing stocks and the best OOS performing stocks from each other. Our equal weight long-short portfolios provide overall return distributions with a concentration shifted towards the right tails and the variances are significantly reduced. The return distributions included in Figure 4 are all from the equal weight portfolios with the similar construction procedure mentioned in Figure 2.

The best OOS long-short portfolio monthly average return in our test is delivered by our neural network model ANN1 128. The long-short portfolio delivers an average of 2.6 % monthly return from 196301:201912. The long-short portfolio based on our tree model, GBM8 100, delivers an average monthly return of 2.2 %. During the same period, the market portfolio delivers an average monthly return of 1.1 %.

Value Weight OOS Portfolio Return Density 196301:201912

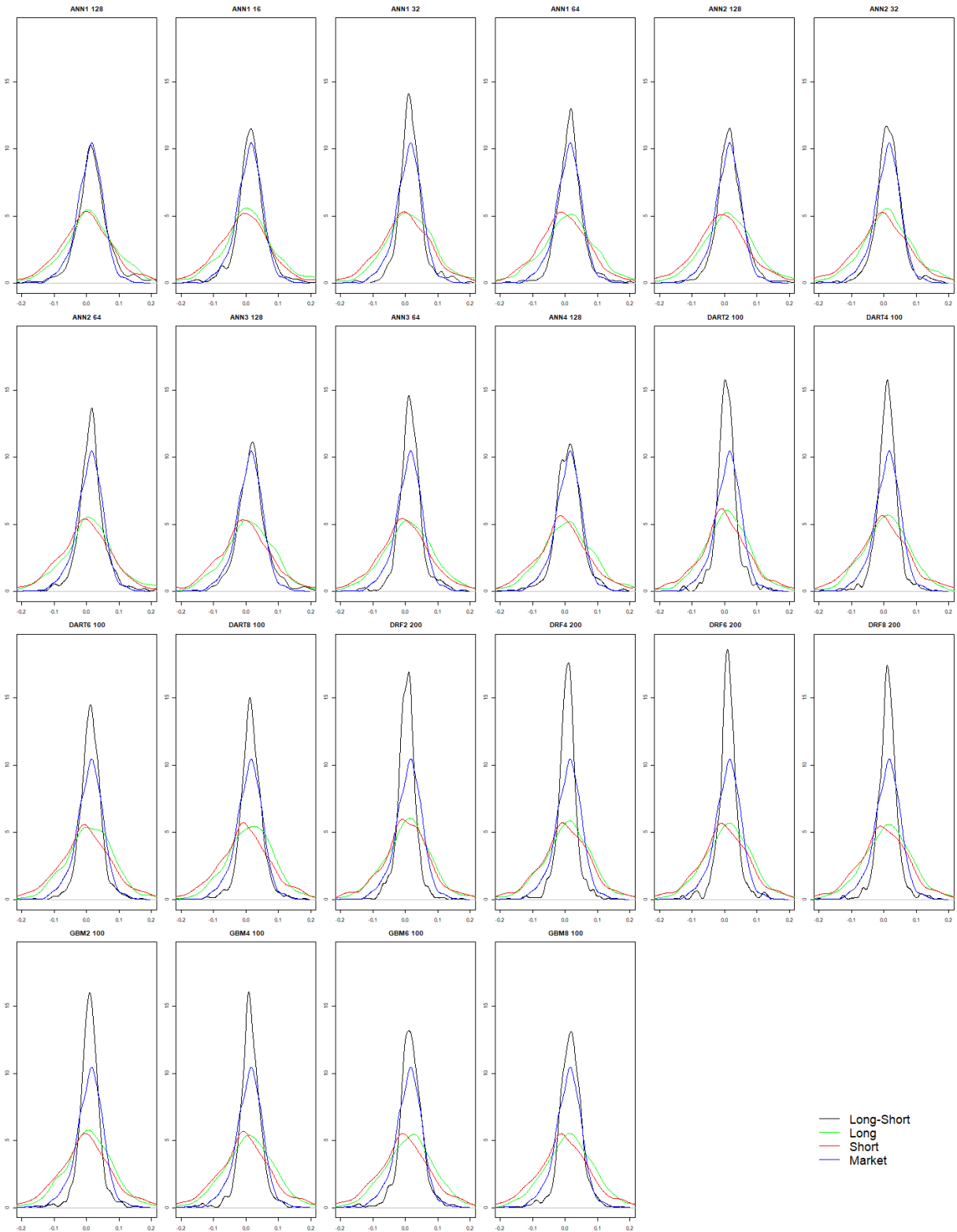


Figure 5 Value Weight OOS Portfolio Return Distributions 196301:201912



### Figure 5 (Continues)

Figure 5 presents the return distribution comparison similar to Figure 4 but of value weight portfolios. The lines in blue, green, red and black represent respectively the OOS return distributions of the market portfolio, the value weight long position in the predicted lowest return state (labeled as “short” in the figure), the value weight long position in the predicted highest return state (labeled as “long” in the figure), and the value weight long-short position. Similar to what is presented in Figure 4, the gaps between the green lines and the red lines indicate that our models are able to distinguish the worst OOS performing stocks and the best OOS performing stocks from each other. Our value weight long-short portfolios also provide overall return distributions with a concentration shifted towards the right tails and the variances are significantly reduced. The return distributions included in Figure 5 are all from the value weight portfolios with the similar construction procedure mentioned in Figure 3.

For the value weight scheme, the long-short portfolio based on our neural network model ANN3 128 delivers an OOS average monthly return of 2 % from 196301:201912, while the market portfolio delivers an average monthly return of 0.9 %.

### Equal Weight OOS Portfolio Net Cumulative Return 196301:201912

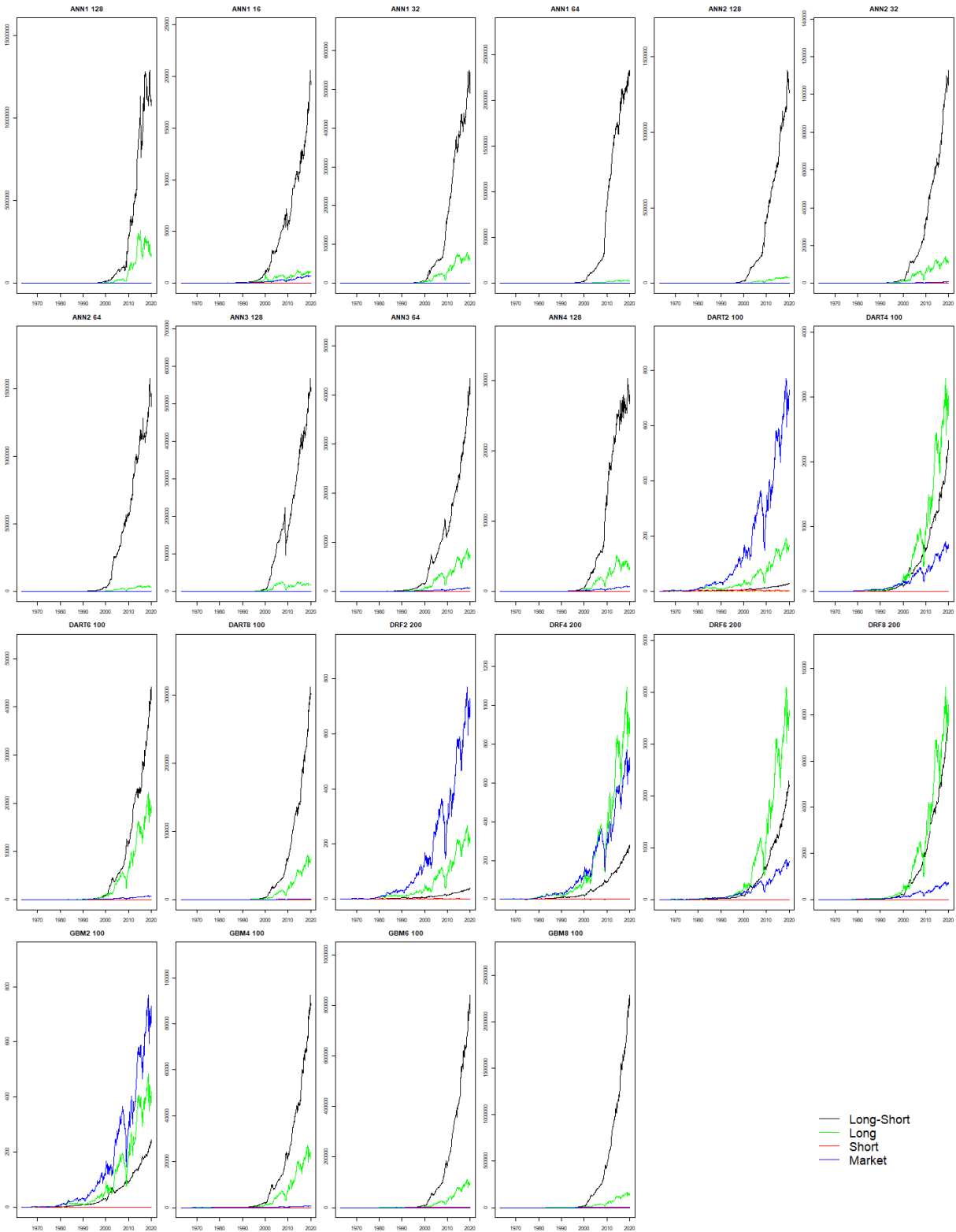


Figure 6 Equal Weight OOS Portfolio Cumulative Returns 196301:201912

### Figure 6 (Continues)

Figure 6 shows the cumulative returns of our equal weight portfolios based on different classification models. The lines in blue, green, red and black represent respectively the OOS cumulative returns of the market portfolio, the equal weight long position in the predicted lowest return state (labeled as “short” in the figure), the equal weight long position in the predicted highest return state (labeled as “long” in the figure), and the equal weight long-short position.

Our equal weight long-short portfolios deliver phenomenal cumulative returns in the OOS test. Comparing to our long-short portfolios, the market portfolios is a horizontal line in the subgraphs as the cumulative return delivered by the market over the same period is too small. Our Neural Network models, ANN1 128 and ANN1 64, deliver OOS cumulative returns of 1,074,286,800 % and 226,665,100% respectively from 196301:201912. Our tree models, GBM8 100 and GBM6 100 deliver OOS cumulative returns of 217,691,000 % and 76,811,200 % respectively during the same investment period. During the same period, the market portfolio achieves a cumulative return of 72,874%, substantially lower than what is delivered by our tree model DRF4 200 by more than 20,000 %, while it is clear that DRF4 200 is underfitted.

Value Weight OOS Portfolio Net Cumulative Return 196301:201912

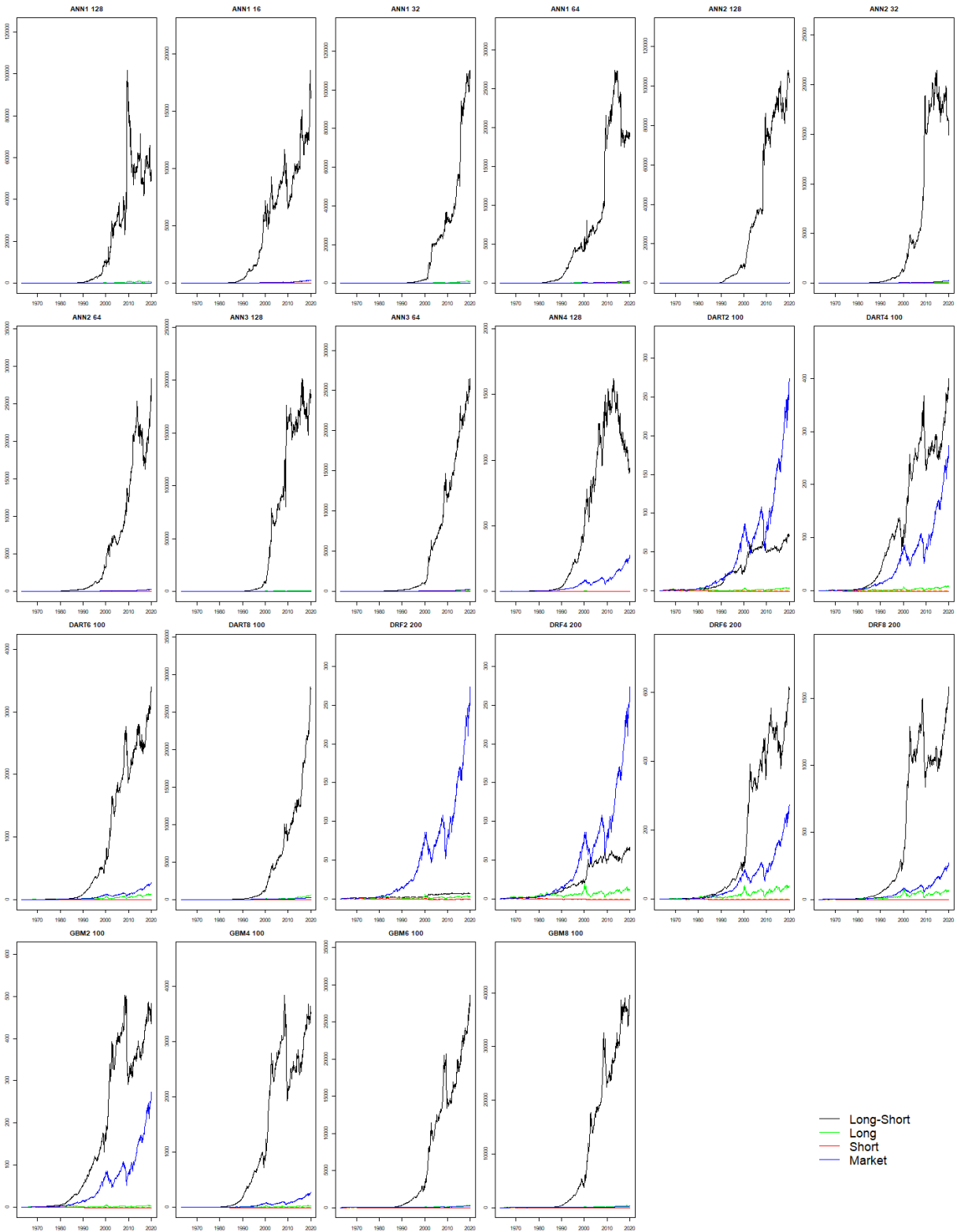


Figure 7 Value Weight OOS Portfolio Cumulative Returns 196301:201912

### Figure 7 (Continues)

Figure 7 is the counterpart of Figure 6 for value weight portfolios. The lines in blue, green, red and black represent respectively the OOS cumulative returns of the market portfolio, the value weight long position in the predicted lowest return state (labeled as “short” in the figure), the value weight long position in the predicted highest return state (labeled as “long” in the figure), and the value weight long-short position.

Our portfolios in value weight also deliver shocking OOS cumulative returns comparing to what the market is capable to achieve. Our long-short portfolios based on neural network models, ANN3 128, ANN1 32, ANN2 128 and ANN1 128, achieve OOS cumulative returns of 18,368,048 %, 10,980,339 %, 10,182,107 % and 5,210,560 % respectively from 196301:201912. The long-short portfolios based on our tree models, GBM8 100, DART8 100 and GBM6 100, deliver cumulative returns of 3,940,967 %, 2,809,332 % and 2,716,935% respectively. During the same period, the value weight market portfolio deliver a cumulative return of 27,396 %.

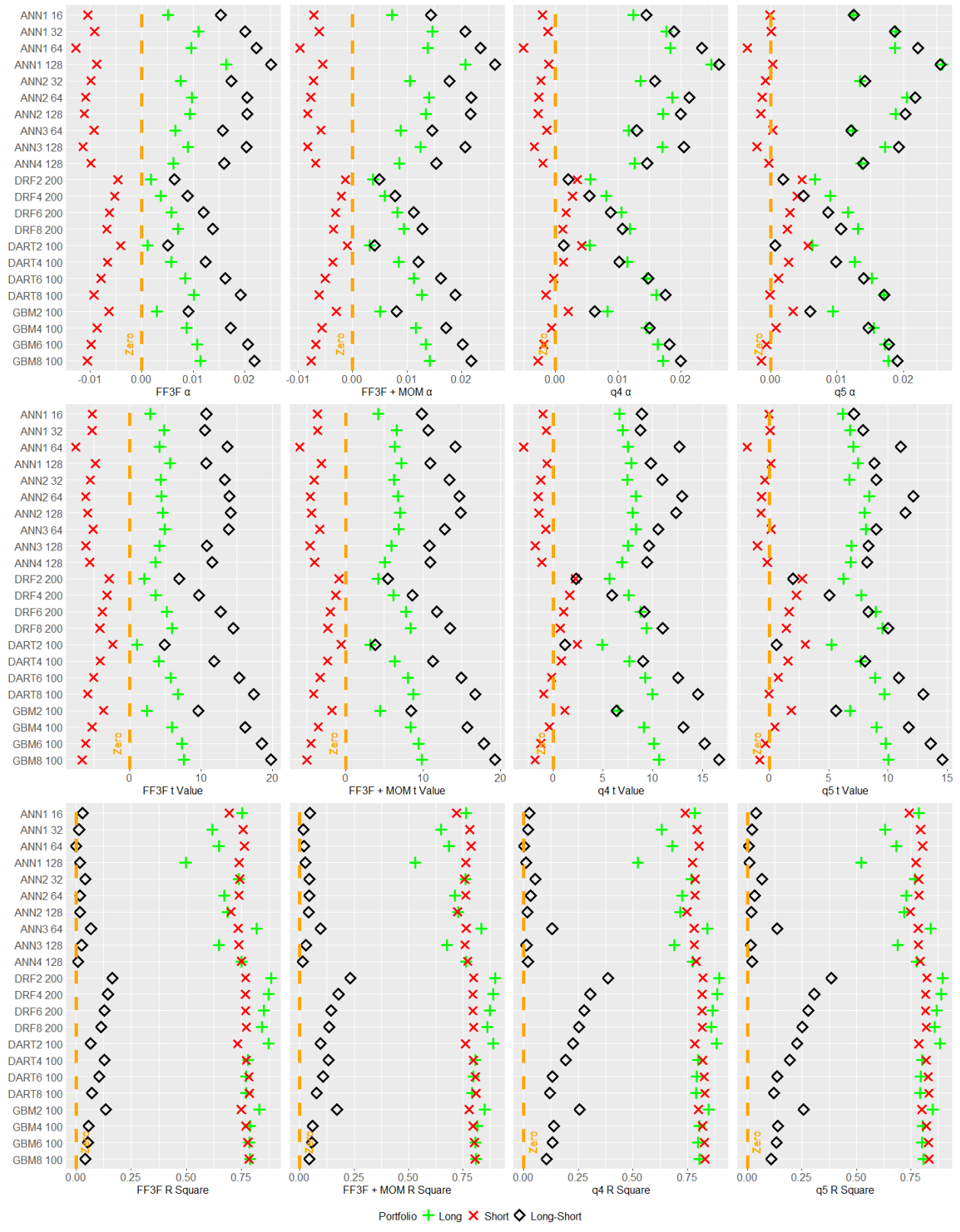


Figure 8 Factor Model Tests on Equal Weight OOS Portfolio Returns 1963:1:2019:12

### Figure 8 (Continues)

Figure 8 demonstrates the factor model tests for the equal weight portfolios. The legends are the same as in Figure 2 and 3 except that the orange dashed line stands for 0. We obtain factor data sets from Kenneth French's data library and Lu Zhang's website. FF3F stands for the Fama French 5 Factor model. MOM stands for the momentum factor. q4 model is the investment CAPM from Hou, Xue, Zhang (2015) and q5 model is the update of the q4 model including the expected growth factor (R\_EG). Note that we do not include Fama French 5 Factor model as we include q4 and we do not include MOM to any of the q-factor models because the q-factor models explains MOM. We present the model alphas, the t statistics of the model alphas and the R square of the models in the 3 lines of subgraphs. It is obvious that the factor models cannot explain the returns achieved by our long-short portfolios. For example, our long-short portfolio based on our neural network model, ANN1 128, has alphas of 2.51 %, 2.62 %, 2.61 % 2.56 % respectively against FF3F, FF3F + MOM, q4 and q5. All R squares of the associated models are below 3 %. The explanatory power of the factor models on the portfolios seems decreasing as the complexity of our classification models increases.

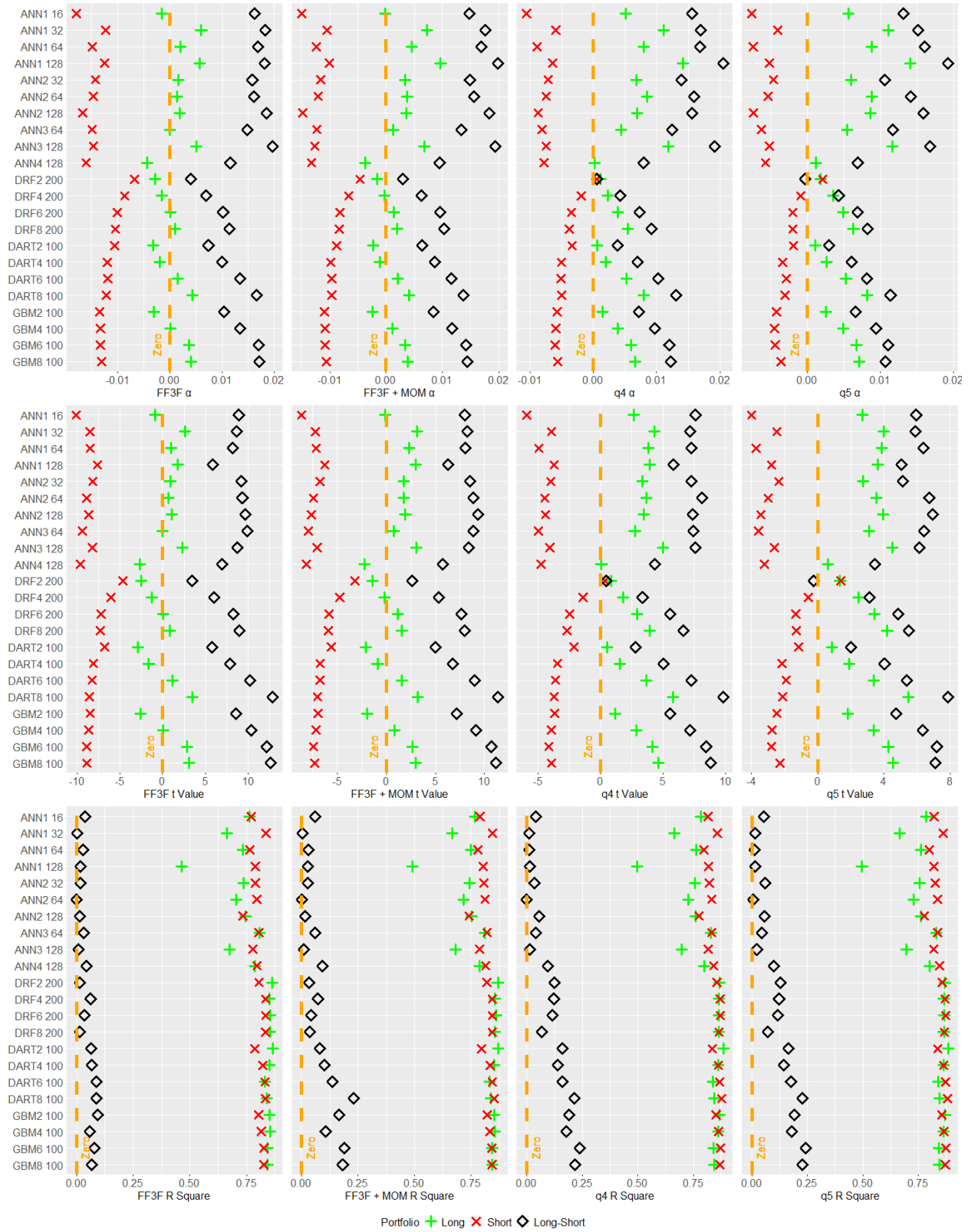


Figure 9 Factor Model Tests on Value Weight OOS Portfolio Returns 1963:01:2019:12



**Figure 9 (Continues)**

Similar to Figure 8, Figure 9 demonstrates the factor model tests for the value weight portfolios. We also present the model alphas, the t statistics of the model alphas and the R square of the models in the 3 lines of subgraphs. Again, it is obvious that the factor models cannot explain the returns achieved by our long-short portfolios. Taking our long-short portfolio based on our neural network model, ANN1 128, the example again, the portfolio has alphas of 1.82 %, 1.99 % 2.05 %, 1.93 % respectively against FF3F, FF3F + MOM, q4 and q5. All R squares of the associated models are below 2 %. The explanatory power of the factor models on the portfolios also seems decreasing as the complexity of our classification models increases.

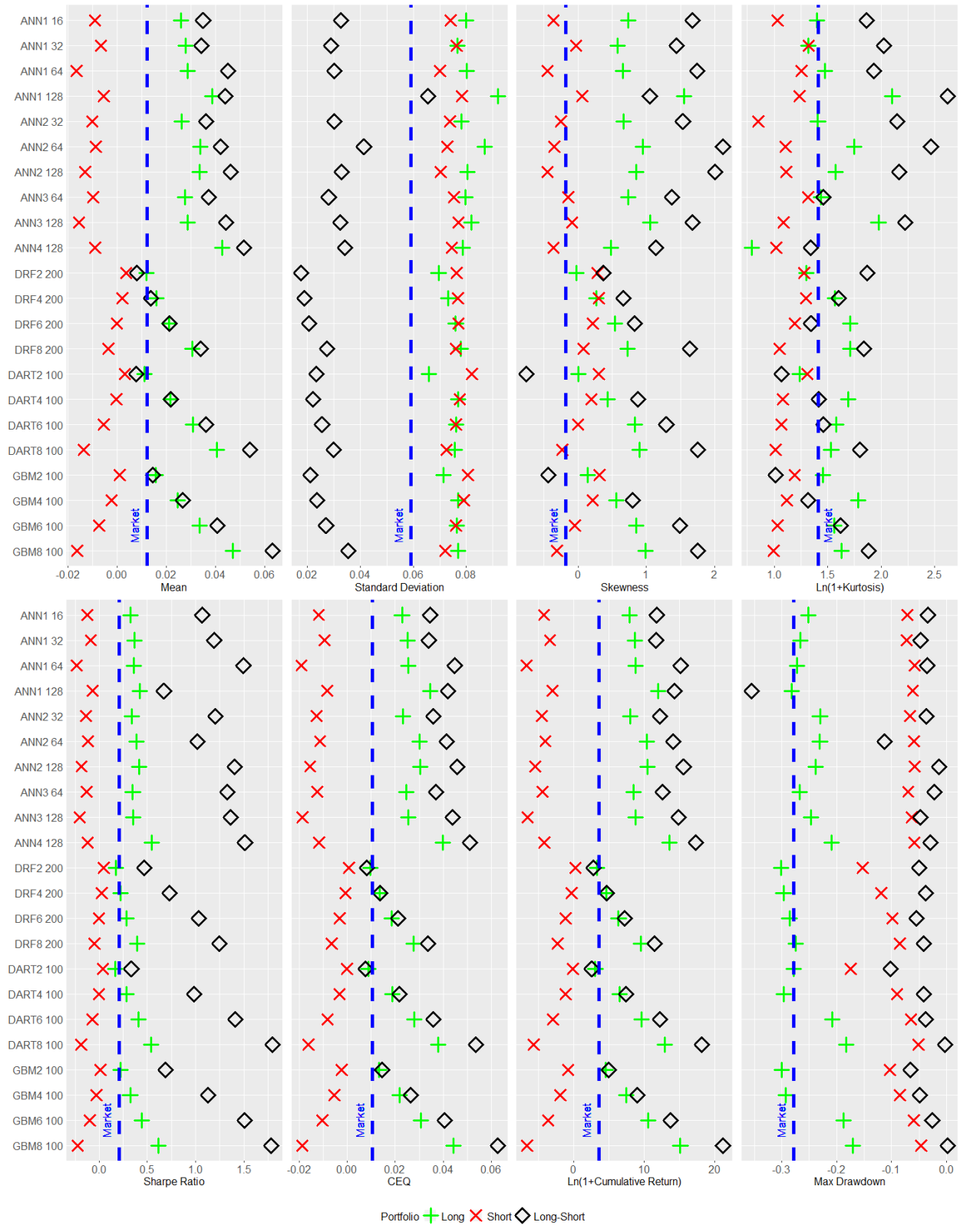


Figure 10 Time Series IS Equal Weight Portfolio Performance 196301:199112

**Figure 10 (Continues)**

Figure 10 presents the economic metrics of the in-sample (IS) performance of the equal weight portfolios based on the fitted classification models from 196301:199112. The setup of the figure is the same as in Figure 2. Comparing the OOS performance figures (Figure 2-5) to Figure 10, we show the consistency between the IS performance and the OOS performance. The sample period from 196301:199112 is used to train the model for OOS evaluation covering 199201:201912. The IS performance of the portfolios is similar to the OOS performance but at a magnified level. For example, during the IS period of 196301:199112, the long-short portfolio based on GBM8 100 achieves a SR of 1.7869 while the market delivers a SR of 0.2078. The best IS model in term of risk-return tradeoff in the period from 196301:199112 is the tree model DART8 100, which delivers a SR of 1.8. Combining Figure 11 and earlier figures on OOS performance, the importance of OOS evaluation is clear.

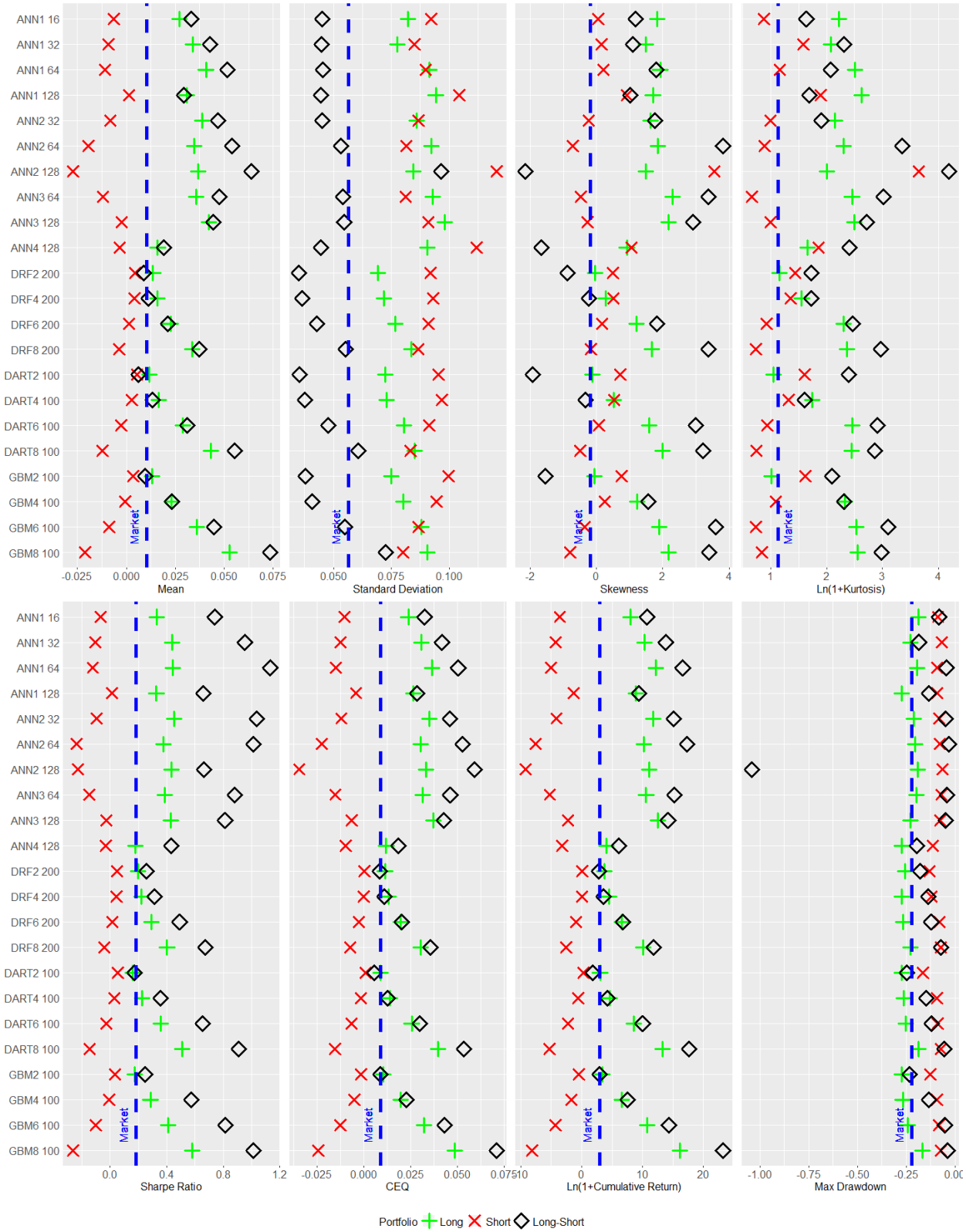


Figure 11 Time Series IS Equal Weight Portfolio Performance 1992:2019

**Figure 11 (Continues)**

Figure 11 demonstrates the IS performance of our portfolios in the complementary time period from 199201:201912. The setup of the figure is the same as in Figure 2. The period 199201:201912 is used to train model for OOS evaluation covering 196301:199112. Similar to what is demonstrated in Figure 10, the IS performance during the period from 199201:201912 is at a magnified level.

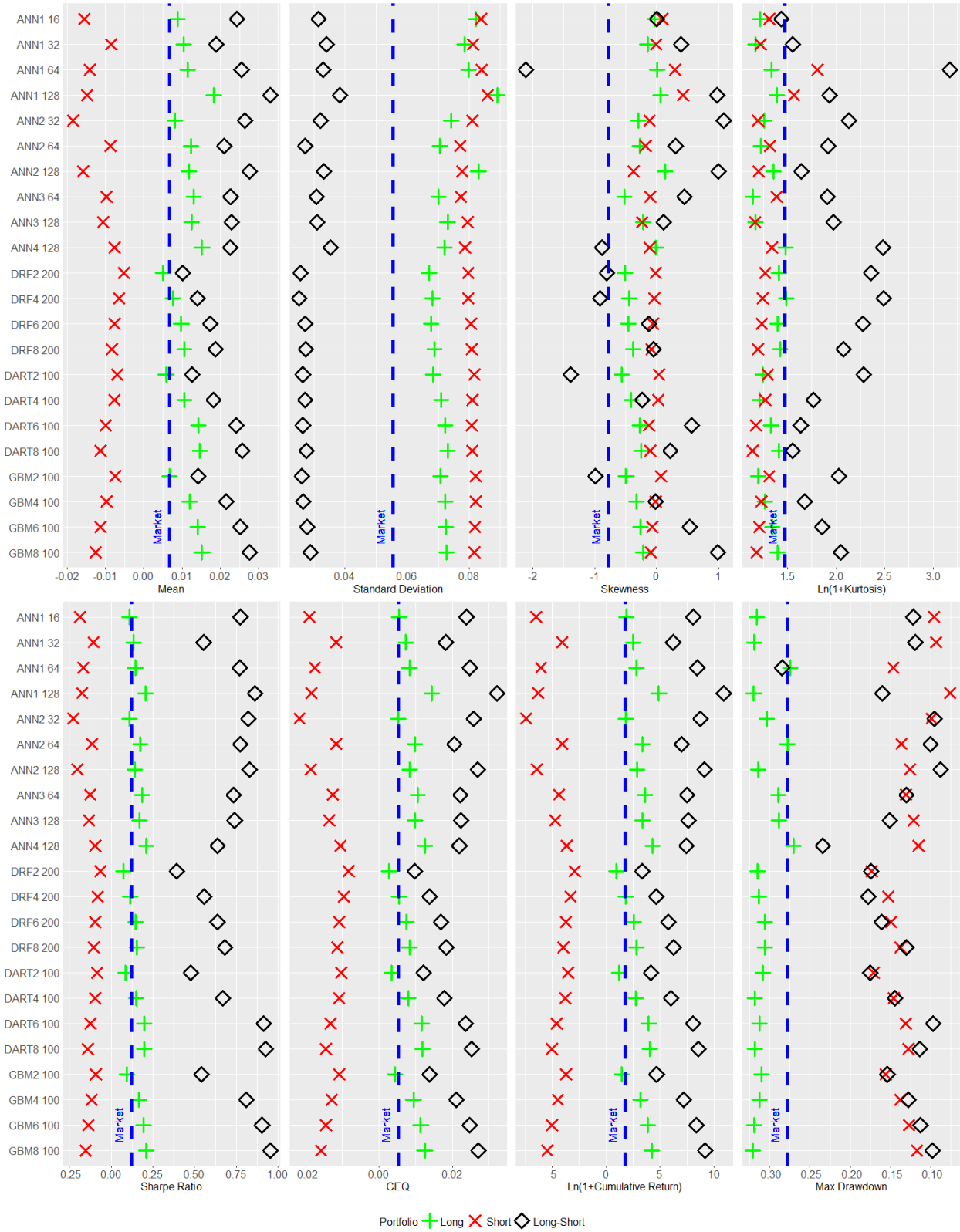


Figure 12 Cross Sectional OOS Equal Weight Portfolio Performance 196302:201912

**Figure 12 (Continues)**

Figure 12 presents the economic metrics of the equal weight portfolios performance evaluated based on our cross-sectional OOS setup. Specifically, we split even number months and odd number months during the sample period of 196301:201912. We train our models based on odd number months and test our models based on even number months. This cross-sectional setup is in the spirit of Fama and French (2018). The cross-sectional OOS evaluation further confirms the superior performance our long-short portfolios based on classification models. The detailed OOS performance of the cross-sectional equal weight portfolios based on the different models mirrors the OOS performance of our equal weight portfolios in the time series test setup as shown in Figure 2.

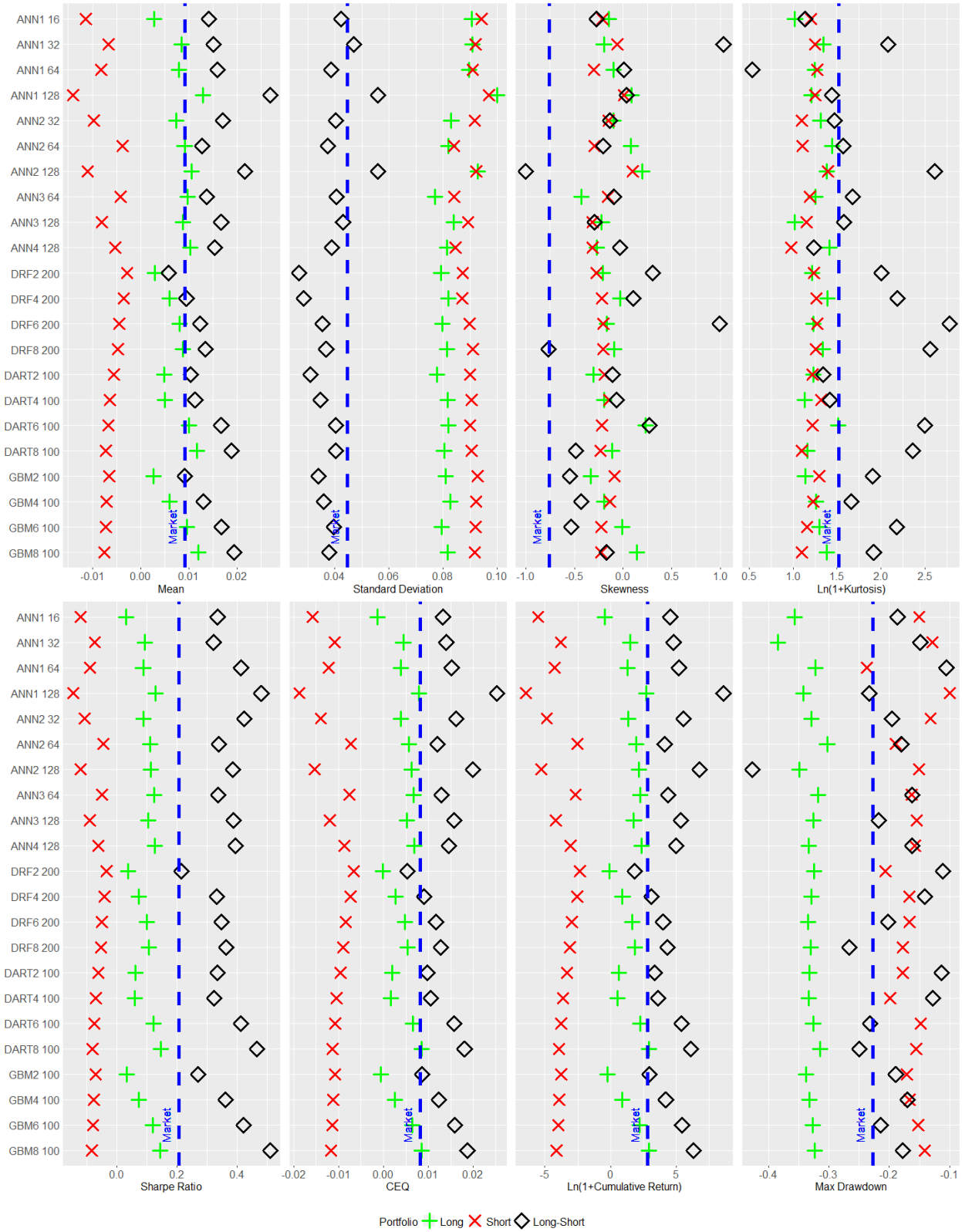


Figure 13 Cross-sectional OOS Value Weight Portfolio Performance 196302:201912



**Figure 13 (Continues)**

Figure 13 presents the cross-sectional performance of the portfolios based on our classification models in value weight scheme. Similar to Figure 12, the cross-sectional OOS evaluation confirms the superior performance of our long-short portfolios based on classification models. The detailed OOS performance of the cross-sectional value weight portfolios based on the different models also seems mirroring the OOS performance of our portfolios based on the time series test setup as in Figure 3.

### Table 1 Specification of Return State Classes

Table 1 describes how we classify return into 10 return states. We cross-sectionally rank individual stock returns by trading month, put them into their corresponding deciles and use the deciles as the classes of return states. For example, if a stock falls into the lowest decile in a trading month, we define the true label of the stock as the class of return state 1. A stock in return state 1 means that the stock delivers a return that is among the worst performing returns of the trading month. A stock in return state 10 indicates that the stock are among the stocks delivering the best performing returns of the trading month.

Specification of Modeling Target	
10 Return States	Criteria
1	Numeric return less than 10 percentile in a month
2	Numeric return less than 20 percentile but greater than or equal to 10 percentile in a month
3	Numeric return less than 30 percentile but greater than or equal to 20 percentile in a month
4	Numeric return less than 40 percentile but greater than or equal to 30 percentile in a month
5	Numeric return less than 50 percentile but greater than or equal to 40 percentile in a month
6	Numeric return less than 60 percentile but greater than or equal to 50 percentile in a month
7	Numeric return less than 70 percentile but greater than or equal to 60 percentile in a month
8	Numeric return less than 80 percentile but greater than or equal to 70 percentile in a month
9	Numeric return less than 90 percentile but greater than or equal to 80 percentile in a month
10	Numeric return greater than or equal to 90 percentile in a month

**Table 2 Model Specification**

Table 2 presents the description of models we apply in this study. Overall, we include 2 modeling architectures and 22 models. Panel A demonstrates the architectural specification for each model. Panel B demonstrates the additional specifications for our neural network models. Panel C demonstrates the hyperparameters that we choose based on cross-validation. Section 2.1.3 discusses the details for each model. Note that we use enum encoding for the categorical variables in our tree models. One-hot-encoding splits the categorical variables into smaller binary choice questions and marks the associated values as positive indicated by 1 and negative indicated by 0. The enum encoding considers the categorical values as nonordinal values. In the tree models, different categories will have different leaves. We do not use one-hot-encoding as it is not the native choice of the tree models and will introduce sparsity which lowers the information ratio for the tree models.

Panel A: Architectural Specifications				
Model	Architecture	Specification	Structural Complexity	Structural Capacity
ANN1 128	Neuron Network	Multilayer Perceptron	1 Hidden Layer	# Neurons = 128
ANN1 16	Neuron Network	Multilayer Perceptron	1 Hidden Layer	# Neurons = 16
ANN1 32	Neuron Network	Multilayer Perceptron	1 Hidden Layer	# Neurons = 32
ANN1 64	Neuron Network	Multilayer Perceptron	1 Hidden Layer	# Neurons = 64
ANN2 128	Neuron Network	Multilayer Perceptron	2 Hidden Layers	# Neurons = {128,64}
ANN2 32	Neuron Network	Multilayer Perceptron	2 Hidden Layers	# Neurons = {32,16}
ANN2 64	Neuron Network	Multilayer Perceptron	2 Hidden Layers	# Neurons = {64,32}
ANN3 128	Neuron Network	Multilayer Perceptron	3 Hidden Layers	# Neurons = {128,64,32}
ANN3 64	Neuron Network	Multilayer Perceptron	3 Hidden Layers	# Neurons = {64,32,16}
ANN4 128	Neuron Network	Multilayer Perceptron	4 Hidden Layers	# Neurons = {128,64,32,16}
DART2 100	Tree	Boosting Tree	Maximum Depth = 2	# Trees = 100
DART4 100	Tree	Boosting Tree	Maximum Depth = 4	# Trees = 100
DART6 100	Tree	Boosting Tree	Maximum Depth = 6	# Trees = 100
DART8 100	Tree	Boosting Tree	Maximum Depth = 8	# Trees = 100
DRF2 200	Tree	Forest	Maximum Depth = 2	# Trees = 200
DRF4 200	Tree	Forest	Maximum Depth = 4	# Trees = 200
DRF6 200	Tree	Forest	Maximum Depth = 6	# Trees = 200
DRF8 200	Tree	Forest	Maximum Depth = 8	# Trees = 200
GBM2 100	Tree	Boosting Tree	Maximum Depth = 2	# Trees = 100
GBM4 100	Tree	Boosting Tree	Maximum Depth = 4	# Trees = 100
GBM6 100	Tree	Boosting Tree	Maximum Depth = 6	# Trees = 100
GBM8 100	Tree	Boosting Tree	Maximum Depth = 8	# Trees = 100

Panel B: ANN Other Specifications				
Hidden Layer Activation	Output Layer Activation	Categorical Variable Encoding	# epochs	Loss
Tanh	Softmax	One hot encoding	50	Cross Entropy

Panel C: Hyperparameters			
Architecture	Model	Parameter	Candidate
Neuron Network	All	L1 Regulation	0.01, 0.001, 0.0001, 0.00001
Tree	All	Sample Rate	0.8, 1
Tree	All	Predictor Sample Rate	0.8, 1

**Table 3 Selection of No Information Benchmark Classifier**

Table 3 presents the Tukey's HSD multiple comparison test with Monte Carlo simulation. The testing samples are generated with the sample covering 199201:201912. Classifier 1 is the random classifier that assigns return states with equal probability. Classifier 2 is the random classifier that assigns return states with IS sample probability mass function observed in the sample 196301:199112. Classifier 3 is the naïve classifier that assigns the most populated IS return state to all OOS observations. Classifier 4 is the random classifier that assigns return states with OOS sample probability mass function with equal probability. Classifier 5 is the naïve classifier that assigns the most populated OOS returns state to all OOS observations. Note that even with minimum information, the naïve classifier which assigns the most populated IS return state to all OOS observations demonstrates higher accuracy than random classifier that uses no information.

	Difference	Lower 95% CI	Upper 95% CI	P Value
2-1	0.00005	-0.00011	0.00022	0.91338
3-1	0.00038	0.00021	0.00054	0.00000
4-1	0.00007	-0.00010	0.00024	0.78601
5-1	0.00060	0.00043	0.00076	0.00000
3-2	0.00033	0.00016	0.00049	0.00000
4-3	0.00002	-0.00015	0.00018	0.99857
5-2	0.00055	0.00038	0.00071	0.00000
4-3	-0.00031	-0.00047	-0.00014	0.00000
5-3	0.00022	0.00005	0.00039	0.00286
5-4	0.00053	0.00036	0.00069	0.00000

**Table 4 Data Construction Summary Statistics**

Table 4 presents the summary statistics of our data with CRSP database as the reference. Panel A presents number of securities in our sample. Panels B and C presents summary statistics and market capitalization in month t-1, respectively.

Panel A: Number of Securities Summary					
Sample	Distinct Total	Mean	Min	Max	Filter
CRSP	33004	6146.905	2069	9366	None
Our Sample	26302	4886.6754	1997	7929	No missing return; EXCHCD and SHRCD

Panel B: Summary Statistics of Returns					
Sample	Mean	SD	Skewness	Kurtosis	Filter
CRSP	0.0102	0.176	20.8963	5165.1519	No missing return
Our Sample	0.0109	0.1883	20.6107	4785.8618	No missing return; EXCHCD and SHRCD

Panel C: Summary Statistics of Market Capitalization at t-1					
Sample	Mean	SD	Skewness	Kurtosis	Filter
CRSP	1601233.289	11003218.89	27.6035	1369.1445	No missing t-1 ME
Our Sample	1723086.927	11923239.69	26.0793	1202.901	No missing t-1 ME; EXCHCD and SHRCD

### Table 5 Sample Splits

Table 5 describes the period we apply for training and testing setup for in-sample (IS) and out-of-sample (OOS). Specifically, we form an overall time series OOS test sample covering 196301:201912 based on splitting the time period into two time series training periods. We use the time period from 196301:199112 to train models for OOS predictions in the period from 199201:201912 and we use the time period from 199201:201912 to train models for OOS predictions in the period from 196301:199112. We combine the OOS predictions for OOS evaluation. In the spirit of Fama and French (2018), we also split the data by even number and odd number months for the cross-sectional OOS evaluation. We train our models cross-sectionally with odd number months from 196301:201911 and test the OOS predictions with the even number months from 196302:201912.

Training and Testing Setup	IS Training	OOS Testing
Time Series Training 1	196301:199112	199201:201912
Time Series Training 2	199201:201912	196301:199112
Time Series Evaluation	We make OOS predictions on 199201:201912 and 196301:199112 with models trained IS covering 196301:199112 and 199201:201912 respectively.	
Cross-sectional Evaluation	Odd Number Months 196301:201911	Even Number Months 196302:201912

**Table 6 Overall Accuracy of Time Series OOS Prediction and Binomial Test 196301:201912**

Table 6 presents the accuracy of each model. The accuracy of a model is the direct evaluation of the correctness of the model predictions. The Kappa statistic measures the level of agreement between the predictions and the actual data and higher Kappa statistic indicates better performance. According to Landis and Koch (1977), a Kappa value greater than 0 but less than 0.2 indicates that the agreement is slight. The confidence interval is associated with the binomial test against the no information accuracy chosen based on Table 3. The P values are associated with the hypothesis test on whether the accuracy is different from the no information accuracy statistically. We discussed our model specifications and the statistical metrics in Section 2. Table 6 shows that all of our models are better than the no information accuracy which is calculated under the assumption that the historical information is useless in terms of prediction future return states.

Model	Accuracy	Kappa	95% Lower Bound	95% Upper Bound	No Information Accuracy	P Value
ANN1 128	0.1476	0.0526	0.1472	0.1480	0.1016	0.0000
ANN1 16	0.1443	0.0488	0.1440	0.1447	0.1016	0.0000
ANN1 32	0.1486	0.0539	0.1482	0.1490	0.1016	0.0000
ANN1 64	0.1458	0.0507	0.1454	0.1461	0.1016	0.0000
ANN2 128	0.1464	0.0512	0.1460	0.1468	0.1016	0.0000
ANN2 32	0.1480	0.0531	0.1476	0.1484	0.1016	0.0000
ANN2 64	0.1476	0.0530	0.1472	0.1480	0.1016	0.0000
ANN3 128	0.1461	0.0513	0.1457	0.1465	0.1016	0.0000
ANN3 64	0.1489	0.0543	0.1485	0.1493	0.1016	0.0000
ANN4 128	0.1509	0.0563	0.1505	0.1513	0.1016	0.0000
DART2 100	0.1516	0.0570	0.1512	0.1520	0.1016	0.0000
DART4 100	0.1552	0.0610	0.1548	0.1556	0.1016	0.0000
DART6 100	0.1563	0.0623	0.1559	0.1567	0.1016	0.0000
DART8 100	0.1559	0.0618	0.1555	0.1563	0.1016	0.0000
DRF2 200	0.1509	0.0562	0.1505	0.1512	0.1016	0.0000
DRF4 200	0.1546	0.0604	0.1542	0.1550	0.1016	0.0000
DRF6 200	0.1561	0.0620	0.1557	0.1565	0.1016	0.0000
DRF8 200	0.1567	0.0627	0.1563	0.1571	0.1016	0.0000
GBM2 100	0.1551	0.0609	0.1548	0.1555	0.1016	0.0000
GBM4 100	0.1572	0.0632	0.1568	0.1576	0.1016	0.0000
GBM6 100	0.1577	0.0638	0.1573	0.1581	0.1016	0.0000
GBM8 100	0.1572	0.0633	0.1568	0.1576	0.1016	0.0000

**Table 7 Return State Transition Probability and Mean Return 196301:201912**

This table presents the return state transition probability and mean return of the transition from the old to new state. A stock in return state 1 means that the stock delivers a return that is among the worst performing returns of the trading month. A stock in return state 10 indicates that the stock are among the stocks delivering the best performing returns of the trading month. Note that the true return state transition probabilities are not evenly distributed. The extreme return states and the middle return states are associated with transition probabilities either substantially greater than 10 % or substantially lower than 10%. These states are thus with higher certainty in the process of state transition. More specifically, return state 3, return state 4 and return state 9 seem the most uncertain states. The return state 1 and return state 10 seem the most certain states. The return states are defined in Table 1.

Panel A: True Return State Transition Probability Matrix 196301:201912										
	New 1	New 2	New 3	New 4	New 5	New 6	New 7	New 8	New 9	New 10
Old 1	0.1741	0.1063	0.0816	0.0686	0.0665	0.0660	0.0719	0.0817	0.1052	0.1782
Old 2	0.1137	0.1073	0.0963	0.0891	0.0879	0.0875	0.0918	0.0993	0.1090	0.1180
Old 3	0.0859	0.0987	0.0997	0.1011	0.1007	0.1033	0.1050	0.1054	0.1059	0.0944
Old 4	0.0713	0.0899	0.1007	0.1073	0.1127	0.1134	0.1122	0.1092	0.1014	0.0817
Old 5	0.0696	0.0860	0.0992	0.1094	0.1128	0.1203	0.1167	0.1098	0.0981	0.0779
Old 6	0.0690	0.0868	0.1002	0.1084	0.1138	0.1177	0.1186	0.1116	0.0970	0.0768
Old 7	0.0675	0.0897	0.1025	0.1083	0.1134	0.1163	0.1164	0.1121	0.0980	0.0758
Old 8	0.0753	0.0973	0.1054	0.1067	0.1102	0.1123	0.1092	0.1058	0.0984	0.0794
Old 9	0.0958	0.1103	0.1061	0.1023	0.0976	0.0974	0.0971	0.1009	0.0999	0.0927
Old 10	0.1742	0.1236	0.0966	0.0825	0.0752	0.0736	0.0743	0.0802	0.0912	0.1284

Panel B: Return State Transition Mean Return 196301:201912										
	New 1	New 2	New 3	New 4	New 5	New 6	New 7	New 8	New 9	New 10
Old 1	-0.2744	-0.1217	-0.0750	-0.0413	-0.0130	0.0114	0.0394	0.0759	0.1350	0.4003
Old 2	-0.2424	-0.1177	-0.0716	-0.0394	-0.0127	0.0127	0.0404	0.0751	0.1292	0.3169
Old 3	-0.2316	-0.1139	-0.0691	-0.0370	-0.0121	0.0126	0.0388	0.0719	0.1243	0.2961
Old 4	-0.2254	-0.1105	-0.0661	-0.0357	-0.0108	0.0119	0.0375	0.0683	0.1193	0.2871
Old 5	-0.2229	-0.1078	-0.0624	-0.0332	-0.0099	0.0120	0.0357	0.0663	0.1174	0.2940
Old 6	-0.2185	-0.1055	-0.0614	-0.0328	-0.0094	0.0127	0.0358	0.0660	0.1165	0.2959
Old 7	-0.2123	-0.1041	-0.0623	-0.0340	-0.0103	0.0114	0.0348	0.0646	0.1128	0.2830
Old 8	-0.2097	-0.1048	-0.0627	-0.0350	-0.0110	0.0113	0.0370	0.0668	0.1168	0.2884
Old 9	-0.2120	-0.1076	-0.0664	-0.0374	-0.0125	0.0113	0.0375	0.0690	0.1207	0.2983
Old 10	-0.2321	-0.1137	-0.0707	-0.0406	-0.0135	0.0115	0.0378	0.0714	0.1276	0.3506



**Table 8 Time Series OOS Prediction Mean Accuracy across Models by Return State Transition 196301:201912**

Table 8 presents our OOS modeling prediction *average* accuracies of return state transitions from the old states to the new states across models. Specifically, we calculate the OOS prediction accuracies of each classification model and form a percentage accuracy table similar to the table below. We then average the numbers across all the models. A stock in return state 1 means that the stock delivers a return that is among the worst performing returns of the trading month. A stock in return state 10 indicates that the stock are among the stocks delivering the best performing returns of the trading month. Details of the return state definition can be found in Table 1.

Combining what is demonstrated in Table 7, Table 8 shows that our models benefit significantly from the most certain return states, i.e., return states 1 and 10. Our models almost give up the most uncertain states, i.e. return states 3, 4, and 9.

	New 1	New 2	New 3	New 4	New 5	New 6	New 7	New 8	New 9	New 10
Old 1	0.5002	0.0334	0.0109	0.0065	0.0246	0.0478	0.0552	0.0521	0.0845	0.4183
Old 2	0.4666	0.0844	0.0246	0.0165	0.0603	0.1290	0.1272	0.1163	0.1318	0.2401
Old 3	0.4102	0.1014	0.0301	0.0225	0.0885	0.2083	0.1754	0.1266	0.1182	0.1697
Old 4	0.3834	0.0962	0.0328	0.0262	0.1085	0.2523	0.1997	0.1334	0.0934	0.1474
Old 5	0.3964	0.0914	0.0311	0.0296	0.1107	0.2675	0.2048	0.1264	0.0804	0.1416
Old 6	0.3961	0.0918	0.0344	0.0296	0.1178	0.2718	0.1997	0.1253	0.0773	0.1413
Old 7	0.3850	0.1001	0.0393	0.0298	0.1179	0.2651	0.1988	0.1306	0.0745	0.1324
Old 8	0.4026	0.1133	0.0450	0.0324	0.1105	0.2419	0.1816	0.1303	0.0794	0.1347
Old 9	0.4634	0.1443	0.0486	0.0335	0.0945	0.1895	0.1322	0.1261	0.0830	0.1398
Old 10	0.7122	0.1223	0.0388	0.0215	0.0532	0.0836	0.0628	0.0716	0.0620	0.1141

**Table 9 OOS Prediction By-Class Accuracy 196301:201912**

This table summarizes the key statistical metrics that measures the performance of our classification models in the OOS predictions covering 196301:201912 by class for the two modeling architectures separately. Our models splits the big question of what is the return state of a stock in the next time period as smaller binary choice questions as describe in Section 2. Specifically, our models ask whether an observation will be in return state 1 or not, whether the observation will be in return state 2 or not, whether the observation will be in return state 3 or not and so on. In short, our models break down the multiclass classification problem into smaller binary classification problems. The evaluation of the accuracy with the statistical metrics by each return state can thus help us understand the modeling performance for each return state. For more details, section 2.2.2 discusses the details of the statistical metrics. Table 1 defines the return states. Table 2 describes model specification.

In the table below, prevalence is the percentage of the associated return state in the population. Sensitivity measures the proportion of correctly predicted positives, while specificity measures the proportion correctly predicted negatives. Sensitivity is also called recall. As an example, for the first line of return state 1, a positive means a prediction of 1 and a negative means a prediction of 0, where the prediction of 1 indicates that the return of the stock in the next period will be in return state 1 and the prediction of 0 indicates that the return of the stock in the next period will not be in return state 1. Precision measures the proportion of correct positives over the sum of correct positives and incorrect negatives. F1 is the balanced F score which is calculated as the harmonic mean of precision and recall. Balanced accuracy is the average of the accuracy between the prediction of positives and the prediction of negatives. In general, greater values in the metrics of sensitivity, specificity, precision, F1 and balanced accuracy indicate better performance.

Panel A presents the by-class metrics in term of average across the neural network models and Panel B presents the by-class metrics in term of average across the tree models. The two panels confirm our findings that our models perform well in return states with high certainty and almost give up the return states with high uncertainty. The balanced accuracy column also shows that our predictions balance the prediction between the positives and negatives.

Panel A: OOS Prediction Average of By-Class Metrics across ANN Models 196301:201912						
Return State	Prevalence	Sensitivity	Specificity	Precision	F1	Balanced Accuracy
1	0.0996	0.4103	0.8382	0.2199	0.2847	0.6243
2	0.0996	0.0849	0.9344	0.1260	0.0940	0.5097
3	0.0988	0.0619	0.9447	0.1109	0.0695	0.5033
4	0.0984	0.0435	0.9636	0.1143	0.0586	0.5035
5	0.0991	0.1252	0.9048	0.1250	0.1196	0.5150
6	0.1008	0.1411	0.8964	0.1316	0.1306	0.5188
7	0.1013	0.1845	0.8637	0.1315	0.1444	0.5241
8	0.1016	0.1080	0.9121	0.1222	0.1075	0.5101
9	0.1004	0.1089	0.9108	0.1203	0.1052	0.5098
10	0.1003	0.2040	0.8838	0.1649	0.1793	0.5439

Panel B: OOS Prediction Average of By-Class Metrics across Tree Models 196301:201912						
Return State	Prevalence	Sensitivity	Specificity	Precision	F1	Balanced Accuracy
1	0.0996	0.5459	0.7697	0.2088	0.3012	0.6578
2	0.0996	0.1101	0.9147	0.1248	0.1159	0.5124
3	0.0988	0.0111	0.9910	0.1198	0.0191	0.5011
4	0.0984	0.0110	0.9908	0.1136	0.0193	0.5009
5	0.0991	0.0674	0.9506	0.1300	0.0880	0.5090
6	0.1008	0.2670	0.8297	0.1496	0.1904	0.5483
7	0.1013	0.1439	0.8958	0.1348	0.1375	0.5199
8	0.1016	0.1238	0.8988	0.1217	0.1218	0.5113
9	0.1004	0.0728	0.9443	0.1273	0.0886	0.5086
10	0.1003	0.1970	0.8758	0.1503	0.1701	0.5364

**Table 10 Time Series IS Training Accuracy and Binomial Test**

Table 10 presents the overall accuracy of the predictions on IS training set across model specifications for the time period from 196301:201912. The table setup is the same as Table 6. Panel A shows overall accuracy of training set across models for 199201:201912. Panel B shows the same information for 199201:201912.

Panel A: Overall Accuracy of Training Set Across Models 196301:199112							
Model	Accuracy	Kappa	95% Lower Bound	95% Upper Bound	No Information Accuracy	P Value	
ANN1 128	0.1641	0.0708	0.1635	0.1647	0.1031	0.0000	
ANN1 16	0.1617	0.0675	0.1611	0.1622	0.1031	0.0000	
ANN1 32	0.1709	0.0788	0.1703	0.1715	0.1031	0.0000	
ANN1 64	0.1702	0.0777	0.1696	0.1708	0.1031	0.0000	
ANN2 128	0.1737	0.0811	0.1731	0.1743	0.1031	0.0000	
ANN2 32	0.1722	0.0794	0.1716	0.1728	0.1031	0.0000	
ANN2 64	0.1742	0.0824	0.1736	0.1748	0.1031	0.0000	
ANN3 128	0.1722	0.0800	0.1716	0.1728	0.1031	0.0000	
ANN3 64	0.1763	0.0848	0.1757	0.1769	0.1031	0.0000	
ANN4 128	0.1735	0.0817	0.1729	0.1741	0.1031	0.0000	
DART2 100	0.1527	0.0578	0.1522	0.1533	0.1031	0.0000	
DART4 100	0.1669	0.0738	0.1663	0.1675	0.1031	0.0000	
DART6 100	0.1841	0.0929	0.1834	0.1847	0.1031	0.0000	
DART8 100	0.2109	0.1228	0.2102	0.2116	0.1031	0.0000	
DRF2 200	0.1540	0.0591	0.1534	0.1546	0.1031	0.0000	
DRF4 200	0.1631	0.0693	0.1625	0.1637	0.1031	0.0000	
DRF6 200	0.1741	0.0817	0.1735	0.1747	0.1031	0.0000	
DRF8 200	0.1920	0.1017	0.1914	0.1927	0.1031	0.0000	
GBM2 100	0.1598	0.0658	0.1592	0.1604	0.1031	0.0000	
GBM4 100	0.1740	0.0816	0.1734	0.1746	0.1031	0.0000	
GBM6 100	0.1945	0.1046	0.1939	0.1951	0.1031	0.0000	
GBM8 100	0.2316	0.1459	0.2309	0.2323	0.1031	0.0000	

**Table 10 (Continues)**

Panel B: Overall Accuracy of Training Set Across Models 199201:201912						
Model	Accuracy	Kappa	95% Lower Bound	95% Upper Bound	No Information Accuracy	P Value
ANN1 128	0.1613	0.0681	0.1608	0.1619	0.1006	0.0000
ANN1 16	0.1662	0.0734	0.1657	0.1668	0.1006	0.0000
ANN1 32	0.1685	0.0760	0.1679	0.1690	0.1006	0.0000
ANN1 64	0.1674	0.0748	0.1668	0.1679	0.1006	0.0000
ANN2 128	0.1652	0.0722	0.1646	0.1657	0.1006	0.0000
ANN2 32	0.1685	0.0760	0.1680	0.1691	0.1006	0.0000
ANN2 64	0.1698	0.0775	0.1693	0.1704	0.1006	0.0000
ANN3 128	0.1660	0.0732	0.1654	0.1665	0.1006	0.0000
ANN3 64	0.1701	0.0778	0.1696	0.1707	0.1006	0.0000
ANN4 128	0.1577	0.0640	0.1572	0.1583	0.1006	0.0000
DART2 100	0.1582	0.0645	0.1577	0.1587	0.1006	0.0000
DART4 100	0.1659	0.0731	0.1654	0.1664	0.1006	0.0000
DART6 100	0.1795	0.0882	0.1790	0.1801	0.1006	0.0000
DART8 100	0.2051	0.1167	0.2045	0.2057	0.1006	0.0000
DRF2 200	0.1589	0.0652	0.1584	0.1594	0.1006	0.0000
DRF4 200	0.1625	0.0692	0.1619	0.1630	0.1006	0.0000
DRF6 200	0.1692	0.0767	0.1687	0.1698	0.1006	0.0000
DRF8 200	0.1836	0.0927	0.1830	0.1842	0.1006	0.0000
GBM2 100	0.1619	0.0686	0.1614	0.1625	0.1006	0.0000
GBM4 100	0.1714	0.0792	0.1708	0.1719	0.1006	0.0000
GBM6 100	0.1895	0.0994	0.1890	0.1901	0.1006	0.0000
GBM8 100	0.2255	0.1393	0.2249	0.2261	0.1006	0.0000

**Table 11 Cross-Sectional OOS Test Accuracy with Even Number Months 196302:201912**

This table presents the overall accuracy of the OOS predictions with even number months across model specifications for our cross-sectional evaluation. The sample split is demonstrated in Table 5 and the table setup is the same as Table 6. We confirm that our models are superior in cross-sectional OOS evaluation than the no information accuracy implied by the efficient market hypothesis.

Model	Accuracy	Kappa	95% Lower Bound	95% Upper Bound	No Information Accuracy	P Value
ANN1 128	0.1476	0.0529	0.1470	0.1481	0.1019	0.0000
ANN1 16	0.1503	0.0559	0.1497	0.1508	0.1019	0.0000
ANN1 32	0.1499	0.0556	0.1494	0.1505	0.1019	0.0000
ANN1 64	0.1543	0.0600	0.1538	0.1548	0.1019	0.0000
ANN2 128	0.1507	0.0566	0.1501	0.1512	0.1019	0.0000
ANN2 32	0.1491	0.0544	0.1485	0.1496	0.1019	0.0000
ANN2 64	0.1539	0.0600	0.1534	0.1545	0.1019	0.0000
ANN3 128	0.1544	0.0602	0.1539	0.1549	0.1019	0.0000
ANN3 64	0.1535	0.0593	0.1530	0.1540	0.1019	0.0000
ANN4 128	0.1516	0.0577	0.1511	0.1522	0.1019	0.0000
DART2 100	0.1545	0.0605	0.1540	0.1551	0.1019	0.0000
DART4 100	0.1579	0.0642	0.1574	0.1585	0.1019	0.0000
DART6 100	0.1601	0.0666	0.1595	0.1606	0.1019	0.0000
DART8 100	0.1608	0.0674	0.1602	0.1614	0.1019	0.0000
DRF2 200	0.1558	0.0616	0.1552	0.1563	0.1019	0.0000
DRF4 200	0.1577	0.0637	0.1571	0.1582	0.1019	0.0000
DRF6 200	0.1594	0.0656	0.1588	0.1599	0.1019	0.0000
DRF8 200	0.1598	0.0661	0.1592	0.1603	0.1019	0.0000
GBM2 100	0.1582	0.0644	0.1577	0.1588	0.1019	0.0000
GBM4 100	0.1600	0.0664	0.1595	0.1606	0.1019	0.0000
GBM6 100	0.1609	0.0674	0.1603	0.1614	0.1019	0.0000
GBM8 100	0.1614	0.0680	0.1608	0.1619	0.1019	0.0000

**Table 12 CS Training Model Average Variable Importance**

Table 12 presents the variable importance of the cross-sectional models based on the training process covering 196301:201911. The variable importance is calculated as the total sum of squared error that the associated variable is able to reduce during the model training process. The sample splits is defined in Table 5. To demonstrate the information in a concise way, we average across the models for the 2 main modeling architectures, ANN and trees, separately as demonstrated respectively in Panel A and B. Specifically, the trading related variables such as lagged idiosyncratic volatility play a part in our models. However, corporate announcement variables such as IPO history in the past 12 months, earning-price ratio, dividend-price ratio and other public information such as number of analysts and SIC classification all play important roles in our models. In addition, historical macro variables make marginal contribution to our models. Note that the suffix number indicates a specific category of the associated indicator variable. For example, ipo.0 stands for the no ipo indicator and sich.-1 stands for the indicator of no SIC information. Because of the difference in encoding the categorical variables (indicators) in our neural network models vs our tree models, our two architectures handle indicator variables differently.

Panel A: Top 50 Average Variable Importance Across Cross Sectional ANN Training Models					
Importance	Variable	Relative Importance	Scaled Importance	Percentage	Rank
1	idiovol	0.9968	0.9968	0.0184	1.1
2	baspread	0.7669	0.7669	0.0137	3.3
3	sich2.60	0.7549	0.7549	0.0136	3.5
4	retvol	0.6986	0.6986	0.0125	3.8
5	ipo.0	0.6538	0.6538	0.0115	4.5
6	ipo.1	0.4564	0.4564	0.0078	8.9
7	label10.0	0.4035	0.4035	0.0071	10.3
8	label10.9	0.3975	0.3975	0.0069	11.6
9	zerotrade	0.3773	0.3773	0.0068	13.5
10	sin.0	0.4056	0.4056	0.007	14.5
11	mom12m	0.3634	0.3634	0.0064	15.2
12	sich2.49	0.372	0.372	0.0067	16.1
13	mom1m	0.3564	0.3564	0.0063	16.5
14	dolvol	0.363	0.363	0.0063	17.6
15	sich2.10	0.3539	0.3539	0.0063	18.1
16	mom6m	0.3252	0.3252	0.0057	20.2
17	beta	0.3298	0.3298	0.0059	21.1
18	sich2.-1	0.3138	0.3138	0.0055	21.3
19	securedind.1	0.3244	0.3244	0.0055	23.6
20	securedind.0	0.3103	0.3103	0.0053	27.9
21	convind.0	0.3086	0.3086	0.0053	28.6
22	dp	0.3107	0.3107	0.0055	28.8
23	rd.1	0.2967	0.2967	0.0051	29.6
24	label10.8	0.3003	0.3003	0.005	31.6
25	label10.1	0.2789	0.2789	0.0048	34.1
26	label10.2	0.2847	0.2847	0.0048	34.1
27	rd.-1	0.2788	0.2788	0.0048	35.3
28	label10.7	0.273	0.273	0.0047	35.7
29	divi.-1	0.2818	0.2818	0.0048	36.1
30	divi.0	0.2839	0.2839	0.0049	36.4
31	age	0.2551	0.2551	0.0045	37.1
32	sich2.56	0.2515	0.2515	0.0045	37.1
33	sich2.28	0.2713	0.2713	0.0046	37.2
34	turn	0.261	0.261	0.0045	37.4
35	divo.0	0.2844	0.2844	0.0049	37.6
36	convind.1	0.2764	0.2764	0.0047	37.7
37	label10.4	0.2692	0.2692	0.0046	38.4
38	label10.6	0.2679	0.2679	0.0046	38.5
39	label10.3	0.2682	0.2682	0.0046	38.9

Panel A: Top 50 Average Variable Importance Across Cross Sectional ANN Training Models ( <b>Continues</b> )						
Importance	Variable	Relative Importance	Scaled Importance	Percentage	Rank	
40	divo.-1	0.2683	0.2683	0.0046	40.2	
41	rd.0	0.2688	0.2688	0.0046	41	
42	invest.y	0.2508	0.2508	0.0043	42.7	
43	sich2.63	0.2359	0.2359	0.0041	42.9	
44	sich2.13	0.2292	0.2292	0.004	51	
45	dy	0.2186	0.2186	0.0039	51.2	
46	sich2.20	0.2172	0.2172	0.0038	51.6	
47	securedind.-1	0.2367	0.2367	0.0041	52.7	
48	ppicmm	0.215	0.215	0.0037	53.5	
49	std_turn	0.2073	0.2073	0.0037	55	
50	label10.5	0.2397	0.2397	0.0041	55.6	

**Table 12 (Continues)**

Panel B: Top 50 Average Variable Importance Across Cross Sectional Tree Training Models					
Importance	Variable	Relative Importance	Scaled Importance	Percentage	Rank
1	idiovol	252020.6634	0.9522	0.3287	1.3333
2	baspread	111399.6761	0.5508	0.1489	1.9167
3	retvol	98918.6611	0.4642	0.1321	2.9167
4	maxret	19223.0792	0.1236	0.0234	6.4167
5	mom6m	14410.1369	0.0873	0.0183	6.5833
6	label10	16302.6913	0.1241	0.0221	6.75
7	mom12m	15108.8604	0.0783	0.0193	7
8	sich2	25173.3684	0.1675	0.039	7.25
9	mom1m	15667.66	0.0752	0.0196	7.5
10	ep	9464.183	0.0616	0.0119	11
11	roaq	6032.3019	0.0377	0.0078	14.9091
12	dy	8893.5185	0.055	0.0105	16.0833
13	betasq	9981.4297	0.0667	0.0118	19.4545
14	beta	6456.432	0.0424	0.0078	21.6364
15	turn	3273.5324	0.0167	0.0037	21.6667
16	zerotrade	2721.8782	0.0105	0.003	23.0833
17	dolvol	5352.8563	0.0335	0.0065	25.9091
18	ill	4037.6283	0.0224	0.0046	27.6364
19	nanalyst	2547.2677	0.0147	0.0031	28.0909
20	label10.0	2044.8495	0.004	0.0017	32.25
21	ill_hi_minus_low	2156.603	0.0074	0.0025	32.4167
22	nonborres	1580.9666	0.0059	0.0019	32.75
23	roeq	5354.326	0.0376	0.0067	35.6364
24	std_turn	1618.6375	0.0077	0.0018	37.5
25	dp	1464.5543	0.006	0.0017	38.1667
26	sich2_ret	1707.6039	0.0057	0.0018	40.25
27	bm.x	1923.647	0.0103	0.0022	41.0909
28	fgr5yr_hi_minus_low	1434.709	0.0053	0.0016	41.8333
29	fgr5yr	1980.004	0.012	0.0025	42.6667
30	ddurrg3m086sbea	1412.1605	0.0048	0.0016	45.6667
31	roavol	8958.6767	0.0642	0.0111	47
32	agr	1188.3276	0.0048	0.0013	49.0909
33	ndmanemp	1114.8946	0.0044	0.0014	49.2727
34	chcsho	1119.6585	0.0046	0.0014	50.4545
35	chmom	1746.2888	0.009	0.002	51.4545
36	turn_hi_minus_low	1065.2049	0.0037	0.0012	55.5
37	roic	4237.5983	0.0303	0.0053	59.4545
38	bm.y	1311.4077	0.0055	0.0015	60.7273
39	cusr0000sas	997.8378	0.0047	0.0011	61.8333
40	bm_hi_minus_low	980.8402	0.0036	0.0011	64.8182
41	chcsho_hi_minus_low	844.7581	0.0031	9.00E-04	65.75
42	age	7391.9657	0.0557	0.0094	68.0909
43	invest.x	978.3864	0.004	0.001	68.2727
44	ms_hi_minus_low	903.585	0.0034	0.001	68.3636
45	pchgm_pchsale_hi_minus_low	841.0044	0.0032	9.00E-04	70.1667
46	realestate_hi_minus_low	672.6382	0.0029	8.00E-04	70.5833
47	nanalyst_hi_minus_low	901.2761	0.0037	0.0011	70.8182
48	dolvol_hi_minus_low	950.7877	0.0034	0.001	76.4545
49	sfe_hi_minus_low	969.5795	0.0034	0.0011	76.6364
50	exszusx	1110.7301	0.0042	0.0013	76.9091



# Appendix

**Table A.1 Firm Characteristics**

This table presents characteristics we reconstructed based on Green, Hand and Zhang (2017).

Name	Description	Min	Max	Mean	Median	Skewness	Kurtosis
acc	Working capital accruals	-1.02	0.58	-0.02	-0.02	-0.88	4.81
aeavol	Abnormal earnings announcement volume	-1.00	21.69	0.87	0.30	3.64	18.51
age	# years since first Compustat coverage	1.00	56.00	12.72	9.00	1.36	1.55
agr	Asset growth	-0.68	5.85	0.15	0.08	4.26	30.02
baspread	Bid-ask spread	0.00	0.91	0.05	0.03	5.15	38.11
beta	Beta	-0.74	3.94	1.08	1.01	0.69	0.69
bm	Book-to-market	-2.35	7.81	0.77	0.60	2.48	11.46
cash	Cash holdings	0.00	0.98	0.16	0.07	1.89	3.15
cashdebt	Cash flow to debt	-7.71	2.23	0.07	0.13	-4.14	25.99
cashpr	Cash productivity	-520.62	600.28	-1.90	-0.73	0.89	29.15
cfp	Cash flow to price ratio	-513.56	156.76	0.05	0.05	-172.20	57390.53
cfp_ia	Industry-adjusted cash flow to price ratio	-449.37	7031.61	13.09	0.00	21.92	479.82
chadv	Change in dividend	-1.59	2.02	0.05	0.03	0.50	8.14
chatoia	Industry-adjusted change in asset turnover	-1.43	1.19	0.00	0.00	-0.15	4.74
chcsho	Change in shares outstanding	-0.89	2.57	0.11	0.01	3.28	13.85
chfeps	Change in forecasted EPS	-6.48	8.25	0.00	0.00	1.29	121.37
chinv	Change in inventory	-0.29	0.37	0.01	0.00	1.10	6.77
chnanalyst	Change in number of analysts	-12.00	9.00	-0.01	0.00	-0.60	9.46
chtx	Change in tax expense	-0.12	0.16	0.00	0.00	0.35	13.06
cinvest	Corporate investment	-26.83	27.87	-0.02	0.00	-2.17	244.24
currat	Current ratio	0.16	60.34	3.16	2.00	5.58	40.35
depr	Depreciation/PP&E	0.01	5.51	0.26	0.15	5.92	49.70
disp	Dispersion in forecasted EPS	0.00	10.00	0.15	0.04	6.48	58.18
dy	Dividend to price	0.00	0.35	0.02	0.00	2.67	10.86
ear	Earnings announcement return	-0.46	0.51	0.00	0.00	0.26	3.17
egr	Growth in common shareholder equity	-3.54	8.19	0.14	0.08	3.32	28.51
ep	Earnings to price	-7.66	0.68	-0.01	0.05	-8.11	107.23
fgr5yr	Forecasted growth in 5-year EPS	-43.50	99.41	16.35	14.50	1.50	5.47
gma	Gross profitability	-0.84	1.78	0.37	0.33	0.81	1.52
grcapx	Growth in capital expenditures	-13.89	55.54	0.89	0.14	5.60	45.95
grltnoa	Growth in long term net operating assets	-0.61	1.18	0.09	0.06	1.64	7.48
herf	Industry sales concentration	0.01	1.00	0.08	0.05	3.10	11.91
hire	Employee growth rate	-0.74	4.00	0.09	0.02	3.81	24.97
idiovol	Idiosyncratic return volatility	0.01	0.26	0.06	0.06	1.47	2.70
ill	Illiquidity	0.00	0.00	0.00	0.00	14.63	355.90

**Table A.1 (Continues)**

Name	Description	Min	Max	Mean	Median	Skewness	Kurtosis
indmom	Industry momentum	-1.00	3.56	0.14	0.12	1.26	5.39
invest	Capital expenditures and inventory	-0.52	2.21	0.08	0.04	2.51	12.80
lev	Leverage	0.00	77.75	2.28	0.69	5.47	43.92
Meanrec	Mean number of analysts	1.00	4.50	2.22	2.20	-0.01	-0.32
mom12m	12-month momentum	-1.00	11.60	0.13	0.06	2.89	21.73
mom1m	1-month momentum	-0.70	2.11	0.01	0.00	1.16	7.77
mom36m	36-month momentum	-0.98	16.20	0.33	0.16	3.08	20.08
ms	Financial statement score	0.00	8.00	3.73	4.00	-0.03	-0.72
mve	Size	6.02	18.90	11.77	11.63	0.29	-0.31
mve_ia	Industry-adjusted size	-	-	-	-	9.17	120.31
nanalyst	Number of analysts covering stock	16608.51	133635.00	158.25	359.12	1.75	2.80
nincr	Number of earnings increases	0.00	34.00	5.17	3.00	2.15	6.35
orgcap	Organizational capital	0.00	8.00	1.00	1.00	2.73	11.60
pchcapx_ia	Industry adjusted % change in capital expenditures	0.00	0.18	0.01	0.01	2.73	11.60
pchcurrat	% change in current ration	-237.42	1640.09	6.50	-0.35	15.17	273.34
pchdepr	% change in depreciation	-0.89	6.72	0.06	-0.01	3.82	23.69
pchgm_pchsale	% change in gross margin - % change in sales	-0.85	7.37	0.10	0.03	4.56	36.69
pchsale_pchinvt	% change in sales - % change in inventory	-12.26	4.77	-0.06	0.00	-5.49	54.90
pchsale_pchrect	% changes in sales - % change in A/R	-11.61	3.02	-0.06	0.01	-5.85	57.26
pchsale_pchxsga	% change in sales - % change in SG&A	-7.93	3.11	-0.04	0.00	-2.90	24.02
pctacc	Percent accruals	-3.50	4.34	0.02	0.00	3.42	31.72
pricedelay	Price delay	-64.75	71.43	-0.65	-0.27	-1.90	34.80
ps	Financial statement score	-15.85	15.52	0.14	0.06	0.09	39.90
rd_mve	R&D to market capitalization	0.00	8.00	4.18	4.00	0.03	-0.55
rd_sale	R&D to sales	0.00	2.23	0.06	0.03	5.12	48.03
retvol	Return volatility	0.00	283.48	0.61	0.03	23.58	733.49
roaq	Return on assets	0.00	0.27	0.03	0.02	2.42	8.82
roavol	Earnings volatility	-0.48	0.16	0.00	0.01	-3.40	16.22
roe	Return on equity	0.00	0.85	0.03	0.01	5.31	41.10
roeq	Quarterly return on equity	-7.05	8.80	0.03	0.10	-2.56	35.80
roic	Return on invested capital	-2.22	1.66	0.00	0.02	-2.80	29.85
rsup	Revenue surprise	-21.24	1.01	-0.08	0.07	-10.40	149.98
salecash	Sales to cash	-4.51	2.33	0.02	0.01	-3.83	64.41
saleinv	Sales to inventory	0.00	2503.48	52.60	10.60	7.68	73.99
salerec	Sales to receivables	0.29	1031.22	25.91	7.59	6.96	62.95
sfe	Scaled earnings forecast	0.00	594.00	11.68	5.94	5.25	31.16
sgr	Sales growth	-36.23	1.09	-0.06	0.04	-14.77	296.41
sp	Sales to price	-0.91	8.50	0.18	0.09	5.68	49.16
		0.00	54.59	2.32	1.10	4.51	29.91

**Table A.1 (Continues)**

Name	Description	Min	Max	Mean	Median	Skewness	Kurtosis
spi	Industry-adjusted sales to price	-0.66	0.19	-0.01	0.00	-5.20	40.89
std_dolvol	Volatility of liquidity (dollar trading volume)	0.18	2.74	0.86	0.79	0.75	0.18
std_turn	Volatility of liquidity (share turnover)	0.02	184.01	3.90	1.90	6.07	64.87
stdcf	Cash flow volatility	0.00	1882.88	9.88	0.14	11.94	178.00
sue	Unexpected quarterly earnings	-5.20	1.70	0.00	0.00	-13.43	550.18
tang	Debt capacity/firm tangibility	0.04	0.98	0.54	0.55	-0.14	0.99
tb	Tax income to book income	-27.70	15.36	-0.10	-0.03	-4.47	66.74
turn	Share turnover	0.00	195.94	1.02	0.52	20.43	1384.13
zerotrade	Zero trading days	0.00	19.95	1.31	0.00	3.03	9.22

**Table A.2 Macroeconomic Variables**

Table A.2 demonstrates the macroeconomic variables we collect from McCracken's page on the website of Federal Reserve Bank of St. Louis. The variables are transformed following the recommended transformation methods. We excluded New Orders for Consumer Goods (ACOGNO), New Orders for Nondefense Capital Goods (ANDENOX) and Trade Weighted USD Index (TWEXMMTH) for data availability issues. In the end, we have 125 macroeconomic variables for our sample period from 1963Q1:2019Q2.

Name	Description	Min	Max	Mean	Median	Skewness	Kurtosis
AAA	Moody's Seasoned Aaa Corporate Bond Yield	-1.18	1.29	0.00	0.00	-0.26	5.52
AAAFFM	Moody's Aaa Corporate Bond Minus FEDFUNDS	-6.27	5.73	2.07	2.21	-0.85	1.17
AMDMNOx	New Orders for Durable Goods	-0.20	0.21	0.00	0.00	-0.14	3.19
AMDMUOX	Unfilled Orders for Durable Goods	-0.03	0.05	0.00	0.00	0.53	1.07
AWHMAN	Avg Weekly Hours : Manufacturing	37.30	42.40	40.81	40.80	-0.41	0.31
AWOTMAN	Avg Weekly Overtime Hours : Manufacturing	-0.90	0.80	0.00	0.00	-0.07	6.76
BAA	Moody's Seasoned Baa Corporate Bond Yield	-1.02	1.57	0.00	0.00	0.55	6.78
BAAFFM	Moody's Baa Corporate Bond Minus FEDFUNDS	-4.05	8.82	3.10	3.18	-0.49	0.31
BOGMBASE	Monetary Base	0.00	0.06	0.00	0.00	15.95	271.36
BUSINVx	Total Business Inventories	-0.02	0.07	0.00	0.00	1.65	22.45
BUSLOANS	Commercial and Industrial Loans	0.00	0.00	0.00	0.00	4.21	26.46
CE16OV	Civilian Employment	-0.01	0.02	0.00	0.00	-0.08	1.85
CES0600000007	Avg Weekly Hours : Goods-Producing	37.20	41.80	40.32	40.30	-0.28	0.28
CES0600000008	Avg Hourly Earnings : Goods-Producing	0.00	0.00	0.00	0.00	5.03	44.16
CES1021000001	All Employees: Mining and Logging: Mining	-0.19	0.20	0.00	0.00	0.42	54.13
CES2000000008	Avg Hourly Earnings : Construction	0.00	0.00	0.00	0.00	5.99	56.58

**Table A.2 (Continues)**

Name	Description	Min	Max	Mean	Median	Skewness	Kurtosis
CES3000000008	Avg Hourly Earnings : Manufacturing	0.00	0.00	0.00	0.00	5.08	36.95
CLAIMSx	Initial Claims	-0.22	0.20	0.00	0.00	0.30	2.18
CLF16OV	Civilian Labor Force	-0.01	0.01	0.00	0.00	0.11	2.62
CMRMTSPLx	Real Manu. and Trade Industries Sales	-0.03	0.05	0.00	0.00	0.01	1.22
COMPAPFFx	3-Month Commercial Paper Minus FEDFUNDS	-2.78	2.22	0.07	0.08	-1.18	8.33
CONSPI	Nonrevolving Consumer Credit to Personal Income	-0.01	0.01	0.00	0.00	-0.34	23.21
CP3Mx	3-Month AA Financial Commercial Paper Rate	-6.29	3.03	0.00	0.00	-2.31	42.14
CPIAPPSL	CPI : Apparel	0.00	0.00	0.00	0.00	4.33	24.18
CPIAUCSL	CPI : All Items	0.00	0.00	0.00	0.00	4.01	23.42
CPIMEDSL	CPI : Medical Care	0.00	0.00	0.00	0.00	3.01	15.77
CPITRNSL	CPI : Transportation	0.00	0.01	0.00	0.00	17.02	358.14
CPIULFSL	CPI : All Items Less Food	0.00	0.00	0.00	0.00	4.66	34.63
CUMFNS	Capacity Utilization: Manufacturing	-3.77	2.22	-0.01	0.03	-0.82	3.76
CUSR0000SA0L2	CPI : All items less shelter	0.00	0.00	0.00	0.00	7.95	98.51
CUSR0000SA0L5	CPI : All items less medical care	0.00	0.00	0.00	0.00	4.24	26.39
CUSR0000SAC	CPI : Commodities	0.00	0.00	0.00	0.00	13.01	233.77
CUSR0000SAD	CPI : Durables	0.00	0.00	0.00	0.00	3.78	16.15
CUSR0000SAS	CPI : Services	0.00	0.00	0.00	0.00	3.74	17.25
DDURRG3M086SBEA	Personal Cons. Exp: Durable goods	0.00	0.00	0.00	0.00	4.73	31.52
DMANEMP	All Employees: Durable goods	-0.05	0.03	0.00	0.00	-1.43	9.48
DNDGRG3M086SBEA	Personal Cons. Exp: Nondurable goods	0.00	0.00	0.00	0.00	12.96	235.86
DPCERA3M086SBEA	Real personal consumption expenditures	-0.03	0.02	0.00	0.00	-0.12	2.94
DSERRG3M086SBEA	Personal Cons. Exp: Services	0.00	0.00	0.00	0.00	2.37	6.50

**Table A.2 (Continues)**

Name	Description	Min	Max	Mean	Median	Skewness	Kurtosis
DTCOLNVHFNM	Consumer Motor Vehicle Loans Outstanding	0.00	0.03	0.00	0.00	9.94	121.28
DTCTHFNM	Total Consumer Loans and Leases Outstanding	0.00	0.06	0.00	0.00	19.57	405.68
EXCAUSx	Canada / U.S. Foreign Exchange Rate	-0.06	0.11	0.00	0.00	0.52	8.22
EXJPUSx	Japan / U.S. Foreign Exchange Rate	-0.11	0.08	0.00	0.00	-0.51	1.80
EXSZUSx	Switzerland / U.S. Foreign Exchange Rate	-0.09	0.12	0.00	0.00	-0.13	1.60
EXUSUKx	U.S. / U.K. Foreign Exchange Rate	-0.11	0.10	0.00	0.00	-0.46	2.81
FEDFUNDS	Effective Federal Funds Rate	-6.63	3.06	0.00	0.01	-2.39	47.47
GS1	1-Year Treasury Rate	-3.91	1.90	0.00	0.00	-1.51	18.10
GS10	10-Year Treasury Rate	-1.76	1.61	0.00	0.00	-0.43	5.94
GS5	5-Year Treasury Rate	-2.03	1.86	0.00	0.00	-0.41	6.68
HOUST	Housing Starts: Total New Privately Owned	6.17	7.82	7.22	7.29	-0.98	1.00
HOUSTMW	Housing Starts, Midwest	4.08	6.38	5.55	5.66	-0.89	0.29
HOUSTNE	Housing Starts, Northeast	3.58	5.98	5.04	5.04	-0.34	-0.11
HOUSTS	Housing Starts, South	5.44	7.08	6.42	6.44	-0.57	0.36
HOUSTW	Housing Starts, West	4.37	6.47	5.78	5.84	-0.93	0.62
HWI	Help-Wanted Index for United States	-633.00	880.00	6.94	6.00	0.07	1.93
HWIURATIO	Ratio of Help Wanted/No. Unemployed	-0.17	0.11	0.00	0.00	-0.42	1.70
INDPRO	IP Index	-0.04	0.03	0.00	0.00	-0.99	5.11
INVEST	Securities in Bank Credit at All Commercial Banks	0.00	0.00	0.00	0.00	5.73	57.39
IPB51222S	IP: Residential Utilities	-0.13	0.14	0.00	0.00	-0.22	1.54

**Table A.2 (Continues)**

Name	Description	Min	Max	Mean	Median	Skewness	Kurtosis
IPBUSEQ	IP: Business Equipment	-0.08	0.05	0.00	0.00	-1.01	5.36
IPCONGD	IP: Consumer Goods	-0.03	0.04	0.00	0.00	-0.02	1.85
IPDCONGD	IP: Durable Consumer Goods	-0.11	0.13	0.00	0.00	0.02	6.60
IPDMAT	IP: Durable Materials	-0.06	0.05	0.00	0.00	-0.82	3.84
IPFINAL	IP: Final Products (Market Group)	-0.03	0.03	0.00	0.00	-0.32	1.90
IPFPNSS	IP: Final Products and Nonindustrial Supplies	-0.03	0.03	0.00	0.00	-0.48	2.31
IPFUELS	IP: Fuels	-0.10	0.15	0.00	0.00	0.65	7.12
IPMANSICS	IP: Manufacturing (SIC)	-0.05	0.03	0.00	0.00	-0.91	4.25
IPMAT	IP: Materials	-0.07	0.03	0.00	0.00	-1.29	7.60
IPNCONGD	IP: Nondurable Consumer Goods	-0.02	0.02	0.00	0.00	-0.08	0.38
IPNMAT	IP: Nondurable Materials	-0.08	0.05	0.00	0.00	-1.40	11.05
ISRATIOx	Total Business: Inventories to Sales Ratio	-0.06	0.11	0.00	0.00	0.46	3.86
M1SL	M1 Money Stock	0.00	0.00	0.00	0.00	10.33	128.42
M2REAL	Real M2 Money Stock	-0.02	0.03	0.00	0.00	0.73	3.58
M2SL	M2 Money Stock	0.00	0.00	0.00	0.00	5.65	50.69
MANEMP	All Employees: Manufacturing	-0.03	0.02	0.00	0.00	-1.47	6.60
MZMSL	MZM Money Stock	0.00	0.01	0.00	0.00	19.28	420.20
NDMANEMP	All Employees: Nondurable goods	-0.02	0.01	0.00	0.00	-1.12	4.60
NONBORRES	Reserves Of Depository Institutions	-195.01	170.56	-0.66	-0.35	-3.28	152.65
NONREVSL	Total Nonrevolving Credit	0.00	0.01	0.00	0.00	16.19	281.62
OILPRICEx	Crude Oil, spliced WTI and Cushing	0.00	0.73	0.01	0.00	19.83	456.21
PAYEMS	All Employees: Total nonfarm	-0.01	0.01	0.00	0.00	-0.48	2.68
PCEPI	Personal Cons. Expend.: Chain Index	0.00	0.00	0.00	0.00	3.09	11.28

**Table A.2 (Continues)**

Name	Description	Min	Max	Mean	Median	Skewness	Kurtosis
PERMIT	New Private Housing Permits (SAAR)	6.24	7.79	7.18	7.21	-0.75	0.43
PERMITMW	New Private Housing Permits, Midwest (SAAR)	4.34	6.23	5.50	5.60	-0.78	-0.02
PERMITNE	New Private Housing Permits, Northeast (SAAR)	4.06	6.05	5.07	5.08	-0.23	-0.26
PERMITS	New Private Housing Permits, South (SAAR)	5.37	7.01	6.32	6.36	-0.30	-0.43
PERMITW	New Private Housing Permits, West (SAAR)	4.57	6.63	5.79	5.85	-0.89	0.57
PPICMM	PPI: Metals and metal products:	0.00	0.02	0.00	0.00	5.68	41.97
REALLN	Real Estate Loans at All Commercial Banks	0.00	0.00	0.00	0.00	8.25	106.81
RETAILx	Retail and Food Services Sales	-0.07	0.06	0.00	0.00	-0.30	4.93
RPI	Real Personal Income	-0.05	0.04	0.00	0.00	-0.64	22.59
S_P__indust	S&P's Common Stock Price Index: Industrials	-0.22	0.11	0.01	0.01	-1.00	4.03
S_P_500	S&P's Common Stock Price Index: Composite	-0.23	0.11	0.01	0.01	-1.02	4.23
S_P_div_yield	S&P's Composite Common Stock: Dividend Yield	-0.64	0.59	0.00	-0.01	0.60	6.06
S_P_PE_ratio	S&P's Composite Common Stock: Price-Earnings Ratio	-0.22	0.24	0.00	0.00	0.07	5.83
SRVPRD	All Employees: Service-Providing Industries	-0.01	0.01	0.00	0.00	0.06	4.34
T10YFFM	10-Year Treasury C Minus FEDFUNDS	-6.51	3.85	1.04	1.21	-1.11	2.12
T1YFFM	1-Year Treasury C Minus FEDFUNDS	-5.00	1.69	0.02	0.13	-2.23	8.91



**Table A.2 (Continues)**

Name	Description	Min	Max	Mean	Median	Skewness	Kurtosis
T5YFFM	5-Year Treasury C Minus FEDFUNDS	-6.31	3.16	0.70	0.87	-1.35	3.27
TB3MS	3-Month Treasury Bill	-4.62	2.61	0.00	0.01	-1.85	28.22
TB3SMFFM	3-Month Treasury C Minus FEDFUNDS	-5.37	0.68	-0.49	-0.25	-2.56	9.28
TB6MS	6-Month Treasury Bill	-4.23	2.17	0.00	0.01	-1.83	23.64
TB6SMFFM	6-Month Treasury C Minus FEDFUNDS	-5.01	1.19	-0.35	-0.13	-2.57	10.00
TOTRESNS	Total Reserves of Depository Institutions	0.00	1.25	0.00	0.00	18.26	364.58
UEMP15OV	Civilians Unemployed - 15 Weeks & Over	-0.18	0.24	0.00	0.00	0.37	1.27
UEMP15T26	Civilians Unemployed for 15-26 Weeks	-0.36	0.29	0.00	0.00	-0.05	0.92
UEMP27OV	Civilians Unemployed for 27 Weeks and Over	-0.21	0.28	0.00	0.00	0.29	1.21
UEMP5TO14	Civilians Unemployed for 5-14 Weeks	-0.22	0.23	0.00	0.00	0.28	1.18
UEMPLT5	Civilians Unemployed - Less Than 5 Weeks	-0.22	0.27	0.00	0.00	-0.01	1.34
UEMPMEAN	Average Duration of Unemployment (Weeks)	-2.70	2.50	0.01	0.00	-0.08	2.25
UMCSENTx	Consumer Sentiment Index	-14.70	17.30	0.01	0.00	0.01	2.06
UNRATE	Civilian Unemployment Rate	-0.70	0.90	0.00	0.00	0.52	1.98
USCONS	All Employees: Construction	-0.04	0.06	0.00	0.00	0.20	5.92
USFIRE	All Employees: Financial Activities	-0.01	0.01	0.00	0.00	-0.51	1.20

**Table A.2 (Continues)**

Name	Description	Min	Max	Mean	Median	Skewness	Kurtosis
USGOOD	All Employees: Goods-Producing Industries	-0.02	0.02	0.00	0.00	-1.26	4.61
USGOVT	All Employees: Government	-0.01	0.02	0.00	0.00	0.67	8.38
USTPU	All Employees: Trade, Transportation & Utilities	-0.01	0.01	0.00	0.00	-0.43	1.33
USTRADE	All Employees: Retail Trade	-0.01	0.01	0.00	0.00	-0.24	2.71
USWTRADE	All Employees: Wholesale Trade	-0.01	0.01	0.00	0.00	-0.57	0.94
VXOCLSx	VXO	8.02	67.15	19.05	17.48	1.96	6.97
W875RX1	Real personal income ex transfer receipts	-0.06	0.04	0.00	0.00	-1.77	32.51
WPSFD49207	Producer Price Index by Commodity: Final Demand: Finished Goods	0.00	0.00	0.00	0.00	6.62	61.76
WPSFD49502	PPI: Final Demand: Personal Consumption Goods (Finished Consumer Goods)	0.00	0.00	0.00	0.00	7.10	70.09
WPSID61	PPI: Intermediate Demand by Commodity Type: Processed Goods for Intermediate Demand	0.00	0.00	0.00	0.00	7.84	77.81
WPSID62	PPI: Intermediate Demand by Commodity Type: Unprocessed Goods for Intermediate Demand	0.00	0.04	0.00	0.00	7.48	69.10