

# Is there an influence of top financial companies on the S&P500?

Simone Pisano

November 18, 2024

## **Abstract**

This study looks at how significant asset management businesses affect the S&P 500 index from 2013 to 2023. We use models like ordinary least squares, random forest regression, and feedforward neural network to analyze data from the biggest listed firms in the industry like Blackrock, SEI and Vanguard. Our findings show that these firms have a prominent impact on the index, which varies according to market conditions. The findings influence investment strategies and policy decisions. What these results suggest could be improving our understanding of market dynamics and stability in different economic conditions.

# 1 Introduction

In 2020 the proportion of active versus passive investment portfolios worldwide was 53% versus 47%. A slightly higher proportion of people's investment portfolios were passive, than the share of active investment portfolios. Passive investment portfolios usually involve less action and are based on long-term investments, such as index portfolios. Active investment portfolios are based on more action following more short-term price fluctuations, where the manager tries to beat the market through research, analysis and own judgement. The dichotomy between active and passive investment portfolios mirrors a profound divergence in investment philosophies, each offering distinct advantages and appealing to different investor profiles. Active investment portfolios are characterized by dynamic decision-making, where fund managers engage in comprehensive research, analysis, and short-term trading to potentially outperform market benchmarks.

The objective of this paper is to gain a more profound understanding of the influence that elite asset managers like Blackrock, Vanguard, SEI investment, Franklin Resources Inc, Blackstone Investments, Brookfield Asset Management and Fidelity International, have on the S&P 500 Index during bull and bear market phases. As we can imagine, the impact these financial firms have on the overall market can vary immensely depending on the state of the market. Investment strategies, systemic stability, risk management, are just some of the practices which make analyzing and studying the differences in impact on the S&P 500 of the top financial firms during bullish and bearish market periods, so crucial. In order to fully understand the topic of discussion, let us delve into what bull and bear markets actually consist of and why they are labelled with these imaginative appellatives. Interestingly enough, the origins of these expressions are still unclear. The most popular version would be related to how each animal is said to attack its prey. A bull will shove its horns vertically into the air, whereas a bear will swipe down. These gestures metaphorically resemble the fluctuations of a market, with bull markets trending up and bear markets trending down. While there is little evidence to suggest that this is the legitimate etymology of the words, the attack strategies can help you recall which way bull and bear markets move. The bull market is regarded as a moment in the economy in which shareholders and investors tend to have very positive sentiment towards the direction and potential behavior of the market, where the stock prices go up, and the general feeling is the one of confidence. In fact, a bull market symbolizes a period of encouraging market conditions where economic indicators such as GDP growth, consumer expenditures, business profits, reveal

overall positive trends. It is exactly in instances like this that less risk-prone investors will pull the trigger and show greater appetite for risk. This will inevitably cause an increase in trade volume and market liquidity. Conversely, in a bear market, the global sentiment is the one of uncertainty. Shareholders and investors feel less confident and become more risk averse. It can be compared to a time of financial shrinkage, where liquidity and trade volume decrease, resulting in phases of recession or stagnation. Metrics like company earnings, investment, investor confidence, shareholder trust, will point towards a negative market landscape. Investigating bear markets requires analysing and studying the factors that cause market crises, assessing the impact on investments, and identifying risk management strategies to reduce potential losses.

Quite wide amount of research has been done with regards to market phases that still has a crucial importance nowadays. For instance, "The Intelligent Investor," (Benjamin Graham, 1949) is widely regarded as one of the most impactful publications on value investing. It offers enduring sagacity and pragmatic guidance for investors managing the fluctuations of the stock market. In his work Graham promotes the idea of oppositional thinking, which can be reflected in challenging the prevailing sentiment and purchasing assets when everyone else is selling (bear markets) and divesting when others are purchasing (bull markets). He holds the belief that markets are susceptible to episodes of irrational exuberance and negativity, which in turn present opportunities for discerning investors to exploit undervalued securities. "Irrational Exuberance," (Robert J. Shiller, 2000) is another research based book that delves into emergence of bullish and bearish markets. Shiller discusses the psychological variables that exert an influence on investor behavior in both bullish and bearish markets. He emphasizes the significance of cognitive biases, such as excessive confidence and the tendency to follow the crowd, in causing fluctuations in market activity. The idea of the book contends that comprehending the psychological foundations of market cycles is crucial for forecasting and controlling volatility in the markets. To put it short, Shiller's study highlights the significance of both investor mood, and social factors in influencing periods of market volatility, offering a complete framework for comprehending and managing the intricacies of the financial system.

Reverting to the relevance of the research question. The comprehension of this topic presents useful knowledge about market behavior dynamics, enabling investors to anticipate and adjust to evolving conditions in the market. Gaining understandings of how financial companies influence market move-

ments during growth periods can provide valuable information for investors looking to maximize the value of their investments. Similarly, knowing their behavior during market downturns can assist investors in minimizing losses and identifying potential investment prospects. In addition, analyzing the impact of leading financial firms on the S&P 500 index during different market stages has significant consequences for policy decisions and regulatory supervision. Precise evaluations of market dynamics are crucial for policy-makers to develop successful strategies that promote market stability and safeguard investors. Having a grasp on the behavior of financial companies in various market conditions can help shape regulatory frameworks and measures that eventually minimize systemic risks and promote the stability of the financial system as a whole. On top of that, this research is of utmost importance for risk management. Financial companies frequently act as early indicators of more general market movements and financial stability. Examining their impact on the S&P 500 during different market phases enables traders and risk managers to evaluate the risk exposure and implement suitable risk reduction strategies. It further enhances the comprehension of systemic risks and vulnerabilities within the economy.

A wide research has been done with relation to various index investments and its performance depending on bullish and bearish trends. In order to thoroughly investigate our presented question and answer it as precise as possible, let's try to examine a number of useful research papers that had actually made a contribution to the financial sector itself. As it was mentioned in the abstract part, prior evidence suggested that financial shocks may have had a significant impact on the real underlying economy, this implies that our research question was at least brought to light by practitioners. But why it is important to investigate economic shocks with relation to S&P500 or any other index if the initial question does not include the exploration of shocks themselves? Due to the fact that both bear and bull markets have quite high correlation with recession and expansion accordingly, it is necessary to understand how potential financial shocks can affect the economic sector. Nonetheless, it should be mentioned that despite causal relation between the above mentioned variables, bullish and bearish phases can actually exist without appearance, or long time after, the specific stages of economic cycle (in our case recession and expansion).

Therefore, the first paper that will help us understand the cause-effect relationship of index performance and different stages of economic development is "The Impact of Stock Market Performance upon Economic Growth" (Najeb M.H. Masoud, 2013; International Journal of Economics and Financial

Issues pp.788-798). As it can be understood from the name, the aim of this research paper is to investigate the dependence of one variable on another. The research provides a deep dive into both empirical and theoretical cases that eventually lead to an agreed conclusion of the causal inference presence. The paper examines eight countries in total (USA, UK, France, Japan, Canada, Australia, Germany, and Switzerland), including both emerging and advanced economies. Authors came up to the conclusion that there exists a significant impact of stock market performance upon economic growth. The outcome has confirmation in the form of both qualitative and quantitative indicators that are being discussed throughout the whole article. Since the topic of correlation between different economic development stages, with the focus on financial shocks, and actual financial and economic sector has been popular for a number of decades, let's take a look at works that gained its glory many years ago. Thus, the first example is the article "Stock Markets, Banks, and Economic Growth" (Ross Levine and Sara Zervos, 1998; The American Economic Review pp. 537-558). Unlike the first mentioned paper, glancing down on the article name is not enough in order to get even remote understanding of the goal that this research paper presents. This research paper focuses on the robust economic growth, rather than recession and stagnation, and its consequences and influences on financial sector. Factors that influence the development of banking sector, market liquidity, and overall economic growth are being thoroughly examined and controlled for.

The sample of approximately 50 countries' markets from the period of 1976 to early 90s has been studied with regards to market size, volatility, accumulation of capital and etc. Authors had also mentioned that in order to get a better understanding on the dependence of the LT growth on the financial and economic systems, there is a need to come up with various models that would depict the simultaneous formation and development of both banks and stock markets, while offering distinct financial services. Eventually, Levine and Zervos prove the existence of quite strong correlation between economics growth and market performance, highlighting the dependence of one variable on another. Article also demonstrates a positive relationship of above mentioned financial factors with economic growth itself. It can be seen that both papers depict quite similar results. However, in order to start answering the question of whether there is a difference in the influence of top financial companies on the S&P500 index between bull and bear market phases we need to investigate additional papers. This time the main focus of the research should be aimed at examining the influence of economic state of the country (or worldwide) on publicly listed companies, which in its turn will have certain effect on both the composition and per-

formance of different indexes.

Before discussing articles that take into consideration the impact of specific companies on indexes, let's first of all look at the influence of external factors on different firms that do actually constitute the indexes themselves. For instance, the work of Bahaaeddin Ahmed Alareeni and Allam Hamdan (2020) "ESG impact on performance of US S&P500 – listed firms". Authors try to understand the behavior of firms depending on the changes in three factors such as environmental, social and governance ones. On top of that, the goal of the research is to find if there is an actual relationship between company's activity (actions and measures) and the performance of the market itself. In order to answer the last question, authors investigate various models, for instance Tobin's Q one, and quantitative factors that depict potential influence, and these are returns on both assets and equity. Article demonstrates the data on all 500 companies that are listed in S&P500 throughout the period of ten years – post global economic crisis and pre Covid19 period. The chosen period was implemented in order to eliminate the possibility of bias occurrence, such as herd bias, panic bias, etc. Eventually Alareeni and Hamdan came to the conclusion that not only the dependence of companies' behavior due to changes in social, governance, and environmental factors, but also that there exists positive relationship between market performance and subsequent actions of the specified firm. Now, once we are faced with the confirmation that external factors do actually influence companies' actions (and most of the time these factors carry out a positive relationship within firms) we can focus on the influence of the companies on indexes depending on their proportionality and other determinants. Therefore, let's have a look at the following research paper "Board Composition: Balancing Family Influence in S&P 500 Firms" (Ronald C. Anderson, David M. Reeb, 2004). In our case, by "other determinants" will be presented the factor of family influence, that, according to the authors, can potential affect firm's operational and development side, which in its turn can influence the performance of the index itself. However, since this factor is not the major target of our research paper, we will not be focusing on it for too much. Moving on to the previously examined factors that can actually help us understand the way top financial companies influence S&P500 index. One of the papers that carries an instructive character in relation to a given topic is "Corporate News Disclosure and Competitive Advantage: What Factors Influence S&P 500 Companies? Competitive Advantage During 2022 Economic Crisis". From the very name it is quite obvious what is the main topic of the article. In order to correctly answer the question, both qualitative and quantitative researches were implemented. The sample size was not drastic, 11 months,

nevertheless a sufficient amount of information was collected in order to carry out the research. The conclusion of the investigation states that the majority of largest firms do in fact publicly demonstrate their competitive benefit. On top of that, the research shows that the larger the company, the easier it is to “stay afloat”, in other words the bigger its competitive privilege during the period of crisis.

## 2 Data

We started looking for some proxies to measure active investment strategy vs passive ones. The main idea behind our decision was to choose the biggest actively managed investment corporations out there that are publicly listed. Among the many potential companies we decided to include the biggest ones: Blackrock, Vanguard, SEI investment, Franklin Resources inc, Blackstone Investments, Brookfield Asset Management and Fidelity International among the others. The passive investment strategy is well proxied by the Standard and Poor's 500, which is a broad market index reflecting the underlying American market very efficiently since it includes approximately 80% of U.S. public companies. The S&P500 index is a free-float capitalization-weighted index, hence its components represent day-by-day the composition of the American market. The data have been downloaded from well-known resources such as Statista, Yahoo Finance and French Data Library. The procedure we followed for data cleaning was simple: merge the documents on date using a brief Python script and store the data in two separate documents, one for each group of regressors. The time horizon falls from 1 January 2013 to 10 December 2023, we included a full decade in order to higher the number of observations for the purpose of lending unbiased estimates, while cross-sectional analyses have their merits, a decade-long dataset provides a richer context for understanding the performance dynamics of passive and active investment strategies. It allows for a more nuanced exploration of how these strategies fare across various market conditions and over different economic cycles, contributing to a more informed and reliable assessment. In the data cleaning process we estimated the difference in returns through logarithmic approximation. This procedure is informative in the context of interpreting the estimates of our regression model. A "log-log" model can be interpreted as the elasticity of one variable in terms of the other. Data on the Fama-French variables such as Small minus Big (SMB) and High minus Low (HML) come from Fama's own website. Each variable has some explanatory power in the model hence we decided to use it as controls to our model. We wanted to to some testing in order to find a potential influential observation in the dataset but we didn't find any which had a great impact. In the first three regression the highest quantity we obtained was ( $e^{-4}$ ) which has undoubtedly a small influence on the regression model.



### 3 Methodology

We want to estimate the baseline model

$$y = \beta X + \epsilon \quad (1)$$

where  $y$  is the return of the passively managed portfolio, proxied by the return of the SPX500 in a decade;  $X$  includes all the regressors, i.e., every actively managed portfolio we decided to include in the regression model as introduced in the Introduction. The  $\beta$  is the vector of coefficients we want to estimate through Ordinary Least Squares, whereas the  $\epsilon$  is the vector of the error term. Ordinary Least Squares gives the Best Linear Unbiased Estimation (BLUE) of the underlying dataset, hence we start looking at the "gross" estimates we get by running this simple regression model. For each regression run we always have to estimate the population model in (1) which is a standard procedure in econometrics. We started analysing the data through Ordinary Least Squares in order to lend unbiased estimates about the true population model, but we still have got to check for unbiasedness of our estimator, which involves both testing and assuming something about our variables of interest and the associated error term.

Random forest is an ensemble learning method mainly utilised for classification, regression and some other tasks that work by building up a multitude of decision trees at training time. Through the usage of multiple learning algorithms, the ensemble learning method enable better predictive performance than we would usually expect from any of the standalone constituent learning algorithms. When it comes to classification tasks, the random forest's output is the class selected by most trees. For regression tasks, the mean or average prediction of the individual trees is returned. Random decision forests correct for decision trees' habit of overfitting to their training set. The Chinese computer scientist and researcher called Tin Kam Ho developed the first random forests' algorithm in 1995. It is based on the random subspace method, which is a method put into action to implement the "stochastic discrimination" approach to classification. An extension of the algorithm was consequently idealized by Leo Breiman and Adele Cutler, who registered the trademark in 2006. The extension in question mixes up Breiman's "Bootstrap aggregating" idea with randomly selected features aiming to construct a series of decision trees with controlled variance. From a historical perspective, Amit and Geman first introduced the idea of searching over a random subset of the available decisions while splitting a node, in the context of growing a single tree. From this idea, Breiman developed the notion of random forests. Also, Ho's idea of random subspace selection directly influenced the design

of random forests. This method involves growing a forest of trees, introducing variation among them by projecting the training data into a randomly chosen subspace before fitting each tree or node. Additionally, the concept of randomized node optimization, whose pioneer was Thomas G. Dietterich, replaced deterministic optimization with a randomized procedure in order to select decisions at each node.

Decision tree is a commonly used method for different machine learning tasks. Tree learning is renowned to almost be an off-the-shelf procedure for data mining, because it is invariant to scaling and other transformations of feature values, it gathers irrelevant features, and produces inspectable models. On the other hand, they are seldomly accurate. In particular, deep grown trees tend to follow highly irregular patterns: they overfit their training sets. Hence they are low biased, but have very high variance. Random forests compute the average of multiple deep decision trees, trained on different parts in the same training set, with the purpose of reducing the variance. The side effect that occurs is a slight surge in the bias and some loss of interpretability, but in general it efficiently boosts the performance in the final model. The training algorithm for random forests applies the bootstrap aggregating technique. Given a training set  $X = x_1, \dots, x_n$  with certain responses  $Y = y_1, \dots, y_n$ , the bootstrap aggregating technique repeatedly selects a random sample with replacement of the training set and fits trees to these samples. For  $b = 1, \dots, B$ :

1. Sample, with replacement,  $n$  training examples from  $X, Y$ ; call these  $X_b, Y_b$ .
2. Train a classification or regression tree  $f_b$  on  $X_b, Y_b$ .

Afterwards, predictions for unseen samples  $x'$  are made by averaging the predictions from all the individual regression trees on the predictions for unseen samples themselves.

$$\hat{f} = \frac{1}{B} \sum_{b=1}^B f_b(x') \quad (2)$$

This bootstrap procedure leads to a more suitable model because it lowers the variance, maintaining the bias at the same level. This implies a high sensitiveness in the predictions of a single tree to noise in its training set, as long as the trees are not correlated. Training many trees on a single training set would lead to a strong correlation among trees (or even a repetition of the

same tree multiple times, if the training algorithm is deterministic). Bootstrap sampling is a way to de-correlate the trees by providing them different training sets. Additionally, an estimate of the uncertainty of the prediction can be made as the standard deviation of the predictions from all the individual regression trees on predictions for unseen samples. On the average, the number of trees used range from few hundred to several thousand, based on the size and nature of the training set. A great way to find the optimal number of trees could be cross-validation, or the observation of out-of-bag error. After fitting a certain number of trees, both the training and test errors stabilize.

Out-of-bag error is a methodology used to measure the predicted error in random forests, boosted decision trees, and other machine learning models that utilize bootstrap aggregating. Bootstrap aggregating uses sub-sampling with replacement, generating training samples from which the model collects information. The out-of-bag error computes the average prediction error for each training sample  $x_i$ , utilizing only the trees that didn't include  $x_i$  in their bootstrap sample. Bootstrap aggregating enables the definition of an out-of-bag estimation of the prediction performance enhancement by assessing predictions on observations not utilized in constructing the subsequent base learner.

Calculating the Out-Of-Bag Error:

Out-of-bag is an efficient way to test the performance of the model. The computation of its error is based on the implementation of the model, but it generally follows four principles:

1. Find all the models (trees if referring to random forest) that are not trained by the out-of-bag.
2. Compare the majority vote of these models' outcomes for the out-of-bag instance with the true value of the OOB case.
3. Compile the out-of-bag error in the dataset for all the cases.
4. The bootstrap aggregating process can be customized to fit different necessities in a model. The sample size of the bootstrap training should be close to the one of the original set in order to ensure a suitable level of accuracy. Moreover, the number of trees in the forest should be taken into consideration to find out the real error. OOB error will, in fact, stabilize over many repetitions. Thus, starting off with a considerable number of iterations is extremely positive.

A neural network is a computational model that, inspired by the workings of the human brain, can develop and recognize complex patterns and correlations in the data it processes.

The neural network model is constructed up of algorithms that interact in the same way that neurons in the human brain do. These algorithms are able to process input and generate output. Neural networks are made up of layers of artificial neurons that are structured hierarchically. Each of these neurons receives input and produces output, which is then transformed mathematically and passed to other neurons in the network. Each neuron includes inputs and outputs, which are linked together by weights that determine the relative importance of each output and input inside the neuron.

Neural networks are trained using a process called deep learning, in the training phase, the network is exposed to an intense number of input and output samples, and to minimize the error present between predicted output and desired output, it updates its weights, through this training process, the network is able to learn and identify patterns and relationships in the data and use that knowledge and resources to make predictions or make decisions.

Artificial neural networks are among the most relevant computational models in artificial intelligence and machine learning processes. Neural networks denote different *architectures of connections* between neurons. In our analysis, we focus on the **Feedforward** architecture, which represents the flow of data from input to output, without feedback or cyclic connections.

**Feedforward** presents a structure according to which, each *layer* of the neuron is connected to the next, information flows forward unidirectionally, through these connections without encountering any form of *recursion*, that is, without recourse within the network structure of cyclic connections, and without processing feedback between input and output.

Feedforward neural networks can be categorized into two types:

1. **Single-layer network:** formed by an input node and an output node, and data propagates in only one direction. There will be only one layer performing data processing since no hidden nodes are present.
2. **Multi-layer network:** Hidden layers are also present but, as in the previous case, the signal propagation is in only one direction.

A feedforward neural network presents a unidirectional flow of information, which is transmitted from the input layer, to the output layer, without drafting feedback or feedback of the results. This model presents a structure with several layers of neurons, in addition to the first layer i.e., the one referring to *input*, and the last one instead of *output*. In between are the so-called

*hidden* or intermediate layers. Data processing takes place sequentially in each of the previously mentioned layers. Transformations are applied to the input data, thanks to the parameters and the activation function. The latter has the role of inputting non-linearity into the network; it is applied to the weighted sum of the inputs of the reference neuron. The goal of neural networks in the training process is to minimize the difference between its outputs and the desired ones, this is done by adjusting its weights and biases.

The activation  $a$  of a neuron is given by the following formula:

$$a_j^{(l)} = f \left( \sum_i w_{ij}^{(l)} x_i + b_j^{(l)} \right) \quad (3)$$

where:

- $a_j^{(l)}$  is the activation of the  $j$ -th neuron in the  $l$ -th layer,
- $w_{ij}^{(l)}$  are the weights connecting the  $i$ -th input to the  $j$ -th neuron in the  $l$ -th layer,
- $x_i$  are the input values from the previous layer (or input data for the first layer),
- $b_j^{(l)}$  is the bias associated with the  $j$ -th neuron in the  $l$ -th layer,
- $f$  is the activation function.

The cost function  $J$  for a neural network is defined as:

$$J(W, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i) + \frac{\lambda}{2m} \sum_{l=1}^L \sum_{i,j} (w_{ij}^{(l)})^2 \quad (4)$$

where:

- $\hat{y}_i$  is the predicted output,
- $y_i$  is the actual output,
- $L$  is the loss function, such as mean squared error for regression,
- $\lambda$  is the regularization parameter to prevent overfitting,
- $m$  is the number of training examples,
- $W, b$  are the collections of all weights and biases in the network,

- $L$  is the total number of layers.

In conclusion we are going to talk about the two tests we decided to employ to check whether there was any major issue with our models: the Diebold-Mariano test and the Fluctuation test from Giacomini et al. (2010). The **Diebold-Mariano** (DM) test is a statistical approach for comparing the prediction accuracy of two different forecasting models. It focuses on the variations in forecasting errors between the two models throughout a given sample. This method enables a direct test of the null hypothesis that the two models have the same prediction accuracy vs the alternative that they do not. The Diebold-Mariano test is widely used in finance to evaluate model performance, particularly time series forecasts.

The **Fluctuation Test** is a statistical approach for evaluating the stability of forecasting models' relative performance over time, particularly in unpredictable circumstances. This test depicts the standardized differences in forecasting errors between two models over time, determining whether and when their performance diverges significantly. It permits the detection of periods when one model outperforms the other, which may be obscured by averaging performance across the full dataset. The Fluctuation Test does not specify an alternative hypothesis, which could indicate that it has lesser power than tests constructed with a specific alternative in mind. However, it is especially useful for emphasizing changes in model efficacy due to changing market or data conditions, making it an excellent choice for examining time-varying model performance.

The **Mean Squared Forecast Error** (MSFE) differences between two models at time  $t$  is calculated as follows:

$$\Delta\text{MSFE}_t = \text{MSFE}_{\text{model 1},t} - \text{MSFE}_{\text{model 2},t} \quad (5)$$

where  $\text{MSFE}_{\text{model},t}$  represents the mean squared forecast error of a model at time  $t$ . This metric is essential for computing the Fluctuation test and is a powerful measure per se.

## 4 Results and Conclusions

The first model is the baseline we chose to start to delve in deeper into our framework of analysis. Its estimates are all significantly different from zero and in particular are all positive. This implies a positive influence of the financial sector on the SPX500 as we would have expected from theory. A special mention is needed for Vanguard Inc. which has a 0.47 beta, which is quite high in relative terms. All the coefficients are statistically significant with a very close to zero p-value, apart from the intercept which is not significant at 10% significance level. The Adjusted  $R^2$  of this model is 0.9339 which is quite high as a result but indicates a strong explanatory power of the independent variables on the dependent variable, this is the main reason why we decided to do out of sample estimation.

The second model includes some factors from the famous Fama-French model suggesting to analyze not only beta for an investment return but also other factors. The factors influencing the stock returns we used were "Small Minus Big," meaning that in a long-run analysis stocks of small caps companies will show higher returns than stocks of big caps companies, and the "High Minus Low," meaning preferring companies that are considered to be undervalued relative to their intrinsic value. Developing the model in MATLAB, we were able to find that the coefficients of the factors taken into account, turn out to be not very positive and both insignificant, this is not relevant since the coefficients themselves are control variables, thus being able to isolate the effect of the factors taken into account.

To address the research question, we constructed a Random Forest model with 16 optimal trees utilizing out-of-bag (OOB) error analysis. The model was developed using the TreeBagger function and employed regression methodology, enabling OOB predictions for performance assessment. Given new data  $X_{\text{out1}}$ , we derived predicted values  $\hat{y}_{\text{rf}}$ . The Mean Square Error (MSE) of our Random Forest model was subsequently computed and gave a comparable result as in the case of other models ( $10^{-5}$ ).

We used a feedforward neural network (FNN) to analyze the predicting accuracy of our dataset. The network architecture comprised of two hidden layers, each with 10 and 8 neurons, and used a sigmoid activation function. We trained the model for up to 1,000 epochs or until early stopping was triggered, which is a mechanism for preventing overfitting that halts training if the mean squared error (MSE) does not improve after 50 consecutive epochs. The MSE for this model is  $1.002 \times 10^{-5}$  which is significantly lower than the initial value we obtained with a plain-vanilla FNN. This result indicate a strong predictive accuracy which is quite remarkable in the case of such a small dataset (compared to Big Data).

Table 1: Regression Model Results

|          | Model 1                             | Model 2                             | Model in 1                          | Model in 2                          |
|----------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| Constant | $-2.86 \times 10^{-5}$<br>(5.4e-05) | $-2.82 \times 10^{-5}$<br>(5.4e-05) | $-3.94 \times 10^{-5}$<br>(5.8e-05) | $-3.93 \times 10^{-5}$<br>(5.8e-05) |
| BLK      | 0.0390<br>(0.007)                   | 0.0390<br>(0.007)                   | 0.0490<br>(0.008)                   | 0.0490<br>(0.008)                   |
| BWG      | 0.0747<br>(0.0089)                  | 0.0750<br>(0.0089)                  | 0.0749<br>(0.0099)                  | 0.0754<br>(0.011)                   |
| BX       | 0.0310<br>(0.012)                   | 0.0310<br>(0.012)                   | 0.0369<br>(0.014)                   | 0.0370<br>(0.014)                   |
| BEN      | 0.0152<br>(0.004)                   | 0.0153<br>(0.004)                   | 0.0200<br>(0.005)                   | 0.0201<br>(0.005)                   |
| FHI      | 0.0460<br>(0.006)                   | 0.0461<br>(0.006)                   | 0.0382<br>(0.006)                   | 0.0383<br>(0.006)                   |
| SEIC     | 0.1908<br>(0.014)                   | 0.1898<br>(0.014)                   | 0.1936<br>(0.017)                   | 0.1920<br>(0.017)                   |
| VGT      | 0.4653<br>(0.009)                   | 0.4654<br>(0.009)                   | 0.4691<br>(0.012)                   | 0.4693<br>(0.012)                   |
| SMB      |                                     | 0.0096<br>(0.0109)                  |                                     | 0.0132<br>(0.013)                   |
| HML      |                                     | 0.0027<br>(0.008)                   |                                     | 0.0034<br>(0.01)                    |

We then decided to run some tests in order to validate our analysis and we chose the Diebold-Mariano test and the Giacomini-Rossi test. The results of the Diebold-Mariano test is that between them (model2, modelRF and modelNN) there is no significant difference in predictive accuracy, which implies the three models are performing at least equivalently.



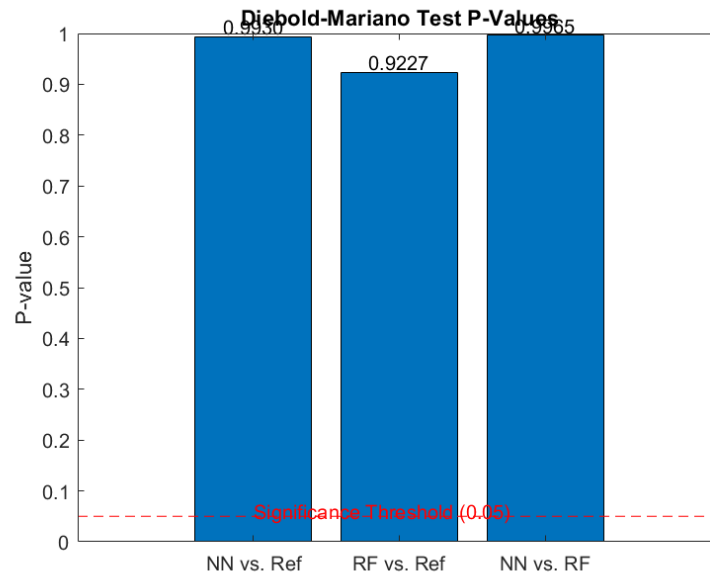


Figure 1: Comparison of p-values for Diebold-Mariano test

The Giacomini-Rossi test, also known as Fluctuation test assesses whether there is a significant time variation in the relative forecast performance of two economic models compared to a random walk model. The results for all the three models is that there is no fluctuation in the series, implying that the accuracy of the forecast does not exhibit any significant changes of variance over observations.

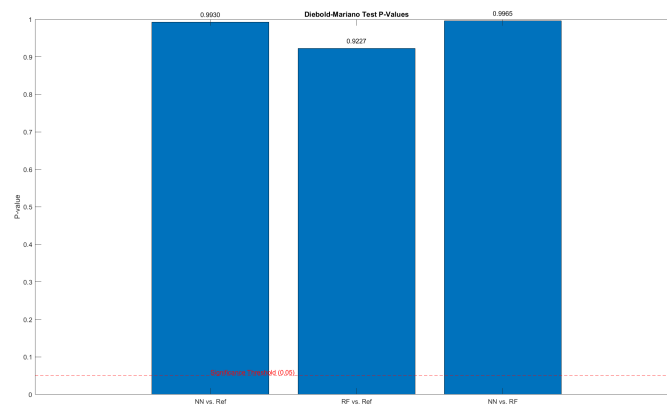


Figure 2: Comparison of p-values for Giacomini-Rossi test

In conclusion, the best model is the linear one. Even though more sophisticated models are able to perform a relatively comparable out of sample estimation, with such a small dataset it is not possible to say whether these models are more useful in predicting over such a short horizon. We want to stick with the basics and choose the simplest model as the final one.

## 5 Bibliography

### References

- [1] Ross Levine and Sara Zervos, *Stock Markets, Banks and Economic Growth*, The American Economic Review, 1998, pp. 537–558.
- [2] Ahmed Alareeni and Allam Hamdam, *ESG impact on performance of US S&P 500 - listed firms*, 2020.
- [3] Ronald C. Anderson and David M. Reeb, *Board Composition: Balancing Family Influence in S&P 500 Firms*, 2004.
- [4] *Corporate News Disclosure and Competitive Advantage: What Factors Influence S&P 500 Companies? Competitive Advantage During 2022 Economic Crisis*, 2022.
- [5] Benjamin Graham, *The Intelligent Investor*, Harper & Brothers, 1949.
- [6] Robert J. Shiller, *Irrational Exuberance*, Princeton University Press, 2000.
- [7] Najeb M.H. Masoud, *The Impact of Stock Market Performance upon Economic Growth*, International Journal of Economics and Financial Issues, 2013, pp. 788–798.
- [8] *Ensemble learning*, [https://en.wikipedia.org/wiki/Ensemble\\_learning](https://en.wikipedia.org/wiki/Ensemble_learning), Accessed: 2024-05-04.
- [9] *Bootstrap aggregating*, [https://en.wikipedia.org/wiki/Bootstrap\\_aggregating](https://en.wikipedia.org/wiki/Bootstrap_aggregating), Accessed: 2024-05-04.
- [10] *Random forest*, [https://en.wikipedia.org/wiki/Random\\_forest](https://en.wikipedia.org/wiki/Random_forest), Accessed: 2024-05-04.
- [11] *Tin Kam Ho*, [https://en.wikipedia.org/wiki/Tin\\_Kam\\_Ho](https://en.wikipedia.org/wiki/Tin_Kam_Ho), Accessed: 2024-05-04.
- [12] *Random forest*, <https://www.ibm.com/topics/random-forest>, Accessed: 2024-05-04.
- [13] *Out-of-bag error*, [https://en.wikipedia.org/wiki/Out-of-bag\\_error](https://en.wikipedia.org/wiki/Out-of-bag_error), Accessed: 2024-05-04.

## 6 Figures

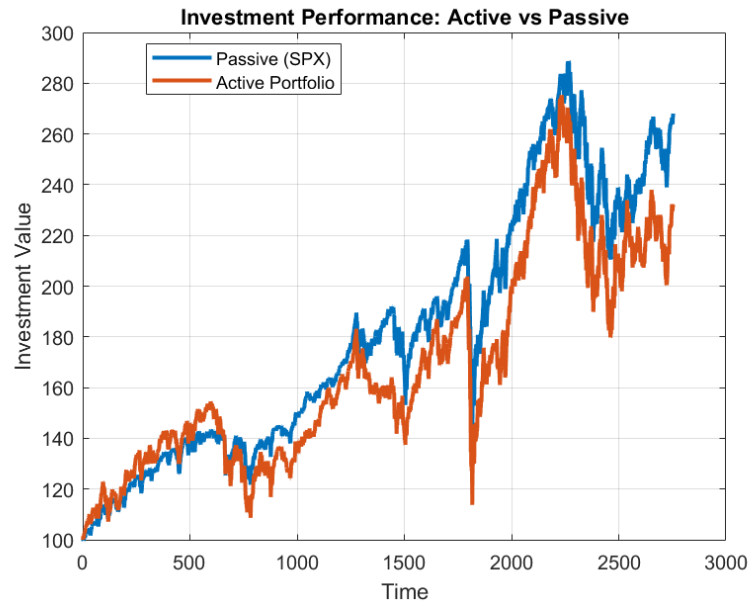


Figure 3: Active vs. Passive investment returns over time

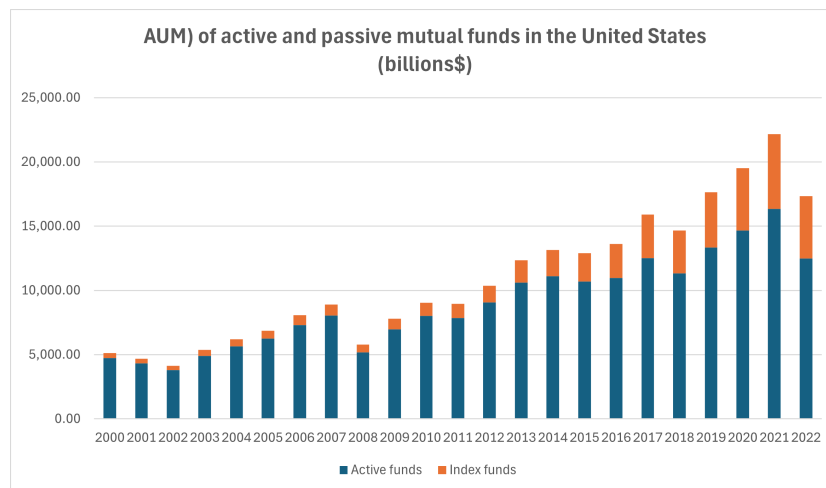


Figure 4: Proportion of Active vs. Passive Investment Portfolios

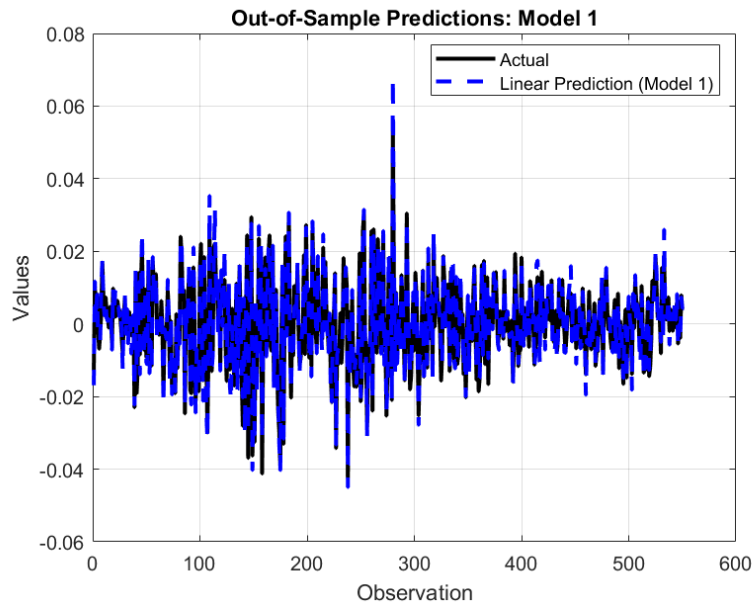


Figure 5: Out-of-Sample Predictions: Model 1

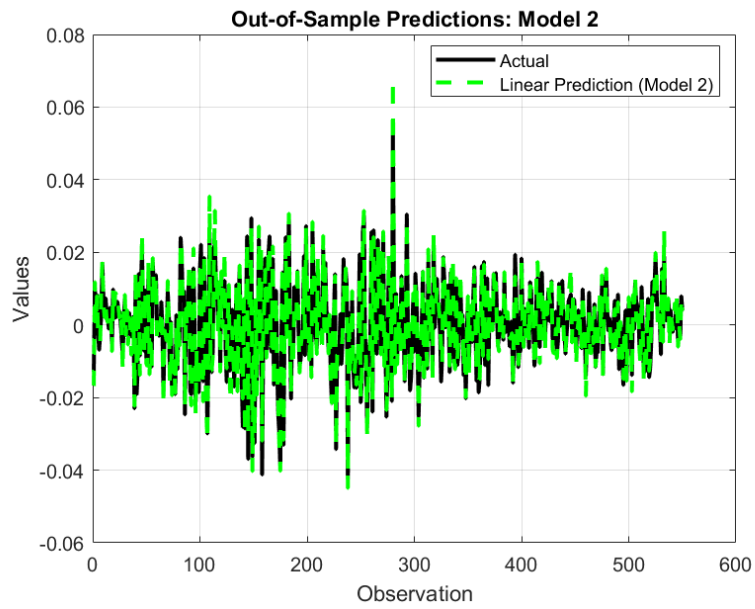


Figure 6: Out-of-Sample Predictions: Model 2

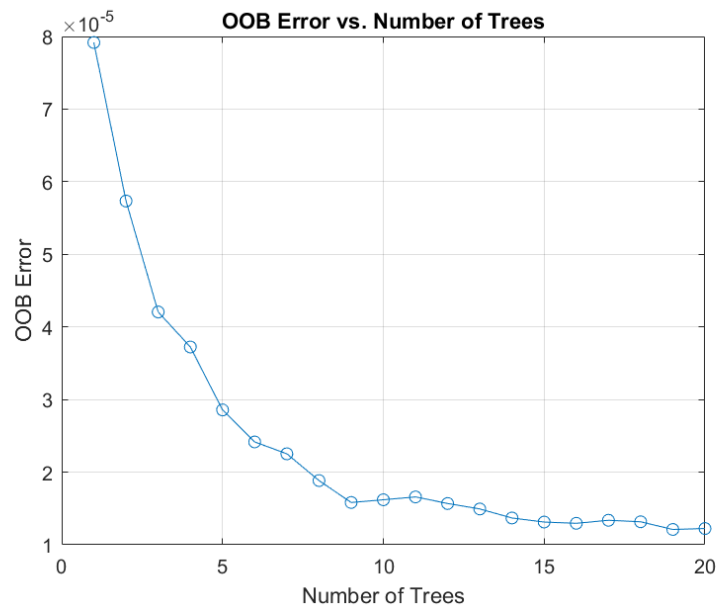


Figure 7: OOB Error vs. Number of Trees

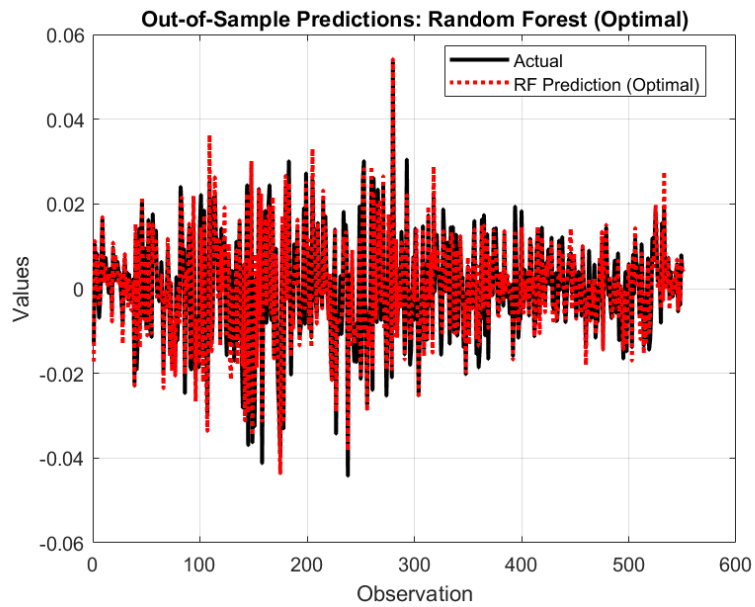


Figure 8: Out-of-Sample Predictions: Random Forest (Optimal)

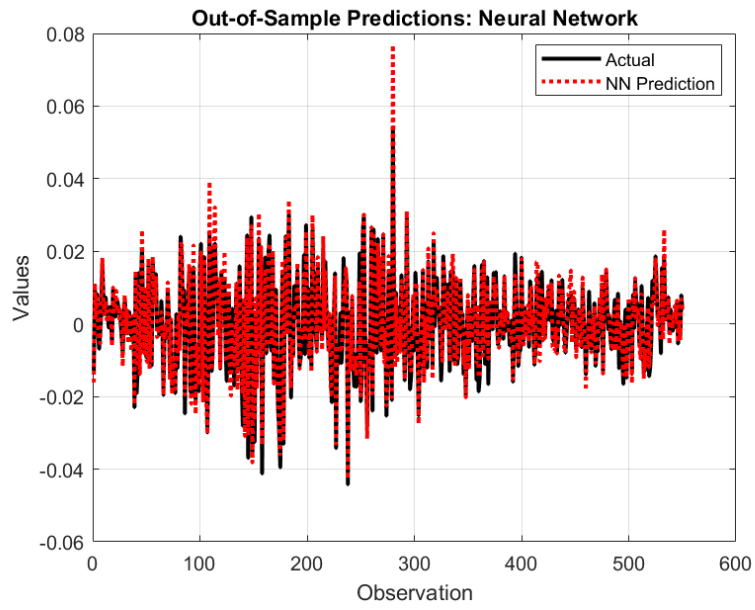


Figure 9: Out-of-Sample Predictions: Feedforward Network Model

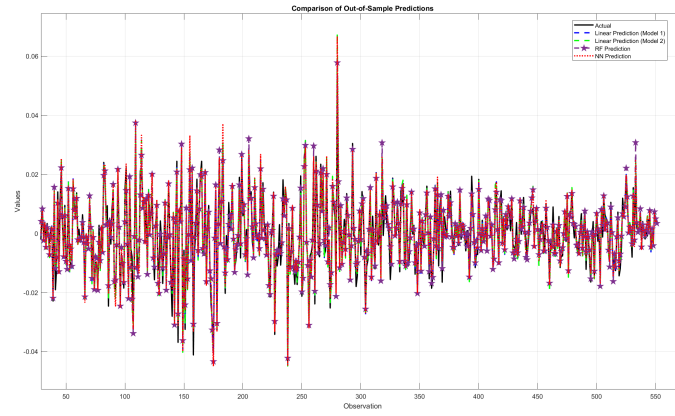


Figure 10: Out-of-Sample Predictions: a comparison

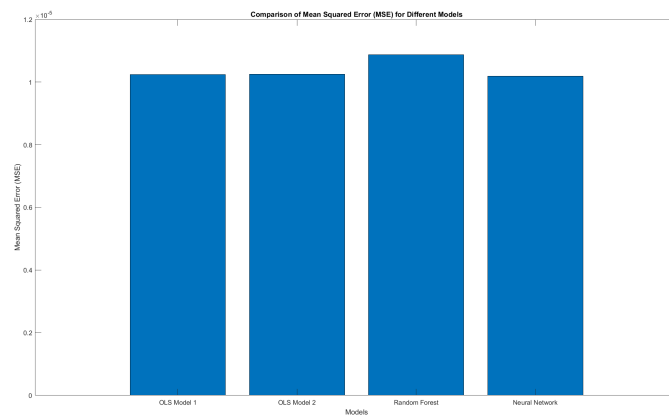


Figure 11: Comparison of MSE for Different Models