

---

# Evaluation of Machine Learning in Empirical Asset Pricing

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        Several recent studies have claimed that machine learning methods provide superior  
2        predictive accuracy of asset returns, relative to simpler modelling approaches, and  
3        can correctly identify factors needed to price portfolio risk. Herein, we demonstrate  
4        that this performance is critically dependent on several features of the data being  
5        analysed; including, the training/test sample split, the frequency at which the data  
6        is observed, and the chosen loss-function. In contrast to existing studies, which  
7        claim that neural nets provide superior predictive accuracy, through a series of  
8        realistic examples that mimics the stylized facts of asset returns, we demonstrate  
9        that neural methods are easily outperformed by simpler methods, such as random  
10       forest and elastic nets.

## 11    1   Introduction

12    The dominance of machine learning methods in terms of predictive accuracy has now begun to filter  
13    into the application and assessment of asset pricing. The most common application of machine  
14    learning methods within finance are for portfolio construction, asset price prediction, and factor  
15    selection.

16    Several studies have now used machine learning techniques to analyze the cross-section of asset  
17    returns and produce portfolios that can capture nonlinear information in the cross-section of asset  
18    returns. Mortiz and Zimmermann (2016) use tree-based methods in an attempt to understand which  
19    firm-level characteristics best predict the cross-section of stock returns, where this information can  
20    then be used within portfolio sorting to help mitigate risk. Similarly, Messemer (2017) uses deep  
21    feedforward neural nets (DFNs) to construct portfolios and predict the returns across a cross-sections  
22    of US asset returns. While Messemer (2017) demonstrates that such DFNs can better capture  
23    nonlinear information, and outperform portfolios generated from linear benchmarks, the author does  
24    not claim that deep learning methods are the best methods to exploit these nonlinear interactions.

25    In addition, several studies have now suggested that machine learning methods can produce better  
26    predictions of asset returns ([?], [?] and [?]). In particular, the results of Gu et al. (2019) suggest that,  
27    in terms of predictive performance, as measured by an out-of-sample  $R^2$ , tree-based methods and  
28    shallow neural nets can provide superior predictive accuracy over other machine learning methods  
29    and simpler model-based approaches. This finding is born out both in terms of simulated data, and  
30    an empirical example with monthly returns data from 1957 to 2016. [?] attribute this to machine  
31    learning’s ability to evaluate and consider non-linear complexities among factors that cannot be  
32    feasibly achieved using traditional techniques.

33    Similarly, work by Kozak et al, (2018), Freyberger et al. (2018), Feng et al., (2019) and Rapach  
34    and Zhou (2013), demonstrate that machine learning methods can “systematically evaluate the  
35    contribution to asset pricing of any new factor” used within an existing linear asset pricing structure.

36 In addition, Gu et al. (2019) use variable importance metrics to quantify the differential impact of  
37 factors across a large set of possible factors available for asset pricing. As such, machine learning  
38 methods can be used, *en masse*, to consistently evaluate the ability of various factors to help price  
39 portfolio risk. Such work is particularly useful given the literature's seeming obsession with the XXX  
40 and constructing such factors: as of 2014, quantitative trading firms were using 81 factor models (Hsu  
41 and Kalesnik, 2014), while Harvey and Liu (2019) currently document that well over 600 different  
42 factors have been suggested in the literature.

43 While the above studies all demonstrate the potential benefits of machine learning methods within  
44 empirical finance, it is unclear whether the findings in these papers are easily generalizable to: one,  
45 different training and validation periods; two, different sampling frequencies, which result in stock  
46 returns with significant different characteristics (e.g., daily volatility is significantly higher than  
47 monthly volatility); and three, different loss-measures of predictive accuracy. The answer to such  
48 questions are particularly pertinent given that the machine learning literature has already documented  
49 the difficulties of certain methods, including those references above, in dealing with data that displays  
50 the stylized facts of asset returns. For instance, methods such as penalized regression and tree-based  
51 models assume a form of conditional independence between observations, which is violated by the  
52 state dependence that exists within, and across, asset returns. In addition, it has already been noted  
53 that training more standard types of neural networks, such as the feed forward kind considered in Gu  
54 et al, becomes particularly difficult when data displays strong dependence, ([?]). In addition, more  
55 complex machine learning approaches require extremely large amounts of data, as well as specialized  
56 sample splitting and cross-validation schemes, to deal with possible model over-fitting.

57 In some ways, existing applications of machine learning to empirical asset pricing have either over-  
58 looked, downplayed, or simply ignored the importance of the above issues. For example, Messemmer  
59 (2017) and [?] use cross validation as part of their model building procedures, thereby destroying  
60 the temporal ordering of data. In addition, [?] and Messemmer (2017) produce models using training  
61 samples that end much earlier than the data sets which they ultimately produce forecasts for: in the  
62 case of Messemmer (1970), the training period ends in 1981, while the which ends in the 1970s to  
63 ultimately produce forecasts for the most recent 30 years; in the case of [?], the training ends in the  
64 1970s, with predictions ultimately produced only for the period of returns from 1987-2016. This is  
65 particularly worrying as the factors driving daily or monthly returns in the 1980s, are starkly different  
66 than those driving returns in, say, 2001 onwards. However, both of these papers suggest that the  
67 training and validation sets used for the various methods does not impact the test set results.

68 While some combination of machine learning methods can undoubtedly lead to better performance  
69 than simpler model-based solutions, a more systematic treatment on the ability of these methods to 1)  
70 accurately detect significant factors; and 2) accurately predict returns according to a range of loss  
71 measures, must be formulated before researchers can rely on such methods in practice. The goal of  
72 this paper is to bridge this gap and thereby provide a systematic, rigorous, realistic, and reproducible  
73 study on the performance of several machine learning methods that have been used in empirical asset  
74 pricing.

75 First, through a rigorous simulation study, which captures the stylized facts of asset returns, we give  
76 an in-depth comparison of several machine learning methods used in the literature. The simulation  
77 study explicitly explores how different aspects of financial data such as persistence in regressors, cross  
78 sectional correlation and different complexities of data generating process can affect a method's ability  
79 to: 1) accurately predict future returns across a range of loss measures; and 2) correctly identify the  
80 significant factors driving returns. In contrast to existing findings, in this realistic simulation design,  
81 we find that neural network procedures, such as feedforward nets, LSTM (CITE), and DeepAR  
82 models (CITE), are among the worst performing methods, while tree-based methods and elastic net  
83 are among the best performing methods. We also demonstrate that this result is consistent across  
84 various levels of volatility, cross-sectional correlation, return signal, and different loss functions. In  
85 addition, we demonstrate that elastic net and tree-based methods also outperform neural net based  
86 approach in terms of correctly identifying significant factors.

87 Next, we validate these findings using a empirical data set of asset returns that considers quarterly  
88 individual price data from CRSP for all firms listed in the NYSE, AMEX and NASDAQ. The starting  
89 period of the data is January first 1957 (starting date of the S&P 500) and the ending date is December  
90 2016, totalling 60 years. A set of 549 possible factors are used to explain the cross-section of returns.  
91 We pay careful attention to the training and test split, and only use the last fourteen years of quarterly

returns to evaluate the different machine learning methods. The results found in the empirical study agree completely with those in the aforementioned simulation study: across all machine learning methods, neural net based procedure perform the worst across various loss functions, while tree-based methods and elastic net perform the best.

The results of this study suggest that great care and diligence is required if one wishes to implement machine learning methods within empirical finance. Indeed, our results suggest that the efficacy of machine learning methods within empirical finance depends are highly-dependent on the samples used for training and testing, the loss functions used for evaluation, and the specific nature of the data series one wishes to predict. As such, while potentially quite useful in empirical finance, machine learning methods are not necessarily a panacea to correctly predict future asset prices or to correctly disentangle which factors are relevant.

The remainder of the paper is organized as follows....

## 2 Model and Methods

### 2.1 Statistical Model

In this section we briefly discuss the statistical model considered for asset returns. Excess monthly returns on asset  $i$ ,  $i = 1, \dots, n$ , at time  $t$ ,  $t = 1, \dots, T$ , are assumed to evolve in an additive fashion:

$$r_{i,t+1} = E(r_{i,t+1}|\mathcal{F}_t) + \epsilon_{i,t+1}, \quad E(\epsilon_{i,t+1}|\mathcal{F}_t) = 0 \quad (1)$$

where  $\mathcal{F}_t$  denotes the observable information at time  $t$ , and  $\epsilon_{i,t+1}$  is a martingale difference sequence (hereafter, mds). We further consider that the conditional mean of returns is an unknown function of a  $P$ -dimensional vector of features, assumed measurable at time  $t$ , such that

$$E(r_{i,t+1}|\mathcal{F}_t) = g(z_{i,t}) \quad (2)$$

The features, or predictors,  $z_{i,t}$  are assumed to be composed of time- $t$  information, and depends only the characteristics of stock  $i$ . It is not assumed that all  $z_{i,t}$  are present within the function  $g(\cdot)$  across all  $i$  units. That is, the function  $g(\cdot)$  need not depend on the same  $z_{i,t}$  as  $i$  varies. The assumption that the information set can be characterized by the variables  $z_{i,t}$  without dependence on the  $j \neq i$  return units, is reasonable given that the collection of  $z_{i,t}$  is rich enough.

In what follows, we represent the space of possible features as the Kronecker product of two pieces

$$z_{i,t} = x_t \otimes c_{i,t} \quad (3)$$

where the variables  $c_{i,t}$  represent a  $P_c \times 1$  vector of individual-level characteristics for return  $i$ , and  $x_t$  represents a  $P_x \times 1$  vector of macroeconomic predictors, and  $\otimes$  represents the Kronecker product. Thus, for  $P = P_c \cdot P_x$ ,  $z_{i,t}$  represents a  $P \times 1$  feature space that can be used to approximate the unknown function  $g(\cdot)$ .

### 2.2 Methods

Given features  $z_{i,t}$ , the goal of any machine learning method is to approximate the unknown function  $g(\cdot)$  in 1. Broadly speaking, how different ML methods choose to approximate this function depends on three components:

1. the model used to make predictions,<sup>1</sup>
2. the regularization mechanism employed to mitigate over-fitting;
3. a loss function that penalized poor predictions.

To ensure the results of ML different methods will be comparable, we fix both the regularization mechanisms and loss functions used within each method, and allow only the models used for prediction to vary. This approach seeks to ensure that performances in one method, relative to another, are based on the model structure and not to some feature of how the models were fit. To this end, we first discuss points 2. and 3. above, and then briefly present the models used for our comparison.

<sup>1</sup>The model used by the ML method need not correspond to the statical models assumed to describe the data. Herein, our goal will not be to asses the ‘‘accuracy’’ of the statistical model, but to determine how different ML methods accurately determine the salient features of this model.

133 **Loss functions:** The choice of loss function used to fit the ML methods is instrumental in the  
 134 methods' ultimate performance. Herein, we consider two separate loss functions: Mean Absolute  
 135 Error (MAE) and Mean Squared Error (MSE):

$$\text{MAE} = \frac{1}{n} \sum_{j=i}^n |y_j - \hat{y}_j| \text{ and } \text{MSE} = \frac{1}{n} \sum_{j=i}^n (y_j - \hat{y}_j)^2,$$

136 We consider both loss functions since MAE is less sensitive to outliers in the data which financial  
 137 returns are known to exhibit, and which are caused by extreme market movements. Given this, we  
 138 expect MAE to produce predictive results that are more robust to such outlier events.

139 **Mitigating over-fitting:** ML methods guard against over-fitting by emphasizing out-of-sample  
 140 performance. To this end, observed data is split into "training", "validation" and "test" sets. Since  
 141 returns data is intrinsically dependent, when constructing such a split we must consider a schema that  
 142 respects this dependence structure.

143 Throughout our experiments/applications, to balance computation and accuracy, we use a hybrid  
 144 "rolling window" and "recursive" approach to training/validation/test splits: for each model refit, the  
 145 training set is increased by one year observations, i.e., 12 monthly observations; the validation set is  
 146 fixed at one year and moves forward (by one year) with each model refit; predictions are generated  
 147 using that model for the subsequent year.

148 **Models** The remaining specification for the ML methods is the chosen model used to generate  
 149 predictions. Herein, we consider a host of different models: including elastic net (Hastie et al.,  
 150 XXX), Random forest (XXX), feed-forward neural nets (XXX), LSTM (XXX), FFORMA (XXX)  
 151 and DeepAR models (XXX). To keep the details as brief as possible, we give full details on each  
 152 model and certain features of its implementation used in this work in the appendix. For each of the  
 153 different methods, we consider two variants, one based on the MAE loss and one based on the MSE  
 154 loss.

## 155 2.3 Model evaluation measures

156 **Predictive accuracy** Predictive performance for individual excess returns are assessed using Mean  
 157 Absolute Error (MAE), Mean Squared Error (MSE) (evaluated over the test set) and an out-of-sample  
 158  $R^2$  measure. While out-of-sample  $R^2$  is a common measure, there is no universally agreed-upon  
 159 definition. As such, we explicitly state the version employed herein as

$$R_{OOS}^2 = 1 - \frac{\sum_{(i,t) \in \mathcal{T}_3} (r_{i,t+1} - \hat{r}_{i,t+1})^2}{\sum_{(i,t) \in \mathcal{T}_3} (r_{i,t+1} - \bar{r}_{i,t+1})^2} \quad (4)$$

160 where  $\mathcal{T}_3$  indicates that the fits are only assessed on the test sub-sample, which is never used for  
 161 training or tuning.

162 Since  $R^2$  is based on in-sample-fit of a linear model, this measure is less meaningful for most of the  
 163 ML methods considered in in this paper. However, we report this measure since this measure has also  
 164 been considered in other applications of ML to empirical finance (see, e.g., Gu et al., 2019).

165 **Factor Selection** An important aspect of empirical finance is the understanding of which features  
 166 drive risk. That is, which features are explicitly represented within  $z_{i,t}$  and can thus be used to help  
 167 price risk using equation 1. To this end, we define a simple variable importance (VI) measure to be  
 168 applied across all ML methods in this research. To this end, we mirror the measure produced in [?]  
 169 and define  $VI_j$  as the reduction in predictive  $R^2$  from setting all values of predictor  $j$  to 0, while  
 170 holding the remaining model estimates fixed. Each  $VI_j$  is then normalized to sum to 1.

171 However, as  $VI_j$  can sometimes be negative, we shift  $VI_j$  by the smallest  $VI_j$  plus a small constant,  
 172 then dividing by this sum to alleviate numerical issues<sup>2</sup>. The resulting VI measure is then.

$$VI_{j,norm} = \frac{VI_j + \min(VI_j) + o}{\sum VI_j + \min(VI_j) + o} \quad ; \quad o = 10^{-100} \quad (5)$$

<sup>2</sup>This mechanism was chosen because the other popular normalization mechanism "softmax" was observed to be unable to preserve the distances between each original  $VI_j$ , making discernment between each  $VI_j$  difficult.

### 3 Simulation study

We begin with the simulation study as a way to explore how machine learning performs with regards to the stylized facts of empirical returns in a controlled environment. We simulate according to a design which incorporates low signal to noise ratio, stochastic volatility in errors, persistence and cross sectional correlation in regressors. Our specification is a latent factor model for excess returns  $r_{t+1}$ , for  $t = 1, \dots, T$ :

$$r_{i,t+1} = g(z_{i,t}) + \beta_{i,t+1}v_{t+1} + e_{i,t+1}; \quad z_{i,t} = (1, x_t)' \otimes c_{i,t}, \quad \beta_{i,t} = (c_{i1,t}, c_{i2,t}, c_{i3,t}) \quad (6)$$

$$e_{i,t+1} = \sigma_{i,t+1}\varepsilon_{i,t+1}; \quad (7)$$

$$\log(\sigma_{i,t+1}^2) = \omega + \gamma \log(\sigma_t^2) + \sigma_u u; \quad u \sim N(0, 1) \quad (8)$$

where  $v_{t+1}$  is a  $3 \times 1$  vector of errors,  $w_{t+1} \sim N(0, 1)$ ,  $\varepsilon_{i,t+1} \sim N(0, 1)$  scalar error terms, matrix  $C_t$  is an  $N \times P_c$  matrix of latent factors, where the first three columns correspond to  $\beta_{i,t}$ , across the  $1 \leq i \leq N$  dimensions, while the remaining  $P_c - 3$  factors do not enter the return equation. The  $P_x \times 1$  vector  $x_t$  is a  $3 \times 1$  multivariate time series, and  $\varepsilon_{t+1}$  is a  $N \times 1$  vector of idiosyncratic errors. The parameters of these were tuned such that the annualized volatility of each return series was approximately 22%, as is often observed empirically.

**Simulating characteristics** We build in correlation across time among factors by drawing normal random numbers for each  $1 \leq i \leq N$  and  $1 \leq j \leq P_c$ , according to :

$$\bar{c}_{ij,t} = \rho_j \bar{c}_{ij,t-1} + \epsilon_{ij,t}; \quad \rho_j \sim \mathcal{U}(0.5, 1) \quad (9)$$

We then build in cross sectional correlation:

$$\hat{C}_t = L\bar{C}_t; \quad B = LL' \quad (10)$$

$$B := \Lambda\Lambda' + \frac{1}{10}\mathbb{I}_n, \quad \Lambda_i = (\lambda_{i1}, \dots, \lambda_{i4}), \quad \lambda_{ik} \sim N(0, \lambda_{sd}), \quad k = 1, \dots, 4 \quad (11)$$

where  $B$  serves as a variance covariance matrix with  $\lambda_{sd}$  controlling the density of the matrix, and  $L$  represents the lower triangle matrix of  $B$  using the Cholesky decomposition.  $\lambda_{sd}$  values of 0.01, 0.1 and 1 were used to explore increasing degrees of cross sectional correlation. Finally, the "observed" characteristics for each  $1 \leq i \leq N$  and for  $j = 1, \dots, P_c$  are constructed according to:

$$c_{ij,t} = \frac{2}{n+1} \text{rank}(\hat{c}_{ij,t}) - 1. \quad (12)$$

with the rank transformation normalizing all predictors to be within  $[-1, 1]$ .

**Simulating macroeconomic series** For simulation of  $x_t$ , a  $3 \times 1$  multivariate time series, we consider a Vector Autoregression (VAR) model <sup>3</sup>:

$$x_t = Ax_{t-1} + u_t; \quad A = 0.95I_3; \quad u_t \sim N(\mu = (0, 0, 0)', \Sigma = I_3)$$

**Simulating return series** We consider three different functions for  $g(z_{i,t})$ :

$$(1) \quad g_1(z_{i,t}) = (c_{i1,t}, c_{i2,t}, c_{i3,t} \times x_t'[3,]) \theta_0 \quad (13)$$

$$(2) \quad g_2(z_{i,t}) = (c_{i1,t}^2, c_{i1,t} \times c_{i2,t}, \text{sgn}(c_{i3,t} \times x_t'[3,])) \theta_0 \quad (14)$$

$$(3) \quad g_3(z_{i,t}) = (1[c_{i3,t} > 0], c_{i2,t}^3, c_{i1,t} \times c_{i2,t} \times 1[c_{i3,t} > 0], \text{logit}(c_{i3,t})) \theta_0 \quad (15)$$

where  $x_t'[3,]$  denotes the third element of the  $x_t'$  vector.

$g_1(z_{i,t})$  allows the characteristics to enter the return equation linearly, and  $g_2(z_{i,t})$  and  $g_3(z_{i,t})$  allow the characteristics to enter the return equation interactively and non-linearly. <sup>4</sup>  $\theta^0$  was tuned such that the predictive  $R^2$  was approximately 5%.

The simulation design results in  $3 \times 3 = 9$  different simulated datasets, each with  $N = 200$  stocks,  $T = 180$  periods and  $P_c = 100$  characteristics. Each design was simulated 10 times to assess the robustness of machine learning algorithms. The number of simulations was kept low for computational feasibility. We employ the hybrid data splitting approach with a training:validation length ratio of approximately 1.5 and a test set that is 1 year in length.

<sup>3</sup>More complex specifications of  $A$  were briefly explored, but these did not have a significant impact on results.

<sup>4</sup>( $g_1, g_2$  correspond to the simulation design used by [?].)

### 3.1 Simulation Study Results

**Prediction Performance** We observe that in general elastic nets are the best performing model, followed closely by random forests, then neural networks. All machine learning models were unaffected by cross sectional correlation in terms of prediction performance, and had better performance when fitted with respect to quantile loss. Random forest models only outperformed the elastic nets on highly non-linear specifications. The neural network models were not observed to outperform any of the machine learning models.

This is in stark contrast to the linear models, whose prediction performance is severely affected by both non-linearities, and increasing cross sectional correlation. This result is consistent across all loss metrics, and is most obvious when looking at the out-of-sample R-squared metrics.

Machine learning models fitted with respect to minimizing MAE (quantile loss) generally perform better, even when evaluated against MSE loss metrics. This is not a surprising result, especially considering the stochastic error design which introduces significant shocks to the returns process. Though the actual difference between the loss metrics between the penalized linear models, random forests and neural networks are very small, when considering the consistency of the results across numerous Monte Carlo simulations, the differences in prediction performance, though small, is robust and significant.

We note that most of these results contradict the sparse literature.

Table 1: Top Models by MAE in Simulation Study

| Corr | model   | Test MAE  |           |           | Test MSE  |           |           |
|------|---------|-----------|-----------|-----------|-----------|-----------|-----------|
|      |         | g1        | g2        | g3        | g1        | g2        | g3        |
| 0.01 | ELN.MAE | 0.0345786 | 0.0361950 | 0.0353345 | 0.0025652 | 0.0026882 | 0.0026210 |
|      | RF.MAE  | 0.0354594 | 0.0354204 | 0.0355399 | 0.0026434 | 0.0026305 | 0.0026446 |
|      | NN2.MAE | 0.0359604 | 0.0369206 | 0.0363047 | 0.0026786 | 0.0027474 | 0.0026996 |
|      | NN1.MAE | 0.0358939 | 0.0368335 | 0.0363352 | 0.0026718 | 0.0027396 | 0.0027028 |
|      | NN3.MAE | 0.0358164 | 0.0369345 | 0.0364712 | 0.0026697 | 0.0027491 | 0.0027181 |
| 1    | ELN.MSE | 0.0346142 | 0.0362761 | 0.0354437 | 0.0025676 | 0.0026980 | 0.0026300 |
|      | RF.MAE  | 0.0359158 | 0.0356434 | 0.0360529 | 0.0026747 | 0.0026445 | 0.0026786 |
|      | NN5.MAE | 0.0370087 | 0.0372705 | 0.0374132 | 0.0027744 | 0.0027832 | 0.0027916 |
|      | NN4.MSE | 0.0373820 | 0.0368966 | 0.0373542 | 0.0028051 | 0.0027505 | 0.0027970 |
|      | NN3.MAE | 0.0372849 | 0.0370382 | 0.0371925 | 0.0027940 | 0.0027652 | 0.0027753 |

**Factor Importance** We observe that the elastic net outperforms all other models consistently in terms of assigning the correct relative importance to the true underlying regressors,<sup>5</sup> even in settings with high cross sectional correlation.

Elastic net models perform the best at identifying the true data generating regressors, and this appears to be mostly robust regardless of the amount of cross sectional correlation, though their performance worsens as the data generating process becomes more non-linear. On more difficult specifications, the elastic net models are conservative and typically identify a single regressor as importance - most apparent on the  $g_2$  specification. Occasionally, the elastic nets identified the incorrect covariates, though the relative importance assigned to them was small.

The random forests and to a lesser extent the neural networks also correctly identified the correct underlying regressors, but struggled with adequately discerning relative importance among correlated regressors. This was became more apparent as the degree of cross sectional correlation increased (see decreasing relative importance of true underlying regressors in Figures ?? and ?? in Appendix).

The linear models unsurprisingly struggled with factor significance analysis with respect to both increasing cross sectional correlation non-linearities. This highlights the non-robustness and ineffectiveness of using traditional linear regression as documented by the literature; linear models

<sup>5</sup>( $c_1$ .constant,  $c_2$ .constant and  $c_3.x_3$  for  $g_1$  and  $g_2$  specifications, and  $c_1$ .constant,  $c_2$ .constant and  $c_3$ .constant for  $g_3$ )

Figure 1: Simulation g1 Variable Importance

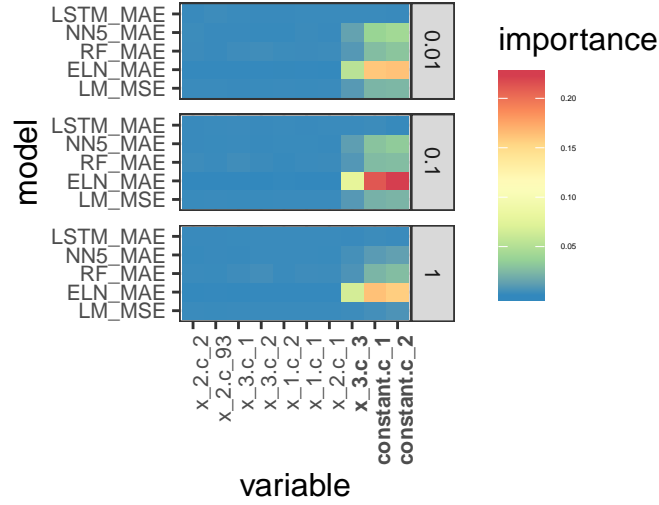
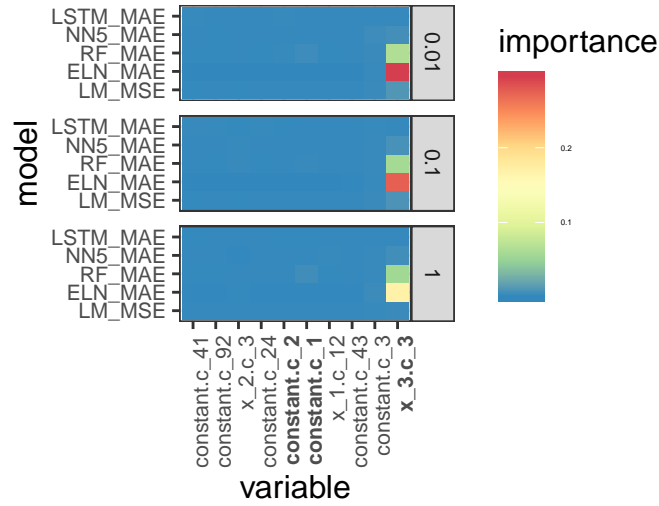


Figure 2: Simulation g2 Variable Importance



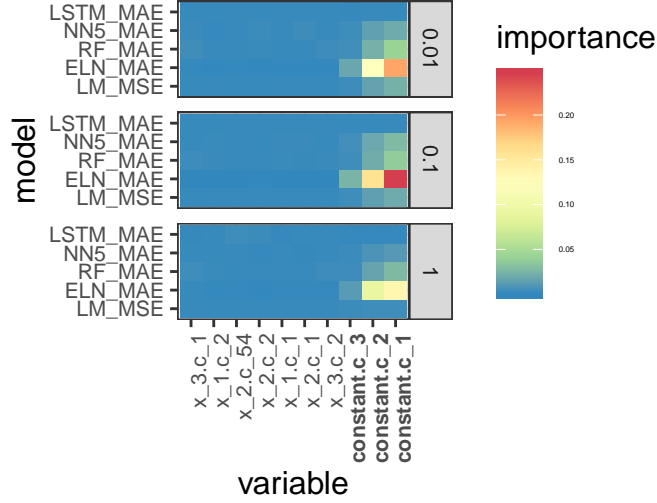
were consistently observed to identify irrelevant regressors as important, especially as the degree of cross sectional correlation increased. Considering that the graphs represent the averaged variable importance metrics over different simulation realisations, this means that on a single simulation realization, the performance of linear models is significantly worse.

#### 4 Empirical analysis

We conduct an empirical study as a final way to corroborate the findings of the properties of machine learning models which we observed in the simulation study. Though our simulation study was aimed at capturing the main features of observed data, the underlying data generating process for empirical returns is unknown. This study thus acts as a robustness check as to how machine learning performs on real world data, which can be significantly more complex and noisy than simulated contexts.

Importantly, we find that our findings from the simulation study are largely corroborated for empirical returns data.

Figure 3: Simulation g3 Variable Importance



#### 4.1 Data

We begin by obtaining monthly individual price data from CRSP for all firms listed in the NYSE, AMEX and NASDAQ, starting from 1957 (starting date of the S&P 500) and ending in December 2016, totalling 60 years. To build individual factors, we construct a factor set based on the cross section of returns literature. This data was sourced from and is the same data used in [?]. Like our initial returns sample, it begins in March 1957 and ends in December 2016, totalling 60 years. It contains 94 stock level characteristics: 61 updated annually, 13 updated quarterly and 20 updated monthly<sup>6</sup>.

We detail our cleaning procedure of this dataset. To reduce the size of the dataset and increase feasibility, we only consider non-penny equities traded primarily on the NASDAQ. To achieve a balance between having a dataset with enough data points and variability among factors, the dataset was converted to a quarterly format. Quarterly returns were then constructed using the PRC variable according to actual returns (ie not logged differences):

$$RET_t = \frac{PRC_t - PRC_{t-1}}{PRC_{t-1}} \quad (16)$$

We allow all stocks which have a quarterly return to enter the dataset, even if they disappear from the dataset for certain periods, as opposed to only keeping stocks which appear continuously throughout the entire period. This was primarily done to reduce survivorship bias in the dataset, and also allows for stocks which were unlisted and relisted again to feature in the dataset.<sup>7</sup>

We then follow [?] and construct eight macroeconomic factors following the variable definitions in [?] (see Table 5). These factors were lagged by one period so as to be used to predict one period ahead quarterly returns. The treasury bill rate was also used from this source to proxy for the risk free rate in order to construct excess quarterly returns.

The two sets of factors were then combined to form a baseline set of covariates, which we define throughout all methods and analysis as:

$$z_{i,t} = (1, x_t)' \otimes c_{i,t} \quad (17)$$

<sup>6</sup>The dataset also included 74 Standard Industrial Classification (SIC) codes, but these were omitted due to their inconsistency, and inadequateness at classifying companies, as noted by WRDS

<sup>7</sup>To deal with missing data, any characteristics that had over 20% of their data missing were omitted. Remaining missing data were then imputed using their cross sectional medians for each year. See Appendix for more details.



Table 2: Empirical Study Loss Statistics

| model   | Sample 1        |                 |                 | Sample 2        |                 |                 | Sample 3       |                 |                 |
|---------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|----------------|-----------------|-----------------|
|         | Test MAE        | Test MSE        | Test $R^2$      | Test MAE        | Test MSE        | Test $R^2$      | Test MAE       | Test MSE        | Test $R^2$      |
| LM.MSE  | 0.229034        | 0.116015        | -1.808481       | 0.397573        | 0.312653        | -6.329935       | 0.566307       | 0.83804         | -17.522476      |
| LM.MAE  | 0.273452        | 0.15894         | -2.8476         | 0.555673        | 0.742223        | -16.400898      | 0.651614       | 1.225121        | -26.077774      |
| ELN.MSE | 0.133887        | 0.039947        | 0.032956        | 0.140402        | 0.04277         | -0.002712       | <b>0.14433</b> | <b>0.043761</b> | <b>0.032789</b> |
| ELN.MAE | 0.131369        | 0.040718        | 0.014306        | <b>0.137092</b> | <b>0.041892</b> | <b>0.017875</b> | 0.146251       | 0.045207        | 0.000835        |
| RF.MSE  | 0.130366        | <b>0.036629</b> | <b>0.113289</b> | 0.195817        | 0.070642        | -0.656158       | 0.157934       | 0.05122         | -0.132066       |
| RF.MAE  | <b>0.126703</b> | 0.036785        | 0.109505        | 0.173721        | 0.057546        | -0.349132       | 0.14692        | 0.046037        | -0.01752        |
| NN1.MSE | 0.169127        | 0.057044        | -0.380909       | 0.207662        | 0.074751        | -0.752494       | 0.192125       | 0.069738        | -0.541369       |
| NN1.MAE | 0.157324        | 0.050418        | -0.22052        | 0.191762        | 0.066746        | -0.564818       | 0.18547        | 0.063053        | -0.393606       |
| NN2.MSE | 0.168773        | 0.059436        | -0.43883        | 0.181808        | 0.063232        | -0.482433       | 0.180584       | 0.062745        | -0.386797       |
| NN2.MAE | 0.162667        | 0.055447        | -0.342256       | 0.194277        | 0.069386        | -0.626702       | 0.185173       | 0.065186        | -0.440746       |
| NN3.MSE | 0.154784        | 0.050152        | -0.21408        | 0.180103        | 0.060193        | -0.411175       | 0.177604       | 0.060404        | -0.335065       |
| NN3.MAE | 0.146411        | 0.044901        | -0.086967       | 0.18499         | 0.06461         | -0.514744       | 0.184986       | 0.063861        | -0.411475       |
| NN4.MSE | 0.153802        | 0.048641        | -0.177503       | 0.193066        | 0.067515        | -0.582833       | 0.172707       | 0.057774        | -0.276929       |
| NN4.MAE | 0.157301        | 0.050286        | -0.217308       | 0.168815        | 0.055711        | -0.306102       | 0.167998       | 0.055129        | -0.218463       |
| NN5.MSE | 0.149436        | 0.047279        | -0.14452        | 0.183584        | 0.064137        | -0.503653       | 0.170238       | 0.056992        | -0.259652       |
| NN5.MAE | 0.140781        | 0.042832        | -0.036882       | 0.181096        | 0.06216         | -0.4573         | 0.164896       | 0.053458        | -0.181528       |

where  $c_{i,t}$  is a  $P_c$  matrix of characteristics for each stock  $i$ , and  $(1, x_t)'$  is a  $P_x \times 1$  vector of macroeconomic predictors,  $\otimes$  represents the Kronecker product.  $z_{i,t}$  is therefore a  $P_x P_c$  vector of features for predicting individual stock returns and includes interactions between stock level characteristics and macroeconomic variables. The total number of covariates in this baseline set is  $61 \times (8 + 1) = 549^8$ . The final dataset spanned from 1993 Q3 to 2016 Q4 with 202, 066 individual observations.<sup>9</sup>

We mimic the procedure used in the simulation study. For the sample splitting procedure, the dataset was split such that the training and validation sets were split such that the training set was approximately 1.5 times the length of the validation set, in order to predict a test set that is one year in length.

## 4.2 Empirical Data Results

In general, the empirical results are in remarkable agreement with the those obtained in the simulation study: the penalized linear models general perform the best, with the random forest models offering slightly worse performance. Machine learning models fitted with respect to median quantile loss were similarly observed to typically offer improvements across all machine learning models across all loss metrics.

**Prediction Accuracy** In terms of prediction accuracy, we can see that in general the results of the simulation study were repeated: the elastic net models perform the best, followed by the random forests, then the neural networks, and finally the linear models. We note that the differences between each model using the MSE and MAE loss metrics are much more pronounced on empirical data. Even so, the predictive performance between the elastic net models and the quantile random forests is not particularly large, and we observe the quantile random forests outperforming the elastic nets in the first data sample. We similarly see that machine learning models perform better when fitted with respect to quantile loss instead of MSE. Most notably, we start to see the neural network models performing poorly on the empirical data, a direct contradiction to what has been reported in the literature.

The non-robustness of DFNs is amplified on the empirical dataset. This was observed to be somewhat more common on neural networks fitted with respect to MSE, suggesting that they are indeed very sensitive to outliers in training data. We similarly observe some evidence that deeper neural networks

<sup>8</sup>As the individual and macroeconomic factors can have similar names, individual and macroeconomic factors were prefixed with ind\_ and macro\_ respectively.

<sup>9</sup>The dataset was not normalized for all methods, as only penalized regression and neural networks are sensitive to normalization. For these two methods, the dataset was normalized such that each predictor column had 0 mean and 1 variance.

303 perform better, though this result is less apparent due to the lower robustness on empirical data (see  
304 ?? in Appendix for results).

305 **Factor Importance** As the data generating process for empirical returns is unknown, the variable  
306 importance results cannot be directly compared with those of the simulation study. Even so, we  
307 see similar results: the elastic net and random forest models tend to agree on the same subset of  
308 predictors, but the random forest struggles to discern between highly correlated regressors. Similar to  
309 the prediction performance results, neural networks perform poorly.

Figure 4: Empirical Individual Factor Importance, averaged across all training samples

Figure 5: Empirical Macroeconomic Factor Importance, faceted by training sample

310 The two top performing models of elastic net and random forest consistently pick out the 1 month and  
311 6 month momentum factors out of the individual characteristics as important, and the book-to-market  
312 and default yield spread factors out of the macroeconomic factors are important. In general, the  
313 variable importance metrics are less consistent for the random forests, and it should be noted in  
314 particular that the random forest tends to determine factors highly correlated with momentum as  
315 important, such as change in momentum, dollar trading volume and return volatility. Within the  
316 macroeconomic factors, penalized linear models tend to identify the average book to market ratio and  
317 the default spread as the most important. The random forests were inconsistent with the elastic nets,  
318 and tended to assign very similar variable importance metrics to most macroeconomic factors.

319 The neural networks tended to believe that the market value factor was the most important among  
320 the individual factors, a result not repeated by any of the other models considered. Within the  
321 macroeconomic factors, the neural networks identified the dividend-price ratio and earnings-price  
322 ratio as the most important among the macroeconomic factors, though these results were non-robust.

323 Interestingly, the linear models assign the controversial dividend price ratio macroeconomic factor  
324 as highly important, a result mirrored only with the neural networks. Their variable importance  
325 for individual factors across different training samples is non-robust, with the important variables  
326 almost completely changing year to year. The linear models consistently identified the controversial  
327 dividend-price ratio as important, a result that was somewhat consistent with the neural networks.

328 The overall results again contradict the results of [?], who conclude that all of the machine methods  
329 agree on the same subset of important factors. Indeed, we only see consistency in variable importance  
330 between the elastic nets and random forests on the individual factors only - all other variable  
331 importance metrics were either inconsistent between different models, or non-robust.

## 332 5 Conclusion

333 Our findings demonstrate that the field of machine learning may offer certain tools to improve stock  
334 prediction and identification of true underlying factors. Penalized linear models and to a lesser extent,  
335 random forests are the best performing methods in the analysis undertaken.

336 Importantly, we find that DFNs fail in the context of stock return prediction, at both prediction  
337 performance and variable importance analysis. This result is consistent across a variety of simulated  
338 datasets, as well as empirical data.

339 Lastly, we find that the top performing models - the elastic nets and random forests, tend to agree  
340 and correctly identify the correct underlying regressors in simulated contexts, and agree on the  
341 same subset of factors which are important in empirical contexts. We find that of all the models  
342 considered, the elastic nets are the most consistent at identifying true underlying regressors through  
343 the simulation study. We find that in the empirical setting, among the individual factors the 1 and  
344 6 month momentum factors are the most powerful predictors of stock returns, according to the  
345 penalized linear models and random forests.

346 The overall findings of this paper differ from the sparse literature on machine learning methods in  
347 empirical finance. However, the performance of the penalized linear models with respect to both out  
348 of sample prediction performance and variable importance analysis is promising, and our findings  
349 show that machine learning provides some tools which may aid in the problems of stock return  
350 prediction and risk factor selection in the financial world.

## 351 5.1 Retrieval of style files

352 The style files for NeurIPS and other conference information are available on the World Wide Web at

353 <http://www.neurips.cc/>

354 The file `neurips_2020.pdf` contains these instructions and illustrates the various formatting re-  
355 quirements your NeurIPS paper must satisfy.

356 The only supported style file for NeurIPS 2020 is `neurips_2020.sty`, rewritten for L<sup>A</sup>T<sub>E</sub>X 2<sub>ε</sub>.  
357 **Previous style files for L<sup>A</sup>T<sub>E</sub>X 2.09, Microsoft Word, and RTF are no longer supported!**

358 The L<sup>A</sup>T<sub>E</sub>X style file contains three optional arguments: `final`, which creates a camera-ready copy,  
359 `preprint`, which creates a preprint for submission to, e.g., arXiv, and `nonatbib`, which will not  
360 load the `natbib` package for you in case of package clash.

361 **Preprint option** If you wish to post a preprint of your work online, e.g., on arXiv, using the  
362 NeurIPS style, please use the `preprint` option. This will create a nonanonymized version of your  
363 work with the text “Preprint. Work in progress.” in the footer. This version may be distributed as  
364 you see fit. Please **do not** use the `final` option, which should **only** be used for papers accepted to  
365 NeurIPS.

366 At submission time, please omit the `final` and `preprint` options. This will anonymize your  
367 submission and add line numbers to aid review. Please do *not* refer to these line numbers in your  
368 paper as they will be removed during generation of camera-ready copies.

369 The file `neurips_2020.tex` may be used as a “shell” for writing your paper. All you have to do is  
370 replace the author, title, abstract, and text of the paper with your own.

371 The formatting instructions contained in these style files are summarized in Sections 6, 7, and 8  
372 below.

## 373 6 General formatting instructions

374 The text must be confined within a rectangle 5.5 inches (33 picas) wide and 9 inches (54 picas) long.  
375 The left margin is 1.5 inch (9 picas). Use 10 point type with a vertical spacing (leading) of 11 points.  
376 Times New Roman is the preferred typeface throughout, and will be selected for you by default.  
377 Paragraphs are separated by 1/2 line space (5.5 points), with no indentation.

378 The paper title should be 17 point, initial caps/lower case, bold, centered between two horizontal  
379 rules. The top rule should be 4 points thick and the bottom rule should be 1 point thick. Allow 1/4 inch  
380 space above and below the title to rules. All pages should start at 1 inch (6 picas) from the top of the  
381 page.

382 For the final version, authors’ names are set in boldface, and each name is centered above the  
383 corresponding address. The lead author’s name is to be listed first (left-most), and the co-authors’  
384 names (if different address) are set to follow. If there is only one co-author, list both author and  
385 co-author side by side.

386 Please pay special attention to the instructions in Section 8 regarding figures, tables, acknowledgments,  
387 and references.

## 388 7 Headings: first level

389 All headings should be lower case (except for first word and proper nouns), flush left, and bold.

390 First-level headings should be in 12-point type.

### 391 7.1 Headings: second level

392 Second-level headings should be in 10-point type.

### 393 7.1.1 Headings: third level

394 Third-level headings should be in 10-point type.

395 **Paragraphs** There is also a `\paragraph` command available, which sets the heading in bold, flush  
396 left, and inline with the text, with the heading followed by 1 em of space.

## 397 8 Citations, figures, tables, references

398 These instructions apply to everyone.

### 399 8.1 Citations within the text

400 The `natbib` package will be loaded for you by default. Citations may be author/year or numeric, as  
401 long as you maintain internal consistency. As to the format of the references themselves, any style is  
402 acceptable as long as it is used consistently.

403 The documentation for `natbib` may be found at

404 <http://mirrors.ctan.org/macros/latex/contrib/natbib/natnotes.pdf>

405 Of note is the command `\citet`, which produces citations appropriate for use in inline text. For  
406 example,

407 `\citet{hasselmo}` investigated\dots

408 produces

409 Hasselmo, et al. (1995) investigated...

410 If you wish to load the `natbib` package with options, you may add the following before loading the  
411 `neurips_2020` package:

412 `\PassOptionsToPackage{options}{natbib}`

413 If `natbib` clashes with another package you load, you can add the optional argument `nonatbib`  
414 when loading the style file:

415 `\usepackage[nonatbib]{neurips_2020}`

416 As submission is double blind, refer to your own published work in the third person. That is, use “In  
417 the previous work of Jones et al. [4],” not “In our previous work [4].” If you cite your other papers  
418 that are not widely available (e.g., a journal paper under review), use anonymous author names in the  
419 citation, e.g., an author of the form “A. Anonymous.”

### 420 8.2 Footnotes

421 Footnotes should be used sparingly. If you do require a footnote, indicate footnotes with a number<sup>10</sup>  
422 in the text. Place the footnotes at the bottom of the page on which they appear. Precede the footnote  
423 with a horizontal rule of 2 inches (12 picas).

424 Note that footnotes are properly typeset *after* punctuation marks.<sup>11</sup>

### 425 8.3 Figures

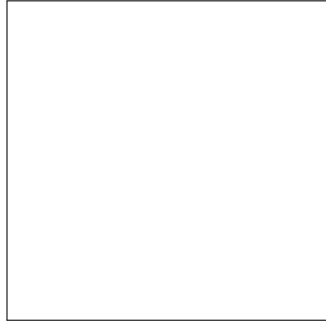
426 All artwork must be neat, clean, and legible. Lines should be dark enough for purposes of reproduction.  
427 The figure number and caption always appear after the figure. Place one line space before the figure  
428 caption and one line space after the figure. The figure caption should be lower case (except for first  
429 word and proper nouns); figures are numbered consecutively.

---

<sup>10</sup>Sample of the first footnote.

<sup>11</sup>As in this example.

Figure 6: Sample figure caption.



Macroeconomic Factors shown on x axis (see Table 5 for definitions)

Table 3: Sample table title

| Part     |                 |                        |
|----------|-----------------|------------------------|
| Name     | Description     | Size ( $\mu\text{m}$ ) |
| Dendrite | Input terminal  | $\sim 100$             |
| Axon     | Output terminal | $\sim 10$              |
| Soma     | Cell body       | up to $10^6$           |

430 You may use color figures. However, it is best for the figure captions and the paper body to be legible  
431 if the paper is printed in either black/white or in color.

## 432 8.4 Tables

433 All tables must be centered, neat, clean and legible. The table number and title always appear before  
434 the table. See Table 3.

435 Place one line space before the table title, one line space after the table title, and one line space after  
436 the table. The table title must be lower case (except for first word and proper nouns); tables are  
437 numbered consecutively.

438 Note that publication-quality tables *do not contain vertical rules*. We strongly suggest the use of the  
439 booktabs package, which allows for typesetting high-quality, professional tables:

440 <https://www.ctan.org/pkg/booktabs>

441 This package was used to typeset Table 3.

## 442 9 Final instructions

443 Do not change any aspects of the formatting parameters in the style files. In particular, do not modify  
444 the width or length of the rectangle the text should fit into, and do not change font sizes (except  
445 perhaps in the **References** section; see below). Please note that pages should be numbered.

## 446 10 Preparing PDF files

447 Please prepare submission files with paper size “US Letter,” and not, for example, “A4.”

448 Fonts were the main cause of problems in the past years. Your PDF file must only contain Type 1 or  
449 Embedded TrueType fonts. Here are a few instructions to achieve this.

- 450 • You should directly generate PDF files using `pdflatex`.

- You can check which fonts a PDF file uses. In Acrobat Reader, select the menu Files>Document Properties>Fonts and select Show All Fonts. You can also use the program `pdffonts` which comes with `xpdf` and is available out-of-the-box on most Linux machines.
- The IEEE has recommendations for generating PDF files whose fonts are also acceptable for NeurIPS. Please see <http://www.emfield.org/icuwb2010/downloads/IEEE-PDF-SpecV32.pdf>
- `xfig` "patterned" shapes are implemented with bitmap fonts. Use "solid" shapes instead.
- The `\bbold` package almost always uses bitmap fonts. You should use the equivalent AMS Fonts:

```
\usepackage{amsfonts}
```

followed by, e.g., `\mathbb{R}`, `\mathbb{N}`, or `\mathbb{C}` for  $\mathbb{R}$ ,  $\mathbb{N}$  or  $\mathbb{C}$ . You can also use the following workaround for reals, natural and complex:

```
\newcommand{\RR}{\mathbb{R}} %real numbers
\newcommand{\Nat}{\mathbb{N}} %natural numbers
\newcommand{\CC}{\mathbb{C}} %complex numbers
```

Note that `amsfonts` is automatically loaded by the `amssymb` package.

If your file contains type 3 fonts or non embedded TrueType fonts, we will ask you to fix it.

## 10.1 Margins in L<sup>A</sup>T<sub>E</sub>X

Most of the margin problems come from figures positioned by hand using `\special` or other commands. We suggest using the command `\includegraphics` from the `graphicx` package. Always specify the figure width as a multiple of the line width as in the example below:

```
\usepackage[pdftex]{graphicx} ...
\includegraphics[width=0.8\linewidth]{myfile.pdf}
```

See Section 4.4 in the `graphics` bundle documentation (<http://mirrors.ctan.org/macros/latex/required/graphics/grfguide.pdf>)

A number of width problems arise when L<sup>A</sup>T<sub>E</sub>X cannot properly hyphenate a line. Please give LaTeX hyphenation hints using the `\-` command when necessary.

## Broader Impact

Authors are required to include a statement of the broader impact of their work, including its ethical aspects and future societal consequences. Authors should discuss both positive and negative outcomes, if any. For instance, authors should discuss a) who may benefit from this research, b) who may be put at disadvantage from this research, c) what are the consequences of failure of the system, and d) whether the task/method leverages biases in the data. If authors believe this is not applicable to them, authors can simply state this.

Use unnumbered first level headings for this section, which should go at the end of the paper. **Note that this section does not count towards the eight pages of content that are allowed.**

## References

References follow the acknowledgments. Use unnumbered first-level heading for the references. Any choice of citation style is acceptable as long as you are consistent. It is permissible to reduce the font size to `small` (9 point) when listing the references. **Note that the Reference section does not count towards the eight pages of content that are allowed.**

[1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In G. Tesauro, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp. 609–616. Cambridge, MA: MIT Press.

- 495 [2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the*  
496 *General NEural Simulation System*. New York: TELOS/Springer-Verlag.
- 497 [3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent  
498 synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.

## 499 **A Additional details: models**

500 In this section, we give a brief overview of all the models considered in the simulation and empirical study.

### 501 **A.1 Linear models**

502 Linear models model the conditional expectation  $g^*(z_{i,t})$  as a linear function of the predictors and the parameter  
503 vector  $\theta$ :

$$g(z_{i,t}; \theta) = z'_{i,t} \theta \quad (18)$$

504 This yields the OLS estimator when optimized w.r.t. MSE, and the LAD estimator when optimized w.r.t. MAE.

### 505 **A.2 Elastic nets**

506 Elastic Nets are similar to linear models but differ via the addition of a penalty term in the loss function:

$$\mathcal{L}(\theta; \cdot) = \underbrace{\mathcal{L}(\theta)}_{\text{Loss Function}} + \underbrace{\phi(\theta; \cdot)}_{\text{Penalty Term}} \quad (19)$$

507 where the elastic net penalty [?] is:

$$\phi(\theta; \lambda, \rho) = \lambda(1 - \rho) \sum_{j=1}^P |\theta_j| + \frac{1}{2} \lambda \rho \sum_{j=1}^P \theta_j^2 \quad (20)$$

508 Further details are given in [?].

### 509 **A.3 Random forests**

510 Further details are given in cite().

### 511 **A.4 Feed forward neural networks**

512 For our application, we considered the following grid of hyperparameters:

513 Further details are given in cite().

### 514 **A.5 Long short term memory networks**

515 Long short term memory (LSTM) networks are

516 For our application, we considered the following grid of hyperparameters:

517 Further details are given in cite().

### 518 **A.6 FFORMA**

519 Feature-based Forecast Model Averaging, cite() is an automated method for obtaining weighted forecast  
520 combinations for time series. We provide a brief overview of the two phases in this methodology.

521 First, we use a collection of

522 To incorporate all regressors in each individual time series model, we applied dimensional reduction techniques  
523 of PCA and UMAP to generate new feature mappings for use in GARCH (1, 1) models (generally the best  
524 performing of the constituent models). It was noted that none of the feature mappings improved fit, however.

525 The constituent models we considered are:

- 526 • Naive
- 527 • Random walk with drift
- 528 • Theta method
- 529 • ARIMA
- 530 • ETS
- 531 • TBATS
- 532 • Neural network auto-regressive model



533       • ARMA (1, 1) with g.e.d. GARCH(1, 1) errors  
534       • ARMA (1, 1) with g.e.d. GARCH(1, 1) errors and UMAP regressors  
535   The time series features used to train the meta-model are detailed in cite(), with the addition of realized volatility.  
536   **A.7   DeepAR**  
537   DeepAR is a generalization of traditional Auto Regressive (AR) models to include additional layers into order to  
538   introduce non-linearities into the model.  
539   Further details are given in cite().

## 540 **A Additional details: simulation design**

541 In this section, we give additional features of the simulation design required to impenent our results. All code  
542 and data can be found at XXXX.

### 543 **A.1 Simulation Design**

544 We simulate a latent factor model with a stochastic volatility process for excess returns  $r_{t+1}$ , for  $t = 1, \dots, T$ :

$$r_{i,t+1} = g(z_{i,t}) + \beta_{i,t+1}v_{t+1} + e_{i,t+1}; \quad z_{i,t} = (1, x_t)' \otimes c_{i,t}, \quad \beta_{i,t} = (c_{i1,t}, c_{i2,t}, c_{i3,t}) \quad (21)$$

$$e_{i,t+1} = \sigma_{i,t+1}\varepsilon_{i,t+1}; \quad (22)$$

$$\log(\sigma_{i,t+1}^2) = \omega + \gamma \log(\sigma_t^2) + \sigma_u u; \quad u \sim N(0, 1) \quad (23)$$

545 Let  $v_{t+1}$  be a  $3 \times 1$  vector of errors, and  $w_{t+1} \sim N(0, 1)$  and  $\varepsilon_{i,t+1} \sim N(0, 1)$  scalar error terms.

546 The matrix  $C_t$  is an  $N \times P_c$  matrix of latent factors, where the first three columns correspond to  $\beta_{i,t}$ , across the  
547  $1 \leq i \leq N$  dimensions, while the remaining  $P_c - 3$  factors do not enter the return equation. The  $P_x \times 1$  vector  
548  $x_t$  is a  $3 \times 1$  multivariate time series, and  $\varepsilon_{t+1}$  is a  $N \times 1$  vector of idiosyncratic errors.

549 The parameters of these were tuned such that the annualized volatility of each return series was approximately  
550 22%, as is often observed empirically.

551 Note that we also reproduce [?]'s error specification as a case where there is no stochastic volatility:

$$v_{t+1} \sim N(0, 0.05^2 \times I_3) \quad (24)$$

$$e_{i,t+1} \sim t_5(0, 0.05^2) \quad (25)$$

#### 552 **A.1.1 Simulating Characteristics**

553 We build in correlation across time among factors by drawing normal random numbers for each  $1 \leq i \leq N$  and  
554  $1 \leq j \leq P_c$ , according to

$$\bar{c}_{ij,t} = \rho_j \bar{c}_{ij,t-1} + \epsilon_{ij,t}; \quad \rho_j \sim \mathcal{U}\left(\frac{1}{2}, 1\right) \quad (26)$$

555 To build in cross sectional correlation, we define the positive-semidefinite matrix  $B$ :

$$B := \Lambda \Lambda' + \frac{1}{10} \mathbb{I}_n, \quad \Lambda_i = (\lambda_{i1}, \dots, \lambda_{i4}), \quad \lambda_{ik} \sim N(0, \lambda_{sd}), \quad k = 1, \dots, 4 \quad (27)$$

556 to serve as a variance covariance matrix with  $\lambda_{sd}$  controlling the density of the matrix, and hence degree of cross  
557 sectional correlation.  $\lambda_{sd}$  values of 0.01, 0.1 and 1 were used to explore increasing degrees of cross sectional  
558 correlation.

559 To build this into our  $N \times P_c$  characteristics matrix  $\bar{C}_t$ , we simulate characteristics according to

$$\hat{C}_t = L \bar{C}_t; \quad B = LL' \quad (28)$$

560 where  $L$  represents the lower triangle matrix of  $B$  using the Cholesky decomposition.

561 Finally, the "observed" characteristics for each  $1 \leq i \leq N$  and for  $j = 1, \dots, P_c$  are constructed according to:

$$c_{ij,t} = \frac{2}{n+1} \text{rank}(\hat{c}_{ij,t}) - 1. \quad (29)$$

562 with the rank transformation normalizing all predictors to be within  $[-1, 1]$ .

#### 563 **A.1.2 Simulating Macroeconomic Series**

564 For simulation of  $x_t$ , a  $3 \times 1$  multivariate time series, we consider a Vector Autoregression (VAR) model <sup>12</sup>:

<sup>12</sup>Other more complex and interactive matrix specifications of  $A$  were briefly explored, but these did not appear to have a significant impact on results. More complex designs were observed to only affect the variable importance metrics, but to an insignificant degree

$$x_t = Ax_{t-1} + u_t; \quad A = \begin{pmatrix} .95 & 0 & 0 \\ 0 & .95 & 0 \\ 0 & 0 & .95 \end{pmatrix} \quad u_t \sim N(\mu = (0, 0, 0)', \Sigma = I_3)$$

### 565 A.1.3 Simulating Return Series

566 We consider three different functions for  $g(z_{i,t})$ :

$$(1) g_1(z_{i,t}) = (c_{i1,t}, c_{i2,t}, c_{i3,t} \times x'_t[3,]) \theta_0 \quad (30)$$

$$(2) g_2(z_{i,t}) = (c_{i1,t}^2, c_{i1,t} \times c_{i2,t}, \text{sgn}(c_{i3,t} \times x'_t[3,])) \theta_0 \quad (31)$$

$$(3) g_3(z_{i,t}) = (1[c_{i3,t} > 0], c_{i2,t}^3, c_{i1,t} \times c_{i2,t} \times 1[c_{i3,t} > 0], \text{logit}(c_{i3,t})) \theta_0 \quad (32)$$

567 where  $x'_t[3,]$  denotes the third element of the  $x'_t$  vector.

568  $g_1(z_{i,t})$  allows the characteristics to enter the return equation linearly, and  $g_2(z_{i,t})$  allows the characteristics to  
 569 enter the return equation interactively and non-linearly. The true underlying regressors for these specifications  
 570 are  $(c_{i1,t}, c_{i2,t}, c_{i3,t} \times x'_t[3,])$ . These two specifications correspond to the simulation design used by [?].

571  $g_3(z_{i,t})$  allows the characteristics to enter in a complex and non-linear fashion. The true underlying regressors  
 572 for this specification are  $(c_{i1,t}, c_{i2,t}, c_{i3,t})$ .

573 It should be noted however, that because  $g_2(z_{i,t})$  has a large part of its signal entering through a  $\text{sgn}$  function,  
 574 this should make it the most difficult to estimate given the regressors and resulting returns process.

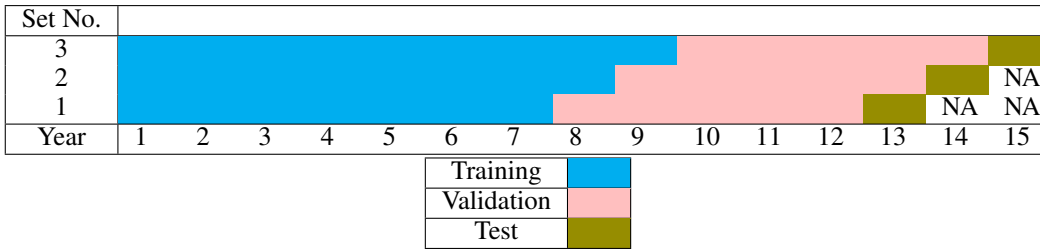
575  $\theta^0$  was tuned such that the predictive  $R^2$  was approximately 5%.

576 The simulation design results in  $3 \times 3 = 12$  different simulated datasets, each with  $N = 200$  stocks,  $T = 180$   
 577 periods and  $P_c = 100$  characteristics. Each design was simulated 10 times to assess the robustness of machine  
 578 learning algorithms. The number of simulations was kept low for computational feasibility.

### 579 A.1.4 Sample Splitting

580 If viewed as monthly periods,  $T = 180$  corresponds to 15 years. A data splitting scheme similar to the scheme  
 581 to be used in the empirical data study was used: a training:validation length ratio of approximately 1.5 to begin,  
 582 and a test set that is 1 year in length. We employ the hybrid growing window approach as described earlier in  
 583 section ?? (see Figure 7 for a graphical representation).

Figure 7: Sample Splitting Procedure



584 Other schemes in the forecasting literature such as using an “inner” rolling window validation loop to find  
 585 the best hyperparameters on average, finally aggregating them in an “outer” loop for a more robust error were  
 586 considered but not implemented for a variety of reasons. Firstly, many of the models were computationally too  
 587 intensive for this to be feasible. More importantly, during the model fitting process it was observed that the  
 588 optimal hyperparameters for the different rolling windows were highly unstable. Thus, this would have made the  
 589 selection of the best hyperparameters on average across all windows significantly less meaningful.

## 590 A.2 Simulation Study Results

### 591 A.2.1 Prediction Performance

Table 4: Simulation Study Loss Statistics

| model   | Corr | g <sup>1</sup> |           |            | g <sup>2</sup> |           |            | g <sup>3</sup> |           |            |
|---------|------|----------------|-----------|------------|----------------|-----------|------------|----------------|-----------|------------|
|         |      | Test MAE       | Test MSE  | Test $R^2$ | Test MAE       | Test MSE  | Test $R^2$ | Test MAE       | Test MSE  | Test $R^2$ |
| LM.MSE  | 0.01 | 0.0366775      | 0.0027400 | 0.0082732  | 0.0382548      | 0.0028801 | -0.1117880 | 0.0373098      | 0.0027954 | -0.0320680 |
|         | 0.10 | 0.0369652      | 0.0027653 | -0.0110198 | 0.0385796      | 0.0029144 | -0.1429443 | 0.0375694      | 0.0028168 | -0.0549404 |
|         | 1.00 | 0.0429486      | 0.0034141 | -0.4387965 | 0.0453765      | 0.0037172 | -0.7809535 | 0.0434339      | 0.0034688 | -0.4887785 |
| LM.MAE  | 0.01 | 0.0366417      | 0.0027373 | 0.0090496  | 0.0383478      | 0.0028862 | -0.1163694 | 0.0373235      | 0.0027967 | -0.0351619 |
|         | 0.10 | 0.0368113      | 0.0027555 | 0.0029188  | 0.0387449      | 0.0029275 | -0.1525797 | 0.0374894      | 0.0028098 | -0.0476746 |
|         | 1.00 | 0.0423399      | 0.0033445 | -0.3930442 | 0.0453420      | 0.0036847 | -0.7699555 | 0.0435349      | 0.0034682 | -0.5445237 |
| ELN.MSE | 0.01 | 0.0345878      | 0.0025663 | 0.1403351  | 0.0362229      | 0.0026898 | 0.0368766  | 0.0353534      | 0.0026227 | 0.0991416  |
|         | 0.10 | 0.0345630      | 0.0025643 | 0.1442376  | 0.0361830      | 0.0026860 | 0.0372585  | 0.0352923      | 0.0026167 | 0.1002410  |
|         | 1.00 | 0.0346142      | 0.0025676 | 0.1671841  | 0.0362761      | 0.0026980 | 0.0378391  | 0.0354437      | 0.0026300 | 0.1198755  |
| ELN.MAE | 0.01 | 0.0345786      | 0.0025652 | 0.1409821  | 0.0361950      | 0.0026882 | 0.0391694  | 0.0353345      | 0.0026210 | 0.1004424  |
|         | 0.10 | 0.0345582      | 0.0025637 | 0.1446272  | 0.0361730      | 0.0026877 | 0.0388747  | 0.0352851      | 0.0026167 | 0.1009186  |
|         | 1.00 | 0.0345989      | 0.0025667 | 0.1677712  | 0.0363047      | 0.0027028 | 0.0365834  | 0.0354652      | 0.0026310 | 0.1180225  |
| RF.MSE  | 0.01 | 0.0357752      | 0.0026710 | 0.0634257  | 0.0357179      | 0.0026571 | 0.0676147  | 0.0358032      | 0.0026613 | 0.0702977  |
|         | 0.10 | 0.0357695      | 0.0026649 | 0.0667382  | 0.0356845      | 0.0026525 | 0.0691389  | 0.0358666      | 0.0026704 | 0.0628386  |
|         | 1.00 | 0.0362325      | 0.0026977 | 0.0687741  | 0.0359893      | 0.0026833 | 0.0571035  | 0.0362129      | 0.0026952 | 0.0698868  |
| RF.MAE  | 0.01 | 0.0354594      | 0.0026434 | 0.0833385  | 0.0354204      | 0.0026305 | 0.0876529  | 0.0355399      | 0.0026446 | 0.0865291  |
|         | 0.10 | 0.0355153      | 0.0026489 | 0.0814253  | 0.0354894      | 0.0026345 | 0.0834048  | 0.0355688      | 0.0026438 | 0.0816426  |
|         | 1.00 | 0.0359158      | 0.0026747 | 0.0870806  | 0.0356434      | 0.0026445 | 0.0809651  | 0.0360529      | 0.0026786 | 0.0753573  |
| NN1.MSE | 0.01 | 0.0364516      | 0.0027219 | 0.0163443  | 0.0367677      | 0.0027319 | -0.0039174 | 0.0366874      | 0.0027384 | 0.0093355  |
|         | 0.10 | 0.0364624      | 0.0027191 | 0.0204223  | 0.0367762      | 0.0027345 | -0.0072588 | 0.0367326      | 0.0027372 | 0.0029550  |
|         | 1.00 | 0.0375452      | 0.0028206 | -0.0144520 | 0.0370492      | 0.0027638 | -0.0146973 | 0.0374589      | 0.0027975 | -0.0124689 |
| NN1.MAE | 0.01 | 0.0359604      | 0.0026786 | 0.0558139  | 0.0369206      | 0.0027474 | -0.0151053 | 0.0363047      | 0.0026996 | 0.0393707  |
|         | 0.10 | 0.0360823      | 0.0026866 | 0.0506976  | 0.0370100      | 0.0027503 | -0.0205616 | 0.0363220      | 0.0027022 | 0.0323034  |
|         | 1.00 | 0.0378894      | 0.0028338 | -0.0431818 | 0.0379790      | 0.0028445 | -0.0840747 | 0.0373056      | 0.0027926 | 0.0021783  |
| NN2.MSE | 0.01 | 0.0370187      | 0.0027850 | -0.0217869 | 0.0373197      | 0.0027752 | -0.0433537 | 0.0370890      | 0.0027745 | -0.0173037 |
|         | 0.10 | 0.0369775      | 0.0027651 | -0.0212763 | 0.0370088      | 0.0027478 | -0.0275384 | 0.0369898      | 0.0027584 | -0.0206446 |
|         | 1.00 | 0.0375360      | 0.0028138 | -0.0139783 | 0.0369035      | 0.0027518 | -0.0058664 | 0.0375157      | 0.0028087 | -0.0169336 |
| NN2.MAE | 0.01 | 0.0358939      | 0.0026718 | 0.0577427  | 0.0368335      | 0.0027396 | -0.0071579 | 0.0363352      | 0.0027028 | 0.0363052  |
|         | 0.10 | 0.0358898      | 0.0026681 | 0.0603096  | 0.0369367      | 0.0027503 | -0.0170774 | 0.0362701      | 0.0026960 | 0.0371567  |
|         | 1.00 | 0.0374795      | 0.0028142 | -0.0095290 | 0.0377146      | 0.0028226 | -0.0653904 | 0.0374711      | 0.0028038 | -0.0101183 |
| NN3.MSE | 0.01 | 0.0367827      | 0.0027568 | -0.0067616 | 0.0368397      | 0.0027379 | -0.0075249 | 0.0370360      | 0.0027644 | -0.0200783 |
|         | 0.10 | 0.0369384      | 0.0027613 | -0.0153994 | 0.0368517      | 0.0027384 | -0.0151060 | 0.0368743      | 0.0027573 | -0.0044063 |
|         | 1.00 | 0.0374242      | 0.0028081 | -0.0129638 | 0.0369376      | 0.0027543 | -0.0063529 | 0.0374202      | 0.0027991 | -0.0103479 |
| NN3.MAE | 0.01 | 0.0358164      | 0.0026697 | 0.0654321  | 0.0369345      | 0.0027491 | -0.0163983 | 0.0364712      | 0.0027181 | 0.0299484  |
|         | 0.10 | 0.0358935      | 0.0026771 | 0.0620017  | 0.0368590      | 0.0027406 | -0.0118497 | 0.0362000      | 0.0026932 | 0.0406114  |
|         | 1.00 | 0.0370087      | 0.0027744 | 0.0213288  | 0.0372705      | 0.0027832 | -0.0296437 | 0.0374132      | 0.0027916 | -0.0083067 |
| NN4.MSE | 0.01 | 0.0368808      | 0.0027586 | -0.0206197 | 0.0368555      | 0.0027423 | -0.0077152 | 0.0371255      | 0.0027752 | -0.0265634 |
|         | 0.10 | 0.0368772      | 0.0027610 | -0.0145791 | 0.0372207      | 0.0027615 | -0.0487112 | 0.0368718      | 0.0027480 | -0.0088940 |
|         | 1.00 | 0.0373820      | 0.0028051 | -0.0064811 | 0.0368966      | 0.0027505 | -0.0053689 | 0.0373542      | 0.0027970 | -0.0077389 |
| NN4.MAE | 0.01 | 0.0359348      | 0.0026782 | 0.0577196  | 0.0368974      | 0.0027487 | -0.0109166 | 0.0367079      | 0.0027376 | 0.0070464  |
|         | 0.10 | 0.0358281      | 0.0026651 | 0.0650415  | 0.0369333      | 0.0027494 | -0.0191117 | 0.0362730      | 0.0026954 | 0.0377039  |
|         | 1.00 | 0.0370948      | 0.0027786 | 0.0198663  | 0.0373230      | 0.0027947 | -0.0293767 | 0.0373013      | 0.0027871 | -0.0018876 |
| NN5.MSE | 0.01 | 0.0372306      | 0.0027846 | -0.0499701 | 0.0369309      | 0.0027474 | -0.0170017 | 0.0371140      | 0.0027720 | -0.0218954 |
|         | 0.10 | 0.0370264      | 0.0027669 | -0.0321897 | 0.0371758      | 0.0027623 | -0.0394362 | 0.0369093      | 0.0027565 | -0.0113522 |
|         | 1.00 | 0.0373642      | 0.0027949 | -0.0104952 | 0.0369277      | 0.0027552 | -0.0053762 | 0.0374751      | 0.0028071 | -0.0149737 |
|         | 0.01 | 0.0358880      | 0.0026693 | 0.0585792  | 0.0368354      | 0.0027380 | -0.0086455 | 0.0366851      | 0.0027371 | 0.0046430  |

Figure 8: g1 BC VIMP

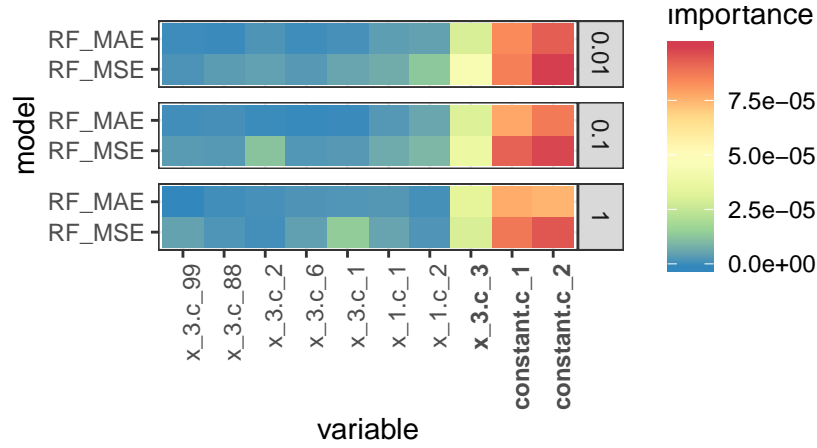


Table 4: Simulation Study Loss Statistics

| model      | Corr | g1        |           |            | g2        |           |            | g3        |           |            |
|------------|------|-----------|-----------|------------|-----------|-----------|------------|-----------|-----------|------------|
|            |      | Test MAE  | Test MSE  | Test $R^2$ | Test MAE  | Test MSE  | Test $R^2$ | Test MAE  | Test MSE  | Test $R^2$ |
| NN5.MAE    | 0.10 | 0.0360381 | 0.0026803 | 0.0509764  | 0.0367451 | 0.0027273 | -0.0049349 | 0.0364843 | 0.0027103 | 0.0181920  |
|            | 1.00 | 0.0372849 | 0.0027940 | 0.0025412  | 0.0370382 | 0.0027652 | -0.0127290 | 0.0371925 | 0.0027753 | 0.0025723  |
| LSTM.MSE   | 0.01 | 0.0372963 | 0.0027982 | -0.0432886 | 0.0372268 | 0.0027764 | -0.0447640 | 0.0375909 | 0.0028180 | -0.0625164 |
|            | 0.10 | 0.0372369 | 0.0027946 | -0.0319550 | 0.0371342 | 0.0027674 | -0.0382547 | 0.0371984 | 0.0027845 | -0.0303936 |
|            | 1.00 | 0.0381282 | 0.0028506 | -0.0820266 | 0.0373821 | 0.0027921 | -0.0442426 | 0.0377803 | 0.0028300 | -0.0443304 |
| LSTM.MAE   | 0.01 | 0.0374310 | 0.0028046 | -0.0564056 | 0.0373372 | 0.0027801 | -0.0518537 | 0.0376270 | 0.0028169 | -0.0674327 |
|            | 0.10 | 0.0374461 | 0.0028036 | -0.0629523 | 0.0371178 | 0.0027679 | -0.0325442 | 0.0372409 | 0.0027931 | -0.0333196 |
|            | 1.00 | 0.0380266 | 0.0028456 | -0.0614833 | 0.0374152 | 0.0027902 | -0.0455057 | 0.0377435 | 0.0028252 | -0.0458837 |
| FFORMA.MSE | 0.01 | 0.0382767 | 0.0028820 | -0.1326717 | 0.0384600 | 0.0028893 | -0.1473902 | 0.0424656 | 0.0033108 | -0.4861451 |
|            | 0.10 | 0.0383581 | 0.0028947 | -0.1407652 | 0.0384795 | 0.0028912 | -0.1600616 | 0.0423231 | 0.0032914 | -0.4739906 |
|            | 1.00 | 0.0388747 | 0.0029647 | -0.1312392 | 0.0388080 | 0.0029331 | -0.1659900 | 0.0430130 | 0.0033713 | -0.4709541 |
| FFORMA.MAE | 0.01 | 0.0387548 | 0.0029387 | -0.1797483 | 0.0387472 | 0.0029178 | -0.1740938 | 0.0429893 | 0.0033651 | -0.5279094 |
|            | 0.10 | 0.0389359 | 0.0029511 | -0.1927930 | 0.0387959 | 0.0029457 | -0.1759939 | 0.0430966 | 0.0034057 | -0.5863752 |
|            | 1.00 | 0.0392468 | 0.0029721 | -0.1636559 | 0.0393873 | 0.0029960 | -0.2116186 | 0.0437090 | 0.0034483 | -0.5260813 |
| DeepAR     | 0.01 | 0.0382993 | 0.0029000 | -0.1289295 | 0.0384895 | 0.0029121 | -0.1325183 | 0.0393898 | 0.0030161 | -0.2049803 |
|            | 0.10 | 0.0388318 | 0.0029353 | -0.1816633 | 0.0384345 | 0.0029045 | -0.1318744 | 0.0391770 | 0.0029932 | -0.1905583 |
|            | 1.00 | 0.0405348 | 0.0031590 | -0.2391417 | 0.0387870 | 0.0029524 | -0.1440285 | 0.0396918 | 0.0030422 | -0.1823646 |

### A.3 Random Forest VIMPs

We note that random forest methods typically have their own built-in ways to calculate variable importance which are different to the variable importance metric presented in the main body of the paper. Here we provide two popular schemes of calculating random forest variable importance metrics - . Importantly, the overall conclusion regarding factor selection does not change with respect to which vimp methodology employed.

Figure 9: g2 BC VIMP

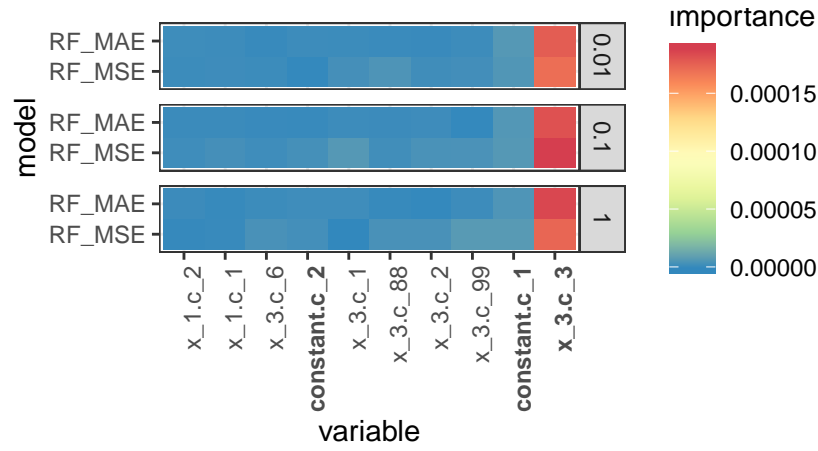


Figure 10: g3 BC VIMP

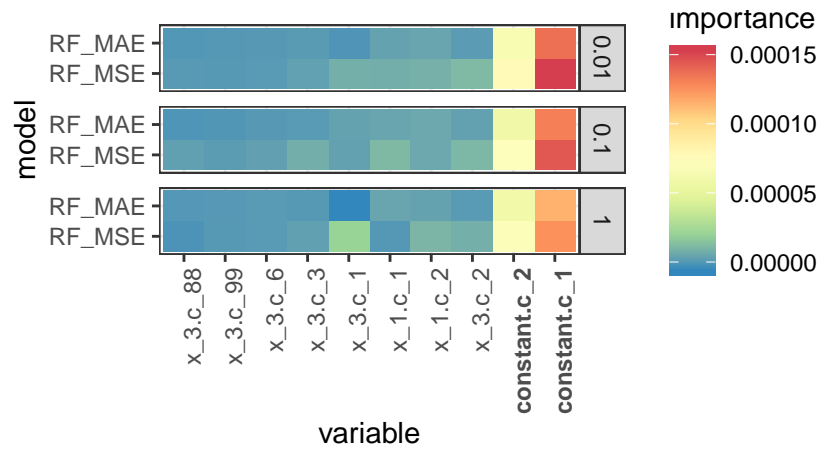


Figure 11: g1 IK VIMP

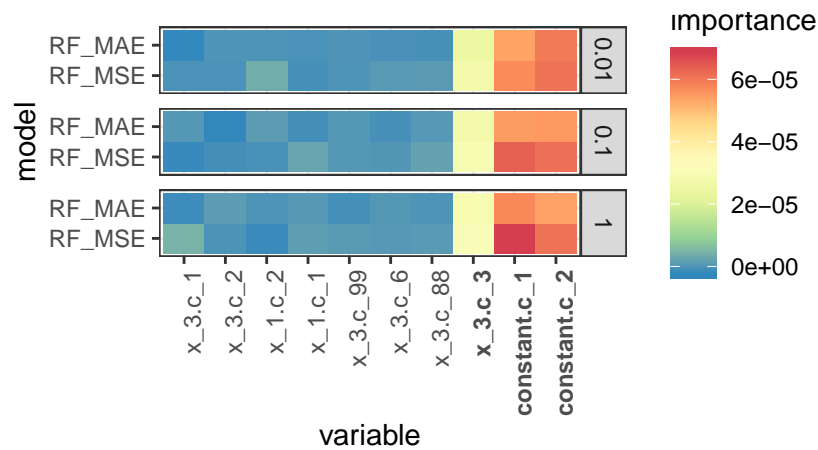


Figure 12: g2 IK VIMP

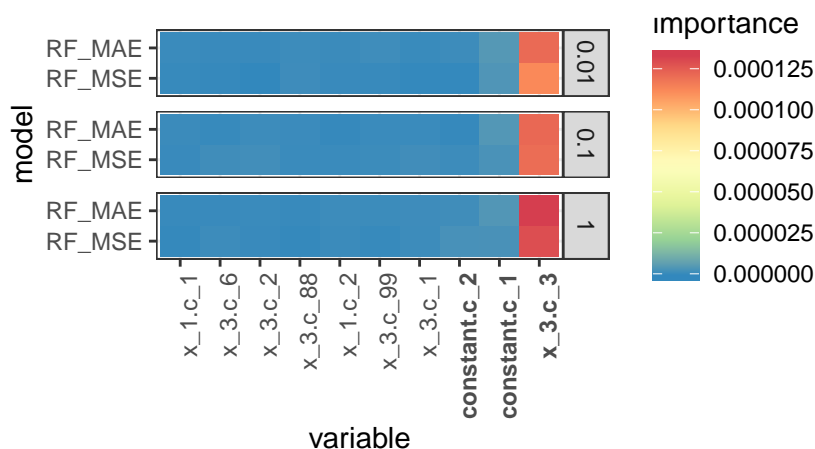
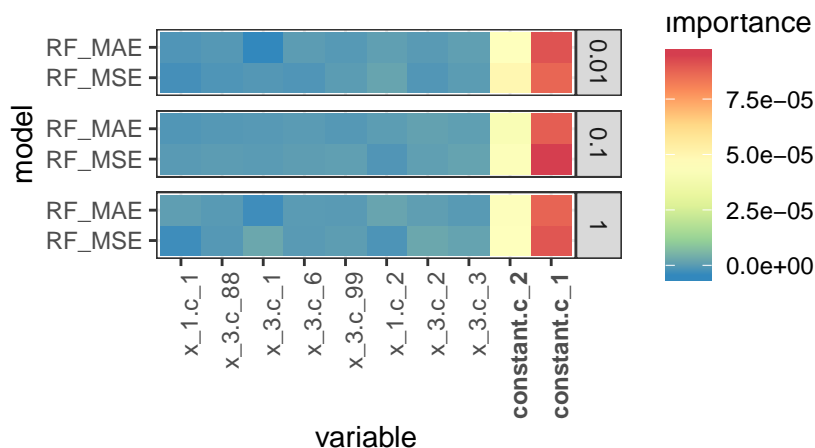


Figure 13: g3 IK VIMP



## A Additional details: Empirical analysis

### A.1 Data & cleaning

We begin by obtaining monthly individual price data from CRSP for all firms listed in the NYSE, AMEX and NASDAQ, starting from 1957 (starting date of the S&P 500) and ending in December 2016, totalling 60 years. To build individual factors, we construct a factor set based on the cross section of returns literature. This data was sourced from and is the same data used in [?]. Like our initial returns sample, it begins in March 1957 and ends in December 2016, totalling 60 years. It contains 94 stock level characteristics: 61 updated annually, 13 updated quarterly and 20 updated monthly, in addition to 74 industry dummies corresponding the the first two digits of the Standard Industrial Classification (SIC) codes. It is noted that this dataset so far contains all securities traded, including those with a CRSP share code other than 10 or 11 and thus includes instruments such as REITs and mutual funds, and those with a share price of less than \$5.

To reduce the size of the dataset and increase feasibility, the dataset was filtered such that only stocks traded primarily on NASDAQ were included (using the PRIMEXCH variable from WRDS). Then, penny stocks (also referred to as microcaps in the literature) with a stock price of less than \$5 were filtered out, as is commonly done in the literature to reduce variability. Stocks without a share code of 10 or 11 (referring to equities) were filtered out, so that securities that are not equities were not included (such as REITs and trust funds). The dataset is provided in a monthly format, which means that many of the factors which are updated only quarterly or annually have very low levels of variability, which can lead to misleading results in the model fitting process. To achieve a balance between having a dataset with enough data points and variability among factors, the dataset

Table 5: Macroeconomic Factors, ([?])

| No. | Acronym    | Macroeconomic Factor |
|-----|------------|----------------------|
| 1   | macro_dp   | Dividend Price Ratio |
| 2   | macro_ep   | Earnings Price Ratio |
| 3   | macro_bm   | Book to Market Ratio |
| 4   | macro_ntis | Net Equity Expansion |
| 5   | macro_tbl  | Treasury Bill Rate   |
| 6   | macro_tms  | Term Spread          |
| 7   | macro_dfy  | Default Spread       |
| 8   | macro_svar | Stock Variance       |

was converted to a quarterly format. Quarterly returns were then constructed using the PRC variable according to actual returns (ie not logged differences):

$$RET_t = \frac{PRC_t - PRC_{t-1}}{PRC_{t-1}} \quad (33)$$

We allow all stocks which have a quarterly return to enter the dataset, even if they disappear from the dataset for certain periods, as opposed to only keeping stocks which appear continuously throughout the entire period. This was primarily done to reduce survivorship bias in the dataset, which can be very prevalent in financial data, and also allows for stocks which were unlisted and relisted again to feature in the dataset.

The sic2 variable, corresponding to the stocks' Standard Industrial Classification (SIC) codes was also dropped. The SIC code system suffers from inconsistent logic in classifying companies, and as a system built for pre-1970s traditional industries has been slow in recognizing new and emerging industries. Indeed, WRDS explicitly cautions the use of SIC codes beyond the use of rough grouping of industries, warning that SIC codes are not strictly enforced by government agencies for accuracy, in addition to most large companies belonging to multiple SIC codes over time. Because of this latter point in particular, there can be inconsistencies on the correct SIC code for the same company depending on the data source. Dropping the sic2 variable also reduced the dimensionality of the dataset by 74 columns, significant increasing computational feasibility.

There existed a significant amount of missing data in the dataset. The dataset's columns were first examined, and any characteristics that had over 20% of their data were removed. However, as the amount of missing data increases dramatically going further back in time, a balance between using more periods at the cost of removing more characteristics versus using less periods but keeping more characteristics was needed. 1993 Q3 was determined to be a reasonable time frame to begin the dataset, as there was a noticeable increase in data availability and quality after this time. Missing characteristics were then imputed using their cross sectional medians for each year.

We then follow [?] and construct eight macroeconomic factors following the variable definitions in [?]. These factors were lagged by one period so as to be used to predict one period ahead quarterly returns. The treasury bill rate was also used from this source to proxy for the risk free rate in order to construct excess quarterly returns.

The two sets of factors were then combined to form a baseline set of covariates, which we define throughout all methods and analysis as:

$$z_{i,t} = (1, x_t)' \otimes c_{i,t} \quad (34)$$

where  $c_{i,t}$  is a  $P_c$  matrix of characteristics for each stock  $i$ , and  $(1, x_t)'$  is a  $P_x \times 1$  vector of macroeconomic predictors, , and  $\otimes$  represents the Kronecker product.  $z_{i,t}$  is therefore a  $P_x P_c$  vector of features for predicting individual stock returns and includes interactions between stock level characteristics and macroeconomic variables. The total number of covariates in this baseline set is  $61 \times (8 + 1) = 549$ <sup>13</sup>.

The dataset was not normalized for all methods, as only penalized regression and neural networks are sensitive to normalization. For these two methods, the dataset was normalized such that each predictor column had 0 mean and 1 variance.

The final dataset spanned from 1993 Q3 to 2016 Q4 with 202, 066 individual observations.

<sup>13</sup>As the individual and macroeconomic factors can have similar names, individual and macroeconomic factors were prefixed with ind\_ and macro\_ respectively.



650 We mimic the procedure used in the simulation study. For the sample splitting procedure, the dataset was split  
 651 such that the training and validation sets were split such that the training set was approximately 1.5 times the  
 652 length of the validation set, in order to predict a test set that is one year in length.

Figure 14: Empirical Data Sample Splitting Procedure

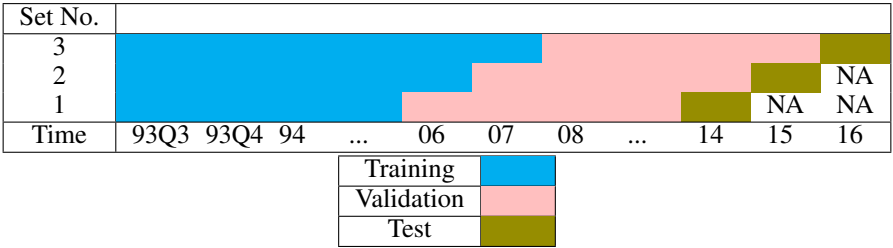


Table 6: Empirical Study Loss Statistics

| model   | Sample 1        |                 |                 | Sample 2        |                 |                 | Sample 3       |                 |                 |
|---------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|----------------|-----------------|-----------------|
|         | Test MAE        | Test MSE        | Test $R^2$      | Test MAE        | Test MSE        | Test $R^2$      | Test MAE       | Test MSE        | Test $R^2$      |
| LM.MSE  | 0.229034        | 0.116015        | -1.808481       | 0.397573        | 0.312653        | -6.329935       | 0.566307       | 0.83804         | -17.522476      |
| LM.MAE  | 0.273452        | 0.15894         | -2.8476         | 0.555673        | 0.742223        | -16.400898      | 0.651614       | 1.225121        | -26.077774      |
| ELN.MSE | 0.133887        | 0.039947        | 0.032956        | 0.140402        | 0.04277         | -0.002712       | <b>0.14433</b> | <b>0.043761</b> | <b>0.032789</b> |
| ELN.MAE | 0.131369        | 0.040718        | 0.014306        | <b>0.137092</b> | <b>0.041892</b> | <b>0.017875</b> | 0.146251       | 0.045207        | 0.000835        |
| RF.MSE  | 0.130366        | <b>0.036629</b> | <b>0.113289</b> | 0.195817        | 0.070642        | -0.656158       | 0.157934       | 0.05122         | -0.132066       |
| RF.MAE  | <b>0.126703</b> | 0.036785        | 0.109505        | 0.173721        | 0.057546        | -0.349132       | 0.14692        | 0.046037        | -0.01752        |
| NN1.MSE | 0.169127        | 0.057044        | -0.380909       | 0.207662        | 0.074751        | -0.752494       | 0.192125       | 0.069738        | -0.541369       |
| NN1.MAE | 0.157324        | 0.050418        | -0.22052        | 0.191762        | 0.066746        | -0.564818       | 0.18547        | 0.063053        | -0.393606       |
| NN2.MSE | 0.168773        | 0.059436        | -0.43883        | 0.181808        | 0.063232        | -0.482433       | 0.180584       | 0.062745        | -0.386797       |
| NN2.MAE | 0.162667        | 0.055447        | -0.342256       | 0.194277        | 0.069386        | -0.626702       | 0.185173       | 0.065186        | -0.440746       |
| NN3.MSE | 0.154784        | 0.050152        | -0.21408        | 0.180103        | 0.060193        | -0.411175       | 0.177604       | 0.060404        | -0.335065       |
| NN3.MAE | 0.146411        | 0.044901        | -0.086967       | 0.18499         | 0.06461         | -0.514744       | 0.184986       | 0.063861        | -0.411475       |
| NN4.MSE | 0.153802        | 0.048641        | -0.177503       | 0.193066        | 0.067515        | -0.582833       | 0.172707       | 0.057774        | -0.276929       |
| NN4.MAE | 0.157301        | 0.050286        | -0.217308       | 0.168815        | 0.055711        | -0.306102       | 0.167998       | 0.055129        | -0.218463       |
| NN5.MSE | 0.149436        | 0.047279        | -0.14452        | 0.183584        | 0.064137        | -0.503653       | 0.170238       | 0.056992        | -0.259652       |
| NN5.MAE | 0.140781        | 0.042832        | -0.036882       | 0.181096        | 0.06216         | -0.4573         | 0.164896       | 0.053458        | -0.181528       |

Table 7: Missing Data Threshold Robustness Check Loss Statistics

| model   | Sample 1        |                 |                 | Sample 2        |                 |                 | Sample 3       |                 |                 |
|---------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|----------------|-----------------|-----------------|
|         | Test MAE        | Test MSE        | Test $R^2$      | Test MAE        | Test MSE        | Test $R^2$      | Test MAE       | Test MSE        | Test $R^2$      |
| LM.MSE  | 0.247457        | 0.130166        | -2.151058       | 0.541089        | 0.700574        | -15.424468      | 0.615714       | 1.188991        | -25.279238      |
| LM.MAE  | 0.214055        | 0.102848        | -1.489727       | 0.372683        | 0.259976        | -5.094954       | 0.507397       | 0.766373        | -15.93847       |
| ELN.MSE | 0.133887        | 0.039947        | 0.032956        | 0.140402        | 0.04277         | -0.002712       | <b>0.14433</b> | <b>0.043761</b> | <b>0.032789</b> |
| ELN.MAE | 0.131338        | 0.040465        | 0.020421        | <b>0.137083</b> | <b>0.041804</b> | <b>0.019938</b> | 0.146589       | 0.045362        | -0.002596       |
| RF.MSE  | 0.129226        | 0.035869        | 0.131692        | 0.198914        | 0.072749        | -0.705542       | 0.168068       | 0.05777         | -0.276838       |
| RF.MAE  | <b>0.124319</b> | <b>0.035103</b> | <b>0.150229</b> | 0.167845        | 0.053578        | -0.256106       | 0.15463        | 0.051594        | -0.140342       |
| NN1.MSE | 0.153785        | 0.048726        | -0.179553       | 0.221019        | 0.084867        | -0.98964        | 0.172557       | 0.058354        | -0.289742       |
| NN1.MAE | 0.154534        | 0.048854        | -0.18266        | 0.199647        | 0.073699        | -0.727823       | 0.176348       | 0.061359        | -0.356155       |
| NN2.MSE | 0.158513        | 0.057061        | -0.381324       | 0.233631        | 0.095004        | -1.227299       | 0.154083       | 0.048353        | -0.068708       |
| NN2.MAE | 0.138489        | 0.043364        | -0.049759       | 0.215253        | 0.078792        | -0.847234       | 0.164459       | 0.055049        | -0.216706       |
| NN3.MSE | 0.167392        | 0.058508        | -0.416345       | 0.19754         | 0.071293        | -0.671422       | 0.156873       | 0.049602        | -0.096299       |
| NN3.MAE | 0.144457        | 0.045293        | -0.096445       | 0.210372        | 0.077747        | -0.822723       | 0.159841       | 0.05152         | -0.138704       |
| NN4.MSE | 0.147989        | 0.047211        | -0.142888       | 0.184277        | 0.064247        | -0.506225       | 0.152214       | 0.048185        | -0.064987       |
| NN4.MAE | 0.15851         | 0.052021        | -0.259326       | 0.18643         | 0.063032        | -0.477746       | 0.177651       | 0.064046        | -0.415562       |
| NN5.MSE | 0.153187        | 0.050053        | -0.211683       | 0.181622        | 0.060313        | -0.413989       | 0.161028       | 0.051221        | -0.132095       |
| NN5.MAE | 0.149496        | 0.050779        | -0.229251       | 0.165726        | 0.053988        | -0.265712       | 0.156151       | 0.049772        | -0.100061       |

## A.2 Empirical study robustness checks

In addition to the main study, we provide four additional robustness checks for our empirical study, with regards to different training/validation splitting schemes, missing data imputation and additional regressors. Importantly, our overall results are consistent across all checks.

We consider training:validation length ratios of 1:1 and 1:2 in addition to 1:1.5 in the main study.

We consider changing the missing data threshold to be 10% - that is, any regressors with over 10% missing data were omitted before being imputed.

We finally consider supplementing our macroeconomic regressor set with the five Fama-French factors.

## A.3 Empirical Data Results

### A.3.1 Prediction Accuracy

Figure 15: Individual Factor Importance

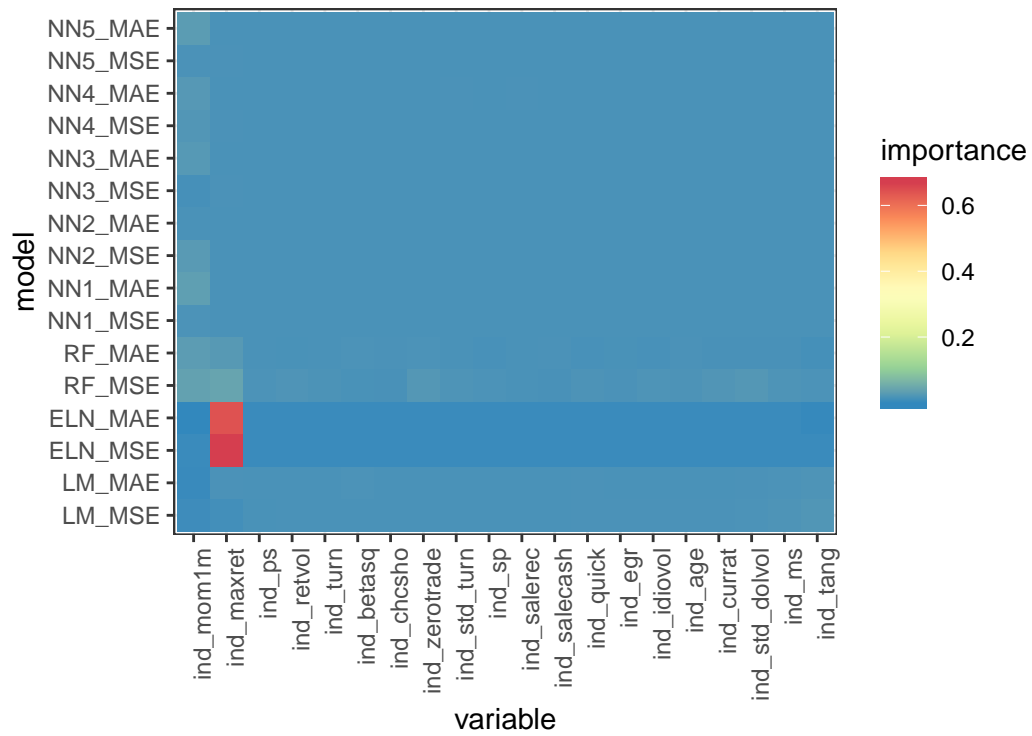


Figure 16: Macroeconomic Factor Importance

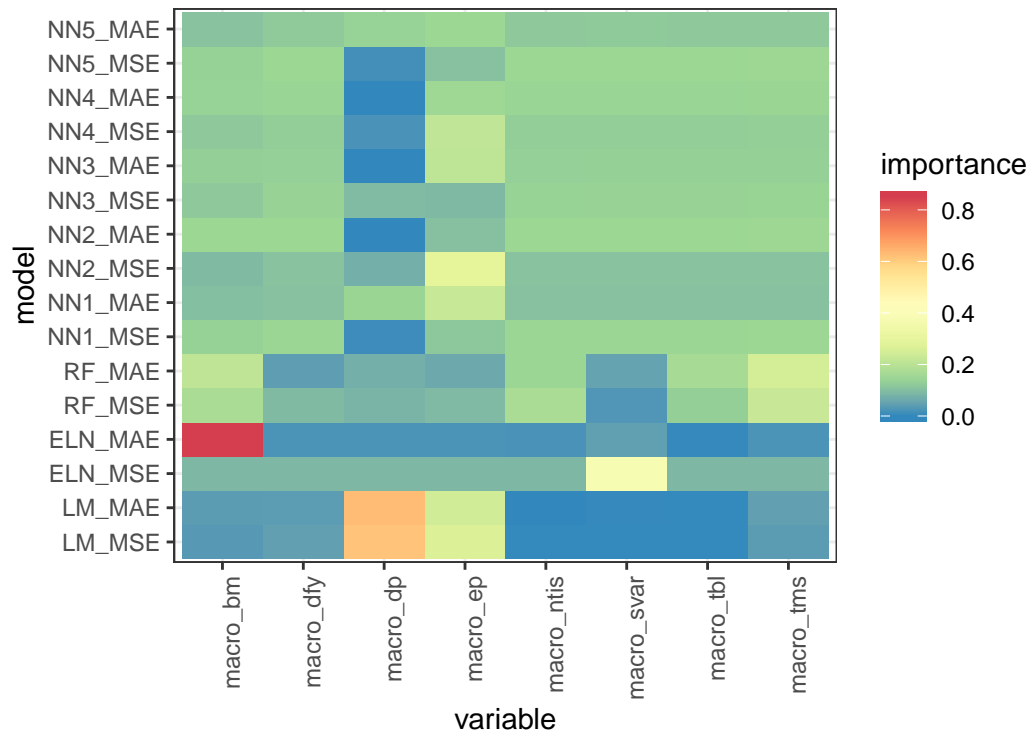


Figure 17: Robustness Check RF VIMP

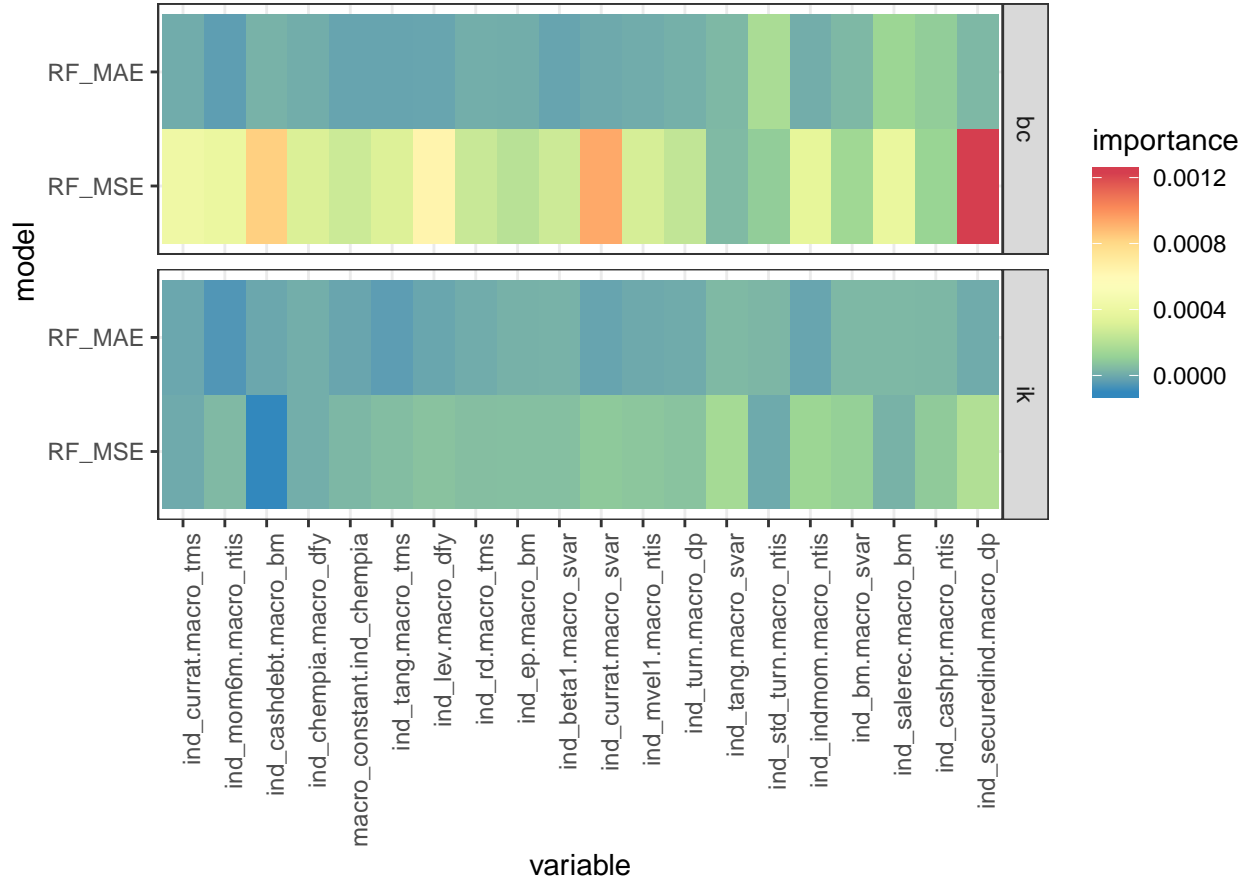


Table 8: Train:Validation 1:1 Robustness Check Loss Statistics

| model   | Sample 1        |                 |                 | Sample 2        |                 |                 | Sample 3        |                 |                 |
|---------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|
|         | Test MAE        | Test MSE        | Test $R^2$      | Test MAE        | Test MSE        | Test $R^2$      | Test MAE        | Test MSE        | Test $R^2$      |
| LM.MSE  | 0.915703        | 2.495094        | -59.401029      | 0.717           | 1.553454        | -35.419641      | 0.451206        | 0.375505        | -7.299459       |
| LM.MAE  | 0.751551        | 1.583265        | -37.32754       | 0.469831        | 0.524686        | -11.300895      | 0.675112        | 1.105759        | -23.43964       |
| ELN.MSE | 0.134609        | <b>0.040072</b> | <b>0.029933</b> | 0.141434        | 0.043169        | -0.012055       | <b>0.144375</b> | <b>0.043705</b> | <b>0.034019</b> |
| ELN.MAE | <b>0.131668</b> | 0.040748        | 0.013583        | <b>0.137494</b> | <b>0.042135</b> | <b>0.012178</b> | 0.146776        | 0.045753        | -0.01123        |
| RF.MSE  | 0.155282        | 0.046655        | -0.129427       | 0.210936        | 0.078006        | -0.828784       | 0.229147        | 0.092622        | -1.047155       |
| RF.MAE  | 0.13882         | 0.04016         | 0.027805        | 0.185338        | 0.063217        | -0.482087       | 0.182753        | 0.063873        | -0.411736       |
| NN1.MSE | 0.218129        | 0.087699        | -1.123002       | 0.238606        | 0.110201        | -1.583582       | 0.260721        | 0.120908        | -1.672321       |
| NN1.MAE | 0.202259        | 0.072844        | -0.763409       | 0.205092        | 0.073567        | -0.724721       | 0.239051        | 0.096477        | -1.132346       |
| NN2.MSE | 0.239446        | 0.101312        | -1.452556       | 0.206109        | 0.078412        | -0.838305       | 0.228591        | 0.095126        | -1.102488       |
| NN2.MAE | 0.19141         | 0.068261        | -0.652455       | 0.184095        | 0.062366        | -0.462125       | 0.220087        | 0.086888        | -0.920403       |
| NN3.MSE | 0.193117        | 0.069206        | -0.675336       | 0.193859        | 0.070747        | -0.658609       | 0.205093        | 0.076497        | -0.690745       |
| NN3.MAE | 0.191596        | 0.066926        | -0.620138       | 0.176555        | 0.060022        | -0.407183       | 0.234768        | 0.091003        | -1.011359       |
| NN4.MSE | 0.191361        | 0.07068         | -0.71101        | 0.175311        | 0.059253        | -0.389136       | 0.18148         | 0.061718        | -0.364096       |
| NN4.MAE | 0.139659        | 0.041096        | 0.005158        | 0.179318        | 0.05976         | -0.401027       | 0.188921        | 0.066144        | -0.461932       |
| NN5.MSE | 0.17209         | 0.056982        | -0.379418       | 0.164756        | 0.054398        | -0.275325       | 0.202012        | 0.074051        | -0.636691       |
| NN5.MAE | 0.170945        | 0.056029        | -0.356356       | 0.180669        | 0.059697        | -0.399552       | 0.189149        | 0.065921        | -0.456988       |

Figure 18: Missing Data Threshold Robustness Check Individual Factor Importance

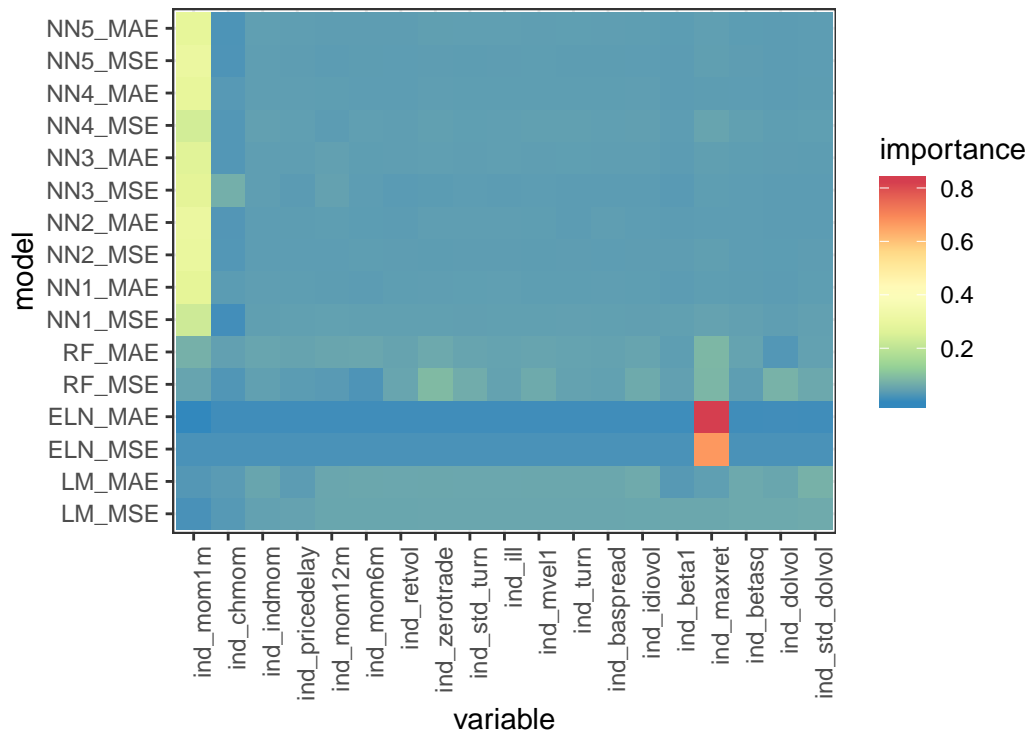


Figure 19: Missing Data Threshold Robustness Check Macroeconomic Factor Importance

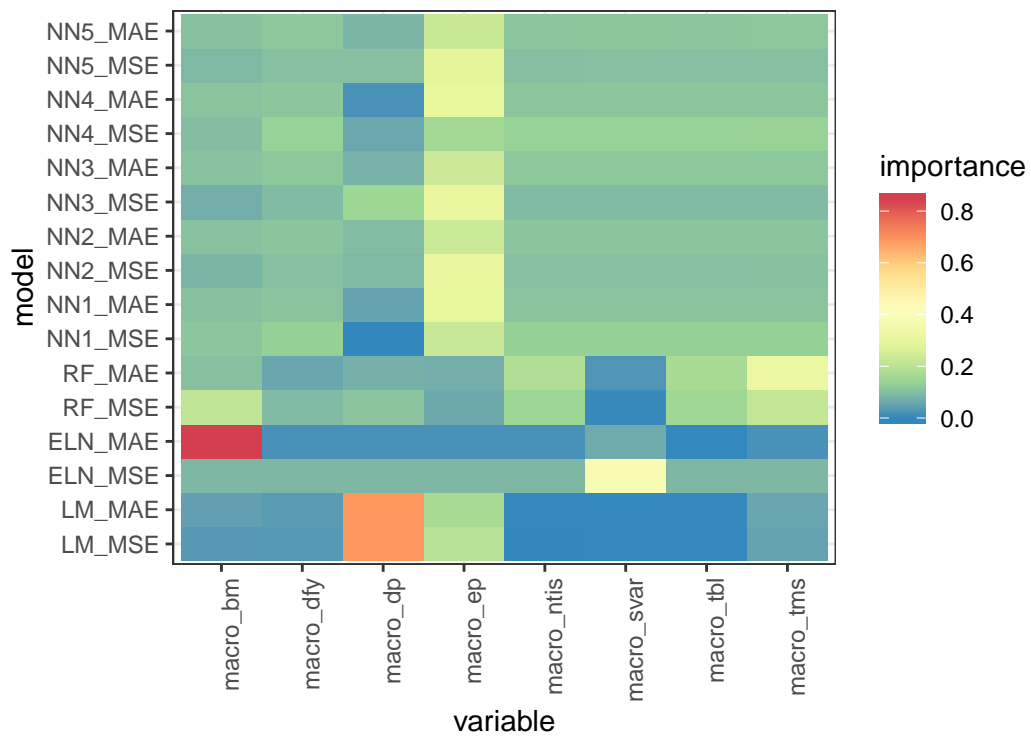


Figure 20: Missing Data Threshold Robustness Check RF VIMP

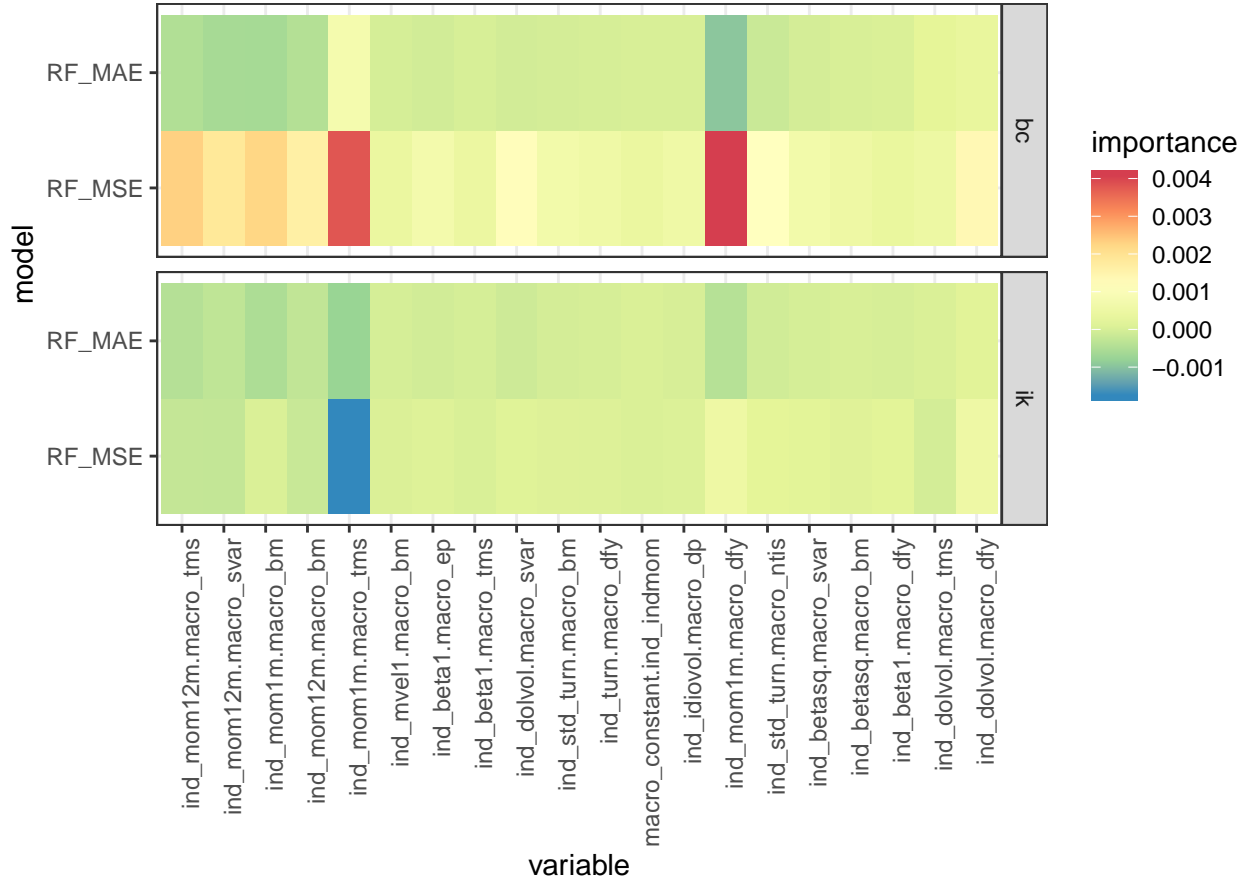


Table 9: Train:Validation 2:1 Robustness Check Loss Statistics

| model   | Sample 1        |                 |                 | Sample 2        |                 |                 | Sample 3        |                 |                 |
|---------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|
|         | Test MAE        | Test MSE        | Test $R^2$      | Test MAE        | Test MSE        | Test $R^2$      | Test MAE        | Test MSE        | Test $R^2$      |
| LM.MSE  | 0.277087        | 0.164599        | -2.98459        | 0.383421        | 0.31299         | -6.337839       | 0.523418        | 0.740288        | -15.361936      |
| LM.MAE  | 0.246936        | 0.147979        | -2.582262       | 0.277044        | 0.161215        | -2.779579       | 0.487285        | 0.631575        | -12.95915       |
| ELN.MSE | 0.133715        | 0.039919        | 0.033647        | 0.139723        | 0.042525        | 0.003028        | <b>0.145034</b> | <b>0.044306</b> | <b>0.020752</b> |
| ELN.MAE | 0.131237        | 0.040361        | 0.022952        | <b>0.137205</b> | <b>0.041858</b> | <b>0.018674</b> | 0.174408        | 0.064513        | -0.425873       |
| RF.MSE  | 0.130808        | 0.036982        | 0.104754        | 0.162762        | 0.051118        | -0.198417       | 0.155264        | 0.048661        | -0.075516       |
| RF.MAE  | <b>0.127013</b> | <b>0.036722</b> | <b>0.111033</b> | 0.146758        | 0.043961        | -0.030633       | 0.168905        | 0.055983        | -0.237348       |
| NN1.MSE | 0.155088        | 0.050284        | -0.217281       | 0.165871        | 0.053459        | -0.253309       | 0.181984        | 0.064621        | -0.428262       |
| NN1.MAE | 0.159797        | 0.050566        | -0.224107       | 0.163397        | 0.052329        | -0.226828       | 0.181636        | 0.062407        | -0.379326       |
| NN2.MSE | 0.155815        | 0.050954        | -0.233492       | 0.168576        | 0.055738        | -0.306745       | 0.170991        | 0.057453        | -0.269824       |
| NN2.MAE | 0.148149        | 0.047617        | -0.152709       | 0.166334        | 0.054058        | -0.26734        | 0.163141        | 0.052639        | -0.163436       |
| NN3.MSE | 0.154141        | 0.04976         | -0.204586       | 0.166218        | 0.053402        | -0.251967       | 0.169539        | 0.05661         | -0.251204       |
| NN3.MAE | 0.142464        | 0.043771        | -0.059594       | 0.154233        | 0.048682        | -0.141321       | 0.184217        | 0.064175        | -0.418401       |
| NN4.MSE | 0.166547        | 0.056184        | -0.360092       | 0.150748        | 0.047566        | -0.115162       | 0.168447        | 0.056575        | -0.250437       |
| NN4.MAE | 0.150167        | 0.046919        | -0.135802       | 0.16197         | 0.05226         | -0.225199       | 0.171676        | 0.057352        | -0.267598       |
| NN5.MSE | 0.155784        | 0.052258        | -0.265047       | 0.139699        | 0.043082        | -0.010018       | 0.166166        | 0.055027        | -0.216219       |
| NN5.MAE | 0.161161        | 0.053216        | -0.28825        | 0.149207        | 0.046344        | -0.086511       | 0.149424        | 0.047544        | -0.050824       |

Figure 21: Train:Validation = 1:1 Robustness Check Individual Factor Importance

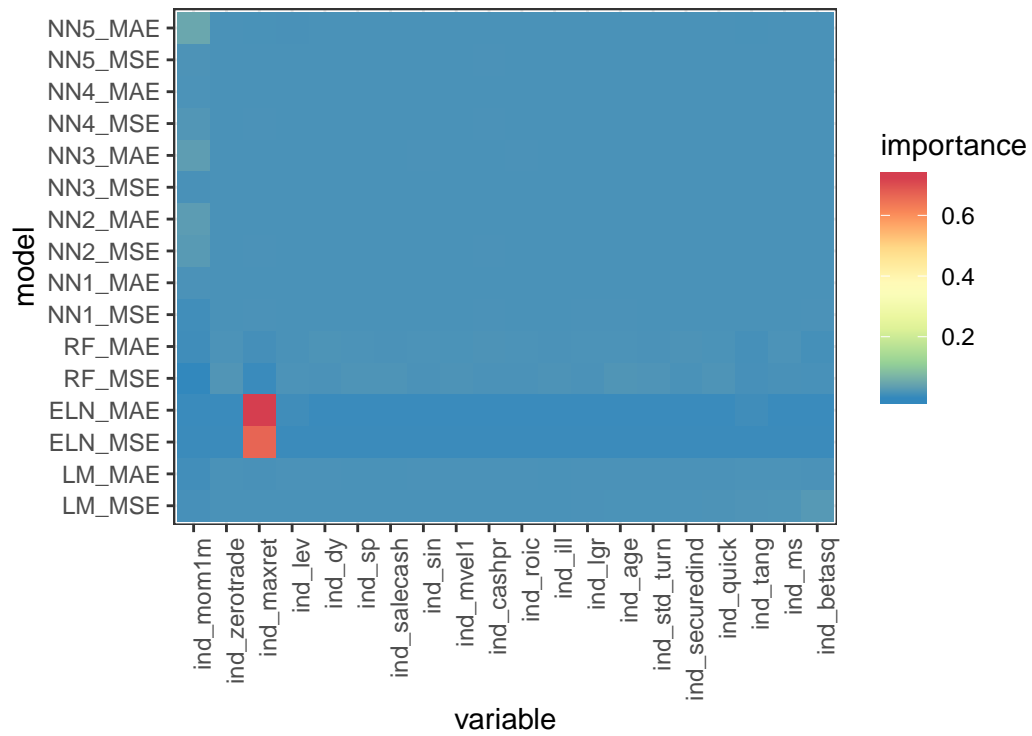


Figure 22: Train:Validation = 1:1 Robustness Check Macroeconomic Factor Importance

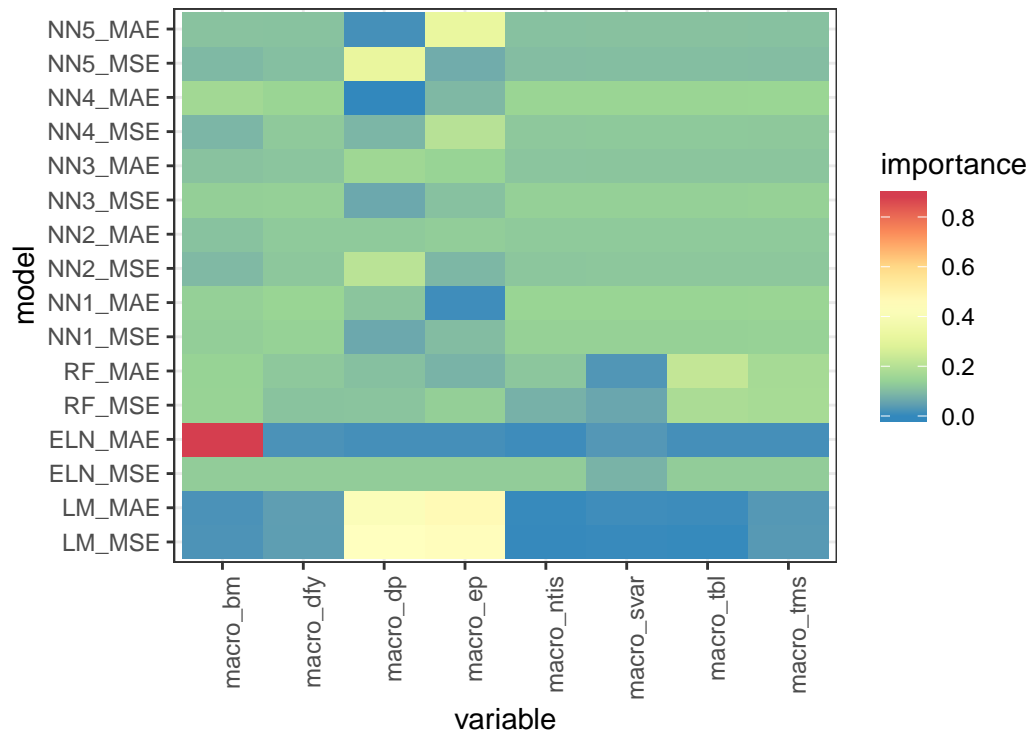


Figure 23: Train:Validation = 1:1 Robustness Check RF VIMP

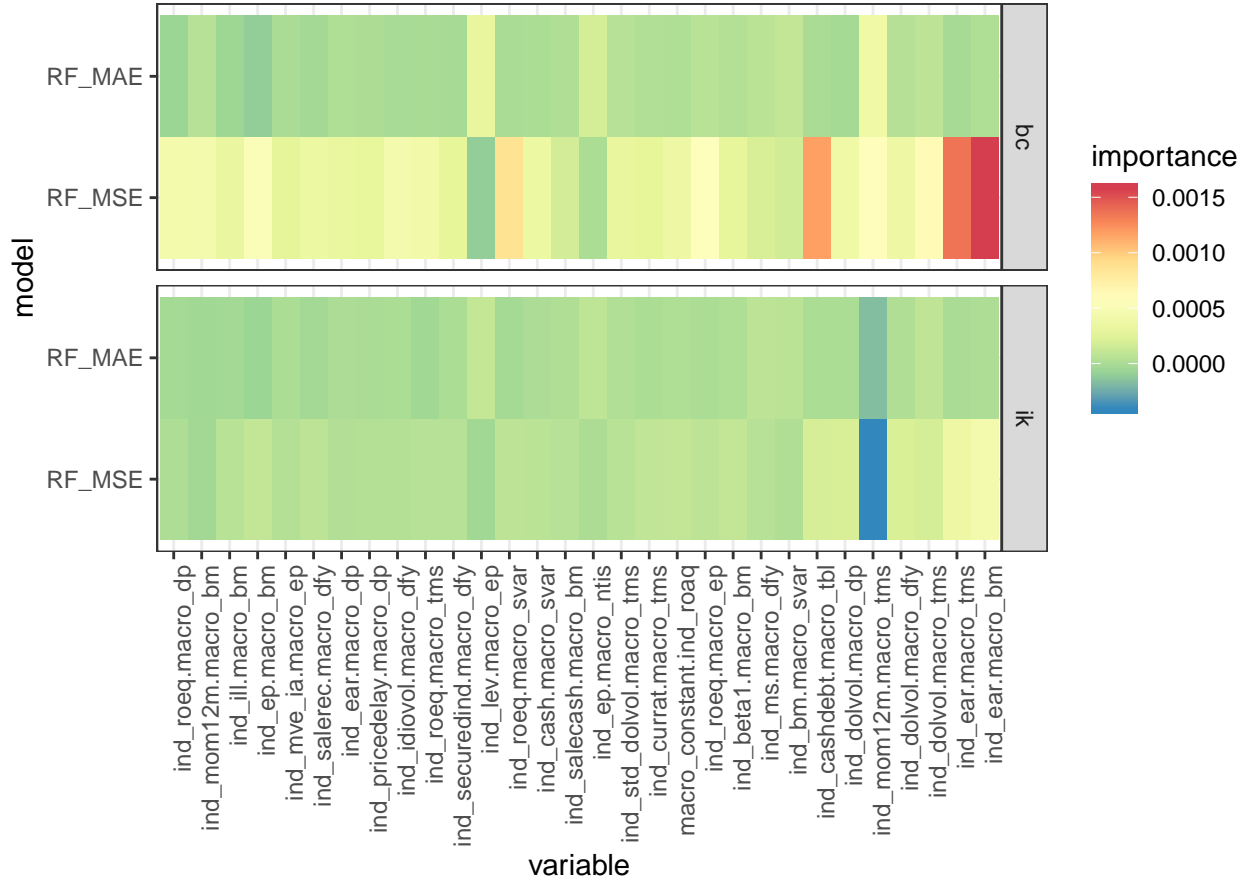


Table 10: Fama French Factor Robustness Check Loss Statistics

| model   | Sample 1        |                 |                 | Sample 2       |                 |                 | Sample 3       |                 |                 |
|---------|-----------------|-----------------|-----------------|----------------|-----------------|-----------------|----------------|-----------------|-----------------|
|         | Test MAE        | Test MSE        | Test $R^2$      | Test MAE       | Test MSE        | Test $R^2$      | Test MAE       | Test MSE        | Test $R^2$      |
| LM.MSE  | 0.288636        | 0.182966        | -3.42923        | 0.367636       | 0.264918        | -5.210825       | 1.101604       | 5.012469        | -109.78624      |
| LM.MAE  | 0.280535        | 0.179777        | -3.352038       | 0.376163       | 0.279476        | -5.552114       | 1.25341        | 7.06036         | -155.048996     |
| ELN.MSE | 0.13383         | 0.039956        | 0.032746        | 0.14022        | 0.0427          | -0.00107        | <b>0.14472</b> | <b>0.043852</b> | <b>0.030769</b> |
| ELN.MAE | <b>0.128936</b> | <b>0.039665</b> | <b>0.039798</b> | <b>0.13716</b> | <b>0.042144</b> | <b>0.011965</b> | 0.172148       | 0.063154        | -0.395841       |
| RF.MSE  | 0.146318        | 0.042607        | -0.031434       | 0.151137       | 0.047091        | -0.104011       | 0.177125       | 0.064664        | -0.429221       |
| RF.MAE  | 0.138266        | 0.04005         | 0.030475        | 0.138714       | 0.042246        | 0.009583        | 0.152068       | 0.048488        | -0.071698       |
| NN1.MSE | 0.168063        | 0.055354        | -0.340017       | 0.192143       | 0.068904        | -0.61541        | 0.275195       | 0.138165        | -2.053731       |
| NN1.MAE | 0.161596        | 0.051507        | -0.246873       | 0.199416       | 0.068181        | -0.598444       | 0.23054        | 0.093434        | -1.065082       |
| NN2.MSE | 0.169842        | 0.056899        | -0.377415       | 0.179733       | 0.058966        | -0.382416       | 0.252929       | 0.117102        | -1.588199       |
| NN2.MAE | 0.155816        | 0.046809        | -0.133147       | 0.185008       | 0.060854        | -0.426679       | 0.219342       | 0.085115        | -0.881213       |
| NN3.MSE | 0.1621          | 0.053165        | -0.287008       | 0.182996       | 0.059643        | -0.398278       | 0.232226       | 0.099353        | -1.195903       |
| NN3.MAE | 0.161255        | 0.050737        | -0.228237       | 0.191625       | 0.064676        | -0.516291       | 0.218355       | 0.085297        | -0.885238       |
| NN4.MSE | 0.166036        | 0.055575        | -0.345349       | 0.191589       | 0.066207        | -0.552182       | 0.23417        | 0.097348        | -1.151607       |
| NN4.MAE | 0.148375        | 0.045227        | -0.094843       | 0.168623       | 0.054176        | -0.270114       | 0.20837        | 0.077667        | -0.7166         |
| NN5.MSE | 0.147379        | 0.044503        | -0.077315       | 0.166006       | 0.054935        | -0.287914       | 0.20667        | 0.077866        | -0.721013       |
| NN5.MAE | 0.150541        | 0.045723        | -0.106868       | 0.172466       | 0.055402        | -0.298865       | 0.218796       | 0.084938        | -0.877301       |



Figure 24: Train:Validation = 2:1 Robustness Check Individual Factor Importance

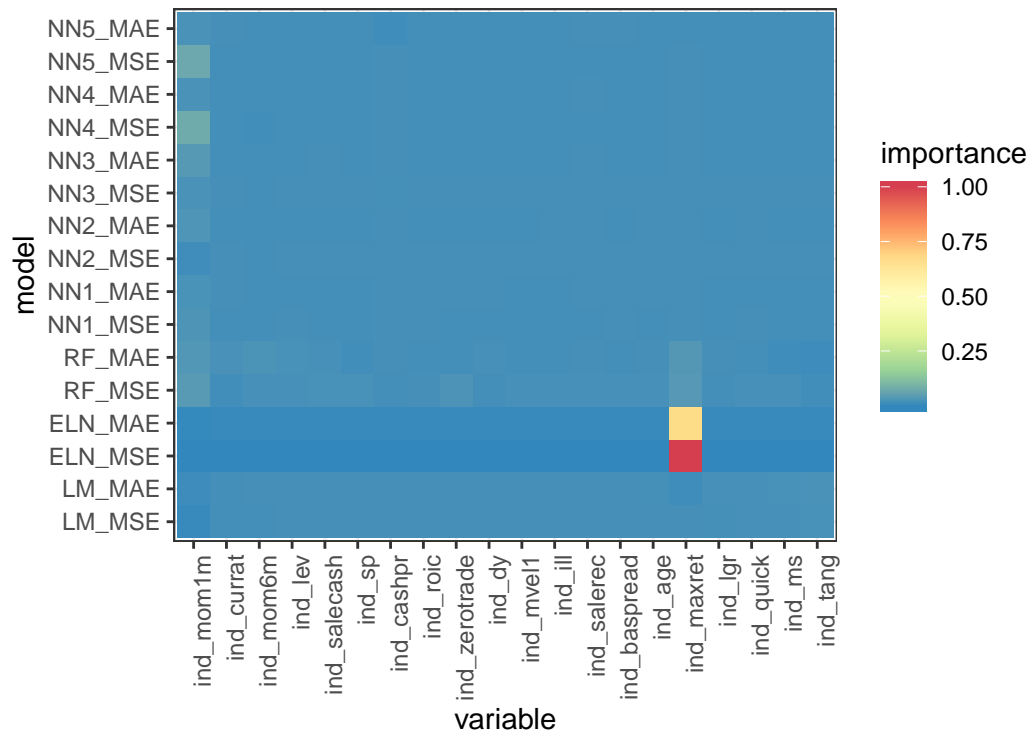


Figure 25: Train:Validation = 2:1 Robustness Check Macroeconomic Factor Importance

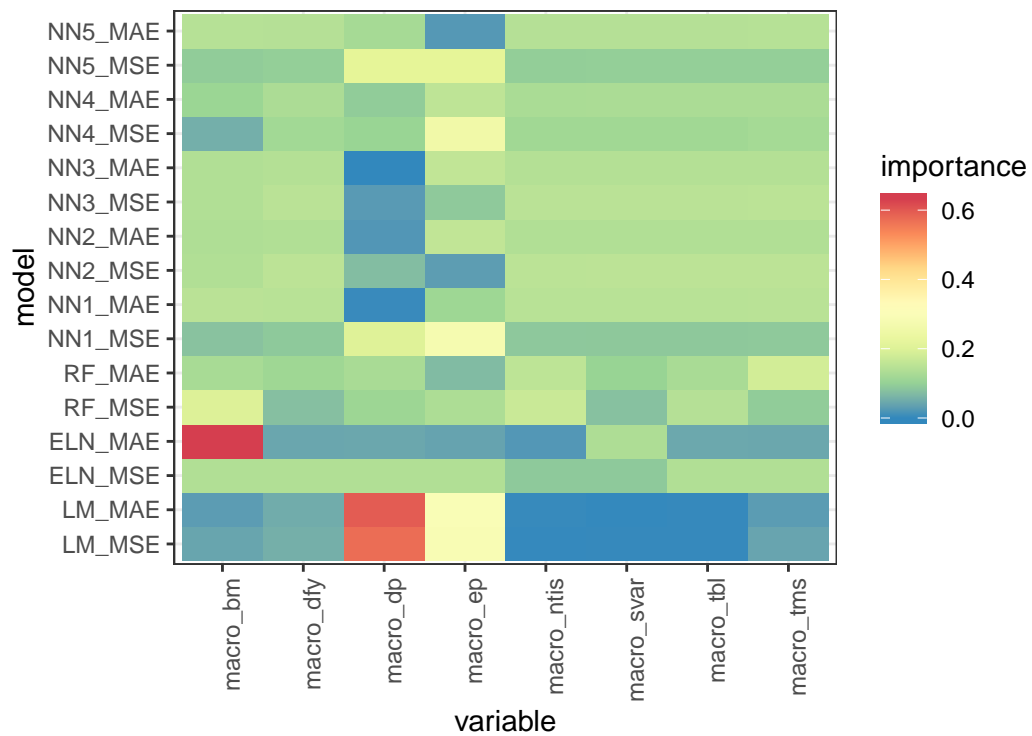


Figure 26: Train:Validation = 2:1 Robustness Check RF VIMP

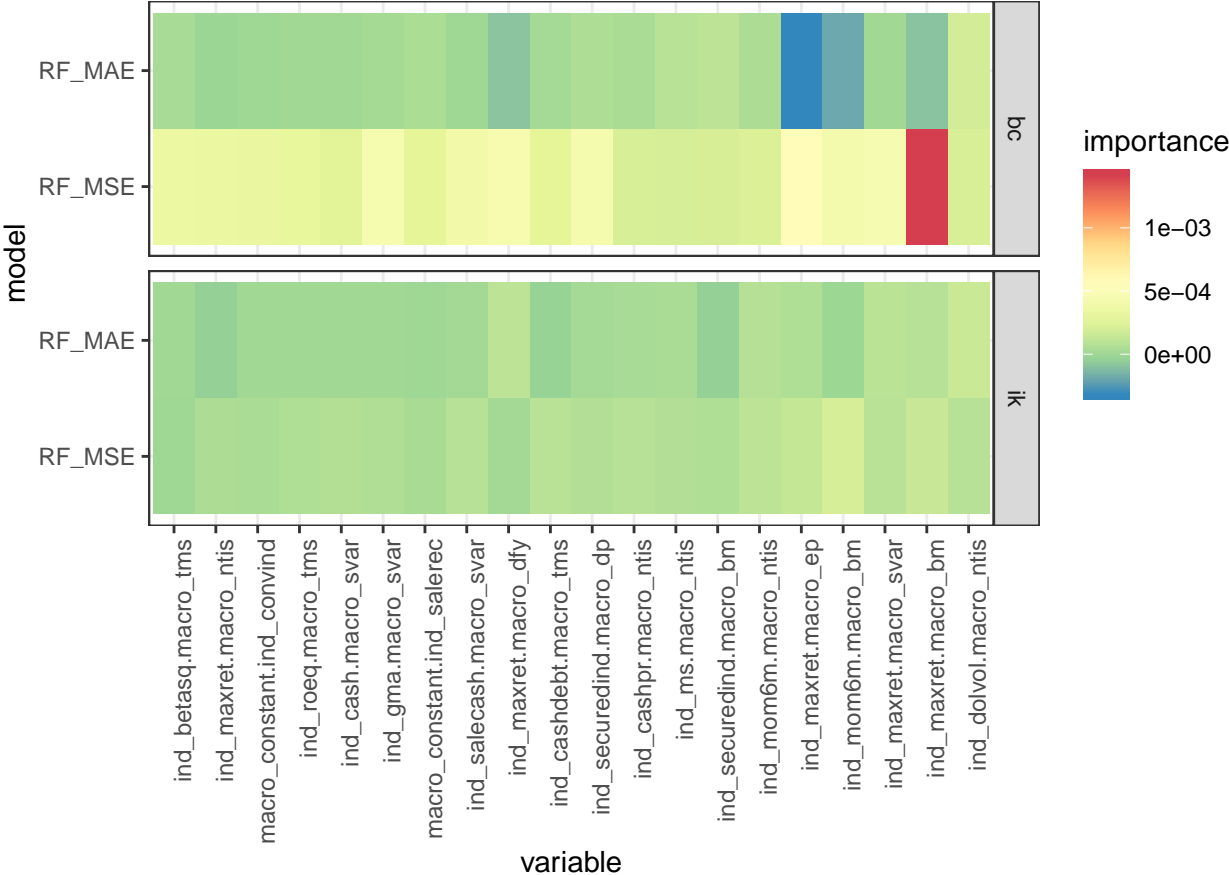


Figure 27: Fama French Factors Robustness Check Individual Factor Importance

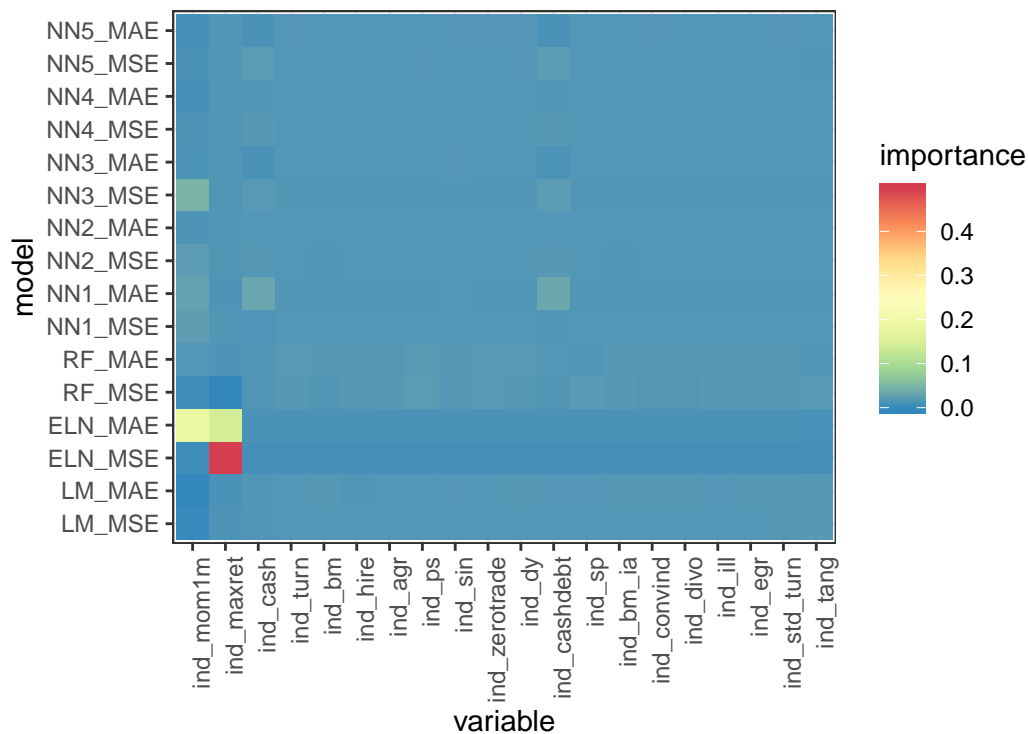


Figure 28: Fama French Factors Robustness Check Macroeconomic Factor Importance

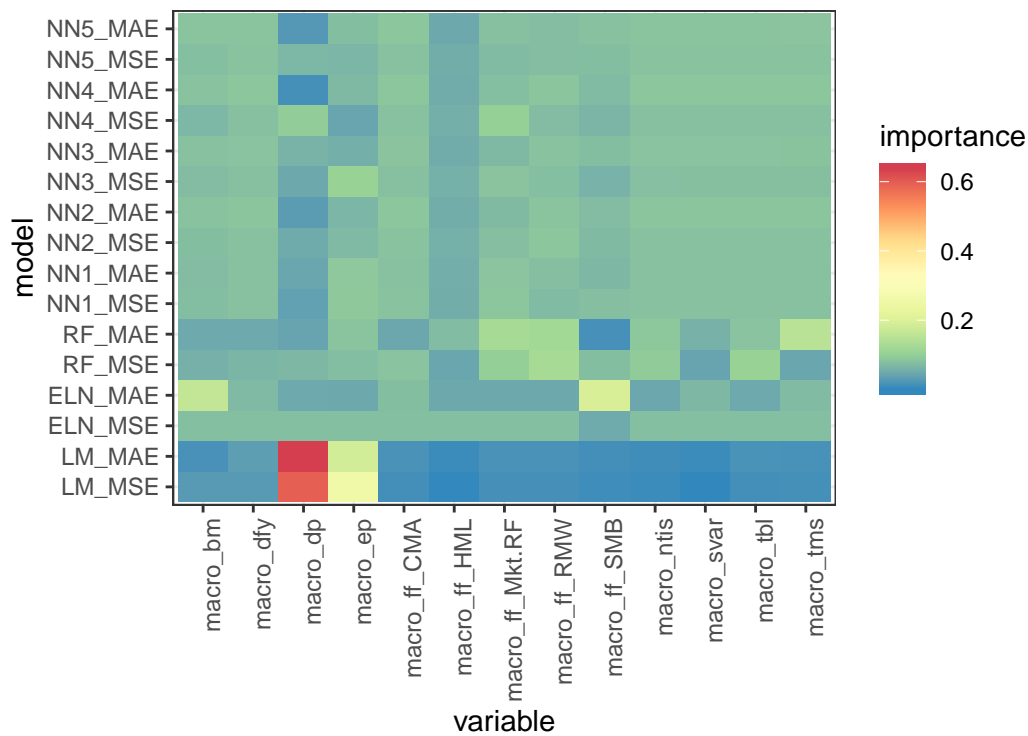


Figure 29: Fama French Factors Robustness Check RF VIMP

