

---

# 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       forests 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:

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

$$B := \Lambda\Lambda' + 0.1\mathbb{I}_n, \quad \Lambda_i = (\lambda_{i1}, \dots, \lambda_{i4}), \quad \lambda_{ik} \sim N(0, \lambda_{sd}), \quad k = 1, \dots, 4 \quad (11)$$

where  $B$  serves as a variance covariance matrix with  $\lambda_{sd}$  its density, and  $L$  represents the lower triangle matrix of  $B$  via the Cholesky decomposition.  $\lambda_{sd}$  values of 0.01, 0.1 and 1 were used to explore increasing degrees of cross sectional correlation. Characteristics are then normalized to be within  $[-1, 1]$  for each  $1 \leq i \leq N$  and for  $j = 1, \dots, P_c$  via:

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

**Simulating macroeconomic series** We consider a Vector Autoregression (VAR) model for  $x_t$ , a  $3 \times 1$  multivariate time series<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, with the number of simulations 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 for  $A$  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** 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 typically had better performance when fitted with respect to quantile loss. Random forests 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, which are severely affected by both increasing non-linearities cross sectional correlation. This result is consistent across all loss 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.

Table 1: Top Models 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
-	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.

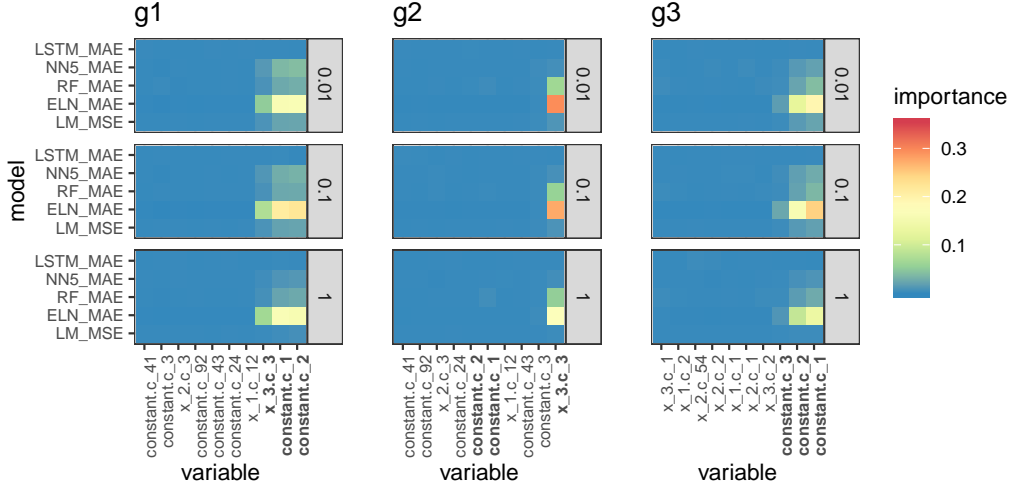
Elastic net models perform the best at identifying the true data generating regressors, and this appears to be mostly robust regardless of cross sectional correlation, though their performance worsens as the data generating process becomes more non-linear. On more difficult specifications, the elastic 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, assigned them low relative importance.

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 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.

<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 variable importance, faceted by simulation specification



## 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.

### 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 is the same data used in [?]. We restrict our dataset to begin from 1993 Q3 and end on 2016 Q4 to alleviate data quality issues. Our individual factor set contains 94 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:

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

We allow all stocks which have a quarterly return to enter the dataset, even if they disappear from the dataset for certain periods. 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

<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: Top 5 models in empirical study

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 $R^2$
ELN.MAE	0.131369	0.040718	0.014306	<b>0.137092</b>	<b>0.041892</b>	<b>0.017875</b>	<b>0.146251</b>	<b>0.045207</b>	<b>0.000835</b>
RF.MAE	<b>0.126703</b>	<b>0.036785</b>	<b>0.109505</b>	0.173721	0.057546	-0.349132	0.14692	0.046037	-0.01752
NN5.MAE	0.146411	0.044901	-0.086967	0.18499	0.06461	-0.514744	0.184986	0.063861	-0.411475
NN4.MAE	0.157301	0.050286	-0.217308	0.168815	0.055711	-0.306102	0.167998	0.055129	-0.218463
NN3.MAE	0.140781	0.042832	-0.036882	0.181096	0.06216	-0.4573	0.164896	0.053458	-0.181528

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:

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

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  $61 \times (8 + 1) = 549^8$ . The final dataset contains 202, 066 individual observations. We note that due to data quality issues, LSTMs, FFORMA and DeepAR are not feasible on empirical data, though the results of the simulation study suggest that even if were to be used, their performance would be underwhelming.<sup>9</sup>

We mimic the sample splitting procedure used in the simulation study: 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 general the results of the simulation study were repeated: the elastic net models perform the best, followed by the random forests, then the DFNs, 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 perform better, though this result is less apparent due to the lower robustness on empirical data (see ?? in Appendix for results).

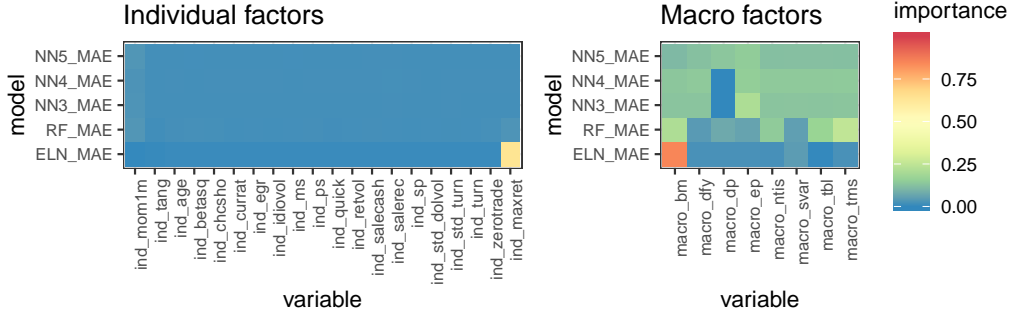
<sup>8</sup>As the individual and macroeconomic factors can have similar names, individual and macroeconomic factors were prefixed with ind\_ and macro\_ respectively.

<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.



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

Figure 2: Empirical individual and macroeconomic factor importance, averaged over all samples



Individual factors shown on x axis (see Table ?? in Appendix for definitions)

The elastic net, random forest and to a lesser extent DFNs tend to pick out the max return and 1 month momentum factors out of the individual characteristics as important, and the book-to-market factor out of the macroeconomic factors are important. In general, the variable importance metrics are less consistent for the random forests, and it should be noted in particular that the random forest tends to determine factors highly correlated with momentum as important, such as change in momentum, dollar trading volume and return volatility. Within the macroeconomic factors, penalized linear models tend to identify the average book to market ratio and the default spread as the most important. The random forests were inconsistent with the elastic nets, and tended to assign very similar variable importance metrics to most macroeconomic factors.

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

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

## 5 Conclusion

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

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

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

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

## 342 5.1 Retrieval of style files

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

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

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

347 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>.  
348 **Previous style files for L<sup>A</sup>T<sub>E</sub>X 2.09, Microsoft Word, and RTF are no longer supported!**

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

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

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

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

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

## 364 6 General formatting instructions

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

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

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

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

## 379 7 Headings: first level

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

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

### 382 7.1 Headings: second level

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

### 384 7.1.1 Headings: third level

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

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

## 388 8 Citations, figures, tables, references

389 These instructions apply to everyone.

### 390 8.1 Citations within the text

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

394 The documentation for `natbib` may be found at

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

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

398 `\citet{hasselmo}` investigated\dots

399 produces

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

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

403 `\PassOptionsToPackage{options}{natbib}`

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

406 `\usepackage[nonatbib]{neurips_2020}`

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

### 411 8.2 Footnotes

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

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

### 416 8.3 Figures

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

---

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

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

Figure 3: Sample figure caption.

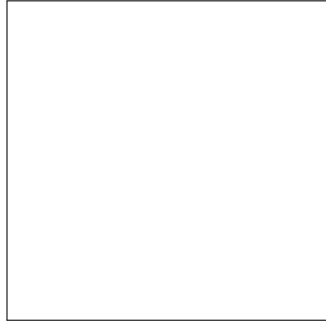


Table 3: Sample table title

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

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

#### 423 8.4 Tables

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

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

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

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

432 This package was used to typeset Table 3.

### 433 9 Final instructions

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

### 437 10 Preparing PDF files

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

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

- 441 • You should directly generate PDF files using `pdflatex`.
- 442 • You can check which fonts a PDF files uses. In Acrobat Reader, select the menu  
443 Files>Document Properties>Fonts and select Show All Fonts. You can also use the program  
444 `pdffonts` which comes with `xpdf` and is available out-of-the-box on most Linux machines.

445 • The IEEE has recommendations for generating PDF files whose fonts are also ac-  
 446 ceptable for NeurIPS. Please see [http://www.emfield.org/icuwb2010/downloads/](http://www.emfield.org/icuwb2010/downloads/IEEE-PDF-SpecV32.pdf)  
 447 [IEEE-PDF-SpecV32.pdf](http://www.emfield.org/icuwb2010/downloads/IEEE-PDF-SpecV32.pdf)

448 • xfig "patterned" shapes are implemented with bitmap fonts. Use "solid" shapes instead.

449 • The `\bbold` package almost always uses bitmap fonts. You should use the equivalent AMS  
 450 Fonts:

451 `\usepackage{amsfonts}`

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

454 `\newcommand{\RR}{\mathbb{R}} %real numbers`

455 `\newcommand{\Nat}{\mathbb{N}} %natural numbers`

456 `\newcommand{\CC}{\mathbb{C}} %complex numbers`

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

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

## 459 10.1 Margins in L<sup>A</sup>T<sub>E</sub>X

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

463 `\usepackage[pdftex]{graphicx} ...`

464 `\includegraphics[width=0.8\linewidth]{myfile.pdf}`

465 See Section 4.4 in the graphics bundle documentation ([http://mirrors.ctan.org/macros/](http://mirrors.ctan.org/macros/latex/required/graphics/grfguide.pdf)  
 466 [latex/required/graphics/grfguide.pdf](http://mirrors.ctan.org/macros/latex/required/graphics/grfguide.pdf))

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

## 469 Broader Impact

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

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

## 478 References

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

483 [1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In  
 484 G. Tesauero, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp.  
 485 609–616. Cambridge, MA: MIT Press.

486 [2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the*  
 487 *General NEural Simulation System*. New York: TELOS/Springer-Verlag.

488 [3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent  
 489 synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.

## A Additional details: models

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

### A.1 Linear models

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

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

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

### A.2 Elastic nets

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

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

where the elastic net penalty [?] is:

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

where the hyperparameter  $\lambda$  controls the overall magnitude of the loss, and hyperparameter  $\rho$  controls the shape of the penalization.  $\rho = 1$  corresponds to ridge regression proposed by [?], which uses  $l_2$  penalty that shrinks all coefficients closer to 0, but not to 0.  $\rho = 0$  case corresponds to the popular LASSO and uses absolute ( $l_1$ ) parameter penalization proposed by [?], which geometrically allows the coefficients to be shrunk to 0. For  $0 < \rho < 1$ , the elastic net aims to produce parsimonious models through both shrinkage and selection by combining the properties of LASSO and ridge regression.

### A.3 Random forests

Random Forests are an extension of Classification and Regression Trees (CART) proposed by [?]. CART are fully non-parametric models that can capture complex multi-way interactions. A tree "grows" in a series of iterations. With each iteration, a split ("branch") is made along one predictor such that it is the best split available at that stage with respect to minimizing the loss function. These steps are continued until each observation is its own node, or more commonly until the stopping criterion is met. The eventual model slices the predictor space into rectangular partitions, and predicts the unknown function  $g^*(z_{i,t})$  with the average value of the outcome variable in each partition. The prediction of a tree,  $\mathcal{T}$ , with  $K$  "leaves" (terminal nodes), and depth  $L$  is

$$g(z_{i,t}; \theta, K, L) = \sum_{k=1}^K \theta_k \mathbf{1}_{z_{i,t} \in C_k(L)} \quad (21)$$

where  $C_k(L)$  is one of the  $K$  partitions in the model. For this study, only recursive binary trees were considered. Though trees were originally proposed and fit with respect to minimizing MSE, they can be grown with respect to a variety of loss functions, including MAE, where the loss within each C partition is denoted by  $H(\theta, C)$ :

$$H(\theta, C) = \frac{1}{|C|} \sum_{z_{i,t} \in C} L(r_{i,t+1} - \theta) \quad (22)$$

where  $|C|$  denotes the number of observations in set C (partition). Given  $C$ , it is clear that the optimal choice for minimising MSE is simply  $\theta = \frac{1}{|C|} \sum_{z_{i,t} \in C} r_{i,t+1}$  i.e. the average of the partition, and the median of the partition when minimising MAE.

Trees, grown to a deep enough level, are highly unbiased and flexible, as each partition can potentially predict a single, or low number of observations. The trade-off is their high variance and instability.

Further details are given in cite().

### A.4 Feed forward neural networks

More specifically, a neural network consists of layers denoted by  $l = 0, 1, \dots, L$ , with  $l = 0$  denoting the input layer and  $l = L$  denoting the output layer. The input layer is defined by the scaled predictor set,  $x^{(0)} = (1, z_1, \dots, z_N)'$ . The model adds complexity through the use of one or more hidden layer, each



containing  $K^{(l)}$  "neurons". Each neuron linearly aggregates the values of the previous layer, and applies some non-linear "activation function" which we denote as  $\alpha$  to its aggregated signal before sending its output to the next layer. The output of neuron  $k$  in layer  $l$  is then  $x_k^{(l)}$ . Next, define the vector of outputs for this layer as  $x^{(l)} = (1, x_1^{(l)}, \dots, x_{K^{(l)}}^{(l)})'$ . The recursive output formula for the neural network at each neuron in layer  $l > 0$  is then:

$$x_k^{(l)} = \alpha(x^{(l-1)'} \theta_k^{l-1}), \quad (23)$$

where  $\alpha(\cdot)$  represents the activation function for that layer with the final output<sup>12</sup>

$$g(z; \theta) = x^{(L-1)'} \theta^{L-1} \quad (24)$$

The neural network's weight and bias parameters for each layer are estimated by minimizing the loss function with respect to the parameters, i.e. by calculating the partial derivative with respect to a specific weight or bias element.

Due to the complexity and hence non-existent analytical form for this solution, this is typically found via backpropagation, an algorithm which exploits the chain rule of the partial derivative and iteratively finds a local optimum using a first order gradient based algorithm, also known as "gradient descent." The gradient descent algorithm minimizes some function (such as the loss function in the context of machine learning) by iteratively moving in the direction of steepest descent, defined as the negative gradient. Formally, for a loss function  $L(x)$  that is defined and has a gradient defined in the neighbourhood of the parameter set  $a$ , the updating algorithm is:

$$a_{n+1} = a_n - \gamma \Delta F(a_n) \quad (25)$$

where  $\gamma$  controls the size of each update. This  $\gamma$  parameter is known as the learning rate in neural network training, and controlling this is critical for good performance. As the loss functions of neural networks can be very complex with many local minima, the learning rate should be high enough such that the parameter updates are large enough to skip or jump over them. Too large of a learning rate however, and the neural may fail to converge to a solution at all. Due to computational limitations, we tune the learning rate manually, and consider a variety of different "optimizers", or algorithms which adapt the learning rate in different ways (see Appendix for computational details).

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

Further details are given in cite().

## A.5 Long short term memory networks

Long short term memory (LSTM) networks are

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

•

Further details are given in cite().

## A.6 FFORMA

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

Phase one involves training the meta learning model. Let  $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N$  denote the different times series of excess returns,  $F$  denote the set of functions for calculating time series features, and  $M$  the set of constituent models considered for each times series. For each time series  $\mathbf{y}_n$ , split  $\mathbf{y}_n$  into training, validation and test sets, calculate the set of features  $f_n \in F$  over the training period, fit each constituent model  $m \in M$  over the training period and generate forecasts over the validation and test periods, and calculate validation losses  $L_{nm}$ . Finally, train the meta-learning model by minimizing:

$$\underset{w}{\operatorname{argmin}} \sum_{n=1}^N \sum_{m=1}^M w(f_n)_m L_{nm} \quad (26)$$

To incorporate all regressors in each individual time series model, we applied dimensional reduction techniques of PCA and UMAP to generate new feature mappings for use in GARCH (1, 1) models (generally the best performing of the constituent models). None of the new external regressors as generated by these feature mappings improved performance, however.

The constituent models we considered are:

<sup>12</sup>Note that the specification of a constant "1" at the beginning of each layer is the same as specifying a bias term as is popular in other parametrizations.

- 569 • Naive
- 570 • Random walk with drift
- 571 • Theta method
- 572 • ARIMA
- 573 • ETS
- 574 • TBATS
- 575 • Neural network auto-regressive model
- 576 • ARMA (1, 1) with g.e.d. GARCH(1, 1) errors
- 577 • ARMA (1, 1) with g.e.d. GARCH(1, 1) errors and UMAP external regressors

578 We follow cite()'s selection of time series features as inputs to the meta-learner. The time series features  
 579 used to train the meta-model are detailed in cite(), with the addition of realized volatility. Note that because  
 580 financial returns data does not typically exhibit seasonality, features and constituent models related which utilized  
 581 seasonality were omitted.

582 Phase two uses the learning model from phase one to produce new, combined forecasts for each time series.  
 583 For each  $\mathbf{y}_n$ , calculate the features  $f_{new}$  by applying  $F$ , use the meta-learner to produce  $w(f_{new})$  an  $M$  vector  
 584 of weights, compute the forecasts from each constituent model in  $M$ , then finally combine them into a final  
 585 forecast using weights  $w$ .

## 586 A.7 DeepAR

587 DeepAR is a generalization of traditional Auto Regressive (AR) models to include additional layers into order to  
 588 introduce non-linearities into the model.

589 DeepAR aims to model the conditional distribution of the

$$P(\mathbf{z}_{i,t_0:T} | \mathbf{z}_{i,1:t_0-1}, \mathbf{x}_{i,1:T})$$

590 of the future of each time series  $[z_{i,t_0}, z_{i,t_0+1}, \dots, z_{i,T}] := \mathbf{z}_{i,t_0:T}$  given its  
 591 past  $[z_{i,1}, \dots, z_{i,t_0-2}, z_{i,t_0-1}] := \mathbf{z}_{i,1:t_0-1}$ , where  $t_0$  denotes the time point from which we assume  
 592  $z_{i,t}$  to be unknown at prediction time, and  $\mathbf{x}_{i,1:T}$  are covariates that are assumed to be known for all time points.  
 593 To prevent confusion we avoid the ambiguous terms “past” and “future” and will refer to time ranges  $[1, t_0 - 1]$   
 594 and  $[t_0, T]$  as the conditioning range and prediction range, respectively. During training, both ranges have to lie  
 595 in the past so that the  $z_{i,t}$  are observed, but during prediction  $z_{i,t}$  is only available in the conditioning range.  
 596 Note that the time index  $t$  is relative, i.e.  $t = 1$  can correspond to a different actual time period for each  $i$ .

597 Our model, summarized in Fig. ??, is based on an autoregressive recurrent network architecture [?, ?]. We  
 598 assume that our model distribution  $Q_{\Theta}(\mathbf{z}_{i,t_0:T} | \mathbf{z}_{i,1:t_0-1}, \mathbf{x}_{i,1:T})$  consists of a product of likelihood factors

$$Q_{\Theta}(\mathbf{z}_{i,t_0:T} | \mathbf{z}_{i,1:t_0-1}, \mathbf{x}_{i,1:T}) = \prod_{t=t_0}^T Q_{\Theta}(z_{i,t} | \mathbf{z}_{i,1:t-1}, \mathbf{x}_{i,1:T}) = \prod_{t=t_0}^T \ell(z_{i,t} | \theta(\mathbf{h}_{i,t}, \Theta))$$

599 parametrized by the output  $\mathbf{h}_{i,t}$  of an autoregressive recurrent network

$$\mathbf{h}_{i,t} = h(\mathbf{h}_{i,t-1}, z_{i,t-1}, \mathbf{x}_{i,t}, \Theta), \quad (27)$$

600 where  $h$  is a function implemented by a multi-layer recurrent neural network with LSTM cells.<sup>13</sup> The model  
 601 is autoregressive, in the sense that it consumes the observation at the last time step  $z_{i,t-1}$  as an input, as well  
 602 as recurrent, i.e. the previous output of the network  $\mathbf{h}_{i,t-1}$  is fed back as an input at the next time step. The  
 603 likelihood  $\ell(z_{i,t} | \theta(\mathbf{h}_{i,t}))$  is a fixed distribution whose parameters are given by a function  $\theta(\mathbf{h}_{i,t}, \Theta)$  of the  
 604 network output  $\mathbf{h}_{i,t}$  (see below).

605 Information about the observations in the conditioning range  $\mathbf{z}_{i,1:t_0-1}$  is transferred to the prediction range  
 606 through the initial state  $\mathbf{h}_{i,t_0-1}$ . In the sequence-to-sequence setup, this initial state is the output of an *encoder*  
 607 *network*. While in general this encoder network can have a different architecture, in our experiments we opt for  
 608 using the same architecture for the model in the conditioning range and the prediction range (corresponding to  
 609 the *encoder* and *decoder* in a sequence-to-sequence model). Further, we share weights between them, so that  
 610 the initial state for the decoder  $\mathbf{h}_{i,t_0-1}$  is obtained by computing (27) for  $t = 1, \dots, t_0 - 1$ , where all required  
 611 quantities are observed. The initial state of the encoder  $\mathbf{h}_{i,0}$  as well as  $z_{i,0}$  are initialized to zero.

612 Given the model parameters  $\Theta$ , we can directly obtain joint samples  $\tilde{\mathbf{z}}_{i,t_0:T} \sim Q_{\Theta}(\mathbf{z}_{i,t_0:T} | \mathbf{z}_{i,1:t_0-1}, \mathbf{x}_{i,1:T})$   
 613 through ancestral sampling: First, we obtain  $\mathbf{h}_{i,t_0-1}$  by computing (27) for  $t = 1, \dots, t_0$ . For  $t = t_0, t_0 +$   
 614  $1, \dots, T$  we sample  $\tilde{z}_{i,t} \sim \ell(\cdot | \theta(\mathbf{h}_{i,t}, \Theta))$  where  $\mathbf{h}_{i,t} = h(\mathbf{h}_{i,t-1}, \tilde{z}_{i,t-1}, \mathbf{x}_{i,t}, \Theta)$  initialized with  $\mathbf{h}_{i,t_0-1} =$   
 615  $\mathbf{h}_{i,t_0-1}$  and  $\tilde{z}_{i,t_0-1} = z_{i,t_0-1}$ . Samples from the model obtained in this way can then be used to compute  
 616 quantities of interest, e.g. quantiles of the distribution of the sum of values for some time range in the future.

617 Further details are given in cite().

<sup>13</sup>Details of the architecture and hyper-parameters are given in the supplementary material.

## A Additional details: simulation design

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

### A.1 Simulation Design

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  $t = 1, \dots, T$ :

$$r_{i,t+1} = g(z_{i,t}) + \beta_{i,t+1}v_{t+1} + e_{i,t+1}; \quad z_{i,t} = (1, x_t)' \otimes c_{i,t}, \quad \beta_{i,t} = (c_{i1,t}, c_{i2,t}, c_{i3,t}) \quad (28)$$

$$e_{i,t+1} = \sigma_{i,t+1}\varepsilon_{i,t+1}; \quad (29)$$

$$\log(\sigma_{i,t+1}^2) = \omega + \gamma \log(\sigma_t^2) + \sigma_u u; \quad u \sim N(0, 1) \quad (30)$$

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 :

$$\bar{c}_{ij,t} = \rho_j \bar{c}_{ij,t-1} + \epsilon_{ij,t}; \quad \rho_j \sim \mathcal{U}(0.5, 1) \quad (31)$$

We then build in cross sectional correlation:

$$\hat{C}_t = L\bar{C}_t; \quad B = LL' \quad (32)$$

$$B := \Lambda\Lambda' + 0.1\mathbb{I}_n, \quad \Lambda_i = (\lambda_{i1}, \dots, \lambda_{i4}), \quad \lambda_{ik} \sim N(0, \lambda_{sd}), \quad k = 1, \dots, 4 \quad (33)$$

where  $B$  serves as a variance covariance matrix with  $\lambda_{sd}$  its density, and  $L$  represents the lower triangle matrix of  $B$  via the Cholesky decomposition.  $\lambda_{sd}$  values of 0.01, 0.1 and 1 were used to explore increasing degrees of cross sectional correlation. Characteristics are then normalized to be within  $[-1, 1]$  for each  $1 \leq i \leq N$  and for  $j = 1, \dots, P_c$  via:

$$c_{ij,t} = \frac{2}{n+1} \text{rank}(\hat{c}_{ij,t}) - 1. \quad (34)$$

**Simulating macroeconomic series** We consider a Vector Autoregression (VAR) model for  $x_t$ , a  $3 \times 1$  multivariate time series<sup>14</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  $g(z_{i,t})$ :

$$(1) \quad g_1(z_{i,t}) = (c_{i1,t}, c_{i2,t}, c_{i3,t} \times x'_t[3,]) \theta_0 \quad (35)$$

$$(2) \quad g_2(z_{i,t}) = (c_{i1,t}^2, c_{i1,t} \times c_{i2,t}, \text{sgn}(c_{i3,t} \times x'_t[3,])) \theta_0 \quad (36)$$

$$(3) \quad g_3(z_{i,t}) = (1[c_{i3,t} > 0], c_{i2,t}^3, c_{i1,t} \times c_{i2,t} \times 1[c_{i3,t} > 0], \text{logit}(c_{i3,t})) \theta_0 \quad (37)$$

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>15</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, with the number of simulations 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.

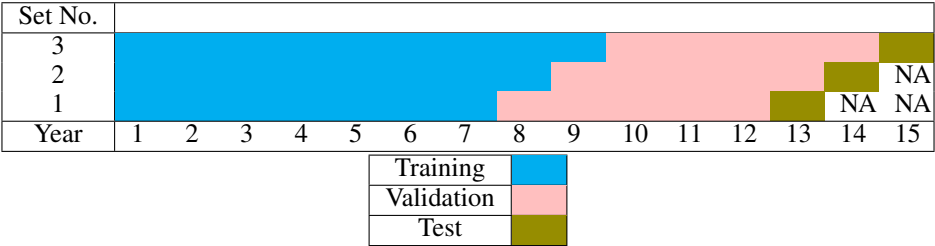
<sup>14</sup>More complex specifications for  $A$  were briefly explored, but these did not have a significant impact on results.

<sup>15</sup>( $g_1, g_2$  correspond to the simulation design used by [?].)

649 **A.1.1 Sample Splitting**

650 If viewed as monthly periods,  $T = 180$  corresponds to 15 years. A data splitting scheme similar to the scheme  
651 to be used in the empirical data study was used: a training:validation length ratio of approximately 1.5 to begin,  
652 and a test set that is 1 year in length. We employ the hybrid growing window approach as described earlier in  
653 section ?? (see Figure 4 for a graphical representation).

Figure 4: Sample Splitting Procedure



654 Other schemes in the forecasting literature such as using an “inner” rolling window validation loop to find  
655 the best hyperparameters on average, finally aggregating them in an “outer” loop for a more robust error were  
656 considered but not implemented due to a) computational feasibility and b) the relative instability of optimal  
657 hyperparameters across different different windows.

## 658 A.2 Simulation Study Results

### 659 A.2.1 Prediction Performance

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 $R^2$
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
	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

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 $R^2$
NN5.MSE	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
	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	0.0025273
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
	0.01	0.0374310	0.0028046	-0.0564056	0.0373372	0.0027801	-0.0518537	0.0376270	0.0028169	-0.0674327
LSTM.MAE	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
	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
FFORMA.MSE	1.00	0.0388747	0.0029647	-0.1312392	0.0388080	0.0029331	-0.1659900	0.0430130	0.0033713	-0.4709541
	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
FFORMA.MAE	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
	0.01	0.0382993	0.0029000	-0.1289295	0.0384895	0.0029121	-0.1325183	0.0393898	0.0030161	-0.2049803

### 660 A.3 Random Forest VIMPs

661 We note that random forest methods typically have their own methodologies to calculate variable importance  
662 which are different to the variable importance metric presented in the main body of the paper. Here we  
663 provide two popular schemes of calculating random forest variable importance metrics - Breiman-cutler VIMP  
664 (traditional) and Ishwaran-Kogalur VIMP, and show that importantly, the overall conclusion regarding factor  
665 selection does not change with respect to which vimp methodology employed.

Figure 5: Simulation Breiman-Cutler vimps

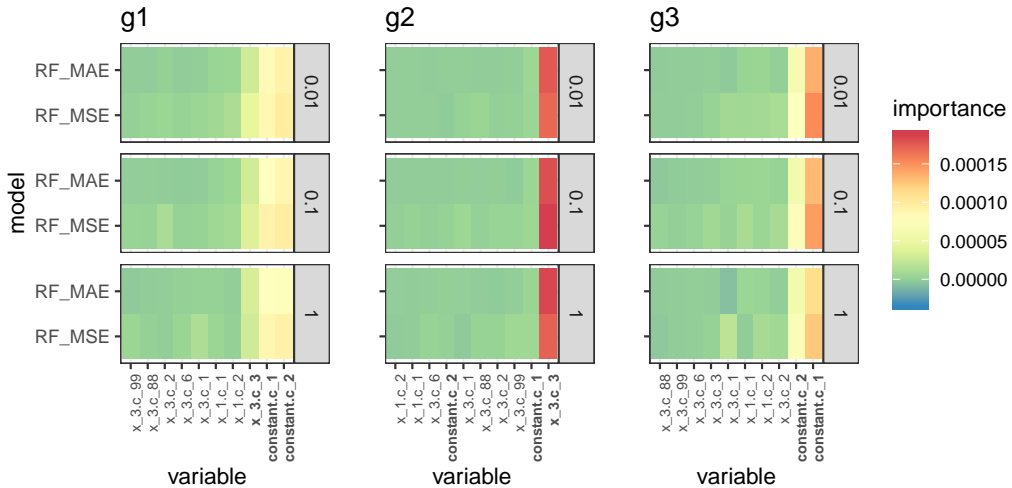
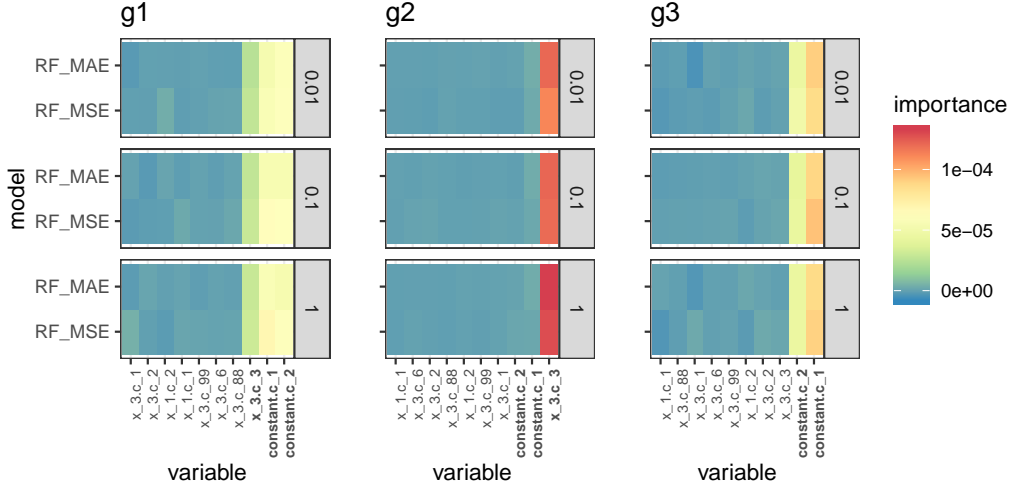


Figure 6: Simulation Ishwaran-Kogalur vimps



## 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. The 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 monthly updated dataset was then converted to a quarterly format, to achieve a balance between having a dataset with enough data points and variability among factors. Quarterly returns were then constructed using the PRC variable according to actual returns:

$$RET_t = (PRC_t - PRC_{t-1}) / PRC_{t-1} \quad (38)$$

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 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, significantly increasing computational feasibility.

There existed a significant amount of missing data in the dataset. For the main empirical study, any characteristics that had over 20% of their data were removed, and remaining missing data points were then imputed with their cross sectional medians. 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



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

701 periods but keeping more characteristics was needed. 1993 Q3 was determined to be a reasonable time frame to  
702 begin the dataset due to a noticeable increase in data quality.

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

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

$$z_{i,t} = (1, x_t)' \otimes c_{i,t} \quad (39)$$

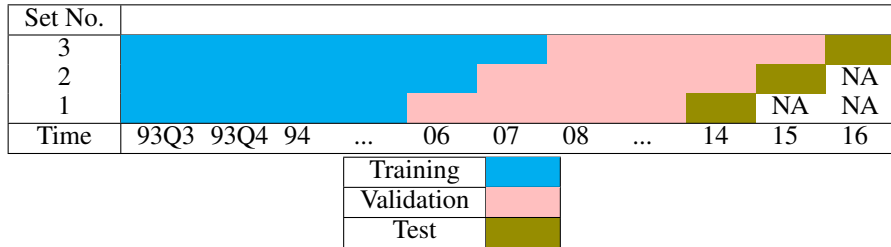
708 where  $c_{i,t}$  is a  $P_c$  matrix of characteristics for each stock  $i$ , and  $(1, x_t)'$  is a  $P_x \times 1$  vector of macroeconomic  
709 predictors, , and  $\otimes$  represents the Kronecker product.  $z_{i,t}$  is therefore a  $P_x P_c$  vector of features for predicting  
710 individual stock returns and includes interactions between stock level characteristics and macroeconomic  
711 variables. The total number of covariates in this baseline set is  $61 \times (8 + 1) = 549$ <sup>16</sup>.

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

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

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

Figure 7: Empirical Data Sample Splitting Procedure



<sup>16</sup>As the individual and macroeconomic factors can have similar names, individual and macroeconomic factors were prefixed with ind\_ and macro\_ respectively.

## 719 A.2 Empirical study results & robustness checks

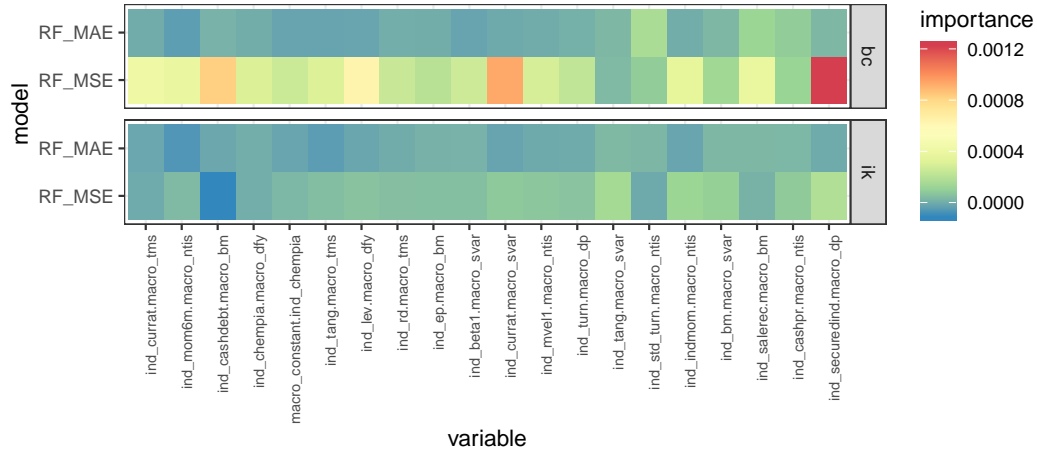
720 In addition to the main study, we provide four additional robustness checks for our empirical study, with regards  
721 to different training/validation splitting schemes, missing data imputation and additional regressors. Importantly,  
722 our overall results are consistent across all checks.

723 **Empirical study results** Here we present the full set of results for our empirical study.

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

Figure 8: Empirical study random forest vimps



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

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

Figure 9: Missing Data Threshold Robustness Check Factor Importance

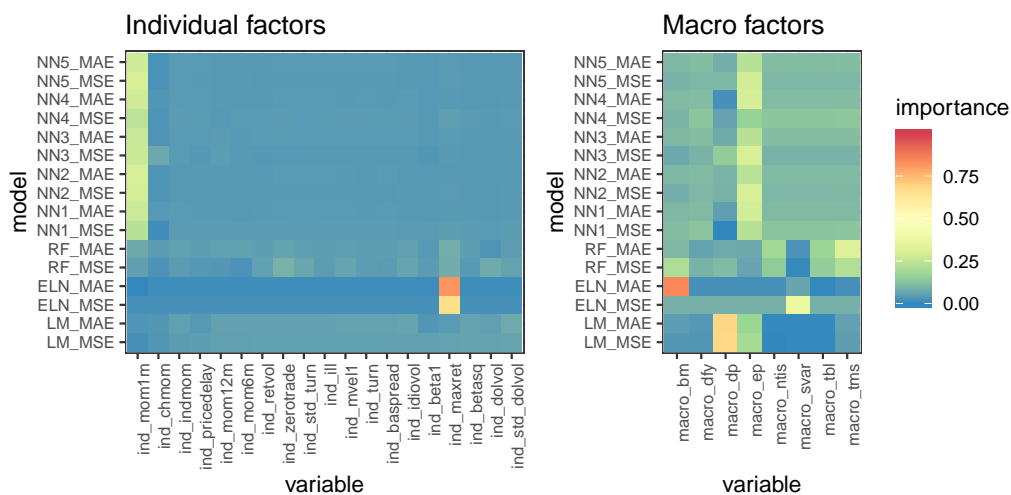
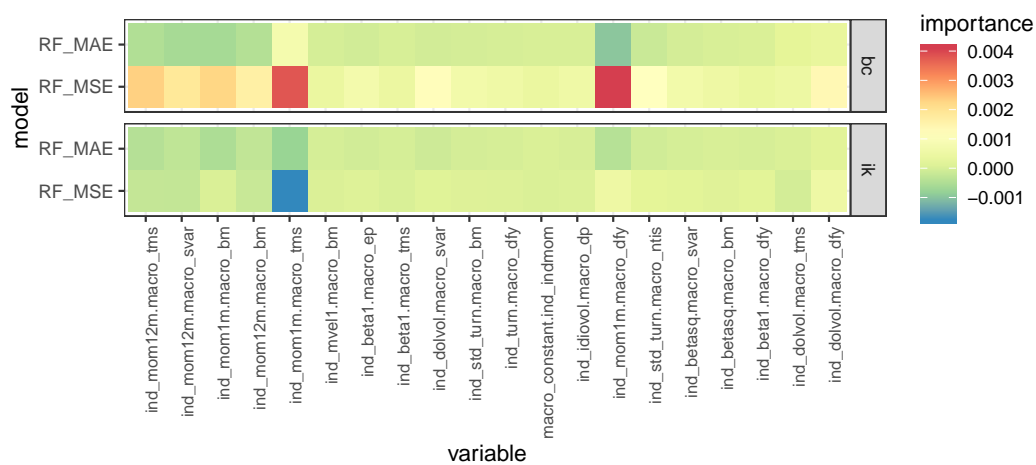


Figure 10: Missing Data Threshold Robustness Check RF VIMP



726 **Train/validation schemes** We consider training:validation length ratios of 1:1 and 1:2 in addition to 1.5:1  
727 in the main study.

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 11: Train:Validation = 1:1 Robustness Check Factor Importance

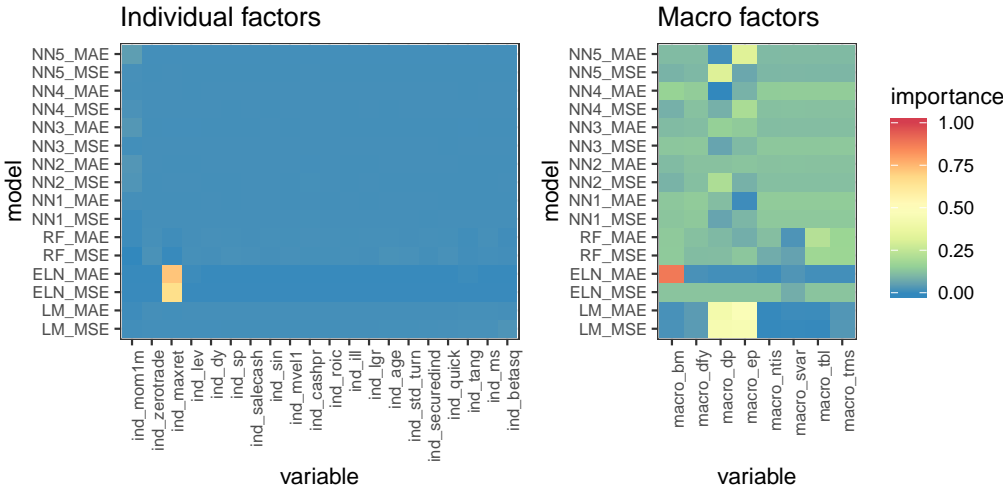


Figure 12: Train:Validation = 1:1 Robustness Check RF VIMP

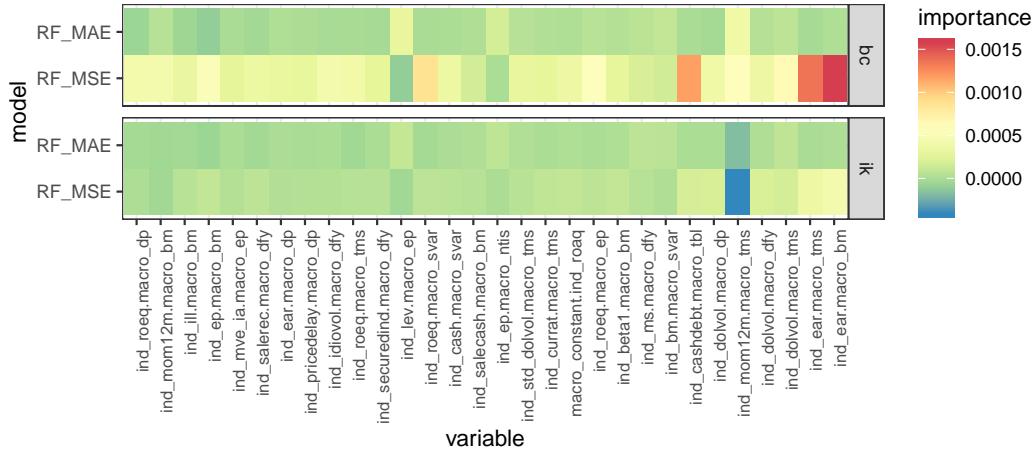


Table 9: Train:Validation 2:1 Robustness Check Loss Statistics

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 13: Train:Validation = 2:1 Robustness Check Factor Importance

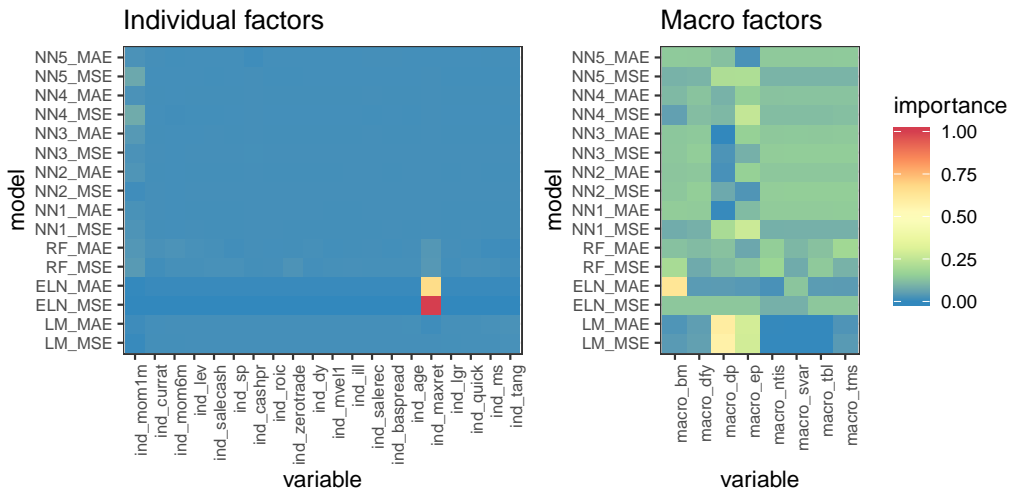


Figure 1: Heatmap of variable importance for RF models. The figure consists of two panels, 'bc' (top) and 'lk' (bottom). Each panel shows a heatmap where rows represent models (RF\_MAE, RF\_MSE) and columns represent variables. The color scale indicates importance, ranging from 0e+00 (blue) to 1e-03 (red). In the 'bc' panel, RF\_MAE shows high importance for 'ind\_maxret\_macro\_ntis' and 'ind\_maxret\_macro\_svar', while RF\_MSE shows high importance for 'ind\_maxret\_macro\_svar'. In the 'lk' panel, RF\_MAE shows high importance for 'ind\_maxret\_macro\_svar' and 'ind\_maxret\_macro\_bm', while RF\_MSE shows high importance for 'ind\_maxret\_macro\_svar' and 'ind\_maxret\_macro\_bm'.

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>0.144472</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	-0.71661
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 1 consists of two heatmaps side-by-side, titled "Individual factors" and "Macro factors". Both heatmaps share the same y-axis, labeled "model", which lists 16 models: NN5\_MAE, NN5\_MSE, NN4\_MAE, NN4\_MSE, NN3\_MAE, NN3\_MSE, NN2\_MAE, NN2\_MSE, NN1\_MAE, NN1\_MSE, RF\_MAE, RF\_MSE, ELN\_MAE, ELN\_MSE, LM\_MAE, and LM\_MSE. The x-axis for both is labeled "variable". The left heatmap's x-axis lists 20 individual factors, and the right heatmap's x-axis lists 15 macro factors. A color scale on the right of the macro factors heatmap indicates "importance" from 0.00 (blue) to 1.00 (red). In the "Individual factors" heatmap, most cells are blue, with a few light green/yellow cells for ELN\_MAE and ELN\_MSE in the first few factors. In the "Macro factors" heatmap, there are more variations, with some cells in red/orange (high importance) and others in blue (low importance), particularly for the first few models and factors.



Figure 16: Fama French Factors Robustness Check RF VIMP

