

Paper Replication: Empirical Asset Pricing via Machine Learning by Gu, Kelly, and Xiu (2018)¹

Zhen HOU, Jianda MAO, Xiaolong WANG

April 29, 2024

¹<https://youtu.be/6NaL8qs5bfM>

Outline

- 1 Introduction
- 2 Dataset and Preprocessing
- 3 Machine Learning Models
- 4 Performance Evaluation
- 5 Variable Importance
- 6 Conclusion

Introduciton

In this work, the excess return of an asset is modeled by:

$$r_{i,t+1} = \mathbb{E}_t(r_{i,t+1}|z_{i,t}) + \epsilon_{i,t+1}, \quad (1)$$

where

$$\mathbb{E}_t(r_{i,t+1}|z_{i,t}) = g^*(z_{i,t}). \quad (2)$$

Notation:

Stocks $i = 1, \dots, N_t$; months $t = 1, \dots, T$.

$r_{i,t+1}$: the excess return, serving as reponse variable.

$z_{i,t}$: vector of predictor variable.

$g^*(\cdot)$: conditional expectation of $r_{i,t+1}$, e.g. risk premium, given $z_{i,t}$.

The main goal: to estimate $g^*(\cdot)$, a typical prediction task and attractive for machine learning methods.

Outline

- 1 Introduction
- 2 Dataset and Preprocessing**
- 3 Machine Learning Models
- 4 Performance Evaluation
- 5 Variable Importance
- 6 Conclusion

Raw Dataset

- Obtained from².
- contains almost 30000 stocks from 1957 to 2020.
- average 6200 stocks per month.
- DATE: the end day of each month t .
- RET: lag-adjusted CRSP returns $r_{i,t+1}$.
- 94 characteristics $c_{i,t}$.
- sic2: first two digits of SIC code.

²<https://www.dropbox.com/s/zzgjdubvv23xkfp/datashare.zip?dl=0>

We impute missing values for each characteristic by the median of the corresponding subgroup of stocks classified by the belonging month. We include 74 industry dummies $s_{i,t}$ using one-hot encoding corresponding to the first two digits of SIC code. Following the original paper, we construct 8 macroeconomic predictors denoted as x_t . The final covariate $z_{i,t}$ in the original work is calculated by

$$z_{i,t} = [[x_t, 1] \otimes c_{i,t}, s_{i,t}], \quad (3)$$

where \otimes denotes the Kronecker product and P -dimensional vector is viewed as a $1 \times P$ matrix. However, limited by the computational resources, we calculated the predictor variable $z_{i,t}$ by

$$z_{i,t} = [x_t, c_{i,t}, s_{i,t}]. \quad (4)$$

in our experiment while the code for (3) is also provided.

Data Splitting

- 18 years of training data(1957-1974), 12 years of validation data(1975-1986) and 34 years of test data(1987-2020).
- Refit models every year while increasing the training data by one year and maintain the same size of validation by rolling it forward until including all the data.
- After each refit, we use the model to predict the next year's excess returns.

Outline

- 1 Introduction
- 2 Dataset and Preprocessing
- 3 Machine Learning Models**
- 4 Performance Evaluation
- 5 Variable Importance
- 6 Conclusion

- OLS, OLS-3, ENet, PCR, PLS, GBRT, NN1, NN2, NN3, NN4 and NN5.
- Below we just discuss some details when conducting OLS-3 and NN models

OLS-3: Linear Regression with 3 Factors

- OLS-3 is the classical empirical asset pricing linear regression model using 3 factors as predictors.
- The three factors are the book-to-market(bm), the size(mvel1), and the momentum(mon1m, mon6m, mon12m, mon36m) while for momentum there are 4 characteristics calculated by data from the past 1, 6, 12, and 36 months, respectively.
- Thus, there are totally 6 predictors in OLS-3.

NN: Neural Networks

	NN1	NN2	NN3	NN4	NN5
neurons	32	32, 16	32,16,8	32,16,8,4	32,16,8,4,2

[Table:](#) Number of neurons in each hidden layer of NN models

- ReLU, Adam optimizer and MSE loss function.

Outline

- 1 Introduction
- 2 Dataset and Preprocessing
- 3 Machine Learning Models
- 4 Performace Evaluation**
- 5 Variable Importance
- 6 Conclusion

Out-of-Sample R^2

For testing sample \mathcal{T}_3 , the out-of-sample R^2 is defined as

$$R_{\text{oos}}^2 = 1 - \frac{\sum_{(i,t) \in \mathcal{T}_3} (r_{i,t+1} - \hat{r}_{i,t+1})^2}{\sum_{(i,t) \in \mathcal{T}_3} r_{i,t+1}^2}. \quad (5)$$

- the denominator here is the sum of squared response variables, not the sum of the squared residuals.

Result

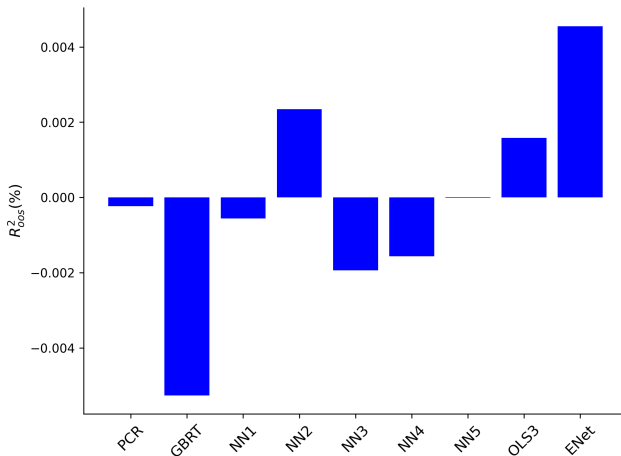
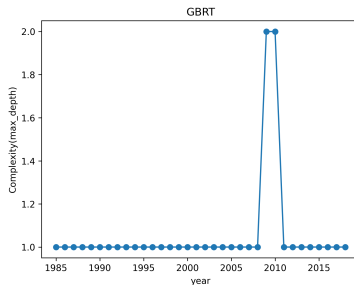
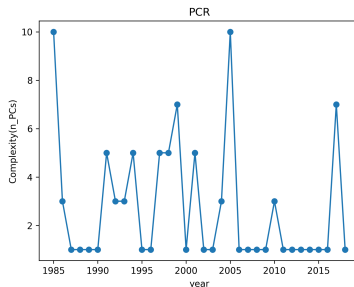
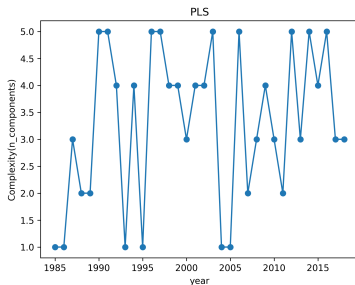


Figure: Out-of-sample R^2 performance of different models except OLS and PLS

Time-varying Model Complexity

Complexity

For PCR and PLS, we use the number of components for evaluation while the maximal depth for GBRT.



Diebold and Mariano test is used to compare the out-of-sample predictive performance of model (1) versus model (2). Denote $\hat{e}_{i,t}^{(1)}$ and $\hat{e}_{i,t}^{(2)}$ as the prediction errors for stock i at testing month t of model (1) and model (2) respectively. Denote $\mathcal{T}_{3,t}$ as the set of stocks in the testing month t . Let

$$d_{12,t} = \frac{1}{|\mathcal{T}_{3,t}|} \sum_{i \in \mathcal{T}_{3,t}} \left((\hat{e}_{i,t}^{(1)})^2 - (\hat{e}_{i,t}^{(2)})^2 \right). \quad (6)$$

The test statistic is defined as $DM_{12} = \bar{d}_{12}/\hat{\sigma}_{12}$, where \bar{d}_{12} is the average of $d_{12,t}$ and $\hat{\sigma}_{12}$ is the Newey-West standard deviation, calculated by

$$\hat{\sigma}_{12} = \frac{1}{T} \sqrt{\sum_{t=1}^T (d_{12,t} - \bar{d}_{12})^2 + 2 \sum_{l=1}^L \sum_{t=l+1}^T w_l (d_{12,t} - \bar{d}_{12})(d_{12,t-l} - \bar{d}_{12})} \quad (7)$$

where $w_l = 1 - \frac{l}{L+1}$ and $L = T - 1$.

And notice that $\hat{\sigma}_{12}$ is the autocorrelation-adjusted estimate of standard deviation of \bar{d}_{12} .

Result

	PCR	PLS	GBRT	NN1	NN2	NN3	NN4	NN5	OLS-3	ENet
OLS	1.73	1.73	1.73	1.73	1.73	1.73	1.73	1.73	1.73	1.73
PCR		-2.52	-1.56	-0.38	3.28*	-0.72	-1.56	0.06	3.09*	6.35*
PLS			2.27	2.54	2.59*	2.66*	2.44	2.39	2.64*	2.86*
GBRT				1.68	2.16	1.20	1.30	1.40	2.26	3.06*
NN1					2.07	-0.55	-2.09	0.29	2.81*	4.90*
NN2						-1.55	-3.20*	-2.31	-0.52	2.26
NN3							-0.04	0.53	1.70	3.22*
NN4								1.17	4.38*	6.28*
NN5									1.34	3.88*
OLS3										3.17*

Figure: Diebold-Mariano test results

A positive value indicates the column model outperforms the row model while a negative value indicates the row model is better. The bold values indicate the difference is significant at 0.05 level. The star symbol indicates significance at 0.05 level for 10-way comparisons. The ENet model outperforms all other models significantly.

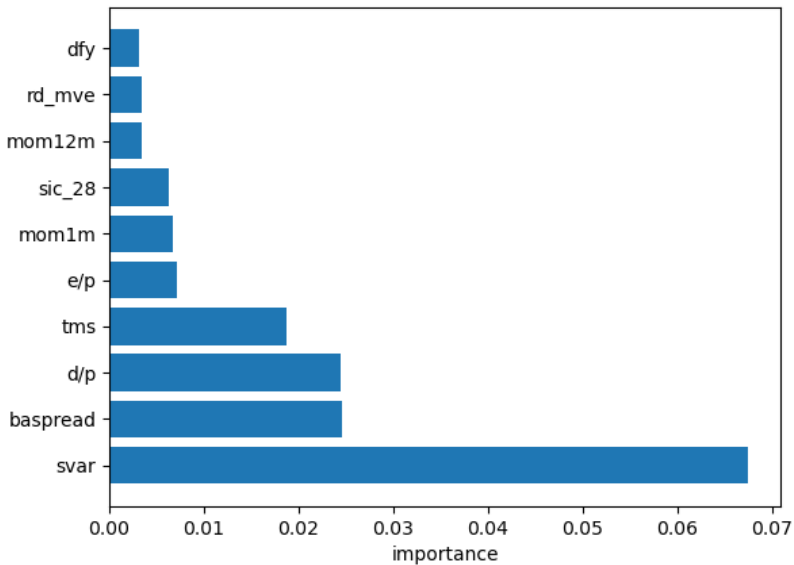
Outline

- 1 Introduction
- 2 Dataset and Preprocessing
- 3 Machine Learning Models
- 4 Performance Evaluation
- 5 Variable Importance**
- 6 Conclusion

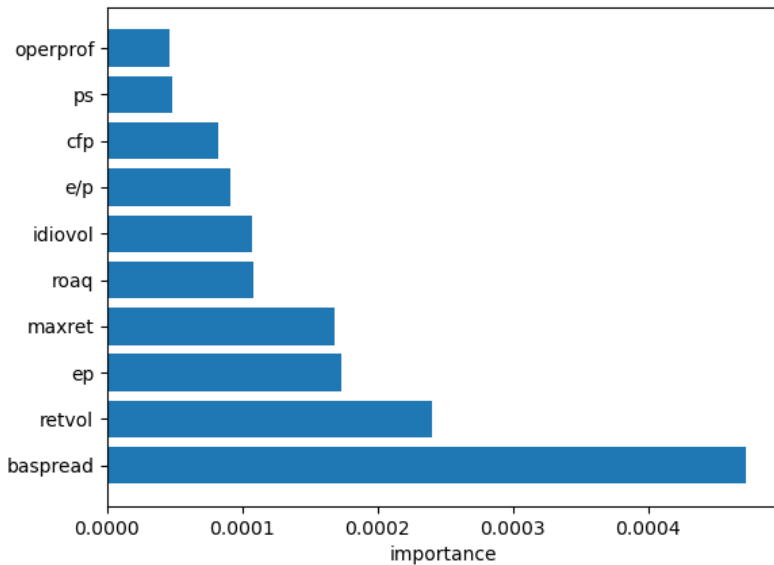
Variable Importance

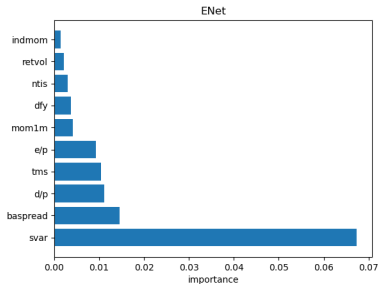
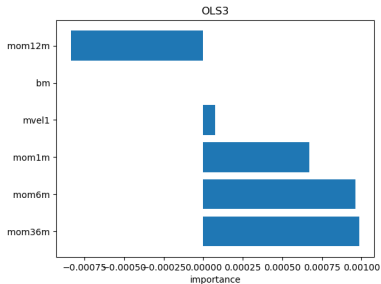
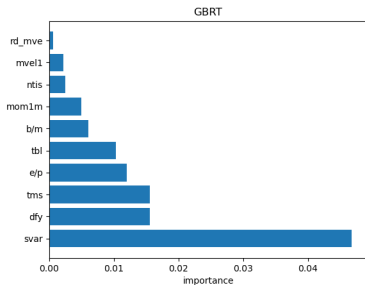
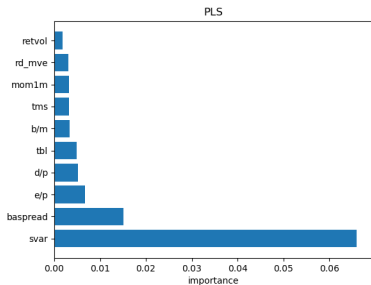
We use permutation feature importance algorithm to study the variable importance. Generally, the bid-ask spread(baspread) and macroeconomic predictors(svar, d/p, etc) are important in our result.

OLS



PCR






Outline

- 1 Introduction
- 2 Dataset and Preprocessing
- 3 Machine Learning Models
- 4 Performance Evaluation
- 5 Variable Importance
- 6 Conclusion**

Conclusion

- replicated the main results of the original paper Empirical Asset Pricing via Machine Learning by Gu, Kelly, and Xiu (2018).
- confirm the findings of the original paper.
- provide more technical details.
- Limited by time and computational resources, we did not conduct the full replication of the original paper.
- More code in our GitHub repository³.

³<https://github.com/hdsfade/math5470-final-project> 

- Zhen HOU is responsible for the code of PCR, performance evaluation, presentation and report writing.
- Jianda MAO is responsible for data preprocessing, training framework implementation, code of ENet, GBRT, NNs and variable importance.
- Xiaolong WANG is responsible for the code of OLS, OLS-3, PLS and running the experiments.