

**Dimensionality Reduction and Machine Learning Approaches to
Factor Investing and Stock Return Prediction:**

**A Comparative Study of SPCA versus LASSO Feature Selection on
Neural Network**

April 19th, 2025

AFM 423, University of Waterloo

Kyle Shi 20933475

Iris Gao 20815552

Table of Contents

<i>Abstract</i>	3
<i>Introduction</i>	3
<i>Application of ML Approach to Factor Investing</i>	3
<i>Research Questions and Objective</i>	4
<i>Data Processing</i>	5
<i>Variables and Measures</i>	6
<i>Experimental Methodology</i>	7
<i>LASSO Feature Selection Process</i>	7
<i>SPCA Implementation</i>	7
<i>Neural Network Architecture</i>	8
LASSO + Neural Network Hyperparameter Tuning	8
SPCA + Neural Network Hyperparameter Tuning	9
<i>Portfolio Construction and Evaluation</i>	9
<i>Results and Discussion</i>	10
<i>LASSO Feature Selection + Neural Networks Result</i>	10
<i>SPCA Implementation + Neural Networks Result</i>	13
<i>Portfolio Performance Comparison</i>	15
<i>Related Work</i>	17
<i>Conclusions</i>	17

Abstract

This project implements two different machine learning methodologies to construct optimal portfolio based on predicted 1-month forward stock returns. The first methodology is LASSO feature selection combined with a Neural Network, and the second methodology is Sparse Principal Component Analysis (SPCA) combined with a Neural Network model. These two machine learning strategies are used to construct top 20% long-only portfolio and compared against the baseline strategy of longing the S&P500 market index. Using a comprehensive financial dataset from 2000 to 2018 that comprises information on 1,207 US listed stocks and their associated 93 firm-specific characteristics, LASSO and SPCA are employed to conduct feature selection, which then is used to train Neural Networks to predict stock returns. The predictive performance is evaluated through out-of-sample testing over the period 2017-2018. Portfolios are constructed monthly by ranking predicted returns and the top 20% of ranked stock are being selected in the long-only strategy portfolio. The experimental result shows that the holding period return for SPCA-NN, LASSO-NN and S&P500 passive portfolio is 53.67%, 32.78% and 10.09%, respectively; while the Sharpe Ratio for SPCA-NN, LASSO-NN and S&P500 passive portfolio is 0.28, 0.20, and 0.14, respectively; which reveals that the top 20% portfolio favors SPCA-NN when comparing against the LASSO-NN and the baseline S&P500 passive strategy.

Introduction

Predicting stock returns has long been a popular subject in academic finance and investment practice. Traditional factor models such as CAPM and Fama-French model, use simple, linear combinations of a few economic drivers, which limits their ability to fully capture non-linear and complex relationships between a wide array of firm-specific and macroeconomic variables. Recent advancement in machine learning models offer enhance predictive performance and model flexibility by allowing the intake of high-dimensional dataset to uncover hidden relationships amongst many predictors.

Application of ML Approach to Factor Investing

This project explores and analyzes various machine learning approaches in factor investing to construct portfolio with optimal return, specifically Sparse principal Component Analysis (SPCA), LASSO regression, and feedforward Neural Networks.

One of the literature papers reviewed in project one introduced a variant of the traditional PCA known as SPCA, which incorporates the *lasso* ($L1$ penalty) and *elastic net* ($L1+L2$ penalties) to modify principal components with sparse loadings, resulting in zero coefficients for many principal component loadings, thereby reducing the number of explicitly used variables. Another literature paper reviewed is based on the Least Absolute Shrinkage and Selection Operator (LASSO) regression, which is a feature selection method that uses $L1$ penalty to shrink some coefficients to exactly zero, thereby reducing the complexity of the model. However, LASSO cannot capture non-linear relationships.

Neural Network is also discussed in one of the literature papers, which predicts stock returns due to the capability of the model to capture non-linear relationship and interaction effects. A simple feed-forward neural network architecture is employed to capture complex relationships between

inputs and outputs, potentially offering improvements in predictability. The NN model consists of three or more layers, where the input layer comprised a set of firm characteristics, then one or more hidden layers capture the interactive effects between different variables and perform non-linear transformations through activation functions on input variables, and the output layer aggregates all information from the hidden layers to generate an output. The literature paper suggests MSE as the preferred objective function and ReLU as the activation function due to its computational simplicity.

In this project, the challenge of high feature dimensionality is addressed by employing two different methods: LASSO regression and SPCA. LASSO performs feature selection across predictors by shrinking some coefficients to exactly zero for those with low predictive power, while SPCA extracts sparse, interpretable latent factors that captures the most variance. These two techniques focus on extracting most informative signals while lower the risk of overfitting. The next step is to implement a feedforward Neural Network for each of the dimensionality reduction technique. After predicting 1-month forward returns for individual stocks, the model outputs are used to rank stocks and construct long-only portfolio based on top 20% predicted performers each month. A baseline passive investing portfolio is also constructed to hold S&P500 over the same investing period. The performance of the machine-learning-based portfolios is evaluated on the test set and then compared against the performance of passive portfolio.

Several caveats must be taken into consideration when applying these models to financial data. Financial returns data is extremely noisy and is characterized by a low signal-to-noise ratio, meaning that the predictive relationships are often weak due to substantial randomness. Financial returns are also subjective to regime shifts and market shocks, which makes model prediction extremely challenging. Since there are many characteristics that can be used as predictors, without appropriate dimensionality reduction, models are highly susceptible to overfitting and fail to generalize to unseen data.

Research Questions and Objective

The goal of this research project is to investigate whether dimensionality reduction combined with a machine learning model can enhance the ability to construct an optimal portfolio that can yield a superior return. Following the in-depth analysis from project one, where several literature papers that deploys different machine learning models to enhance the factor selection and portfolio construction process were compared in terms of each of their strengths and weaknesses; this research project replicates the models from those literature papers by generally following the implementation plan developed in project one report, with some tweaks and improvements in the implementation process that will be discussed in details in the methodology section. The different combinations of dimensionality reduction and predictive modelling approaches is compared to answer the following research questions that guides the report:

- Can machine learning models built on SPCA or LASSO selected features generate superior investment performance compared to passive investment on S&P500 index over the same period?
- Does dimensionality reduction using SPCA improve stock return predictability relative to direct feature selection via LASSO?
- Does Neural Network provide added value to portfolio construction?

By addressing the first question, it provides direct evidence of whether feature selection and nonlinear modelling technique offers additional value beyond simple market exposure. With the industry being increasingly dominated by the passive investing strategies, leveraging machine learning models that demonstrate an edge over passive benchmarks is valuable and useful. By examining the relative advantages of SPCA versus LASSO in enhancing portfolio performance addresses the question of which method is best at reducing dimensionality in financial machine learning application, which may provide a pathway toward more robust predictive models. The third question explores whether Neural Networks' ability to capture nonlinearities can provide additional value in constructing a more profitable portfolio. All three questions contribute to the broader understanding of how machine learning techniques can be adapted to improve portfolio construction, as well as addressing this research project's objective of comparing which dimensionality reduction model is more effective and whether the combination with Neural Network can enhance factoring investing and improve stock return prediction.

Data Processing

This project involves two data sources: the first data source is the `data_ml` dataset provided in this course, and the second data source is the S&P500 index (^GSPC) from Yahoo Finance. The `data_ml` dataset comprises information on 1,207 stocks listed in the US with time periods starting from 1998/11/11 to 2019/03/31 for a total of 244 monthly observations. Besides the stock id and date columns which are used as identifiers, this dataset also contains 93 firm characteristics associated with each of the stocks and 4 labels which are the 1-month, 3-month, 6-month and 12-month forward returns. The returns are total returns that incorporate potential dividend payments over the considered periods.

The data wrangling process begins with selecting the data from period 1999-12-31 to 2019-01-01 and dropping potential target leaking predictors. Extracting data from the period 2000 to 2018, inclusive, to ensure reliability and completeness of the dataset. This timeframe was chosen because several stocks exhibit missing information beyond 2019. The response variable of interest in the project is "`R1M_Usd`", which is the 1-month forward return. Besides the target variable, the target leaking variables are identified to be "`R3M_Usd`", "`R6M_Usd`", and "`R12M_Usd`". These predictors introduce forward-looking information, in this case, these three variables all represent forward returns, which must be removed to ensure that the model relies only on information available at the time of prediction to avoid target leakage.

Observations (rows) with missing values across predictors are dropped. Duplicated rows are removed as well. Note that columns with zero variance needs to be removed as they carry no predictive information and may cause instability during standardization. Fortunately, no constant columns have been identified, therefore all columns are preserved. The next step is to check whether the predictors need to be standardized. The calculation shows that the mean range of the dataset is between 0.0127 and 0.506, while the standard deviation range is between 0.176 and 0.289. Prior to modeling, the predictors exhibit a wide range of means and standard deviations vary across features. Ideally, the means should be centered around 0 and standard deviations should be around 1, since models like LASSO, SPCA and NN are sensitive to the scale of input variables. Without standardization, predictors with larger magnitudes can disproportionately

influence model outcomes, biasing feature selection in LASSO, dominating the identification of principal components and destabilizing gradient optimization in Neural Networks, therefore standardization is applied to all predictor variables except the target variable.

To avoid lookahead bias, the dataset is split into training, validation and test sets chronologically. The training set contains data from 2000 to 2013, which is used to fit the models. The validation set contains data from 2014 to 2016, which is used to tune hyperparameters. The test set from 2017 to 2018 is used for final evaluation of models and portfolio performance on unseen data. This splitting ensures that only past information is used in the training set to predict future returns, while the test set only contains truly unseen data.

The S&P500 data is used separately as a baseline passive investing strategy. From the tidyquant package, the monthly adjusted S&P500 prices from 2017 to 2018 are pulled from Yahoo Finance. The adjusted price is used because it accounts for stock splits, dividends and corporate actions, which reflects the true economic return. The time frame chosen is the same as the time frame for the test set, since ultimately the portfolios being compared holds assets over the period 2017-2018.

Variables and Measures

The 93 predictor variables in this study consist exclusively of firm-specific characteristics. These features capture different economic and financial aspects of individual stocks which include the following:

- Valuation metrics such as price-to-book ratio, price-to-earnings ratio, EBITDA and enterprise value to EBITDA that describes how the market values the company's assets and earnings
- Size and liquidity indicators such as market capitalization, trading volume, volatility measures that reflect the company size and stock liquidity.
- Profitability metrics such as return on assets, return on equity, operating margins, and free cash flow margins, which assess a company's ability to generate profit and cash flows.
- Leverage and capital structure metrics such as debt-to-equity ratio, total debt to capital and net debt, which measures financial risk.
- Momentum indicators such as past 11-month momentum, which captures price trends.

All these predictors are standardized and fed into the LASSO and SPCA model. Many of these variables are highly correlated, and how these two techniques handle highly correlated variables will be discussed in later sections. The response variable is the 1-month forward stock return (R1M_Usd), which already accounts for dividends, stock splits and other corporate actions.

The ultimate objective of this project is not just to predict stock returns, but to construct optimal portfolios based these predictions, several portfolio performance metrics are used to evaluate the effectiveness of each model.

- The average monthly return is the simple mean of all predicted monthly portfolio returns over the test period

- Holding period return is a compounded return over the test period and it is calculated as $HPR = \prod_{i=1}^n (1 + r_i) - 1$ where i is the predicted return in month i . This metric captures the total growth in stock return since the beginning of the test period
- Standard deviation of monthly returns measures the volatility of monthly returns. The higher the volatility, the greater the risk.
- Monthly Sharpe Ratio assesses the risk-adjusted returns, which is calculated as $SharpeRatio = (E[R_p] - R_f) / \sigma_p$. Note that the calculation of Sharpe Ratio in this project simplifies to $SharpeRatio = E[R_p] / \sigma_p$ because of two reasons. The first reason is that the US risk-free rates were very low during 2017-2018, roughly around 2%-2.5% annually, which translates to roughly 0.2% per month. This small amount is insignificant and is negligible. The second reason is that dropping this risk-free rate is a common simplification technique in research projects like this when the focus is on relative performance, so it is acceptable to just assume the risk-free rate is 0.

Experimental Methodology

LASSO Feature Selection Process

After chronologically splitting the full dataset into training, validation and test sets, the target returns variable forms the y_{train} , y_{val} , and y_{test} dataset, while the rest of predictor variables form the x_{train} , x_{val} , and x_{test} dataset. Prior to predicting using machine learning models, LASSO regression was applied using 'glmnet' to perform feature selection with 'R1M_Usd' as the response variable. The LASSO model imposes an L1 penalty on the coefficient to shrink those with low predictive power to exactly zero, requires the input of the hyperparameter value λ , which controls the degree of penalization. In order to determine the optimal λ value, cross-validation with a default 10-folds was performed on the training set by using 'cv.glmnet' function. Upon obtaining the optimal λ value, use this value to perform LASSO regression, which would result in a model that retained a number of non-zero coefficients out of the initial 93 characteristics.

In this study, LASSO was applied solely for feature selection to reduce dimensionality, and stock return predictions were subsequently conducted using a feedforward Neural Network model trained on these LASSO selected features. Neural Network model training procedure using LASSO selected feature inputs is discussed in the Neural Networks section below.

SPCA Implementation

The second dimensionality reduction or feature generation technique deployed is SPCA, which constructs new and compressed feature representations for the NN model. The SPCA input involve all 93 standardized firm-specific features, then it seeks to extract a set of orthogonal principal components that depends only on a subset of input variables. Unlike traditional PCA, sparsity is introduced into the loading vectors to enhance interpretability and reduces overfitting risks, therefore SPCA was chosen over traditional PCA to achieve a sparser set of principal components with heavy loading only on a subset of significant variables.

To construct the SPCA, the function `sPCA()` is used from the PMA package, with hyperparameters being pre-specified. The number of components extracted K is set to 5, which balances the need for sufficient information while maintaining a manageable number of features. The sparsity level hyperparameter (`para`) encourages greater sparsity for lower values. Several values of `para` including `para` value of 0.2, 0.4, 0.6, and 0.9 are each used as the input to perform SPCA dimensionality reduction to assess the impact on the resulting SPCA loadings and component structure. Based on the resulting loadings, the most optimal `para` value is chosen that can balance between stability and sparsity. This chosen `para` value is then used in the SPCA model to generate a new set of features.

In this study, SPCA was applied solely for feature generation, and stock return predictions were subsequently conducted using a feedforward Neural Network model trained on these SPCA generated features. Neural Network model training procedure using SPCA generated features is discussed in the Neural Networks section below.

Neural Network Architecture

LASSO + Neural Network Hyperparameter Tuning

After selecting a subset of features from LASSO regression, a feedforward Neural Network (NN) model was constructed to capture non-linear interactions amongst variables. To achieve the best model performance, a hyperparameter tuning process is required to determine the optimal parameter values to be inputted into the model:

- 1st hidden layer units represent the number of neurons in the first hidden layer, which controls the model's ability to learn complex feature interactions
- 2nd hidden layer units represent the number of neurons in the second hidden layer, which compresses and refines the information learned from the first layer
- Learning rate controls the step size in optimization process. A smaller value leads to slower model weights updates, meaning that the model learns slower and results in more stable convergence.
- batch size is the number of samples used in each smaller batches during training.

The first step of the tuning process was to set a reasonable range for each of the above hyperparameters. For instance, units for 1st layer may take on any value in (16, 32, 64), units for 2nd layer may take on any value in (8, 16, 32), learning rate can be 0.001 or 0.0005, and batch sizes can be either 64 or 128. For all possible combinations within the reasonable value range, the NN model is performed on training data with features selected by LASSO. Each candidate model was trained for 100 epochs, and performance was evaluated on the validation set based on the mean squared error (MSE). After obtaining the optimal hyperparameter combination, the NN model is trained based on these hyperparameters.

As for training control parameter, each epoch represents one full pass through the entire training dataset, the NN model needs to train for enough epochs to gradually learn and minimize loss function. The model is trained for 100 epochs, which is a common value because it is large enough to provide enough room for learning and converge toward a minimum, and at the same time, 100 epochs is short enough to avoid overfitting.

The training and validation loss curves generated by the Neural Network model is plotted to assess the model learning capability, and how well the model is trained on the validation set comparing to the training set. Then the model prediction performance is assessed on both training and test set, where training and test RMSE values are compared to assess the predictive power of the model and whether overfitting issue occurs. A scatter plot between stock returns in test data and predicted stock returns is produced for a sanity check to evaluate whether the model prediction makes sense in a financial modeling setting.

SPCA + Neural Network Hyperparameter Tuning

Following the feature generation step using SPCA, the new features are fed into a feedforward Neural Network for stock return prediction. Although hyperparameters were tuned for the LASSO model, it is necessary to tune again using SPCA features. This is because LASSO uses original predictors, but SPCA generates transformed features, the change in data inputs would affect the training dynamics as the nature and distribution of the inputs are different, therefore, tuning prior to training NN is required. Following a similar tuning process as before such that the units for 1st layer may take on any value in (16, 32, 64), units for 2nd layer may take on any value in (8, 16, 32), learning rate can be 0.001 or 0.0005, and batch sizes can be either 64 or 128. For all possible combinations within the reasonable value range, the NN model is performed on training data with features generated by SPCA. Each candidate model was trained for 100 epochs, and performance was evaluated on the validation set based on the mean squared error (MSE). After obtaining the optimal hyperparameter combination, the NN model is trained based on these hyperparameters.

The training and validation loss curves are plotted to assess model generalization and risk of overfitting. Then the model prediction performance is assessed on both training and test set, where training and test RMSE values are compared to assess the predictive power of the model and whether overfitting issue occurs. A scatter plot between stock returns in test data and predicted stock returns is produced for a sanity check to evaluate whether the model prediction makes sense in a financial modeling setting.

Portfolio Construction and Evaluation

Three long-only portfolios are constructed in this project to compare which method yields the best portfolio performance. For the baseline passive portfolio, the S&P500 index is held over the test period (2017-2018) using financial market data. For the other two portfolios, after training the NN models, for the predicted returns generated by both LASSO+NN and SPCA+NN models, stocks were ranked based on their predicted returns using a percentile ranking. To construct long-only portfolios, select stocks in the top 20% percentile each month and take an equal weight within selected stocks to calculate the average realized returns. The monthly long-only portfolio returns are aggregated over the test period (2017-2018).

Performance metrics such as holding period return, monthly return standard deviation and Sharpe Ratio are then calculated across three portfolios to determine whether machine-learning models have an advantage over the passive strategy, as well as comparing which machine-learning model has a superior performance.

Results and Discussion

LASSO Feature Selection + Neural Networks Result

LASSO Feature Selection Result

The cross-validation result is shown below:

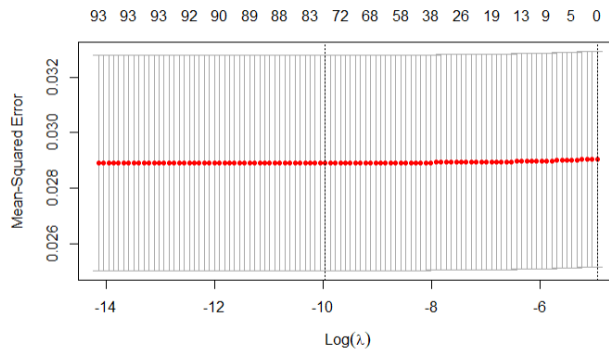


Figure 1 lambda cross-validation result

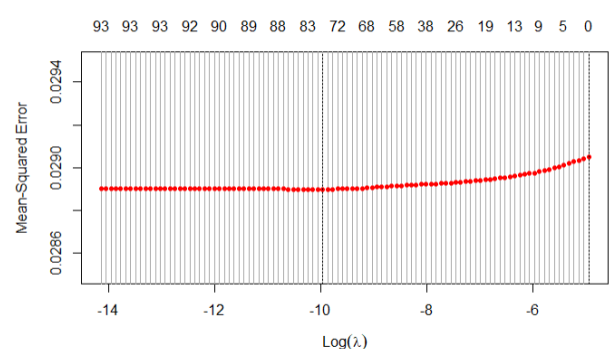


Figure 2 lambda cross-validation result (zoomed-in)

This plot depicts the cross-validation MSE on y-axis and $\log(\lambda)$ on x-axis. The red dots plot the MSE of each corresponding $\log(\lambda)$, while the vertical dashed line indicates the $\log(\lambda)$ value that correspond to the minimum MSE. Note that the above plots are the same except that the right plot is zoomed in, with y-axis between 0.0285 and 0.0295, which can better display the change in MSE. The left plot shows that the MSE curve appears to be relative flat across a wide range of λ values, suggesting that the model is not highly sensitive to the choice of λ . From the right plot, there is an obvious trend that the MSE curve starts to increase for $\log(\lambda)$ greater than -10, and the MSE stays flat for $\log(\lambda)$ less than -10. Therefore, the optimal λ value selected would be 0.00004720529, which yields the lowest cross-validation error.

Using the optimal λ , the LASSO regression was performed to conduct feature selection, which results in a model that retained 81 non-zero coefficients out of the initial 93 characteristics. The LASSO retained a majority of features with similar coefficients, suggesting that these firm-specific variables provide incremental predictive information about the target return but none of which shows significant predictive power. After sorting the coefficient value, the market capitalization and free cash flow on book value have a relatively higher coefficient, which in standardized LASSO, may indicate greater predictive contribution to the model, but it is by no means a direct measure of predictive power. Typically, LASSO would select one representative variable from each correlated group of predictors, however, the selected list retains many variables that are correlated. This might be because these variables contribute to weak but complementary signals. The relatively large set of selected predictors suggests that the financial data is every noisy, a broad predictor set may be needed to extract useful predictive signals.

LASSO + NN Model Prediction Result

The resulting optimal hyperparameter combination is obtained to be 32 units for 1st layer, 16 units for 2nd layer, a learning rate of 0.001, and batch size of 128, which corresponds to the lowest validation MSE of 0.01064.

Using the tuned hyperparameters, the NN model was re-trained on the training set with a model architecture consists of the following:

- an input layer corresponding to the 81 LASSO selected features
- a 1st hidden layer with 32 neurons and ReLU activation as recommended in the literature paper reviewed in project one
- a 2nd hidden layer with 16 neurons and ReLU activation as recommended in the literature paper reviewed in project one
- an output layer with a single neural representing the predicted forward stock return.

The training and validation loss curves for the Neural Network plot is shown below:

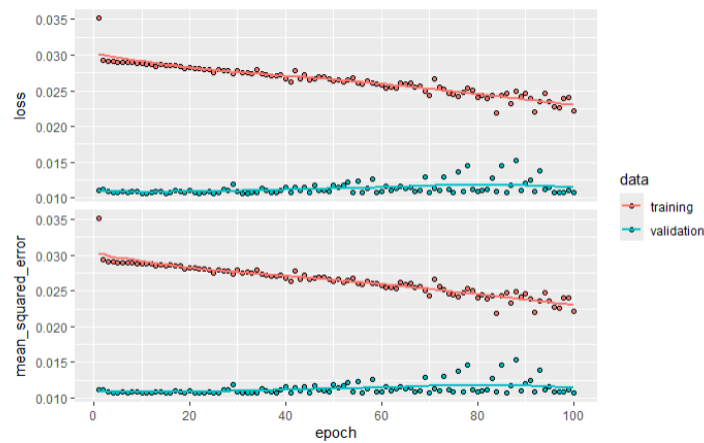


Figure 3 - LASSO+NN Training and Validation Loss Curve

The y-axis represents the mean squared error, and the x-axis represents the number of epochs. The red dots represent training loss curve over epochs while the green dots represent validation loss curve over epochs. Starting from 0 epochs, both training and validation loss steadily decrease across epochs, indicating that the model was picking up meaningful patterns. Note that the initial decrease in training MSE was more stable for less than 40 epochs, then as epoch increases, the training MSE decreases relative stable but with mild fluctuations. The validation MSE remains relatively flat until epoch reaches 70, mild fluctuations (higher MSE) in validation loss toward later epochs, reflecting inherent noise in stock return prediction. On a positive note, the validation MSE fluctuations are contained and do not show a specific upward trend, suggesting that the model is not overfitted. Although the training loss curve is above the validation loss curve across all epochs, the gap between training and validation loss remains relatively small. The small difference of roughly 0.01-0.02 suggests that the model generalizes well to unseen data within the training period. Although it is uncommon for validation MSE to be consistently below the training MSE, it is still explainable as the training set over 2000-2013 period is subject to dramatic economic regime shifts (great financial crisis) whereas the validation set over 2014-2016 may be statistically slightly easier to predict due to stable bull market. Overall, using LASSO selected features and tuned hyperparameters to train the NN achieves a stable learning behavior and a reasonable generalization to the validation set without severe overfitting.

Following NN model training, the model's predictive performance needs to be evaluated separately on the training set and the test set. The validation set is not used for model evaluation because it is used before training to tune the parameters and during training to guide model selection and prevent overfitting. Using the tuned NN model to predict the 1-month forward stock returns on the training and test set. The resulting training RMSE is 0.1515377 while the test RMSE is 0.2806537. The training RMSE measures how well the model predict the training data, and the test RMSE assess the model's ability to generalize to unseen data. The training RMSE of 0.1515377 indicates that the model can capture patterns with reasonable accuracy and without underfitting as the error is not huge. As expected, the test RMSE is higher than the training RMSE, which reflects the challenges of predicting stock returns. The difference in RMSE is noticeable but not extreme, suggesting potential mild overfitting as the test error is almost double the training error; however, since the financial returns are inherently noisy and difficult to predict, the test RMSE is still within reasonable range.

To complement the RMSE evaluation, the actual test returns are plotted against the predicted test returns.

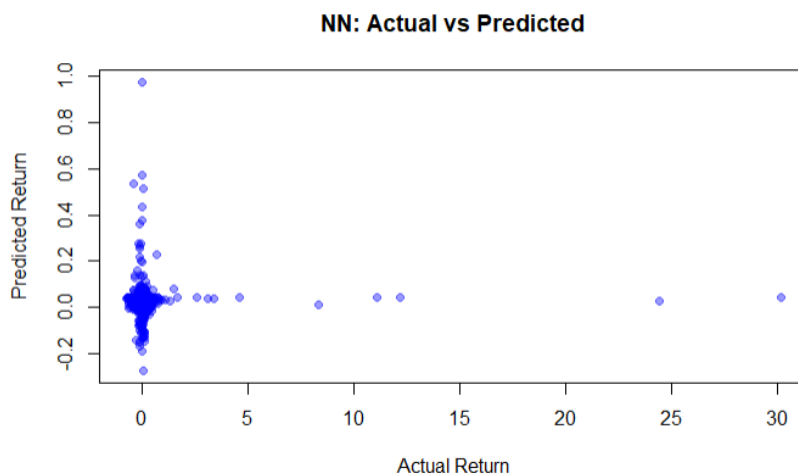


Figure 4 - LASSO+NN Actual vs Predicted Stock Returns

This scattered plot visually displays how well the NN model predicts the return and captures the variability in the target variable. Most points are concentrated and clustered around zero regardless of actual return values, with some negative predictions around -0.2% and some positive predictions between 0.2% and 1%. In contrast, the actual stock returns are heavily right-skewed, with a couple of extreme outliers being around 10%, 25% and 30%. The NN model underestimates the magnitude of returns. This conservative prediction is common in financial modeling and is expected because NN model tends to avoid overreacting to extreme observations. Since the focus is on portfolio construction, and less focused on being able to predict precise returns, the absolute accuracy of the predicted returns in this case is less critical than the relative ranking amongst stocks. In the later portfolio construction step, correctly predicting the ranking of the most promising stocks is crucial to the long-only portfolio strategy. Even though the graph shows more compressed and conservative predictions with a limitation in predicting extreme returns, the model's ability to perform stable ranking differentiation can still drive successful stock selection.

SPCA Implementation + Neural Networks Result

SPCA Feature Generation Result

Several values of para (0.2, 0.4, 0.6, 0.9) were tested to evaluate the impact on loadings and the resulting component structure (note that the testing result on different para values are not included in the r-code). The result shows that across all para values, the SPCA loadings remain identical up to the fourth decimal place, suggesting that changes in penalty strength have little impact on extracted components due to highly correlated variables. Based on this finding, a para setting of 0.4 was adopted because this value provides a balance between sparsity and stability. Since para value does not alter the result significantly, then hyperparameter tuning for SPCA is not necessary.

The final loadings for five principal components are listed in the R output, with approximately 62% of loadings being less than 0.05, suggesting that SPCA can successfully shrink the influence of many weaker predictors toward zero. Although none of the SPCA coefficients are forced to exact zeros, this method still effectively prioritized meaningful predictors and suppressed noise from over half of the less informative variables.

Following the extraction of principal components, the original standardized features are multiplied by the SPCA loading matrix to transform the high-dimensional data with 93 predictors into a lower-dimensioned representation based on 5 principal components. The resulting transformed datasets consists of 5 features, each feature being the linear combination of the original predictors with many weaker predictors having minimal influence due to near zero coefficients. By working with the new SPCA components, multicollinearity is handled, risk of overfitting is reduced, and interpretability is improved. Now the transformed features serve as the new inputs for training the NN model, which is discussed in detail in the NN section below.

SPCA + NN Model Prediction Result

The resulting optimal hyperparameter combination is obtained to be 64 units for 1st layer, 8 units for 2nd layer, a learning rate of 0.001, and batch size of 64, which corresponds to the lowest validation MSE of 0.0106.

Using the tuned hyperparameters, the NN model was re-trained on the training set with a model architecture consists of the following:

- an input layer corresponding to 5 SPCA generated features
- a 1st hidden layer with 64 neurons and ReLU activation as recommended in the literature paper reviewed in project one
- a 2nd hidden layer with 8 neurons and ReLU activation as recommended in the literature paper reviewed in project one
- an output layer with a single neural representing the predicted forward stock return.

The training and validation loss curves for the Neural Network plot is shown below:

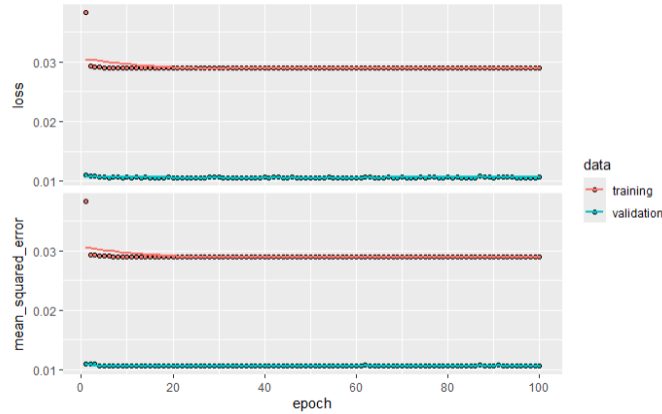


Figure 5 - SPCA+NN Training and Validation Loss Curve

The red dots represent training loss curve while the green dots represent the validation loss curve. Both training and validation loss decreases initially for small epoch, then remained relatively flat throughout the epochs. Furthermore, training loss is roughly 0.03 whereas the validation loss is roughly 0.01, the loss difference of 0.02 remains relatively constant across epochs. This shows that the model generalizes well to unseen data and does not suffer from overfitting. The plot does not show any spikes or oscillations amongst dots, suggesting that the learning rate was appropriately tuned for the SPCA features.

Then the SPCA+NN model is used to predict stock returns based on training and test sets. The RMSE is used as performance metrics to evaluate the model performance. The resulting training RMSE and test RMSE is 0.1702 and 0.2803, respectively. The training RMSE is relatively low, indicating a good training fit. However, the test RMSE is somewhat higher, which is within expectation due to the inherent noise and high unpredictability of financial returns. Nonetheless, the gap between test and training RMSE is not huge and within reasonable range, suggesting that the model generalizes reasonably well.

The outputs below provide a direct comparison between LASSO+NN and SPCA+NN model performance. As can be seen by the LASSO+NN result on the left side and SPCA+NN result on the right side, the LASSO+NN training RMSE is lower than SPCA+NN by 0.03, while the LASSO+NN test RMSE is higher than SPCA+NN by 0.0004. Although both models produce similar in-sample and out-of-sample RMSE, SPCA+NN is slightly better due to the slightly better SPCA test error. It is alright to accept a slightly higher training error, so the model generalizes better on unseen data. On the other hand, LASSO+NN is relatively overfitted compared to SPCA because LASSO's training error is lower, but test error is higher than SPCA, which is not desirable.

```
6192/6192 [=====] - 3s 402us/step
853/853 [=====] - 0s 394us/step
LASSO + NN Train RMSE: 0.1515377
LASSO + NN Test RMSE: 0.2806537
```

Table 1 - LASSO+NN Training vs Test RMSE

```
6192/6192 [=====] - 3s 417us/step
853/853 [=====] - 0s 407us/step
SPCA + NN Train RMSE: 0.1702111
SPCA + NN Test RMSE: 0.2802628
```

Table 2 - SPCA+NN Training vs Test RMSE

To further investigate model behavior, a scatter plot of actual return versus predicted return on test set is generated. The resulting plot shows that majority of the predicted returns are clustered near zero, while the actual returns are right-skewed, with some outliers around 10%, 25% and 30%. The predicted returns range from -0.04% to 0.04% which underestimates the actual returns. This conservative prediction is normal since the model is trained to minimize MSE, thus it tends to predict values close to the mean zero rather than matching the extreme outliers. Same as the explanation provided for LASSO, the relative ranking of predicted returns is more critical than the accuracy of return predictions. Overall, the plot makes sense given the noisy nature of the financial data.

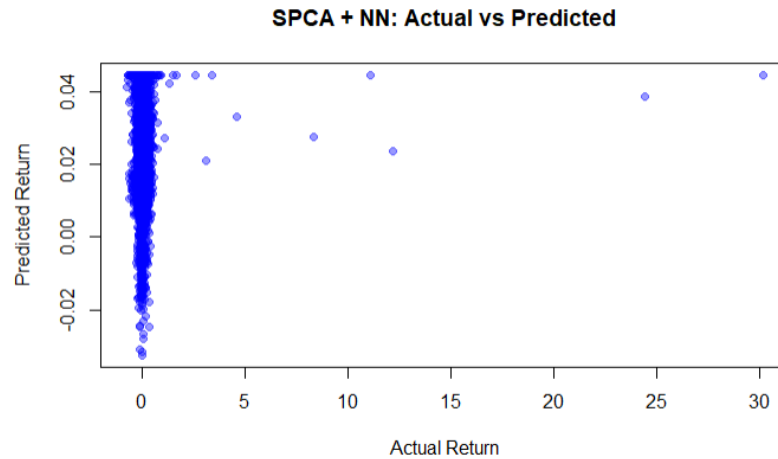


Figure 6 - SPCA+NN Actual vs Predicted Stock Returns

Portfolio Performance Comparison

Three portfolios have been constructed in total. To assess the practicality of machine learning models, two long-only portfolios were constructed on the predicted stock rankings from LASSO+NN and SPCA+NN models. Each month, the portfolio was formed by investing equally in the top 20% ranked stocks. The baseline portfolio with passive investing strategy is also constructed to hold S&P500 over the same period. To compare the performance of three portfolios, the average monthly returns, holding period returns, standard deviation of monthly returns and Sharpe ratios were calculated over the test period 2017-2018.

The SPCA+NN model achieves the highest performance across all metrics. It yields an average monthly return of 2.15%, a holding period return of 58.54% and a Sharpe ratio of 0.322. The LASSO+NN model generates a slightly lower average monthly return of 1.55%, a lower holding period return of 36.40% and a lower Sharpe ratio of 21.47%. Clearly, the SPCA+NN model outperforms LASSO+NN, which aligns with the previous analysis on training and test RMSE, where the SPCA+NN was deemed a slightly better model as it generalizes better on unseen data. In comparison against the passive strategy, the S&P500 investment achieves the lowest average monthly return of 0.46%, its holding period return is only 10.09%, significantly lower than the returns generated by the machine-learning-based models. Its Sharpe ratio of 0.135 is also significantly lower compared to the other two portfolios.

model <chr>	avg_monthly_return <dbl>	holding_period_return <dbl>	sd_monthly_return <dbl>	sharpe_ratio <dbl>
LASSO_NN	0.01550896	0.3639586	0.07225157	0.2146522
SPCA_NN	0.02146822	0.5854073	0.06660703	0.3223116

Table 3 - LASSO+NN vs SPCA+NN Portfolio Performance

model <chr>	avg_monthly_return <dbl>	holding_period_return <dbl>	sd_monthly_return <dbl>	sharpe_ratio <dbl>
SP500	0.004584832	0.100942	0.03392919	0.1351294

Table 4 - Passive Investing on S&P500 Portfolio Performance

The cumulative return plot further visualizes the change in portfolio cumulative return over the test period (2017-2018) across three portfolios. Both machine-learning-based strategies significantly outperforms the S&P500 index, particularly during the bull market period in 2018. The blue line represents SPCA+NN strategy while the red line represents LASSO+NN strategy. Not only does SPCA-based strategy consistently generate higher return than the LASSO-based strategy, it also delivers better risk-adjusted performance as indicated by its relative lower monthly volatility and higher Sharpe ratio.

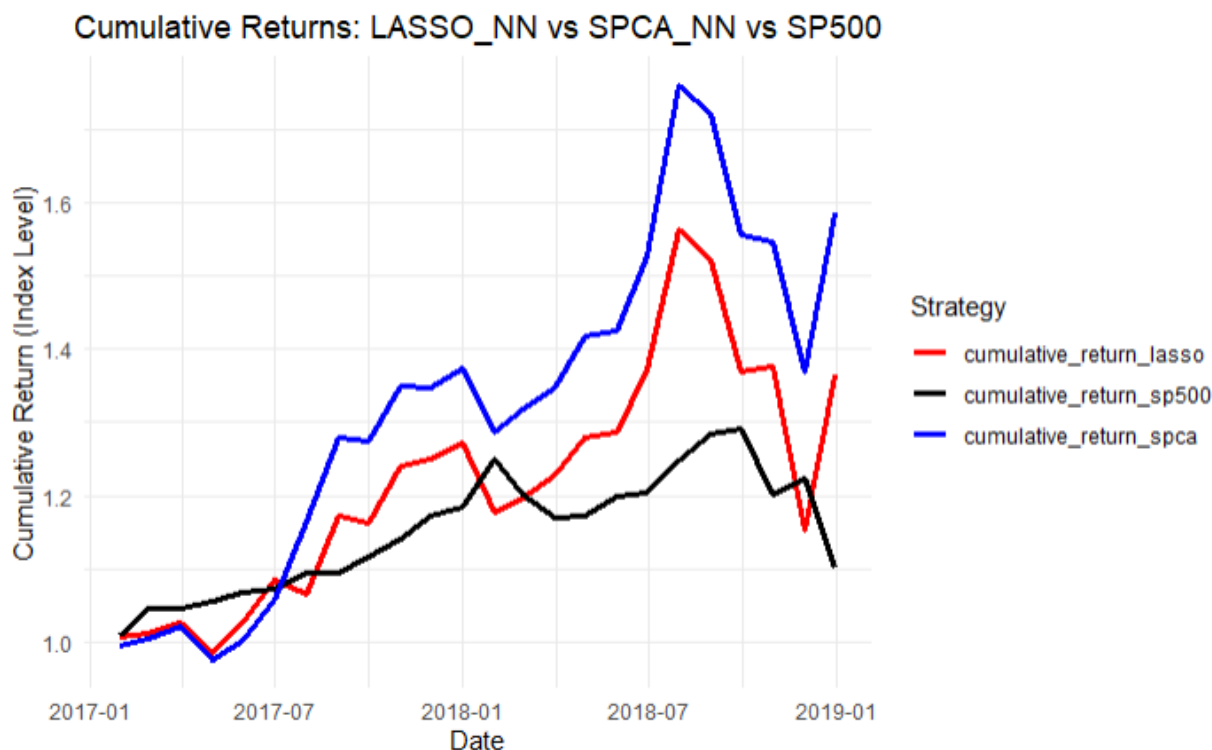


Figure 7 - Cumulative Returns Comparison between Three Portfolios

Related Work

“Projects 1 and 2 are highly related, therefore the reader is referred to our GPR #1 report for related work. Below is a quick summary of GRP #1 and what are the differences.

In GPR #1, we focused on the application of machine learning to factor investing, especially on SPCA, LASSO Regression, and neural networks as key tools for feature extraction and return prediction. In GPR #2, we successfully implemented both SPCA and LASSO, using each as a separate pipeline to generate inputs for our neural network model.

While SPCA is often applied to macroeconomic variables in literature, in our case, we applied it to firm-specific variables only, since our datasets contains solely of firm level financial and fundamental features. Separately we use LASSO regression to select most predictive factors directly. The features from each method were then used as inputs to feedforward neural networks, so we can compare predictive power and downstream portfolio performance.

Beyond what was covered in GPR #1, we did not do a long-short strategy (i.e.. Long top 10% and short bottom 10%), instead we constructed a long-only top 20% portfolio based on predicted returns and evaluated its performance using cumulative return and Sharpe ratio against the S&P 500 benchmark.

Conclusions

This research project explores the application machine learning techniques, especially Sparse Principal Component Analysis (SPCA) combined with Neural Networks and LASSO combined with Neural Networks, to stock return prediction and portfolio construction. The primary research objective is to determine whether machine-learning-based models can deliver superior investment performance relative to passive investing strategy, and whether feature generation via SPCA can improve predictability compared to dimensionality reduction via LASSO.

The portfolio performance results based on test data provide strong evidence that machine-learning-based models can enhance portfolio construction. Both LASSO-based and SPCA-based Neural Networks strategy significantly outperforms the passive S&P500 investment over 2017-2018 period in terms of holding period returns and Sharpe ratios. This suggests that properly regularize machine learning models can extract useful predictive signals to enhance investment outcomes. When comparing between machine-learning-based models, the higher holding period returns, and Sharpe ratio generated by SPCA-based strategy indicates that feature generation via SPCA offers advantages over direct feature selection using LASSO. The SPCA transformation can reduce noise and compress more information, which contributes to higher returns and lower volatility.

The use of Neural Networks in both strategies allow the non-linear complex relationships to be captured. Although the stock return predictions remain conservative, the model’s ability to rank stocks was sufficient to construct effective portfolios that can outperform passive investing strategy. Overall, combining dimensionality reduction or feature selections techniques with Neural Networks can improve stock selection and portfolio returns beyond passive benchmarks.