

# Stock Market Trend Prediction Using Regression Errors

Omair Sandhu

Department of Electrical, Computer, and Biomedical Engineering, Ryerson University

## ABSTRACT

Stock exchanges are one of the major areas of investment because of the possibility of high returns and big winners. They are affected by a variety of factors making it difficult to get consistent returns and accurate predictions when using systematic forecasting techniques.

We consider a portfolio formation problem by comparison of the trend strengths of multiple assets. The trend strength determined by the slope and errors from the regression line provides a useful method for cross-sectional comparison of stocks.

We use weekly and monthly data from 1965 to 2018 from the CRSP US Stocks Database to test the performance of these factors when used to predict the direction of movement for an asset in the future. We investigate the feasibility of this two factor model and various methods of combination to determine the optimal stock trend forecasting model.

## INTRODUCTION

The stock market is widely known to be difficult to predict because it is affected by a very large number of factors. A popular approach is to use machine learning techniques to try and predict changes in the stock market based on a verity of features such as fundamental and technical indicators. However, highly efficient stock exchange prediction models have yet to be designed due to the high volatility in price of assets.

This study aims to use linear regression techniques to predict stock price trends by training on past data and testing on out of sample data to evaluate the performance. Once we have our predictions, we will create a portfolio with the best stocks to test our technique and also check the performance of different portfolio weighting techniques.

This research uses the closing price of stocks at the end of the stock exchanges business day and as such we will not investigate the effect on day trading and high frequency trading even though there may be opportunities for profit there.

## RELATED WORK

Using natural language processing techniques on text data to predict returns has proven quite effective achieving a sharpe ratio higher than 4. The random forest of decision trees also generates good results when considered on individual stocks and longer holding periods

## METHOD

### Data Collection

Data from CRSP US Stock Database was used containing monthly and weekly close, dollar liquidity, and market cap data for the top 2000 market cap stocks. The data is collected from 1965 to 2017 and includes delisted companies to remove survivor bias.

### Ranking Based On Regression

A portfolio of stocks was formed at time  $t$  by regressing linearly on the previous 12 price points,  $P_{t-1}$  to  $P_{t-12}$ , and ranking stocks. Before we can send the stock price data to be regressed we need to perform some data transformations such that we produce plausible and comparable results. To do this, we normalize the inputs by scaling them to the first price point,  $P_{t-12}$ , so that a broad range of prices for a stock will be proportionally comparable to that of another stocks. This will ensure that trend line slope and errors are comparable across stocks with various price points.

The transformed data is then linearly regressed. This gives us the trend line slope which we save to use later on. The slope and y-intercept (bias) are also used to calculate 12 fitted price points. Then using each of the 12 fitted price points,  $FP$ , we calculate the inverse root mean square error,  $E$ , using Equation 1 below.

$$E_{i,t} = \sqrt{\frac{12}{\sum_{j=1}^{12} (P_{i,t-j} - FP_{i,t-j})^2}}$$

Equation 1: The inverse root mean square error for stock  $i$  at time  $t$

We create 2 different scoring systems for portfolio formation.

$S_{i,t} = b_{i,t}/E_{i,t}$ , Where  $S$  is the score of stock  $i$  at time  $t$ .  
 $S_{i,t} = \alpha E_{i,t} - \beta b_{i,t}$ , Where  $S$  is the score of stock  $i$  at time  $t$  and  $\alpha$  and  $\beta$  sum to 1 and are found using a grid search on the performance of the first half of the dataset.

The scores for each stock will then be sorted from lowest to highest and the 10 highest scoring stocks will be longed while the 10 lowest scoring stocks will be shorted. Furthermore, we will simulate 3 different portfolio weighting methods on each scoring system; equal weighting, standard deviation weighting, and weighting by the sharpe ratio divided by the standard deviation.

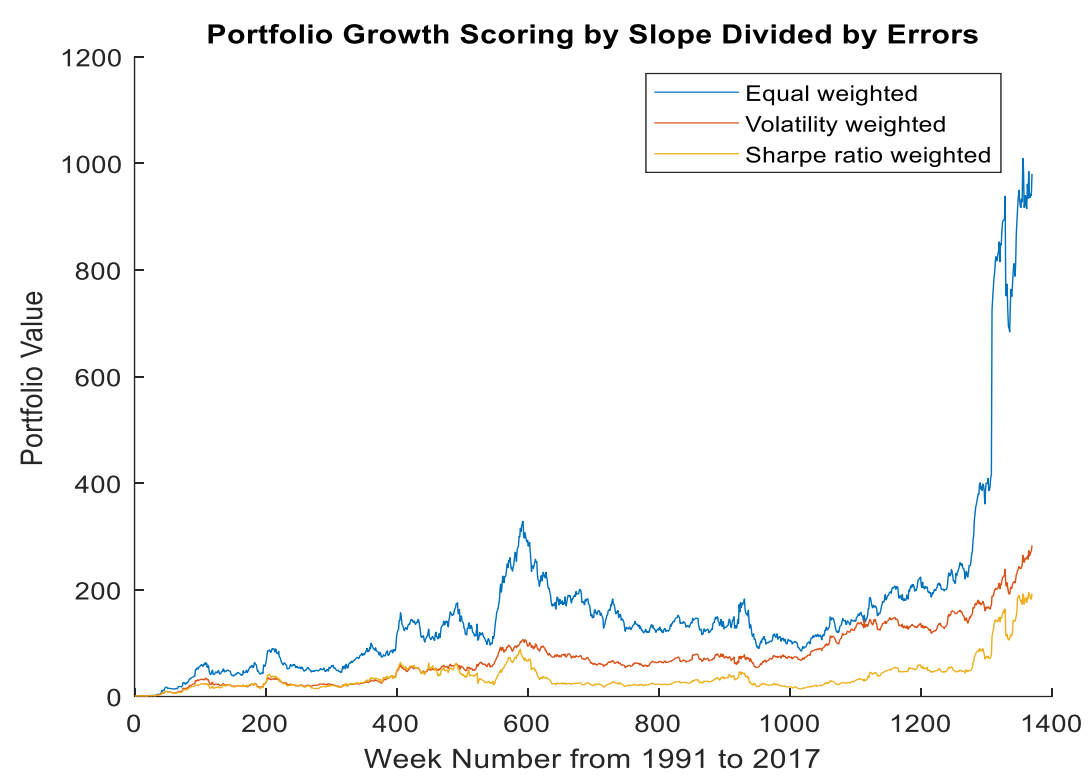
### Model Evaluation

The different models will be evaluated by comparing the annualized sharpe ratio of each model against the others. The sharpe ratio is calculated by dividing the annualized returns by the annualized volatility. We will compare the advantages and disadvantages of returns, volatility, and sharpe ratio when considering the results.

## RESULTS

### Scoring stocks based on regressed slope divided by inverse mean squared errors

The portfolio value curve when scoring and selecting stocks based on slope divided by inverse mean squared errors after linearly regressing on the previous 12 weeks prices as shown below suggests a consistent gradual increase in portfolio value over time for all weighting methods.

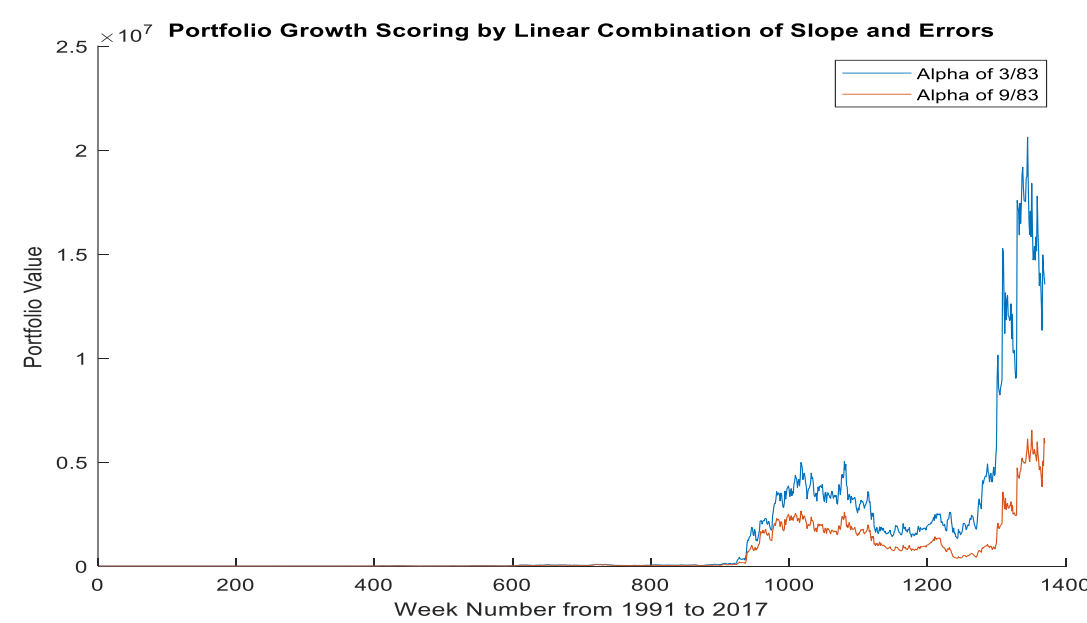


Volatility and sharpe ratio weighted portfolios do not perform as well as an equal weighted portfolio for this scoring method but in this case they still perform acceptably when returns and volatility are considered. We can see this in in the table below when we compare the sharpe ratios of the three weighting methods.

Weighting Method	Annualized returns	Annualized Volatility	Annualized Sharpe Ratio
Equal Weighted	0.3429	0.4217	0.8130
Volatility Weighted	0.2734	0.3520	0.7767
Sharpe Ratio Weighted	0.3137	0.4848	0.6470

### Scoring stocks based on a linear combination of slope and errors

We apply weights of  $\alpha$  and  $\beta$  to the errors and slope, respectively, such that  $\alpha$  and  $\beta$  sum to 1 or effectively  $\beta = 1 - \alpha$ . The optimal value for weights  $\alpha$  and  $\beta$  are found using a grid search on the dataset from 1965-01-08 to 1991-03-28 and their consistency is then tested from 1991-04-05 to 2017-06-30. The optimal value for  $\alpha$  on the back test is 9/83 with annualized returns of 0.9938, annualized volatility of 0.5615, and annualized sharpe ratio of 1.7700. The figure below shows the portfolio value over time for the back test optimal value as well as the optimal value of 3/83 found when testing from 1991-04-05 to 2017-06-30 for comparison.



Test Date Range	$\alpha = 9/83$			$\alpha = 3/83$		
	Annualized returns	Annualized Volatility	Annualized Sharpe Ratio	Annualized returns	Annualized Volatility	Annualized Sharpe Ratio
1965-01-08 to 1991-03-28	0.9938	0.5615	1.7700	1.013	0.5780	1.7517
1991-04-05 to 2017-06-30	0.9113	0.8423	1.0820	0.9521	0.8511	1.1187

With the optimal  $\alpha$  value being so close to 0, we can conclude that a linear combination of regressed slope inverse mean squared errors is not needed to significantly improve performance.

## CONCLUSION AND FUTURE WORK

The paper addresses the successful use of linearity as a scoring method with both regressed slope and inverse mean squared errors performing well individually. A fair investigation into a combination of the factors was also conducted by normalizing and standardizing the data before combination.

It is interesting to see that picking stocks with the largest negative regressed slope to long and stock with the largest positive slope to short in scoring method  $A$  produced the best results. This is likely because these stocks have a more pronounced reversal as they have the most drastic price movements when we conduct a cross sectional comparison of all stocks in our dataset. Another interesting result is that in scoring method  $B$ ; the equal weighted portfolio performs quite well while volatility and sharpe ratio weighted portfolios do not do well. In general, it appears that equal weighted portfolios perform best for linearity scoring.

Future research includes testing the use of linearity as a feature for random forest and neural network models and observing the performance compared to simple scoring methods. Additionally, the usefulness of linearity's properties (regressed slope and inverse errors) can be tested by running different feature importance tests and observing which contributes most to the results and how their contribution compares to other commonly used features such as momentum and volatility.

## ACKNOWLEDGEMENTS

The CRSP dataset was generously provided by Dr. Xiao Ping Zhang at Ryerson University. I would also like to thank Dr. Tuan Quoc Tuan for suggesting the idea of linearity as a scoring method.