

Forecasting Constituents of the MSCI Minimum Volatility Index Through Logistic
Regression

A Thesis
Presented to
The Division of Statistics
Harvard College

In Partial Fulfillment
of the Requirements for the Degree
Bachelor of Arts in Statistics (Honors)

John A. Gilheany

November 6, 2017

Approved for the Division
(Statistics)

Professor Michael Parzen

David Kane

Acknowledgements

I want to thank Prof. Parzen and David Kane for all of their help.

Preface

This thesis explores a way of predicting index constituents using logistic regression.

Table of Contents

Chapter 1: thesisdown::thesis_gitbook: default	1
Chapter 2: Introduction	3
2.1 Background	4
2.1.1 Exchange Traded Funds (ETFs)	4
2.1.2 iShares MSCI Min Vol USA ETF	4
2.1.3 Purpose	4
2.1.4 Logistic Regression Model	4
2.2 Literature Review	4
2.2.1 Overview of the Low-Risk Anomaly	4
2.2.2 Evidence of the Low-Risk Anomaly	4
2.2.3 Possible Explanations for the Low-Risk Anomaly	4
2.2.4 Further Decomposition of the Low-Risk Anomaly into Micro and Macro Effects	4
2.2.5 Real-World Applications of the Low-Risk Anomaly	4
Chapter 3: Data Gathering Process	5
3.1 Data Aggregation	5
3.2 Data Cleaning	5
3.2.1 Non-US Exchanges	5
3.2.2 Unrecognized Tickers	5
3.2.3 Price Discrepancies	5
3.3 Data Overview	5
3.4 Data Check	5
3.4.1 Weights	5
3.4.2 Comparing ETF returns to Constructed Index returns	5
Chapter 4: Tables, Graphics, References, and Labels	7
4.1 Change in 5 largest holdings by average weight for EUSA and USMV	7
4.2 Tables	8
4.3 Figures	10
4.4 Footnotes and Endnotes	12
4.5 Bibliographies	12
4.6 Anything else?	14

Conclusion	15
Chapter 5: The First Appendix	17
References	19

List of Tables

4.1	Correlation of Inheritance Factors for Parents and Child	9
-----	--	---

List of Figures

4.1	For USMV, the 5 largest holdings were VZ, T, ADP, JNJ, and MCD. Their change in weights are shown below. As we can see below, with the exception of Verizon, the holdings generally remain between 1 and 1.6 percent of the overall portfolio.	8
4.2	Reed logo	10
4.3	Mean Delays by Airline	11
4.4	Subdiv. graph	12
4.5	A Larger Figure, Flipped Upside Down	12

Abstract

The low-risk anomaly has created opportunities for arbitrage in the financial markets. As Baker et al. discuss in “Benchmarks as Limits to Arbitrage: Understanding the Low-Volatility Anomaly,” low-volatility and low-beta portfolios outperform and high-volatility and high-beta portfolios by a factor of several times due to benchmarking and lottery-preferences. The iShares MSCI USA Minimum Volatility (USMV) is an ETF tracking a minimum volatility index that was used to find data and will be used for trading arbitrage. Frazzini et al. discuss arbitrage opportunities by quantitative focused funds like AQR in “Betting Against Beta”, and this thesis explores a more advanced type of index front-running as a potential arbitrage opportunity. Data was collected from USMV from its inception in October 2011, and from EUSA, the parent ETF of USMV, from the same period until December 2016. 52-week trailing beta, 52-week trailing volatility, lagged price/book, and current index membership were calculated, and a regression model was run to quantify the relationship between current index membership and these four variables. In the model, a probabilities of index membership were calculated and an optimal cutoff was calculated to which the model would be 95% accurate of its findings of a stock to be in or out of USMV, given the historical data. Backtesting with prior data showed with a model accuracy of 95%, arbitrage opportunities of X% could be collected after each rebalancing.

Chapter 1

thesisdown::thesis_gitbook:
default

Placeholder

Chapter 2

Introduction

Placeholder

2.1 Background

2.1.1 Exchange Traded Funds (ETFs)

2.1.2 iShares MSCI Min Vol USA ETF

2.1.3 Purpose

2.1.4 Logistic Regression Model

2.2 Literature Review

2.2.1 Overview of the Low-Risk Anomaly

2.2.2 Evidence of the Low-Risk Anomaly

Measures of Risk

1929-2015

1968-2008

Critique of the Capital Asset Pricing Model

2.2.3 Possible Explanations for the Low-Risk Anomaly

Compounding

Benchmarking

Single-Period Returns

Psychological and Behavioral Factors

Profitability and Value

2.2.4 Further Decomposition of the Low-Risk Anomaly into Micro and Macro Effects

2.2.5 Real-World Applications of the Low-Risk Anomaly

Betting Against Beta (BaB)

Betting Against Correlation

Stock Price Response to Index Rebalancing

Chapter 3

Data Gathering Process

Placeholder

3.1 Data Aggregation

3.2 Data Cleaning

3.2.1 Non-US Exchanges

Mislabeled Exchanges

Mislabeled Tickers

3.2.2 Unrecognized Tickers

3.2.3 Price Discrepancies

3.3 Data Overview

3.4 Data Check

3.4.1 Weights

3.4.2 Comparing ETF returns to Constructed Index returns

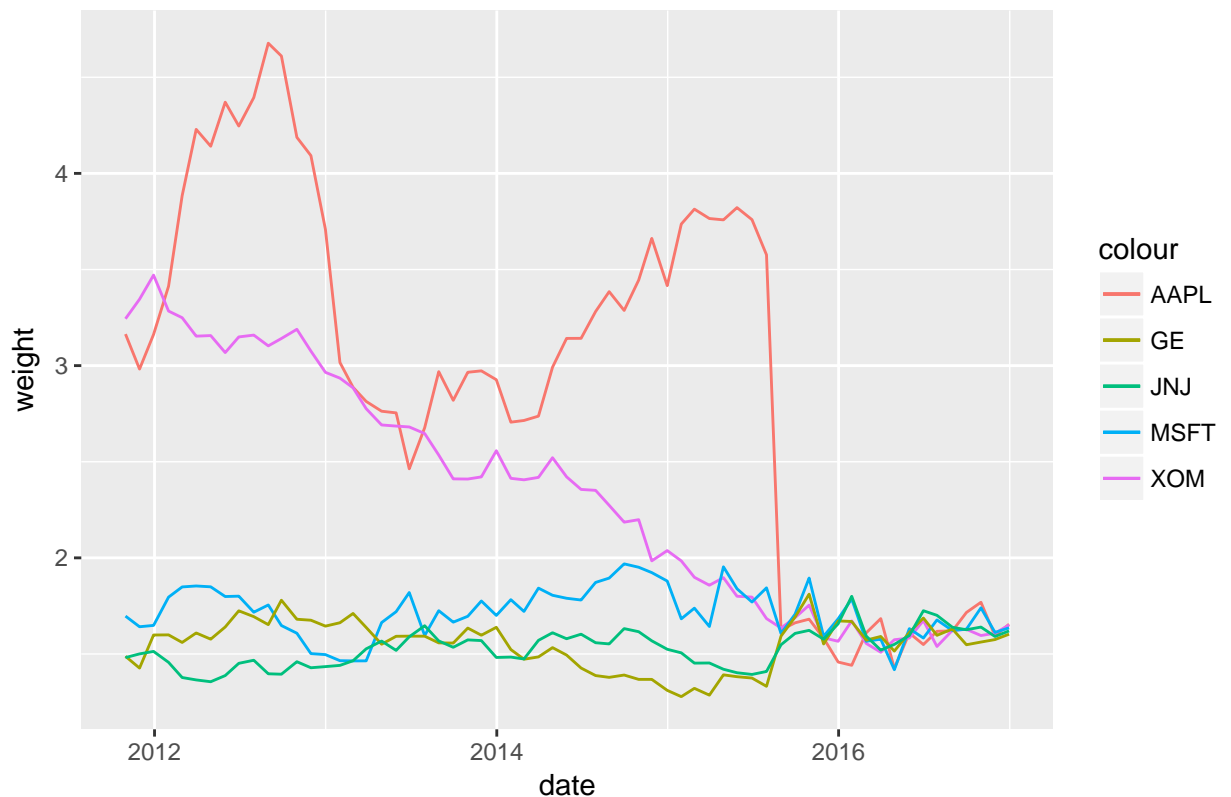
Chapter 4

Tables, Graphics, References, and Labels

4.1 Change in 5 largest holdings by average weight for EUSA and USMV

The next thing we want to see is how the top 5 largest holdings, by average weight, in each index have changed in weighting over time. For EUSA, the 5 largest holdings were AAPL, XOM, MSFT, GE, and JNJ. Their change in weights are shown below.

Change in Weights of Top 5 EUSA Holdings



Shown above, for EUSA, are some very interesting findings. The weights of the

5 companies are all very high, then suddenly all spike. Verifying this in the data, showed that for all 5 companies, holdings dropped significantly between 2015-07-31 and 2015-08-31. The reason for this is not entirely clear, but the general ETF started performing poorly around this time too. In July of 2015 the price per share was 45.20, then it dropped to 42.60 the following month, and dropped again to 40.50 in August 2015. Perhaps these large companies were doing poorly, and MSCI decided to try underweighting them.

Change in Weights of Top 5 USMV Holdings

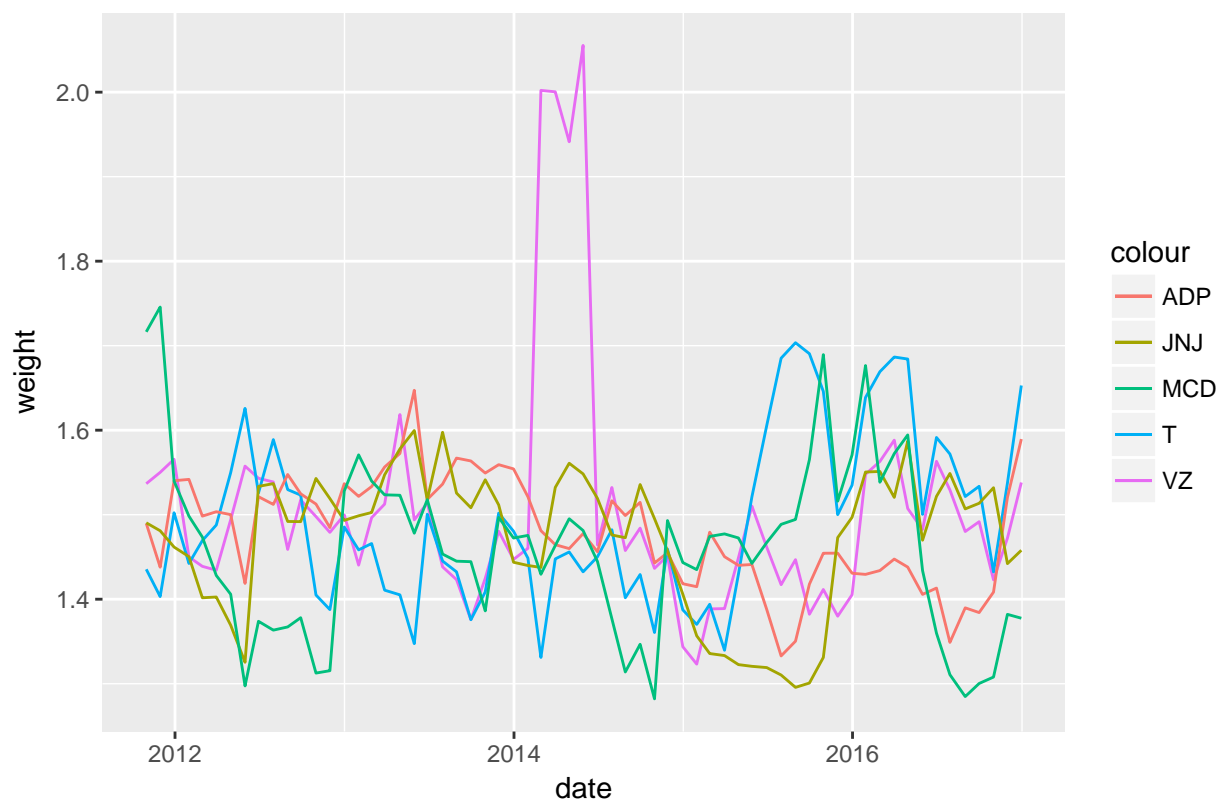


Figure 4.1: For USMV, the 5 largest holdings were VZ, T, ADP, JNJ, and MCD. Their change in weights are shown below. As we can see below, with the exception of Verizon, the holdings generally remain between 1 and 1.6 percent of the overall portfolio.

4.2 Tables

In addition to the tables that can be automatically generated from a data frame in **R** that you saw in [R Markdown Basics] using the `kable` function, you can also create tables using *pandoc*. (More information is available at <http://pandoc.org/README.html#tables>.) This might be useful if you don't have values specifically stored in **R**, but you'd like to display them in table form. Below is an example. Pay careful attention to the alignment in the table and hyphens to create the rows and columns.

Table 4.1: Correlation of Inheritance Factors for Parents and Child

Factors	Correlation between Parents & Child	Inherited
Education	-0.49	Yes
Socio-Economic Status	0.28	Slight
Income	0.08	No
Family Size	0.18	Slight
Occupational Prestige	0.21	Slight

We can also create a link to the table by doing the following: Table 4.1. If you go back to [Loading and exploring data] and look at the `kable` table, we can create a reference to this max delays table too: Table `??`. The addition of the `(\#tab:inher)` option to the end of the table caption allows us to then make a reference to Table `\@ref(tab:label)`. Note that this reference could appear anywhere throughout the document after the table has appeared.

4.3 Figures

If your thesis has a lot of figures, *R Markdown* might behave better for you than that other word processor. One perk is that it will automatically number the figures accordingly in each chapter. You'll also be able to create a label for each figure, add a caption, and then reference the figure in a way similar to what we saw with tables earlier. If you label your figures, you can move the figures around and *R Markdown* will automatically adjust the numbering for you. No need for you to remember! So that you don't have to get too far into LaTeX to do this, a couple **R** functions have been created for you to assist. You'll see their use below.

In the **R** chunk below, we will load in a picture stored as `reed.jpg` in our main directory. We then give it the caption of "Reed logo", the label of "reedlogo", and specify that this is a figure. Make note of the different **R** chunk options that are given in the R Markdown file (not shown in the knitted document).

```
include_graphics(path = "figure/reed.jpg")
```



Figure 4.2: Reed logo

Here is a reference to the Reed logo: Figure 4.2. Note the use of the `fig:` code here. By naming the **R** chunk that contains the figure, we can then reference that figure later as done in the first sentence here. We can also specify the caption for the figure via the R chunk option `fig.cap`.

Below we will investigate how to save the output of an **R** plot and label it in a way similar to that done above. Recall the `flights` dataset from Chapter ?? (Note that we've shown a different way to reference a section or chapter here.) We will next explore a bar graph with the mean flight departure delays by airline from Portland for 2014. Note also the use of the `scale` parameter which is discussed on the next page.

```
flights %>% group_by(carrier) %>%  
  summarize(mean_dep_delay = mean(dep_delay)) %>%  
  ggplot(aes(x = carrier, y = mean_dep_delay)) +  
  geom_bar(position = "identity", stat = "identity", fill = "red")
```

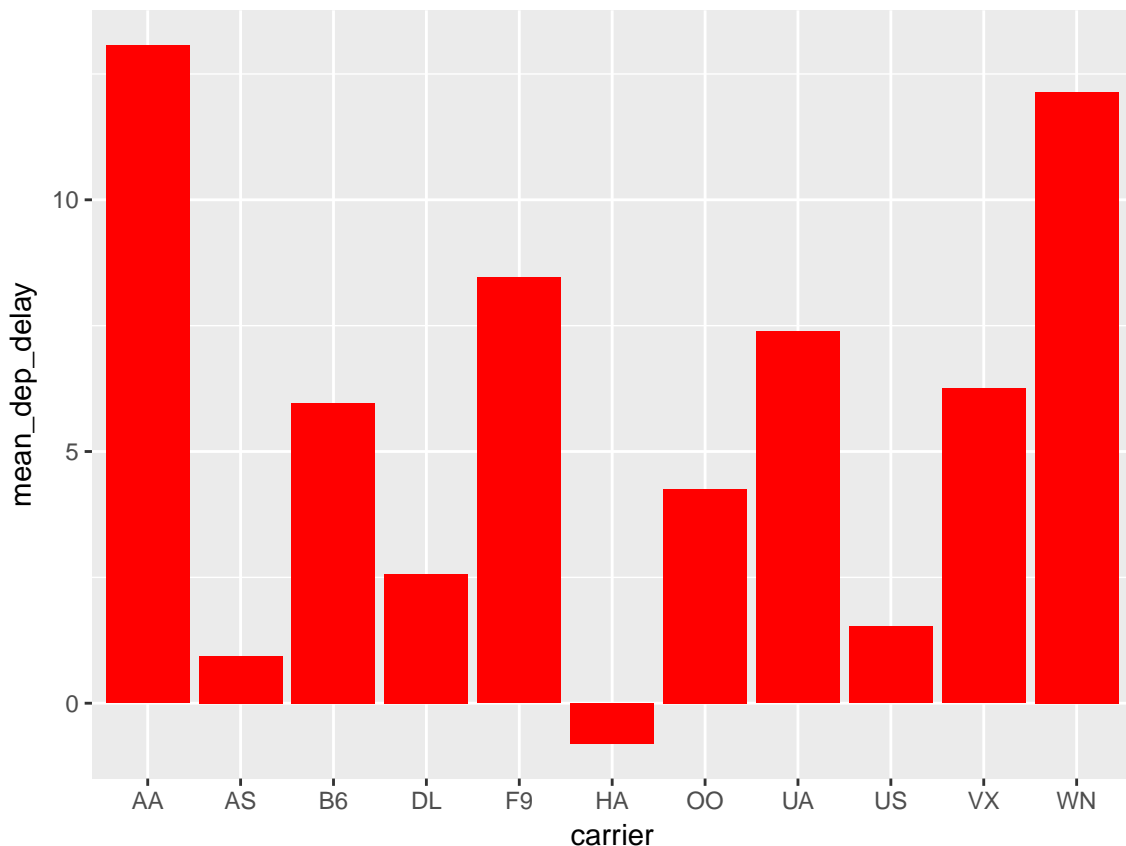


Figure 4.3: Mean Delays by Airline

Here is a reference to this image: Figure 4.3.

A table linking these carrier codes to airline names is available at <https://github.com/ismayc/pnwflights14/blob/master/data/airlines.csv>.

Next, we will explore the use of the `out.extra` chunk option, which can be used to shrink or expand an image loaded from a file by specifying `"scale= "`. Here we use the mathematical graph stored in the “subdivision.pdf” file.

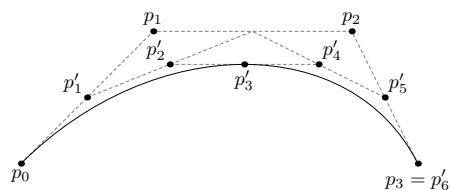


Figure 4.4: Subdiv. graph

Here is a reference to this image: Figure 4.4. Note that `echo=FALSE` is specified so that the **R** code is hidden in the document.

More Figure Stuff

Lastly, we will explore how to rotate and enlarge figures using the `out.extra` chunk option. (Currently this only works in the PDF version of the book.)

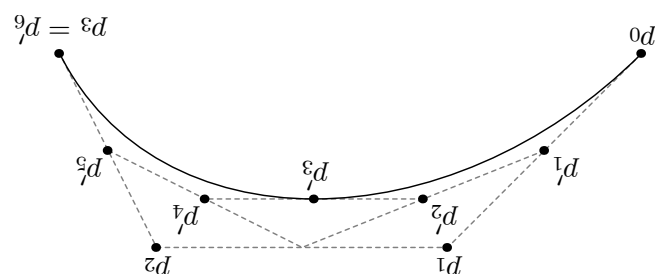


Figure 4.5: A Larger Figure, Flipped Upside Down

As another example, here is a reference: Figure 4.5.

4.4 Footnotes and Endnotes

You might want to footnote something.¹ The footnote will be in a smaller font and placed appropriately. Endnotes work in much the same way. More information can be found about both on the CUS site or feel free to reach out to `data@reed.edu`.

4.5 Bibliographies

Of course you will need to cite things, and you will probably accumulate an armful of sources. There are a variety of tools available for creating a bibliography database (stored with the `.bib` extension). In addition to BibTeX suggested below, you may want to consider using the free and easy-to-use tool called Zotero. The Reed librarians have created Zotero documentation at <http://libguides.reed.edu/>

¹footnote text

citation/zotero. In addition, a tutorial is available from Middlebury College at <http://sites.middlebury.edu/zoteromiddlebury/>.

R Markdown uses *pandoc* (<http://pandoc.org/>) to build its bibliographies. One nice caveat of this is that you won't have to do a second compile to load in references as standard LaTeX requires. To cite references in your thesis (after creating your bibliography database), place the reference name inside square brackets and precede it by the “at” symbol. For example, here's a reference to a book about worrying: (Molina & Borkovec, 1994). This `Molina1994` entry appears in a file called `thesis.bib` in the `bib` folder. This bibliography database file was created by a program called BibTeX. You can call this file something else if you like (look at the YAML header in the main `.Rmd` file) and, by default, is to placed in the `bib` folder.

For more information about BibTeX and bibliographies, see our CUS site (<http://web.reed.edu/cis/help/latex/index.html>)². There are three pages on this topic: *bibtex* (which talks about using BibTeX, at <http://web.reed.edu/cis/help/latex/bibtex.html>), *bibtexstyles* (about how to find and use the bibliography style that best suits your needs, at <http://web.reed.edu/cis/help/latex/bibtexstyles.html>) and *bibman* (which covers how to make and maintain a bibliography by hand, without BibTeX, at <http://web.reed.edu/cis/help/latex/bibman.html>). The last page will not be useful unless you have only a few sources.

If you look at the YAML header at the top of the main `.Rmd` file you can see that we can specify the style of the bibliography by referencing the appropriate csl file. You can download a variety of different style files at <https://www.zotero.org/styles>. Make sure to download the file into the `csl` folder.

Tips for Bibliographies

- Like with thesis formatting, the sooner you start compiling your bibliography for something as large as thesis, the better. Typing in source after source is mind-numbing enough; do you really want to do it for hours on end in late April? Think of it as procrastination.
- The cite key (a citation's label) needs to be unique from the other entries.
- When you have more than one author or editor, you need to separate each author's name by the word “and” e.g. `Author = {Noble, Sam and Youngberg, Jessica},.`
- Bibliographies made using BibTeX (whether manually or using a manager) accept LaTeX markup, so you can italicize and add symbols as necessary.
- To force capitalization in an article title or where all lowercase is generally used, bracket the capital letter in curly braces.
- You can add a Reed Thesis citation³ option. The best way to do this is to use the `phdthesis` type of citation, and use the optional “type” field to enter “Reed thesis” or “Undergraduate thesis.”

²Reed College (2007)

³Noble (2002)

4.6 Anything else?

If you'd like to see examples of other things in this template, please contact the Data @ Reed team (email data@reed.edu) with your suggestions. We love to see people using *R Markdown* for their theses, and are happy to help.

Conclusion

Placeholder

Chapter 5

The First Appendix

Placeholder

References

Placeholder

Molina, S. T., & Borkovec, T. D. (1994). The Penn State worry questionnaire: Psychometric properties and associated characteristics. In G. C. L. Davey & F. Tallis (Eds.), *Worrying: Perspectives on theory, assessment and treatment* (pp. 265–283). New York: Wiley.

Noble, S. G. (2002). *Turning images into simple line-art* (Undergraduate thesis). Reed College.

Reed College. (2007, march). LaTeX your document. Retrieved from <http://web.reed.edu/cis/help/LaTeX/index.html>