

My Final College Paper

A Thesis
Presented to
The Division of Mathematics and Natural Sciences
Reed College

In Partial Fulfillment
of the Requirements for the Degree
Bachelor of Arts

Your R. Name

May 20xx

Approved for the Division
(Mathematics)

Advisor F. Name

Acknowledgements

I want to thank a few people.

Preface

This is an example of a thesis setup to use the reed thesis document class.

Table of Contents

Introduction	1
Chapter 1: Chapter 1	3
1.1 1.1 Trees and Random Forests	3
1.1.1 Trees	3
1.1.2 Random Forests	4
1.2 1.x What We Mean When We Talk About Inference	4
1.3 1.x Permutations and Populations	5
1.4 1.x Inference on Random Forests	5
1.4.1 The Problem	5
1.5 Permuatations Tests Theory and Application to Conditional Variable Importance	6
Chapter 2: Tables, Graphics, References, and Labels	7
2.1 Tables	7
2.2 Figures	8
2.3 Footnotes and Endnotes	11
2.4 Bibliographies	11
2.5 Anything else?	13
Conclusion	15
Appendix A: The First Appendix	17
Appendix B: The Second Appendix, for Fun	19
References	21

List of Tables

2.1	Correlation of Inheritance Factors for Parents and Child	7
-----	--	---

List of Figures

2.1	Reed logo	8
2.2	Mean Delays by Airline	9
2.3	Subdiv. graph	10
2.4	A Larger Figure, Flipped Upside Down	10
2.5	Subdivision of arc segments	11

Abstract

The preface pretty much says it all.

Second paragraph of abstract starts here.

Dedication

You can have a dedication here if you wish.

Introduction

Welcome to the *R Markdown* thesis template. This template is based on (and in many places copied directly from) the L^AT_EX template, but hopefully it will provide a nicer interface for those that have never used T_EX or L^AT_EX before. Using *R Markdown* will also allow you to easily keep track of your analyses in **R** chunks of code, with the resulting plots and output included as well. The hope is this *R Markdown* template gets you in the habit of doing reproducible research, which benefits you long-term as a researcher, but also will greatly help anyone that is trying to reproduce or build onto your results down the road.

Hopefully, you won't have much of a learning period to go through and you will reap the benefits of a nicely formatted thesis. The use of L^AT_EX in combination with *Markdown* is more consistent than the output of a word processor, much less prone to corruption or crashing, and the resulting file is smaller than a Word file. While you may have never had problems using Word in the past, your thesis is likely going to be about twice as large and complex as anything you've written before, taxing Word's capabilities. After working with *Markdown* and **R** together for a few weeks, we are confident this will be your reporting style of choice going forward.

Why use it?

R Markdown creates a simple and straightforward way to interface with the beauty of L^AT_EX. Packages have been written in **R** to work directly with L^AT_EX to produce nicely formatting tables and paragraphs. In addition to creating a user friendly interface to L^AT_EX, *R Markdown* also allows you to read in your data, to analyze it and to visualize it using **R** functions, and also to provide the documentation and commentary on the results of your project. Further, it allows for **R** results to be passed inline to the commentary of your results. You'll see more on this later.

Who should use it?

Anyone who needs to use data analysis, math, tables, a lot of figures, complex cross-references, or who just cares about the final appearance of their document should use *R Markdown*. Of particular use should be anyone in the sciences, but the user-friendly nature of *Markdown* and its ability to keep track of and easily include figures, automatically generate a table of contents, index, references, table of figures, etc. should make it of great benefit to nearly anyone writing a thesis project.

Chapter 1

Chapter 1

1.1 1.1 Trees and Random Forests

1.1.1 Trees

Decision trees may be familiar to many with a background in the social or medical sciences as convenient ways to represent data and can assist in decision making. Morgan and Sonquist (1963) derived a way for constructing trees motivated by the specific feature space of data collected from interviews and surveys. Unlike, say agricultural data which involves mostly numerical variables like rainfall, the data collected from interviews is mostly categorical. On top of this issue, the datasets Morgan and Sonquist dealt with had few participants (n) and much data collected on them (p). To continue with their list of difficulties, there was reason to believe that there were lurking errors in the variables that would be hard to identify and quantify. Lastly, many of the predictors were correlated and Morgan and Sonquist doubted that the additive assumptions of many models would be appropriate for this data. Morgan and Sonquist noted that while many statistical methods would have a difficult time accurately parsing this data, a clever researcher with quite a lot of time could create a suitable model simply by grouping values in the feature space and predicting that the response corresponding to these values would be the mean of the observed responses given the grouped conditions. Their formalization of this procedure in terms of “decision rules” laid the ground work for future research on decision trees.

In 1984, Breiman et al introduces a revolutionary new algorithm for trees. **Need to acquire *Classification and Regression Trees* to make sure the method discussed in MASS is the same that Breiman uses/is used in randomForest**

Tree Algorithm CART?

Begin by considering the entire feature space X_1, \dots, X_n . Then:

1. Consider every possible pair of partitions of this feature space, P_1, P_2 , so that if $X_1 = x_1, X_2 = x_2, \dots, X_n = x_n$ where $x_1, \dots, x_n \in P_1$ then our prediction is the mean value of y given $x_1, \dots, x_n \in P_1$.
2. Choose the partitions that minimize RSS

3. For each new partition, repeat steps 1 and 2 until some stopping condition is reached.

An alternative to this method is conditional inference trees. Torsten Hothorn, Kurt Hornik, Achim Zeileis argue in their 2006 paper **Unbiased Recursive Partitioning: A Conditional Inference Framework**, CART has a selection bias toward variables with either missing values or a great number of possible splits. This bias can effect the interpretability of all tree models fit using this method. As an alternative to CART and other algorithms, Hothorn et al propose a new method, conditional inference trees.

1. For case weights w test the global null hypothesis of independence between any of the m covariates and the response. Stop if this hypothesis cannot be rejected. Otherwise select the j_{th} covariate X_j with strongest association to Y .
2. Choose a set $A \subset X_j$ in order to split X_j into two disjoint sets A and $X_j \setminus A$. The case weights w_{left} and w_{right} determine the two subgroups with $w_{left,i} = w_i I(X_{j,i} \in A)$ and $w_{right,i} = w_i I(X_{j,i} \notin A)$ for all $i = 1, \dots, n$ ($I(\cdot)$ denotes the indicator function).
3. Recursively repeat steps 1 and 2 with modified case weights w_{left} and w_{right} , respectively.

from <https://eeecon.uibk.ac.at/~zeileis/papers/Hothorn+Hornik+Zeileis-2006.pdf>

After step 1 is completed, any goodness of fit method can be used to generate the split and choose the set A . Note that in this method the splitting is done separately from the variable selection.

1.1.2 Random Forests

There is a limit to the predictive capabilities of a single tree; they suffer from high variance. To alleviate this, random forests are often used instead. They function by enlisting the help of many trees, and then by aggregating the responses over all of them.

- history
- algorithm
- uses

1.2 1.x What We Mean When We Talk About Inference

- Inferential vs Descriptive
- Frequentist vs Bayesian

1.3 1.x Permutations and Populations

As stated in the introduction of the *Chronical of Permutations Statistical Methods* by KJ Berry et al, 2014, there are two models of statistical inference. One is the population model, where we assume that the data was randomly sampled from one (or more) populations. Under this model, we assume that the data generated follows some known distribution. “Under the population model, the level of statistical significance that results from applying a statistical test to the results of an experiment or a survey corresponds to the frequency with which the null hypothesis would be rejected in repeated random samplings from the same specified population(s)”, (Berry et al, 2014).

The permutation family of methods, on the other hand, only assumes that the observed result was caused by experimental variability.

1.4 1.x Inference on Random Forests

1.4.1 The Problem

Random forests create models with great predictive-, but poor inferential capabilities. A single tree is simple to i ### Proposed solutions to this problem

Statisticians Leo Breiman and ____ Cutler proposed a method of permuted variable importance that hoped to answer this problem. Their method compares the variable importance for each variable in a tree-wise manner. For each tree, the permuted variable importance of the variable X_j is:

$$PV^t(x_j) = \frac{\sum_{i \in |B|} y - \hat{y}^t}{|B|} - \frac{\sum_{i \in |*B|} y - *\hat{y}^t}{|*B|}$$

Where B is the matrix representing the feature space, $|B|$ is the number of observations, $*B$ is the matrix of predictors but with X_j permuted, \hat{y} is the predicted outcome, and $*\hat{y}^t$ is the predicted outcomes after variable X_j has been permuted. This value is averaged over all the trees. It’s important to note that if the variable X_j is not split on in the tree t , the tree-wise variable importance will be 0.

Creating a permutation-based method is certainly an attractive solution to our problem. One, it allows us to estimate the distribution of variable importance and generate a Z score under the null hypothesis that $PV = 0$.

$$PV(x_j) = \frac{\sum_1^n treePV^t(x_j)}{\frac{\hat{\sigma}}{\sqrt{ntree}}}$$

Strobl et al from the University of Munich criticize this method in their 2008 technical report, *Danger: High Power! – Exploring the Statistical Properties of a Test for Random Forest Variable Importance*. One, this method has the downside of increasing power with increasing numbers of trees in the forest. This is a more or less arbitrary parameter which we would hope would not affect our importance estimates. Secondly, the null hypothesis under Breiman and Cutler’s strategy is that the variable

importance V for any variable X_j is not equal to zero given Y , the response. Because random forests are most often used in situations with multicollinearity that would make other methods like the linear model difficult, Strobl argues that any variable importance measure worth its salt should not be misled by correlation within the predictors.

The researchers at the University of Munich published a fully fleshed response to the Breiman and Cutler method in 2008, titled *Conditional Variable Importance for Random Forests* that address these issues. Strobl et al propose restructuring the Breiman and Cutler algorithm to account for conditional dependence among the predictors. Their algorithm looks like this:

1. Fit a random forest to the model, R_0 , and calculate base variable importance for each variable V
2. For every predictor $X_j \in X_1, \dots, X_n$:
 - 2a. Conditionally permute X_j given the splits found in R_0
 - 2b. Fit a new random forest R_j with the permuted data
 - 2c. Calculate a new variable importance \hat{V}_j
3. For every variable X_1, \dots, X_n ,

$$CV(X_j) = \hat{V}_j - V_j$$

The null hypothesis is that $CV(X_j) = 0$ given the predictor Y and all other predictors X_1, \dots, X_n . This accounts for interactions between X_j and the other predictors. Using the simulated data from the previous example, here's an implementation of the algorithm outlined here as it is in the **party** package.

This paper aims to provide a response to this method. One the conditional permutation algorithm is notoriously slow with any dataset of a size that is appropriate for a random forest. Two, the partitions are made from the random forest corresponding to the formula of $Y \sim X_1, \dots, X_n$ instead of a model of $X_j \sim X_1, \dots, X_n$.

1.5 Permutatations Tests Theory and Application to Conditional Variable Importance

Chapter 2

Tables, Graphics, References, and Labels

2.1 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 2.1: 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 2.1. If you go back to [Loading and exploring data] and look at the `kable` function code, you'll see that I added in a similar `\\label` to be able to reference that table later. (The extra backslash there is a way that *Markdown* interfaces with *L^AT_EX*.) We can create a reference to the max delays table: ??.

The addition of the `\\label{}` option to the end of the table caption allows us to then use the *L^AT_EX* `autoref` function to produce the link. The `ref` function in **R** allows for tables and figures to be referenced in the document easily without having to directly use the `autoref` function. It will automatically add "Table" before your number if you add the "tab:" prefix to your label. Note that this reference could appear anywhere throughout the document.

2.2 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 L^AT_EX 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 "reed", and specify that this is a figure. Note again the use of the `results = "asis"` specification to automatically include and compile the L^AT_EX code.

```
label(path = "figure/reed.jpg", caption = "Reed logo",  
      label = "reed", type = "figure")
```



Figure 2.1: Reed logo

Here is a reference to the Reed logo: Figure 2.1. Note the use of the inline **R** code here. By default "figure" is specified as the type. For clarity, we could have also added the `label` and `type` to the parameter specifications and this would give us the same result: Figure 2.1.

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 . (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.

```
#delay_airline <- 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")
#ggsave("figure/delays.pdf", plot = delay_airline,
#   # height = 3, width = 6)
```

```
label(path = "figure/delays.pdf",
      caption = "Mean Delays by Airline",
      label = "delays", type = "figure")
```

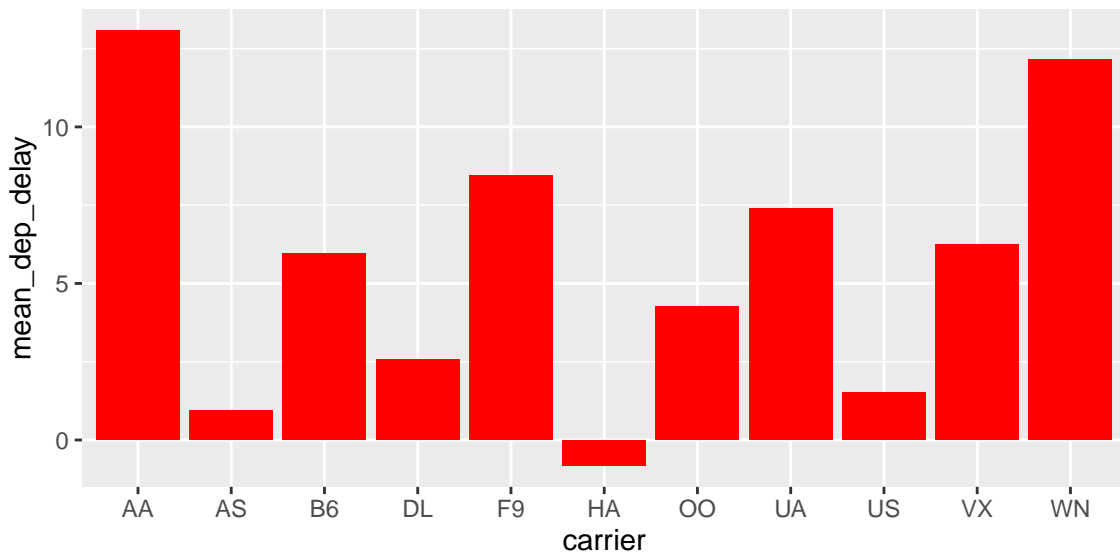


Figure 2.2: Mean Delays by Airline

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 `scale` parameter which can be used to shrink or expand an image. Here we use the mathematical graph stored in the “subdivision.pdf” file. Note that we didn’t specify the `caption =` or `label =` here, but we could have.

```
label("figure/subdivision.pdf", "Subdiv. graph", "subd",
      scale = 0.75)
```

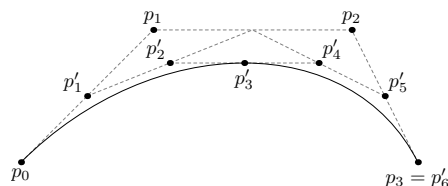


Figure 2.3: Subdiv. graph

Here is a reference to this image: Figure 2.3. (Move this around throughout the document as you wish.)

More Figure Stuff

Lastly, we will explore how to rotate figures using the `angle` parameter.

```
label("figure/subdivision.pdf",
      "A Larger Figure, Flipped Upside Down",
      scale = 1.5,
      angle = 180,
      label = "subd2")
```

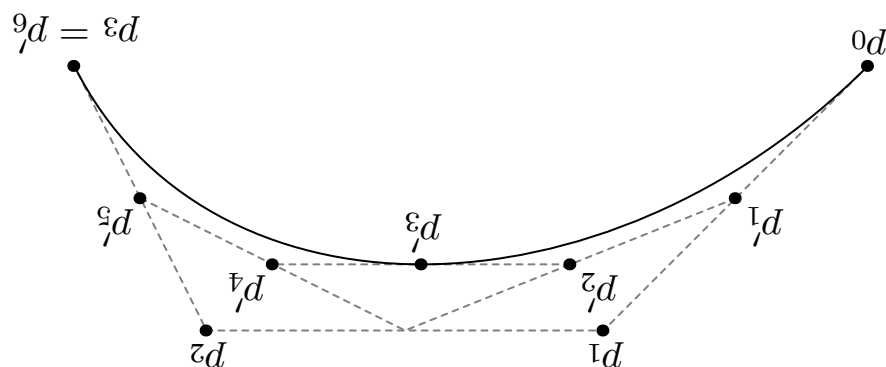


Figure 2.4: A Larger Figure, Flipped Upside Down

As another example, here is a reference to this figure: Figure 2.4.

Common Modifications

The following figure features the more popular changes thesis students want to make to their figures. We can add math to the caption that displays below the picture, specify the size of our caption to display below the figure (list of sizes available at this link), and also specify that a different caption `alt.cap` be what appears in the Table of Figures for this figure.

If you'd like to make further tweaks to figures, you might need to invoke some \LaTeX code. Please email us at `data@reed.edu` if you need assistance.

```
label("figure/subdivision.pdf",
      caption = "Subdivision of arc segments",
      alt.cap = "You can see that  $p_3 = p_6^{\prime}$ ",
      cap.size = "footnotesize",
      label = "subd3")
```

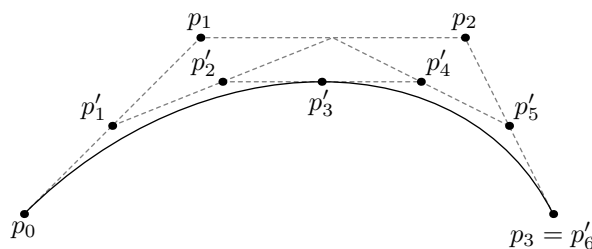


Figure 2.5: You can see that $p_3 = p_6'$

2.3 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`.

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

¹footnote text

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 L^AT_EX 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 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 L^AT_EX 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)

2.5 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

If we don't want Conclusion to have a chapter number next to it, we can add the `{.unnumbered}` attribute. This has an unintended consequence of the sections being labeled as 3.6 for example though instead of 4.1. The \LaTeX commands immediately following the Conclusion declaration get things back on track.

More info

And here's some other random info: the first paragraph after a chapter title or section head *shouldn't be* indented, because indents are to tell the reader that you're starting a new paragraph. Since that's obvious after a chapter or section title, proper typesetting doesn't add an indent there.

Appendix A

The First Appendix

This first appendix includes all of the R chunks of code that were hidden throughout the document (using the `include = FALSE` chunk tag) to help with readability and/or setup.

In the main Rmd file:

```
# This chunk ensures that the reedtemplates package is  
# installed and loaded. This reedtemplates package includes  
# the template files for the thesis and also two functions  
# used for labeling and referencing  
if(!require(devtools))  
  install.packages("devtools", repos = "http://cran.rstudio.com")  
if(!require(reedtemplates)){  
  library(devtools)  
  devtools::install_github("ismayc/reedtemplates")  
}  
library(reedtemplates)
```

In :

```
# This chunk ensures that the reedtemplates package is  
# installed and loaded. This reedtemplates package includes  
# the template files for the thesis and also two functions  
# used for labeling and referencing  
if(!require(devtools))  
  install.packages("devtools", repos = "http://cran.rstudio.com")  
if(!require(dplyr))  
  install.packages("dplyr", repos = "http://cran.rstudio.com")  
if(!require(ggplot2))  
  install.packages("ggplot2", repos = "http://cran.rstudio.com")
```

```
if(!require(reedtemplates)){  
  library(devtools)  
  devtools::install_github("ismayc/reedtemplates")  
}  
library(reedtemplates)  
#flights <- read.csv("data/flights.csv")
```

Appendix B

The Second Appendix, for Fun

References

- Angel, E. (2000). *Interactive computer graphics : A top-down approach with opengl*. Boston, MA: Addison Wesley Longman.
- Angel, E. (2001a). *Batch-file computer graphics : A bottom-up approach with quicktime*. Boston, MA: Wesley Addison Longman.
- Angel, E. (2001b). *Test second book by angel*. Boston, MA: Wesley Addison Longman.
- 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>