



# Introduction to Data Analysis

*(First Edition)*

**By Matt Teachout**

**College of the Canyons  
Santa Clarita, CA, USA**



*Introduction to Data Analysis, first edition by M. Teachout, College of the Canyons, Santa Clarita, CA, USA, and is licensed under a [Creative Commons Attribution 4.0 International license](https://creativecommons.org/licenses/by/4.0/) – 10/1/17*

**Special thanks to all of the people that made this book possible.**

**Most especially all of the *Intermediate Algebra for Statistics* students and teachers at College of the Canyons. Your work to pioneer this material and improve statistics education continues to inspire education reform.**

**I would also like to thank Kathy Kubo, Joe Gerda, Dustin Silva, Ambika Silva, Ralph (Randy) Ades, Udani Ranasinghe, and Rupa Sinha, Myra Snell and Katie Hern at the California Acceleration Project, James Glapa-Grossklag and Brian Weston and the COC OER office staff, The incredible COC math department, reprographics, and the rest of the COC faculty and staff.**

**Your help and support made this possible.**

## Introduction

We live in the age of computers and the internet. We are exposed to huge volumes of data every day. How do we make sense of this massive amount of information? How can we tell the difference between helpful and misleading information? How can businesses know what their customers want and need, or hospitals analyze various types of infections and which treatments are working and which are not? All of these questions revolve around the study of data and statistics. A good understanding of statistics is vital to anyone living in the modern world, however very few people understand how to analyze data. The shortage of trained statisticians, data analysts, and data scientists is a huge problem worldwide.

There are many fabulous books on statistics and analyzing data. Unfortunately, they are extremely expensive and most people cannot afford the cost. I wrote this book to help people learn to analyze data. It is free to use the material in this book, update it, add to it, print it or just read it. It is an open educational resource (OER) and so anyone can use it.

Many college students struggle to balance work and family with their education. One of the biggest roadblocks for many students is the cost of textbooks. Students today cannot afford the cost of textbooks and so chose to attend classes without purchasing books and materials needed for the class. It goes without saying, that this is a major impediment to passing their classes, but the students have no choice. They simply cannot afford \$150-\$200 textbooks. For this reason, I believe strongly in open educational resources (OER). Open source materials like this book are available and are virtually free for students.

## Notes about OER and Creative Commons Licensing

This textbook is licensed through Creative Commons as “Attribution CC-BY”. Creative Commons describes this license as follows: “This license lets others distribute, remix, tweak, and build upon (the author’s) work, even commercially, as long as they (give) credit (to the author) for the original creation.” This is the most accommodating of licenses offered. Recommended for maximum dissemination and use of licensed materials.” If you need to see the license deed or legal code, please go to <https://creativecommons.org/licenses/> and look under the “CC-By” section.

## Pre-Statistics or Intermediate Algebra for Statistics

I tell my beginning statistics students all the time that the study of statistics is a deep well of knowledge, and they are only playing in the puddle. Statisticians, data analysts and data scientists are life-long learners and spend years and years studying this subject.

This is an introduction to some very basic data analysis techniques. It is a book designed for anyone new to statistics. It can be used with a pre-statistics class or an intermediate algebra for statistics class.

Pre-statistics classes focus on helping students understand and analyze categorical and quantitative data sets.

Intermediate algebra for statistics has the same information as a pre-statistics class but often includes some intermediate algebra curve analysis and regression techniques. Many statisticians and statistics educators feel that curve analysis and regression is a topic better addressed in more advanced level statistics classes since this is a topic explored by many graduate level statistics students.

If your college requires intermediate algebra for statistics, I have included that material in chapter 6. If your college is using a pre-statistics class, then chapters 1-5 should suffice.

## Important Note about Technology

We live in the age of computers, internet and a huge volume of data. No practicing statistician or data scientist uses a calculator or tables to analyze data. You cannot even begin to analyze a data set of 100,000 values by hand with a calculator. You need high-powered computer software. There are many statistics software programs on the market, but very few of them are free.

If you read the history of statistics, you will find brilliant scientists, mathematicians and people in business who had to try to figure out data, but had no access to a computer. (Computers had not been invented yet.) Our pioneers of statistics dreamed of the day that they could compute statistics and graphs and analyze data with the touch of a button. They invented complicated techniques for analyzing data because they had no choice. Today, computers can calculate statistics and graphs directly.

Here is the problem. Most statistics classes taught in high schools, community colleges and even some universities are teaching statistics as if computers have not been invented yet. They are teaching the techniques developed by our pioneers of statistics before the computer age. This is a terrible approach to the subject, especially for the thousands of students that actually want to work in the field. A statistics class should be a study of how to practically collect and analyze data with a computer, not a class on what to do if computers have not been invented yet.

Are formulas important in statistics? Yes. We look at formulas to understand what they tell us about the data and the world around us. The pioneers of statistics did an amazing job of addressing the major ideas of statistics with formulas and inventive calculations. However, we should not use a formula and a calculator to calculate a statistic for a data set with 10,000 values or use charts that list critical values and P-values. High-powered computers with statistics software can calculate the statistic and make graphs directly. Then students can focus on the analysis part, and explore and discover the meaning behind the data.

This book will show students how to use statistics software to calculate statistics and graphs. I want students to learn to analyze the data and not spend all their time just trying to calculate something. Remember, no one pays a data analyst to calculate something a computer can already do. They are paid to explore and explain what the data may be telling us.

## Statcato

Teaching statistics with computer software is very important, but many schools and students cannot afford to pay for software. For this reason, I used the free statistics software Statcato throughout the book. Statcato is a JAVA based software, which can be saved to any computer (MAC or PC) and is easy to use.

Statcato Website: [www.statcato.org](http://www.statcato.org)

Statcato Download for MAC: <http://mac.softpedia.com/get/Utilities/Statcato.shtml>

Statcato Download for PC: <https://sourceforge.net/projects/statcato/>

You can of course use the book with any statistics software. Most basic statistics software programs are very similar to Statcato. You will just need to supplement software directions.

## Data Sets

The national (GAISE) guidelines for teaching statistics recommend that you use real data. Allowing students to learn statistics principles through analysis of real data is key. With that being said, there are many places where raw data can found and used. The key data sets throughout this book are located at my website under “Int Alg for Stats” and “Data Sets”.

Matt Teachout’s Website: [www.matt-teachout.org](http://www.matt-teachout.org)

## The Computer Dilemma

A statistics or pre-statistics class should be taught in a computer lab. It is important to allow the computers to do the difficult calculations. Students need to focus on interpretation and discovering the meaning behind the data. They cannot do that if they spend all their time trying to calculate with a formula or making graphs by hand.

If your school wants to teach statistics or pre-statistics, but you cannot teach in a computer lab, here are some suggestions for you.

1. Reserve unused computer labs. Some schools may have a couple computer labs that are not always in use. Schedule your statistics and pre-statistics classes in order to use the computer lab. Even if you can only reserve the lab once a week or once every two weeks, it will be a huge help to students.
2. Have groups of students share computers. If you do have a few computers in your classroom, you can divide the class up into groups and share computers. This has many advantages like encouraging explanations to one another and teamwork.
3. Teachers can use their own computer or laptop to project statistics software on a screen and have the class analyze the data with you. Teachers without any computer can make printed copies of the software printouts for your class and for exams. It is a poor substitute for a computer lab, but it is much better than teaching statistics as if computers have not been invented yet.

<b>INTRODUCTION TO DATA ANALYSIS .....</b>	<b>1</b>
<b>Chapter 1 – Categorical Data Analysis .....</b>	<b>9-43</b>
Section 1A – Two Types of Data – Categorical and Quantitative .....	10-11
Section 1B – Proportions and Percentages .....	12-16
Section 1C – Pie Charts and a Bar Charts with Technology .....	17-26
Section 1D – Comparing Percentages from Multiple Groups.....	27-34
Section 1E – Using Percentage Data.....	35-38
Chapter 1 Review .....	38-42
Project Chapter 1 - Categorical Data Analysis Group Poster .....	42-43
<b>Chapter 2 – Relationships between Categorical Variables .....</b>	<b>44-85</b>
Section 2A – Two-Way Tables with Technology .....	45-52
Section 2B – Using Bar Charts and Pie Charts to Summarize Two-Way Tables .....	53-60
Section 2C – Marginal and Joint Percentages from Two-Way Tables .....	61-65
Section 2D – Conditional Percentages and Categorical Relationships .....	66-79
Chapter 2 Review .....	80-83
Project Chapter 2 - Two-Way Tables & Categorical Relationships .....	84-85
<b>Chapter 3 – Analyzing Normal Quantitative Data .....</b>	<b>86-142</b>
Section 3A – Finding the Shape of a Quantitative Data Set with Dot Plots and Histograms .....	87-94
Section 3B – Shapes and Centers .....	95-101
Section 3C – Understanding the Mean Average.....	102-107
Section 3D – Introduction to Spread, Standard Deviation, and Typical Values for Normal Data .....	108-121
Section 3E – Unusual Values in Normal Data, Using the Dot Plot, and Summarizing Quantitative Data .....	122-134
Chapter 3 Review .....	135-140
Project Chapter 3 - Quantitative Data Analysis Poster for Bell Shaped Data .....	141-142
<b>Chapter 4 – Analyzing Skewed Quantitative Data .....</b>	<b>143-194</b>
Section 4A – Review of Shapes and Centers, Creating Histograms and Dot Plots with Technology .....	144-146
Section 4B – Understanding the Median Average.....	147-152
Section 4C – Understanding Spread for Skewed Data, Quartiles, Interquartile Range (IQR), and the Five Number Summary .....	153-161
Section 4D – Box Plots, Typical and Unusual Values for Skewed Data .....	162-174
Section 4E – Summary Report Paragraph for Skewed Data .....	175-184
Section 4F – Measures of Center, Spread and Position.....	185-187
Chapter 4 Review.....	188-193
Project Chapter 4 – Skewed Data Analysis Group Poster .....	193-194

<b>Chapter 5 – Linear Quantitative Relationships .....</b>	<b>195-278</b>
Section 5A – Introduction to Quantitative Relationships, Explanatory and Response Variables, Scatterplots with Technology .....	196-206
Section 5B – Strength and Direction of Linear Relationships and the Correlation Coefficient “r” .....	207-215
Section 5C – Confounding Variables, r-squared, Correlation is not Causation, and Multivariable Studies...	216-227
Section 5D – Best Fit Regression Line with Technology with Slope and Y-intercept Interpretation .....	228-247
Section 5E – Residuals, Standard Deviation of the Residual Errors (Se), Residual Plots and Histogram of the Residuals with Technology .....	248-270
Chapter 5 Review.....	270-277
Project Chapter 5 - Linear Quantitative Relationships .....	277-278
<b>Chapter 6 – Curved Quantitative Relationships .....</b>	<b>279-326</b>
Section 6A – Exponential Relationships with Technology .....	282-293
Section 6B – Logarithmic Relationships with Technology .....	294-303
Section 6C – Quadratic Relationships with Technology .....	304-317
Chapter 6 Review.....	318-324
Project Chapter 6 - Curved Quantitative Relationships.....	325-326