

Psychological Statistics and Psychometrics Using Stata

SCOTT A. BALDWIN
Brigham Young University



STATA® *Press*

A Stata Press Publication
StataCorp LLC
College Station, Texas



Copyright © 2019 StataCorp LLC
All rights reserved. First edition 2019

Published by Stata Press, 4905 Lakeway Drive, College Station, Texas 77845

Typeset in L^AT_EX 2_ε

Printed in the United States of America

10 9 8 7 6 5 4 3 2 1

Print ISBN-10: 1-59718-303-2

Print ISBN-13: 978-1-59718-303-1

ePub ISBN-10: 1-59718-304-0

ePub ISBN-13: 978-1-59718-304-8

Mobi ISBN-10: 1-59718-305-9

Mobi ISBN-13: 978-1-59718-305-5

Library of Congress Control Number: 2019935130

No part of this book may be reproduced, stored in a retrieval system, or transcribed, in any form or by any means—electronic, mechanical, photocopy, recording, or otherwise—without the prior written permission of StataCorp LLC.

Stata, **stata**, Stata Press, Mata, **mata**, and NetCourse are registered trademarks of StataCorp LLC.

Stata and Stata Press are registered trademarks with the World Intellectual Property Organization of the United Nations.

NetCourseNow is a trademark of StataCorp LLC.

L^AT_EX 2_ε is a trademark of the American Mathematical Society.

To Autumn

Contents

List of figures	xvii	
List of tables	xxiii	
Acknowledgments	xxv	
Notation and typography	xxvii	
I	Getting oriented to Stata	1
1	Introduction	3
1.1	Structure of the book	3
1.2	Benefits of Stata	4
1.3	Scientific context	5
2	Introduction to Stata	9
2.1	Point-and-click versus writing commands	9
2.2	The Stata interface	10
2.3	Getting data in Stata	12
2.4	Viewing and describing data	14
2.4.1	list, in, and if	15
2.5	Creating new variables	17
2.5.1	Missing data	18
2.5.2	Labels	19
2.6	Summarizing data	20
2.6.1	summarize	20
2.6.2	table and tabulate	21
2.7	Graphing data	22
2.7.1	Histograms	23
2.7.2	Box plots	24

2.7.3	Scatterplots	25
2.8	Reproducible analysis	27
2.8.1	Do-files	28
2.8.2	Log files	30
2.8.3	Project Manager	30
2.8.4	Workflow	31
2.9	Getting help	32
2.9.1	Help documents	32
2.9.2	PDF documentation	33
2.10	Extending Stata	33
2.10.1	Statistical Software Components	33
2.10.2	Writing your own programs	33
II	Understanding relationships between variables	35
3	Regression with continuous predictors	37
3.1	Data	38
3.2	Exploration	38
3.2.1	Demonstration	40
	Simulation program	40
3.3	Bivariate regression	42
3.3.1	Lines	43
3.3.2	Regression equation	43
3.3.3	Estimation	45
3.3.4	Interpretation	46
	Slope	46
	Intercept	47
3.3.5	Residuals and predicted values	48
3.3.6	Partitioning variance	51
3.3.7	Confidence intervals	54
3.3.8	Null hypothesis significance testing	60

<i>Contents</i>	ix
3.3.9	Additional methods for understanding models 63
	Using predicted scores to understand model implications . . 64
	Composite contrasts 73
3.4	Conclusions 74
4	Regression with categorical and continuous predictors 75
4.1	Data for this chapter 76
4.2	Why categorical predictors need special care 77
4.3	Dummy coding 78
4.3.1	Example: Incorrect use of categorical variable 85
4.4	Multiple predictors 86
4.4.1	Interpretation 86
	Model fit 86
	Intercept 87
	Slopes 88
4.4.2	Unique variance 89
4.5	Interactions 90
4.5.1	Categorical by continuous interactions 91
	Dichotomous by continuous interactions 91
	Polytomous by continuous interactions 101
	Joint test for interactions with polytomous variables 105
4.5.2	Continuous by continuous interactions 107
4.6	Summary 109
5	t tests and one-way ANOVA 111
5.1	Data 112
5.2	Comparing two means 112
5.2.1	t test 114
5.2.2	Effect size 115
5.3	Comparing three or more means 116
5.3.1	Analysis of variance 116

5.3.2	Multiple comparisons	120
	Planned comparisons	122
	Direct adjustment for multiple comparisons	126
5.4	Summary	129
6	Factorial ANOVA	131
6.1	Data for this chapter	132
6.2	Factorial design with two factors	134
6.2.1	Examining and visualizing the data	134
6.2.2	Main effects	138
	Testing the null hypothesis	139
6.2.3	Interactions	139
6.2.4	Partitioning the variance	140
6.2.5	2 x 2 source table	142
6.2.6	Using anova to estimate a factorial ANOVA	144
6.2.7	Simple effects	146
6.2.8	Effect size	149
6.3	Factorial design with three factors	151
6.3.1	Examining and visualizing the data	152
6.3.2	Marginal means	154
6.3.3	Main effects and interactions	155
6.3.4	Three-way interaction	158
6.3.5	Fitting the model with anova	159
6.3.6	Interpreting the interaction	160
6.3.7	A note about effect size	165
6.4	Conclusion	166
7	Repeated-measures models	167
7.1	Data for this chapter	169
7.2	Basic model	172

7.3	Using mixed to fit a repeated-measures model	175
7.3.1	Covariance structures	176
	Compound symmetry (exchangeable)	177
	First-order autoregressive	180
	Toeplitz	183
	Unstructured	186
7.3.2	Degrees of freedom	189
7.3.3	Pairwise comparisons	190
7.4	Models with multiple factors	192
7.5	Estimating heteroskedastic residuals	197
7.6	Summary	200
8	Planning studies: Power and sample-size calculations	201
8.1	Foundational ideas	202
8.1.1	Null and alternative distributions	202
8.1.2	Simulating draws out of the null and alternative distributions	204
8.2	Computing power manually	210
8.3	Stata's commands	214
8.3.1	Two-sample z test	214
8.3.2	Two-sample t test	215
8.3.3	Correlation	220
8.3.4	One-way ANOVA	223
8.3.5	Factorial ANOVA	226
8.4	The central importance of power	229
8.4.1	Type M and S errors	230
	Type S errors	231
	Type M errors	233
8.5	Summary	235
9	Multilevel models for cross-sectional data	237
9.1	Data used in this chapter	239

9.2	Why clustered data structures matter	239
9.2.1	Statistical issues	239
9.2.2	Conceptual issues	242
9.3	Basics of a multilevel model	244
9.3.1	Partitioning sources of variance	244
9.3.2	Random intercepts	246
9.3.3	Estimating random intercepts	248
9.3.4	Intraclass correlations	250
9.3.5	Estimating cluster means	252
	Comparing pooled and unpooled means	259
9.3.6	Adding a predictor	262
9.4	Between-clusters and within-cluster relationships	264
9.4.1	Partitioning variance in the predictor	265
9.4.2	Total- versus level-specific relationships	266
9.4.3	Exploring the between-clusters and within-cluster relationships	267
9.4.4	Estimating the between-clusters and within-cluster effects	270
9.5	Random slopes	273
9.6	Summary	283
10	Multilevel models for longitudinal data	285
10.1	Data used in this chapter	286
10.2	Basic growth model	287
10.2.1	Multilevel model	295
10.3	Adding a level-2 predictor	300
10.4	Adding a level-1 predictor	307
10.5	Summary	310
III	Psychometrics through the lens of factor analysis	313
11	Factor analysis: Reliability	315
11.1	What you will learn in this chapter	316
11.2	Example data	316

11.3	Common versus unique variance	317
11.4	One-factor model	320
11.4.1	Parts of a path model	321
11.4.2	Where do the latent variables come from?	321
11.5	Prediction equation	323
11.6	Using sem to estimate CFA models	326
11.7	Model fit	328
11.7.1	Computing χ^2	330
11.8	Obtaining σ_C^2 and σ_U^2	334
11.8.1	Computing R^2 for an item	335
11.8.2	Computing σ_C^2 and σ_U^2 for all items	338
11.8.3	Computing reliability— ω	340
11.8.4	Bootstrapping the standard error and 95% confidence interval for ω	341
11.9	Comparing ω with α	343
11.9.1	Evaluating the assumption of tau-equivalence	346
11.9.2	Parallel items	349
11.10	Correlated residuals	350
11.11	Summary	353
12	Factor analysis: Factorial validity	355
12.1	Data for this chapter	357
12.2	Exploratory factor analysis	358
12.2.1	Common factor model	359
12.2.2	Extraction methods	360
12.2.3	Interpreting loadings	362
12.2.4	Eigenvalues	363
12.2.5	Communality and uniqueness	364
12.2.6	Factor analysis versus principal-component analysis	365

12.2.7	Choosing factors and rotation	366
	How many factors should we extract?	366
	Eigenvalue-greater-than-one rule	366
	Scree plots	367
	Parallel analysis	368
	Orthogonal rotation—varimax	370
	Oblique rotation—promax	374
12.3	Confirmatory factor analysis	377
12.3.1	EFA versus CFA	377
12.3.2	Estimating a CFA with sem	380
12.3.3	Mean structure versus variance structure	381
12.3.4	Identifying models	383
	Imposing constraints for identification	384
	How much information is needed to identify a model?	385
12.3.5	Refitting the model with constrained latent variables	386
12.3.6	Standardized solutions	389
12.3.7	Global fit	392
	RMSEA	393
	TLI	394
	CFI	394
	SRMR	394
	A summary and a caution	395
12.3.8	Refining models further	395
12.4	Summary	396
13	Measurement invariance	397
13.1	Data	398
13.2	Measurement invariance	398
13.3	Measurement invariance across groups	400
13.3.1	Configural invariance	400
13.3.2	Metric invariance	407

13.3.3	Scalar invariance	407
13.3.4	Residual invariance	408
13.3.5	Using the comparative fit index to evaluate invariance . . .	409
13.4	Structural invariance	412
13.4.1	Invariant factor variances	412
13.4.2	Invariant factor means	412
13.5	Measurement invariance across time	413
13.5.1	Configural invariance	413
	Effects coding for identification	416
	Effects-coding constraints in Stata	417
13.5.2	Metric invariance	424
13.5.3	Scalar invariance	425
13.5.4	Residual invariance	426
13.6	Structural invariance	427
13.7	Summary	428
	References	429
	Author index	441
	Subject index	445

(Pages omitted)

1 Introduction

Statistics is a major part of research in many fields, often serving as the primary method for establishing whether results support hypotheses. Indeed, it is unusual for a quantitative study in psychology to not use statistical analysis. Consequently, it's essential for psychologists to be competent users of statistics, and that is the primary aim of this book.

That is a lofty goal. Statistics is a huge, technical field. Psychologists who spend their careers studying statistics only master small portions, and most psychologists study psychological content, not statistics. Furthermore, studying a single book will not make you a competent user of statistics. Competence comes from trial and error and from applying statistical methods and ideas to your own research. Competence comes from years of building your skills. Competence comes from collaborating with more competent people and learning from them. My hope is that this book serves as a step toward competence for those starting their training and as a useful reference or development tool for those already using statistics.

1.1 Structure of the book

This book covers some foundational topics in psychological statistics and psychometrics. Topics include t tests, analysis of variance (ANOVA), regression, power analysis, multilevel models, structural equation modeling, and factor analysis. I do not provide an exhaustive treatment of each topic because many books already do that. Instead, I provide an introduction to key concepts regarding how to use these statistical models to answer research questions.

Competent statistical analysis requires the use of statistical software. Consequently, I weave the statistical content with Stata code that illustrates how to fit models and make sense of the output. Most of the code examples illustrate how to use Stata commands, such as `regress` (for regression) or `sem` (for structural equation modeling). I also illustrate how to use some of Stata's graphical commands, such as `histogram` and `twoway scatter` for scatterplots, and data management commands, such as `generate` and `egen` (to create new variables). I show how to program Stata to run simulations to help you learn a concept or as a method for understanding your models.

The book is divided into three parts. Part I provides an introduction to Stata, including the interface, loading data, do-files and Stata syntax, descriptive statistics, graphics, and help and documentation. Part II discusses regression, ANOVA, multilevel

models, and statistical power. I begin with regression because it provides the foundation for the discussion of ANOVA, multilevel models, and power. Part III covers psychometrics, including a discussion of the issues of reliability and validity. In these chapters, I discuss these concepts from the perspective of factor analysis, because I believe this provides a coherent framework for thinking about these measurement concepts. Looking at these concepts from the perspective of factor analysis will also help you make connections between the psychometric concepts of reliability and validity and regression models.

I *strongly* recommend that you work through the Stata code in each chapter rather than just reading the code and seeing the results in the book. You will learn statistical concepts better by toying with code and seeing what happens, especially as you change the specification or the options. My students tell me that it is more effective to practice using Stata and to run models during class rather than just watch me use Stata or read lecture notes. Furthermore, as you get more comfortable with Stata, you will find it much easier to learn how to use other Stata commands because you will start to see the connections between the various commands.

1.2 Benefits of Stata

Statistical software is necessary, but why choose Stata? Some treat statistical software kind of like a sports team, showing undying devotion to the software and viewing criticisms of the software as personal attacks. Others are pragmatic users, simply using what is available to them or what seems useful. Some want only open source software; others want a company to back the development and certification of the software. Many use a specific software because that's what their adviser used or what was taught in their program. I started using Stata because I was a longtime Mac user. As I was finishing my graduate degree in 2003, SAS did not support Mac at all, and SPSS's Mac software was buggy and slow. Stata's Mac support was excellent, and the software did most of what I needed. I had been introduced to Stata as an undergraduate and was happy to return. I have never looked back. Stata has remained my primary statistical package.¹

Although my reasons for choosing Stata had little to do with statistics itself, I believe Stata is an excellent choice for five reasons:

1. **Consistent syntax.** Learning Stata, like learning any programming language, is challenging. However, learning Stata is manageable because Stata's syntax is expressive and consistent. It is expressive because it is easy to understand once you have the basics down. It is consistent because command syntax does not fundamentally change from one command to another. Thus, once you get the basic structure down, it is easy to learn how to use new commands: the structure will be similar to commands you are already familiar with.

1. That's not to say that I do not use other software. Sometimes, collaboration requires that I do. Sometimes, I need analysis routines not available in Stata. And sometimes, I just like to learn new stuff, so I try software that is available to me.

2. **Aids the replicability of analyses.** Statistical analyses should be reproducible (Long 2009), which means that if I asked you to run an analysis a second time (or third, fourth, and so on), you will get the same answer. Likewise, if I ran your code on my computer, I would get the same results. Stata includes the ability to specify the Stata version number for your analysis. For example, you may start your analysis with the command `version 14.2`, which tells Stata to run the analysis using the code that was present for Stata version 14.2. Therefore, if something changed between version 14.2 and 15.0, you will not get different answers.²
3. **Comprehensive documentation.** Stata’s support documentation is comprehensive and ships as part of the software. Additionally, all documentation is freely available on the web. In my opinion, Stata documentation is second-to-none in the statistical software world. The documentation provides a readable explanation of each command, including options. Furthermore, the documentation includes worked examples and discussion of output to aid your learning. Finally, the documentation provides technical details regarding the mathematical and statistical underpinnings of commands. You can learn a lot by studying the documentation.
4. **Data management.** Most analyses require data management: cleaning data, generating variables, labeling variables and values, reshaping the data into a specific format, and so on. Stata’s data management commands and utilities are excellent and make data preparation straightforward and replicable. Even when I need to use another software package for analysis, I nearly always prepare and manage the data in Stata.
5. **Graphics.** Graphics are an essential part of data analysis and are often superior to tables of numbers when it comes to communicating results (Cox 2004; Gelman, Pasarica, and Dodhia 2002; Gelman 2011). Stata includes comprehensive graphical tools to aid in exploring your data and interpreting the results of your models. The graphics are customizable and can be quite beautiful. Nearly all figures in this book were created with Stata. I include the code for creating these figures so that you can reproduce the graphs. Study this code. You will appreciate how flexible Stata graphs can be, and you will learn about Stata programming.

1.3 Scientific context

As I write this introduction, psychology, and science generally, has some problems. Pressure to publish and to ensure it is something exciting and novel combined with bad methodological practices means that a lot of research is not replicable (Ioannidis 2005, 2008, 2012). Indeed, many psychological studies simply do not replicate (Open Science Collaboration 2015), leading some to call the situation a “replication crisis”

2. Of course, if the Stata developers caught a bug in the code that results in a different answer in version 15 than 14.2, you may want the different answer. Nevertheless, by making the version number explicit, you can ensure that the results will only change when you expect them to.

(Pashler and Harris 2012). Psychology is not alone in this, with some pointing to problems in other disciplines, such as economics, biology, and medicine (Ioannidis 2005, 2013, 2014). Do these problems mean that many (or even most) theories are not worthwhile and the scientific literature cannot be trusted? I hope it is not that bad. Regardless of the answer to that question, I think the replication crisis does suggest we ought to step back and think about why we are facing these problems.

Some argue that a major reason for these problems is the incentives that influence scientists (Baldwin 2017; Ioannidis 2014; Nelson, Simmons, and Simonsohn 2012; Simmons, Nelson, and Simonsohn 2011). For example, in universities across the world, promotion and tenure depend upon publications—hence the phrase “publish or perish”. Getting your first academic job often depends upon having many publications, including some in prestigious journals. Some research positions are “soft money” jobs, which means that salary and research support comes from grant money rather than the university itself. Getting grants requires that you are productive and that previous grants worked out. Sometimes we joke that to get a grant, we have to do all the research that the grant is proposing to prove that the research will work. When your salary and reputation depend upon getting papers published and the research turning out in a specific way, the incentive is to make sure the research works out as predicted.

Given these incentives, publishing becomes the end goal of research. That is, rather than publishing being the means to communicate scientific observation, publishing and adding lines to your vita become tantamount to science itself. The book *The Compleat Academic*³ provides advice to researchers in the early stages of their careers on things like graduate school, applying for postdocs and jobs, submitting grants, and teaching classes. The advice to graduate students states:

The information that we need to arrive at a short list of applicants is contained in the letters of recommendation and, primarily, in the academic vita. Wise graduate students, therefore, will start at day one of their first year in a PhD program to develop a strong vita. [...] Alter your perspective so that you derive your professional self-respect entirely from what is on that document. From the start of graduate school on, throughout what we hope will be a long and productive career, you *are* your vita. (Lord 2004, 10, emphasis in original)

Given such advice, combined with the incentives for getting and keeping a job (including securing your own salary!), it is not difficult to see why publishing became equivalent to doing science.

I learned statistics in this context, as did most researchers before me. Consequently, a number of problematic research and statistical practices evolved that ultimately helped publication rates but did not improve the quality of the science. For example, consider the use of p -values to evaluate statistical significance. There are many criticisms of

3. The dictionary on my computer defines the word “Compleat” as “archaic spelling of complete”. Leave it to academics to take the simple word *complete* and make it snooty.

p -values in the scientific literature (for example, Meehl [1978]). Many of the criticisms are about how p -values are used, not so much about p -values themselves. McElreath (2016) says it well:

This audience accepts that there is something vaguely wrong about typical statistical practice in the early 21st century, dominated as it is by p -values and a confusing menagerie of testing procedures. [...] The problem in my opinion is not so much p -values as the set of odd rituals that have evolved around them, in the wilds of the sciences, as well as the exclusion of so many other useful tools. (pp. xi–xii)

In short, p -values became the primary source of evidence that a result is publishable. Consequently, the goal of analysis becomes finding a p -value that is less than 0.05. If you complete a study and find null results, there's a good chance you will not even try to publish it.

Researcher flexibility, especially with respect to design, analysis, and reporting, means that finding significant effects probably required torturing the data. More advice from *The Compleat Academic*⁴ explains this well:

To compensate for this remoteness from our participants, let us at least become familiar with the record of their behavior: the data. Examine them from every angle. Analyze the sexes separately. Make up composite indexes. If a datum suggests a new hypothesis, try to find additional evidence for it elsewhere in the data. If you see dim traces of interesting patterns, try to reorganize the data to bring them into bolder relief. If there are participants you do not like, or trials, observers, or interviewers who gave you anomalous results, drop them (temporarily). Go on a fishing expedition for something—anything—interesting. (Bem 2004, 187)

“Interesting patterns” here usually means small p -values. Such flexibility in analysis is sometimes called researcher degrees of freedom or p -hacking (Simmons, Nelson, and Simonsohn 2011) or the garden of forking paths (Gelman and Loken 2014). Examples of such analyses include running multiple experiments and only reporting results from those with significant results. If main effects are not significant, test interactions; dropping observations when doing so takes the p -value from $p = 0.09$ to $p = 0.03$. Change a continuous variable to a categorical variable because the categorical variable produces significant results. Ignore problems with estimation and fit because a result is statistically significant. Fail to look at the raw data to see what a model (which is a reduction of the data) implies about the data because the analysis is statistically significant (Simmons, Nelson, and Simonsohn 2011). I think you get the idea.

4. Honestly, I cannot get over that name.

Fortunately, changes are in process. For example, the Open Science Framework (<https://osf.io/>) provides tools for researchers to register hypotheses before seeing the data and to create a website for hosting data and analysis files. I have seen job postings where an emphasis on improving rigor and replicability of science is a job qualification. Journals are accepting registered reports, wherein studies are reviewed prior to data collection and evaluated solely on the basis of the research question and quality of the design and proposed analyses (<https://cos.io/rr/>).

I hope this book contributes to the positive changes. My goal is to teach statistical concepts and software in a way that helps researchers a) address their research question transparently and openly, b) better understand their data, and c) better understand the models they use and what the models imply about their theory or research area. To be sure, I teach and use p -values—they can be useful. Stata includes many tools that supplement the information provided by p -values, and as a consequence, Stata can be used in a way that improves how statistics are used in psychology specifically and in science generally. So buckle up! We have a lot of great stuff to discuss.

(Pages omitted)

3.3 Bivariate regression

The mean was our best guess regarding the value of `new_attitude`. Another term meaning best guess is *expected value*. That is, in the long run, what do we expect the value of `attitude` to be? For a single variable, such as `attitude`, the expected value is the mean. Regression, on the other hand, allows us to determine an expected value for `attitude` that incorporates information from other variables. Our best guess about `attitude` will be based on variables such as the respondent's level of education and mental health symptoms.

Figure 3.2 is a scatterplot of `educ` and `attitude`.

```
. use http://www.stata-press.com/data/pspus/gss_attitude, clear
. twoway scatter attitude educ, jitter(2) scheme(lean2)
```

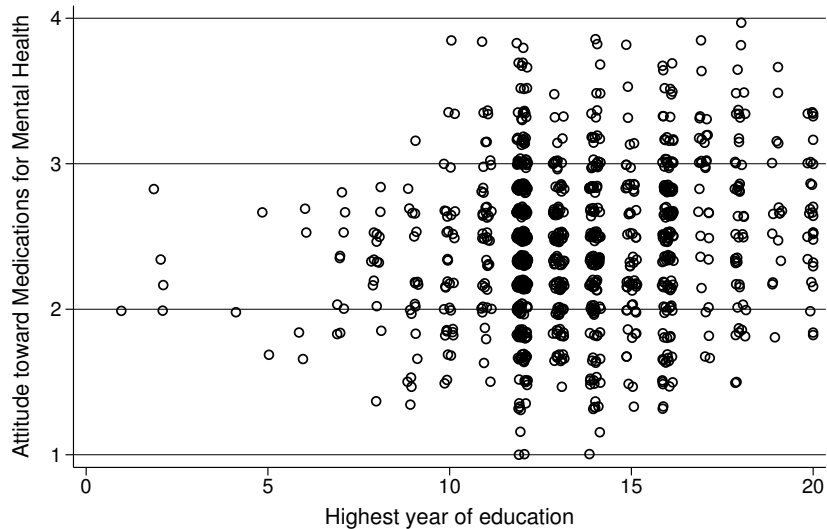


Figure 3.2. Relationship between `educ` and `attitude`

Based on figure 3.2, what do you think is the expected value for `attitude`? Does that expected value differ for someone who has 10 years of education compared with someone who has 20 years of education? That is, if we consider `educ`, does that affect what we expect `attitude` to be? If there is a relationship between education and attitude toward psychotropic medication, then expected values of `attitude` will vary as a function of `educ`. Regression provides information about the nature of that relationship in the form of a line, which is why regression is also called linear regression.

3.3.1 Lines

Let's go back in time and revisit junior high school algebra. You may recall the formula for a line:

$$y = mx + b$$

where x and y are variables, m is the slope, and b is the y intercept. We say that y is a function of two quantities: i) the product of m and x and ii) b . The slope describes changes in y with respect to changes in x (that is, "rise over run"), and the y intercept provides the value of y when $x = 0$. If we know m , x , and b , then we know y (see figure 3.3). Regression produces slope and intercept values that describe the relationship between an outcome and predictors.

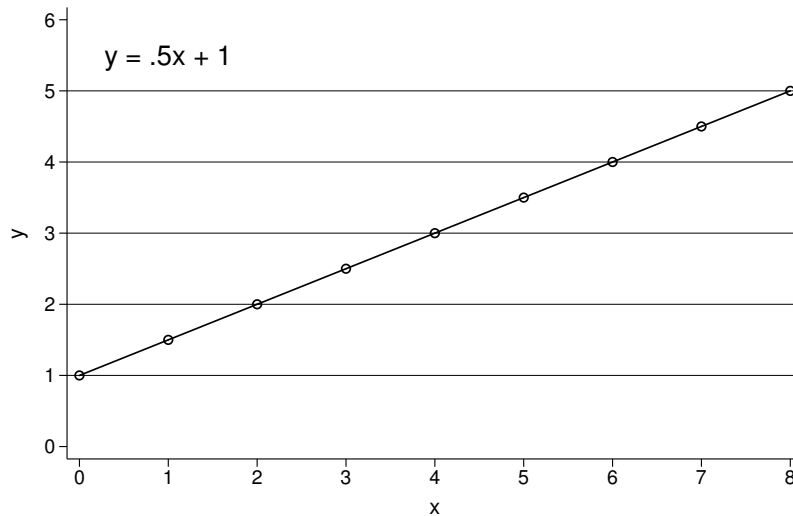


Figure 3.3. Illustration of $y = mx + b$

3.3.2 Regression equation

Compare figures 3.2 and 3.3. In figure 3.3, it is obvious what the line should be because the points all fall on a straight line. In figure 3.2, you could draw any number of lines through the data—a horizontal line, a vertical line, a line moving up from left to right, or a line moving down from left to right. All lines would provide some indication of the relationship between **attitude** and **educ**, and all lines will provide an intercept and slope. Which line is best? The best line is the line that provides the closest correspondence between the expected value of y and the actual value of y . In other words, the regression line will provide the best prediction of y given the x values across the entire dataset.

A regression equation for the relationship between `attitude` and `educ` is

$$\text{attitude}_i = \beta_0 + \beta_1 \text{educ}_i + \epsilon_i \quad (3.1)$$

This equation says that the observed value of `attitude` for person i is a function of an intercept (β_0), a slope (β_1) relating `educ` to `attitude`, and error in prediction (ϵ_i). Error is included in (3.1) because the observed value of `attitude` is not perfectly captured by β_0 and β_1 (that is, the data do fall on a straight line).

Commonly, the regression equation is written in terms of the expected value of y . For the example model, the equation is

$$E(\text{attitude}_i | \text{educ}_i) = \beta_0 + \beta_1 \text{educ}_i \quad (3.2)$$

which is read as, “The expected value of `attitude` given `educ` is equal to β_0 plus β_1 times `educ`.” Alternatively, we can substitute $E(\text{attitude}_i | \text{educ}_i)$ with $\widehat{\text{attitude}}_i$:

$$\widehat{\text{attitude}}_i = \beta_0 + \beta_1 \text{educ}_i \quad (3.3)$$

The “hat” over attitude_i denotes the predicted value of `attitude`. Equations (3.2) and (3.3) omit ϵ_i because these equations deal with the expected or predicted values of `attitude`—errors come into play when comparing the expected value to the actual value of y .

A general bivariate regression (one y and one x) equation is

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad (3.4)$$

Further, the general form of (3.2) and (3.3) is

$$E(y_i | x_i) = \beta_0 + \beta_1 x_i \quad (3.5)$$

and

$$\widehat{y}_i = \beta_0 + \beta_1 x_i \quad (3.6)$$

Population parameters versus sample-based estimates

In research, we collect data on samples to learn about populations. In regression, we estimate slopes and intercepts based on samples to learn about what the slopes and intercepts are in the broader population. Slopes and intercepts in the population are called population parameters. Parameters are typically symbolized using Greek letters, such as β in the case of regression coefficients [see (3.4)]. In this book, sample-based estimates of those equations are symbolized using lowercase, Roman letters, such as b in the case of regression coefficients. Thus, the population slope is β_1 and the sample-based estimate of the slope is b_1 . Other examples include σ and s for the population and sample-based standard deviation, and ρ and r for the population and sample-based correlation.

3.3.3 Estimation

The slope is computed as

$$b_1 = r_{yx} \frac{s_y}{s_x} \quad (3.7)$$

where r_{yx} is the correlation between x and y , and s_y and s_x are the standard deviations for y and x , respectively. The intercept is computed as

$$b_0 = \bar{y} - b_1 \bar{x} \quad (3.8)$$

Of course, we rarely—if ever—compute these quantities by hand. Stata’s `regress` command does this for us. The syntax for `regress` is straightforward. Following the keyword `regress`, you type the dependent variable followed by the independent variables. As discussed in section 2.4.1, you can use `if` and `in` with `regress`.

`regress depvar [indepvars] [if] [in] [weight] [, options]`

Estimating (3.3) with `regress` is done as follows:

```
. regress attitude educ
```

Source	SS	df	MS	Number of obs	=	1,006
Model	8.81344558	1	8.81344558	F(1, 1004)	=	33.02
Residual	267.958141	1,004	.266890579	Prob > F	=	0.0000
Total	276.771587	1,005	.275394614	R-squared	=	0.0318
				Adj R-squared	=	0.0309
				Root MSE	=	.51661

attitude	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
educ	.0335253	.005834	5.75	0.000	.0220771 .0449736
_cons	1.991708	.0804314	24.76	0.000	1.833876 2.149541

Focus on the bottom portion of the output for now. Stata calls the estimated intercept `_cons` (that is, constant). Thus, $b_0 = 1.99$. The slope for `educ` is $b_1 = 0.03$. We can now fill in the details of (3.3).

$$\widehat{\text{attitude}}_i = 1.99 + 0.03\text{educ}_i \quad (3.9)$$

We can make predictions based on the coefficients. For example, if we want to know the predicted `attitude` for someone who has 15 years of education, we simply plug 15 into (3.9).

$$2.44 = 1.99 + 0.03 \times 15 \quad (3.10)$$

The predicted `attitude` for someone who has 15 years of education is 2.44. Another way to say that is, “The expected value of `attitude` given that `educ` equals 15 is 2.44, or $E(\text{attitude}|\text{educ} = 15) = 2.44$.”