

LINEAR MODELS AND ŘURITÀ KRÀLOVSTVÌ

USING THE KINGDOM FOR GREATER INSIGHT

Version: 0.704442ζ

OLE J. FORSBERG

*Department of Mathematics
Knox College
Galesburg, IL, USA*

December 2023

ALL RIGHTS RESERVED. No part of this work covered by copyright herein may be reproduced, transmitted, stored, or used in any form or by any means graphic, electronic, or mechanical, including, but not limited to photocopying, recording, scanning, digitizing, taping, Internet distribution, information networks, or information storage and retrieval systems, except as permitted under Section 107 or 108 of the 1976 US Copyright Act, without the prior written permission of the author.

The current draft version of this document is free (without cost). This document is distributed in the hope that it will be useful, but without any warranty; without even the implied warranty of merchantability or fitness for a particular use or for a particular purpose.

This document was typeset using the $\text{\LaTeX} 2_{\epsilon}$ environment. All graphics were created by the author using the \R statistical environment (and referenced packages). All data is either public domain or specifically generated for this text. The following are the photograph credits.

Page	Site	Source/Holder	CC Type
Cover	Own AI Artwork	Ole J. Forsberg	CC0
2	King Harald V of Norway	Government of Norway	CC0
3	Ruritanian Flag	Ole J. Forsberg	BY-NC-SA
??	Ambassador Georgette Mosbacher	US Embassy to Poland	CC0
11	King Olav V of Norway	Government of Norway	CC0
15	St. Canute's Cathedral, Odense, Denmark	Ole J. Forsberg	BY-NC-SA
57	Street View, Karlstad, Sweden	Ole J. Forsberg	BY-NC-SA
79	A Steadfast Tin Soldier statue, Odense, Denmark	Ole J. Forsberg	BY-NC-SA
159	A bridge in Karlstad, Sweden	Ole J. Forsberg	BY-NC-SA
195	A canal in The Hague, Netherlands	Ole J. Forsberg	BY-NC-SA
208	A map of San Marino	Wikipedia User Aotearoa	BY-SA
223	The Peace Palace in The Hague, Netherlands	Ole J. Forsberg	BY-NC-SA
241	Sunne Church in Sunne, Sweden	Ole J. Forsberg	BY-NC-SA
265	The Vigeland Park in Oslo, Norway	Ole J. Forsberg	BY-NC-SA
289	Bergen, Norway	Ole J. Forsberg	BY-NC-SA
329	Bergenshus Fortress in Bergen, Norway	Ole J. Forsberg	BY-NC-SA
357	Windmill scene in Malmö, Sweden	Ole J. Forsberg	BY-NC-SA
387	Street scene in Malmö, Sweden	Ole J. Forsberg	BY-NC-SA
421	Street scene in Karlstad, Sweden	Ole J. Forsberg	BY-NC-SA
443	Royal Guard in Oslo, Norway	Ole J. Forsberg	BY-NC-SA
465	William Sealy Gosset, 1908	Wikipedia	CC0
466	Ronald Fisher, 1913	Wikipedia	CC0
467	Augustin-Louis Cauchy, 1901	Wikipedia	CC0

Contents

List of Figures	ix
List of Tables	xiii
Preface	xv
1 Introduction to Ruruitania	1
1.1 Background of Ruruitania	2
1.2 Economics	4
1.3 US–Ruritanian Relations	5
1.4 Illustrations of Analyses	5
1.5 Conclusion	11
I Least Squares	13
2 Introduction to Linear Regression	15
2.1 Scalar Representation	18
2.2 Matrix Representation	30
2.3 Predictions and the Hat Matrix	37
2.4 The PRE Measures	42
2.5 Multicollinearity and Categorical Independent Variables	45
2.6 Conclusion	53
2.7 End-of-Chapter Materials	54
3 Improved! Now with Probabilities	57
3.1 Probability Distributions	59
3.2 Test Statistics and Hypothesis Testing	68
3.3 Confidence Intervals	70
3.4 The Working-Hotelling Bands	75

3.5	Conclusion	75
3.6	End-of-Chapter Materials	76
4	Dood! Check the Requirements	79
4.1	Normality	81
4.2	Constant Expected Value	95
4.3	Constant Variance	101
4.4	Multicollinearity	110
4.5	Conclusion	116
4.6	End-of-Chapter Materials	117
5	A Time for Some Examples	123
5.1	Full Example: Violent Crime	125
5.2	Full Example: Violent Crime, Wealth, Region	130
5.3	Full Example: Cows in the City of Děčín	138
5.4	Conclusion	153
5.5	End-of-Chapter Materials	154
6	Fixing the Violations	159
6.1	The Issue of Boundedness	161
6.2	Full Example: The South Sudanese Referendum	176
6.3	Heteroskedastic Adjustments	182
6.4	Conclusion	184
6.5	End-of-Chapter Materials	185
II	Beyond the Ordinary	193
7	Other Least Squares	195
7.1	Ordinary Least Squares	197
7.2	Weighted Least Squares	198
7.3	Generalized Least Squares	205
7.4	Full Example: May the Strong Force be with You	210
7.5	Full Example: Elections in Ruritania	212
7.6	Conclusion	217
7.7	End-of-Chapter Materials	218
8	Quantile Regression	223
8.1	Parameter Estimation	225
8.2	Quantile Regression	230
8.3	Conclusion	235

8.4	End-of-Chapter Materials	237
9	Maximizing the Likelihood	241
9.1	The Likelihood	243
9.2	The MLE and the CLM	251
9.3	Conclusion	258
9.4	End-of-Chapter Materials	259
III	Beyond the Classical Model	263
10	Generalized Linear Models	265
10.1	The CLM and the GLM	267
10.2	The Requirements for GLMs	268
10.3	Assumptions of GLMs.	272
10.4	The Gaussian Distribution	273
10.5	Generalized Linear Models in R	275
10.6	Conclusion	282
10.7	End-of-Chapter Materials	283
11	Binary Dependent Variables	289
11.1	Binary Dependent Variables	291
11.2	Latent Variable Modeling	294
11.3	The Mathematics.	296
11.4	Modeling with the Logit	303
11.5	Prediction Accuracy	306
11.6	Modeling with Other Links.	313
11.7	Model Selection	316
11.8	Conclusion	321
11.9	End-of-Chapter Materials	322
12	Binomial Dependent Variables	329
12.1	Binomial Distribution	331
12.2	The Mathematics.	333
12.3	Full Example: O Canada!	335
12.4	Full Example: Sri Lanka in 2010.	345
12.5	Conclusion	350
12.6	End-of-Chapter Materials	351
13	Count Dependent Variables	357
13.1	Linear or Poisson Regression?.	359

13.2	The Mathematics	361
13.3	Overdispersion	368
13.4	Full Example: Body counts	373
13.5	The Bias-Variance Trade-Off	379
13.6	Conclusion	381
13.7	End of Chapter Materials	382
14	Nominal and Ordinal Dependent Variables	387
14.1	Nominal Dependent Variable	390
14.2	Ordinal Dependent Variable	401
14.3	Extended Example: Cattle Feed	405
14.4	Extended Example: The State University of Ruritania	410
14.5	Conclusion	413
14.6	End-of-Chapter Materials	414
IV	The Appendices	419
Appendix M	The Appendix of Matrices	421
M.1	Matrix Basics	422
M.2	Addition	424
M.3	Multiplication	425
M.4	Other Matrix Terms.	434
M.5	Consequences	436
M.6	Statistics in Matrices	438
M.7	End-of-Appendix Materials.	441
Appendix S	The Appendix of Statistics	443
S.1	Importantly Confusing Points.	445
S.2	Sample Statistics	446
S.3	Population Parameters	452
S.4	Probability Distributions	457
S.5	Distributions of Sample Statistics	469
S.6	Other Topics	474
S.7	End-of-Appendix Materials.	496
Index		499

List of Figures

1.1	His Majesty Rudolph II, King of Ruritania.	2
1.2	The Flag of Ruritania, 1954	3
1.3	Economic Output	4
1.4	Mark Brzezinski	5
1.5	Rudolph II (1952)	11
2.1	Illustrating a line of best fit	17
2.2	Example an OLS line of best fit	23
2.3	Comparison of \mathbf{Y} and $\hat{\mathbf{Y}}$ spaces	39
3.1	Illustrating residuals	59
3.2	Comparing confidence interval and prediction interval widths	66
3.3	A visual of a confidence interval	71
3.4	Central (or symmetric) confidence interval	72
3.5	Minimum-width confidence interval	73
4.1	Example quantile-quantile plot	82
4.2	Default residuals histograms	84
4.3	Enhanced residuals histograms	84
4.4	Residuals plot of misspecified model	96
4.5	Residuals plot of properly-specified model	97
4.6	Illustration of homo- and heteroskedasticity	102
4.7	Diagram of multicollinearity	115
5.1	Plot of the violent crime rate in 2000 against that in 1990	128
5.2	Plot of the 2000 violent crime rate against the 1990 GSP per capita	136
5.3	Correlation plots between the three independent variables.	141
5.4	Prediction graphs of our <code>cows</code> model.	149
5.5	Plot of the predicted vote outcomes from the MC experiment.	152
6.1	Schematic of the variable transformation procedure.	162

6.2	Schematic of the transformation procedure for Example 6.1.	166
6.3	Histogram of the MC experiment for Děčín.	167
6.4	Histogram of the MC experiment for Ruritania.	172
6.5	A scatterplot of the results of the South Sudan referendum.	177
6.6	The results of the South Sudan referendum with predictions.	181
7.1	A map of Ruritania	208
7.2	Invalidation plot for the Parliamentary election	216
8.1	Line of “best” fit	224
8.2	Comparing OLS and medreg estimators	230
8.3	Quantile effects of crime on crime	233
8.4	Quantile effects of wealth on crime	235
10.1	Graphs comparing CLM and GLM transformed results.	278
10.2	Comparison of CLM and GLM predictions.	281
11.1	Residual scatterplot for example 11.1.	292
11.2	Schematic of logistic regression.	295
11.3	Three symmetric link functions.	301
11.4	Two asymmetric link functions.	302
11.5	Plot of predicted coin weightings.	306
11.6	The model accuracy against various thresholds.	309
11.7	The ROC curve for the coin flipping model.	312
11.8	The ROC curve using the <code>ROC</code> command.	313
11.9	Plot of logit and complementary log-log functions.	314
11.10	Plot of logit and log-log functions.	316
12.1	A map of the world and Canada	335
12.2	The ‘O Canada’ data and estimates	343
12.3	A map of the world and Sri Lanka	345
12.4	Sri Lankan Presidential returns from 2010, without postal	348
12.5	Sri Lankan Presidential returns from 2010	349
13.1	Poisson data plot with regression curves.	360
13.2	Plot of initiative use against state population.	367
13.3	Plot of the number of deaths due to terrorism.	379
14.1	Explanatory schematic for thresholding.	402
14.2	Feed type probabilities for RUR ranch.	409
14.3	Class level probabilities.	412

M.1	A schematic of commensurate matrices for multiplication	426
S.1	William Sealy Gosset, 1908	465
S.2	Ronald Fisher, 1913	466
S.3	Augustin-Louis Cauchy, 1901	467
S.4	Abraham Wald	480
S.5	Distribution of the number of runs	482

List of Tables

1	The Greek alphabet	xx
2.1	Toy data for Example 2.1.	21
2.2	Data for Example 2.2.	24
2.3	Raw data for Example 2.5.	45
2.4	Structured data for Example 2.5.	47
3.1	P-value calculations	69
5.1	Statistics on the <code>cows</code> data set.	139
5.2	Correlations between the variables in the <code>cows</code> data	140
5.3	The grammar of formulas	143
5.4	Results table for the <code>cows</code> example.	147
6.1	Results table for the <code>cows</code> vote model (logit).	165
6.2	Results table for the GDP per capita model.	170
6.3	Results table for the South Sudan referendum model.	179
6.4	Regression table using White standard errors	183
7.1	Data for the Strong Force example	211
8.1	Example of median regression	226
10.1	GLM distributions and links.	269
10.2	Results table for the <code>cows</code> example.	277
10.3	Results table for the <code>cows</code> vote model (logit).	278
10.4	Results table from fitting the GDP data with GLMs.	279
11.1	Insurance data to accompany Example 11.1.	292
11.2	Binary dependent variable link functions.	300
11.3	Results table for logit regression on the coin flip data.	304
11.4	Results table for clog-log regression on the coin flip data.	315

13.1 Initiative model, with adjusted standard errors.	369
13.2 Initiative model, fitted using QMLE.	370
13.3 Results table for initiatives using the Negative Binomial.	372
13.4 Three terrorism models, days as independent variable.	375
13.5 Three terrorism models, using days as offset.	376
13.6 Two terrorism models with higher degrees.	378
14.1 Correlation matrix.	396
14.2 Results table for the logit model.	397
14.3 Results table for the multinomial regression model.	399
14.4 Result of ordinal regression in \mathbb{R}	403

Preface

I sat there, as I was wont to do, in the Spillane Reading Room, drinking coffee in the early morning hour, trying to find my wakefulness. The air was peaceful, with two first-years discussing and sharing insights about international politics and world events with each other and with me. It was life as I expected it to be in academia.

Frequently, more so as the term wore on, 'Simon' would rush into the room and explode, "I don't know what this *thing* is telling me!" Then, he would throw down five pages of printout from a well-known statistical package and throw up his hands, as if beseeching the gods of statistics to send down an answer to him.

After allowing the situation to calm, and the first-years to start breathing again, I would ask the question "What did you do to get the printout?" (Show your work.) For some reason, I always expected the answer to differ from all the times previous; I expected him to tell me what he intended to do, what specific actions he performed, what information those commands were supposed to give him, and why he needed that information.

"I clicked on some menu things and this came out." As expected.

He and I would then sit down and go over the five pages of printout, examining what each of the tables and numbers meant in relation to his research. Eventually, after dealing with several "Why is it giving me *that* information?!" questions, Simon would be vaguely satisfied with the printout and could select several statistics from the printout that would provide the information he sought.

However, there were many more questions I wanted to ask him. Most centered on questions about the validity of the tests performed. I knew, however, that such a line of questioning would be moot with the statistical package he (and his class) used. Either the company that owned the package had the test available, or it did not. There was no (easy) way to add tests and procedures. Thus, Bartlett's test of equal variances was not an option, even when the data was such that it should be analyzed using it. Furthermore, the number of available tests was quite limited.

In addition to the extensibility issue, there was the issue of clicking your way to an analysis. If Simon needed to repeat the exact analysis, except

for a tweak or two, he would have to start from scratch and repeat it all, clicking all the same menu items, hoping that he did not make a mistake along the way. This repeat analysis happens quite frequently in life.



One of the many pleasures I have as a statistician is my exposure to many different disciplines. I consulted for a woman doing research in science education. Her specific problem was to determine if a specific Science Teaching Unit (STU) had the students like science more. She performed her experiment on a class of fifth and sixth graders in rural North Dakota (hardly representative). She gave them a 50-question pre-test, taught the Unit, then gave the same students a 50-question post-test (what about reliability?!?!).

She then contacted me and sent the data.

On the data she sent to me, I spent about three hours analyzing the data, coming to some interesting (and counter-intuitive) conclusions. In my experience, researchers have a sufficient feel for their discipline that surprising results are frequently a result of analysis error. Thus, I checked my analysis for errors.

I was actually able to check the analysis because I wrote a script — a series of commands — detailing every bit of the analysis I performed. Mouse-clicking my way through the analysis would make it all but impossible to check my work (something my math teachers in grade school always emphasized). Thus, I was confident that the analysis I returned to her were correct ...

... conditional on the data being correct.

The next day, she emailed me and let me know that she found several serious errors in her data. Relying on mouse-clicks would have made those original three hours a waste of time. However, once she sent the corrected version of the data, the analysis took 90 seconds. Clicking my way to re-analysis would have taken the same amount of time as the first analysis — time we did not have (we were facing a deadline). Re-running the script only took processing time.

While I *am* a fan of mouse-clicks under many circumstances, I am not a fan when it comes to statistical analysis. Scripting provides three definite advantages over mouse clicks: You get what you request, you can check your work easily, and you can re-run the analysis with little effort.



One benefit of being in academia is that I get to travel the world sharing my research and its results. One June found me standing in front of a formidable group in Denmark's second-city (Odense) giving such a presentation. After I had discussed the current literature on the causes of terrorism, my statistical model, the results, and my conclusions, I opened it up for questions. After about five minutes of questions on the validity of my statistical model, one professor asked how my analysis would change were I to add an offset to the model.

After a brief panic, I decided to actually run the regression with the offset in front of the audience. I apologized for not knowing the answer to his question, but I would be happy to hypothesize the effect of the variable offset while running the analysis. The professor smiled as I tried to open my analysis script to modify and run it. It turned out that R was not installed on the presentation computer. No worries. I opened the R folder on my USB drive, double-clicked on the R program, and proceeded with the altered analysis — all the while discussing the theoretical effects of using such an offset under these circumstances. Before I could finish hypothesizing, R gave me the answer (which, thankfully, agreed with my hypotheses).

Now that I had my model laid bare before all, many more questions arose about different alterations I could (or should) make to the model. All of which I was able to perform in front of the now hyper-interested crowd.



These three vignettes illustrate many strengths of a statistical environment like R. First, it encourages one to write out the analysis and “show the work.” This makes it easier to see the entire scope and sequence of the analysis. It also makes it easier to check for errors. Second, it is extensible. If there is a cutting-edge test or procedure you wish to run, there is probably a package that contains it. If not, you are quite free to write it yourself. Finally, one can carry R around on a USB drive, allowing anyone to perform analyses whenever there is a computer, like in Denmark.

Oh yeah, that R is free is also a nice feature, especially as statistical packages can run from \$600 to \$6000 and up, *and* can have licenses requiring annual payments. As budgets get tighter, an ability to work successfully with a free (and powerful) statistical environment is invaluable.

Prerequisites

For any book (or course), there is a necessary assumption made about the background of the reader (or student). For the material in this book, I assume that you have had experiences with elementary statistics, matrices, and calculus.

In **differential calculus**, you will need to understand how to optimize (minimize or maximize) functions. In **integral calculus**, you should understand how to calculate areas under a curve (probabilities related to density functions). Beyond that, there is little calculus needed.

The **matrix topics** you need consist of being able to perform algebra on matrices. Beyond that, anything you remember from a typical linear algebra course will help things make a bit more sense. Matrix rank, invertibility, idempotency, projection matrices, orthogonality, etc., are all important in ordinary least squares. So, if you remember those topics, you will be ahead of the curve. If you do not remember them, then you will need to (re-)learn them in this course. Appendix M will help with remembering and learning the important matrix topics.

Finally, I wrote this book to be a second course in statistics, one that started where a typical **introductory course** ended. Because of this, I also assume you remember many topics from such a course. These topics definitely include the meanings of confidence intervals and p-values. They also include probability distributions, t-tests, issues with multiple testing, and the Central Limit Theorem (CLT). To help refresh your memory, work through Appendix S. Note that Appendix S also introduces you to some (optional) advanced items. These topics were included at the request of past students who wanted to actually see a proof of the Central Limit Theorem. Rest assured that understanding the CLT is more important than being able to prove it. Furthermore, the proof offers little in the way of a deeper understanding.

A Note on Notation

Sadly, notation varies across the discipline. This is a result of the history of statistics: Many of the methods came from disciplines that *used* statistics, rather than from statistics itself. Different disciplines use different notation for the same idea. Thus, any discussion of Survival Analysis needs to include Event History Analysis and Reliability Analysis, as they all study the same phenomena but from different disciplines (medicine, social sciences, and engineering, respectively).¹

Even within a discipline, there is often a variety of notation used to indicate the same ideas. For instance, probability functions are often parametrized in different ways. The parameter of the Exponential distribution can be the rate λ or the expected value θ ; the second parameter of the Normal distribution (Gaussian distribution, Gauss-Laplace distribution) may be the variance σ^2 , the standard deviation σ , or the precision as $\tau = 1/\sigma$ or as $\tau = 1/\sigma^2$. The symbol for the average rate in the Poisson distribution can be μ or λ . In this volume, I will (try to) keep consistent with notation, and I will explain the notation before I use it.

To that end, population parameters will be signified using Greek minuscules. Sets from which the population parameters can belong (parameter spaces) will be Greek majuscules. Both are included in Table 1. All random variables are Roman majuscules. All realized random variables (data) are Roman minuscules. Violations of these rules will exist, but should be kept to a minimum.

Thus, if we theorize that our measurements come from a population that is Normally distributed, with mean μ and standard deviation σ , we would specify that $\mu \in M$ and that $\sigma \in \Sigma$, where $M = \mathbb{R}$ and $\Sigma = (0, \infty)$.²

Now, if we know that the mean is 15 and the standard deviation is 10, I would write this as

$$X \sim \mathcal{N}(\mu = 15; \sigma = 10), \quad (0.1)$$

where the mean of the population is denoted by μ , the standard deviation by σ , and the Normal distribution by \mathcal{N} .

¹Furthermore, the term “reliability analysis” means different things in different areas. It could mean studies of how long until a part or a machine breaks. It could also mean how robust conclusions are to changes in model assumptions.

²Be aware of the difference between Σ , \sum , and \mathfrak{Y} . The first is the set of possible values of σ . The second is the symbol indicating the sum of what follows. The third is a generic covariance matrix. Frequently, the difference will be obvious: first, from context; second, because the sum will be performed over an index and symbolized as \sum_i or as $\sum_{i=1}^n$.

Minuscule Letter	Majuscule Letter	Name
α	A	alpha
β	B	beta
γ	Γ	gamma
δ	Δ	delta
ϵ, ε	E	epsilon
ζ	Z	zeta
η	H	eta
θ	Θ	theta
ι	I	iota
κ	K	kappa
λ	L	lambda
μ	M	mu
ν	N	nu
ξ	Ξ	xi
\omicron	O	omicron
π	Π	pi
ρ	R	rho
σ	Σ	sigma
τ	T	tau
υ	Y	upsilon
ϕ	Φ	phi
χ	X	chi
ψ	Ψ	psi
ω	Ω	omega

Table 1: The usual Greek alphabet in the canonical order. Being familiar with the letters will make it easier to recognize the implied meaning behind the letter.

Once we take those measurements, we would call the variable x . The difference between random variables and *realizations* of those random variables is that the random variable has a probability distribution associated with it; the realized data are just numbers.

By the way, if we wish to specify “parameter” in general, we use θ . As a result, its parameter space is Θ .

Matrices (and vectors) will be indicated with **bold-faced** letters. Thus, \mathbf{x} is the data matrix (observed values) and \mathbf{X} is the data matrix (theoretical values).

Between these two, it *does* make sense to say something like

$$\mathbf{X} \sim \mathcal{N}(\mu = 15; \sigma = 10)$$

It does NOT make sense to say something similar about \mathbf{x} . The observed data do *not* have a theoretical distribution.

Conclusion

And so, with all of this said, turn the page and begin your trek through linear models. The first chapter introduces you to the topics of both linear models and the Kingdom of Ruritania. The former is the purpose of the book. The latter is a common theme and source of examples. Since the Kingdom of Ruritania does not exist, think of it as a generic country with no real information about it beyond what is given.

Had I used a real country, it would be perfectly defensible for the student to bring in real information about that country. This may cloud the intended statistical lesson.

Also, using Ruritania allows me the ability to be creative in my storytelling.

I hope you enjoy the journey.

Vyčkej Času
Adolph Heyduk

Nespěchej k štětí, tiše jen,
vše nemůž' najednou býti;
nejprv se travou zjeví len,
pozděj jak modravé kvítí.

Poupě, jež touží růží být,
obalu hledí se zbavit;
dříve než můžeš v nebi žít,
musíš se v očistci stavít.

Translation:
Bide Your Time
Adolph Heyduk

Don't rush to fortune, ease your tracks,
all can't at once be present;
You would for grass mistake the flax
ere blue-flower'd iridescent.

A bud, so keen to be a rose,
is for its calyx sorry;
But ere in heaven you'll repose,
you'll bide in purgatory.

Source: <http://www.vz.jp.cz/basne.htm#Heyduk>

