# Statistical Modeling: A Fresh Approach

*Daniel T Kaplan*

*E-book version of Second Edition ©2017*

# Contents

# Preface to this electronic version

Placeholder

# Chapter 1

# Introduction

Placeholder

## Example: Applying to Law School

## Example: Nitrogen Fixing

## Example: Sex Discrimination

## 1.1  Models and their Purposes

## 1.2  Observation and Knowledge

## 1.3  The Main Points of this Book

## Reading Questions

# Chapter 2

# Data: Cases, Variables, Samples

Placeholder

## 2.1 Kinds of Variables

## 2.2 Data Frames and the Unit of Analysis

## 2.3 Populations and Samples

## 2.4 Longitudinal and Cross-Sectional Samples

# Chapter 3

# Describing Variation

Placeholder

## 3.1 Coverage Intervals

## Aside: What's Normal?

## 3.2 The Variance and Standard Deviation

## 3.3 Displaying Variation

## 3.4 Shapes of Distributions

### 3.4.1 Categorical Variables

### 3.4.2 Quantifying Categorical Variation

# Chapter 4

# Groupwise Models

Placeholder

## 4.1   Grand and Group-wise Models

## 4.2   Accounting for Variation

## Aside:  The geometry of partitioning

## 4.3   Group-wise Proportions

## 4.4   What's the Precision?

## 4.5   Misleading Group-wise Models

# Chapter 5

# Confidence Intervals

Placeholder

## 5.1  The Sampling Distribution

## Aside: Precision and Sample Size

## 5.2  The Resampling Distribution & Bootstrapping

## 5.3  Re-sampling

## 5.4  The Re-Sampling Distribution

## Example: The Precision of Grades

## 5.5  The Confidence Level

## 5.6  Interpreting Confidence Intervals

## 5.7  Confidence Intervals from Census Data

# Chapter 6

# Language of Models

Placeholder

**6.1   Models as Functions**

**6.2   Model Functions with Multiple Explanatory
         Variables**

**6.3   Reading a Model**

**6.4   Choices in Model Design**

The Data

The Response Variable

Explanatory Variables

**6.5   Model Terms**

The Intercept Term (and no other terms)

Intercept and Main Terms

Interaction Terms

Transformation Terms

Aside: Are swimmers slowing down?

Main Effects without the Intercept

**6.6   Notation for Describing Model Design**

# Chapter 7

# Model Formulas and Coefficients

Placeholder

## 7.1   The Linear Model Formula

## 7.2   Linear Models with Multiple Terms

## Aside: Interpreting Interaction Terms

## 7.3   Formulas with Categorical Variables

## 7.4   Coefficients and Relationships

## 7.5   Model Values and Residuals

## 7.6   Coefficients of Basic Model Designs

## 7.7   Coefficients have Units

## Aside: Comparing Coefficients

## 7.8   Untangling Explanatory Variables

## Aside: Interaction terms and partial derivatives

## 7.9   Why Linear Models?

# Chapter 8

# Fitting Models to Data

Placeholder

## 8.1 The Least Squares Criterion

## Aside: Why Square Residuals?

## 8.2 Partitioning Variation

## Aside: Partitioning and the Sum of Squares

## 8.3 The Geometry of Least Squares Fitting

## 8.4 Redundancy

## Example: Almost redundant

## Aside: Redundancy – (Almost) Anything Goes

# Chapter 9

# Correlation and Partitioning of Variance

Placeholder

## 9.1   Properties of R²

Example: Quantifying the capture of variation

## 9.2   Simple Correlation

Example: Relationships without Correlation

Example: R versus R²

## 9.3   The Geometry of Correlation

## 9.4   Nested Models

Example: R² Out of the Headlines

## 9.5   The Geometry of R²

# Chapter 10

# Total and Partial Change

Placeholder

**10.1   Total and Partial Relationships**

**10.2   Example: Covariates and Death**

Example: Used Car Prices

**10.3   Models and Partial Relationships**

Aside: Partial change and partial derivatives

**10.4   Adjustment**

**10.5   Simpson's Paradox**

Example: Cancer Rates Increasing?

**10.6   Explicitly Holding Covariates Constant**

Example: SAT Scores and School Spending

Aside: Divide and Be Conquered!

**10.7   Adjustment and Truth**

**10.8   The Geometry of Covariates and Adjustment**

Aside: Interaction terms and partial derivatives

# Chapter 11

# Modeling Randomness

Placeholder

## 11.1   Describing Pure Randomness

## Example: Rolling a Die

## Example: The Chance of Rain

## Example: Flipping two coins.

## 11.2   Settings for Probability Models

## 11.3   Models of Counts

**The Binomial Model**

## Example: Multiple Coin Flips

## Example: Houses for Sale

## The Poisson Model

## Example: The Rate of Highway Accidents

## 11.4   Common Probability Calculations

## Example: A Political Poll

## Example: A Normal Year?

## 11.5   Models of Continuous Outcomes

## The Uniform Model

## Example: Equally Likely

## The Normal Model

## Example: IQ Test Scores

## The Log-normal Model

# Chapter 12

# Confidence in Models

Placeholder

**12.1  The Sampling Distribution & Model Coefficients**

**12.2  Standard Errors and the Regression Report**

**12.3  Confidence Intervals**

Example: Wage discrimination in trucking?

Example: SAT Scores and Spending, revisited

**12.4  Confidence in Predictions**

Example: Catastrophe in Grand Forks

**12.5  A Formula for the Standard Error**

**12.6  Confidence and Collinearity**

Aside: Redundancy and Multi-collinearity

**12.7  Confidence and Bias**

# Chapter 13

# The Logic of Hypothesis Testing

Placeholder

## 13.1 Example: Ups and downs in the stock market

## 13.2 An Example of a Hypothesis Test

## 13.3 Inductive and Deductive Reasoning

## Deductive Reasoning

## Inductive Reasoning

## 13.4 The Null Hypothesis

## 13.5 The p-value

## 13.6 Rejecting by Mistake

## 13.7 Failing to Reject

## Aside: Calculating a Power

## 13.8 A Glossary of Hypothesis Testing

## 13.9 Update on Stock Prices

# Chapter 14

# Hypothesis Testing on Whole Models

Placeholder

## 14.1   The Permutation Test

## 14.2   R² and the F Statistic

## Example: Marriage and Astrology

## 14.3   The ANOVA Report

## Aside: F and R²

## Example: Is height genetically determined?

## Aside: The shape of F

## Example: F and Astrology

## 14.4   Interpreting the p-value

## Multiple Comparisons

## Example: Multiple Jeopardy

## Significance vs Substance

## Example: The Significance of Finger Lengths

# Chapter 15

# Hypothesis Testing on Parts of Models

Placeholder

## 15.1   The Term-by-Term ANOVA Table

## 15.2   Covariates Soak Up Variance

Example: Wages and Race

## 15.3   Measuring the Sum of Squares

## 15.4   ANOVA, Collinearity, and Multi-Collinearity

Example: Height and Siblings

## 15.5   Choosing the Order of Terms in ANOVA

Example: Wages and Race: Part 2

Example: Wages and Race: Part 3

Example: Testing Universities

## 15.6   Non-Parametric Statistics

# Chapter 16

# Models of Yes/No Variables

Placeholder

## 16.1 The 0-1 Encoding

## 16.2 Inference on Logistic Models

## 16.3 Model Probabilities

## Example: Log-odds ratios of Prostate Cancer

# Chapter 17

# Causation

Placeholder

# 17.1   Interpreting Models Causally

# Example: Greenhouse Gases and Global Warming

# 17.2   Causation and Correlation

# 17.3   Hypothetical Causal Networks

# 17.4   Networks and Covariates

# 17.5   Pathways

# 17.6   Pathways and the Choice of Covariates

# Example: Learning about Learning

# 17.7   Sampling Variables

# 17.8   Disagreements about Networks

# Example: Sex Discrimination in Salary

# Update on Global Warming

# Chapter 18

# Experiment

Placeholder

## 18.1 Experiments

## 18.2 Experimental Variables and Experimental Units

## Example: Virtues of Doing Surgery while Blind

## Example: Oops! An accidental correlation!

## 18.3 Choosing levels for the experimental variables

## 18.4 Replication

## 18.5 Experiments vs Observations

## 18.6 Creating Orthogonality

### 18.6.1 Random Assignment

## Blocking in Experimental Assignment

## 18.7 When Experiments are Impossible

## 18.8 Intent to Treat

## 18.9 Destroying Associations

## Instrumental Variables

## Matched Sampling

## Example: Returning to Campaign Spending ...

## 18.10 Conclusion