

A Minimal Book Example

Yihui Xie

2020-01-22

Contents

1	Prerequisites	5
2	Sampling, Experiments, and Exploratory Data Analysis	7
2.1	Data in the Wild	7
2.2	Statistical Testing Ideas	11
3	Point Estimates	13
3.1	Estimate with means	13
3.2	Estimate with quantiles	13
4	Accounting for Uncertainty	17
4.1	Example one	17
4.2	Example two	17
5	Inference via Hypothesis Tests for One Sample	19
6	Inference via Confidence Intervals for One Sample	21
7	Inference for Two Categorical Variables	23
8	One-Way ANOVA	25
9	Motivating example	27
10	Simple model for the data	29
11	exploratory analysis	31

12 sources of variation	33
13 statistical model and analysis	35
14 compare analyses	37
15 Multi-way ANOVA	39
16 Block Designs	41
17 Regression Models	43
18 The General Linear Model	45
19 Mixed Models	47
20 Split Plot and Repeated Measures Designs	49
21 Logistic Regression and Generalized Linear Models	51
21.1 Stuff here	51
22 Generalized Linear Mixed Models	53
23 Appendix - Learning Objectives	55
23.1 Book-level	55
23.2 Topic-level	59
23.3 From ST512	59
23.4 For Point Estimates Chapter	63
24 Appendix - Notation	65
24.1 Standard notation	65
24.2 Mixed models	66
24.3 Effects model representation	66
24.4 Estimators vs. Estimates	66

Chapter 1

Prerequisites

Stuff

Chapter 2

Sampling, Experiments, and Exploratory Data Analysis

2.1 Data in the Wild

Data is a collection of information about a group, which may include both quantitative and qualitative variables. Data is ubiquitous in today's society. Healthcare, marketing, history, biology, ... basically every field has a quantitative aspect. The quality of data varies greatly from study to study.

2.1.1 Data from Experiments

Some data comes from a well-designed experiment where a researcher uses sound principles to select units and conduct interventions.

For example, a mechanical engineer wants to determine which variables influence overall gas mileage of a certain year and model of a car. Gas mileage would be referred to as the **response** variable for this study.

After careful consideration, the engineer chooses to investigate the following **factors** that may affect the overall gas mileage:

- Tire pressure (low, standard)
- Octane rating of fuel (regular, midgrade, premium)
- Type of driving (defensive, aggressive)

They also choose to **control** or hold constant the following variables during the implementation of the study:

- Weather conditions
- Route
- Tire type
- Past car usage

The engineer randomly selects 24 cars from the assembly line for that year and model of car (we'll learn more about the importance of selecting a representative sample of cars shortly). Software is used to randomly assign a **treatment** or combination of the factors to each car of the 24 cars. For instance, low tire pressure, regular octane fuel, and defensive driving would be a treatment. The cars would be called the **experimental units** or (EUs) as they are the unit the treatments are assigned to.

The experiment is run and the gas mileage found for each car. As the car is being measured we'd refer to the car as the **observational unit**.

This short description exhibits three important concepts in experimental design that we'll come back to many times.

Definition 2.1. Experiment - researchers manipulate the conditions in which the study is done.

Pillars of experimental design:

- Randomization - treatments are randomly assigned to the experimental units
- Replication - multiple (independent) experimental units are assigned the same treatment
- Control - study conditions are held constant where possible to reduce variability in the response

2.1.2 Data from Observational Studies

Other data comes from observations studies where ... For example, ...

<>

The implications for the conclusions that can be made from a set of data varies greatly with the quality of the data and study design.

<>

Both studies have some things in common - popu sample...

Statistics is the science of learning from data. (Other definition from some book but not sure which - the science of designing studies or experiments, collecting data and modeling/analyzing data for the purpose of decisions making and scientific discovery when the available information is both limited and variable.)

«call out about statistics, general usage vs this type of definition»

Statistical methods are needed because data is variable. For example... The full statistical process begins with ... - Define the objective of the experiment

- Select appropriate response variables
- Identify sources of variation
- Choose experimental design
- Perform the test
- Statistically analyze the data
- Draw conclusions

We'll focus on the entire process of a study and mostly investigate designed experiments. We attempt to tackle each topic in this text with a problem-based approach. That is, we identify a real-world problem and discuss the relevant statistical ideas in context. Summaries at the end of each chapter recap the main statistical ideas.

2.1.3 Experiment Background

Marketing example. Goal to describe the customers, how they tend to purchase/shop, and maybe find some shared qualities in order to advertise curated packages to folks.

Define basic things like population, parameters, statistics, and sample.

Discuss conceptual vs actual populations and when we might care about one or the other. Our “sample” is really a bit of data from the conceptual population. Or we could consider it as the population and we just want to describe it.

2.1.4 Selecting Response Variables

Marketing example with data such as Clicks, Impressions, Total Revenue, Total Spent, Average Order Value, Sport, Time of visit/purchase, Campaigns running, etc.

2.1.5 Identifying Sources of Variation

Consider variables linked to the user. Age, other accounts, etc.

2.1.6 Choose an Experimental Design

Discuss our “sampling” scheme vs a random sample. This seems like a case where we aren’t doing a “good” scheme but not much else could be done...

Maybe talk about how in the future you could do alternate email ads or something and do an AB type study.

2.1.7 Perform the Test

Get the data from google analytics or whatever, have a plan for updating each month?

2.1.8 Look at the Data

Careful discussion of not selecting a modeling technique based on this unless it is a pilot study or an exploratory study else we have increased our nominal type I error rate...

(sometimes EDA sometimes data validation only/cleaning - more formal experiments)

Spend a lot of time here talking about graphs of different types. Sample means, sample variances, etc.

Discuss population curves vs sample histograms and the relationship.

2.1.9 Statistically Analyze the Data

New variables as functions of old?

Not a formal test here but comparisons of interest etc.

2.1.10 Draw conclusions

What actionable things have we found? Likely some trends to investigate further. Perhaps run an experiment to formally see if some alteration can be effective.

What can we conclude realistically from this data? To what population are we talking?

2.2 Statistical Testing Ideas

2.2.1 Experiment Background

This example would lend itself to a reasonably easy randomization test or simulation based test. Maybe an AB type study where we swap labels and do that with a nice visual.

Maybe third example with simulation test.

2.2.2 Selecting Response Variables

2.2.3 Identifying Sources of Variation

2.2.4 Choose an Experimental Design

Good discussion of what makes a good sampling design. Maybe a stratified example like the river and selecting houses example as a quick expose of the issues with not doing a truly random sampling technique.

Basics of experimental design (randomization, replication, error control ideas).

Recap benefits of doing an experiment vs an observational study.

2.2.5 Perform the Test

2.2.6 Explore the Data

NHST paradigm with false discovery?

2.2.7 Statistically Analyze the Data

2.2.8 Draw conclusions

Chapter 3

Point Estimates

Learning objectives for this lesson: - How to estimate a mean - Definition of “convenience sample” - Definition of “systematic sample” - Benefits/drawbacks to both approaches - Understand how to estimate a mean - Understand how to estimate a quantile - Understand implicit assumptions for these approaches

3.1 Estimate with means

3.1.1 Experiment background

Someone wants to know how much of something they need to satisfy some population. To get a good estimate of this, we can use the average amount for each one and then multiply by the whole population.

3.2 Estimate with quantiles

3.2.1 Experiment background

Big Deborah's is making new packaging for their cookies. The engineer responsible for the new design needs to make sure that the packaging fits the new cookies. While the cookie manufacturing process is standardized, there's inevitably some degree of variation in cookie size. After discussing the issue with corporate, the engineer decides that the new cookie sleeves should be large enough to fit 95% of cookies that are baked. (The largest five percent will be marketed separately as “JUMBO” cookies.)

3.2.2 Define the object of the experiment

The Engineer is tasked with determining how large the cookie sleeve needs to be. There's no way for her to know the size of every cookie that Big Deborah's has made (or will make going forward!), so she'll need to collect data on existing cookies to inform her cookie sleeve size determination.

3.2.3 Select appropriate response variables

If the maximum distance from any one point on the (round) cookie's perimeter to any other point is smaller than the diameter of the cookie sleeve, then the cookie will fit. This makes "cookie diameter" a good measure for this test. It is easy to measure for each cookie and is directly relevant to the experiment's objective.

[probably have something in here about]

3.2.4 Identify sources of variation

While the manufacturing process is standardized, there is variation in size from one cookie to the next. This is one source of variation. The engineer isn't sure of any others. However, she knows that cookies are made in multiple factories, and that each factory has multiple ovens. Ovens and factories could also be sources of variation.

3.2.5 Choose an experimental design

The Engineer knows that she needs to look at multiple cookies, since she knows that there is variation in diameter from one cookie to the next. One option would be to just use the remaining cookies in the box she has in her office (22 of the 25-count box remain). [something about convenience sample] However, she knows that cookies from the same oven are typically packaged together. If there is variation from one oven to the next, looking at the cookies she has in her office may not tell the whole story.

Instead, she chooses to take every 20th cookie manufactured off the assembly line until she gets 500 cookies. [something about systematic sample]

3.2.6 Perform the test

The day of the test comes, and the Engineer starts collecting cookies. However, problems arise! The plan has to shut down half-way through, so she only gets 431 cookies instead of the 500 she thought she would. However, she measures the diameters of each cookie and records the data in a spreadsheet.

3.2.7 Statistically analyze the data

The initial plan had been to rank-order the 500 cookies and estimate the 95th percentile using the diameter of the 475th largest cookie. Since we didn't get all of our data, we have to improvise. 431 doesn't neatly yield a value such that exactly 95% are less than or equal and 5% are greater than or equal. One option is to choose the 410th largest cookie to estimate our percentile. Slightly more than 95% of cookies will have smaller diameters than this. Alternatively, we could interpolate between the 409th and 410th cookies. [reasons and logic and math for each of these]

3.2.8 Draw conclusions

Based on this study, the Engineer concludes that a cookie sleeve large enough for a cookie of diameter XX will be big enough to contain 95% of Big Deborah cookies.

3.2.9 Discussion

- pros and cons to the approach chosen
- generalizing to other types of point estimates

Chapter 4

Accounting for Uncertainty

Some *significant* applications are demonstrated in this chapter.

4.1 Example one

4.2 Example two

Chapter 5

Inference via Hypothesis Tests for One Sample

We have finished a nice book.

Chapter 6

Inference via Confidence Intervals for One Sample

We have finished a nice book.

Chapter 7

Inference for Two Categorical Variables

We have finished a nice book.

Chapter 8

One-Way ANOVA

Learning objectives for this lesson: - Write one-way ANOVA model - Define terms - state assumptions - interpret results - Interpret ANOVA table - Describe SSE, SST, MSE - F-statistic - degrees of freedom - understand how all of these interrelate - Understand how to compare multiple group means how ANOVA is similar/different to t-tests - Understand partitioning of variation and coefficient of determination

Chapter 9

Motivating example

The United States Air Force Academy has 24 sections of Calculus I, taught by three different types of instructors: In-uniform instructors, full-time civilian instructors, and visiting faculty. The Dean of Students wants to give students the best experience possible and make sure that all three types of instructors are doing a good job. There are plausible reasons why any one of the three could be doing well: In-uniform instructors are all members of the Air Force, and students may be extra attentive in these classes because they know that these instructors rank above them in their chain of command. On the other hand, full-time instructors have been around the Academy for many years and understand the Cadets and their workloads. Alternatively, visiting faculty tend to come from prestigious institutions and may be familiar with more recently-developed pedagogical techniques. Regardless, the Dean wants to understand if there is any variation in end-of-semester grades of classes taught by these three types of instructors. At the end of the semester, she collects the average grades from each of the 24 sections. How can she go about investigating this question?

Recall from Chapter 6 that we can use t-tests to compare two group means. In this case, we'd like to do a comparison across three groups, and instead of looking at a direct comparison of one group to another, what the Dean is interested in is whether there's an *overall* difference across the three groups.

One option might be to just do a bunch of different t-tests. We could first compare classes taught by in-uniform instructors to classes taught by full-time civilians, then compare the classes taught by the in-uniform instructors to the classes taught by the visiting instructors, and then finally compare the classes taught by the full-time civilians with the classes taught by the visiting faculty. We'd end up with three p-values, each addressing different questions than the one we initially set out to answer.

We could do the same thing, except comparing courses taught by one type of instructor to the combined group of courses taught by the other two, and this

gets a bit closer to the mark. But we're still doing three tests that individually fail to answer the Dean's question.

What we'd like instead is a single hypothesis that we could test that directly gets at the Dean's concern about whether the three types of instructors were producing end-of-semester grades that were, on average, the same. [Need to make that motivation clearer above.]

Chapter 10

Simple model for the data

Narrative explanation that instructor type might matter, there should be some variation from class to class. - write some things in greek, including model without any difference by instructor type - write model with differences by instructor type - note that we can use Gaussian errors b/c Academy grades do actually tend to be centered around a C, particularly for classes like Calc - discuss model assumptions in general sense

Chapter 11

exploratory analysis

- course-to-course variability is expected
- maybe show a plot of it or something
- visualize groups using box-and-whisker plots

Chapter 12

sources of variation

Things like student population, time of day, etc. But we'll throw this all into an error term and focus on the main one, instructor type

Chapter 13

statistical model and analysis

- ANVOA model explicit w/ assumptions
- variation around overall mean w/ no groups
- variation around group means
- introduce idea of reference level

Chapter 14

compare analyses

- t-test methods from above
- ANOVA method
- compare and contrast results, interpretations, etc.

Chapter 15

Multi-way ANOVA

We have finished a nice book.

Chapter 16

Block Designs

We have finished a nice book.

Chapter 17

Regression Models

We have finished a nice book.

Chapter 18

The General Linear Model

We have finished a nice book.

Chapter 19

Mixed Models

We have finished a nice book.

Chapter 20

Split Plot and Repeated Measures Designs

We have finished a nice book.

Chapter 21

Logistic Regression and Generalized Linear Models

21.1 Stuff here

We have finished a nice book.

Chapter 22

Generalized Linear Mixed Models

We have finished a nice book.

Chapter 23

Appendix - Learning Objectives

23.1 Book-level

After reading this book you will be able to:

- identify relevant sources of variability for a potential study and, if applicable, utilize principles of design to plan a reasonable experiment to help answer questions of interest
 - covariates
 - noise variables
 - random effects
 - variance of individual observations vice variance of summary statistics
 - randomization
 - systematic variation of factors/covariates
 - factor identifiability
 - understand issues surrounding multiple comparisons
 - * Bonferroni correction
 - * at least one other method (Tukey?)
 - tradeoffs from replication within groups vice getting more groups
 - compare and contrast methods for designing an experiment when the goal of a study is prediction versus when the goal is statistical inference
- explain the general concept of point estimation and how to account for sampling variability

- definition
 - identify the right point estimate for your response variable of interest
 - estimating uncertainty for point estimates
 - * normal approximation
 - * bootstrap CI
 - * others?
 - Types of point estimates:
 - * means
 - Simple effects
 - interaction effects
 - main effects
 - * standard deviations/variance components
 - * correlation coefficients
 - * quantiles/percentiles from distributions
 - * probabilities
 - * parameters of a distribution
 - * model parameters
- describe relevant properties of random variables and probabilities
 - Distinguish between mutually exclusive and independent events.
 - Calculate probability for a given scenario, either numerically or using a Venn diagram.
 - Apply the General Addition Rule to solve probability problems.
 - Apply the Rules for Probability Distributions to create a probability distribution for a given scenario.
 - Use the complement of an event to solve probability problems.
 - Apply the Multiplication Rule for Independent Processes to solve probability problems.
 - random variables
 - * have a defined set of possible outcomes (“sample space”)
 - * Discrete vs. continuous RVs
 - * others???
 - probabilities/PDFs
 - * between 0 and 1 inclusive
 - * sum of probability of all possible events is 1
 - * $P(A) + P(A^c) = 1$, where A is an event and A^c is the complement of A
- explain the importance of statistical distributions when conducting statistical inference
 - normal distribution and approximations plus properties
 - * robustness
 - * generality

* CLT

- costs and benefits of using nonparametric approaches
- describe the fundamental inferential techniques of hypothesis testing and confidence intervals as well as compare and contrast their uses and interpretations
- identify a null and alternative for a given problem - interpret hypotheses
 - characterize the test statistic under the null - explain what a rejection region and be able to identify one - define statistical power - calculate statistical power for one- and two-sample tests of continuous and binary random variables - define statistical confidence
 - identify when using a CI and NSHT will result in the same conclusion - explain when you can use a confidence interval to test for differences (e.g., comparing a single point estimate to a threshold) and when you can't (e.g., when you have CIs for two different means)
- choose appropriate numerical summaries and graphical displays for a set of data and create these using software
 - when to use tables vs. a picture
 - types of graphical displays
 - * bar charts
 - * pie charts
 - * plotting data vice just predictions/conclusions
 - * when to include uncertainty bounds
 - * five-number summaries
 - * means vs. medians
 - * general plotting recommendations
 - * use of colors in you plots (discrete vs. divergent vs. continuous color scales, gray-scale, color-blind-friendly scales)
 - use of annotations
 - general graphical design philosophy (building a chart to illustrate a conclusion)
 - trade-offs between detail and interpretability
 - not screwing up your axes
- fit statistical models in software and interpret their output
 - Which PROCs from SAS? REG, GLM, MIXED, GLIMMIX, others??
 - `lm()`, `glm()`, `anova()` `broom`? `modelr`? `ciTools`?
 - p-values, point estimates, standard errors, f-statistics, chi-square-statistics, degrees of freedom, SS/MS, residual plots
- connect common statistical methods under the linear model framework
 - Write statistical models using matrix representaiton

- identify models written in matrix representation with their representation in software
 - identify when models written in different notation are the same or different
 - describe when specific models will give you the same results
 - * ANOVA w/ 2 factors and a t-test or a SLR
 - * ANCOVA and MLR
 - * random effects vs. fixed effects
 - * split plots vs. more general mixed models
 - * logistic regression w/ categorical factors vice contingency table analysis
 - discuss differences in assumptions associated with ANOVA vice SLR/MLR
-
- articulate the scope of inferential conclusions in light of the method of data collection, the experimental design used, the assumptions made, and the statistical analysis applied
-
- limitations due to sampling/sample frame
 - missing data
 - modeling assumptions
 - sampling assumptions
 - requirements for causal inference

23.2 Topic-level

23.2.1 Chapter 2 - Sampling, Design, and Exploratory Data Analysis

23.2.2 Chapter 3 - Point Estimation

23.2.3 Chapter 4 - Accounting for Uncertainty in Estimation

23.2.4 Chapter 5 - Inference via Hypothesis Testing for a Proportion or Mean

23.2.5 Chapter 6 - Inference via Confidence Intervals for a Proportion or Mean

23.2.6 Chapter 7 - Inference on Two Categorical Variables

23.2.7 Chapter 8 - Inference for Multiple Means

23.2.8 Chapter 9 - Multiway ANOVA

23.2.9 Chapter 10 - Block Designs

23.2.10 Chapter 11 - Regression

23.2.11 Chapter 12 - The General Linear Model

23.2.12 Chapter 13 - Mixed Models

23.2.13 Chapter 14 - Repeated Measures and Split Plot Designs

23.2.14 Chapter 15 - Logistic Regression and Generalized Linear Models

23.2.15 Chapter 16 - Generalized Linear Mixed Models

23.3 From ST512

WE NEED TO ORGANIZE THESE UNDER DIFFERENT CHAPTERS AT SOME POINT Learning Objectives

1. Recognize a completely randomized design with one treatment factor and write the corresponding one-way analysis of variance model, with assumptions
2. Estimate treatment means
3. Estimate the variance among replicates within a treatment
4. Construct the analysis of variance table for a one factor analysis of variance, including computing degrees of freedom, sums of squares, mean squares, and F-ratios
5. Interpret results and draw conclusions from a one-factor analysis of variance
6. Estimate differences between two treatment means in a one factor analysis of variance
7. Test differences between two treatment means in a one factor analysis of variance
8. Construct a contrast to estimate or test a linear combination of treatment means
9. Estimate the standard error of a linear combination of treatment means
10. Make inferences about linear combinations of treatment means, including contrasts.
11. Obtain and understand SAS output for linear combinations of treatment means, including contrasts.
12. Explain when and why corrections for multiple comparisons are needed
13. Know when and how to use Tukey's correction for all pairwise comparisons
14. Compute Bonferroni confidence intervals
15. Create and interpret orthogonal contrasts.
16. Define main effects and interactions
17. Write contrasts to estimate main effects and interactions
18. Estimate these contrasts and their standard errors
19. Compute sums of squares associated with these contrasts
20. Test hypotheses about the main effects and interactions.
21. Identify and define simple effects.
22. Identify and define interaction effects.

23. Identify and define main effects.
24. Understand when to use simple, interaction, and main effects when drawing inferences in a two-way ANOVA.
25. Write the analysis of variance model and SAS code for a completely randomized design with two factors
26. Test hypotheses and interpret the analysis of variance for a factorial experiment.
27. Explain the appropriate use of correlations and compute the correlation coefficient
28. Read and interpret a scatterplot and guess the correlation coefficient by examination of a scatter plot
29. Interpret the strength and direction of association indicated by the correlation coefficient and judge when a correlation coefficient provides an appropriate summary of a bivariate relationship
30. Test the hypothesis that the correlation coefficient is zero using either a t-test or the Fisher z transformation, Compute confidence intervals using Fisher's z transformation
31. Write a statistical model for a straight line regression or a multiple regression and explain what all the terms of the model represent
32. Explain the assumptions underlying regression models, evaluate whether the assumptions are met
33. Estimate the intercept, slope and variance for a simple linear regression model
34. Fit a multiple regression model in SAS and interpret the output, use the coefficient of determination to evaluate model fit
35. Use a regression model to predict Y for new values of X
36. Estimate the variance and standard error of parameters in regression models, test hypotheses about the parameters, and construct confidence intervals for the parameters.
37. Explain the difference between a confidence interval and a prediction interval and know when to use each of them
38. Construct a confidence interval for the expected value of Y at a given value of X
39. Construct a prediction interval for a new value of Y at a given value of X
40. Write a linear model in matrix notation

41. Find the expectation and variance of a linear combination of random variables, $a'Y$
42. Set up the expressions to calculate parameter estimates and predicted values using the matrix form of the model
43. Estimate standard errors for parameter estimates and predicted values
44. Use extra sums of squares to test hypotheses about subsets of parameters
45. Construct indicator variables for including categorical regressor variables in a linear model
46. Understand how to interpret parameters of a general linear model with indicator variables
47. Estimate contrasts of treatment means and their standard errors using the general linear model notation and matrix form of the model
48. Compare nested models with a lack of fit test to select a model
49. Explain what a covariate is and how they are used
50. Explain the assumptions of the analysis of covariance model and determine when these assumptions are met
51. Fit an analysis of covariance model in SAS and conduct appropriate tests for treatment effects
52. Estimate and interpret treatment means and their standard errors adjusted for covariates using SAS, Construct confidence intervals for adjusted treatment means
53. Construct and estimate contrasts of treatment means adjusted for covariates and estimate the standard errors and confidence intervals of such contrasts.

Analysis of variance and design of experiments Recognize each of the following types of experimental designs and determine when each type would be advantageous. 1. completely randomized design 2. randomized complete block design 3. split plot design Recognize whether factors should be considered fixed effects or random effects and explain the scope of inference for each case. Recognize whether factors are crossed or nested. For all of the designs listed and for experiments with crossed and/or nested fixed factors, random factors, or a combination of fixed and random effects, be able to 1. Write the corresponding analysis of variance model, with assumptions, and define all terms 2. Estimate treatment means and their standard errors 3. Construct the analysis of variance table, including computing degrees of freedom, sums of squares, mean squares, and F-ratios 4. Determine whether the assumptions of the model are satisfied

5. Interpret results and draw conclusions 6. Construct and estimate linear combinations of treatment means and their standard errors 7. Test hypotheses and construct confidence intervals about linear combinations of treatment means 8. Explain when and why corrections for multiple comparisons are needed, know when and how to use Tukey's correction for all pairwise comparisons, compute Bonferroni confidence intervals 9. Create and interpret orthogonal contrasts. 10. Define and interpret main effects, simple effects and interactions 11. Use a table of expected mean squares to estimate variance components and determine appropriate F-statistics for testing effects in the analysis of variance 12. Interpret variance components and estimate and interpret the intraclass correlation coefficient. Regression and correlation Explain the appropriate use of correlations and compute the correlation coefficient, read and interpret a scatterplot and guess the correlation coefficient by examination of a scatter plot, test the hypothesis that the correlation coefficient is zero using either a t-test or the Fisher z transformation, compute confidence intervals using Fisher's z transformation You should be able to do the following for fitting models to describe the relationships of one or several variables to a response variable. The regressor variables may be continuous or categorical or a mix of the two (e.g., analysis of covariance models) 1. Write a general linear model, including assumptions, in standard or matrix notation, and explain what all the terms and assumptions represent. Be able to handle models that contain interaction terms, polynomial terms, and dummy variables. 2. Evaluate whether the model assumptions are met 3. Fit a general linear model in SAS and interpret the output 4. Work with the general linear model in matrix form, including finding the expectation and variance of a linear combination of regression coefficients or treatment means 5. Test hypotheses and construct confidence intervals for linear combinations of the parameters 6. Construct and interpret a confidence interval for the expected value of Y at a given value of X 7. Construct and interpret a prediction interval for a new value of Y at a given value of X 8. Use extra sums of squares to test hypotheses about subsets of parameters. 9. Explain what a covariate is and how covariates are used

23.4 For Point Estimates Chapter

- Definitions for Mean, Median, Quantile, Percentile
- Explain uses for the above
- Identify the correct point estimate to use for a given test
- Define Systematic Random Sample and Convenience Sample
- Explain strengths and weaknesses of each
- Identify conditions when Systematic and Convenience Sampling may not provide representative samples

Chapter 24

Appendix - Notation

24.1 Standard notation

Vectors of variables are denoted with Roman letters, such as x and Y . Capital letters denote random variables while lower case letters denote fixed variables. Note that these vectors may be of length 1 depending on context. Bolded values (x) denote matrices, and in the case of Y , possibly single-column matrices.

Unknown parameters are denoted with Greek letters, with boldface font indicating matrices.

In most models, Y will denote the univariate response, x will describe a matrix of predictor variables, and E a vector of random errors. The Greek letter β will be commonly used for regression parameters (either with subscripts for each values as in $\beta_0 + \beta_1 X_1$ or as a vector (as in $X\beta$). The letters i, j, k , and l will be most commonly used as subscripts or indices. N will typically denote a sample size (not a random vector), with subscripted versions (n_i) describing the number of observations in a group, and p describing the number of parameters in a model beyond the intercept.

We may therefore describe a simple linear regression model as:

$$Y = x\beta + E$$

In this model, Y is a $N \times 1$ random vector, x is a $N \times (p + 1)$ matrix of fixed values, and E is a $N \times 1$ vector.

π is typically used to describe probability parameters, as in Bernoulli or binomial random variables.

24.2 Mixed models

Still need to add something for this

24.3 Effects model representation

In the effects formulation of ANOVA models, additional greek letters (α , γ , etc.) will appear as parameter effects, as will μ , which will typically represent the grand mean. Group-specific means will be denoted via subscripts: μ_{ij} . When using this representation, it is convenient to describe a single observation as Y_{ijk} , which is the k th observation from the group with the i th level of the first factor and the j th level of the second factor. In the main effects version of this model, we have:

$$Y_{ijk} = \mu + \alpha_i + \gamma_j + E_{ijk}$$

We can therefore estimate μ_{ij} as $\hat{\mu}_{ij} = \frac{1}{n} \sum_{k=1}^n Y_{ijk} = \bar{Y}_{ij\cdot}$. This “dot” notation can be extended to any subscript and indicates summing over the index that has been replaced by the dot. Further note that the “hat” over a parameter value denotes the estimator for that parameter value, and the “bar” indicates an average. These features are used generally throughout this book.

24.4 Estimators vs. Estimates

If we want to get pedantic, we can differentiate between estimates and estimators in our notation. Estimators are functions of random variables used to estimate parameters. Estimates are realized values of estimators. To differentiate these, we use Roman letters with hats to represent estimators ($\hat{B} = (x'x)^{-1}x'Y$) and Greek letters with hats to represent estimates ($\hat{\beta} = 1.52$).