

# Statistics 500=Psychology 611=Biostatistics 550 Introduction to Regression and Anova

Paul R. Rosenbaum  
Professor, Statistics Department, Wharton School

## Description

Statistics 500/Psychology 611 is a second course in statistics for PhD students in the social, biological and business sciences. It covers multiple linear regression and analysis of variance. Students should have taken an undergraduate course in statistics prior to Statistics 500.

## Topics

- 1-Review of basic statistics.
- 2-Simple regression.
- 3-Multiple regression.
- 4-General linear hypothesis.
- 5-Woes of Regression Coefficients.
- 6-Transformations.
- 7-Polynomials.
- 8-Coded variables.
- 9-Diagnostics.
- 10-Variable selection.
- 11-One-way anova.
- 12-Two-way and factorial anova.

How do I get R for free? <http://cran.r-project.org/>

Final exam date: <http://www.upenn.edu/registrar/>

Holidays, breaks, last class: <http://www.upenn.edu/almanac/3yearcal.html>

My web page: <http://www-stat.wharton.upenn.edu/~rosenbap/index.html>

Email: [rosenbaum@wharton.upenn.edu](mailto:rosenbaum@wharton.upenn.edu)

Phone: 215-898-3120

Office: 473 Huntsman Hall (in the tower, 4<sup>th</sup> floor)

Office Hours: Tuesday 1:30-2:30 and by appointment.

**The bulk pack and course data in R are on my web page.**

## Overview

### Review of Basic Statistics

Descriptive statistics, graphs, probability, confidence intervals, hypothesis tests.

### Simple Regression

Simple regression uses a line with one predictor to predict one outcome.

### Multiple Regression

Multiple regression uses several predictors in a linear way to predict one outcome.

### General Linear Hypothesis

The general linear hypothesis asks whether several variables may be dropped from a multiple regression.

### Woes of Regression Coefficients

Discussion of the difficulties of interpreting regression coefficients and what you can do.

### Transformations

A simple way to fit curves or nonlinear models: transform the variables.

### Polynomials

Another way to fit curves: include quadratics and interactions.

### Coded Variables

Using nominal data (NY vs Philly vs LA) as predictors in regression.

### Diagnostics

How to find problems in your regression model: residual, leverage and influence.

### Variable Selection

Picking which predictors to use when many variables are available.

### One-Way Anova

Simplest analysis of variance: Do several groups differ, and if so, how?

### Two-Way Anova

Study two sources of variation at the same time.

### Factorial Anova

Study two or more treatments at once, including their interactions.

## Common Questions

Statistics Department Courses (times, rooms)

<http://www.upenn.edu/registrar/roster/stat.html>

Final Exams (dates, rules)

[http://www.upenn.edu/registrar/finals/spring05\\_index.html](http://www.upenn.edu/registrar/finals/spring05_index.html)

Computing and related help at Wharton

<http://inside.wharton.upenn.edu/>

Statistical Computing in the Psychology Department

<http://www.psych.upenn.edu>

When does the the course start? When does it end? Holidays?

<http://www.upenn.edu/almanac/3yearcal.html>

Does anybody have any record of this?

<http://www.upenn.edu/registrar/>

Huntsman Hall

[http://www.facilities.upenn.edu/mapsBldgs/view\\_bldg.php3?id=146](http://www.facilities.upenn.edu/mapsBldgs/view_bldg.php3?id=146)

[http://www.facilities.upenn.edu/mapsBldgs/view\\_map.php3?id=393](http://www.facilities.upenn.edu/mapsBldgs/view_map.php3?id=393)

Suggested reading

Box, G. E. P. (1966) Use and Abuse of Regression, *Technometrics*, 8, 625-629.

<http://www.jstor.org/> or

<http://www.jstor.org/stable/1266635?&Search=yes&term=abuse&term=box&list=hide&searchUri=%2Faction%2FdoAdvancedSearch%3Fq0%3Dbox%26f0%3Dau%26c0%3DAND%26q1%3Dabuse%26f1%3Dti%26c1%3DAND%26q2%3D%26f2%3Dall%26c2%3DAND%26q3%3D%26f3%3Dall%26wc%3Don%26sd%3D%26ed%3D%26la%3D%26jo%3D%26dc.Statistics%3DStatistics%26Search%3DSearch&item=1&ttl=1&returnArticleService=showArticle>

Helpful articles from JSTOR <http://www.jstor.org/>

1. The Analysis of Repeated Measures: A Practical Review with Examples  
B. S. Everitt *The Statistician*, Vol. 44, No. 1. (1995), pp. 113-135.
2. The hat matrix in regression and anova. D. Hoaglin and R. Welsh, *American Statistician*, Vol 32, (1978), pp. 17-22.
3. The Use of Nonparametric Methods in the Statistical Analysis of the Two-Period Change-Over Design Gary G. Koch  
*Biometrics*, Vol. 28, No. 2. (Jun., 1972), pp. 577-584.

## Some Web Addresses

Web page for Sheather's text

<http://www.stat.tamu.edu/~sheather/>

Amazon for Sheather's text (required)

[http://www.amazon.com/Modern-Approach-Regression-Springer-Statistics/dp/0387096078/ref=tmm\\_hrd\\_title\\_0/186-7302133-0606755?ie=UTF8&qid=1315493088&sr=1-1](http://www.amazon.com/Modern-Approach-Regression-Springer-Statistics/dp/0387096078/ref=tmm_hrd_title_0/186-7302133-0606755?ie=UTF8&qid=1315493088&sr=1-1)

Alternative text used several years ago (optional alternative, not suggested)

[http://www.amazon.com/Applied-Regression-Analysis-Multivariable-Methods/dp/0495384968/ref=sr\\_1\\_1?s=books&ie=UTF8&qid=1315493363&sr=1-1](http://www.amazon.com/Applied-Regression-Analysis-Multivariable-Methods/dp/0495384968/ref=sr_1_1?s=books&ie=UTF8&qid=1315493363&sr=1-1)

Good supplement about R (optional, suggested)

[http://www.amazon.com/Data-Analysis-Graphics-Using-Example-Based/dp/0521762936/ref=sr\\_1\\_1?s=books&ie=UTF8&qid=1315493138&sr=1-1](http://www.amazon.com/Data-Analysis-Graphics-Using-Example-Based/dp/0521762936/ref=sr_1_1?s=books&ie=UTF8&qid=1315493138&sr=1-1)

Review basic statistics, learn basic R (optional, use if you need it)

[http://www.amazon.com/Introductory-Statistics-R-Computing/dp/0387790535/ref=sr\\_1\\_1?s=books&ie=UTF8&qid=1315493184&sr=1-1](http://www.amazon.com/Introductory-Statistics-R-Computing/dp/0387790535/ref=sr_1_1?s=books&ie=UTF8&qid=1315493184&sr=1-1)

Excellent text, alternative to Sheather, more difficult than Sheather

[http://www.amazon.com/Applied-Regression-Analysis-Probability-Statistics/dp/0471170828/ref=sr\\_1\\_1?s=books&ie=UTF8&qid=1315493220&sr=1-1](http://www.amazon.com/Applied-Regression-Analysis-Probability-Statistics/dp/0471170828/ref=sr_1_1?s=books&ie=UTF8&qid=1315493220&sr=1-1)

Good text, alternative/supplement to Sheather, easier than Sheather

[http://www.amazon.com/Regression-Analysis-Example-Probability-Statistics/dp/0471746967/ref=tmm\\_hrd\\_title\\_0?ie=UTF8&qid=1315493316&sr=1-1](http://www.amazon.com/Regression-Analysis-Example-Probability-Statistics/dp/0471746967/ref=tmm_hrd_title_0?ie=UTF8&qid=1315493316&sr=1-1)

Free R manuals at R home page. Start with "An Introduction to R"

<http://cran.r-project.org/>

--> Manuals --> An Introduction to R

--> Search --> Paradis --> R for Beginners

My web page (bulk pack, course data)

<http://www-stat.wharton.upenn.edu/~rosenbap/index.html>

## Computing

**How do I get R for free?** <http://cran.r-project.org/>

After you have installed R, you can get the **course data** in the R-workspace on my web page: <http://www-stat.wharton.upenn.edu/~rosenbap/index.html>

I will probably add things to the R-workspace during the semester. So you will have to go back to my web page to **get the latest version**.

**A common problem:** You go to my web page and download the latest R-workspace, but it looks the same as the one you had before – the new stuff isn't there. This happens when your web browser thinks it has downloaded the file before and will save you time by not downloading it again. Bad web browser. You need to clear the cache; then it will get the new version.

**Most people find an R book helpful.** I recommend Maindonald and Braun, *Data Analysis and Graphics Using R*, published by Cambridge. A more basic book is Dalgaard, *Introductory Statistics with R*, published by Springer.

---

At <http://cran.r-project.org/>, click on **manuals** to get free documentation. "An Introduction to R" is there, and it is useful. When you get good at R, do a search at the site for Paradis' "R for Beginners," which is very helpful, but not for beginners.

---

## Textbook

My sense is that students need a textbook, not just the lectures and the bulk pack.

The 'required' textbook for the course is Sheather (2009) *A Modern Approach to Regression with R*, NY: Springer. There is a little matrix algebra in the book, but there is none in the course. Sheather replaces the old text, Kleinbaum, Kupper, Muller and Nizam, *Applied Regression and other Multivariable Methods*, largely because this book has become very expensive. An old used edition of Kleinbaum is a possible alternative to Sheather – it's up to you. Kleinbaum does more with anova for experiments. A book review by Gudmund R. Iversen of Swathmore College is available

at: <http://www.jstor.org/stable/2289682?&Search=yes&term=kleinbaum&term=kupper&list=hide&searchUri=%2Faction%2FdoAdvancedSearch%3Fq0%3Dkleinbaum%26f0%3Dau%26c0%3DAND%26q1%3Dkupper%26f1%3Dau%26c1%3DAND%26q2%3D%26f2%3Dall%26c2%3DAND%26q3%3D%26f3%3Dall%26wc%3Don%26re%3Don%26sd%3D%26ed%3D%26la%3D%26jo%3D%26dc.Statics%3DStatistics%26Search%3DSearch&item=6&ttl=7&returnArticleService=showArticle>

Some students might prefer one of the textbooks below, and they are fine substitutes.

If you would prefer an easier, less technical textbook, you might consider *Regression by Example* by Chatterjee and Hadi. The book has a nice chapter on transformations, but it barely covers anova. An earlier book, now out of print, with the same title by Chatterjee and Price is very similar, and probably available inexpensively used.

[http://www.amazon.com/Regression-Analysis-Example-Probability-Statistics/dp/0471746967/ref=sr\\_1\\_2?ie=UTF8&s=books&qid=1252524629&sr=1-2](http://www.amazon.com/Regression-Analysis-Example-Probability-Statistics/dp/0471746967/ref=sr_1_2?ie=UTF8&s=books&qid=1252524629&sr=1-2)

If you know matrix algebra, you might prefer the text *Applied Regression Analysis* by Draper and Smith. It is only slightly more difficult than Kleinbaum, and you can read around the matrix algebra.

[http://www.amazon.com/Applied-Regression-Analysis-Probability-Statistics/dp/0471170828/ref=sr\\_1\\_1?ie=UTF8&s=books&qid=1252524403&sr=1-1](http://www.amazon.com/Applied-Regression-Analysis-Probability-Statistics/dp/0471170828/ref=sr_1_1?ie=UTF8&s=books&qid=1252524403&sr=1-1)

If you use R, then as noted previously, I recommend the additional text Maindonald and Braun, *Data Analysis and Graphics Using R*, published by Cambridge. It is in its third edition, which is a tad more up to date than the first

or second editions, but you might prefer an inexpensive used earlier edition if you can find one.

## **Graded Work**

**Your grade is based on three exams.** Copies of old exams are at the end of this bulkpack. The first two exams are take-homes in which you do a data-analysis project. They are exams, so you do the work by yourself. The first exam covers the basics of multiple regression. The second exam covers diagnostics, model building and variable selection. The final exam is sometimes in-class, sometimes take home. The date of the final exam is determined by the registrar – see the page above for Common Questions. The decision about whether the final is in-class or take-home will be made after the first take-home is graded. That will be in the middle of the semester. If you need to make travel arrangements before the middle of the semester, you will need to plan around an in-class final.

**The best way to learn the material is to practice using the old exams.** There are three graded exams. If for each graded exam, you did two practice exams, then you would do nine exams in total, which means doing nine data analysis projects. With nine projects behind you, regression will start to be familiar.

## Review of Basic Statistics – Some Statistics

- The review of basic statistics is a quick review of ideas from your first course in statistics.

- $n$  measurements:  $X_1, X_2, \dots, X_n$

- **mean** (or average):  $\bar{X} = \frac{\sum_{i=1}^n X_i}{n} = \frac{X_1 + X_2 + \dots + X_n}{n}$

- **order statistics** (or data sorted from smallest to largest): Sort  $X_1, X_2, \dots, X_n$  placing the smallest first, the largest last, and write  $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ , so the smallest value is the first order statistic,  $X_{(1)}$ , and the largest is the  $n^{\text{th}}$  order statistic,  $X_{(n)}$ . If there are  $n=4$  observations, with values  $X_1 = 5, X_2 = 4, X_3 = 9, X_4 = 5$ , then the  $n=4$  order statistics are  $X_{(1)} = 4, X_{(2)} = 5, X_{(3)} = 5, X_{(4)} = 9$ .

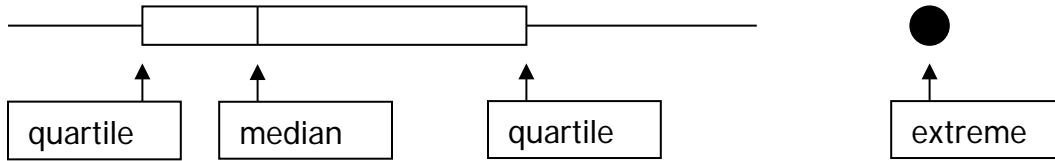
- **median** (or middle value): If  $n$  is odd, the median is the middle order statistic – e.g.,  $X_{(3)}$  if  $n=5$ . If  $n$  is even, there is no middle order statistic, and the median is the average of the two order statistics closest to the middle – e.g.,  $\frac{X_{(2)} + X_{(3)}}{2}$  if  $n=4$ . Depth of median is  $\frac{n+1}{2}$  where a “half” tells you to average two order statistics – for  $n=5$ ,  $\frac{n+1}{2} = \frac{5+1}{2} = 3$ , so the median is  $X_{(3)}$ , but for  $n=4$ ,  $\frac{n+1}{2} = \frac{4+1}{2} = 2.5$ , so the median is  $\frac{X_{(2)} + X_{(3)}}{2}$ .

The median cuts the data in half – half above, half below.

- **quartiles**: Cut the data in quarters – a quarter above the upper quartile, a quarter below the lower quartile, a quarter between the lower quartile and the median, a quarter between the median and the upper quartile. The **interquartile range** is the upper quartile minus the lower quartile.



- **boxplot:** Plots median and quartiles as a box, calls attention to extreme observations.



- **sample standard deviation:** square root of the typical squared deviation from the mean, sorta,

$$s = \sqrt{\frac{(X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + \dots + (X_n - \bar{X})^2}{n - 1}}$$

however, you don't have to remember this ugly formula.

- **location:** if I add a constant to every data value, a measure of location goes up by the addition of that constant.
- **scale:** if I multiply every data value by a constant, a measure of scale is multiplied by that constant, but a measure of scale does not change when I add a constant to every data value.

**Check your understanding:** What happens to the mean if I drag the biggest data value to infinity? What happens to the median? To a quartile? To the interquartile range? To the standard deviation? Which of the following are measures of location, of scale or neither: median, quartile, interquartile range, mean, standard deviation? In a boxplot, what would it mean if the median is closer to the lower quartile than to the upper quartile?

### Topic: Review of Basic Statistics – Probability

- **probability space:** the set of everything that can happen,  $\Omega$ . Flip two coins, dime and quarter, and the sample space is  $\Omega = \{HH, HT, TH, TT\}$  where HT means “head on dime, tail on quarter”, etc.
- **probability:** each element of the sample space has a probability attached, where each probability is between 0 and 1 and the total probability over the sample space is 1. If I flip two fair coins:  $\text{prob}(HH) = \text{prob}(HT) = \text{prob}(TH) = \text{prob}(TT) = \frac{1}{4}$ .
- **random variable:** a rule  $\mathbf{X}$  that assigns a number to each element of a sample space. Flip two coins, and the number of heads is a random variable: it assigns the number  $\mathbf{X}=2$  to HH, the number  $\mathbf{X}=1$  to both HT and TH, and the number  $\mathbf{X}=0$  to TT.
- **distribution of a random variable:** The chance the random variable  $\mathbf{X}$  takes on each possible value,  $x$ , written  $\text{prob}(\mathbf{X}=x)$ . Example: flip two fair coins, and let  $\mathbf{X}$  be the number of heads; then  $\text{prob}(\mathbf{X}=2) = \frac{1}{4}$ ,  $\text{prob}(\mathbf{X}=1) = \frac{1}{2}$ ,  $\text{prob}(\mathbf{X}=0) = \frac{1}{4}$ .
- **cumulative distribution of a random variable:** The chance the random variable  $\mathbf{X}$  is less than or equal to each possible value,  $x$ , written  $\text{prob}(\mathbf{X} \leq x)$ . Example: flip two fair coins, and let  $\mathbf{X}$  be the number of heads; then  $\text{prob}(\mathbf{X} \leq 0) = \frac{1}{4}$ ,  $\text{prob}(\mathbf{X} \leq 1) = \frac{3}{4}$ ,  $\text{prob}(\mathbf{X} \leq 2) = 1$ . Tables at the back of statistics books are often cumulative distributions.
- **independence of random variables:** Captures the idea that two random variables are unrelated, that neither predicts the other. The formal definition which follows is not intuitive – you get to like it by trying many intuitive examples, like unrelated coins and taped coins, and finding the definition always works. Two random variables,  $\mathbf{X}$  and  $\mathbf{Y}$ , are independent if the chance that simultaneously  $\mathbf{X}=x$  and  $\mathbf{Y}=y$  can be found by multiplying the separate probabilities

$$\text{prob}(\mathbf{X}=x \text{ and } \mathbf{Y}=y) = \text{prob}(\mathbf{X}=x) \text{ prob}(\mathbf{Y}=y) \quad \text{for every choice of } x,y.$$

**Check your understanding:** Can you tell exactly what happened in the sample space from the value of a random variable? Pick one: Always, sometimes, never. For people, do you think  $\mathbf{X}$ =height and  $\mathbf{Y}$ =weight are independent? For undergraduates, might  $\mathbf{X}$ =age and  $\mathbf{Y}$ =gender (1=female, 2=male) be independent? If I flip two fair coins, a dime and a quarter, so that  $\text{prob}(\text{HH}) = \text{prob}(\text{HT}) = \text{prob}(\text{TH}) = \text{prob}(\text{TT}) = 1/4$ , then is it true or false that getting a head on the dime is independent of getting a head on the quarter?

### Topic: Review of Basics – Expectation and Variance

- **Expectation:** The expectation of a random variable  $\mathbf{X}$  is the sum of its possible values weighted by their probabilities,

$$E(\mathbf{X}) = \sum_x x \cdot \text{prob}(\mathbf{X} = x)$$

- **Example:** I flip two fair coins, getting  $\mathbf{X}=0$  heads with probability  $1/4$ ,  $\mathbf{X}=1$  head with probability  $1/2$ , and  $\mathbf{X}=2$  heads with probability  $1/4$ ; then the expected number of heads is  $E(\mathbf{X}) = 0 \cdot \frac{1}{4} + 1 \cdot \frac{1}{2} + 2 \cdot \frac{1}{4} = 1$ , so I expect 1 head when I flip two fair coins. Might actually get 0 heads, might get 2 heads, but 1 head is what is typical, or expected, on average.
- **Variance and Standard Deviation:** The standard deviation of a random variable  $\mathbf{X}$  measures how far  $\mathbf{X}$  typically is from its expectation  $E(\mathbf{X})$ . Being too high is as bad as being too low – we care about errors, and don't care about their signs. So we look at the squared difference between  $\mathbf{X}$  and  $E(\mathbf{X})$ , namely  $\mathbf{D} = \{\mathbf{X} - E(\mathbf{X})\}^2$ , which is, itself, a random variable. The variance of  $\mathbf{X}$  is the expected value of  $\mathbf{D}$  and the standard deviation is the square root of the variance,  $\text{var}(\mathbf{X}) = E(\mathbf{D})$  and  $\text{st. dev.}(\mathbf{X}) = \sqrt{\text{var}(\mathbf{X})}$ .
- **Example:** I independently flip two fair coins, getting  $\mathbf{X}=0$  heads with probability  $1/4$ ,  $\mathbf{X}=1$  head with probability  $1/2$ , and  $\mathbf{X}=2$  heads with probability  $1/4$ . Then  $E(\mathbf{X})=1$ , as noted above. So  $\mathbf{D} = \{\mathbf{X} - E(\mathbf{X})\}^2$  takes the value  $\mathbf{D} =$

$(0 - 1)^2 = 1$  with probability  $\frac{1}{4}$ , the value  $\mathbf{D} = (1 - 1)^2 = 0$  with probability  $\frac{1}{2}$ , and the value  $\mathbf{D} = (2 - 1)^2 = 1$  with probability  $\frac{1}{4}$ . The variance of  $\mathbf{X}$  is the expected value of  $\mathbf{D}$  namely:  $\text{var}(\mathbf{X}) = E(\mathbf{D}) = 1 \cdot \frac{1}{4} + 0 \cdot \frac{1}{2} + 1 \cdot \frac{1}{4} = \frac{1}{2}$ . So the standard deviation is  $st.dev.(\mathbf{X}) = \sqrt{\text{var}(\mathbf{X})} = \sqrt{\frac{1}{2}} = 0.707$ . So when I flip two fair coins, I expect one head, but often I get 0 or 2 heads instead, and the typical deviation from what I expect is 0.707 heads. This 0.707 reflects the fact that I get exactly what I expect, namely 1 head, half the time, but I get 1 more than I expect a quarter of the time, and one less than I expect a quarter of the time.

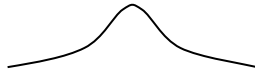
**Check your understanding:** If a random variable has zero variance, how often does it differ from its expectation? Consider the height  $\mathbf{X}$  of male adults in the US. What is a reasonable number for  $E(\mathbf{X})$ ? Pick one: 4 feet, 5'9", 7 feet. What is a reasonable number for  $st.dev.(\mathbf{X})$ ? Pick one: 1 inch, 4 inches, 3 feet. If I independently flip three fair coins, what is the expected number of heads? What is the standard deviation?

## Topic: Review of Basics – Normal Distribution

- Continuous random variable:** A continuous random variable can take values with any number of decimals, like 1.2361248912. Weight measured perfectly, with all the decimals and no rounding, is a continuous random variable. Because it can take so many different values, each value winds up having probability zero. If I ask you to guess someone's weight, not approximately to the nearest millionth of a gram, but rather exactly to all the decimals, there is no way you can guess correctly – each value with all the decimals has probability zero. But for an interval, say the nearest kilogram,

there is a nonzero chance you can guess correctly. This idea is captured in by the density function.

- **Density Functions:** A density function defines probability for a continuous random variable. It attaches zero probability to every number, but positive probability to ranges (e.g., nearest kilogram). The probability that the random variable  $\mathbf{X}$  takes values between 3.9 and 6.2 is the area under the density function between 3.9 and 6.2. The total area under the density function is 1.
- **Normal density:** The Normal density is the familiar “bell shaped curve”.



The standard Normal distribution has expectation zero, variance 1, standard deviation  $1 = \sqrt{1}$ . About 2/3 of the area under the Normal density is between  $-1$  and  $1$ , so the probability that a standard Normal random variable takes values between  $-1$  and  $1$  is about 2/3. About 95% of the area under the Normal density is between  $-2$  and  $2$ , so the probability that a standard Normal random variable takes values between  $-2$  and  $2$  is about .95. (To be more precise, there is a 95% chance that a standard Normal random variable will be between  $-1.96$  and  $1.96$ .) If  $\mathbf{X}$  is a standard Normal random variable, and  $\mu$  and  $\sigma > 0$  are two numbers, then  $\mathbf{Y} = \mu + \sigma\mathbf{X}$  has the Normal distribution with expectation  $\mu$ , variance  $\sigma^2$  and standard deviation  $\sigma$ , which we write  $N(\mu, \sigma^2)$ . For example,  $\mathbf{Y} = 3 + 2\mathbf{X}$  has expectation 3, variance 4, standard deviation 2, and is  $N(3,4)$ .

- **Normal Plot:** To check whether or not data,  $X_1, \dots, X_n$  look like they came from a Normal distribution, we do a Normal plot. We get the order statistics – just the data sorted into order – or  $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$  and plot this ordered data against what ordered data from a standard Normal distribution should look like. The computer takes care of the details. A straight line in a

Normal plot means the data look Normal. A straight line with a couple of strange points off the lines suggests a Normal with a couple of strange points (called outliers). Outliers are extremely rare if the data are truly Normal, but real data often exhibit outliers. A curve suggest data that are not Normal. Real data wiggle, so nothing is ever perfectly straight. In time, you develop an eye for Normal plots, and can distinguish wiggles from data that are not Normal.

### Topic: Review of Basics – Confidence Intervals

- Let  $X_1, \dots, X_n$  be  $n$  independent observations from a Normal distribution with expectation  $\mu$  and variance  $\sigma^2$ . A compact way of writing this is to say  $X_1, \dots, X_n$  are iid from  $N(\mu, \sigma^2)$ . Here, iid means independent and identically distributed, that is, unrelated to each other and all having the same distribution.
- How do we know  $X_1, \dots, X_n$  are iid from  $N(\mu, \sigma^2)$ ? We don't! But we check as best we can. We do a boxplot to check on the shape of the distribution. We do a Normal plot to see if the distribution looks Normal. Checking independence is harder, and we don't do it as well as we would like. We do look to see if measurements from related people look more similar than measurements from unrelated people. This would indicate a violation of independence. We do look to see if measurements taken close together in time are more similar than measurements taken far apart in time. This would indicate a violation of independence. Remember that statistical methods come with a warrantee of good performance if certain assumptions are true, assumptions like  $X_1, \dots, X_n$  are iid from  $N(\mu, \sigma^2)$ . We check the assumptions to make sure we get the promised good performance of statistical methods. Using statistical methods when the assumptions are not

true is like putting your CD player in washing machine – it voids the warrantee.

- To begin again, having checked every way we can, finding no problems, assume  $X_1, \dots, X_n$  are iid from  $N(\mu, \sigma^2)$ . We want to estimate the expectation  $\mu$ . We want an interval that in most studies winds up covering the true value of  $\mu$ . Typically we want an interval that covers  $\mu$  in 95% of studies, or a **95% confidence interval**. Notice that the promise is about what happens in most studies, not what happened in the current study. If you use the interval in thousands of unrelated studies, it covers  $\mu$  in 95% of these studies and misses in 5%. You cannot tell from your data whether this current study is one of the 95% or one of the 5%. All you can say is the interval usually works, so I have confidence in it.
- If  $X_1, \dots, X_n$  are iid from  $N(\mu, \sigma^2)$ , then the confidence interval uses the sample mean,  $\bar{X}$ , the sample standard deviation,  $s$ , the sample size,  $n$ , and a critical value obtained from the t-distribution with  $n-1$  degrees of freedom, namely the value,  $t_{0.025}$ , such that the chance a random variable with a t-distribution is above  $t_{0.025}$  is 0.025. If  $n$  is not very small, say  $n > 10$ , then  $t_{0.025}$  is near 2. The 95% confidence interval is:

$$\bar{X} \pm (\text{allowance for error}) = \bar{X} \pm \frac{t_{0.025} \cdot s}{\sqrt{n}}$$

## Topic: Review of Basics – Hypothesis Tests

- Null Hypothesis:** Let  $X_1, \dots, X_n$  be  $n$  independent observations from a Normal distribution with expectation  $\mu$  and variance  $\sigma^2$ . We have a particular value of  $\mu$  in mind, say  $\mu_0$ , and we want to ask if the data contradict this value. It means something special to us if  $\mu_0$  is the correct value – perhaps it means the treatment has no effect, so the treatment should be discarded. We wish to test the null hypothesis,  $H_0: \mu = \mu_0$ . Is the null hypothesis plausible? Or do the data force us to abandon the null hypothesis?
- Logic of Hypothesis Tests:** A hypothesis test has a long-winded logic, but not an unreasonable one. We say: Suppose, just for the sake of argument, not because we believe it, that the null hypothesis is true. As is always true when we suppose something for the sake of argument, what we mean is: Let's suppose it and see if what follows logically from supposing it is believable. If not, we doubt our supposition. So suppose  $\mu_0$  is the true value after all. Is the data we got, namely  $X_1, \dots, X_n$ , the sort of data you would usually see if the null hypothesis were true? If it is, if  $X_1, \dots, X_n$  are a common sort of data when the null hypothesis is true, then the null hypothesis looks sorta ok, and we *accept* it. Otherwise, if there is no way in the world you'd ever see data anything remotely like our data,  $X_1, \dots, X_n$ , if the null hypothesis is true, then we can't really believe the null hypothesis having seen  $X_1, \dots, X_n$ , and we *reject* it. So the basic question is: Is data like the data we got commonly seen when the null hypothesis is true? If not, the null hypothesis has gotta go.
- P-values or significance levels:** We measure whether the data are commonly seen when the null hypothesis is true using something called the P-value or significance level. Supposing the null hypothesis to be true, the P-value is the chance of data at least as inconsistent with the null hypothesis as



the observed data. If the P-value is  $\frac{1}{2}$ , then half the time you get data as or more inconsistent with the null hypothesis as the observed data – it happens half the time by chance – so there is no reason to doubt the null hypothesis. But if the P-value is 0.000001, then data like ours, or data more extreme than ours, would happen only one time in a million by chance if the null hypothesis were true, so you gotta be having some doubts about this null hypothesis.

- **The magic 0.05 level:** A convention is that we “reject” the null hypothesis when the P-value is less than 0.05, and in this case we say we are testing at **level** 0.05. Scientific journals and law courts often take this convention seriously. It is, however, only a convention. In particular, sensible people realize that a P-value of 0.049 is not very different from a P-value of 0.051, and both are very different from P-values of 0.00001 and 0.3. It is best to report the P-value itself, rather than just saying the null hypothesis was rejected or accepted.
- **Example:** You are playing 5-card stud poker and the dealer sits down and gets 3 royal straight flushes in a row, winning each time. The null hypothesis is that this is a fair poker game and the dealer is not cheating. Now, there are or 2,598,960 five-card stud poker hands, and 4 of these are royal straight flushes, so the chance of a royal straight flush in a fair game is

$$\frac{4}{2,598,960} = 0.000001539. \text{ In a fair game, the chance of three royal straight}$$

flushes in a row is  $0.000001539 \times 0.000001539 \times 0.000001539 = 3.6 \times 10^{-18}$ .

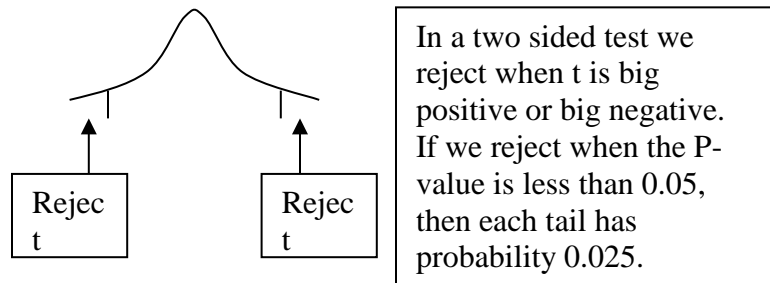
(Why do we multiply probabilities here?) Assuming the null hypothesis, for the sake of argument, that is assuming he is not cheating, the chance he will get three royal straight flushes in a row is very, very small – that is the P-value or significance level. The data we see is highly improbable if the null hypothesis were true, so we doubt it is true. Either the dealer got very, very lucky, or he cheated. This is the logic of all hypothesis tests.

- One sample t-test:** Let  $X_1, \dots, X_n$  be  $n$  independent observations from a Normal distribution with expectation  $\mu$  and variance  $\sigma^2$ . We wish to test the null hypothesis,  $H_0: \mu = \mu_0$ . We do this using the one-sample t-test:

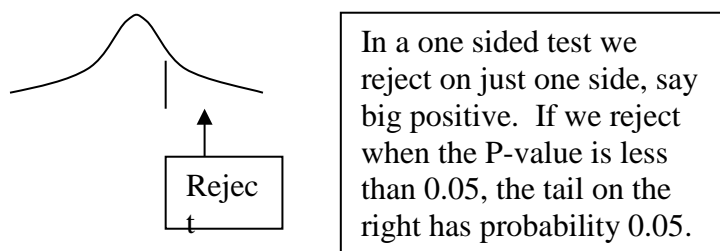
$$t = \frac{\sqrt{n}(\bar{X} - \mu_0)}{s}$$

looking this up in tables of the t-distribution with  $n-1$  degrees of freedom to get the P-value.

- One-sided vs Two-sided tests:** In a two-sided test, we don't care whether  $\bar{X}$  is bigger than or smaller than  $\mu_0$ , so we reject at the 5% level when  $|t|$  is one of the 5% largest values of  $|t|$ . This means we reject for 2.5% of  $t$ 's that are very positive and 2.5% of  $t$ 's that are very negative:



In a one sided test, we do care, and only want to reject when  $\bar{X}$  is on one particular side of  $\mu_0$ , say when  $\bar{X}$  is bigger than  $\mu_0$ , so we reject at the 5% level when  $t$  is one of the 5% largest values of  $t$ . This means we reject for the 5% of  $t$ 's that are very positive:



- Should I do a one-sided or a two-sided test:** Scientists mostly report two-sided tests.

## Obtaining a Confidence Interval for a Parameter by Inverting a Test

1. We have a valid way of testing the hypothesis  $H_0 : \theta = \theta_0$  that works for any real number  $\theta_0$ . For example, we can test  $H_0 : \theta = 0$  and  $H_0 : \theta = 1$  and  $H_0 : \theta = 1.263$ , etc.
2. A valid level  $\alpha$  test of  $H_0 : \theta = \theta_0$  falsely rejects  $H_0$  when it is true with probability at most  $\alpha$ . That is the definition of “valid” in the phrase “valid test”. The popular  $\alpha$  is  $\alpha = 0.05$ , but that is merely a convention.
3. Suppose we test every possible value  $\theta_0$  at level  $\alpha = 0.05$ , keeping the values we do not reject in a set  $\mathcal{C}$ . So  $\theta_0$  is in  $\mathcal{C}$ , or equivalently  $\theta_0 \in \mathcal{C}$ , if we did not reject  $\theta_0$  when we tested it because its  $P$ -value was  $> 0.05$ .
4. The set  $\mathcal{C}$  is a random set. We computed it from the data, and the data are random, so  $\mathcal{C}$  is random.
5. There is one true value,  $\theta^*$ , of  $\theta$ . This true value  $\theta^*$  is a number, not a random variable, a number we do not know but a number nonetheless. What is the chance that the random set  $\mathcal{C}$  fails to cover the fixed number  $\theta^*$ ?
6. We know that we will, sooner or later, test the true value,  $H_0 : \theta = \theta^*$ , because we test every value. We are happy to reject false values of  $\theta$  — after all, they are false — but when we test the one true value,  $\theta^*$ , we don’t want to reject it.
7. When we test the true value,  $H_0 : \theta = \theta^*$ , the chance that we falsely reject it is at most  $\alpha = 0.05$ . That is, again, what it means to use a valid test.
8. Therefore, the probability that the random interval  $\mathcal{C}$  fails to cover  $\theta^*$  is at most  $\alpha = 0.05$ . We call  $\mathcal{C}$  a 95% confidence interval.

## REGRESSION ASSUMPTIONS

| Assumption                     | If untrue:  | How to detect:   |
|--------------------------------|---|--|
| Independent errors             | 95% confidence intervals may cover much less than 95% of the time. Tests that reject with $p < 0.05$ may reject true hypotheses more than 5% of the time. You may think you have much more information than you do. | Often hard to detect. Questions to ask yourself: (i) Are the observations clustered into groups, such as several measurements on the same person? (ii) Are observations repeated over time?      |
| Normal errors                  | Thick tails and outliers may distort estimates, and they may inflate the estimated error variance, so that confidence intervals are too long, and hypothesis tests rarely reject false hypotheses.                  | Do a Normal quantile plot. This is the one use of the Normal quantile plot. A more or less straight line in the plot suggests the data are approximately Normal.                                 |
| Errors have constant variance. | Least squares gives equal weight to all observations, but if some observations are much more stable than others, it is not sensible to give equal weight to all observations.                                       | Plot the residuals against the predicted values. A fan shape in the plot – narrow on one end, wide on the other – suggests unequal variances. Can also plot residuals against individual $x$ 's. |
| Model is linear.               | Linear model may not fit, or may give the wrong interpretation of the data.   | Plot the residuals against the predicted values. Curves, such as a U-shape, suggest the relationship is not linear. Can also plot residuals against individual $x$ 's.                           |

### Statistics 500: Basic Statistics Review

- **Reading:** In Kleinbaum, read chapter 3.
- **Practice:** The blood pressure data we discussed in class is given below. It is from MacGregor, et. al. (1979) British Medical Journal, 2, 1106-9. It is the change in systolic blood pressure two hours after taking Captopril, in mm Hg, after-before, so a negative number means a decline in blood pressure. Use JMP or another package to do a Normal plot, a boxplot and a t-test. Think about how you would describe what you see.

| Patient # | Change in bp |
|-----------|--------------|
| 1         | -9           |
| 2         | -4           |
| 3         | -21          |
| 4         | -3           |
| 5         | -20          |
| 6         | -31          |
| 7         | -17          |
| 8         | -26          |
| 9         | -26          |
| 10        | -10          |
| 11        | -23          |
| 12        | -33          |
| 13        | -19          |
| 14        | -19          |
| 15        | -23          |

**Homework:** The following data are from Kaneto, Kosaka and Nakao (1969) *Endocrinology*, 80, 530-536. It is an experiment on 7 dogs. Question is whether stimulation of the vagus nerve increases levels of immunoreactive insulin in the blood. Two measurements were taken on each dog, one before, one five minutes after stimulation. The measurements are blood lead levels of immunoreactive insulin ( $\mu U / ml$ ).

| <b><i>Dog</i></b> | Before | After |
|-------------------|--------|-------|
| 1                 | 350    | 480   |
| 2                 | 200    | 130   |
| 3                 | 240    | 250   |
| 4                 | 290    | 310   |
| 5                 | 90     | 280   |
| 6                 | 370    | 1450  |
| 7                 | 240    | 280   |

Do an appropriate analysis.

---

## Topic: Simple Regression

- **Simple regression:** Fitting a line a response  $Y$  using one predictor  $X$ .
- **Data:** 48 contiguous states in 1972,  $i=1, \dots, 48$ .  $Y = \text{FUEL} =$  motor fuel consumption per person in gallons per person.  $X = \text{TAX} =$  motor fuel tax rate in cents per gallon.
- **First thing you do:** Plot the data.
- **Least squares:** Fit the line  $\hat{Y}_i = \hat{\alpha} + \hat{\beta}X_i$  by minimizing the sum of the squares of the residuals  $Y_i - \hat{Y}_i$  around the line,  $\sum_{i=1}^n (Y_i - \hat{Y}_i)^2$ .
- **Plot the residuals:** After you fit a line, you plot the residuals,  $Y_i - \hat{Y}_i$ . They tell you where and how the line fits poorly. The minimum is: (i) a boxplot of residuals, (ii) a Normal plot of residuals, (iii) a plot of residuals vs predicted values,  $Y_i - \hat{Y}_i$  vs  $\hat{Y}_i$ .
- **Statistical Model:** The statistical model says:
 
$$Y_i = \alpha + \beta X_i + \varepsilon_i \text{ where the } \varepsilon_i \text{ are iid } N(0, \sigma^2),$$
 so the  $Y$ 's were generated by a true line,  $\alpha + \beta X_i$ , which we do not know, plus errors  $\varepsilon_i$  that are independent of each other and Normal with mean zero and constant variance  $\sigma^2$ . We use the residual plots to check whether the model is a reasonable description of the data. The line fitted by least squares,  $\hat{Y}_i = \hat{\alpha} + \hat{\beta}X_i$ , is our estimate of the true line  $\alpha + \beta X_i$ .
- **Properties of least squares estimates:** The least squares estimators are great estimators – the best there are – when the model is correct, and not so great when the model is wrong. Checking the model is checking whether we are getting good estimates. When the model is true, least squares estimates are **unbiased**, that is, correct in expectation or on average, and they have **minimum variance** among all unbiased estimates, so they are the most stable, most accurate unbiased estimates (but only if the model is correct!).

They are not robust to outliers – one weird observation can move the fitted line anywhere it wants.

- **Basic Regression Output**

| Variable | Estimated Coefficient | Estimated Standard Error of Estimated Coefficient | t-ratio                                 |
|----------|-----------------------|---|---|
| Constant | $\hat{\alpha}$        | $se(\hat{\alpha})$                                | $\frac{\hat{\alpha}}{se(\hat{\alpha})}$ |
| X        | $\hat{\beta}$         | $se(\hat{\beta})$                                 | $\frac{\hat{\beta}}{se(\hat{\beta})}$   |

- **Hypothesis tests:** Use the t-ratio to test the null hypothesis  $H_0: \beta = 0$ . Under the model, the hypothesis  $H_0: \beta = 0$  implies X and Y are unrelated.
- **Confidence intervals:** Under the model, a 95% confidence interval for  $\beta$  is:

$$estimate \pm allowance = \hat{\beta} \pm t_{0.025} \cdot se(\hat{\beta})$$

where  $t_{0.025}$  is the upper 2.5% point of the t-distribution with n-2 degrees of freedom. When n-2 is not small, the interval is almost (but not quite)  $\hat{\beta} \pm 2 \cdot se(\hat{\beta})$ .

- **Points on a line vs Predictions:** Two problems look almost the same, but really are very different. One asks: Where is the line at X=8.5? That is, what is  $\alpha + \beta 8.5$ ? That problem gets easier as I collect more data and learn where the line really is. The other asks: Where will a new observation on Y be if X=8.5? That is, what is  $\alpha + \beta 8.5 + \varepsilon_{new}$ ? That problem always stays pretty hard, no matter how much data I collect, because I can't predict the new error,  $\varepsilon_{new}$ , for this new observation no matter how well I know where the line is. Important thing is to make sure you know which answer you



want and to use the right method for that answer. They look similar, but they're not.

- **Regression Anova Table:** Partitions the total variation (or sum of squares)

in the data about the mean, namely  $\sum_{i=1}^n (Y_i - \bar{Y})^2$  into two parts that add back

to the total, namely the variation fitted by the regression,  $\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$ , and

the variation in the residuals,  $\sum_{i=1}^n (Y_i - \hat{Y}_i)^2$ . Degrees of freedom measure

keep track of how many distinct numbers are really being described by a sum of squares. In simple regression, the variation fitted by the regression is just fitted by the slope,  $\hat{\beta}$ , which is just one number, so this sum of squares has

1 degree of freedom. A mean square is the ratio  $\frac{\text{sum of squares}}{\text{degrees of freedom}}$ . The

F-ratio is the ratio of two mean squares, a signal to noise ratio. The F-ratio is used to test that all the slopes are zero.

- **Simple correlation:** If the data fall perfectly on a line tilted up, the correlation  $r$  is 1. If the data fall perfectly on a line tilted down, the correlation  $r$  is  $-1$ . If a line is not useful for predicting  $Y$  from  $X$ , the correlation  $r$  is 0. Correlation is always between  $-1$  and 1. The correlation between  $Y$  and  $X$  is the regression coefficient of standardized  $Y$  on

standardized  $X$ , that is, the regression of  $\frac{Y_i - \bar{Y}}{\text{st. dev}(Y)}$  on  $\frac{X_i - \bar{X}}{\text{st. dev}(X)}$ . In

simple, one-predictor regression, the square of the correlation,  $r^2$ , is the percent of variation fitted by the regression, so it summarizes the anova table. Correlation discards the units of measurement, which limits its usefulness.

## Homework: Vocabulary Data

**Homework:** The following data are from M. E. Smith, (1926), "An investigation of the development of the sentence and the extent of vocabulary in young children." It relates the  $X$ =age of children in years to their  $Y$ =vocabulary size in words. I would like you to do a regression of  $Y$  and  $X$ , look closely at what you've done, and comment on what it all means. You should turn in (1) one paragraph of text, (2) linear regression output, (3) at most two plots you find interesting and helpful in thinking about what is special about these data. This is real data, so it is not a "trick question", but it does require some real thought about what makes sense and what is happening.

| <b>X=age</b> | <b>Y=vocabulary</b> |
|--------------|---------------------|
| 0.67         | 0                   |
| 0.83         | 1                   |
| 1            | 3                   |
| 1.25         | 19                  |
| 1.5          | 22                  |
| 1.75         | 118                 |
| 2            | 272                 |
| 2.5          | 446                 |
| 3            | 896                 |
| 3.5          | 1,222               |
| 4            | 1,540               |
| 4.5          | 1,870               |
| 5            | 2,072               |
| 5.5          | 2,289               |
| 6            | 2,562               |

## Topic: Multiple Regression

- **Multiple regression:** Uses several predictor variables  $X_1, X_2, \dots, X_k$  to fit a single response variable  $Y$ .
- **FUEL DATA:** Trying to predict  $Y = \text{FUEL}$  from  $X_1 = \text{TAX}$  and a second predictor,  $X_2 = \text{LICENSES}$ .
- **Least squares fit:** Multiple regression fits a plane

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \dots + \hat{\beta}_k X_{ki} \text{ making the residuals } Y_i - \hat{Y}_i \text{ small, in}$$

the sense that the sum of the squares of the residuals, namely,  $\sum_{i=1}^n (Y_i - \hat{Y}_i)^2$ ,

is minimized.

- **Multiple Correlation:** The multiple correlation,  $R$ , is the ordinary correlation between the observed  $Y_i$  and the fitted  $\hat{Y}_i$ . The square of the multiple correlation,  $R^2$ , is the percent of variation fitted by the regression, that is, regression sum of squares in the ANOVA table divided by the total sum of squares of  $Y$  around its mean.
- **Fit vs Prediction:** Fit refers to how close the model is to the observed data. Prediction refers to how close the model is to new data one might collect. They are not the same. Adding variables, even junk variables, always improves the fit, but the predictions may get better or worse.  $R^2$  is a measure of fit, not of prediction. We will develop a measure of prediction,  $C_p$ , later in the course.
- **Statistical Model:** The model underlying multiple regression says:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + \varepsilon_i$$

where the  $\varepsilon_i$  are independent  $N(0, \sigma^2)$

The true model is unknown, but the least squares fit is an estimate.

- **Hypothesis Tests and Confidence Intervals for a Coefficient:** Testing a hypothesis about a regression coefficient, say  $H_0: \beta_5 = 0$ , is done using the t-

statistic as in simple regression. Confidence intervals are also done as in simple regression.

- **Testing that all coefficients are zero:** The F-test from the ANOVA table is used to test  $H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$ .
- **Residual analysis:** One checks the model by plotting the residuals. The minimum is a plot of residuals against predicted, a boxplot of residuals, and a Normal plot of residuals, as in simple regression.

## Topic: General Linear Hypothesis

- What is it? In model,

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + \varepsilon_i$$

where the  $\varepsilon_i$  are independent  $N(0, \sigma^2)$ ,

we know how to test a hypothesis about one coefficient, say  $H_0: \beta_5 = 0$ , (t-test) and we know how to test that all of the variables are unneeded,

$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$  (F-test from regression anova table). The general linear hypothesis says that a particular subset of the coefficients is zero. For example, the hypothesis might say that the last  $k-J$  variables are not needed,

$$H_0: \beta_{J+1} = \beta_{J+2} = \dots = \beta_k = 0.$$

- Why do this? Generally, a hypothesis expresses an idea. Some ideas need to be expressed using more than one variable. For example, in the FUEL data, the 48 states might be divided into five regions, Northeast, Southeast, Midwest, Mountain, and Pacific, say. Later on, we will see how to code region into several variables in a regression. Testing whether “REGION” matters is testing whether all of these variables can be dropped from the model.
- Comparing Two Models: The test involves comparing two models, a reduced model which assumes the hypothesis is true, and a full model which assumes it is false. To test  $H_0: \beta_{J+1} = \beta_{J+2} = \dots = \beta_k = 0$ , one fits the full model:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + \varepsilon_i$$

and the reduced model without variables  $X_{J+1}, \dots, X_k$ ,

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_J X_{Ji} + \varepsilon_i,$$

and the test is based on comparing the ANOVA tables for these two models. Details in the textbook.

## Topic: Woes of Regression Coefficients

- The phrase, “the woes of regression coefficients” is due to Fred Mosteller and John Tukey in a standard text, *Data Analysis and Regression*. The bulk pack contains another standard reading: George Box’s paper “The use and abuse of regression”. (The paper is mostly easy to read, but contains some technical material – just skip the technical material. The main points are not technical and not difficult.)
- The issue concerns the interpretation of regression coefficients. The bottom line is that it is hard to interpret regression coefficients. The reason is that whenever you add (or delete) a variable from a regression model, all of the other coefficients change to reflect the added (or deleted) variable. People want (but they can’t have) a way of speaking of THE coefficient of a variable, but actually the coefficient of a variable always depends on what other variables are in the model with it. People want to say that  $\beta_j$  is the change in Y expected from a one unit change in  $X_j$ , but it simply isn’t true. It can’t be true, since  $\beta_j$  keeps changing as variables are added or deleted from a model, whereas changing  $X_j$  out in the world has nothing to do with which variables I put in the model.
- The bottom line is this: Whenever you hear people say that changing  $X_j$  will produce a particular change in Y, and they say they know this solely because they did a regression, you should be a little skeptical. There is more to knowing something like this than just running regressions.

## Topic: Transformations

- **Key Idea:** Fit many kinds of curved (i.e., nonlinear) models by transforming the variables and fitting a linear model to the transformed variables.

- **Logs.** For  $b > 0$ , the base  $b$  log has the property that:

$$y = b^a \quad \text{is the same as} \quad \log_b(y) = a.$$

Common choices of the base  $b$  are  $b=10$ ,  $b=2$  and  $b=e=2.71828\dots$  for natural logs. Outside high school, if no base is mentioned (e.g.,  $\log(y)$ ) it usually means base  $e$  or natural logs. Two properties we use often are:  $\log(xy)=\log(x)+\log(y)$  and  $\log(y^a) = a \cdot \log(y)$ .

- **Why transform?** (i) You plot the data and it is curved, so you can't fit a line. (ii) The  $Y$ 's have boundaries (e.g.,  $Y$  must be  $>0$  or  $Y$  must be between 0 and 1), but linear regression knows nothing of the boundaries and overshoots them, producing impossible  $\hat{Y}$ 's. (iii) The original data violate the linear regression assumptions (such as Normal errors, symmetry, constant variance), but perhaps the transformed variables satisfy the assumptions. (iv) If some  $Y$ 's are enormously bigger than others, it may not make sense to compare them directly. If  $Y$  is the number of people who work at a restaurant business, the  $Y$  for McDonald's is very, very big, so much so that it can't be compared to the  $Y$  for Genji's (4002 Spruce & 1720 Samson). But you could compare  $\log(Y)$ .

- **Family of transformations:** Organizes search for a good transformation. Family is  $\frac{Y^p - 1}{p}$  which tends to  $\log(y)$  as  $p$  gets near 0.

Often we drop the shift of 1 and the scaling of  $1/p$ , using just  $\text{sign}(p) \cdot Y^p$  for  $p \neq 0$  and  $\log(y)$  for  $p=0$ . Important members of this family are: (i)

$p=1$  for no transformation or  $Y$ , (ii)  $p=1/2$  for  $\sqrt{Y}$ , (iii)  $p=1/3$  for  $\sqrt[3]{Y}$ , (iv)  $p=0$  for  $\log(y)$ , (v)  $p = -1$  for  $1/Y$ .

- **Straightening a scatterplot:** Plot  $Y$  vs  $X$ . If the plot looks curved, then do the following. Divide the data into thirds based on  $X$ , low, middle, high. In each third, find median  $Y$  and median  $X$ . Gives you three  $(X,Y)$  points. Transform  $Y$  and/or  $X$  by adjusting  $p$  until the slope between low and middle equals the slope between middle and high. Then plot the transformed data and see if it looks ok. You want it to look straight, with constant variance around a line.
- **Logit:**  $\text{logit}(a) = \log\{a/(1-a)\}$  when  $a$  is between 0 and 1. If the data are between 0 and 1, their logits are unconstrained.
- **Picking Curves that Make Sense:** Sometimes we let the data tell us which curve to fit because we have no idea where to start. Other times, we approach the data with a clear idea what we are looking for. Sometimes we know what a sensible curve should look like. Some principles – (i) If the residuals show a fan pattern, with greater instability for larger  $Y$ 's, then a log transformation may shift things to constant variance. (ii) If there is a naïve model based on a (too) simple theory (e.g., weight is proportional to volume), then consider models which include the naïve theory as a very special case. (iii) If outcomes  $Y$  must satisfy certain constraints (e.g., percents must be between 0% and 100%), consider families of models that respect those constraints.
- **Interpretable transformations:** Some transformations have simple interpretations, so they are easy to think and write about. Base 2 logs, i.e.,  $\log_2(y)$  can be interpreted in terms of doublings. Reciprocals,  $1/Y$ , are often interpretable if  $Y$  is a ratio (like density) or a time. Squares and



square roots often suggest a relationship between area and length or diameter. Cubes and cube roots suggest a relationship between volume and diameter.

- **Transformations to constant variance:** A very old idea, which still turns up in things you read now and then. Idea is that certain transformations – often strange ones like the arcsin of the square root – make the variance nearly constant, and that is an assumption of regression.

## Topic: Polynomials

**Why fit polynomials?** The transformations we talked about all keep the order of Y intact – big Y’s have big transformed Y’s. Often that is just what we want. Sometimes, however, we see a curve that goes down and comes back up, like a  $\cup$ , or goes up and comes back down, like a  $\cap$ , and the transformations we looked at don’t help at all. Polynomials can fit curves like this, and many other wiggles. They’re also good if you want to find the X that maximizes Y, the top point of the curve  $\cap$ .

- **Quadratic:**  $y = a + bx + cx^2$  has a  $\cup$  shape if  $c > 0$  and a  $\cap$  shape if  $c < 0$  (why?) and is a line if  $c = 0$ . Top of hill or bottom of valley is at  $x = \frac{-b}{2c}$ .
- **Fitting a Quadratic:** Easy – put two variables in the model, namely X and  $X^2$ .
- **Centering:** If  $X > 0$ , then X is big at the same time  $X^2$ , so these two variables are highly correlated. Often a good idea to center, using X and  $(X - \bar{X})^2$  instead of X and  $X^2$ . Fits the same curve, but is more stable as a computing algorithm.
- **Orthogonal polynomials:** Typically used in anova rather than in regression. Transforms  $X^2$  so it is uncorrelated with X. Does this by regressing  $X^2$  on X and using the residuals in place of  $X^2$ .
- **Cubics:** Can fit cubics using X,  $X^2$  and  $X^3$ . Usually don’t go beyond cubics. Usually center.
- **Polynomials in several predictors:** If I have two predictors, say x and w, the quadratic in x and w has squared terms,  $x^2$  and  $w^2$ , but it adds something new, their crossproduct or interaction, xw:

$$y = a + b \cdot x + c \cdot w + d \cdot x^2 + f \cdot w^2 + h \cdot w \cdot x$$

- **Are quadratic terms needed?** You can judge whether you need several quadratic terms using a general linear hypothesis and its avova table.

## Topic: Coded Variables (i.e., Dummy Variables)

- **Why use coded variables?** Coded or dummy variables let you incorporate nominal data (Philly vs New York vs LA) as predictors in regression.
- **Two categories:** If there are just two categories, say male and female, you include a single coded variable, say  $C=1$  for female and  $C=0$  for male. Fits a parallel line model. If you add interactions with a continuous variable,  $X$ , then you are fitting a two-line model, no longer a parallel line model.
- **More than Two Categories:** If there are 3 categories (Philly vs New York vs LA) then you need two coded variables to describe them ( $C=1, D=0$  for New York;  $C=0, D=1$  for LA;  $C=0, D=0$  for Philly). Such a model compares each group to the group left out, the group without its own variable (here, Philly). When there are more than two categories – hence more than one coded variable – interesting hypotheses often involve several variables and are tested with the general linear hypothesis. Does it matter which group you leave out? Yes and no. Had you left out NY rather than Philly, you get the same fitted values, the same residuals, the same overall F-test, etc. However, since a particular coefficient multiplies a particular variable, changing the definition of a variable changes the value of the coefficient.

---

## Topic: Diagnostics -- Residuals

- **Why do we need better residuals?:** We look at residuals to see if the model fits ok – a key concern for any model. But the residuals we have been looking at are not great. The problem is that least squares works very hard to fit data points with extreme X's – unusual predictors – so it makes the residuals small in those cases. A data point with unusual X's is called a high leverage point, and we will think about them in detail a little later. A single outlier (weird Y) at a high leverage point can pull the whole regression towards itself, so this point looks well fitted and the rest of the data looks poorly fitted. We need ways of finding outliers like this. We want regression to tell us what is typical for most points – we don't want one point to run the whole show.

- The model is:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + \varepsilon_i$$

where the  $\varepsilon_i$  are independent  $N(0, \sigma^2)$

- By least squares, we estimate the model to be:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \dots + \hat{\beta}_k X_{ki} \quad \text{with residuals } E_i = Y_i - \hat{Y}_i$$

- Although the true errors,  $\varepsilon_i$  have constant variance,  $\text{var}(\varepsilon_i) = \sigma^2$ , the same for every unit i, the residuals have different variances,  $\text{var}(E_i) = \sigma^2(1 - h_i)$  where  $h_i$  is called the leverage.
- The standardized residual we like best does two things: (i) it uses the leverages  $h_i$  to give each residual the right variance for that one residual, and (ii) it removes observation i when estimating the variance of the residual  $E_i = Y_i - \hat{Y}_i$  for observation i. That is,  $\text{var}(E_i) = \sigma^2(1 - h_i)$  is

estimated by  $\hat{\sigma}_{[-i]}^2(1 - h_i)$ , where  $\hat{\sigma}_{[-i]}^2$  is the estimate of the residual variance we get by setting observation  $i$  aside and fitting the regression without it.

- **The residual we like has several names** – studentized, deleted, jackknife – no one of which is used by everybody. It is the residual divided by its estimated standard error:

$$r_{[i]} = \frac{E_i}{\hat{\sigma}_{[-i]} \sqrt{1 - h_i}}$$

- Another way to get this residual is to create a coded variable that is 1 for observation  $i$  and 0 for all other observations. Add this variable to your regression. The t-statistic for its coefficient equals  $r_{[i]}$ .
- We can test for outliers as follows. The null hypothesis says there are no outliers. If there are  $n$  observations, there are  $n$  deleted residuals  $r_{[i]}$ . Find the largest one in absolute value. To test for outliers at level 0.05, compute  $0.025/n$ , and reject the hypothesis of no outliers if the largest absolute deleted residual is beyond the  $0.025/n$  percentage point of the t-distribution with one less degree of freedom than the error line in the anova table for the regression. (You lose one degree of freedom for the extra coded variable mentioned in the last paragraph.)

---

## Topic: Diagnostics -- Leverage

- **Three very distinct concepts:** An outlier is an observation that is poorly fitted by the regression – it has a response Y that is not where the other data points suggest its Y should be. A high leverage point has predictors, X's, which are unusual, so at these X's, least squares relies very heavily on this one point to decide where the regression plane should go – a high leverage point has X's that allow it to move the regression if it wants to. A high influence point is one that did move the regression – typically, such a point has fairly high leverage (weird X's) and is fairly poorly fitted (weird Y for these X's); however, it may not be the one point with the weirdest X or the one point with the weirdest Y. People often mix these ideas up without realizing it. Talk about a weird Y is outlier talk; talk about a weird X is leverage talk; talk about a weird Y for these X's is influence talk. We now will measure leverage and later influence.
- **Measuring Leverage:** Leverage is measured using the leverages  $h_i$  we encountered when we looked at the variance of the residuals. The leverages are always between 0 and 1, and higher values signify more pull on the regression.
- **When is leverage large?** If a model has k predictors and a constant term, using n observations, then the average leverage, averaging over the n observations is always  $\frac{k+1}{n} = \frac{1}{n} \sum_{i=1}^n h_i$ . A rule a thumb that works well is that leverage is large if it is at least twice the average,  $h_i \geq \frac{2(k+1)}{n}$ .
- **What do you do if the leverage is large?** You look closely. You think. Hard. You find the one or two or three points with  $h_i \geq \frac{2(k+1)}{n}$  and you

look closely at their data. What is it about their X's that made the leverage large? How, specifically, are they unusual? Is there a mistake in the data? If not, do the X's for these points make sense? Do these points belong in the same regression with the other points? Or should they be described separately? Regression gives high leverage points a great deal of weight. Sometimes that makes sense, sometimes not. If you were looking at big objects in our solar system, and  $X$ =mass of object, you would find the sun is a high leverage point. After thinking about it, you might reasonably decide that the regression should describe the planets and the sun should be described separately as something unique. With the solar system, you knew this before you looked at the data. Sometimes, you use regression in a context where such a high leverage point is a discovery. (If you remove a part of your data from the analysis, you must tell people you did this, and you must tell them why you did it.)

- **Interpretation of leverage:** Leverage values  $h_i$  have several interpretations. You can think of them as the distance between the predictors  $X$  for observation  $i$  and the mean predictor. You can think of them as the weight that observation  $i$  gets in forming the predicted value  $\hat{Y}_i$  for observation  $i$ . You can think of leverages as the fraction of the variance of  $Y_i$  that is variance of  $\hat{Y}_i$ . We will discuss these interpretations in class.



---

**Topic: Diagnostics -- Influence**

- **What is influence?** A measure of influence asks whether observation  $i$  *did* move the regression. Would the regression change a great deal if this one observation were removed? Not whether it *could* move the regression – that’s leverage. Not whether it fits poorly – that’s an outlier.
- **Measures of influence.** There are several measures of influence. They are all about the same, but no one has become the unique standard. Two common choices are DFFITS and Cook’s Distance. Cook’s distance is (almost) a constant times the square of DFFITS, so it makes little difference which one you use. It is easier to say what DFFITS does.
- **What is DFFITS?** Roughly speaking, DFFITS measures the change in the predicted value for observation  $i$  when observation  $i$  is removed from the regression. Let  $\hat{Y}_i$  be the predicted value for observation  $i$  using all the data, and let  $\hat{Y}_{i[-i]}$  be the predicted value for observation  $i$  if we fit the regression without this one observation. Is  $\hat{Y}_i$  close to  $\hat{Y}_{i[-i]}$ ? If yes, then this observation does not have much influence. If no, then it does have influence. DFFITS divides the difference,  $\hat{Y}_i - \hat{Y}_{i[-i]}$  by an estimate of the standard error of  $\hat{Y}_i$ , so a value of 1 means a movement of one standard error. Recall that  $\hat{\sigma}_{[-i]}^2$  is the estimated residual variance when observation  $i$  is removed from the regression. Then DFFITS is:

$$DFFITS_i = \frac{\hat{Y}_i - \hat{Y}_{i[-i]}}{\hat{\sigma}_{[-i]} \sqrt{h_i}}.$$

- **DFBETAS:** A related quantity is DFBETAS which looks at the standardized change in the regression coefficient  $\hat{\beta}_j$  when observation  $i$  is removed. There is one DFBETAS for each observation and for each coefficient. DFFITS is always bigger than the largest DFBETAS, and there is only one DFFITS per observation, so many people look at DFFITS instead of all  $k$  DFBETAS.

---

## Topic: Variable Selection

- **What is variable selection?** You have a regression model with many predictor variables. The model looks ok – you’ve done the diagnostic checking and things look fine. But there are too many predictor variables. You wonder if you might do just as well with fewer variables. Deciding which variables to keep and which to get rid of is variable selection.
- **Bad methods.** There are two bad methods you should not use. One bad method is to drop all the variables with small t-statistics. The problem is the t-statistic asks whether to drop a variable *providing you keep all the others*. The t-statistic tells you little about whether you can drop two variables at the same time. It might be you could drop either one but not both, and t can’t tell you this. Another bad method uses the squared multiple correlation,  $R^2$ . The problem is  $R^2$  always goes up when you add variables, and the size of the increase in  $R^2$  is not a great guide about what to do. Fortunately, there is something better.
- **A good method.** The good method uses a quantity called  $C_p$  which is a redesigned  $R^2$  built for variable selection. Suppose the model is, as before,

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + \varepsilon_i$$

where the  $\varepsilon_i$  are independent  $N(0, \sigma^2)$ , but now  $k$  is large (many predictors) and we think some  $\beta$ 's might be zero. We fit this model and get the usual estimate  $\hat{\sigma}^2$  of  $\sigma^2$ . A submodel has some of the  $k$  variables but not all of them, and we name the submodel by the set  $P$  of variables it contains. So the name of the model  $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_3 X_{3i} + \beta_9 X_{9i} + \varepsilon_i$  is

$P=\{1,3,9\}$ , and it has residual sum of squares  $SSE_p$  from the residual line in its Anova table, and  $p=3$  variables plus one constant term or 4 parameters. (Note carefully – I let  $p=\#$ variables, but a few people let  $p=\#$ parameters.) We have  $n$  observations. Then the strange looking but simple formula for  $C_p$  is:  $C_p = \frac{SSE_p}{\hat{\sigma}^2} - [n - 2(p + 1)]$ . Then  $C_p$  compares the model with all variable to the model with just the variables in  $P$  and asks whether the extra variables are worth it.

- Using  $C_p$ : The quantity  $C_p$  estimates the standardized total squared error of prediction when using model  $P$  in place of the model with all the variables. We like a model  $P$  with a small  $C_p$ . If a model  $P$  contains all the variables with nonzero coefficients, then  $C_p$  tends on average to estimate  $p+1$ , the number of variables plus 1 for the constant, so a good value of  $C_p$  is not much bigger than  $p+1$ . For instance, if  $C_{\{1,3,9\}} = 8$ , then that is much bigger than  $p+1=3+1=4$ , so the model seems to be missing important variables, but if  $C_{\{1,3,9,11\}} = 5.1$  then that is close to  $p+1=4+1=5$  and smaller than 8, so that model predicts better and might have all important variables.
- **Searching:** If a model has  $k$  variables, then there are  $2^k$  submodels formed by dropping variables, or about a billion models for  $k=30$  variables. There are various strategies for considering these models: forward selection, backward elimination, stepwise, all subsets, best subsets.
- **Cautions:** Variable selection is an exploratory method, one that looks for interesting things, but because it searches so extensively, it may find some things that don't replicate. If we reject hypotheses when  $P$ -

value  $< 0.05$ , then only 1 time in 20 do we reject a true hypothesis. But if we fit billions for regressions, calculating billions of P-values, then we reject many true hypotheses and make many mistakes. The results of variable selection need to be examined with caution avoiding overstatement. A good strategy is to split the sample, perform variable selection on one half, and confirm the results on the other. This is a simple type of cross-validation.

---

## Topic: One Way Analysis of Variance

- **What is ANOVA?** Anova, or analysis of variance, is the decomposition of data into parts that add back up to the original data, and the summary of the parts in terms of their sizes measured by summing and squaring their numerical entries. At an abstract level, in statistical theory, anova and regression are not really different. In practice, however, they look very different. Most computer programs have separate routines for regression and anova. Center questions, issues and methods arise in anova that don't arise in regression. Anova tends to be used with structured data sets, often from carefully designed experiments, while regression is often used with data that arises naturally. However, by running enough regressions, knowing exactly what you are doing, and putting together the pieces very carefully, you can do even a complex anova using a regression program – it's easy to use an anova program. Anova has a nice geometry.
- **What is one-way anova?** One-way anova is the very simplest case. People fall into one of several groups and we want to understand the difference between the groups. Basic questions are: Do the groups differ? (F-test.) If so, how? (Multiple comparisons.) If I anticipate a specific pattern of differences between the groups, is this pattern confirmed by the data? (Contrasts.) Notation: There are  $I$  groups,  $i=1, \dots, I$ , and  $n_i$  people in group  $i$ , with  $n = n_1 + \dots + n_I$  people in total. Each person is in just one group and people in different groups have nothing to do with each other. Person  $j$  in group  $i$ , has response  $y_{ij}$ . The

mean response in group  $i$  is  $\bar{y}_{i\cdot}$  and the mean response for everyone is  $\bar{y}_{\cdot\cdot}$ . The anova decomposition is:

$$y_{ij} = \bar{y}_{\cdot\cdot} + (\bar{y}_{i\cdot} - \bar{y}_{\cdot\cdot}) + (y_{ij} - \bar{y}_{i\cdot})$$

data = (grand mean) + (group difference) + (residual)

Model for One-Way Anova: The model for one-way anova says the observations are Normal with the same variance, are independent of each other, and have different means in the different groups. Specifically:

$$y_{ij} = \mu + \alpha_i + \varepsilon_{ij} \text{ where the } \varepsilon_{ij} \text{ are iid } N(0, \sigma^2) \text{ and } 0 = \alpha_1 + \alpha_2 + \dots + \alpha_I.$$

- **Do the groups differ?** We test the hypothesis that the groups do not differ using the F-ratio from the one-way analysis of variance table. If F is large enough, judged by the F-table, we reject the null hypotheses and conclude there is strong evidence the groups differ. Otherwise, we conclude that we lack strong evidence that the groups differ (which is not the same thing as saying we know for certain they are the same).
- **If the groups differ, how do they differ?** It is not enough to know the groups differ – we need to understand what differences are present. There are two cases: (1) we have no idea what we are looking for, or (2) we have a clear and specific idea what we are looking for. Case 2 is the better case – if you know what you are looking for in statistics, then you can find it with a smaller sample size. We handle case 1 using multiple comparisons and case 2 using contrasts, described later. In multiple comparison, every group is compared to every other group. If there are  $I=10$  groups, there are 45 comparisons of two groups, 1 with 2, 1 with 3, ..., 9 with 10. If you did 45 t-tests to compare the groups, rejected for P-values  $< 0.05$ , then you would falsely reject a true hypothesis of no

difference in one out of 20 tests. This means that with  $I=10$  groups and 45 comparisons, if the groups are really the same, you expect to get  $45 \times 0.05 = 2.25$  significant (P-value  $< 0.05$ ) difference by chance alone. That's a problem – a big problem. It means you expect 2 mistakes – you expect to say a treatment worked when it didn't. Gotta do something to prevent this. There are many, many solutions to this problem – there are whole books of multiple comparison procedures. One of the first is due to John Tukey of Princeton. You can think of it as using essentially a t-statistic to compare groups in pairs, but as the number of groups,  $I$ , gets bigger, so more tests are being done, the procedure requires a bigger value of the t-statistic before declaring the difference significant. If you did this with  $I=10$  groups and 45 comparisons of two groups, the procedure promises that if the groups are really the same, the chance of a significant difference anywhere in the 45 comparisons is less than 0.05 – if you find anything, you can believe it. For example, with just two groups and 30 degrees of freedom, we would reject at the 0.05 level if  $|t| > 2.04$ , but using Tukey's method with  $I=10$  groups, 45 comparisons, we would reject if  $|t| > 3.41$ , and if  $t$  is that big, then it is unlikely to happen by chance even if you did 45 comparisons.

- **Planned Contrasts for Specific Hypotheses:** Ideally, when you do research, you have a clear idea what you are looking for and why. When this is true, you can build a test tailored to your specific hypothesis. You do this with contrasts among group means. You express your hypothesis using a set of contrast weights you pick, one weight for each group mean, summing to zero:  $c_1, c_2, \dots, c_I$  with  $c_1 + c_2 + \dots + c_I = 0$ . For instance, consider a study with  $I=3$  groups, with  $n_1 = n_2 = n_3 = 100$  people in each



group. The groups are two different drug treatments, A and B, and a placebo control. Then the contrast “drug vs placebo” is:

|                           |               |               |
|---------------------------|---------------|---------------|
| Contrast: Drug vs Placebo |               |               |
| Placebo                   | Drug A        | Drug B        |
| $c_1$                     | $c_2$         | $c_3$         |
| -1                        | $\frac{1}{2}$ | $\frac{1}{2}$ |

whereas the contrast “drug A vs drug B” is:

|                            |        |        |
|----------------------------|--------|--------|
| Contrast: Drug A vs Drug B |        |        |
| Placebo                    | Drug A | Drug B |
| $d_1$                      | $d_2$  | $d_3$  |
| 0                          | 1      | -1     |

- The value of contrast applies the contrast weights to the group means,

$$L = \sum_{i=1}^I c_i \cdot \bar{y}_{i\cdot}, \text{ so for “Drug vs Placebo” it is } L = -1 \cdot \bar{y}_{1\cdot} + \frac{1}{2} \cdot \bar{y}_{2\cdot} + \frac{1}{2} \cdot \bar{y}_{3\cdot}.$$

- The t-test for a contrast tests the null hypothesis  $H_0: 0 = \sum_{i=1}^I c_i \cdot \alpha_i$ . Let

$\hat{\sigma}^2$  be the residual mean square from the anova table, which estimates

$\sigma^2$ . The t-statistic is  $t = \frac{L}{\sqrt{\hat{\sigma}^2 \cdot \sum \frac{c_i^2}{n_i}}}$  and the degrees of freedom are

from the residual line in the anova table.

- The sum of squares for a contrast is  $\frac{L^2}{\sum \frac{c_i^2}{n_i}}$ . Two contrasts,  $c_1, c_2, \dots, c_I$

and  $d_1, d_2, \dots, d_I$  are orthogonal if  $0 = \sum_{i=1}^I \frac{c_i \cdot d_i}{n_i}$ . Example: “Drug vs

Placebo” is orthogonal to “Drug A vs Drug B” because

$$\sum_{i=1}^I \frac{c_i \cdot d_i}{n_i} = \frac{-1 \times 0}{100} + \frac{\frac{1}{2} \times 1}{100} + \frac{\frac{1}{2} \times -1}{100} = 0. \text{ When contrasts are orthogonal, the}$$

sum of squares between groups may be partitioned into separate parts, one for each contrast. If there are  $I$  groups, then there are  $I-1$  degrees of freedom between groups, and each degree of freedom can have its own contrast. Both of these formulas are mostly used in balanced designs where the sample sizes in the groups are the same,  $n_1 = n_2 = \dots = n_I$ .

---

## Topic: Two Way Analysis of Variance

- What is two-way ANOVA? In two-way anova, each measurement is classified into groups in two different ways, as in the rows and columns of a table. In the social sciences, the most common situation is to measure the same unit or person under several different treatments – this is the very simplest case of what is know as repeated measurements. Each person is a row, each treatment is a column, and each person gives a response under each treatment. The two-way's are person and treatment. Some people give higher responses than others. Some treatments are better than others. The anova measures both sources of variation. The units might be businesses or schools or prisons instead of people.

Notation: There are I people,  $i=1,\dots,I$ , and J treatments,  $j=1,\dots,J$ , and person i gives response  $y_{ij}$  under treatment j. The mean for person i is

$$\bar{y}_{i\cdot} = \frac{1}{J} \sum_{j=1}^J y_{ij} \text{ and the mean for treatment j is } \bar{y}_{\cdot j} = \frac{1}{I} \sum_{i=1}^I y_{ij}, \text{ and the}$$

mean of everyone is  $\bar{y}_{\cdot\cdot}$ . The anova decomposition is:

$$y_{ij} = \bar{y}_{\cdot\cdot} + (\bar{y}_{i\cdot} - \bar{y}_{\cdot\cdot}) + (\bar{y}_{\cdot j} - \bar{y}_{\cdot\cdot}) + (y_{ij} - \bar{y}_{i\cdot} - \bar{y}_{\cdot j} + \bar{y}_{\cdot\cdot}).$$

- **Anova table:** The anova table now has “between rows”, “between columns” and “residual”, so the variation in the data is partitioned more finely.
- **Normal model:** The model is  $y_{ij} = \mu + \alpha_i + \beta_j + \varepsilon_{ij}$  where the errors are independent Normals with mean zero and variance  $\sigma^2$ . Under this model, F-statistics from the anova table may be used to test the hypotheses of no difference between rows and no difference between columns. Can do multiple comparisons and contrasts using the residual line from the anova table to obtain the estimate  $\hat{\sigma}^2$ .

**Error Rates When Performing More Than One Hypothesis Test**

Table Counts Null Hypotheses

|                       | Accepted or untested null hypotheses | Rejected null hypotheses | Total     |
|-----------------------|--------------------------------------|--------------------------|-----------|
| True null hypotheses  | U                                    | V                        | $m_0$     |
| False null hypotheses | T                                    | S                        | $m - m_0$ |
| Total                 | $m - R = m - (U + T)$                | R                        | m         |

Family-wise error rate:  $\Pr(V \geq 1)$ , the probability of at least one false rejection in  $m$  tests.

The family-wise error rate is weakly controlled at  $\alpha = 0.05$  if  $\alpha = 0.05 \geq \Pr(V \geq 1)$  whenever  $m = m_0$ , that is, whenever all  $m$  null hypotheses are true.

The family-wise error rate is strongly controlled at  $\alpha = 0.05$  if  $\alpha = 0.05 \geq \Pr(V \geq 1)$  for all values of  $m_0$ , that is, no matter how many null hypotheses are true.

Weak control is not enough. Weak control means you are unlikely to find something when there is nothing, but you are still likely to find too much when there is something.

False discovery rate (FDR) is the expected number of false rejections,  $E(V/R)$  where  $E/R = 0/0$  is defined to be 0 (i.e., no false rejections if no rejections). This is a more lenient standard than the family-wise error rate, rejecting more true hypotheses.

If you do  $m$  tests at level  $\alpha = 0.05$ , you expect to falsely reject  $0.05 \times m_0$  hypotheses, and if all hypotheses are true, this might be as high as  $0.05 \times m$ . The expected ratio of false rejections to tests,  $E(V/m)$ , is called the per comparison error rate.

Example  
Table Counts Null Hypotheses

|                       | Accepted or untested null hypotheses | Rejected null hypotheses | Total |
|-----------------------|--------------------------------------|--------------------------|-------|
| True null hypotheses  | U                                    | V                        | 100   |
| False null hypotheses | T                                    | S                        | 1     |
| Total                 | $m-R=m-(U+T)$                        | R                        | 101   |

In this example, there are 101 hypotheses and 100 are true.

If you test each hypothesis at level  $\alpha=0.05$ , you expect  $0.05 \times 100 = 20$  false rejections of true null hypotheses, plus if you are lucky a rejection of the one false null hypothesis, so you expect most rejections to be false rejections.

If you strongly control the family-wise error rate at  $\alpha=0.05$ , then the chance of at least one false rejection is at most 5%.

If you weakly control the family-wise error rate at  $\alpha=0.05$ , then there are no promises about false rejections in this case, as one null hypothesis is false.

---

**What are adjusted P-values?** (e.g. as produced by `pairwise.t.test()`)

A test of null hypothesis  $H_0$  either rejects  $H_0$  or it does not.

The level,  $\alpha$ , of the test is such that  $\alpha \geq \Pr(\text{Reject } H_0)$  when  $H_0$  is true.

The P-value is the smallest  $\alpha$  such that we reject  $H_0$ .

This definition of a P-value continues to work with multiple hypothesis testing.

---

## Topic: Factorial Analysis of Variance

- **Two Factor Factorial Anova:** The simplest case of factorial anova involves just two factors – similar principles apply with more than two factors, but things get large quickly. Suppose you have two drugs, A and B – then “drug” is the first factor, and it has two levels, namely A and B. Suppose each drug has two dose levels, low and high – then “dose” is the second factor, and it too has two levels, low and high. A person gets one combination, perhaps drug B at low dose. Maybe I give 50 people each drug at each level, so I have 200 people total, 100 on A, 100 on B, 100 at low dose, 100 at high dose.

**Main effects and interactions:** We are familiar with main effects – we saw them in two-way anova. Perhaps drug A is better than drug B – that’s a main effect. Perhaps high dose is more effective than low dose – that’s a main effect. But suppose instead that drug A is better than drug B at high dose, but drug A is inferior to drug B at low dose – that’s an interaction. In an interaction, the effect of one factor changes with the level of the other.

- **Anova table:** The anova table has an extra row beyond that in two-way anova, namely a row for interaction. Again, it is possible to do contrasts and multiple comparisons.
- **More Complex Anova:** Anova goes on and on. The idea is to pull apart the variation in the data into meaningful parts, each part having its own row in the anova table. There may be many factors, many groupings, etc.



## Some Aspects of R

*Script is my commentary to you.* **Bold Courier is what I type in R.** Regular Courier is what R answered.

*What is R?*

*R is a close relative of Splus, but R is available for free. You can download R from <http://cran.r-project.org/>. R is very powerful and is a favorite (if not the favorite) of statisticians; however, it is not easiest package to use. It is command driven, not menu driven, so you have to remember things or look them up - that's the only thing that makes it hard. You can add things to R that R doesn't yet know how to do by writing a little program. R gives you fine control over graphics. Most people need a book to help them, and so Mainland & Braun's *Data Analysis and Graphics Using R*, Cambridge University Press. Another book is Dalgaard's *Introductory Statistics with R*, NY: Springer. Dalgaard's book is better at teaching basic statistics, and it is good if you need a review of basic statistics to go with an introduction to R. R is similar to Splus, and there are many good books about Splus. One is: Venables and Ripley *Modern Applied Statistics with S-Plus* (NY: Springer-Verlag).*

*Who should use R?*

*If computers terrify you, if they cause insomnia, cold sweats, and anxiety attacks, perhaps you should stay away from R. On the other hand, if you want a very powerful package for free, one you won't outgrow, then R worth a try. If you find you need lots of help to install R or make R work, then R isn't for you. Alternatives for Statistics 500 are JMP-IN, SPSS, Systat, Stata, SAS and many others. For Statistics 501, beyond the basics, R is clearly best.*

*You need to download R the first time from the webpage above.*

*You need to get the "Rst500" workspace for the course from*

*<http://www-stat.wharton.upenn.edu/>*

*going to "Course downloads" and the most recent Fall semester, Statistics 500, or in one step to*

*<http://download.wharton.upenn.edu/download/pub/stat/Fall-2006/STAT-500/>*

*For Statistics 501,*

*<http://stat.wharton.upenn.edu/statweb/course/Spring-2007/stat501/>*

---

Start R.

From the File Menu, select "Load Workspace".

Select "Rst500"

To see what is in a workspace, type

```
ls()
```

or type

```
objects()
```

```
> ls()
```

```
[1] "fuel"
```

To display an object, type its name

```
> fuel
```

|    | ID | state | Fuel | Tax  | License | Inc   | Road  |
|----|----|-------|------|------|---------|-------|-------|
| 1  | 1  | ME    | 541  | 9.00 | 52.5    | 3.571 | 1.976 |
| 2  | 2  | NH    | 524  | 9.00 | 57.2    | 4.092 | 1.250 |
| 3  | 3  | VT    | 561  | 9.00 | 58.0    | 3.865 | 1.586 |
|    |    |       |      |      |         |       |       |
|    |    |       |      |      |         |       |       |
|    |    |       |      |      |         |       |       |
|    |    |       |      |      |         |       |       |
|    |    |       |      |      |         |       |       |
|    |    |       |      |      |         |       |       |
|    |    |       |      |      |         |       |       |
| 46 | 46 | WN    | 510  | 9.00 | 57.1    | 4.476 | 3.942 |
| 47 | 47 | OR    | 610  | 7.00 | 62.3    | 4.296 | 4.083 |
| 48 | 48 | CA    | 524  | 7.00 | 59.3    | 5.002 | 9.794 |

Fuel is a data frame.

```
> is.data.frame(fuel)
```

```
[1] TRUE
```

You can refer to a variable in a data frame as `fuel$Tax`, etc. It returns one column of fuel.

```
> fuel$Tax
```

```
[1] 9.00 9.00 9.00 7.50 8.00 10.00 8.00 8.00 8.00 7.00 8.00 7.50
[13] 7.00 7.00 7.00 7.00 7.00 7.00 7.00 8.50 7.00 8.00 9.00 9.00
[25] 8.50 9.00 8.00 7.50 8.00 9.00 7.00 7.00 8.00 7.50 8.00 6.58
[37] 5.00 7.00 8.50 7.00 7.00 7.00 7.00 7.00 6.00 9.00 7.00 7.00
```

`length()` and `dim()` tell you how big things are. There are 48 states and seven variables.

```
> length(fuel$Tax)
```

```
[1] 48
```

```
> dim(fuel)
```

```
[1] 48 7
```

*To get a summary of a variable, type `summary(variable)`*

```
> summary(fuel$Tax)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 5.000  7.000   7.500   7.668  8.125  10.000
```

*R has very good graphics. You can make a boxplot with*

```
boxplot(fuel$Fuel)
```

*or dress it up with*

```
boxplot(fuel$Fuel,ylab="gallons per person",main="Figure 1:
Motor Fuel Consumption")
```

*To learn about a command, type `help(command)`*

```
help(boxplot)
```

```
help(plot)
```

```
help(t.test)
```

```
help(lm)
```

### *Optional Trick*

*It can get tiresome typing `fuel$Tax`, `fuel$Licenses`, etc. If you type `attach(data.frame)` then you don't have to mention the data frame. Type `detach(data.frame)` when you are done.*

```
> summary(fuel$Tax)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 5.000  7.000   7.500   7.668  8.125  10.000
> summary(Tax)
Error in summary(Tax) : Object "Tax" not found
> attach(fuel)
> summary(Tax)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 5.000  7.000   7.500   7.668  8.125  10.000
> summary(License)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 45.10  52.98   56.45   57.03  59.52   72.40
> detach(fuel)
```

## HELP

*R contains several kinds of help. Use `help(keyword)` to get documentation about keyword.*

```
> help(boxplot)
```

*Use `help("key")` to find the keywords that contain "key". The quotes are needed.*

```
> apropos("box")
```

```
[1] "box"           "boxplot"       "boxplot.default"  
     "boxplot.stats"
```

*Use `help.search("keyword")` to search the web for R functions that you can download related to keyword. Quotes are needed.*

```
> help.search("box")
```

```
> help.search("fullmatch")
```

*At <http://cran.r-project.org/> there is free documentation, some of which is useful, but perhaps not for first-time users. To begin, books are better.*

## Some R

*A variable, "change" in a data.frame bloodpressure.*

```
> bloodpressure$change
[1] -9 -4 -21 -3 -20 -31 -17 -26 -26 -10 -23 -33 -19 -19 -23
```

*It doesn't know what "change" is.*

```
> change
Error: Object "change" not found
```

*Try attaching the data.frame*

```
> attach(bloodpressure)
```

*Now it knows what "change" is.*

```
> change
[1] -9 -4 -21 -3 -20 -31 -17 -26 -26 -10 -23 -33 -19 -19 -23
```

```
> mean(change)
```

```
[1] -18.93333
```

```
> sd(change)
```

```
[1] 9.027471
```

```
> summary(change)
```

```
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
-33.00 -24.50  -20.00  -18.93  -13.50   -3.00
```

```
> stem(change)
```

The decimal point is 1 digit(s) to the right of the |

```
-3 | 31
-2 | 663310
-1 | 9970
-0 | 943
```

```
> hist(change)
```

```
> boxplot(change)
```

```
> boxplot(change,main="Change in Blood Pressure After
Captopril",ylab="Change mmHg")
```

```
> boxplot(change,main="Change in Blood Pressure After
Captopril",ylab="Change mmHg",ylim=c(-40,40))
```

```
> abline(0,0,lty=2)
```

```
> t.test(change)
```

One Sample t-test

```
data: change
```

```
t = -8.1228, df = 14, p-value = 1.146e-06
```

```
alternative hypothesis: true mean is not equal to 0
```

```
95 percent confidence interval:
```

```
-23.93258 -13.93409
```

```
sample estimates:
```

```
mean of x
```

```
-18.93333
```

## Are the Data Normal?

```
> attach(bloodpressure)
> change
[1] -9 -4 -21 -3 -20 -31 -17 -26 -26 -10 -23 -33 -19 -19 -23
> par(mfrow=c(1,2))
> boxplot(change)
> qqnorm(change)
```

*A straight line in a Normal quantile plot is consistent with a Normal distribution.*

*You can also do a Shapiro-Wilk test. A small p-value suggests the data are not Normal.*

```
> shapiro.test(change)

Shapiro-Wilk normality test

data:  change
W = 0.9472, p-value = 0.4821
```

*The steps below show what the qqnorm() function is plotting*

```
> round(ppoints(change), 3)
[1] 0.033 0.100 0.167 0.233 0.300 0.367 0.433 0.500 0.567 0.633
[11] 0.700 0.767 0.833 0.900 0.967
```

*The plotting positions in the normal plot:*

```
> round(qnorm(ppoints(change)), 3)
[1] -1.834 -1.282 -0.967 -0.728 -0.524 -0.341 -0.168 0.000 0.168
[10] 0.341 0.524 0.728 0.967 1.282 1.834
```

*qqnorm(change) is short for*

```
> plot(qnorm(ppoints(change)), sort(change))
```

*Here are Normal quantile plots of several Normal and non-Normal distributions.*

*Can you tell from the plot which are Normal?*

```
> qqnorm(rnorm(10))
> qqnorm(rnorm(100))
> qqnorm(rnorm(1000))
> qqnorm(rcauchy(100))
> qqnorm(rlogis(100))
> qqnorm(rexp(100))
```

## Regression in R

*Script is my commentary to you. Bold Courier is what I type in R. Regular Courier is what R answered.*

```
> ls()
[1] "fuel"
```

*To display an object, type its name*

```
> fuel
  ID state Fuel  Tax License  Inc  Road
1  1  ME  541  9.00   52.5 3.571 1.976
2  2  NH  524  9.00   57.2 4.092 1.250
3  3  VT  561  9.00   58.0 3.865 1.586
      .
      .
      .
46 46  WN  510  9.00   57.1 4.476 3.942
47 47  OR  610  7.00   62.3 4.296 4.083
48 48  CA  524  7.00   59.3 5.002 9.794
```

*To do regression, use lm. lm stands for linear model.*

*To fit  $Fuel = \alpha + \beta Tax + \epsilon$  type*

```
> lm(Fuel~Tax)
```

```
Call:
lm(formula = Fuel ~ Tax)
```

```
Coefficients:
(Intercept)          Tax
   984.01         -53.11
```

*To fit  $Fuel = \alpha + \beta_1 Tax + \beta_2 License + \epsilon$  type*

```
> lm(Fuel~Tax+License)
```

```
Call:
lm(formula = Fuel ~ Tax + License)
```

```
Coefficients:
(Intercept)          Tax          License
   108.97         -32.08           12.51
```

To see more output, type  
**> summary(lm(Fuel~Tax))**

Call:  
 lm(formula = Fuel ~ Tax)

Residuals:

|  | Min      | 1Q      | Median | 3Q     | Max     |
|--|----------|---------|--------|--------|---------|
|  | -215.157 | -72.269 | 6.744  | 41.284 | 355.736 |

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t ) |     |
|-------------|----------|------------|---------|----------|-----|
| (Intercept) | 984.01   | 119.62     | 8.226   | 1.38e-10 | *** |
| Tax         | -53.11   | 15.48      | -3.430  | 0.00128  | **  |

---  
 Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 100.9 on 46 degrees of freedom  
 Multiple R-Squared: 0.2037, Adjusted R-squared: 0.1863  
 F-statistic: 11.76 on 1 and 46 DF, p-value: 0.001285

*You can save the regression in an object and then refer to it:*

**> reg1<-lm(Fuel~Tax+License)**

*Now the workspace has a new object, namely reg1:*

**> ls()**  
 [1] "fuel" "reg1"

*To see reg1, type its name:*

**> reg1**

Call:  
 lm(formula = Fuel ~ Tax + License)

Coefficients:

|             | Tax    | License |
|-------------|--------|---------|
| (Intercept) | 108.97 | 12.51   |
|             | -32.08 |         |

*To get residuals, type*

**> reg1\$residuals**

*This works only because I defined reg1 above. To boxplot residuals, type:*

**> boxplot(reg1\$residuals)**



*To plot residuals against predicted values, type*

```
> plot(reg1$fitted.values,reg1$residuals)
```

*To do a normal plot of residuals, type*

```
> qqnorm(reg1$residuals)
```

*To get deleted or jackknife residuals, type*

```
> rstudent(reg1)
```

*To get leverages or hats, type*

```
> hatvalues(reg1)
```

*To get dffits*

```
> dffits(reg1)
```

*To get Cook's distance*

```
> cooks.distance(reg1)
```

*Clean up after yourself. To remove reg1, type rm(reg1)*

```
> ls()
```

```
[1] "fuel" "reg1"
```

```
> rm(reg1)
```

```
> ls()
```

```
[1] "fuel"
```

## Predictions

*Fit a linear model and save it.*

```
> mod<-lm(Fuel~Tax)
```

*A confidence interval for the line at Tax = 8.5*

```
> predict(mod,data.frame(Tax=8.5),interval="confidence")
      fit      lwr      upr
[1,] 532.6041 493.4677 571.7405
```

*A prediction interval for a new observation at Tax = 8.5*

```
> predict(mod,data.frame(Tax=8.5),interval="prediction")
      fit      lwr      upr
[1,] 532.6041 325.7185 739.4897
```

*Same point estimate, 532.6 gallons, but a very different interval, because the prediction interval has to allow for a new error for the new observation.*

## Multiple Regression Anova in R

*The standard summary output from a linear model in R contains the key elements of the anova table, which are underlined.*

```
> summary(lm(Fuel~Tax+License))
Call:
lm(formula = Fuel ~ Tax + License)

Residuals:
    Min       1Q   Median       3Q      Max
-123.177  -60.172   -2.908   45.032  242.558

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  108.971     171.786   0.634   0.5291
Tax          -32.075     12.197  -2.630   0.0117 *
License       12.515      2.091   5.986 3.27e-07 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 76.13 on 45 degrees of freedom
Multiple R-Squared: 0.5567, Adjusted R-squared: 0.537
F-statistic: 28.25 on 2 and 45 DF, p-value: 1.125e-08
```

*More explicitly, the model `lm(Fuel~1)` fits just the constant term, and the F test compares that model (with just the constant term) to the model with all the variables (here Tax & License).*

```
> anova(lm(Fuel~1),lm(Fuel~Tax+License))
Analysis of Variance Table

Model 1: Fuel ~ 1
Model 2: Fuel ~ Tax + License
  Res.Df    RSS Df Sum of Sq    F    Pr(>F)
  1      47 588366
  2      45 260834  2    327532 28.253 1.125e-08 ***
```

*Most regression programs present an explicit anova table, similar to that above, rather than just the F-test.*

## Partial Correlation Example

Here are the first two lines of data from a simulated data set. We are interested in the relationship between  $y$  and  $x_2$ , taking account of  $x_1$ .

```
> partialcorEG[1:2,]
      y          x1          x2
1 -3.8185777 -0.8356356 -1.0121903
2  0.3219982  0.1491024  0.0853746
```

Plot the data. Always plot the data.

```
> pairs(partialcorEG)
```

Notice that  $y$  and  $x_2$  have a positive correlation.

```
> cor(partialcorEG)
      y          x1          x2
y  1.0000000  0.9899676  0.9535053
x1  0.9899676  1.0000000  0.9725382
x2  0.9535053  0.9725382  1.0000000
```

The partial correlation is the correlation between the residuals. Notice that  $y$  and  $x_2$  have a negative partial correlation adjusting for  $x_1$ .

```
> cor(lm(y~x1)$residual,lm(x2~x1)$residual)
[1] -0.2820687
```

Notice that the multiple regression coefficient has the same sign as the partial correlation.

```
> summary(lm(y~x1+x2))
Call:
lm(formula = y ~ x1 + x2)
Residuals:
      Min       1Q   Median       3Q      Max
-1.13326 -0.27423 -0.02018  0.32216  1.07808
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.007177   0.048662   0.147  0.88305
x1           4.768486   0.243833  19.556 < 2e-16 ***
x2          -0.720948   0.248978  -2.896  0.00468 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 0.4866 on 97 degrees of freedom
Multiple R-Squared: 0.9816, Adjusted R-squared: 0.9812
F-statistic: 2591 on 2 and 97 DF, p-value: < 2.2e-16
```

### Added Variable Plots

You have fit a model, `mod1` say,  $Y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \epsilon$  where  $\epsilon$  are iid  $N(0, \sigma^2)$  and now want to ask about adding a new variable,  $x_{k+1}$ , to this model.

It can't hurt to plot  $Y$  against  $x_{k+1}$ . However, that plot does not tell you what  $x_{k+1}$  will do in the model above. It could happen that  $Y$  increases with  $x_{k+1}$  but  $\beta_{k+1} < 0$ .

The added variable plot uses the idea of regression by stages. In regression by stages, you estimate  $\beta_{k+1}$  by regressing the residuals from `mod1` on the residuals of  $x_{k+1}$  when regressed on  $x_1, \dots, x_k$ . The added variable plot is simply the plot of these two sets of residuals, residuals of  $Y$  versus residuals of  $x_{k+1}$ . The slope in that plot estimates  $\beta_{k+1}$ . So the added variable plot lets you see what happens when  $x_{k+1}$  is added to `mod1`.

You can calculate the two set of residuals and plot them. That works fine. Or you can use `addedvarplot` in the course workspace.

```
> attach(fuel)
> head(fuel)
  ID state Fuel Tax License Inc Road
1  1    ME  541  9.0    52.5 3.571 1.976
2  2    NH  524  9.0    57.2 4.092 1.250
3  3    VT  561  9.0    58.0 3.865 1.586
4  4    MA  414  7.5    52.9 4.870 2.351
5  5    RI  410  8.0    54.4 4.399 0.431
6  6    CN  457 10.0    57.1 5.342 1.333
> mod1<-lm(Fuel~Tax+License)
> addedvarplot(mod1,Inc)
The same plot is produced directly by:
> plot(lm(Inc~Tax+License)$resid,mod1$resid)
```

**ADDED VARIABLE PLOTS IN THE car Package**

```
> attach(fuel)
> pairs(cbind(Fuel,Tax,License))
> library(car)
> help(avPlots)
> avPlots(lm(Fuel~Tax+License))
> avPlots(lm(Fuel~Tax+License+Inc),term=~Inc)
> avPlots(lm(Fuel~Tax+License+Inc))
> summary(lm(Fuel~Tax+License+Inc))
```

car stands for "Companion to Applied Regression". The book, *An R Companion to Applied Regression* by John Fox and Sanford Weisberg discusses regression using this package.

Vocabulary Homework

```
> vocabulary
  Age Vocab
1  0.67     0
2  0.83     1
3  1.00     3
4  1.25    19
5  1.50    22
6  1.75   118
7  2.00   272
8  2.50   446
9  3.00   896
10 3.50  1222
11 4.00  1540
12 4.50  1870
13 5.00  2072
14 5.50  2289
15 6.00  2562
```

```
> attach(vocabulary)
```

*Fit linear model (a line) and store results in "mod".*

```
> mod<-lm(Vocab~Age)
```

*Summary output for mod.*

```
> summary(mod)
```

Call:

```
lm(formula = Vocab ~ Age)
```

Residuals:

| Min     | 1Q      | Median | 3Q    | Max    |
|---------|---------|--------|-------|--------|
| -249.67 | -104.98 | 13.14  | 78.47 | 268.25 |

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t )     |
|-------------|----------|------------|---------|--------------|
| (Intercept) | -621.16  | 74.04      | -8.389  | 1.32e-06 *** |
| Age         | 526.73   | 22.12      | 23.808  | 4.17e-12 *** |

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 148 on 13 degrees of freedom  
 Multiple R-Squared: 0.9776, Adjusted R-squared: 0.9759  
 F-statistic: 566.8 on 1 and 13 DF, p-value: 4.170e-12

*Plot the data. Does a line look appropriate?*

```
> plot(Age,Vocab,ylim=c(-1000,3000))  
> abline(mod)
```

*Plot residuals vs predicted. Is there a pattern?*

```
> plot(mod$fitted.values,mod$residuals)
```

*Boxplot residuals. Unusual points? Skewness?*

```
> boxplot(mod$residuals)
```

*Normal plot of residuals. Do the residuals look Normal? (Is it a line?)*

```
> qqnorm(mod$residuals)
```

*Test of the null hypothesis that the residuals are Normal.*

```
> shapiro.test(mod$residuals)
```

Shapiro-Wilk normality test

```
data: mod$residuals
```

```
W = 0.9801, p-value = 0.9703
```



### General Linear Hypothesis

```
> help(anova.lm)
```

```
> attach(fuel)
```

```
> fuel[1:2,]
```

|   | ID | state | Fuel | Tax | License | Inc   | Road  |
|---|----|-------|------|-----|---------|-------|-------|
| 1 | 1  | ME    | 541  | 9   | 52.5    | 3.571 | 1.976 |
| 2 | 2  | NH    | 524  | 9   | 57.2    | 4.092 | 1.250 |

*Fit the full model.*

```
> mod<-lm(Fuel~Tax+License+Inc)
```

```
> anova(mod) Optional step - for your education only.
```

Analysis of Variance Table

Response: Fuel

|           | Df | Sum Sq | Mean Sq | F value | Pr(>F)        |
|-----------|----|--------|---------|---------|---------------|
| Tax       | 1  | 119823 | 119823  | 27.560  | 4.209e-06 *** |
| License   | 1  | 207709 | 207709  | 47.774  | 1.539e-08 *** |
| Inc       | 1  | 69532  | 69532   | 15.992  | 0.0002397 *** |
| Residuals | 44 | 191302 | 4348    |         |               |

---

*Fit the reduced model.*

```
> mod2<-lm(Fuel~Tax)
```

```
> anova(mod2) Optional step - for your education only.
```

Analysis of Variance Table

Response: Fuel

|           | Df | Sum Sq | Mean Sq | F value | Pr(>F)      |
|-----------|----|--------|---------|---------|-------------|
| Tax       | 1  | 119823 | 119823  | 11.764  | 0.001285 ** |
| Residuals | 46 | 468543 | 10186   |         |             |

*Compare the models*

```
> anova(mod2,mod)
```

Analysis of Variance Table

Model 1: Fuel ~ Tax

Model 2: Fuel ~ Tax + License + Inc

|   | Res.Df | RSS    | Df | Sum of Sq | F      | Pr(>F)        |
|---|--------|--------|----|-----------|--------|---------------|
| 1 | 46     | 468543 |    |           |        |               |
| 2 | 44     | 191302 | 2  | 277241    | 31.883 | 2.763e-09 *** |

*Notice the residual sum of squares and degrees of freedom in the three anova tables!*

## Polynomial Regression

```
> attach(cars)
```

*Quadratic in size*  $y = \beta_0 + \beta_1 x + \beta_2 x^2$

```
> lm(mpg~size+I(size^2))
```

Call:

```
lm(formula = mpg ~ size + I(size^2))
```

Coefficients:

|             |            |           |
|-------------|------------|-----------|
| (Intercept) | size       | I(size^2) |
| 39.3848313  | -0.1485722 | 0.0002286 |

*Centered quadratic in size*  $y = \beta_0 + \beta_1 x + \beta_2 \{x - \text{mean}(x)\}^2$

```
> lm(mpg~size+I((size-mean(size))^2))
```

Call:

```
lm(formula = mpg ~ size + I((size - mean(size))^2))
```

Coefficients:

|             |      |                          |
|-------------|------|--------------------------|
| (Intercept) | size | I((size - mean(size))^2) |
| 28.8129567  |      | -0.0502460               |
| 0.0002286   |      |                          |

*Orthogonal Polynomial Quadratic in size*

```
> lm(mpg~poly(size,2))
```

Call:

```
lm(formula = mpg ~ poly(size, 2))
```

Coefficients:

|             |                |                |
|-------------|----------------|----------------|
| (Intercept) | poly(size, 2)1 | poly(size, 2)2 |
| 20.74       | -24.67         | 12.33          |

*To gain understanding:*

- do all three regressions
- look at t-test for  $\beta_2$
- type `poly(size,2)`
- plot `poly(size,2)[,1]` and `poly(size,2)[,2]` against size

### Centered Polynomial with Interaction

```
> fuel[1:2,]
  ID state Fuel Tax License  Inc  Road
1  1    ME  541   9    52.5 3.571 1.976
2  2    NH  524   9    57.2 4.092 1.250

> attach(fuel)
Construct the squared and crossproduct terms. Alternatives: use "*" or ":" in model formula.
> TaxC<-Tax-mean(Tax)
> LicC<-License-mean(License)
> TaxC2<-TaxC*TaxC
> LicC2<-LicC*LicC
> TaxLicC<-TaxC*LicC

> modfull<-lm(Fuel~Tax+License+TaxC2+LicC2+TaxLicC)
> summary(modfull)
```

Call:  
 lm(formula = Fuel ~ Tax + License + TaxC2 + LicC2 + TaxLicC)  
 Residuals:

|  | Min        | 1Q        | Median   | 3Q       | Max       |
|--|------------|-----------|----------|----------|-----------|
|  | -121.52425 | -51.08809 | -0.01205 | 46.27051 | 223.28655 |

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t )     |
|-------------|----------|------------|---------|--------------|
| (Intercept) | 169.7242 | 179.6332   | 0.945   | 0.3501       |
| Tax         | -32.4465 | 12.2906    | -2.640  | 0.0116 *     |
| License     | 11.2776  | 2.3087     | 4.885   | 1.55e-05 *** |
| TaxC2       | 1.3171   | 8.6638     | 0.152   | 0.8799       |
| LicC2       | 0.2575   | 0.2868     | 0.898   | 0.3743       |
| TaxLicC     | -2.5096  | 2.7343     | -0.918  | 0.3640       |

---  
 Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1  
 Residual standard error: 76.42 on 42 degrees of freedom  
 Multiple R-Squared: 0.5831, Adjusted R-squared: 0.5335  
 F-statistic: 11.75 on 5 and 42 DF, p-value: 3.865e-07

*Test whether the three squared and interaction terms are needed:*

```
> modred<-lm(Fuel~Tax+License)
> anova(modred,modfull)
Analysis of Variance Table

Model 1: Fuel ~ Tax + License
Model 2: Fuel ~ Tax + License + TaxC2 + LicC2 + TaxLicC
  Res.Df  RSS Df Sum of Sq    F Pr(>F)
1     45 260834
2     42 245279  3     15555 0.8879 0.4552
```

## Understanding Linear Models with Interactions or Polynomials

*NIDA data (DC\*MADS) on birth weight of babies in DC and attributes of mom.*

```
> DCBabyCig[1:2,]
  Age Married CIGS  BW
1  17      0    0 2385
2  23      1    0 4175
```

*Age x Cigarettes interaction*

```
> AC<-Age*CIGS
```

*Model with interaction*

```
> lm(BW~Age+CIGS+AC)
```

Call:

```
lm(formula = BW ~ Age + CIGS + AC)
```

Coefficients:

```
(Intercept)      Age      CIGS      AC
  2714.81      13.99     562.66    -28.04
```

*How do you understand a model with interactions?*

*Let's create a new data.frame with 6 moms in it. Three moms are 18, three are 35. Some smoke 0, 1 or 2 packs.*

```
> new[,1]<-c(18,35,18,35,18,35)
> new[,2]<-c(0,0,1,1,2,2)
> new[,3]<-new[,1]*new[,2]
> colnames(new)<-c("Age", "CIGS", "AC")
> new<-data.frame(new)
```

```
> new
  Age CIGS AC
1  18    0  0
2  35    1 35
3  18    0  0
4  35    1 35
5  18    0  0
6  35    1 35
```

*Now, for these six moms, let's predict birth weight of junior. It is usually easier to talk about people than about coefficients, and that is what this table does: it talks about 6 moms.*

```
> round(cbind(new,predict(lm(BW~Age+CIGS+AC),new,interval="confidence"))
  Age CIGS AC  fit  lwr  upr
1  18    0  0 2967 2865 3068
2  35    0  0 3204 3073 3336
3  18    1 18 3024 2719 3330
4  35    1 35 2786 2558 3013
5  18    2 36 3082 2474 3691
6  35    2 70 2367 1919 2814
```

**Interpretation of an Interaction**

```
> DCBabyCig[1:6,]
```

|   | Age | Married | CIGS | BW   |
|---|-----|---------|------|------|
| 1 | 17  | 0       | 0    | 2385 |
| 2 | 23  | 1       | 0    | 4175 |
| 3 | 25  | 0       | 0    | 3655 |
| 4 | 18  | 0       | 0    | 1855 |
| 5 | 20  | 0       | 0    | 3600 |
| 6 | 24  | 0       | 0    | 2820 |

Age = mother's age

Married, 1=yes, 0=no

CIGS = packs per day, 0, 1, 2.

BW = birth weight in grams

```
> dim(DCBabyCig)
```

```
[1] 449 4
```

```
> mod<-lm(BW~Age+Married+CIGS+I(Married*CIGS))
```

```
> summary(mod)
```

Coefficients:

|                   | Estimate  | Std. Error | t value | Pr(> t ) |     |
|-------------------|-----------|------------|---------|----------|-----|
| (Intercept)       | 2973.1866 | 152.5467   | 19.490  | < 2e-16  | *** |
| Age               | 0.1699    | 6.5387     | 0.026   | 0.97928  |     |
| Married           | 274.0662  | 89.2913    | 3.069   | 0.00228  | **  |
| CIGS              | -88.4957  | 81.7163    | -1.083  | 0.27941  |     |
| I(Married * CIGS) | -415.1501 | 160.4540   | -2.587  | 0.00999  | **  |

Residual standard error: 687.8 on 444 degrees of freedom

Multiple R-squared: 0.05337, Adjusted R-squared: 0.04484

F-statistic: 6.258 on 4 and 444 DF, p-value: 6.618e-05

Plot the data

```
> boxplot(BW~Married:CIGS)
```

A 25 year old mom in all combinations of Married and CIGS.

```
> DCBabyCigInter
```

|   | Age | Married | CIGS |
|---|-----|---------|------|
| 1 | 25  | 0       | 0    |
| 2 | 25  | 0       | 1    |
| 3 | 25  | 0       | 2    |
| 4 | 25  | 1       | 0    |
| 5 | 25  | 1       | 1    |
| 6 | 25  | 1       | 2    |

Predicted birth weights for this mom, with confidence intervals.

```
> predict(mod,DCBabyCigInter,interval="conf")
      fit      lwr      upr
1 2977.434 2890.180 3064.688
2 2888.938 2738.900 3038.977
3 2800.443 2502.124 3098.761
4 3251.500 3114.163 3388.838
5 2747.854 2476.364 3019.345
6 2244.209 1719.423 2768.995
```

Let's clean it up, converting to pounds (2.2 pounds per kilogram), and add the predictors:

```
> pr<-predict(mod,DCBabyCigInter,interval="conf")

> round(cbind(DCBabyCigInter,2.2*pr/1000),1)
  Age Married CIGS fit lwr upr
1  25      0     0 6.6 6.4 6.7
2  25      0     1 6.4 6.0 6.7
3  25      0     2 6.2 5.5 6.8
4  25      1     0 7.2 6.9 7.5
5  25      1     1 6.0 5.4 6.6
6  25      1     2 4.9 3.8 6.1
```

### Using Restricted Cubic Splines (aka Natural Splines)

```

> library(Hmisc)
> head(cars)
      car  size  mpg  group
1 ToyotaC  71.1 33.9     1
2  HondaC  75.7 30.4     1

> x<-rcspline.eval(size,nk=3) #Three knots, one additional
variable
> plot(size,x) #What does the new variable look like?
> plot(size,mpg)
> m<-lm(mpg~size+x) #Add the new variable to the model.
> points(size,m$fit,pch=16,col="red") #What does the fit
look like?
> summary(m)
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 38.81737    1.91493  20.271 < 2e-16 ***
size        -0.11667    0.01419  -8.225 7.86e-09 ***
x           0.14618    0.02642   5.533 7.29e-06 ***
Residual standard error: 2.346 on 27 degrees of freedom
Multiple R-squared:  0.8395,    Adjusted R-squared:  0.8276
F-statistic: 70.62 on 2 and 27 DF,  p-value: 1.877e-11

> x<-rcspline.eval(size,nk=5) #Five knots, three additional
variables
> m<-lm(mpg~size+x)
> summary(m)
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 46.28095    3.61684  12.796 1.79e-12 ***
size        -0.19689    0.03655  -5.387 1.37e-05 ***
x1          2.17972    0.97629   2.233  0.0348 *
x2         -3.52907    1.75277  -2.013  0.0550 .
x3          1.45720    0.90924   1.603  0.1216
Residual standard error: 2.199 on 25 degrees of freedom
Multiple R-squared:  0.8694,    Adjusted R-squared:  0.8485
F-statistic: 41.61 on 4 and 25 DF,  p-value: 1.054e-10

> points(size,m$fit,pch=16,col="purple")

```

Reference: Harrell, F. (2015) *Regression Modeling Strategies*, New York: Springer, section 2.4.5.

Comment: There are many types of splines. Natural splines are linear beyond the final knots, so they wiggle less at the ends.

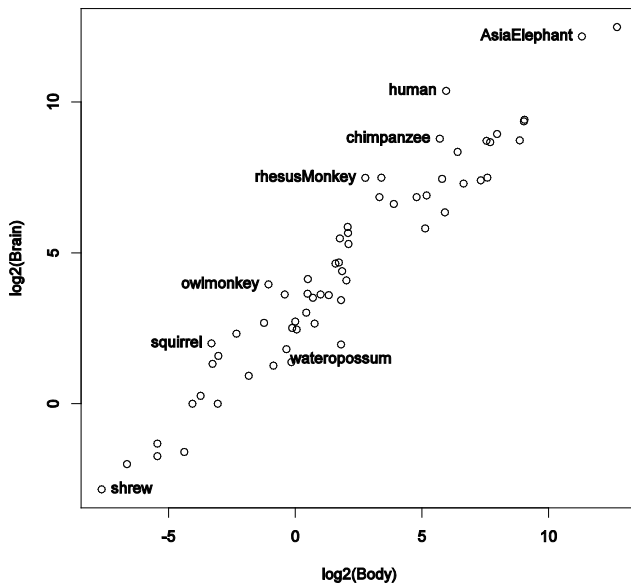
### Dummy Variable in Brains Data

First two rows of "brains" data.

```
> brains[1:2,]
  Body Brain Animal Primate Human
1 3.385 44.500 articfox      0     0
2 0.480 15.499 owlmonkey  1     0

> attach(brains)

> plot(log2(Body), log2(Brain))
> identify(log2(Body), log2(Brain), labels=Animal)
```



```
> mod<-lm(log2(Brain)~log2(Body)+Primate)
> mod
```

Call:  
lm(formula = log2(Brain) ~ log2(Body) + Primate)

Coefficients:  
(Intercept)    log2(Body)       Primate  
      2.8394        0.7402        1.6280

$$\log_2(\text{Brain}) \sim \log_2(\text{Body}) + \text{Primate}$$

$$\text{is } 2^{\text{Brain}} = 2^{(\alpha + \beta \log_2(\text{Body}) + \gamma \text{Primate} + \epsilon)} = (2^\alpha) (\text{Body}^\beta) (2^{\gamma \text{Primate}}) (2^\epsilon)$$

$$2^{1.628 \text{Primate}} = 3.1 \text{ for a primate, } = 1 \text{ for a nonprimate}$$



### Computing the Diagnostics in the Rat Data

```

> ratdata[1:3,]
  BodyWgt LiverWgt Dose Percent Rat3
1    176     6.5 0.88  0.42    0
2    176     9.5 0.88  0.25    0
3    190     9.0 1.00  0.56    1

> attach(ratdata)
> mod<-lm(Percent~BodyWgt+LiverWgt+Dose)

Standardized residuals (first 5)
> rstandard(mod)[1:5]
      1          2          3          4          5
1.766047 -1.273040  0.807154 -1.377232 -1.123099

Deleted or jackknife or "studentized" residuals (first 5)
> rstudent(mod)[1:5]
      1          2          3          4          5
1.9170719 -1.3022313  0.7972915 -1.4235804 -1.1337306

dffits (first 5)
> dffits(mod)[1:5]
      1          2          3          4          5
0.8920451 -0.6087606  1.9047699 -0.4943610 -0.9094531

Cook's distance (first 5)
> cooks.distance(mod)[1:5]
      1          2          3          4          5
0.16882682 0.08854024 0.92961596 0.05718456 0.20291617

Leverages or 'hats' (first 5)
> hatvalues(mod)[1:5]
      1          2          3          4          5
0.1779827 0.1793410 0.8509146 0.1076158 0.3915382

> dfbeta(mod)[1:3,]
  (Intercept)   BodyWgt   LiverWgt   Dose
1 -0.006874698  0.0023134055 -0.011171761 -0.3419002
2  0.027118946 -0.0007619302 -0.008108905  0.1869729
3 -0.045505614 -0.0134632770  0.005308722  2.6932347

> dfbetas(mod)[1:3,]
  (Intercept)   BodyWgt   LiverWgt   Dose
1 -0.03835128  0.31491627 -0.7043633 -0.2437488
2  0.14256373 -0.09773917 -0.4817784  0.1256122
3 -0.23100202 -1.66770314  0.3045718  1.7471972

```

### High Leverage Example

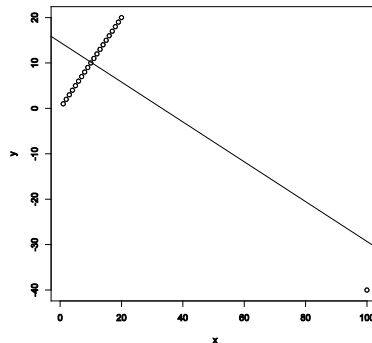
```
> t(highlev)  t() is transpose - make rows into columns and columns into rows - compact printing
  1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21
x 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 100
y 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 -40
```

```
> mod<-lm(y~x)
> summary(mod)
Residuals:
    Min       1Q   Median       3Q      Max
-13.1343  -7.3790  -0.1849   7.0092  14.2034

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 14.57312     2.41264   6.040 8.24e-06 ***
x           -0.43882     0.09746  -4.503 0.000244 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Residual standard error: 8.875 on 19 degrees of freedom  
 Multiple R-Squared: 0.5162, Adjusted R-squared: 0.4908  
 F-statistic: 20.27 on 1 and 19 DF, p-value: 0.0002437

```
> plot(x,y)
> abline(mod)  Puts the fitted line on the plot — What a dumb model!
```



*The bad guy, #21, doesn't have the biggest residual!*

```
> mod$residual[21]
      21
-10.69070
```

```
Residuals:
    Min       1Q   Median       3Q      Max
-13.1343  -7.3790  -0.1849   7.0092  14.2034
```

*But our diagnostics find him!*

```
> rstudent(mod)[21]
      21
-125137800
> hatvalues(mod)[21]
      21
0.9236378
> dffits(mod)[21]
      21
-435211362
```

## Outlier Testing

Use the Bonferroni inequality with the deleted/jackknife/"studentized" residuals.

Example uses random data - should not contain true outliers

```
> x<-rnorm(1000)
> y<-rnorm(1000)
> plot(x,y)
> summary(lm(y~x))
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) -0.005469   0.031685  -0.173   0.863
x            -0.044202   0.031877  -1.387   0.166
```

```
Residual standard error: 1.002 on 998 degrees of freedom
Multiple R-Squared: 0.001923, Adjusted R-squared: 0.0009228
F-statistic: 1.923 on 1 and 998 DF, p-value: 0.1659
```

Look at the deleted residuals (from rstudent)

```
> summary(rstudent(lm(y~x)))
      Min.      1st Qu.      Median      Mean      3rd Qu.      Max.
-3.189e+00 -6.596e-01  2.186e-02 -7.077e-06  6.517e-01  3.457e+00
```

The big residual in absolute value is 3.457. Is that big for the biggest of 1000 residuals?

The `pt(value,df)` command looks up value in the *t*-table with *df* degrees of freedom, returning  $Pr(t < \text{value})$ . You need the other tail,  $Pr(t > \text{value})$ , and you need to double it for a 2-tailed test. The degrees of freedom are *one less than the degrees of freedom in the error for the regression*, here 997.

```
> 2*(1-pt(3.457,997))
[1] 0.0005692793
```

This is uncorrected *p*-value. Multiply by the number of tests, here 1000, to correct for multiple testing. (It's an inequality, so it can give a value bigger than 1.)

```
> 1000* 2*(1-pt(3.457,997))
[1] 0.5692793
```

As this is bigger than 0.05, the null hypothesis of no outliers is not rejected - it is plausible there are no outliers present.



## Testing Whether a Transformation of Y is Needed

### Tukey's One Degree of Freedom for Nonadditivity

Tukey (1949) proposed testing whether a transformation of  $y$  is needed in the model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon \quad \varepsilon \sim \text{iid } N(0, \sigma^2)$$

by adding a scaled centered version of  $\hat{y}^2$  to the model, specifically  $\frac{(\hat{y} - \bar{y})^2}{2\bar{y}}$ ; see

Atkinson (1985, p. 157). The function `tukey1df(mod)` in the class workspace does this, but you could easily do it yourself.

```
> mod<-lm(BW~Age+Married+CIGS)
> summary(mod)
```

Call:

```
lm(formula = BW ~ Age + Married + CIGS)
```

Residuals:

| Min      | 1Q      | Median | 3Q     | Max     |
|----------|---------|--------|--------|---------|
| -2408.30 | -358.49 | 99.69  | 453.34 | 1952.79 |

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t )    |
|-------------|----------|------------|---------|-------------|
| (Intercept) | 2936.618 | 152.859    | 19.211  | < 2e-16 *** |
| Age         | 2.557    | 6.515      | 0.392   | 0.69488     |
| Married     | 200.615  | 85.198     | 2.355   | 0.01897 *   |
| CIGS        | -196.644 | 70.665     | -2.783  | 0.00562 **  |

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 692.2 on 445 degrees of freedom  
 Multiple R-squared: 0.0391, Adjusted R-squared: 0.03262  
 F-statistic: 6.036 on 3 and 445 DF, p-value: 0.0004910

```
> boxplot(mod$resid)
> qqnorm(mod$resid)
> shapiro.test(mod$resid)
```

Shapiro-Wilk normality test

```
data: mod$resid
W = 0.9553, p-value = 1.996e-10
```

```
> plot(mod$fit,mod$resid)
> lines(lowess(mod$fit,mod$resid))
```

To do the test, add the transformed variable to the model and look at its t-statistic.

```
> summary(lm(BW~Age+Married+CIGS+tukey1df(mod)))
```

Call:

```
lm(formula = BW ~ Age + Married + CIGS + tukey1df(mod))
```

Residuals:

| Min     | 1Q     | Median | 3Q    | Max    |
|---------|--------|--------|-------|--------|
| -2301.3 | -334.5 | 107.7  | 420.8 | 1981.2 |

Coefficients:

|               | Estimate | Std. Error | t value | Pr(> t )     |
|---------------|----------|------------|---------|--------------|
| (Intercept)   | 2962.751 | 151.775    | 19.521  | < 2e-16 ***  |
| Age           | 0.508    | 6.494      | 0.078   | 0.93769      |
| Married       | 104.750  | 90.366     | 1.159   | 0.24701      |
| CIGS          | -489.699 | 120.693    | -4.057  | 5.86e-05 *** |
| tukey1df(mod) | 31.507   | 10.567     | 2.982   | 0.00302 **   |

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 686.2 on 444 degrees of freedom  
 Multiple R-squared: 0.05796, Adjusted R-squared: 0.04947  
 F-statistic: 6.83 on 4 and 444 DF, p-value: 2.428e-05

### Box – Cox Method

An alternative approach is due to Box and Cox (1964).

```
> library(MASS)
> help(boxcox)
> boxcox(mod)
```

Andrews, D. F. (1971) A note on the selection of data transformations. *Biometrika*, 58, 249-254.

Atkinson, A. C. (1985) *Plots, Transformations and Regression*. NY: Oxford.

Box, G. E. P. and Cox, D. R. (1964) An analysis of transformations (with discussion). *Journal of the Royal Statistical Society B*, 26, 211–252.

Tukey, J. W. (1949) One degree of freedom for nonadditivity. *Biometrics*, 5, 232-242.

Venables, W. N. and Ripley, B. D. (2002) *Modern Applied Statistics with S*. Fourth edition. Springer.

## Tukey Tests in the `car` Package in R

- The function `residualPlots()` in the `car` package automates some checks for nonlinearity using the general form of Tukey's test for nonadditivity.
- If `m` is a linear model, then type `residualPlots(m)`
- For example, using the cars data:

```
> library(car)
> residualPlots(lm(mpg~size))
      Test stat Pr(>|t|)
size          4.99      0
Tukey test    4.99      0
```

- Because there is only one `x` in the cars data, there is only one test. It asks whether a quadratic in `size` would predict the residuals. The t-statistic is 4.99 with a tiny P-value, so the answer is yes.
- For example, using the fuel data:

```
> attach(fuel)
> residualPlots(lm(Fuel~Tax+License))
      Test stat Pr(>|t|)
Tax          0.202  0.841
License      1.321  0.193
Tukey test   1.633  0.103
```

- Now there are 3 tests. One is about whether  $\hat{y}^2$  predicts the residuals, the "Tukey test". Another is whether  $\text{Tax}^2$  predicts the residuals. Another is whether  $\text{License}^2$  predicts the residuals. In all three cases, the t-statistic is not large and the P-value is large, so it is plausible that no quadratic term is needed.

## Transformations in the car package

```
> attach(cars)
> plot(size,mpg)
> library(car)
```

This is looking for a transformation of size that would improve the fit. Because  $y=mpg$  is not transformed, only  $x=size$ , the procedure can compare residual sums of squares of  $y$  (RSS) for different transformations of  $x$ . Here,  $\lambda$  is the power transformation, what we called  $p$  in class.

```
> invTranPlot(mpg~size)
      lambda      RSS
1 -1.262328 123.4974
2 -1.000000 126.3604
3  0.000000 192.1809
4  1.000000 317.0362
```

It likes the  $-1.26$  power of  $x$  as a transformation, but  $-1$  and  $-1.5$  are in the confidence interval.

```
> invTranEstimate(size,mpg)
```

```
$lambda
[1] -1.262328
$lowerCI
[1] -1.737292
$upperCI
[1] -0.8237128
```

This is trying to transform  $y=mpg$ , not  $x=size$ . It uses the Box-Cox likelihood method. It likes  $y^{(-0.99)}$  or approximately  $1/y$ .

```
> summary(powerTransform(lm(mpg~size)))
```

```
bcPower Transformation to Normality
      Est.Power Std.Err. Wald Lower Bound Wald Upper Bound
Y1   -0.9878   0.5553      -2.0763           0.1007
Likelihood ratio tests about transformation parameters
              LRT df          pval
LR test, lambda = (0)  3.194518  1 0.0738855550
LR test, lambda = (1) 12.478121  1 0.0004117462
```

A graphical version.

```
> boxCox(lm(mpg~size))
```



The car package has an alternative to Tukey's test for a transformation.

Atkinson' method

```
> v<-boxCoxVariable(mpg)
> summary(lm(mpg~size+v))
```

Call:

```
lm(formula = mpg ~ size + v)
```

Residuals:

|  | Min     | 1Q      | Median | 3Q     | Max    |
|--|---------|---------|--------|--------|--------|
|  | -4.3645 | -1.3910 | 0.0462 | 1.0028 | 6.4195 |

Coefficients:

|             | Estimate  | Std. Error | t value | Pr(> t ) |     |
|-------------|-----------|------------|---------|----------|-----|
| (Intercept) | 76.651791 | 9.132471   | 8.393   | 5.27e-09 | *** |
| size        | -0.033455 | 0.004324   | -7.737  | 2.55e-08 | *** |
| v           | 2.519663  | 0.486271   | 5.182   | 1.87e-05 | *** |

Andrews/Tukey method

```
> summary(lm(mpg~size+tukey1df(md)))
```

Coefficients:

|              | Estimate  | Std. Error | t value | Pr(> t ) |     |
|--------------|-----------|------------|---------|----------|-----|
| (Intercept)  | 28.812957 | 1.008748   | 28.56   | < 2e-16  | *** |
| size         | -0.050246 | 0.004508   | -11.15  | 1.32e-11 | *** |
| tukey1df(md) | 5.577240  | 1.117584   | 4.99    | 3.12e-05 | *** |

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.471 on 27 degrees of freedom  
 Multiple R-squared: 0.8218, Adjusted R-squared: 0.8087  
 F-statistic: 62.28 on 2 and 27 DF, p-value: 7.686e-11

References:

For Tukey's method and related methods  
 Andrews, D. F. (1971) A note on the selection of data transformations. *Biometrika* 58, 249-254.  
 For the method in the car package:  
 Atkinson, A. C. (1973) Testing transformations to Normality. *JRSS-B* 473-479.

## Checking for Non-constant Variance

The usual linear model assumes that the errors are independent and identically distributed (iid) with a Normal distribution having expectation 0 and constant variance  $\sigma^2$ .

Is there evidence in the data that the variance is not constant?

There is a graphical aid and a test.

In the car data, the relationship between mile-per-gallon and engine size was curved.

```
> attach(cars)
> head(cars)
      car  size  mpg group
1 ToyotaC  71.1 33.9     1
2  HondaC  75.7 30.4     1
3  Fiat128  78.7 32.4     1
4  FiatX19  79.0 27.3     1
5 LotusEur  95.1 30.4     1
6  Datsun 108.0 22.8     1
> plot(size,mpg)
```

In the car data, our first attempt to straighten the curve was to raise mpg to the -4.5 power. That fixed one problem, but created another: it was straight, but the plot of residuals versus predicted had a fan shape, suggesting non-constant variance.

```
> tmpg<-mpg^(-4.5)
> plot(size,tmpg)
```

The third plot you get from `plot(model)` is a substitute for looking for a fan shape. It plots the square root of the absolute value of the standardized residuals against the fitted values. A trend, up or down, in this plot suggests non-constant variance.

```
> m<-lm(tmpg~size)
> plot(m)
```

The third plot you get is the plot for non-constant variance. In this case, you see a clear trend, indicating non-constant variance. You could make this plot yourself by typing:

```
> v<-sqrt(abs(rstandard(m)))
> plot(m$fit,v)
```

Why is this better than plotting residuals versus predicted, the first plot you get from `plot(model)`? There are three reasons:

- (i) Even when the variance of the errors is constant, the usual residuals, `m$residual`, do not have constant variance. However, in this case, the standardized residuals, `rstandard(m)`, do have constant variance.
- (ii) By taking absolute values, you make the standardized residuals positive, so instead of looking for a fan, you are looking for a trend, up or down. It is often easier to see the trend.
- (iii) The absolute value of iid Normal data tends to be skewed right. This is distracting. Taking the square root removes much of skew.

In the `car` package, there is a test for non-constant variance. The null hypothesis is that the variance is constant, and a small P-value raises doubt about the null hypothesis, suggesting the variance changes with one or more of the predictors.

```
> library(car)
> m<-lm(tmpg~size)
> ncvTest(m)
Non-constant Variance Score Test
Variance formula: ~ fitted.values
Chisquare = 20.54211    Df = 1    p = 5.833364e-06
```

This test is similar to, but not quite the same as, regressing the squared residuals on the predicted values from the model. The residuals themselves are uncorrelated with the predicted values. Can you predict the squared residuals?

The test is due to Breusch and Pagan (1979) and Cook and Weisberg (1983).

It is easy to confuse an outlier, non-constant variance, and nonlinearity. A test of each might find all three. Graphs help you recognize which one is actually present.

Cook, R. D. and Weisberg, S. (1983) Diagnostics for heteroscedasticity in regression. *Biometrika* 70, 1-10.

## Calculating $C_p$ for the Cathedral Data

```

> attach(cathedral)
> mod<-lm(length~height+gothic+GH)
> summary(mod)
Call:
lm(formula = length ~ height + gothic + GH)
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  241.833     336.471   0.719   0.480
height         3.138         4.506   0.696   0.494
gothic       -204.722     347.207  -0.590   0.562
GH             1.669         4.641   0.360   0.723
Residual standard error: 79.11 on 21 degrees of freedom
Multiple R-Squared:  0.5412,    Adjusted R-squared:  0.4757
F-statistic: 8.257 on 3 and 21 DF,  p-value: 0.0008072

> drop1(lm(length~height+gothic+GH),scale=79.11^2)
Single term deletions

Model:
length ~ height + gothic + GH

scale:  6258.392

      Df Sum of Sq    RSS    Cp
<none>          131413 3.9979
height  1         3035 134448 2.4829
gothic  1         2176 133589 2.3455
GH      1          810 132223 2.1273

> drop1(lm(length~height+gothic),scale=79.11^2)
Single term deletions

Model:
length ~ height + gothic

scale:  6258.392

      Df Sum of Sq    RSS    Cp
<none>          132223 2.1273
height  1     119103 251326 19.1582
gothic  1       37217 169440  6.0740

```

## Variable Selection

Highway data. First two rows of 39 rows. More details in the Variable Selection section of this bulkpack.

```
> highway[1:2,]
  ID rate  len adt trks slim lwid shld  itg sigs acpt lane fai pa ma
1  1 4.58  4.99  69   8  55  12  10 1.20   0  4.6   8  1  0  0
2  2 2.86 16.11  73   8  60  12  10 1.43   0  4.4   4  1  0  0
```

Highway data has 39 rows, 15 columns, of which y=rate, and columns 3 to 15 or 3:15 are predictors. Want to select predictors.

```
> dim(highway)
[1] 39 15
```

```
> attach(highway)
```

To use "leaps" for best subsets regression, need to get it from the library. To get documentation, type help!

```
> library(leaps)
> help(leaps)
```

Easiest if you put the x's in a separate variable. These are columns 3:15 of highway, including all the rows.

```
> x<-highway[,3:15]
```

First three rows of 39 rows of x. Notices that the first two columns of highway are gone.

```
> x[1:3,]
  len adt trks slim lwid shld  itg sigs acpt lane fai pa ma
1  4.99  69   8  55  12  10 1.20   0  4.6   8  1  0  0
2 16.11  73   8  60  12  10 1.43   0  4.4   4  1  0  0
3  9.75  49  10  60  12  10 1.54   0  4.7   4  1  0  0
```

There are 13 predictors, hence  $2^{13} = 8,192$  possible models formed by including each variable or not.

```
> dim(x)
[1] 39 13
> 2^13
[1] 8192
```

Look at the names of your predictors: len = length of segment, ..., slim = speed limit, ..., acpt = number of access points per mile, ...

```
> colnames(x)
[1] "len"  "adt"  "trks" "slim" "lwid" "shld" "itg"  "sigs" "acpt"
"lane" "fai"  "pa"   "ma"
```

A quick and easy, but not very complete, answer is obtained from regsubsets. Here, it gives the best model with 1 variable, the best with 2 variables, etc. Look for the \*'s. The best 3 variable model is len, slim, acpt.

```
> summary(regsubsets(x=x,y=rate))
1 subsets of each size up to 8 Selection Algorithm: exhaustive
  len adt trks slim lwid shld itg sigs acpt lane fai pa ma
1 ( 1 ) " " " " " " " " " " " " " " " " " " " " " " " "
2 ( 1 ) "*" " " " " " " " " " " " " " " " " " " " " " "
3 ( 1 ) "*" " " " " " "*" " " " " " " " " " " " " " " " "
4 ( 1 ) "*" " " " " " "*" " " " " " " "*" "*" " " " " " " "
5 ( 1 ) "*" " " " " " "*" " " " " " " "*" "*" " " " " "*" " "
6 ( 1 ) "*" " " "*" "*" " " " " " " " " "*" "*" " " " " "*" " "
7 ( 1 ) "*" " " "*" "*" " " " " " " " " "*" "*" " " " " "*" "*"
8 ( 1 ) "*" " " "*" "*" " " " " " " " " "*" "*" " " " " "*" "*"

```

To get the two best models of each size, type:

```
> summary(regsubsets(x=x,y=rate,nbest=2))
```

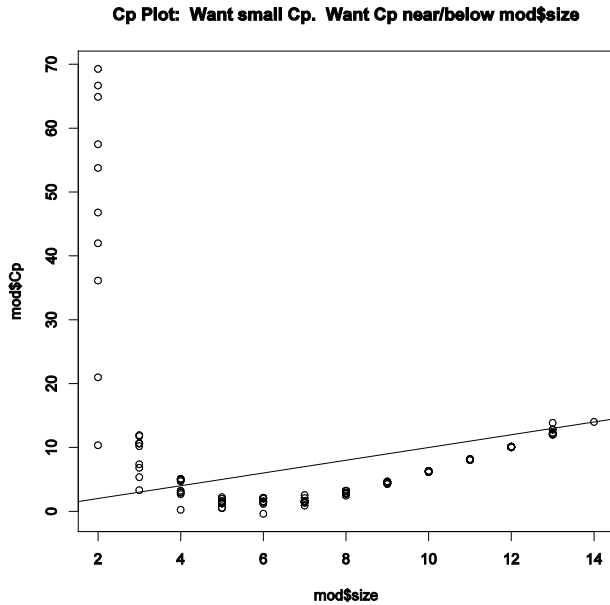




## Variable Selection, Continued

This is the  $C_p$  plot.

```
> plot(mod$size,mod$Cp)
> abline(0,1)
```



There is one pretty good 2 variable model (size=3), with  $C_p$  near the  $x=y$  line, and one very good 3 variable model (size=4), with  $C_p$  way below the line. The best model has 5 variables (size =6) but is only trivially better than the 3 variable model.  $R^2$  is highest for the 14 variable model, but chances are it won't predict as well as the 3 variable model.

Let's put together the pieces.

```
> join<-cbind(mod$which,mod$Cp,mod$size)
```

Let's look at the 3 variable models (size=4). The best has  $C_p = 0.236$  and variables len, slim, and acpt.

```
> join[mod$size==4,]
  len adt trks slim lwid shld itg sigs acpt lane fai pa ma
3  1  0  0  1  0  0  0  0  1  0  0  0  0 0.2356971 4
3  0  0  1  1  0  0  0  0  1  0  0  0  0 2.6805672 4
3  1  0  0  0  0  0  0  1  1  0  0  0  0 2.8975068 4
3  1  0  0  0  0  1  0  0  1  0  0  0  0 3.0404482 4
3  1  0  1  0  0  0  0  0  1  0  0  0  0 3.2366902 4
3  1  0  0  0  1  0  0  0  1  0  0  0  0 4.7193511 4
3  0  0  0  1  0  0  0  1  1  0  0  0  0 4.8847460 4
3  1  0  0  0  0  0  0  0  1  0  0  1  0 4.9933327 4
3  1  0  0  0  0  0  0  0  1  1  0  0  0 5.0489720 4
3  1  1  0  0  0  0  0  0  1  0  0  0  0 5.1013513 4
```

The full model has  $C_p = 14$ .

```
> join[mod$size==14,]
  len adt trks slim lwid shld itg sigs acpt lane fai pa ma
  1  1  1  1  1  1  1  1  1  1  1  1  1  1 14 14
```

$C_p$  thinks that the 14 variable model will have squared errors 59 times greater than the 3 variable model with len, slim, and acpt.

```
> 14/0.2356971
[1] 59.39827
```

A key problem is that variable selection procedures overfit. Need to cross-validate!



## Variable Selection, Continued (O2Uptake Example)

*Load libraries*

```
> library(leaps)
> library(car)

> O2Uptake[1:3,]
  Day Bod TKN  TS  TVS  COD O2UP LogO2Up
1   0 1125 232 7160 85.9 8905 36.0  1.5563
2   7  920 268 8804 86.5 7388  7.9  0.8976
3  15  835 271 8108 85.2 5348  5.6  0.7482

> dim(O2Uptake)
[1] 20  8
```

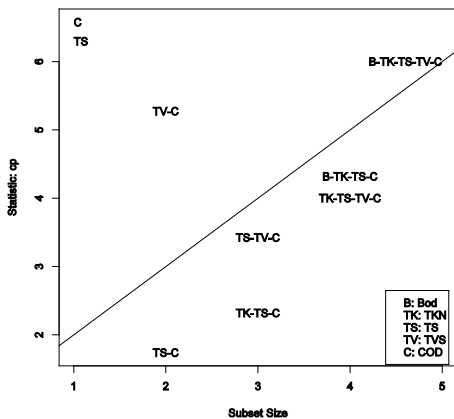
*Find best 2 models of each size .*

```
> mod<-regsubsets(x=O2Uptake[,2:6],y=O2Uptake$LogO2Up,nbest=2)
> summary(mod)
2 subsets of each size up to 5
Selection Algorithm: exhaustive
      Bod TKN TS  TVS COD
1 ( 1 ) " " " " "*" " " " "
1 ( 2 ) " " " " " " " " "*"
2 ( 1 ) " " " " "*" " " " "*"
2 ( 2 ) " " " " " " " "*" "*"
3 ( 1 ) " " "*" "*" " " " "*"
3 ( 2 ) " " " " "*" "*" "*"
4 ( 1 ) " " "*" "*" "*" "*"
4 ( 2 ) "*" "*" "*" " " "*"
5 ( 1 ) "*" "*" "*" "*" "*"

```

*C<sub>p</sub> plot*

```
> subsets(mod,stat="cp")
> abline(1,1)
```



## PRESS (and writing little programs in R)

We have seen many times in many ways that ordinary residuals, say  $E_i$ , tend to be too small, because  $Y_i$  was used in fitting the model, so the model is too close to  $Y_i$ . Predicting  $Y_i$  having fitted the model using  $Y_i$  is called "in-sample prediction," and it tends to suggest that a model is better than it is, because it says you are making progress getting close to your current  $Y_i$ 's, even if you could not do well in predicting a new  $Y_i$ .

If you left  $i$  out of the regression, and tried to predict  $Y_i$  from the regression without  $i$ , the error you would make is:

$$Y_i - \hat{Y}_{i[-i]} = V_i \text{ say.}$$

Here,  $V_i$  is an "out-of-sample prediction," a true effort to predict a "new" observation, because  $i$  did not get used in fitting this equation. It gives me a fair idea as to how well a model can predict an observation not used in fitting the model.

The predicted error sum of squares or PRESS is

$$\text{PRESS} = \sum V_i^2.$$

It turns out that  $V_i = E_i / (1 - h_i)$  where  $E_i$  is the residual and  $h_i$  is the leverage or hatvalue.

```
> fuel[1:2,]
  ID state Fuel Tax License Inc Road
1  1  ME  541   9   52.5 3.571 1.976
2  2  NH  524   9   57.2 4.092 1.250
> attach(fuel)
> modMAX<-lm(Fuel~Tax+License+Inc+Road)
```

These are the out of sample prediction errors or  $V_i$ 's:

```
> v<-modMAX$residual / (1-hatvalues(modMAX))
```

Let's look at Wyoming, WY. It's residual is about 235 gallons:

```
> modMAX$residual[state=="WY"]
 40
234.9472
```

but it's out of sample prediction error is about 26 gallons larger:

```
> v[state=="WY"]
 40
260.9721
```

**PRESS (and writing little programs in R), continued**

*PRESS is the sum of the squares of the  $V_i$ :*

```
> sum(V^2)
[1] 235401.1
```

*How does PRESS compare to  $R^2$ ? Well  $R^2$  is an in-sample measure, while press is an out-of-sample measure. For modMAX,  $R^2$  is:*

```
> summary(modMAX)$r.squared
[1] 0.6786867
```

*Let's take Road out of the model, and see what happens to  $R^2$  and PRESS.*

```
> modSmall<-lm(Fuel~Tax+License+Inc)
```

```
> summary(modSmall)$r.squared
[1] 0.6748583
```

*So  $R^2$  went down, "got worse," which it always does when you delete variables;*

```
> Vsmall<-modSmall$residual/(1-hatvalues(modSmall))
```

```
> sum(Vsmall^2)
[1] 229998.9
```

*however, PRESS went down too, or "got better." In other words, adding Road to the model makes the residuals smaller, as adding variables always does, but it makes the prediction errors bigger. Sometimes adding a variable makes prediction errors smaller, sometimes it makes them bigger, and PRESS tells which is true in your model.*

*You could compute PRESS as above each time you fit a model, but it is easier to add a little program to R. Here is how you write a program called PRESS that computes PRESS.*

```
> PRESS<-function(mod){
+ V<-mod$residual/(1-hatvalues(mod))
+ sum(V^2)}
```

*If you type in the name of your program, here PRESS, it prints the program for you to look at.*

```
> PRESS
function(mod){
  V<-mod$residual/(1-hatvalues(mod))
  sum(V^2)}
```

*Your new program will compute PRESS for you:*

```
> PRESS(modMAX)
[1] 235401.1
> PRESS(modSmall)
[1] 229998.9
```

**Variance Inflation Factor (VIF)**

*Need library DAAG. You may have to install it the first time.*

> **library(DAAG)**

> **fuel[1:2,]**

|   | ID | state | Fuel | Tax | License | Inc   | Road  |
|---|----|-------|------|-----|---------|-------|-------|
| 1 | 1  | ME    | 541  | 9   | 52.5    | 3.571 | 1.976 |
| 2 | 2  | NH    | 524  | 9   | 57.2    | 4.092 | 1.250 |

> **attach(fuel)**

*Run a regression, saving results.*

> **mod<-lm(Fuel~Tax+License+Inc+Road)**

*Here are the VIF's*

> **vif(mod)**

|  | Tax    | License | Inc    | Road   |
|--|--------|---------|--------|--------|
|  | 1.6257 | 1.2164  | 1.0433 | 1.4969 |

*You can convert the VIF's to R<sup>2</sup>*

> **1-1/vif(mod)**

|  | Tax               | License    | Inc        | Road       |
|--|-------------------|------------|------------|------------|
|  | <u>0.38488036</u> | 0.17790201 | 0.04150292 | 0.33195270 |

*This says: If you predict Tax from License, Inc and Road, the R<sup>2</sup> is 0.3849. You could do the regression and get the same answer; see below.*

> **summary(lm(Tax~License+Inc+Road))**

Call:

lm(formula = Tax ~ License + Inc + Road)

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t ) |     |
|-------------|----------|------------|---------|----------|-----|
| (Intercept) | 11.14672 | 1.35167    | 8.247   | 1.79e-10 | *** |
| License     | -0.05791 | 0.02058    | -2.814  | 0.00728  | **  |
| Inc         | 0.15455  | 0.19881    | 0.777   | 0.44109  |     |
| Road        | -0.14935 | 0.03232    | -4.621  | 3.34e-05 | *** |

Residual standard error: 0.7707 on 44 degrees of freedom  
Multiple R-Squared: 0.3849, Adjusted R-squared: 0.3429  
 F-statistic: 9.177 on 3 and 44 DF, p-value: 7.857e-05

### Spjotvoll's Method in Variable Selection

The maximum model has variables  $\{1, 2, \dots, k\}$ ,  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon$  where the errors  $\varepsilon$  are independent and  $N(0, \sigma^2)$ . Suppose that  $T$  is the true model, where  $T$  is a subset of  $\{1, 2, \dots, k\}$ . That is,  $T$  is the model with exactly the nonzero  $\beta_j$ s, so  $\beta_j = 0$  if and only if  $j$  is not in  $T$ . We could do a general linear hypothesis F-test to test model  $T$  against the maximum model, and if we did just this one test, the chance that we would falsely reject model  $T$  at level  $\alpha = 0.05$  would be 5%. The problem is that we don't know  $T$ , so we end up testing many models, and might make many mistakes in all those tests.

Spjotvoll (1977) defines an inadequate model as any model that omits a variable in  $T$ , and an adequate model as any model that includes all of the variables in  $T$ , perhaps including some extra variables whose coefficients are 0. If  $k=5$  and  $T=\{1, 2\}$ , then  $\{1, 2\}$  and  $\{1, 2, 3\}$  are adequate models, but  $\{1, 3\}$  and  $\{2, 3, 4, 5\}$  are inadequate models. Spjotvoll wants to reject some models as inadequate. He is not worried about having too many variables, and is only worried about omitting a needed variable.

Spjotvoll (1977) declares a model  $Q$  to be inadequate at level  $\alpha = 0.05$  if and only if the F-test rejects both  $Q$  and every model contained in  $Q$ . With  $k=5$ , to reject  $Q=\{1, 3\}$  as inadequate, you would have to reject  $Q=\{1, 3\}$ , and also  $\{1\}$ ,  $\{3\}$  and  $\{\}$ , where  $\{\}$  is the model with no variables, that is,  $\{\}$  is  $y = \beta_0 + \varepsilon$ . So to reject  $Q=\{1, 3\}$ , four F-tests have to reject  $\alpha = 0.05$ . A different way of saying the same thing is that the maximum of the four p-values must be less than or equal to  $\alpha = 0.05$ ; that is, the maximum of the F-test p-values for  $\{1, 3\}$ ,  $\{1\}$ ,  $\{3\}$ , and  $\{\}$  must be less than or equal to 0.05. The chance that Spjotvoll's makes at least one mistake, saying that an adequate model is inadequate, is  $\alpha = 0.05$ , despite doing lots of tests.

It is easy to see why this works. Model  $Q$  is adequate if and only if the true model  $T$  is a subset of  $Q$ , possibly  $T=Q$ ; that's the definition of "adequate". Suppose  $Q$  is adequate, so it contains (or equals)  $T$ . To declare  $Q$  inadequate at the 0.05 level, you have to reject every model formed as a subset of  $Q$  - that's Spjotvoll's method - but as  $T$  is one of those subsets, you have to reject  $T$ , and

the chance that the F-test falsely rejects the true model  $T$  is  $\alpha=0.05$ .

Instead of focusing on  $\alpha=0.05$ , we could do this for any  $\alpha$  in just the same way. We can define an adjusted p-value for model  $Q$  as the maximum F-test p-value for all of the models that are subsets of  $Q$ , including  $Q$  itself and the empty model. This adjusted p-value rejects  $Q$  as inadequate at level  $\alpha$  if and only if the adjusted p-value is at most  $\alpha$ .

**Spjotvoll's Method in Variable Selection, continued.**

```
> attach(O2Uptake)
> y<-LogO2Up
> x<-O2Uptake[,2:6]
```

**Test of model lm(LogO2Up~TS+COD). Compare with row 7.**

```
>anova(lm(LogO2Up~TS+COD),lm(LogO2Up~Bod+TKN+TS+TVS+COD))
```

|   | Res.Df | RSS     | Df | Sum of Sq | F             | Pr(>F)        |
|---|--------|---------|----|-----------|---------------|---------------|
| 1 | 17     | 1.08502 |    |           |               |               |
| 2 | 14     | 0.96512 | 3  | 0.1199    | <b>0.5797</b> | <b>0.6379</b> |

The p-value from the F-test is 0.638, whereas the adjusted p-value is the maximum of the p-values for the models contained in {TS,COD}, namely 0.638 for {TS,COD}, 0.139 for {TS}, 0.129 for {COD} and 0.000 for the empty model {}, so the adjusted p-value is 0.638. So this model is not declared inadequate.

Model #16, {TKN, TVS}, is declared inadequate, but its adjusted p-value comes from model #5, {TVS}, because the largest p-value for a model contained in {TKN,TVS} is from model {TVS}.

```
> spjotvoll(x,y)
```

|    | p | Cp     | Fp           | pval         | adjusted.pval | inadequate | Bod | TKN | TS | TVS | COD |
|----|---|--------|--------------|--------------|---------------|------------|-----|-----|----|-----|-----|
| 1  | 1 | 55.463 | 11.893       | 0.000        | 0.000         | TRUE       | 0   | 0   | 0  | 0   | 0   |
| 2  | 2 | 6.297  | 2.074        | 0.139        | 0.139         | FALSE      | 0   | 0   | 1  | 0   | 0   |
| 3  | 2 | 6.576  | 2.144        | 0.129        | 0.129         | FALSE      | 0   | 0   | 0  | 0   | 1   |
| 4  | 2 | 13.505 | 3.876        | 0.025        | 0.025         | TRUE       | 1   | 0   | 0  | 0   | 0   |
| 5  | 2 | 20.331 | 5.583        | 0.007        | 0.007         | TRUE       | 0   | 0   | 0  | 1   | 0   |
| 6  | 2 | 56.861 | 14.715       | 0.000        | 0.000         | TRUE       | 0   | 1   | 0  | 0   | 0   |
| 7  | 3 | 1.739  | <b>0.580</b> | <b>0.638</b> | 0.638         | FALSE      | 0   | 0   | 1  | 0   | 1   |
| 8  | 3 | 5.274  | 1.758        | 0.201        | 0.201         | FALSE      | 0   | 0   | 0  | 1   | 1   |
| 9  | 3 | 6.872  | 2.291        | 0.123        | 0.129         | FALSE      | 0   | 1   | 0  | 0   | 1   |
| 10 | 3 | 6.885  | 2.295        | 0.122        | 0.139         | FALSE      | 1   | 0   | 1  | 0   | 0   |
| 11 | 3 | 7.165  | 2.388        | 0.113        | 0.139         | FALSE      | 0   | 0   | 1  | 1   | 0   |
| 12 | 3 | 7.336  | 2.445        | 0.107        | 0.139         | FALSE      | 0   | 1   | 1  | 0   | 0   |
| 13 | 3 | 7.705  | 2.568        | 0.096        | 0.129         | FALSE      | 1   | 0   | 0  | 0   | 1   |
| 14 | 3 | 9.097  | 3.032        | 0.065        | 0.065         | FALSE      | 1   | 1   | 0  | 0   | 0   |
| 15 | 3 | 11.331 | 3.777        | 0.036        | 0.036         | TRUE       | 1   | 0   | 0  | 1   | 0   |
| 16 | 3 | 21.369 | 7.123        | 0.004        | 0.007         | TRUE       | 0   | 1   | 0  | 1   | 0   |
| 17 | 4 | 2.319  | 0.160        | 0.854        | 0.854         | FALSE      | 0   | 1   | 1  | 0   | 1   |
| 18 | 4 | 3.424  | 0.712        | 0.508        | 0.638         | FALSE      | 0   | 0   | 1  | 1   | 1   |
| 19 | 4 | 3.439  | 0.720        | 0.504        | 0.638         | FALSE      | 1   | 0   | 1  | 0   | 1   |
| 20 | 4 | 5.665  | 1.833        | 0.196        | 0.201         | FALSE      | 0   | 1   | 0  | 1   | 1   |
| 21 | 4 | 6.253  | 2.126        | 0.156        | 0.156         | FALSE      | 1   | 1   | 1  | 0   | 0   |
| 22 | 4 | 6.515  | 2.258        | 0.141        | 0.141         | FALSE      | 1   | 1   | 0  | 0   | 1   |
| 23 | 4 | 7.152  | 2.576        | 0.112        | 0.201         | FALSE      | 1   | 0   | 0  | 1   | 1   |
| 24 | 4 | 8.155  | 3.077        | 0.078        | 0.139         | FALSE      | 1   | 0   | 1  | 1   | 0   |
| 25 | 4 | 8.165  | 3.082        | 0.078        | 0.139         | FALSE      | 0   | 1   | 1  | 1   | 0   |
| 26 | 4 | 8.681  | 3.341        | 0.065        | 0.065         | FALSE      | 1   | 1   | 0  | 1   | 0   |
| 27 | 5 | 4.001  | 0.001        | 0.972        | 0.972         | FALSE      | 0   | 1   | 1  | 1   | 1   |
| 28 | 5 | 4.319  | 0.319        | 0.581        | 0.854         | FALSE      | 1   | 1   | 1  | 0   | 1   |
| 29 | 5 | 5.068  | 1.068        | 0.319        | 0.638         | FALSE      | 1   | 0   | 1  | 1   | 1   |
| 30 | 5 | 6.776  | 2.776        | 0.118        | 0.201         | FALSE      | 1   | 1   | 0  | 1   | 1   |
| 31 | 5 | 7.697  | 3.697        | 0.075        | 0.156         | FALSE      | 1   | 1   | 1  | 1   | 0   |
| 32 | 6 | 6.000  | NA           | 1.000        | 1.000         | FALSE      | 1   | 1   | 1  | 1   | 1   |

Spjotvoll's method is a case of clostet testing; see Marcus et al. (1976)

Spjotvoll, E. (1977) Alternatives to plotting  $C_p$  in multiple regression. *Biometrika* 64, 1-8. Correction: page 241.

Marcus R, Peritz E, Gabriel KR. (1976). On closed testing procedures with special reference to ordered analysis of variance. *Biometrika* 63, 655-60.

**ANOVA**

*Memory data. 36 kids randomized to form 3 groups of 12, which were given different treatments. The 'data' are columns 1 and 2, for group and 4=words. It is a "balanced design" because every group has the same sample size. The rest of memory consists of various ways of coding the 2 degrees of freedom between the three groups into two coded variables. The variables ten and five are "dummy variables" for two categories, leaving out the third category. The variables five\_ten and nh\_ten are used to produce effects that are deviations from a mean for all three groups. The best coding is hier and info which involve "orthogonal contrasts," discussed below. It is only with orthogonal contrasts that you partition the sum of squares between groups into single degree of freedom parts that add back to the total.*

```

> memory[1:3, ]
> memory
      group words five_ten nh_ten ten five hier info
1      Ten   50      -1     -1   1    0  0.5   1
2      Ten   49      -1     -1   1    0  0.5   1
3      Ten   44      -1     -1   1    0  0.5   1
4      Ten   31      -1     -1   1    0  0.5   1
5      Ten   47      -1     -1   1    0  0.5   1
6      Ten   38      -1     -1   1    0  0.5   1
7      Ten   38      -1     -1   1    0  0.5   1
8      Ten   48      -1     -1   1    0  0.5   1
9      Ten   45      -1     -1   1    0  0.5   1
10     Ten   48      -1     -1   1    0  0.5   1
11     Ten   35      -1     -1   1    0  0.5   1
12     Ten   33      -1     -1   1    0  0.5   1
13     Five  44       1      0   0    1  0.5  -1
14     Five  41       1      0   0    1  0.5  -1
15     Five  34       1      0   0    1  0.5  -1
16     Five  35       1      0   0    1  0.5  -1
17     Five  40       1      0   0    1  0.5  -1
18     Five  44       1      0   0    1  0.5  -1
19     Five  39       1      0   0    1  0.5  -1
20     Five  39       1      0   0    1  0.5  -1
21     Five  45       1      0   0    1  0.5  -1
22     Five  41       1      0   0    1  0.5  -1
23     Five  46       1      0   0    1  0.5  -1
24     Five  32       1      0   0    1  0.5  -1
25 NoHier  33       0      1   0    0 -1.0   0
26 NoHier  36       0      1   0    0 -1.0   0
27 NoHier  37       0      1   0    0 -1.0   0
28 NoHier  42       0      1   0    0 -1.0   0
29 NoHier  33       0      1   0    0 -1.0   0
30 NoHier  33       0      1   0    0 -1.0   0
31 NoHier  41       0      1   0    0 -1.0   0
32 NoHier  33       0      1   0    0 -1.0   0
33 NoHier  38       0      1   0    0 -1.0   0
34 NoHier  39       0      1   0    0 -1.0   0
35 NoHier  28       0      1   0    0 -1.0   0
36 NoHier  42       0      1   0    0 -1.0   0

```



**ANOVA**

```
> attach(memory)
```

*The anova can be done as a linear model with a factor as the predictor.*

```
> anova(lm(words~group))
```

Analysis of Variance Table

Response: words

|           | Df | Sum Sq | Mean Sq | F value | Pr(>F)    |
|-----------|----|--------|---------|---------|-----------|
| group     | 2  | 215.06 | 107.53  | 3.7833  | 0.03317 * |
| Residuals | 33 | 937.92 | 28.42   |         |           |

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

*Or you can use the aov command. You get the same answer.*

```
> summary(aov(words~group))
```

|           | Df | Sum Sq | Mean Sq | F value | Pr(>F)    |
|-----------|----|--------|---------|---------|-----------|
| group     | 2  | 215.06 | 107.53  | 3.7833  | 0.03317 * |
| Residuals | 33 | 937.92 | 28.42   |         |           |

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

**Multiple Comparisons using Tukey's Method**

```
> TukeyHSD(aov(words~group))
```

Tukey multiple comparisons of means  
95% family-wise confidence level

Fit: aov(formula = words ~ group)

| \$group     |  | diff      | lwr        | upr       |
|-------------|--|-----------|------------|-----------|
| NoHier-Five |  | -3.750000 | -9.0905714 | 1.590571  |
| Ten-Five    |  | 2.166667  | -3.1739047 | 7.507238  |
| Ten-NoHier  |  | 5.916667  | 0.5760953  | 11.257238 |

*These are simultaneous 95% confidence intervals for the difference in means between two groups. The promise is that all 3 confidence intervals will cover their population differences in 95% of experiments. This is a better promise than that each one, by itself, covers in 95% of uses, because then the first interval would have a 5% chance of error, and so would the second, and so would the third, and the chance of at least one error would be greater than 5%. If the interval includes zero, as the first two intervals do, then you can't declare the two groups significantly different. If the interval excludes zero, as the third interval does, you can declare the two groups significantly different.*

### Tukey, Bonferroni and Holm

```

> help(pairwise.t.test)
> help(p.adjust)

> TukeyHSD(aov(words~group))
  Tukey multiple comparisons of means
    95% family-wise confidence level
Fit: aov(formula = words ~ group)
$group
              diff            lwr            upr
NoHier-Five -3.750000 -9.0905714  1.590571
Ten-Five     2.166667 -3.1739047  7.507238
Ten-NoHier   5.916667  0.5760953 11.257238

> pairwise.t.test(words,group,p.adj = "none")
  Pairwise comparisons using t tests with pooled SD
data:  words and group
       Five NoHier
NoHier 0.094 -
Ten     0.327 0.010
P value adjustment method: none

> pairwise.t.test(words,group,p.adj = "bonf")
  Pairwise comparisons using t tests with pooled SD
data:  words and group
       Five NoHier
NoHier 0.283 -
Ten     0.980 0.031
P value adjustment method: bonferroni

> pairwise.t.test(words,group,p.adj = "holm")
  Pairwise comparisons using t tests with pooled SD
data:  words and group
       Five NoHier
NoHier 0.189 -
Ten     0.327 0.031

Holm, S. (1979) A simple sequentially rejective multiple test
procedure. Scandinavian Journal of Statistics, 6, 65-
70. http://www.jstor.org/

Wright, S. P. (1992). Adjusted P-values for simultaneous
inference. Biometrics, 48, 1005-1013. http://www.jstor.org/

```

**ANOVA: Many Ways to Code the Same Anova**

*Here are three different codings with the same Anova table. Notice that much is the same, but some things differ. >*

**summary(lm(words~ten+five))**

Call:

lm(formula = words ~ ten + five)

Residuals:

| Min     | 1Q     | Median | 3Q    | Max   |
|---------|--------|--------|-------|-------|
| -11.167 | -3.479 | 0.875  | 4.771 | 7.833 |

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t )   |
|-------------|----------|------------|---------|------------|
| (Intercept) | 36.250   | 1.539      | 23.554  | <2e-16 *** |
| ten         | 5.917    | 2.176      | 2.718   | 0.0104 *   |
| five        | 3.750    | 2.176      | 1.723   | 0.0943 .   |

Residual standard error: 5.331 on 33 degrees of freedom  
 Multiple R-Squared: 0.1865, Adjusted R-squared: 0.1372  
 F-statistic: 3.783 on 2 and 33 DF, p-value: 0.03317

**> summary(lm(words~five\_ten+nh\_ten))**

Call:

lm(formula = words ~ five\_ten + nh\_ten)

Residuals:

| Min     | 1Q     | Median | 3Q    | Max   |
|---------|--------|--------|-------|-------|
| -11.167 | -3.479 | 0.875  | 4.771 | 7.833 |

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t )   |
|-------------|----------|------------|---------|------------|
| (Intercept) | 39.4722  | 0.8885     | 44.424  | <2e-16 *** |
| five_ten    | 0.5278   | 1.2566     | 0.420   | 0.6772     |
| nh_ten      | -3.2222  | 1.2566     | -2.564  | 0.0151 *   |

Residual standard error: 5.331 on 33 degrees of freedom  
 Multiple R-Squared: 0.1865, Adjusted R-squared: 0.1372  
 F-statistic: 3.783 on 2 and 33 DF, p-value: 0.03317

**> summary(lm(words~hier+info))**

Call:

lm(formula = words ~ hier + info)

Residuals:

| Min     | 1Q     | Median | 3Q    | Max   |
|---------|--------|--------|-------|-------|
| -11.167 | -3.479 | 0.875  | 4.771 | 7.833 |

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t )   |
|-------------|----------|------------|---------|------------|
| (Intercept) | 39.4722  | 0.8885     | 44.424  | <2e-16 *** |
| hier        | 3.2222   | 1.2566     | 2.564   | 0.0151 *   |
| info        | 1.0833   | 1.0882     | 0.996   | 0.3267     |

---  
 Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 5.331 on 33 degrees of freedom  
 Multiple R-Squared: 0.1865, Adjusted R-squared: 0.1372  
 F-statistic: 3.783 on 2 and 33 DF, p-value: 0.03317

**Contrasts in ANOVA: Better Coding of Nominal Variables**

*The third coding is best, because the predictors (contrasts) are uncorrelated, so the sums of squares partition.*

```
> cor(memory[,3:4])
      five_ten nh_ten
five_ten  1.0  0.5
nh_ten    0.5  1.0
```

```
> cor(memory[,5:6])
      ten five
ten  1.0 -0.5
five -0.5  1.0
```

```
> cor(memory[,7:8])
      hier info
hier  1  0
info  0  1
```

*Notice that hier and info have zero correlation: they are orthogonal. Because of this, you can partition the two degrees of freedom between groups into separate sums of squares.*

```
> anova(lm(words~hier+info))
Analysis of Variance Table

      Df Sum Sq Mean Sq F value Pr(>F)
hier    1 186.89  186.89  6.5756 0.01508 *
info    1  28.17   28.17  0.9910 0.32674
Residuals 33 937.92   28.42
---
```

*Reverse the order of info and hier, and you get the same answer.*

```
> anova(lm(words~info+hier))
Analysis of Variance Table

      Df Sum Sq Mean Sq F value Pr(>F)
info    1  28.17   28.17  0.9910 0.32674
hier    1 186.89  186.89  6.5756 0.01508 *
Residuals 33 937.92   28.42
---
```

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

*You can't do this with correlated predictors, because they overlap, and the order of the variables changes the sum of squares for the variable, so one can't really say that what portion of the sum of squares belongs to the variable.*

```
> anova(lm(words~ten+five))
Analysis of Variance Table

      Df Sum Sq Mean Sq F value Pr(>F)
ten    1 130.68  130.68  4.5979 0.03947 *
five   1  84.38   84.38  2.9687 0.09425 .
Residuals 33 937.92   28.42
```

```
> anova(lm(words~five+ten))
Analysis of Variance Table

      Df Sum Sq Mean Sq F value Pr(>F)
five   1   5.01    5.01  0.1764 0.67720
ten    1 210.04  210.04  7.3902 0.01037 *
Residuals 33 937.92   28.42
```

### Coding the Contrasts in ANOVA in R

*This is about shortcuts to get R to convert a nominal variable into several contrasts. There's no new statistics here; just R details.*

*The group variable, memory\$group, is a factor.*

```
> is.factor(memory$group)
[1] TRUE
```

*This factor has 3 levels. Notice that the levels are ordered and the order matters.*

```
> levels(memory$group)
[1] "Five" "NoHier" "Ten"
```

```
> memory$group
 [1] Ten      Ten      Ten      Ten      Ten      Ten      Ten      Ten      Ten
[10] Ten      Ten      Ten      Five     Five     Five     Five     Five     Five
[19] Five     Five     Five     Five     Five     Five     NoHier  NoHier  NoHier
[28] NoHier  NoHier  NoHier  NoHier  NoHier  NoHier  NoHier  NoHier  NoHier
Levels: Five NoHier Ten
```

*If you do nothing, R codes a factor in a linear model using 'dummy coding'.*

```
> contrasts(memory$group)
      NoHier Ten
Five      0  0
NoHier    1  0
Ten       0  1
```

*You can change the coding. Essentially, you can replace the little table above by whatever you want. We will build a new 3x2 table and redefine the contrasts to be this new table.*

```
> hier2<-c(.5,-1,.5)
> hier2
[1] 0.5 -1.0 0.5
```

```
> info2<-c(-1,0,1)
> info2
[1] -1 0 1
```

```
> cm<-cbind(hier2,info2)
> cm
      hier2 info2
[1,]  0.5   -1
[2,] -1.0    0
[3,]  0.5    1
```

*So cm is our new table, and we redefine the contrasts for memory\$group.*

```
> contrasts(memory$group)<-cm
```

*This replaces the 'dummy coding' by our new coding.*

```
> contrasts(memory$group)
      hier2 info2
Five      0.5   -1
NoHier   -1.0    0
Ten       0.5    1
```

**Coding the Contrasts in ANOVA, Continued**

*If you ask R to extend the contrasts into variables, it will do this with "model.matrix". Notice that this is the coding in the original data matrix, but R is happy to generate it for you using the contrasts you specified.*

```
> m<-model.matrix(memory$words~memory$group)
> m
      (Intercept) memory$grouphier2 memory$groupinfo2
1             1             0.5             1
2             1             0.5             1
3             1             0.5             1
4             1             0.5             1
5             1             0.5             1
6             1             0.5             1
7             1             0.5             1
8             1             0.5             1
9             1             0.5             1
10            1             0.5             1
11            1             0.5             1
12            1             0.5             1
13            1             0.5             -1
14            1             0.5             -1
15            1             0.5             -1
16            1             0.5             -1
17            1             0.5             -1
18            1             0.5             -1
19            1             0.5             -1
20            1             0.5             -1
21            1             0.5             -1
22            1             0.5             -1
23            1             0.5             -1
24            1             0.5             -1
25            1             -1.0             0
26            1             -1.0             0
27            1             -1.0             0
28            1             -1.0             0
29            1             -1.0             0
30            1             -1.0             0
31            1             -1.0             0
32            1             -1.0             0
33            1             -1.0             0
34            1             -1.0             0
35            1             -1.0             0
36            1             -1.0             0
```

```
> hcontrast<-m[,2]
> icontrast<-m[,3]
```

*We now do the anova with single degree of freedom contrasts.*

```
> anova(lm(memory$words~hcontrast+icontrast))
```

Analysis of Variance Table

|           | Df | Sum Sq | Mean Sq | F value | Pr(>F)    |
|-----------|----|--------|---------|---------|-----------|
| hcontrast | 1  | 186.89 | 186.89  | 6.5756  | 0.01508 * |
| icontrast | 1  | 28.17  | 28.17   | 0.9910  | 0.32674   |
| Residuals | 33 | 937.92 | 28.42   |         |           |

**ANOVA DECOMPOSITION**

```

> mod<-aov(words~group,projections=T)

> mod$projections
  (Intercept)      group      Residuals
1    39.47222  2.6944444  7.833333e+00
2    39.47222  2.6944444  6.833333e+00
3    39.47222  2.6944444  1.833333e+00
4    39.47222  2.6944444 -1.116667e+01
5    39.47222  2.6944444  4.833333e+00
6    39.47222  2.6944444 -4.166667e+00
7    39.47222  2.6944444 -4.166667e+00
8    39.47222  2.6944444  5.833333e+00
9    39.47222  2.6944444  2.833333e+00
10   39.47222  2.6944444  5.833333e+00
11   39.47222  2.6944444 -7.166667e+00
12   39.47222  2.6944444 -9.166667e+00
13   39.47222  0.5277778  4.000000e+00
14   39.47222  0.5277778  1.000000e+00
15   39.47222  0.5277778 -6.000000e+00
16   39.47222  0.5277778 -5.000000e+00
17   39.47222  0.5277778 -9.503032e-16
18   39.47222  0.5277778  4.000000e+00
19   39.47222  0.5277778 -1.000000e+00
20   39.47222  0.5277778 -1.000000e+00
21   39.47222  0.5277778  5.000000e+00
22   39.47222  0.5277778  1.000000e+00
23   39.47222  0.5277778  6.000000e+00
24   39.47222  0.5277778 -8.000000e+00
25   39.47222 -3.2222222 -3.250000e+00
26   39.47222 -3.2222222 -2.500000e-01
27   39.47222 -3.2222222  7.500000e-01
28   39.47222 -3.2222222  5.750000e+00
29   39.47222 -3.2222222 -3.250000e+00
30   39.47222 -3.2222222 -3.250000e+00
31   39.47222 -3.2222222  4.750000e+00
32   39.47222 -3.2222222 -3.250000e+00
33   39.47222 -3.2222222  1.750000e+00
34   39.47222 -3.2222222  2.750000e+00
35   39.47222 -3.2222222 -8.250000e+00
36   39.47222 -3.2222222  5.750000e+00
attr(,"df")
  (Intercept)      group      Residuals
            1            2            33

```

### Orthogonal and Non-orthogonal Predictors

```
> attach(memory)
> summary(aov(words~group))
```

|           | Df | Sum Sq | Mean Sq | F value | Pr(>F)    |
|-----------|----|--------|---------|---------|-----------|
| group     | 2  | 215.06 | 107.53  | 3.7833  | 0.03317 * |
| Residuals | 33 | 937.92 | 28.42   |         |           |

---

*ten and five are not orthogonal predictors - so there is not a unique sum of squares for each*

```
> anova(lm(words~ten+five))
```

Analysis of Variance Table

Response: words

|           | Df | Sum Sq | Mean Sq | F value | Pr(>F)    |
|-----------|----|--------|---------|---------|-----------|
| ten       | 1  | 130.68 | 130.68  | 4.5979  | 0.03947 * |
| five      | 1  | 84.38  | 84.38   | 2.9687  | 0.09425 . |
| Residuals | 33 | 937.92 | 28.42   |         |           |

---

```
> anova(lm(words~five+ten))
```

Analysis of Variance Table

Response: words

|           | Df | Sum Sq | Mean Sq | F value | Pr(>F)    |
|-----------|----|--------|---------|---------|-----------|
| five      | 1  | 5.01   | 5.01    | 0.1764  | 0.67720   |
| ten       | 1  | 210.04 | 210.04  | 7.3902  | 0.01037 * |
| Residuals | 33 | 937.92 | 28.42   |         |           |

*her and info are orthogonal predictors - so there is a unique sum of squares for each*

```
> anova(lm(words~hier+info))
```

Analysis of Variance Table

Response: words

|           | Df | Sum Sq | Mean Sq | F value | Pr(>F)    |
|-----------|----|--------|---------|---------|-----------|
| hier      | 1  | 186.89 | 186.89  | 6.5756  | 0.01508 * |
| info      | 1  | 28.17  | 28.17   | 0.9910  | 0.32674   |
| Residuals | 33 | 937.92 | 28.42   |         |           |

```
> anova(lm(words~info+hier))
```

Analysis of Variance Table

Response: words

|           | Df | Sum Sq | Mean Sq | F value | Pr(>F)    |
|-----------|----|--------|---------|---------|-----------|
| info      | 1  | 28.17  | 28.17   | 0.9910  | 0.32674   |
| hier      | 1  | 186.89 | 186.89  | 6.5756  | 0.01508 * |
| Residuals | 33 | 937.92 | 28.42   |         |           |



**Simulating in R**

*Ten observations from the standard Normal distribution:*

```
> rnorm(10)
[1] 0.8542301 -1.3331572 1.4522862 0.8980641 0.1456334
[6] 0.4926661 -0.4366962 0.6204263 -0.1582319 -0.6444449
```

*Fixed integer sequences*

```
> 1:2
[1] 1 2
```

```
> 1:5
[1] 1 2 3 4 5
```

```
> 0:1
[1] 0 1
```

*20 coin flips*

```
> sample(0:1,20,r=T)
[1] 0 1 0 0 1 1 0 0 0 1 0 1 0 0 0 1 0 1 1 1
```

*10 random numbers from 1 to 5*

```
> sample(1:5,10,r=T)
[1] 5 2 3 5 2 2 1 1 3 2
```

*More information:*

```
help(sample)
help(rnorm)
```

PROBLEM SET #1 STATISTICS 500 Fall 2017: DATA PAGE 1

**Due at noon on Tuesday, 24 Oct 2017 in class.**

**This is an exam. Do not discuss it with anyone.**

The data are from NHANES 2009-2010. You can obtain the original data at <https://www.cdc.gov/nchs/nhanes/>, but there is no reason to do this unless you want to. The data relate lung function to smoking. The data are in an object **smokelung** in the course workspace. You will have to download the workspace again, and may need to clear your web browser's cache to get the new version. The data are also available as a csv file on my web page using the button data.cv. There are 2360 people, of whom 1842 never smoked and 518 are daily smokers, meaning that they smoked at least 5 cigarettes per day every day of the last 30 days. For an explanation of the lung function measures, fvc, fev1, and ratio = fev1/fvc, see [http://oac.med.jhmi.edu/res\\_phys/Encyclopedia/ForcedExpiration/ForcedExpiration.HTML](http://oac.med.jhmi.edu/res_phys/Encyclopedia/ForcedExpiration/ForcedExpiration.HTML) For an explanation of bmi, see [https://www.nhlbi.nih.gov/health/educational/lose\\_wt/BMI/bmicalc.htm](https://www.nhlbi.nih.gov/health/educational/lose_wt/BMI/bmicalc.htm) The variables are

**SEQN** = NHANES id number

**fvc** = forced vital capacity in ml

**fev1** = forced expiratory volume in 1 second, in ml

**ratio** = fev1/fvc

**cigsperday** = cigarettes smoked per day, 0 for never smokers

**smoke** = daily or never

**female** = 1 for female, 0 for male

**age** = age in years, >=20

**bmi** = body mass index

**educ, educf** = education, 1=<9<sup>th</sup> grade, 2=9-11<sup>th</sup> grade, 3=high school or equivalent, 4="some college", say a 2 year associates degree, 5="college", >= 4 year BA degree.

**income, incomef** = ratio of family income to the poverty level, capped at 5xPoverty.

**cotinine** = cotinine in blood, ng/ml, a marker for recent tobacco exposure

**lead** = lead in blood, ug/dL

**cadmium** = cadmium in blood ug/L

```
> dim(smokelung)
```

```
[1] 2360 16
```

```
> attach(smokelung)
```

You should plot the data in various ways, such as

```
plot(cigsperday,ratio)
```

```
lines(lowess(cigsperday,ratio))
```

```
boxplot(ratio~(bmi>30))
```

```
boxplot(ratio~educf)
```

```
boxplot(ratio~incomef)
```

**This is an exam. Do not discuss it with anyone.**

**Due at noon on Tuesday, 24 Oct 2017 in class.**

Model 1:

$$\text{ratio} = \beta_0 + \beta_1 \text{cigsperday} + \beta_2 \text{bmi} + \beta_3 \text{female} + \beta_4 \text{age} + \varepsilon$$

where  $\varepsilon$  is iid  $N(0, \sigma^2)$ .

**Important:** Write your name on both sides of the answer page, **last name first**. **Turn in only the answer page**. Do not turn in plots. Brief answers suffice. **Circle** the correct answer, but do not cross out an answer. A circled answer may be correct or incorrect, but every crossed out answer is incorrect. Do not give one answer adding a note explaining why a different answer is correct.

**This is an exam. Do not discuss the exam with anyone.** If you discuss the exam, you have cheated on an exam. The single dumbest thing a PhD student at Penn can do is cheat on an exam.

Name (**Last**, First): \_\_\_\_\_ ID# \_\_\_\_\_

PROBLEM SET #1 STATISTICS 500 Fall 2017: ANSWER PAGE 1

**This is an exam. Do not discuss it. Due noon Oct 24.**

|   |  |
|---|--|
| Part 1.   | Fill in or <b>circle</b> the correct answer.                                   |
| 1.1 What is the value of the smallest ratio? How old is this person? Does this person smoke?                                | ratio = _____ age = _____<br><b>Circle one</b><br>Daily-smoker    Never-smoker |
| 1.2 What is the value of the largest bmi? How old is this person? Does this person smoke?                                   | bmi = _____ age = _____<br><b>Circle one</b><br>Daily-smoker    Never-smoker   |
| 1.3 The distribution of ratio is skewed right, with a longer tail to the right (high ratios) than to the left (low ratios). | <b>Circle one</b><br>True    False   |

|  |   |
|--|---|
| Fit <b>model 1</b> on the data page. Use model 1 for the questions in part 2.  | Fill in or <b>circle</b> the correct answer.  |
| 2.1 Test the null hypothesis $H_0: \beta_1 = 0$ , that the coefficient of cigspersday is zero. What is the <b>name</b> of the test, the <b>value</b> of the test statistic, the 2-sided <b>P-value</b> . Is $H_0$ <b>plausible</b> ? | Name: _____ Value: _____<br>P-value: _____<br><b>Circle one</b><br>Plausible    Not plausible |
| 2.2 Give the 2-sided 95% confidence interval for $\beta_1$ , the coefficient of cigspersday.   | [ _____ , _____ ]   |
| 2.3 Test the null hypothesis $H_0: \beta_1 = \beta_2 = \beta_3 = \beta_4 = 0$ . What is the <b>name</b> of the test, the <b>value</b> of the test statistic, the 2-sided <b>P-value</b> . Is $H_0$ <b>plausible</b> ?                | Name: _____ Value: _____<br>P-value: _____<br><b>Circle one</b><br>Plausible    Not plausible |
| 2.4 What is the value of the correlation between ratio and its fitted values from model 1? What is the square of this correlation?   | Correlation = _____<br>Correlation <sup>2</sup> = _____                                       |

Name (**Last**, First): \_\_\_\_\_ ID# \_\_\_\_\_

PROBLEM SET #1 STATISTICS 500 Fall 2017: ANSWER PAGE 1

**This is an exam. Do not discuss it. Due noon Oct 24.**

3.1 In model 1, test the null hypothesis  $H_0: \beta_2 = \beta_3 = \beta_4 = 0$ . Fill in the anova table.

| Source of Variation           | Sums of square | Degrees of freedom | Mean square | F-ratio                      |
|-------------------------------|----------------|--------------------|-------------|------------------------------|
| Full Model                    |                |                    |             |                              |
| cigsperday Alone              |                |                    |             | XXXXXXXXXXXX<br>XXXXXXXXXXXX |
| Added by bmi, age, and female |                |                    |             |                              |
| Residual from full model      |                |                    |             | XXXXXXXXXXXX<br>XXXXXXXXXXXX |

|  |  |
|--|--|
| <p>3.2 In model 1, test the null hypothesis <math>H_0: \beta_2 = \beta_3 = \beta_4 = 0</math>. Fill in the anova table. What is the <b>name</b> of the test, the <b>value</b> of the test statistic, the 2-sided <b>P-value</b>. Is <math>H_0</math> <b>plausible</b>?</p> | <p>Name: _____ Value: _____</p> <p>P-value: _____</p> <p style="text-align: center;"><b>Circle one</b></p> <p>Plausible      Not plausible</p> |
|--|--|

|  |   |
|--|---|
| <p>Use model 1 to answer the following questions.</p>  | <p>Fill in or <b>circle</b> the correct answer</p>  |
| <p>4.1 The plot of residuals against fitted values shows a clear U shape, with the largest positive residuals at the largest and smallest fitted values.</p> | <p style="text-align: center;"><b>Circle one</b></p> <p style="text-align: center;">True    False</p> |
| <p>4.2 The Normal quantile plot of residuals indicates that they look Normal.</p>  | <p style="text-align: center;"><b>Circle one</b></p> <p style="text-align: center;">True    False</p> |
| <p>4.3 There are 8 residuals <math>\leq -0.3</math>, and no residuals <math>\geq 0.3</math></p>  | <p style="text-align: center;"><b>Circle one</b></p> <p style="text-align: center;">True    False</p> |
| <p>4.4 The Shapiro-Wilk test applied to the residuals from model 1 accepts, at the 0.05 level, the null hypothesis that the residuals are Normal.</p>        | <p style="text-align: center;"><b>Circle one</b></p> <p style="text-align: center;">True    False</p> |

ANSWERS

PROBLEM SET #1 STATISTICS 500 Fall 2017: ANSWER PAGE 1

**This is an exam. Do not discuss it. Due noon Oct 24.**

|  |  |
|--|--|
| Part 1. 7 points each, except as noted.  | Fill in or <b>circle</b> the correct answer.   |
| 1.1 What is the value of the smallest ratio? How old is this person? Does this person smoke?   | ratio = 0.3130 age = 69<br><b>Circle one</b><br>Daily-smoker Never-smoker                                  |
| 1.2 What is the value of the largest bmi? How old is this person? Does this person smoke?  | bmi = 82.1 age = 64<br><b>Circle one</b><br>Daily-smoker Never-smoker                                      |
| 1.3 The distribution of ratio is skewed right, with a longer tail to the right (high ratios) than to the left (low ratios).  | <b>Circle one</b><br>True False  |
| Fit <b>model 1</b> on the data page. Use model 1 for the questions in part 2.  | Fill in or <b>circle</b> the correct answer.   |
| 2.1 Test the null hypothesis $H_0: \beta_1 = 0$ , that the coefficient of cigspersday is zero. What is the <b>name</b> of the test, the <b>value</b> of the test statistic, the 2-sided <b>P-value</b> . Is $H_0$ <b>plausible</b> ? | Name: t-test Value: -13.85<br>P-value: $2 \times 10^{-16}$<br><b>Circle one</b><br>Plausible Not plausible |
| 2.2 Give the 2-sided 95% confidence interval for $\beta_1$ , the coefficient of cigspersday.   | $[-.00306, -0.00230]$  |
| 2.3 Test the null hypothesis $H_0: \beta_1 = \beta_2 = \beta_3 = \beta_4 = 0$ . What is the <b>name</b> of the test, the <b>value</b> of the test statistic, the 2-sided <b>P-value</b> . Is $H_0$ <b>plausible</b> ?                | Name: F-test Value: 198.6<br>P-value: $2 \times 10^{-16}$<br><b>Circle one</b><br>Plausible Not plausible  |
| 2.4 What is the value of the correlation between ratio and its fitted values from model 1? What is the square of this correlation?   | Correlation = 0.502<br>Correlation <sup>2</sup> = 0.252  |

ANSWERS

PROBLEM SET #1 STATISTICS 500 Fall 2017: ANSWER PAGE 1

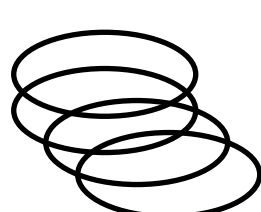
**This is an exam. Do not discuss it. Due noon Oct 24.**

3.1 In model 1, test the null hypothesis  $H_0: \beta_2 = \beta_3 = \beta_4 = 0$ . Fill in the anova table. (16 points)

| Source of Variation           | Sums of square | Degrees of freedom | Mean square | F-ratio                    |
|-------------------------------|----------------|--------------------|-------------|----------------------------|
| Full Model                    | 3.8825         | 4                  | 0.970625    | 198.6                      |
| cigsperday Alone              | 1.1541         | 1                  | 1.1541      | xxxxxxxxxxx<br>xxxxxxxxxxx |
| Added by bmi, age, and female | 2.7284         | 3                  | 0.9094667   | 186.12                     |
| Residual from full model      | 11.508         | 2355               | 0.00489     | xxxxxxxxxxx<br>xxxxxxxxxxx |

|  |   |
|--|---|
| 3.2 In model 1, test the null hypothesis $H_0: \beta_2 = \beta_3 = \beta_4 = 0$ . Fill in the anova table. What is the <b>name</b> of the test, the <b>value</b> of the test statistic, the 2-sided <b>P-value</b> . Is $H_0$ <b>plausible</b> ? | Name:F-test Value:186.12<br>P-value: $2 \times 10^{-16}$<br><b>Circle one</b><br>Plausible <b>Not plausible</b> |
|--|---|

|   |   |
|---|---|
| Use model 1 to answer the following questions.  | Fill in or <b>circle</b> the correct answer |
| 4.1 The plot of residuals against fitted values shows a clear U shape, with the largest positive residuals at the largest and smallest fitted values. | <b>Circle one</b><br>True <b>False</b>      |
| 4.2 The Normal quantile plot of residuals indicates that they look Normal.  | <b>Circle one</b><br>True <b>False</b>      |
| 4.3 There are 8 residuals $\leq -0.3$ , and no residuals $\geq 0.3$   | <b>Circle one</b><br><b>True</b> False      |
| 4.4 The Shapiro-Wilk test applied to the residuals from model 1 accepts, at the 0.05 level, the null hypothesis that the residuals are Normal.        | <b>Circle one</b><br>True <b>False</b>      |



PROBLEM SET #1 STATISTICS 500 Fall 2017

DOING THE PROBLEM SET IN R

```
attach(smokelung)
plot(cigsperday,ratio)
lines(lowess(cigsperday,ratio))
boxplot(ratio~(bmi>30))
boxplot(ratio~educf)
boxplot(ratio~incomef)
boxplot(bmi)

#Part 1
which.min(ratio)
smokelung[1744,]
which.max(bmi)
smokelung[1449,]

#Model 1
m<-lm(ratio~cigsperday+bmi+female+age)

#Part 2
summary(m)
confint(m)
cor(m$fitted.values,ratio)
cor(m$fitted.values,ratio)^2

#Part 3
mr<-lm(ratio~cigsperday)
anova(mr,m)

#Part 4
plot(m$fitted.values,m$residuals)
lines(lowess(m$fitted.values,m$residuals),col="red")
qqnorm(m$residuals)
qqline(m$residuals)
sum(m$residuals<=-.3)
sum(m$residuals>=.3)
shapiro.test(m$residuals)
detach(smokelung)
```



PROBLEM SET #2 STATISTICS 500 Fall 2017: DATA PAGE 1

**Due at noon on Tuesday, 24 Oct 2017 in class.**

**This is an exam. Do not discuss it with anyone.**

The data are the same as Problem Set 1, from NHANES 2009-2010. The data relate lung function to smoking. The data are in an object **smokelung** in the course workspace. The data are also available briefly as a csv file on my web page using the button data.cv. There are 2360 people, of whom 1842 never smoked and 518 are daily smokers, meaning that they smoked at least 5 cigarettes per day every day of the last 30 days. For an explanation of the lung function measures, fvc, fev1, and ratio = fev1/fvc, see [http://oac.med.jhmi.edu/res\\_phys/Encyclopedia/ForcedExpiration/ForcedExpiration.HTML](http://oac.med.jhmi.edu/res_phys/Encyclopedia/ForcedExpiration/ForcedExpiration.HTML) For an explanation of bmi, see [https://www.nhlbi.nih.gov/health/educational/lose\\_wt/BMI/bmicalc.htm](https://www.nhlbi.nih.gov/health/educational/lose_wt/BMI/bmicalc.htm) The variables are

**SEQN** = NHANES id number

**fvc** = forced vital capacity in ml

**fev1** = forced expiratory volume in 1 second, in ml

**ratio** = fev1/fvc

**cigsperday** = cigarettes smoked per day, 0 for never smokers

**smoke** = daily or never

**female** = 1 for female, 0 for male

**age** = age in years, >=20

**bmi** = body mass index

**educ, educf** = education, 1=<9<sup>th</sup> grade, 2=9-11<sup>th</sup> grade, 3=high school or equivalent, 4="some college", say a 2 year associates degree, 5="college", >= 4 year BA degree.

**income, incomef** = ratio of family income to the poverty level, capped at 5xPoverty.

**cotinine** = cotinine in blood, ng/ml, a marker for recent tobacco exposure

**lead** = lead in blood, ug/dL

**cadmium** = cadmium in blood ug/L

```
> dim(smokelung)
```

```
[1] 2360 16
```

```
> attach(smokelun
```

As always, you should plot the data in various ways.

**IMPORTANT:** If you look at the interaction of a binary (1 or 0) variable and a continuous variable, then **do not** center either variable before multiplying. You may mess up several questions if you do this incorrectly. Also, make sure you construct **mini** as described below.

**This is an exam. Do not discuss it with anyone.**

**Due at noon on Tuesday, 24 Oct 2017 in class.**

**You will need to construct several new variables.**

```
smoker<-1*(smokelung$smoke=="Daily")
smokerAge<-smoker*smokelung$age
d<-cbind(smokelung,smoker,smokerAge)
rm(smoker,smokerAge)
attach(d)
mini<-d[c(1781,1726,554,1394),]
mini<-mini[,c(8,9,18,19)]
mini
```

|      | female | age | smoker | smokerAge |
|------|--------|-----|--------|-----------|
| 1781 | 1      | 25  | 0      | 0         |
| 1726 | 1      | 25  | 1      | 25        |
| 554  | 1      | 60  | 0      | 0         |
| 1394 | 1      | 60  | 1      | 60        |

**So**, mini is a data.frame with just 4 people. Look at mini carefully, so you understand it. It has 4 females, two aged 25, two aged 60, two daily smokers, two never smokers.

**Model 1:** (based on d)

ratio=  $\beta_0 + \beta_1 \text{age} + \beta_2 \text{female} + \beta_3 \text{smoker} + \varepsilon$   
 where  $\varepsilon$  is iid  $N(0, \sigma^2)$ .

**Model 2:** (based on d)

ratio=  $\gamma_0 + \gamma_1 \text{age} + \gamma_2 \text{female} + \gamma_3 \text{smoker} + \gamma_4 \text{smokerAge} + \varepsilon$   
 where  $\varepsilon$  is iid  $N(0, \omega^2)$ .

**For question 2.2**, use the **restricted cubic spline** function, rcspline.eval(), in the Hmisc package, using the default to let the function decide the spacing of the 5 knots. The question asks you test for curvature in **age** in model 2. For understanding, you might plot the fitted values against age for the two models in 2.2.

**For question 3.1**, the **studentized residual** is obtained using rstudent() in R. It uses the deleted estimate of  $\omega^2$ .

**Important:** Write your name on both sides of the answer page, **last name first**. **Turn in only the answer page.**

Brief answers suffice. **Circle** the correct answer, but do not cross out an answer. A circled answer may be correct or incorrect, but every crossed out answer is incorrect.

**This is an exam. Do not discuss the exam with anyone.** If you discuss the exam, you have cheated on an exam. The single dumbest thing a PhD student at Penn can do is cheat on an exam.

Name (Last, First): \_\_\_\_\_ ID# \_\_\_\_\_

PROBLEM SET #2 STATISTICS 500 Fall 2017: ANSWER PAGE 1

**This is an exam. Do not discuss it. Due noon Oct 24.**

|   |  | Fill in/ <b>CIRCLE</b> the answer |        |
|---|--|-----------------------------------|--------|
| 1.1 Fit model 1 and use Tukey's test to see if a transformation of y would be needed if this model were used. Give the value of the t-statistic and the P-value. Does the test indicate that a transformation is needed? <b>Repeat</b> the calculation for model 2, giving just the P-value.  | <b>Model 1:</b><br>t-statistic: _____<br><br>P-value: _____<br>Transformation is:<br>(CIRCLE ONE)<br>NEEDED                  NOT NEEDED<br><br><b>Model 2:</b><br>P-value: _____           |                                   |        |
| 1.2 In models 1 and 2, is it plausible that the relationship between the relationship between ratio and age is parallel for smokers and nonsmokers? Give the name of the test statistic, its numerical value, its degrees of freedom, its P-value and state whether parallelism is plausible. | Name of test: _____<br>Value of the test statistic: _____<br><br>Degrees of freedom: _____<br><br>P-value: _____<br>Parallelism is:<br>(CIRCLE ONE)<br>PLAUSIBLE          NOT PLAUSIBLE    |                                   |        |
| 1.3 The estimated coefficient of smoker in model 2 is positive and not significantly different from zero. Recall that y=ratio is a measure of lung function.  | From the fact on the left, we can reasonably conclude that model 2 provides no indication that smoking is associated with lower values of ratio. CIRCLE ONE<br>TRUE                  FALSE |                                   |        |
| 1.4 Use the predict function applied to model 2 and to mini to estimate the expected ratio for female smokers and nonsmokers aged 25 and 60. Put the four estimates in the table.   |  | age=25                            | age=60 |
|   | smoke=0  |                                   |        |
| 1.5 Use the predict function applied to model 2 and to mini to obtain two-sided 95% confidence intervals for the expected ratio in part 1.4.  |  | age=25                            | age=60 |
|   | smoke=0  |                                   |        |
|   | smoke=1  |                                   |        |
|   |  |                                   |        |

(Round to 2 digits after the decimal.)

(Round to 2 digits).

Name (**Last**, First): \_\_\_\_\_ ID# \_\_\_\_\_

PROBLEM SET #2 STATISTICS 500 Fall 2017: ANSWER PAGE 2

**This is an exam. Do not discuss it. Due noon Oct 24.**

|   |  |
|---|--|
|   | Fill in/ <b>CIRCLE</b> the answer  |
| 2.1 Who said: "To find out what happens to a system when you interfere with it you have to interfere with it"? Circle <b>ONE</b> .  | Gauss Cauchy Sheather<br>Box Cox Fox Cook<br>Fisher Markov Tukey<br>Harrell Everitt Weisberg<br>Ripley Titterington Trump      |
| 2.2 Add to model #2 <b>restricted cubic splines</b> in <b>age</b> with 5 knots. <b>See the data page.</b> Test that model 2 is adequate against the alternative that the splines in age are needed. Give the <b>name</b> of the test, the <b>value</b> of the test statistic, the <b>degrees of freedom (df)</b> , the <b>P-value</b> , and state whether the test <b>indicates curvature</b> in age. | Name: _____ Value: _____<br>Degrees of Freedom: _____<br>P-value: _____<br>Indicates curvature in age?<br>CIRCLE ONE<br>YES NO |

|   |   |
|---|---|
| Use <b>model 2</b> for part 3.  | Fill in/ <b>CIRCLE</b> the answer   |
| 3.1 Which observation in model 2 has the largest <b>absolute</b> studentized residual? Give <b>signed</b> studentized residual, the row number and SEQN, and circle TRUE or FALSE about "This person is an outlier at the 0.05 level, taking account of the number of tests done." Give <b>DF</b> and adjusted <b>P-value</b> for test. | Row: _____ SEQN: _____<br>Studentized Residual: _____<br>This person has the <b>lowest ratio</b> .<br>CIRCLE: TRUE FALSE<br><b>Outlier:</b><br>CIRCLE: TRUE FALSE<br>DF=degrees of freedom:<br>DF= _____<br>Adjusted P-value: _____ |
| 3.2 Which observation has the largest <b>absolute</b> dffits? Give <b>signed</b> dffits, the row number and SEQN, and circle TRUE or FALSE. Does this observation move its $\hat{y}$ up 4.48 standard errors.   | Row: _____ SEQN: _____<br>dffits: _____<br>CIRCLE: YES NO   |
| 3.3 Which observation in model 2 has the largest <b>leverage</b> ? Give <b>hatvalue</b> , row # and SEQN.   | hatvalue: _____<br>Row: _____ SEQN: _____   |

**This is an exam. Do not discuss it.**

| Fill in/ <b>CIRCLE</b> the answer  |  |            |        |        |         |            |            |         |            |            |
|--|--|------------|--------|--------|---------|------------|------------|---------|------------|------------|
| <p>1.1 Fit model 1 and use Tukey's test to see if a transformation of y would be needed if this model were used. Give the value of the t-statistic and the P-value. Does the test indicate that a transformation is needed? <b>Repeat</b> the calculation for model 2, giving just the P-value.</p>  | <p><b>Model 1:</b><br/>                     t-statistic: -4.195<br/><br/>                     P-value: &lt;0.001<br/>                     Transformation is:<br/>                     (CIRCLE ONE)<br/> <input checked="" type="radio"/> <b>NEEDED</b>      <input type="radio"/> NOT NEEDED</p> <p><b>Model 2:</b><br/>                     P-value: 0.596</p>  |            |        |        |         |            |            |         |            |            |
| <p>1.2 In models 1 and 2, is it plausible that the relationship between the relationship between ratio and age is parallel for smokers and nonsmokers? Give the name of the test statistic, its numerical value, its degrees of freedom, its P-value and state whether parallelism is plausible.</p> | <p>Name of test: t-test<br/>                     Value of the test statistic: -5.530<br/><br/>                     Degrees of freedom: 2355<br/><br/>                     P-value: <math>3.56 \times 10^{-8}</math><br/>                     Parallelism is:<br/>                     (CIRCLE ONE)<br/> <input type="radio"/> PLAUSIBLE      <input checked="" type="radio"/> <b>NOT PLAUSIBLE</b></p> |            |        |        |         |            |            |         |            |            |
| <p>1.3 The estimated coefficient of smoker in model 2 is positive and not significantly different from zero. Recall that y=ratio is a measure of lung function.</p>  | <p>From the fact on the left, we can reasonably conclude that model 2 provides no indication that smoking is associated with lower values of ratio. <b>CIRCLE ONE</b><br/> <input type="radio"/> TRUE      <input checked="" type="radio"/> <b>FALSE</b></p>   |            |        |        |         |            |            |         |            |            |
| <p>1.4 Use the predict function applied to model 2 and to mini to estimate the expected ratio for female smokers and nonsmokers aged 25 and 60. Put the four estimates in the table.</p>   | <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th></th> <th>age=25</th> <th>age=60</th> </tr> </thead> <tbody> <tr> <td>smoke=0</td> <td>0.85</td> <td>0.78</td> </tr> <tr> <td>smoke=1</td> <td>0.82</td> <td>0.71</td> </tr> </tbody> </table> <p style="text-align: center;">(Round to 2 digits)</p>   |            | age=25 | age=60 | smoke=0 | 0.85       | 0.78       | smoke=1 | 0.82       | 0.71       |
|  | age=25   | age=60     |        |        |         |            |            |         |            |            |
| smoke=0  | 0.85   | 0.78       |        |        |         |            |            |         |            |            |
| smoke=1  | 0.82   | 0.71       |        |        |         |            |            |         |            |            |
| <p>1.5 Use the predict function applied to model 2 and to mini to obtain two-sided 95% <b>confidence intervals for the expected ratio</b> in part 1.4.</p>   | <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th></th> <th>age=25</th> <th>age=60</th> </tr> </thead> <tbody> <tr> <td>smoke=0</td> <td>[.84, .85]</td> <td>[.78, .79]</td> </tr> <tr> <td>smoke=1</td> <td>[.81, .83]</td> <td>[.70, .72]</td> </tr> </tbody> </table> <p style="text-align: center;">(Round to 2 digits).</p>                                      |            | age=25 | age=60 | smoke=0 | [.84, .85] | [.78, .79] | smoke=1 | [.81, .83] | [.70, .72] |
|  | age=25   | age=60     |        |        |         |            |            |         |            |            |
| smoke=0  | [.84, .85]   | [.78, .79] |        |        |         |            |            |         |            |            |
| smoke=1  | [.81, .83]   | [.70, .72] |        |        |         |            |            |         |            |            |

PROBLEM SET #2 STATISTICS 500 Fall 2017: ANSWER PAGE 2  
**This is an exam. Do not discuss it. Due noon Oct 24.**

|   |  |
|---|--|
|   | Fill in/ <b>CIRCLE</b> the answer  |
| 2.1 Who said: "To find out what happens to a system when you interfere with it you have to interfere with it"? Circle <b>ONE</b> .  | Gauss Cauchy Sheather<br><input checked="" type="radio"/> Box Cox Fox Cook<br>Fisher Markov Tukey<br>Harrell Everitt Weisberg<br>Ripley Titterington Trump   |
| 2.2 Add to model #2 <b>restricted cubic splines</b> in <b>age</b> with 5 knots. <b>See the data page.</b> Test that model 2 is adequate against the alternative that the splines in age are needed. Give the <b>name</b> of the test, the <b>value</b> of the test statistic, the <b>degrees of freedom (df)</b> , the <b>P-value</b> , and state whether the test <b>indicates curvature</b> in age. | Name: F-test Value:4.04<br>Degrees of Freedom: 3 & 2352<br><i>F-tests have degrees of freedom for numerator &amp; denominator</i><br>P-value: 0.007079<br><br>Indicates curvature in age?<br>CIRCLE ONE<br><br><input checked="" type="radio"/> YES <input type="radio"/> NO   |
| Use <b>model 2</b> for part 3.  | Fill in/ <b>CIRCLE</b> the answer  |
| 3.1 Which observation in model 2 has the largest <b>absolute</b> studentized residual? Give <b>signed</b> studentized residual, the row number and SEQN, and circle TRUE or FALSE about "This person is an outlier at the 0.05 level, taking account of the number of tests done." Give <b>DF</b> and adjusted <b>P-value</b> for test.   | Row: 1212 SEQN: 67220<br><br>Studentized Residual: -6.49<br>This person has the <b>lowest ratio</b> .<br>CIRCLE: TRUE <input checked="" type="radio"/> FALSE<br><b>Outlier:</b><br>CIRCLE: <input checked="" type="radio"/> TRUE FALSE<br>DF=degrees of freedom:<br>DF= 2354<br>Adusted P-value: 2.4591x10 <sup>-7</sup> |
| 3.2 Which observation has the largest <b>absolute</b> dffits? Give <b>signed</b> dffits, the row number and SEQN, and circle TRUE or FALSE. Does this observation move its yhat up 4.48 standard errors.  | Row: 1744 SEQN: 69468<br><br>dffits: -0.4484<br><br>CIRCLE: YES <input checked="" type="radio"/> NO  |
| 3.3 Which observation in model 2 has the largest <b>leverage</b> ? Give <b>hatvalue</b> , row # and SEQN.   | hatvalue: 0.013067<br><br>Row: 674 SEQN: 65071   |

## Doing the Problem Set in R

```
smoker<-1*(smokelung$smoke=="Daily")
smokerAge<-smoker*smokelung$age
d<-cbind(smokelung,smoker,smokerAge)
rm(smoker,smokerAge)
attach(d)
mini<-d[c(1781,1726,554,1394),]
mini<-mini[,c(8,9,18,19)]

m1<-lm(ratio~age+female+smoker)
m2<-lm(ratio~age+female+smoker+smokerAge)
library(car)
#1.1
residualPlots(m1)
residualPlots(m2)
#1.2
summary(m2)
#1.4 and 1.5
predict(m2,mini,interval="confidence")
#2.1 Last sentence of Box's Use and Abuse of Regression.
(Assigned reading.)
#2.2
library(Hmisc)
sp<-rcspline.eval(age)
m3<-lm(ratio~age+female+smoker+smokerAge+sp)
anova(m2,m3)
#For understanding
plot(age,m2$fitted.values)
plot(age,m3$fitted.values)
#3.1
outlierTest(m2)
which.max(abs(rstudent(m2)))
rstudent(m2)[1212]
d[1212,]
m2$df.residual-1
#3.2
which.max(abs(dffits(m2)))
dffits(m2)[1744]
d[1744,]
#3.3
which.max(hatvalues(m2))
hatvalues(m2)[674]
d[674,]
```

PROBLEM SET #3 STATISTICS 500 Fall 2017: DATA PAGE 1  
**Due at noon on Friday 15 December 2017, Statistics Dept.**

**This is an exam. Do not discuss it with anyone.**

The first data set expands the object **smokelung** in the course workspace. **You must download the workspace again.** The data are also available briefly as a csv file on my web page using the button data.csv. There are 2360 people, of whom 1842 never smoked and 518 are daily smokers, meaning that they smoked at least 5 cigarettes per day every day of the last 30 days. For an explanation of the lung function measures, fvc, fev1, and ratio = fev1/fvc, see [http://oac.med.jhmi.edu/res\\_phys/Encyclopedia/ForcedExpiration/ForcedExpiration.HTML](http://oac.med.jhmi.edu/res_phys/Encyclopedia/ForcedExpiration/ForcedExpiration.HTML) For an explanation of bmi, see [https://www.nhlbi.nih.gov/health/educational/lose\\_wt/BMI/bmicalc.htm](https://www.nhlbi.nih.gov/health/educational/lose_wt/BMI/bmicalc.htm) The variables are

**SEQN** = NHANES id number

**fvc** = forced vital capacity in ml

**fev1** = forced expiratory volume in 1 second, in ml

**ratio** = fev1/fvc

**cigsperday** = cigarettes smoked per day, 0 for never smokers

**smoke** = daily or never

**female** = 1 for female, 0 for male

**age** = age in years, >=20

**bmi** = body mass index

**educ** = education, 1=<9<sup>th</sup> grade, 2=9-11<sup>th</sup> grade, 3=high school or equivalent, 4="some college", say a 2 year associates degree, 5="college", >= 4 year BA degree.

**income** = ratio of family income to the poverty level, capped at 5xPoverty.

**You will need to construct several new variables.**

```
smoker<-1*(smokelung$smoke=="Daily")
```

```
smokerAge<-smoker*smokelung$age
```

```
d<-cbind(smokelung,smoker,smokerAge)
```

```
rm(smoker,smokerAge)
```

**Use randomhalf in d to split the data into two parts.**

**Use the random division in d; DO NOT make a new one.**

**One half is used for exploration, dex. The other half is used for validation, dva.**

```
dex<-d[d$randomhalf=="Explore",]
```

```
dva<-d[d$randomhalf=="Validate",]
```

```
dim(d)
```

```
[1] 2360    20
```

```
dim(dex)
```

```
[1] 1180    20
```

```
dim(dva)
```

```
[1] 1180    20
```



**This is an exam. Do not discuss it with anyone.**

**Due at noon on Friday 15 December 2017, Statistics Dept.**

**Model 1:**

$$\text{ratio} = \beta_0 + \beta_1 \text{age} + \beta_2 \text{female} + \beta_3 \text{smoker} + \beta_4 \text{smokerAge} \\ + \beta_5 \text{bmi} + \beta_6 \text{educ} + \beta_7 \text{income} + \varepsilon$$

where  $\varepsilon$  is iid  $N(0, \sigma^2)$ .

You will work with Model 1 and its  $2^7=128$  submodels (including Model 1 and the model with no predictors). In **question 1**, you use the exploratory random half, dex. In **question 2**, you check what you found using the validation half, dva. Spjotvoll's method corrects for multiple testing, so you don't have to split the data to use it. In **question 3**, you use all of the data in  $d$ , 2360 observations, and all  $2^7=128$  submodels in Spjotvoll's method. Use Spjotvoll's method at simultaneous level 0.05, so the chance of rejecting a true model as inadequate is at most 5% despite looking at  $2^7=128$  submodels.

The **second data set** is **nls500**. It is from the National Longitudinal Study (NLS). Decades ago, before charter schools, people used the NLS comparing Catholic high schools and public high schools. The data you have are simplified in several ways, so don't use these data to decide about the education of your children. There are three variables. **income50** indicates whether family income was  $<$  or  $\geq$  \$50,000. **school** is Catholic or public. **math** is the change in a math test score from before high school to the end of high school. Define **insec**  $\leftarrow$  **income50:school**.

```
> attach(nls500) > table(income50,school)
> insec <- income50:school > boxplot(math~insec)
```

**Important:** In **question 4**, refer to " $<50000$ :Catholic" as " $<C$ ", ..., " $\geq 50000$ :Public" as " $>P$ ".

**Model 2:**  $\text{math}_{ij} = \mu + \tau_j + \varepsilon_{ij}$ ,  $i=1, \dots, 163$ ,  $j=1, \dots, 4$ ,  $\varepsilon$  iid  $N(0, \sigma^2)$

**Important:** Write your name on both sides of the answer page, **last name first**. **Turn in only the answer page**. Brief answers suffice. **Circle** the correct answer. A circled answer may be correct or incorrect, but every crossed out answer is incorrect. **Turn in the exam** at noon on Friday 15 December 2017 in the Statistics Dept, 4<sup>th</sup> floor Huntsman. Give to Noelle at the front desk or place in an envelope addressed to me in my mailbox in Statistics. **Make and keep a photocopy of your answer page**. The answer key will be posted on-line in the revised bulk-pack.

**This is an exam. Do not discuss the exam with anyone.** If you discuss the exam, you have cheated on an exam. The single dumbest thing a PhD student at Penn can do is cheat on an exam.

**Have a great holiday!**

Name (**Last**, First): \_\_\_\_\_ ID# \_\_\_\_\_

PROBLEM SET #3 STATISTICS 500 Fall 2017: ANSWER PAGE 1

**This is an exam. Do not discuss it. Due noon December 15**

|   |  |
|---|--|
| <p>1. For questions in part 1, <b>use only the exploratory half of the data in dex</b>, n=1180. Do not use dva.</p>   | <p>FILL IN OR CIRCLE THE CORRECT ANSWER</p>  |
| <p>1.1 Fit model 1 to the 1180 observations in dex. Which predictor (x variable) has the largest variance-inflation-factor (VIF)? What is its numerical value? What is the <math>R^2</math> for this predictor as derived from VIF. CIRCLE True or False.</p> | <p>Predictor: _____<br/> VIF value: _____<br/> <math>R^2</math> Value: _____<br/> This VIF value means that this predictor is a poor predictor of y=ratio.<br/> CIRCLE ONE<br/> TRUE FALSE</p> |
| <p>1.2 When fitted to the 1180 observations in dex, which model P of the <math>2^7</math> submodels of model 1 has the smallest <math>C_p</math>? List the variables in P, give the value of <math>C_p</math>, the "size" of the model.</p>                   | <p>Predictors in P:<br/> <hr/> <math>C_p</math>=_____ size=_____</p>   |
| <p>1.3 <math>C_p</math> estimates of <math>J_p</math>. If the true <math>J_p</math> equaled its estimate, <math>C_p</math>, then the total expected squared error for the model in 1.2 would be more than twice as large as model 1 with 7 predictors.</p>    | <p>CIRCLE ONE<br/> TRUE FALSE</p>  |
| <p>1.4 The value of <math>C_p</math> in 1.2 suggests that the model in 1.2 omits at least one important predictor.</p>  | <p>CIRCLE ONE<br/> TRUE FALSE</p>  |

|   |                                     |
|---|-------------------------------------|
| <p><b>In 2, use dva, not dex.</b></p>   | <p>Fill in or CIRCLE the answer</p> |
| <p>2.1 Fit the model in 1.2 to the 1180 validation observations in dva. List the <b>predictors</b> and their <b>t-statistics</b>.</p> | <p>Variable names t-statistics</p>  |
| <p>2.2 Using dva, what is the correlation between the fitted values in 2.1 and the fitted values for model 1?</p>                     | <p>Correlation: _____</p>           |

Name (**Last**, First): \_\_\_\_\_ ID# \_\_\_\_\_

PROBLEM SET #3 STATISTICS 500 Fall 2017: ANSWER PAGE 2

**This is an exam. Do not discuss it. Due noon December 15**

|   |   |
|---|---|
| <b>In 3 use d, not dex nor dva.</b>   | Fill in or CIRCLE the answer  |
| <p>3.1 Use Spjotvoll's method and all 2360 observations in d to examine the <math>2^7=128</math> submodels of model 1. At the 0.05 level, how many of the 128 models are <b>not</b> judged "inadequate"? Which of these "not inadequate" models has the fewest predictors? (List them.)</p> | <p><b>How many?</b> _____<br/>         List the predictors in the one "not inadequate" model with the fewest predictors:<br/> <b>Predictor names:</b></p> |

| In 4, use the <b>nls500</b> data.   | Fill in or CIRCLE the answer  |    |    |    |    |    |        |  |  |  |  |        |  |  |  |  |              |  |  |  |  |
|---|---|----|----|----|----|----|--------|--|--|--|--|--------|--|--|--|--|--------------|--|--|--|--|
| <p>4.1 Under <b>model 2</b>, do a 4-group one-way anova of <b>math</b> score changes by <b>insc</b> group. Test the null hypothesis <math>H_0: \tau_1 = \tau_2 = \tau_3 = \tau_4 = 0</math>. Give the name of the test, the value of the test statistic, the P-value, degrees of freedom (DF) and circle an answer.</p> | <p>Name: _____ Value: _____<br/>         P-value: _____ DF: _____<br/>         CIRCLE ONE<br/>         H<sub>0</sub> IS<br/>         PLAUSIBLE NOT PLAUSIBLE</p>  |    |    |    |    |    |        |  |  |  |  |        |  |  |  |  |              |  |  |  |  |
| <p>4.2 Use Holm's method to test 6 hypotheses <math>H_0: \tau_j = \tau_k</math> controlling the familywise error rate at 0.05. Use the <b>notation from the data page</b> (e.g., "&lt;C") to indicate <b>pairs</b> of that differ significantly (eg "&lt;C, &gt;C")</p>   | <p>List all <b>pairs</b> of groups that differ significantly. If none, write "none".</p>  |    |    |    |    |    |        |  |  |  |  |        |  |  |  |  |              |  |  |  |  |
| <p>4.3 Give three orthogonal contrasts with integer weights for income (&lt;50,000, &gt;50,000), school (Catholic, Public) and their <b>interaction</b>.</p>  | <table border="1"> <thead> <tr> <th></th> <th>&lt;C</th> <th>&lt;P</th> <th>&gt;C</th> <th>&gt;P</th> </tr> </thead> <tbody> <tr> <td>Income</td> <td></td> <td></td> <td></td> <td></td> </tr> <tr> <td>School</td> <td></td> <td></td> <td></td> <td></td> </tr> <tr> <td>Inter-action</td> <td></td> <td></td> <td></td> <td></td> </tr> </tbody> </table> |    | <C | <P | >C | >P | Income |  |  |  |  | School |  |  |  |  | Inter-action |  |  |  |  |
|   | <C  | <P | >C | >P |    |    |        |  |  |  |  |        |  |  |  |  |              |  |  |  |  |
| Income  |   |    |    |    |    |    |        |  |  |  |  |        |  |  |  |  |              |  |  |  |  |
| School  |   |    |    |    |    |    |        |  |  |  |  |        |  |  |  |  |              |  |  |  |  |
| Inter-action  |   |    |    |    |    |    |        |  |  |  |  |        |  |  |  |  |              |  |  |  |  |
| <p>4.4 Test the hypothesis that the <b>interaction contrast</b> in the <math>\tau_j</math> does not differ significantly from zero.</p>   | <p>P-value: _____<br/>         CIRCLE ONE: H<sub>0</sub> IS<br/>         PLAUSIBLE NOT PLAUSIBLE</p>  |    |    |    |    |    |        |  |  |  |  |        |  |  |  |  |              |  |  |  |  |

**This is an exam. Do not discuss it.**

|   |  |
|---|--|
| <p>1. For questions in part 1, <b>use only the exploratory half of the data in dex</b>, n=1180. Do not use dva.</p>   | <p>FILL IN OR CIRCLE THE CORRECT ANSWER</p>  |
| <p>1.1 Fit model 1 to the 1180 observations in dex. Which predictor (x variable) has the largest variance-inflation-factor (VIF)? What is its numerical value? What is the <math>R^2</math> for this predictor as derived from VIF. CIRCLE True or False.</p> | <p>Predictor: smokerAge<br/> VIF value: 10.1410<br/> <math>R^2</math> Value: 0.90139<br/> This VIF value means that this predictor is a poor predictor of y=ratio.<br/> CIRCLE ONE<br/> TRUE      <b>FALSE</b></p> |
| <p>1.2 When fitted to the 1180 observations in dex, which model P of the <math>2^7</math> submodels of model 1 has the smallest <math>C_p</math>? List the variables in P, give the value of <math>C_p</math>, the "size" of the model.</p>                   | <p>Predictors in P:<br/> age, female, smokerAge</p> <hr/> <p><math>C_p=1.92</math> size=3</p>  |
| <p>1.3 <math>C_p</math> estimates of <math>J_p</math>. If the true <math>J_p</math> equaled its estimate, <math>C_p</math>, then the total expected squared error for the model in 1.2 would be more than twice as large as model 1 with 7 predictors.</p>    | <p>CIRCLE ONE<br/> TRUE      <b>FALSE</b></p>  |
| <p>1.4 The value of <math>C_p</math> in 1.2 suggests that the model in 1.2 omits at least one important predictor.</p>  | <p>CIRCLE ONE<br/> TRUE      <b>FALSE</b></p>  |

|   |   |                |              |     |        |        |      |           |        |
|---|---|----------------|--------------|-----|--------|--------|------|-----------|--------|
| <p><b>In 2, use dva, not dex.</b></p>   | <p>Fill in or CIRCLE the answer</p>   |                |              |     |        |        |      |           |        |
| <p>2.1 Fit the model in 1.2 to the 1180 validation observations in dva. List the <b>predictors</b> and their <b>t-statistics</b>.</p> | <table border="0"> <tr> <td>Variable names</td> <td>t-statistics</td> </tr> <tr> <td>age</td> <td>-16.47</td> </tr> <tr> <td>female</td> <td>4.36</td> </tr> <tr> <td>smokerAge</td> <td>-11.62</td> </tr> </table> | Variable names | t-statistics | age | -16.47 | female | 4.36 | smokerAge | -11.62 |
| Variable names  | t-statistics  |                |              |     |        |        |      |           |        |
| age   | -16.47  |                |              |     |        |        |      |           |        |
| female  | 4.36  |                |              |     |        |        |      |           |        |
| smokerAge   | -11.62  |                |              |     |        |        |      |           |        |
| <p>2.2 Using dva, what is the correlation between the fitted values in 2.1 and the fitted values for model 1?</p>                     | <p>Correlation: 0.9897</p>  |                |              |     |        |        |      |           |        |

**This is an exam. Do not discuss it.**

| <p><b>In 3 use d, not dex nor dva.</b></p>  | <p>Fill in or CIRCLE the answer</p>   |    |    |    |  |  |    |    |    |    |        |   |   |    |    |        |   |    |   |    |              |   |    |    |   |
|---|---|----|----|----|--|--|----|----|----|----|--------|---|---|----|----|--------|---|----|---|----|--------------|---|----|----|---|
| <p>3.1 Use Spjotvoll's method and all 2360 observations in d to examine the <math>2^7=128</math> submodels of model 1. At the 0.05 level, how many of the 128 models are <b>not</b> judged "inadequate"? Which of these "not inadequate" models has the fewest predictors? (List them.)</p>                             | <p><b>How many?</b> 16<br/>                 List the predictors in the one "not inadequate" model with the fewest predictors:<br/> <b>Predictor names:</b><br/>                 age, female, smokerAge<br/>                 (Every adequate model includes these three!)</p>  |    |    |    |  |  |    |    |    |    |        |   |   |    |    |        |   |    |   |    |              |   |    |    |   |
| <p>In 4, use the <b>nls500</b> data.</p>  | <p>Fill in or CIRCLE the answer</p>   |    |    |    |  |  |    |    |    |    |        |   |   |    |    |        |   |    |   |    |              |   |    |    |   |
| <p>4.1 Under <b>model 2</b>, do a 4-group one-way anova of <b>math</b> score changes by <b>insc</b> group. Test the null hypothesis <math>H_0: \tau_1 = \tau_2 = \tau_3 = \tau_4 = 0</math>. Give the name of the test, the value of the test statistic, the P-value, degrees of freedom (DF) and circle an answer.</p> | <p>Name: F-test Value: 5.085<br/>                 P-value: 0.00174 DF: 3, 648<br/>                 CIRCLE ONE<br/>                 H<sub>0</sub> IS<br/>                 PLAUSIBLE <input type="radio"/> NOT PLAUSIBLE <input checked="" type="radio"/></p>   |    |    |    |  |  |    |    |    |    |        |   |   |    |    |        |   |    |   |    |              |   |    |    |   |
| <p>4.2 Use Holm's method to test 6 hypotheses <math>H_0: \tau_j = \tau_k</math> controlling the familywise error rate at 0.05. Use the <b>notation from the data page</b> (e.g., "&lt;C") to indicate <b>pairs</b> of that differ significantly (eg "&lt;C, &gt;C")</p>   | <p>List all <b>pairs</b> of groups that differ significantly. If none, write "none".<br/>                 (&lt;C, &lt;P)<br/>                 (&gt;C, &lt;P)</p>  |    |    |    |  |  |    |    |    |    |        |   |   |    |    |        |   |    |   |    |              |   |    |    |   |
| <p>4.3 Give three orthogonal contrasts with integer weights for income (&lt;50,000, &gt;50,000), school (Catholic, Public) and their <b>interaction</b>.</p>  | <table border="1"> <thead> <tr> <th></th> <th>&lt;C</th> <th>&lt;P</th> <th>&gt;C</th> <th>&gt;P</th> </tr> </thead> <tbody> <tr> <td>Income</td> <td>1</td> <td>1</td> <td>-1</td> <td>-1</td> </tr> <tr> <td>School</td> <td>1</td> <td>-1</td> <td>1</td> <td>-1</td> </tr> <tr> <td>Inter-action</td> <td>1</td> <td>-1</td> <td>-1</td> <td>1</td> </tr> </tbody> </table> |    |    |    |  |  | <C | <P | >C | >P | Income | 1 | 1 | -1 | -1 | School | 1 | -1 | 1 | -1 | Inter-action | 1 | -1 | -1 | 1 |
|   | <C  | <P | >C | >P |  |  |    |    |    |    |        |   |   |    |    |        |   |    |   |    |              |   |    |    |   |
| Income  | 1   | 1  | -1 | -1 |  |  |    |    |    |    |        |   |   |    |    |        |   |    |   |    |              |   |    |    |   |
| School  | 1   | -1 | 1  | -1 |  |  |    |    |    |    |        |   |   |    |    |        |   |    |   |    |              |   |    |    |   |
| Inter-action  | 1   | -1 | -1 | 1  |  |  |    |    |    |    |        |   |   |    |    |        |   |    |   |    |              |   |    |    |   |
| <p>4.4 Test the hypothesis that the <b>interaction contrast</b> in the <math>\tau_j</math> does not differ significantly from zero.</p>   | <p>P-value: 0.31830<br/>                 CIRCLE ONE: H<sub>0</sub> IS<br/>                 PLAUSIBLE <input checked="" type="radio"/> NOT PLAUSIBLE <input type="radio"/></p>   |    |    |    |  |  |    |    |    |    |        |   |   |    |    |        |   |    |   |    |              |   |    |    |   |

### Doing the Problem Set in R: Problem 3, Fall 2017

```
smoker<-1*(smokelung$smoke=="Daily")
smokerAge<-smoker*smokelung$age
d<-cbind(smokelung,smoker,smokerAge)
rm(smoker,smokerAge)
dex<-d[d$randomhalf=="Explore",]
dva<-d[d$randomhalf=="Validate",]

attach(dex)
maxm<-lm(ratio~age+female+smoker+smokerAge+bmi+educ+income)

#1.1
library(car)
vif(maxm)
1-1/vif(maxm)

#1.2
x<-data.frame(age,female,smoker,smokerAge,bmi,educ,income)
library(leaps)
mex<-leaps(x=x,y=ratio,names=colnames(x))
cbind(mex$which,mex$Cp,mex$size)

#2.1
detach(dex)
attach(dva)
summary(lm(ratio~age+female+smokerAge))
#2.2
cor(lm(ratio~age+female+smokerAge)$fitted,lm(ratio~age+female+smoker+smokerAge+bmi+educ+income)$fitted)

#3
sp<-spjotvoll(x,ratio)
dim(sp[!sp$inadequate,])
sp[!sp$inadequate,]
```

```
#4.1
summary(aov(math~insc))

#4.2
pairwise.t.test(math,insc)

#4.3
cpschool<-c(1,-1,1,-1)
inc<-c(1,1,-1,-1)
interact<-cpschool*inc
contrasts(insc)<-cbind(cpschool,inc,interact)
contrasts(insc)

#4.4
x<-model.matrix(lm(math~insc))
head(x)
x<-as.data.frame(x)
m<-lm(math~x$insccpschool+x$inscinc+x$inscinteract)
summary(m)
anova(m)
```