

MAT 167: STATISTICS

FINAL EXAM

INSTRUCTOR: ANTHONY TANBAKUCHI

SPRING 2008

Name: _____

Computer / Seat Number: _____

No books, notes, or friends. **Show your work.** You may use the attached equation sheet, R, and a calculator. No other materials. If you choose to use R, write what you typed on the test or copy and paste your work into a word document labeling the question number it corresponds to. When you are done with the test print out the document. Be sure to save often on a memory stick just in case. Using any other program or having any other documents open on the computer will constitute cheating.

You have until the end of class to finish the exam, manage your time wisely.

If something is unclear quietly come up and ask me.

If the question is legitimate I will inform the whole class.

Express all final answers to 3 significant digits. Probabilities should be given as a decimal number unless a percent is requested. Circle final answers, ambiguous or multiple answers will not be accepted. Show steps where appropriate.

The exam consists of 12 questions for a total of 70 points on 9 pages.

This Exam is being given under the guidelines of our institution's **Code of Academic Ethics**. You are expected to respect those guidelines.

Points Earned: _____ out of 70 total points

Exam Score: _____

1. The following is a partial list of statistical methods that we have discussed:

- | | |
|---|---|
| 1. mean | 15. confidence interval for a proportion |
| 2. median | 16. confidence interval for difference in proportions |
| 3. mode | 17. one sample mean test |
| 4. standard deviation | 18. two independent sample mean test |
| 5. z-score | 19. match pair test |
| 6. percentile | 20. one sample proportion test |
| 7. coefficient of variation | 21. two sample proportion test |
| 8. scatter plot | 22. test of homogeneity |
| 9. histogram | 23. test of independence |
| 10. pareto chart | 24. linear correlation coefficient & test |
| 11. box plot | 25. regression |
| 12. normal-quantile plot | 26. 1-way ANOVA |
| 13. confidence interval for a mean | |
| 14. confidence interval for difference in means | |

For each situation below, which method is most applicable?

- If it's a hypothesis test, **also state what the null and alternative hypothesis are.**
 - If it's a graphical method, **also describe what you would be looking for.**
 - If it's a statistic, how susceptible to outliers is it?
- (a) (2 points) Ten pairs of chicks were selected to test the effect of a vitamin supplement on early growth. The chicks in each pair were siblings of high birth weight. One chick in each pair was given the supplement and the other was not. After two weeks, the weight of each chick was recorded. The researcher would like to test the research hypothesis that the supplement increases the growth rate of chicks in the first weeks after hatching against the null hypothesis that it has no effect.
- (b) (2 points) A researcher would like to compare the distribution of incomes of individuals in the four major regions of the US: the west, the south, the midwest, the east.
- (c) (2 points) A researcher wants determine if the mean income in the four major regions of

the US — the west, the south, the midwest, the east — are not all the same.

- (d) (2 points) An investigator is interested in the success of a job training program for current welfare recipients. If fewer than 30% of participants in the program are able to find work within three months, the program will be discontinued.
- (e) (2 points) A high school principal is interested in how well she can predict the number of days that her students miss school as a function of their GPA.
- (f) (2 points) A manufacturer needs to measure how consistently a new machine can cut nails to the desired length.
2. (1 point) 1-Way ANOVA can be thought of as a generalization of what two sample test?
3. (1 point) If the mean, median, and mode for a data set are all the same, what can you conclude about the data's distribution?
4. The following questions regard hypothesis testing in general.
- (a) (2 points) When we conduct a hypothesis test, we assume something is true and calculate the probability of observing the sample data under this assumption. What do we assume is true?

- (b) (2 points) Do we use the population distribution or the sampling distribution when calculating the p -value?
- (c) (2 points) If you reject H_0 but H_a is false. What type of error has occurred? (Type I or Type II)
- (d) (2 points) In the one sample proportion test, what is the distribution of the test statistic?
- (e) (2 points) In the one sample proportion test, what requirements must be met so that the test statistic's distribution is valid?
- (f) (2 points) Why is it important to use random sampling?
- (g) (2 points) A two sample mean hypothesis test was conducted with $H_0 : \Delta\mu = 0$ and $H_a : \Delta\mu > 0$. The first and second samples had respective sample sizes of 18 and 20. The test statistic calculated from the sample data is $t = 1.89$. Find the p -value.
5. Nine students were randomly selected who had taken the SAT twice. A researcher would like to test the claim that students who take the SAT test a second time score higher than their first test.

| Student | A | B | C | D | E | F | G | H | I |
|------------------|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| First SAT Score | 480 | 510 | 530 | 540 | 550 | 560 | 600 | 620 | 660 |
| Second SAT Score | 460 | 500 | 530 | 520 | 580 | 580 | 560 | 640 | 690 |

- (a) (1 point) What type of hypothesis test will you use?
- (b) (2 points) What are the test's requirements?

- (c) (2 points) What are the hypothesis H_0 and H_a ?
- (d) (1 point) What α will you use?
- (e) (2 points) Conduct the hypothesis test. What is the p -value?
- (f) (1 point) What is your formal decision?
- (g) (2 points) State your final conclusion in words.

6. The following table lists the the fuel consumption (in miles/gallon) and weight (in lbs) of a vehicle.

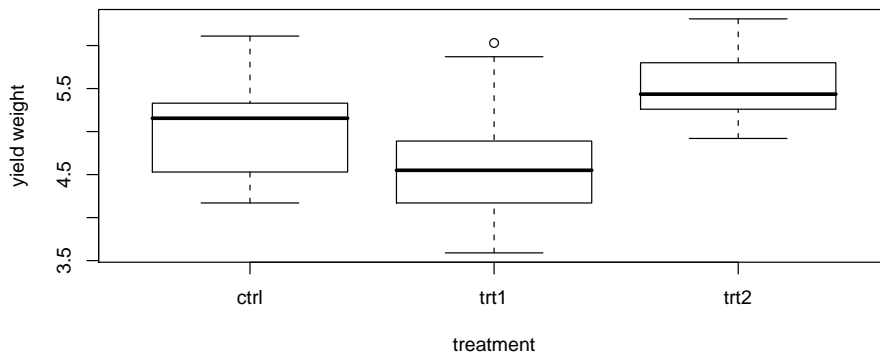
| | | | | | | | |
|--------|------|------|------|------|------|------|------|
| Weight | 3175 | 3450 | 3225 | 3985 | 2440 | 2500 | 2290 |
| MPG | 27 | 29 | 27 | 24 | 37 | 34 | 37 |

- (a) (2 points) Upon looking at the scatter plot of the data, the relationship of fuel consumption and milage looks linear. Is the linear relationship statistically significant? (**Justify your answer with an analysis.**)
- (b) (2 points) What percent of a vehicle's fuel consumption can be explained by its weight?
- (c) (2 points) You are designing a new vehicle and would like to be able to predict its fuel consumption. Write the equation for fitted model (with the actual values of the coeffi-

cients).

- (d) (2 points) What range of vehicle weights is the model valid for making predictions of fuel efficiency?
- (e) (2 points) What is the best predicted fuel consumption for a new vehicle that weights 2800 lbs?
- (f) (2 points) If the liner relationship had not been statistically significant, what would the best predicted fuel consumption for a new vehicle that weights 2800 lbs be?

7. Results from a randomized experiment to compare yields (as measured by dried weight of plants in grams) obtained under a control and two different treatment conditions are shown with a box plot of the data. The researcher who has developed the two new treatments hopes that at least one increases crop yield as compared to the control group.



- (a) (1 point) From the above box plot, the researcher notices that there do appear to be differences in the crop yield depending on the treatment. The researcher concludes from the box plots that the sample data supports the claim that the treatment type a plant receives affects the crop yield.
The researcher's thought process has a serious error in forming the conclusion. What has

the researcher forgotten to consider?

(b) (1 point) What type of hypothesis test should the researcher use to test her claim?

(c) (2 points) State the null and alternative hypothesis for this study **in words**.

(d) (1 point) If the researcher runs your suggested hypothesis test and the p -value is 0.18, what should her final conclusion be?

8. (2 points) You would like to conduct a study to estimate (at the 95% confidence level) the mean waist size of men with a margin of error of 1 in. Assuming that the standard deviation of waist sizes is $\sigma = 2.3$ in, what sample size should you use for this study?

9. (2 points) A random sample of 5 people was conducted to determine the mean length of index fingers. Below is the study data in inches.

3.2, 3.9, 3, 3.7, 3.7

Construct a 90% confidence interval for the true population mean index finger length using the

above data. (**Assume σ is unknown.**)

10. A bag of M&M's contains 18 red, 12 blue, 8 green, and 7 brown candies.

(a) (2 points) What is the probability of randomly selecting a red or brown M&M?

(b) (2 points) If 10 M&M's are randomly selected with replacement, what is the probability of getting exactly 4 green M&M's?

11. Engineers must consider the breadths of male heads when designing motorcycle helmets. Men have head breadths that are normally distributed with a mean of 6.0 in and a standard deviation of 1.0 in (based on anthropometric survey data from Gordon, Churchill, et al.).

(a) (2 points) If 1 man is randomly selected, find the probability that his head breadth is greater than 6.1 in.

(b) (2 points) If 100 men are randomly selected, find the probability that their mean head breadth is greater than 6.1 in.

12. (2 points) Given $x = \{4c, 2c, -2c\}$, where c is a constant, completely simplify the following expression:

$$\sqrt{\frac{\sum(x_i^2 - 2c)}{6c}}$$

End of exam. Reference sheets follow.

Statistics Quick Reference

Card & R Commands

by Anthony Tanbakuchi. Version 1.8.2
<http://www.tanbakuchi.com>
 ANTHONY@TANBAKUCHI.COM
 Get R at: <http://www.r-project.org>
 R commands: bold typewriter text

1 Misc R

To make a vector v store data: $x=c(x1, x2, \dots)$
 Help: general `RSiteSearch("Search Phrase")`
 Get: function `?functionName`
 Get column of data from table:
`tableName$columnName`
 List all variables: `ls()`
 Delete all variables: `rm(list=ls())`

$$\sqrt{x} = \text{sqrt}(x) \quad (1)$$

$$x^n = x^n \quad (2)$$

$$n = \text{length}(x) \quad (3)$$

$$T = \text{table}(x) \quad (4)$$

2 Descriptive Statistics

2.1 NUMERICAL

Let $x=c(x1, x2, x3, \dots)$

$$\text{total} = \sum_{i=1}^n x_i = \text{sum}(x) \quad (5)$$

$$\text{min} = \text{min}(x) \quad (6)$$

$$\text{max} = \text{max}(x) \quad (7)$$

$$\text{six number summary} = \text{summary}(x) \quad (8)$$

$$\mu = \frac{\sum x_i}{N} = \text{mean}(x) \quad (9)$$

$$\bar{x} = \frac{\sum x_i}{N} = \text{mean}(x) \quad (10)$$

$$\bar{x} = P_{50} = \text{median}(x) \quad (11)$$

$$\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}} \quad (12)$$

$$s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} = \text{sd}(x) \quad (13)$$

$$CV = \frac{\sigma}{\mu} = \frac{s}{\bar{x}} \quad (14)$$

2.2 RELATIVE STANDING

$$z = \frac{x - \mu}{\sigma} = \frac{x - \bar{x}}{s} \quad (15)$$

Percentiles:

$$P_k = x_i, \text{ (sorted } x) \quad (16)$$

$$k = \frac{i-0.5}{n} \cdot 100\%$$

To find x_i given P_k , i is:

- $L = (k/100)n$
- if L is an integer: $i = L + 0.5$; otherwise $i = L$ and round up.

2.3 VISUAL

All plots have optional arguments:

- `main=""` sets title
- `xlab=""`, `ylab=""` sets x/y-axis label
- `type="p"` for point plot
- `type="l"` for line plot
- `type="b"` for both points and lines

Ex: plot(x, y, type="b", main="My Plot")

Plot Types:

`hist(x)` histogram
`stem(x)` stem & leaf
`boxplot(x)` box plot
`plot(T)` bar plot, `T=table(x)`
`plot(x, y)` scatter plot, x, y are ordered vectors
`plot(t, y)` time series plot, t, y are ordered vectors
`curve(expr, xmin, xmax)` plot expr involving x

2.4 ASSESSING NORMALITY

Q-Q plot: `qqnorm(x)`; `qqline(x)`

3 Probability

Number of successes x with n possible outcomes. (Don't double count!)

$$P(A) = \frac{x}{n} \quad (17)$$

$$P(\bar{A}) = 1 - P(A) \quad (18)$$

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B) \quad (19)$$

$$P(A \text{ or } B) = P(A) + P(B) \text{ if } A, B \text{ mut. excl.} \quad (20)$$

$$P(A \text{ and } B) = P(A) \cdot P(B|A) \quad (21)$$

$$P(A \text{ and } B) = P(A) \cdot P(B) \text{ if } A, B \text{ independent} \quad (22)$$

$$n! = n(n-1) \cdots 1 = \text{factorial}(n) \quad (23)$$

$${}_n P_k = \frac{n!}{(n-k)!} \text{ Perm. no elem. alike} \quad (24)$$

$${}_n C_k = \frac{n!}{n!k! \cdots n_k!} \text{ Perm. } n_1 \text{ alike, } \dots \quad (25)$$

$${}_n C_k = \frac{n!}{(n-k)!k!} = \text{choose}(n, k) \quad (26)$$

4 Discrete Random Variables

$$P(x_i) : \text{probability distribution} \quad (27)$$

$$E = \mu = \sum x_i \cdot P(x_i) \quad (28)$$

$$\sigma = \sqrt{\sum (x_i - \mu)^2 \cdot P(x_i)} \quad (29)$$

4.1 BINOMIAL DISTRIBUTION

$$\mu = n \cdot p \quad (30)$$

$$\sigma = \sqrt{n \cdot p \cdot q} \quad (31)$$

$$P(x) = {}_n C_x p^x q^{n-x} = \text{dbinom}(x, n, p) \quad (32)$$

4.2 POISSON DISTRIBUTION

$$P(x) = \frac{\mu^x \cdot e^{-\mu}}{x!} = \text{dpois}(x, \mu) \quad (33)$$

5 Continuous random variables

CDF $F(x)$ gives area to the left of x , $F^{-1}(p)$ expects p is area to the left.

$$f(x) : \text{probability density} \quad (34)$$

$$E = \mu = \int_{-\infty}^{\infty} x \cdot f(x) dx \quad (35)$$

$$\sigma = \sqrt{\int_{-\infty}^{\infty} (x - \mu)^2 \cdot f(x) dx} \quad (36)$$

$$F(x) : \text{cumulative prob. density (CDF)} \quad (37)$$

$$F^{-1}(x) : \text{inv. cumulative prob. density} \quad (38)$$

$$F(x) = \int_{-\infty}^x f(x') dx' \quad (39)$$

$$p = P(x < x') = F(x') \quad (40)$$

$$x' = F^{-1}(p) \quad (41)$$

$$p = P(x > a) = 1 - F(a) \quad (42)$$

$$p = P(a < x < b) = F(b) - F(a) \quad (43)$$

5.1 UNIFORM DISTRIBUTION

$$p = P(a < u') = F(u') \quad (44)$$

$$u' = F^{-1}(p) = \text{qunif}(p, \text{min}=0, \text{max}=1) \quad (45)$$

5.2 NORMAL DISTRIBUTION

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} \frac{(x-\mu)^2}{\sigma^2}} \quad (46)$$

$$p = P(z < z') = F(z') = \text{pnorm}(z') \quad (47)$$

$$z' = F^{-1}(p) = \text{qnorm}(p) \quad (48)$$

$$p = P(x < x') = F(x') \quad (49)$$

$$= \text{pnorm}(x', \text{mean}=\mu, \text{sd}=\sigma) \quad (49)$$

$$x' = F^{-1}(p) \quad (50)$$

$$= \text{qnorm}(p, \text{mean}=\mu, \text{sd}=\sigma) \quad (50)$$

5.3 t-DISTRIBUTION

$$p = P(t < t') = F(t') = \text{pt}(t', \text{df}) \quad (51)$$

$$t' = F^{-1}(p) = \text{qt}(p, \text{df}) \quad (52)$$

5.4 χ^2 -DISTRIBUTION

$$p = P(\chi^2 < \chi'^2) = F(\chi'^2) \quad (53)$$

$$= \text{pchisq}(\chi'^2, \text{df}) \quad (53)$$

$$\chi'^2 = F^{-1}(p) = \text{qchisq}(p, \text{df}) \quad (54)$$

5.5 F-DISTRIBUTION

$$p = P(F < F') = F(F') \quad (55)$$

$$= \text{pf}(F', \text{df1}, \text{df2}) \quad (55)$$

$$F' = F^{-1}(p) = \text{qf}(p, \text{df1}, \text{df2}) \quad (56)$$

6 Sampling distributions

$$\mu_{\bar{x}} = \mu \quad \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \quad (57)$$

$$\mu_{\bar{p}} = p \quad \sigma_{\bar{p}} = \sqrt{\frac{pq}{n}} \quad (58)$$

7 Estimation

7.1 CONFIDENCE INTERVALS

$$\text{proportion: } \hat{p} \pm E, \quad E = z_{\alpha/2} \cdot \sigma_{\bar{p}} \quad (59)$$

$$\text{mean } (\sigma \text{ known): } \bar{x} \pm E, \quad E = z_{\alpha/2} \cdot \sigma_{\bar{x}} \quad (60)$$

$$\text{mean } (\sigma \text{ unknown, use } s): \bar{x} \pm E, \quad E = t_{\alpha/2} \cdot \sigma_{\bar{x}} \quad (61)$$

$$df = n - 1$$

$$\text{variance: } \frac{(n-1)s^2}{\chi^2_R} < \sigma^2 < \frac{(n-1)s^2}{\chi^2_L} \quad (62)$$

$$df = n - 1$$

$$2 \text{ proportions: } \Delta \hat{p} \pm z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}} \quad (63)$$

$$2 \text{ means (indep): } \Delta \bar{x} \pm t_{\alpha/2} \cdot \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \quad (64)$$

$$df \approx \text{min}(n_1 - 1, n_2 - 1)$$

$$\text{matched pairs: } \bar{d} \pm t_{\alpha/2} \cdot \frac{s_d}{\sqrt{n}}, \quad d_i = x_i - y_i, \quad (65)$$

$$df = n - 1$$

7.2 CI CRITICAL VALUES (TWO SIDED)

$$z_{\alpha/2} = F_z^{-1}(1 - \alpha/2) = \text{qnorm}(1 - \alpha/2) \quad (66)$$

$$t_{\alpha/2} = F_t^{-1}(1 - \alpha/2) = \text{qt}(1 - \alpha/2, \text{df}) \quad (67)$$

$$\chi^2_L = F_{\chi^2}^{-1}(\alpha/2) = \text{qchisq}(\alpha/2, \text{df}) \quad (68)$$

$$\chi^2_R = F_{\chi^2}^{-1}(1 - \alpha/2) = \text{qchisq}(1 - \alpha/2, \text{df}) \quad (69)$$

7.3 REQUIRED SAMPLE SIZE

$$\text{proportion: } n = \hat{p} \hat{q} \left(\frac{z_{\alpha/2}}{E} \right)^2, \quad (70)$$

$$(\hat{p} = \hat{q} = 0.5 \text{ if unknown})$$

$$\text{mean: } n = \left(\frac{z_{\alpha/2} \cdot \hat{\sigma}}{E} \right)^2 \quad (71)$$

8 Hypothesis Tests

Test statistic and R function (when available) are listed for each.

Optional arguments for hypothesis tests:

alternative="two.sided" can be:
"two.sided", "less", "greater"

conf.level=0.95 constructs a 95% confidence interval. Standard CI only when **alternative="two.sided"**.

Optional arguments for power calculations & Type II error:

alternative="two.sided" can be:
"two.sided" or "one.sided"

sig.level=0.05 sets the significance level α .

8.1 1-SAMPLE PROPORTION

$H_0: p = p_0$

prop.test(x, n, p=p₀, alternative="two.sided")

$$z = \frac{\hat{p} - p_0}{\sqrt{p_0q_0/n}} \quad (72)$$

8.2 1-SAMPLE MEAN (σ KNOWN)

$H_0: \mu = \mu_0$

$$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} \quad (73)$$

8.3 1-SAMPLE MEAN (σ UNKNOWN)

$H_0: \mu = \mu_0$

t.test(x, mu=μ₀, alternative="two.sided")

Where **x** is a vector of sample data.

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}, \quad df = n - 1 \quad (74)$$

Required Sample size:

power.t.test(delta=h, sd=σ, sig.level=α, power=1 - β, type="one.sample", alternative="two.sided")

8.4 2-SAMPLE PROPORTION TEST

$H_0: p_1 = p_2$ or equivalently $H_0: \Delta p = 0$

prop.test(x, n, alternative="two.sided")

where: **x=c(x₁, x₂)** and **n=c(n₁, n₂)**

$$z = \frac{\hat{\Delta p} - \Delta p_0}{\sqrt{\frac{\hat{p}_1\hat{q}_1}{n_1} + \frac{\hat{p}_2\hat{q}_2}{n_2}}}, \quad \hat{\Delta p} = \hat{p}_1 - \hat{p}_2 \quad (75)$$

$$\hat{p} = \frac{x_1 + x_2}{n_1 + n_2}, \quad \hat{q} = 1 - \hat{p} \quad (76)$$

Required Sample size:

power.prop.test(p1=p₁, p2=p₂, power=1 - β, sig.level=α, alternative="two.sided")

8.5 2-SAMPLE MEAN TEST

$H_0: \mu_1 = \mu_2$ or equivalently $H_0: \Delta \mu = 0$

t.test(x1, x2, alternative="two.sided")

where: **x1** and **x2** are vectors of sample 1 and sample 2 data.

$$t = \frac{\Delta \bar{x} - \Delta \mu_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}, \quad df \approx \min(n_1 - 1, n_2 - 1), \quad \Delta \bar{x} = \bar{x}_1 - \bar{x}_2 \quad (77)$$

Required Sample size:

power.t.test(delta=h, sd=σ, sig.level=α, power=1 - β, type="two.sample", alternative="two.sided")

8.6 2-SAMPLE MATCHED PAIRS TEST

$H_0: \mu_d = 0$

t.test(x, y, paired=TRUE, alternative="two.sided")

where: **x** and **y** are ordered vectors of sample 1 and sample 2 data.

$$t = \frac{\bar{d} - \mu_{d0}}{s_d/\sqrt{n}}, \quad d_i = x_i - y_i, \quad df = n - 1 \quad (78)$$

Required Sample size:

power.t.test(delta=h, sd=σ, sig.level=α, power=1 - β, type="paired", alternative="two.sided")

8.7 TEST OF HOMOGENEITY, TEST OF INDEPENDENCE

$H_0: p_1 = p_2 = \dots = p_k$ (homogeneity)

$H_0: X$ and Y are independent (independence)

chisq.test(D)

Enter table: **D=matrix(c1, c2, ...)**, where **c1, c2, ...** are column data vectors.

Or generate table: **D=table(x1, x2)**, where **x1, x2** are ordered vectors of raw categorical data.

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}, \quad df = (\text{num rows} - 1)(\text{num cols} - 1) \quad (79)$$

$$E_i = \frac{(\text{row total})(\text{column total})}{(\text{grand total})} = np_i \quad (80)$$

For 2×2 contingency tables, you can use the Fisher Exact Test:

fisher.test(D, alternative="greater")

(must specify alternative as greater)

9 Linear Regression

9.1 LINEAR CORRELATION

$H_0: \rho = 0$

cor.test(x, y)

where: **x** and **y** are ordered vectors.

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{(n-1)s_x s_y}, \quad t = \frac{r - 0}{\sqrt{\frac{1-r^2}{n-2}}}, \quad df = n - 2 \quad (81)$$

9.2 MODELS IN R

| MODEL TYPE | EQUATION | R MODEL |
|---------------------|--|--------------------------------|
| linear 1 indep var | $y = b_0 + b_1x_1$ | $y \sim x_1$ |
| ... 0 intercept | $y = 0 + b_1x_1$ | $y \sim 0 + x_1$ |
| linear 2 indep vars | $y = b_0 + b_1x_1 + b_2x_2$ | $y \sim x_1 + x_2$ |
| ... interaction | $y = b_0 + b_1x_1 + b_2x_2 + b_{12}x_1x_2$ | $y \sim x_1 + x_2 + x_1 * x_2$ |
| polynomial | $y = b_0 + b_1x_1 + b_2x_1^2$ | $y \sim x_1 + 1(x_1^2)^2$ |

9.3 REGRESSION

Simple linear regression steps:

1. Make sure there is a significant linear correlation.
2. **results=lm(y~x)** Linear regression of y on x vectors
3. **results** View the results
4. **plot(x, y); abline(results)** Plot regression line on data
5. **plot(x, results\$residuals)** Plot residuals

$$y = b_0 + b_1x_1 \quad (82)$$

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad (83)$$

$$b_0 = \bar{y} - b_1\bar{x} \quad (84)$$

9.4 PREDICTION INTERVALS

To predict y when $x = 5$ and show the 95% prediction interval with regression model in results:

predict(results, newdata=data.frame(x=5), int="pred")

10 ANOVA

10.1 ONE WAY ANOVA

1. **results=aov(depVarColName~indepVarColName, data=tableName)** Run ANOVA with data in **tableName**, factor data in **indepVarColName** column, and response data in **depVarColName** column.
2. **summary(results)** Summarize results
3. **boxplot(depVarColName~indepVarColName, data=tableName)** Boxplot of levels for factor

$$F = \frac{MS(\text{treatment})}{MS(\text{error})}, \quad df_1 = k - 1, df_2 = N - k \quad (85)$$

To find required sample size and power see **power.anova.test(...)**

11 Loading and using external data and tables

11.1 LOADING EXCEL DATA

1. Export your table as a CSV file (comma separated file) from Excel.
2. Import your table into R using:
MyTable=read.csv(file.choose())

11.2 LOADING AN .RDATA FILE

You can either double click on the .RData file or use the menu:

- Windows: **File→Load Workspace...**
- Mac: **Workspace→Load Workspace File...**

11.3 USING TABLES OF DATA

1. To see all the available variables type: **ls()**
2. To see what's inside a variable, type its name.
3. If the variable **tableName** is a table, you can also type **names(tableName)** to see the column names or type **head(tableName)** to see the first few rows of data.
4. To access a column of data type **tableName\$columnName**

An example demonstrating how to get the women's height data and find the mean:

```
> ls() # See what variables are defined
[1] "women" "x"
> head(women) # Look at the first few entries
  height weight
1    58   115
2    59   117
3    60   120
> names(women) # Just get the column names
[1] "height" "weight"
> women$height # Display the height data
[1] 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72
> mean(women$height) # Find the mean of the heights
[1] 65
```