

MAT 167: STATISTICS

TEST II: CHAPTERS 4-7

INSTRUCTOR: ANTHONY TANBAKUCHI

FALL 2007

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, 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 4 questions for a total of 25 points on 7 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 25 total points

**Exam Score:** \_\_\_\_\_



2. The clothing manufacturer's association (CMA) publishes data that manufacture's use to determine what sizes of clothing they should make. As mentioned before, the CMA states that men's waists are normally distributed with  $\mu = 35$  in and  $\sigma = 2.3$  in. Lately, you are getting a lot of returns on your one size fit's all sweat pants (that you designed in a previous question) because they are too small.

(a) (1 point) 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?

(b) (1 point) A study was conducted (and they ignored your recommendation of sample size!) of 5 randomly selected men and the following waist sizes were measured:

36.8 38.3 37.8 40 38.3

Construct a 95% confidence interval for the true population mean waist size using the above data. (**Assume  $\sigma$  is unknown.**)

(c) (1 point) Why is the margin of error in your calculated confidence interval greater than our original desired margin of error of 1 in?

3. As mentioned before, the CMA states that men's waists are normally distributed with  $\mu = 35$  in. Lately, you are getting allot of returns on your one size fit's all sweat pants (that you designed in a previous question) because they are too small.

You believe that the mean waist size of men is actually greater than 35 in. Using the same study data from the previous question of 5 randomly sampled men (shown below again), conduct a hypothesis test to test your claim. Use a significance level of 0.01 and **assume  $\sigma$  is unknown**.

36.8 38.3 37.8 40 38.3

- (a) (1 point) What type of hypothesis test will you use?
- (b) (1 point) What are the test's requirements?
- (c) (1 point) Are the requirements satisfied? **State how they are satisfied.**
- (d) (1 point) What are the hypothesis?
- (e) (1 point) What  $\alpha$  will you use?
- (f) (1 point) Conduct the hypothesis test. What is the  $p$ -value?
- (g) (1 point) What is your formal decision?

- (h) (1 point) State your final conclusion in words.
- (i) (1 point) What is the *actual* probability of a Type I error for this study data?
- (j) (1 point) Assume a Type I error has occurred, state what the formal decision was and what the error is.
- (k) (1 point) If the researcher had an  $\alpha$  of 0.005 and failed reject  $H_0$ , what have they proven?
- (l) (1 point) Assume a Type II error has occurred, state what the formal decision was and what the error is.
- (m) (1 point) In general, 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?

4. A newly married couple does not want to get pregnant during their first year of marriage and decides to only use condoms as a contraceptive device. Studies<sup>1</sup> have shown that the probability of pregnancy per use of a condom is 0.15% (individual probability of pregnancy). During the first year the couple use a condom every time for a total of 162 times.
- (a) (1 point) Find the probability of 0 pregnancies over 1 year.
  
  
  
  
  
  
  
  
  
  
  - (b) (1 point) Find the probability of 1 or more pregnancies over 1 year.
  
  
  
  
  
  
  
  
  
  
  - (c) (1 point) Would it be unusual to become pregnant during 1 year of use when only using this form of contraceptive?
  
  
  
  
  
  
  
  
  
  
  - (d) (1 point) What would be the average number of pregnancies that the couple should expect to have given the above information?
  
  
  
  
  
  
  
  
  
  
  - (e) (1 point) Using the equation for the mean of the binomial distribution, solve for the number of uses  $n$  required to have an average (mean) of 1 pregnancy.
  
  
  
  
  
  
  
  
  
  
  - (f) (1 point) Based on your above results, if you were the primary care physician for this couple, what would you tell them about their plan to prevent pregnancy? Would it be effective?

---

<sup>1</sup>The data presented in this problem are approximate values derived from actual published data based on *typical* use effectiveness.

# Basic Statistics: Quick Reference & R Commands

by Anthony Tambakuchi. Version 1.4  
<http://www.tambakuchi.com>  
 ANTHONY@TANBAKUCHI.COM  
 Get R at: <http://www.r-project.org>  
 R commands: bold typewriter text

## 1 Misc R

To make a vector / store data: `x=c(x1, x2, ...)`  
 Get help on function: `?functionName`  
 Get column of data from table: `tableName$columnName`  
 List all variables: `ls()`  
 Delete all variables: `rm(list=ls())`

$$\begin{aligned} \sqrt{x} &= \text{sqrt}(x) & (1) \\ x^n &= x^n & (2) \\ n &= \text{length}(x) & (3) \\ T &= \text{table}(x) & (4) \end{aligned}$$

## 2 Descriptive Statistics

### 2.1 NUMERICAL

Let `x=c(x1, x2, x3, ...)`

$$\begin{aligned} \text{total} &= \sum_{i=1}^n x_i = \text{sum}(x) & (5) \\ \text{min} &= \text{min}(x) & (6) \\ \text{max} &= \text{max}(x) & (7) \\ \text{six number summary} &: \text{summary}(x) & (8) \\ \bar{\mu} &= \frac{\sum x_i}{N} = \text{mean}(x) & (9) \\ \bar{x} &= \frac{\sum x_i}{n} = \text{mean}(x) & (10) \\ \bar{x} &= P_{50} = \text{median}(x) & (11) \\ \sigma &= \sqrt{\frac{\sum (x_i - \bar{\mu})^2}{N}} & (12) \\ s &= \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} = \text{sd}(x) & (13) \\ CV &= \frac{\sigma}{\bar{\mu}} = \frac{s}{\bar{x}} & (14) \\ z &= \frac{x - \bar{x}}{\sigma} = \frac{x - \bar{x}}{s} & (15) \end{aligned}$$

Percentiles:

$$\begin{aligned} P_k &= x_{(k)} \quad (\text{sorted } x) & (16) \\ k &= \frac{i-0.5}{n} \cdot 100\% \end{aligned}$$

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!)

$$\begin{aligned} P(A) &= \frac{x}{n} & (17) \\ P(\bar{A}) &= 1 - P(A) & (18) \\ P(A \text{ or } B) &= P(A) + P(B) - P(A \text{ and } B) & (19) \\ P(A \text{ or } B) &= P(A) + P(B) \quad \text{if } A, B \text{ mut. excl.} & (20) \\ P(A \text{ and } B) &= P(A) \cdot P(B|A) & (21) \\ P(A \text{ and } B) &= P(A) \cdot P(B) \quad \text{if } A, B \text{ independent} & (22) \\ n! &= n(n-1) \cdots 1 = \text{factorial}(n) & (23) \\ n!_k &= \frac{n!}{(n-k)!} \quad \text{Perm. no elem. alike} & (24) \\ n!_k &= \frac{n!}{n_1! n_2! \cdots n_k!} \quad \text{Perm. } n_1 \text{ alike, } \dots & (25) \\ nC_k &= \frac{n!}{(n-k)! k!} = \text{choose}(n, k) & (26) \end{aligned}$$

## 4 Discrete Random Variables

$$\begin{aligned} P(x_i) &: \text{probability distribution} & (27) \\ E = \mu &= \sum x_i \cdot P(x_i) & (28) \\ \sigma &= \sqrt{\sum (x_i - \mu)^2 \cdot P(x_i)} & (29) \end{aligned}$$

### 4.1 BINOMIAL DISTRIBUTION

$$\begin{aligned} \mu &= n \cdot p & (30) \\ \sigma &= \sqrt{n \cdot p \cdot q} & (31) \\ P(x) &= {}^n C_x p^x q^{n-x} = \text{dbinom}(x, n, p) & (32) \end{aligned}$$

### 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.

$$\begin{aligned} f(x) &: \text{probability density} & (34) \\ E = \mu &= \int_{-\infty}^{\infty} x \cdot f(x) dx & (35) \\ \sigma &= \sqrt{\int_{-\infty}^{\infty} (x - \mu)^2 \cdot f(x) dx} & (36) \\ F(x) &: \text{cumulative prob. density (CDF)} & (37) \\ F^{-1}(p) &: \text{inv. cumulative prob. density} & (38) \\ F(x) &= \int_{-\infty}^x f(x') dx' & (39) \\ p &= P(x < x') = F(x') & (40) \\ x' &= F^{-1}(p) & (41) \\ p &= P(x > a) = 1 - F(a) & (42) \\ p &= P(a < x < b) = F(b) - F(a) & (43) \end{aligned}$$

### 5.1 UNIFORM DISTRIBUTION

$$\begin{aligned} p &= P(a < u') = F(u') \\ &= \text{punif}(u', \text{min}=0, \text{max}=1) & (44) \\ u' &= F^{-1}(p) = \text{qunif}(p, \text{min}=0, \text{max}=1) & (45) \end{aligned}$$

### 5.2 NORMAL DISTRIBUTION

$$\begin{aligned} f(x) &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} & (46) \\ p &= P(z < z') = F(z') = \text{pnorm}(z') & (47) \\ z' &= F^{-1}(p) = \text{qnorm}(p) & (48) \\ p &= P(x < x') = F(x') \\ &= \text{pnorm}(x', \text{mean}=\mu, \text{sd}=\sigma) & (49) \\ x' &= F^{-1}(p) \\ &= \text{qnorm}(p, \text{mean}=\mu, \text{sd}=\sigma) & (50) \end{aligned}$$

### 5.3 t-DISTRIBUTION

$$\begin{aligned} p &= P(t < t') = F(t') = \text{pt}(t', \text{df}) & (51) \\ p &= F^{-1}(p) = \text{qt}(p, \text{df}) & (52) \end{aligned}$$

### 5.4 $\chi^2$ -DISTRIBUTION

$$\begin{aligned} p &= P(\chi^2 < \chi'^2) = F(\chi'^2) \\ &= \text{pchisq}(\chi'^2, \text{df}) & (53) \\ \chi'^2 &= F^{-1}(p) = \text{qchisq}(p, \text{df}) & (54) \end{aligned}$$

### 5.5 F-DISTRIBUTION

$$\begin{aligned} p &= P(F < F') = F(F') \\ &= \text{pf}(F', \text{df1}, \text{df2}) & (55) \\ F' &= F^{-1}(p) = \text{qf}(p, \text{df1}, \text{df2}) & (56) \end{aligned}$$

## 6 Sampling distributions

$$\mu = \mu \quad \sigma_x = \frac{\sigma}{\sqrt{n}} \quad (57)$$

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

## 7 Estimation

### 7.1 CONFIDENCE INTERVALS

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

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

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

$$df = n - 1$$

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

$$df = n - 1$$

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

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

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

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

$$df = n - 1$$

### 7.2 CI CRITICAL VALUES (TWO SIDED)

$$z_{\alpha/2} = P(z > \alpha) = \text{qnorm}(1 - \alpha/2) \quad (66)$$

$$t_{\alpha/2} = P(t > \alpha) = \text{qt}(1 - \alpha/2, \text{df}) \quad (67)$$

$$\chi^2_L = P(\chi^2 < \alpha) = \text{qchisq}(\alpha/2, \text{df}) \quad (68)$$

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

### 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 Test

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

Optional arguments:

**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"**.

### 8.1 1-SAMPLE PROPORTION

$H_0: p = p_0$

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

$$z = \frac{\hat{p} - p_0}{\sqrt{p_0 q_0 / n}} \quad (72)$$

### 8.2 1-SAMPLE MEAN ( $\sigma$ KNOWN)

$H_0: \mu = \mu_0$

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

### 8.3 1-SAMPLE MEAN ( $\sigma$ UNKNOWN)

$H_0: \mu = \mu_0$

**t.test(x, mu=mu0, alternative="two.sided")**

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

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

### 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)**

$$z = \frac{\Delta \hat{p} - \Delta p_0}{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}}, \quad \Delta \hat{p} = \hat{p}_1 - \hat{p}_2 \quad (75)$$

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

### 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.

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

### 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.

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

### 8.7 TEST OF HOMOGENEITY, TEST OF INDEPENDENCE

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

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

### chisq.test(D)

Enter table: **D=data.frame(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  $2 \times 2$  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.

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{(n-1)s_x s_y}, \quad t = \frac{r - \rho_0}{\sqrt{\frac{1-r^2}{n-2}}}, \quad df = n - 2 \quad (81)$$

### 9.2 MODELS IN R

MODEL TYPE	EQUATION	R MODEL
linear 1 indep var	$y = b_0 + b_1 x_1$	$y \sim x_1$
... 0 intercept	$y = 0 + b_1 x_1$	$y \sim 0 + x_1$
linear 2 indep vars	$y = b_0 + b_1 x_1 + b_2 x_2$	$y \sim x_1 + x_2$
... interaction	$y = b_0 + b_1 x_1 + b_2 x_2 + b_{12} x_1 x_2$	$y \sim x_1 + x_2 + x_1 * x_2$

### 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

$$y = b_0 + b_1 x_1 \quad (82)$$

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

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

### 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

## 11 Loading external data

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