

SOLUTIONS  
MAT 167: STATISTICS

TEST I: CHAPTERS 1-4

INSTRUCTOR: ANTHONY TANBAKUCHI

FALL 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. Write your work in the provided space for each problem (including any R work if appropriate). You may not use personal computers, only use the classroom computer at your desk. 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 8 questions for a total of 40 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.

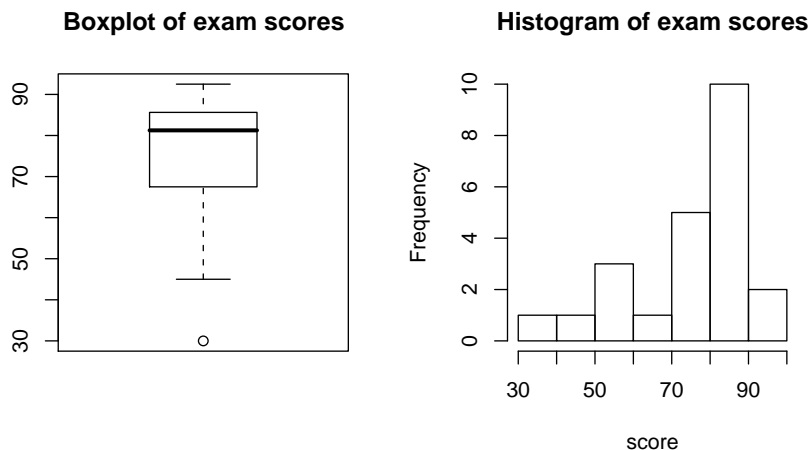
**Points Earned:** \_\_\_\_\_ out of 40 total points

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

**Solution:**

Exam Results. Below are the summary statistics.

```
> summary(score)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
30.00  67.50   81.25   75.11  85.62   92.50
> par(mfrow = c(1, 2))
> boxplot(score, main = "Boxplot of exam scores")
> hist(score, main = "Histogram of exam scores")
```



1. Provide **short succinct** written answers to the following conceptual questions.

- (a) (1 point) Would ethnicity be classified as a nominal, ordinal, interval, or ratio level of measurement?

**Solution:** Ethnicity would be categorical, with no meaningful order, therefore we would classify it as nominal.

- (b) (1 point) Which of the following measures of center is least susceptible to outliers:  
**median, mean, midrange, mode**

**Solution:** The mode is the least susceptible.

- (c) (1 point) What percent of data is greater than  $Q_1$ ?

**Solution:** 75%

- (d) (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?

**Solution:** If all three measures of center are the same, the distribution is symmetrical. (Not necessarily a normal distribution, all we know is that it is symmetrical.)

- (e) (1 point) If the mean is less than the mode for a data set, what can you conclude about the data's distribution?

**Solution:** The data set is negatively skewed.

- (f) (1 point) What does the standard deviation represent conceptually **in words**? (Be concise but don't simply state the equation in words verbatim.)

**Solution:** The standard deviation represents the average variation of the data from the mean.

- (g) (1 point) To determine the proportion of support in Tucson for Vice Presidential candidate Palin, a researcher randomly samples 10 people (using an ideal simple random sample). From the sample data, 70% expressed their support. The researcher was suspicious of the result so she repeated the study randomly sampling 10 people 3 more times, the resultant statistics were 30% 70% and 0% respectively. The researcher thinks some sort of mistake has occurred causing the numbers to change each time the study is repeated. However, after carefully checking the procedures used, no mistake occurred. What type of error could we attribute to the variation between the sample results?

**Solution:** This result should not be surprising, sampling error can cause fluctuations in the statistic from sample to sample.

- (h) (2 points) A histogram is a useful tool that can quickly communicate many traits about a set of data. List 4 useful pieces of information that an observer can easily assess using a histogram.

**Solution:** A histogram can be used to get an approximation of:

1. central tendency
2. variation in the data
3. shape of the data
4. assess if outliers exist
5. min
6. max

- (i) (1 point) Your SAT results state that you scored in the 95<sup>th</sup> percentile. What does this mean?

**Solution:** This is a good score, it means that 95% of the the people taking the test had a lower score than you.

- (j) (1 point) Why would a SAT percentile be preferred over a raw SAT score for college admissions committees?

**Solution:** The percentile compares how the applicant did to their peers who took the test (a measure of relative standing). A raw score doesn't give information as to how this score compared to others taking the test, making it hard to determine if a 1100 is easy or hard to get.

2. A survey conducted in our class asked 24 students how many credits they were enrolled in this semester. Use the R output below to answer the following questions.

There are 24 data points stored in the variable  $x$ , below is the sorted data:

```
> sort(x)
[1] 6 10 10 10 11 11 11 12 12 12 12 12 12 13 13 13 14 15 15 16 16 16 17 19
```

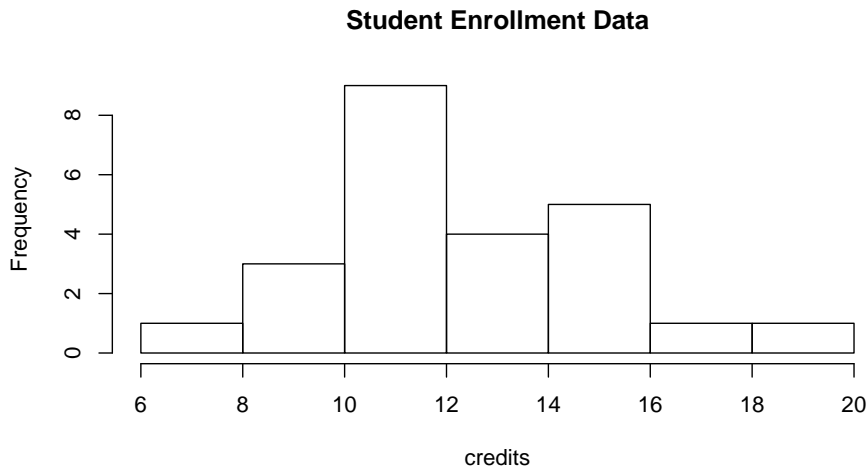
The basic descriptive statistical analysis is as follows:

```
> summary(x)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
  6.00   11.00   12.00   12.83   15.00   19.00

> var(x)
[1] 7.884058

> sd(x)
[1] 2.807856
```

```
> hist(x, xlab = "credits", main = "Student Enrollment Data")
```



(a) (1 point) Use the range rule of thumb to estimate the standard deviation. Is it close to the actual standard deviation?

**Solution:**

```
> s.est = (max(x) - min(x))/4
> signif(s.est, 3)
[1] 3.25
```

The range rule is a rough estimate of  $\sigma$ , it is relatively close to the actual standard deviation shown in the output.

- (b) (1 point) What is  $P_{25}$  equal to?

**Solution:** You can find the value in the summary output for the 1st quarter, it is equal to 11.

- (c) (1 point) What is the IQR (inter quartile range) equal to?

**Solution:**  $IQR = Q_3 - Q_1$ , using the summary output, the IQR is 4.

- (d) (1 point) What percent of the data falls within the IQR?

**Solution:** The IQR contains the data between  $Q_1$  and  $Q_3$  or equivalently  $P_{25}$  and  $P_{75}$ , thus 50% of the data falls inside the IQR.

- (e) (1 point) What is the percentile for the student who is taking 14 credits?

**Solution:** The percentile is the percent of data points less than the given data point. The student taking 14 credits is at position  $i = 17$ , thus:

$$\begin{aligned} P_k &= x_i, \quad (\text{sorted } x) \\ k &= \frac{i - 0.5}{n} \cdot 100\% \\ &= \frac{17 - 0.5}{24} \cdot 100\% \\ &= 68.8\% \end{aligned}$$

A quick approximation would be  $16/24 \cdot 100\%$  since there are 16 data points less than 14. However,  $17/24 \cdot 100\%$  would not be correct.

- (f) (1 point) What is the z-score for the student who is taking 10 credits?

**Solution:**

```
> x.bar
[1] 12.83333
> s
[1] 2.807856
> z = (10 - x.bar)/s
> signif(z, 3)
[1] -1.01
```

- (g) (1 point) Is 6 credits an unusual (outlier) value based on its  $z$  score? (Why)

**Solution:**

```
> x.bar
[1] 12.83333
> s
[1] 2.807856
> z = (6 - x.bar)/s
> signif(z, 3)
[1] -2.43

Yes, since  $|z| > 2$ .
```

- (h) (1 point) Is the data positively skewed, negatively skewed, or symmetrical?

**Solution:** Since the histogram is close to symmetrical, it's difficult to tell by just looking at it. However since the mean  $>$  median we know that it is positively skewed.

- (i) (1 point) Construct an interval using the Empirical Rule which you would expect 99.7% of the data to fall within.

**Solution:** 99.7% of the data falls within  $\mu \pm 3\sigma$  which for this data set is:  $12.8 \pm 8.42 = (4.41, 21.3)$ .

3. A researcher conducts a study in which 1,000 individuals in Tucson are randomly selected and asked if they prefer orange juice over grape juice. Eighty percent of the respondents preferred orange juice. The researcher concludes that "Study data indicates that eighty percent of Americans prefer orange juice over grape juice."

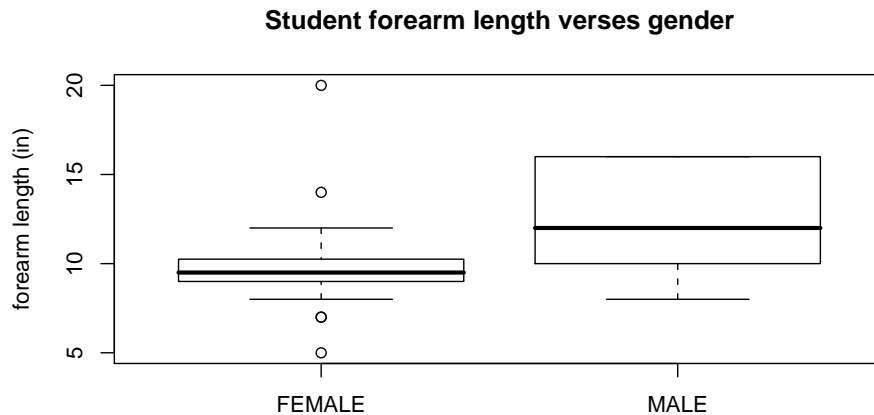
- (a) (1 point) What is wrong with this conclusion?

**Solution:** The study is not representative of all Americans. It only represents people living in Tucson. (We cannot comment as to whether or not the sample size is large enough since we haven't specified the desired margin of error needed — soon we will learn how to do that.)

- (b) (1 point) Restate an appropriate conclusion that can be reached from the study.

**Solution:** "Study data indicates that eighty percent of Tucson residents prefer orange juice over grape juice." (Shouldn't say simply "Eighty percent of Tucson residents prefer orange juice" because it sounds like you have determined the true value — the parameter — when we have only measured a statistic)

4. Use the below box plot to answer the following questions.



- (a) (1 point) Which gender has a higher median forearm length?

**Solution:** Males

- (b) (1 point) What is the approximate median forearm length for the females?

**Solution:** Just below 10 in.

- (c) (1 point) Which gender has a larger variation in forearm length for the middle 50% of individuals?

**Solution:** Males

- (d) (1 point) How many outliers are there in this data set as indicated by the box plots?

**Solution:** The circles indicate outliers, there are 4.

5. Using the below table for our class to answer the following questions.

	BLACK	BLOND	BROWN
FEMALE	2	5	11
MALE	3	0	3

- (a) (1 point) Find the probability of selecting a person with brown hair.

**Solution:**

$$P(\text{brown}) = \frac{14}{24} = 0.583$$

- (b) (1 point) Would it be unusual to randomly select a person with brown hair?

**Solution:** No since  $P \not\leq 0.05$ .

- (c) (1 point) Find the probability of randomly selecting three males without replacement.

**Solution:**

$$P(\text{male \& male \& male}) = \frac{6}{24} \cdot \frac{5}{23} \cdot \frac{4}{22} = 0.00988$$

- (d) (1 point) If you randomly select 6 people with replacement, what is the probability that at least one has brown hair?

**Solution:**

$$P(\text{at least one brown}) = 1 - P(\text{none brown}) = 1 - P(\text{not brown})^6 = 1 - \left(\frac{10}{24}\right)^6 = 0.995$$

- (e) (1 point) Find the probability of selecting a female student or a student with brown hair.

**Solution:**

$$P(\text{brown or female}) = \frac{21}{24} = 0.875$$

- (f) (1 point) Find the probability of selecting a person with brown hair given that they are female.

**Solution:**

$$P(\text{brown}|\text{female}) = \frac{11}{18} = 0.611$$



6. In the 1964 movie *Dr. Strangelove or: How I Learned to Stop Worrying and Love the Bomb*, Brigadier General Jack D. Ripper, a delusional commander of the United States Air Force, orders nuclear armed B-52 bomber planes to attack Russia. The planes use the CRM-114<sup>1</sup> discriminator, which, to prevent false or misleading orders from being received, is designed not to receive at all, unless the message is preceded by a three-letter code prefix.

- (a) (1 point) To recall the planes before they bomb Russia, the US government must issue the correct recall code. How many possible codes are there if the code is composed of three characters, each character being one of the 26 letters in the alphabet?

**Solution:**  $N = n_1 \cdot n_2 \cdot n_3 = 26 \cdot 26 \cdot 26 = 17,576$

- (b) (1 point) If there are only 20 minutes before the planes drop the bombs and it takes 2 seconds to issue an individual code, how many random codes can be issued to guess the correct code?

**Solution:**

$$x = 20\text{min} \cdot \frac{60 \text{ sec}}{\text{min}} \cdot \frac{1 \text{ code}}{2 \text{ sec}} = 600 \text{ codes}$$

- (c) (1 point) What is the probability that the US government will issue the correct recall code — by randomly guessing codes — before it's too late?

**Solution:** All we need is one correct guess out of the 600 tries we have. So we need to find the probability of **at least one correct guess**.

$$P(\text{at least 1 correct}) = 1 - P(\text{none correct for all 600 guesses})$$

We can treat this as independent since we are randomly selecting  $n = 600$  codes out of  $N = 17,576$  satisfies  $n/N \leq 0.05$ . Therefore:

$$P(\text{at least 1 correct}) = 1 - P(\text{none})^{600} = 1 - (1 - 1/17,576)^{600} = 0.0336$$

Caution, it would not be correct to calculate the probability as  $600/17,576 = 0.0341$  because this represents the probability that you would be happy with any of the 600 guesses out of 17,576, or more formally,  $P(\text{guess 1 or guess 2 or } \dots \text{ guess 600}) = 600/17,576$ . We are only happy — ie. the world survives — if we guess the 1 correct possibility of the 17,576 during our 600 trials.

- (d) (1 point) Would it be unusual to issue the recall code in 20 minutes by guessing? (WHY)

**Solution:** Yes, it would be unusual since  $p \leq 0.05$ !

<sup>1</sup>CRM114 is also the name of a computer program which uses a statistical approach for classifying data and is especially utilized for filtering email spam. It was named after the fictional device.

7. (2 points) When doing blood testing for HIV infections, the procedure can be made more efficient and less expensive by combining samples of blood specimens. If samples from five people are combined and the mixture tests negative, we know that all five individual samples are negative. Find the probability of a positive result for five samples combined into one mixture, assuming the probability of an individual blood sample testing positive is 10%. (Based on data from the NY State Health Department)

**Solution:** If multiple individual blood samples are mixed together, the mixture will test positive if 1 or more individuals are positive. Therefore,

$$\begin{aligned}
 P(\text{Positive Mixture}) &= P(1 \text{ or more samples pos}) \\
 &= 1 - P(\text{None positive}) \\
 &= 1 - P(\text{NEG \& NEG \& NEG \& NEG \& NEG}) \\
 &= 1 - P(\text{NEG})^5 && \text{prop. of independence} \\
 &= 1 - (1 - P(\text{POS}))^5 && P(\text{NEG}) = 1 - P(\text{POS}) \\
 &= 1 - (1 - 0.1)^5 \\
 &= 0.41
 \end{aligned}$$

8. (2 points) Given  $y = \{a, -2a, 4a\}$ , completely simplify the following expression. Assume  $a$  is an unknown constant.

$$\left(\sum(y_i - 2a)\right)^2$$

**Solution:**

$$\begin{aligned}
 \left(\sum(y_i - 2a)\right)^2 &= (a - 2a + -2a - 2a + 4a - 2a)^2 \\
 &= (-3a)^2 \\
 &= 9a^2
 \end{aligned}$$

# Basic Statistics: Quick Reference & R Commands

by Anthony Tanbakuchi. Version 1.8  
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:  $x=c(x1, x2, ...)$   
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  $x=c(x1, x2, x3, ...)$

$$\begin{aligned}\text{total} &= \sum_{i=1}^n x_i = \text{sum}(x) & (5) \\ \text{min} &= \text{min}(x) & (6) \\ \text{max} &= \text{max}(x) & (7) \\ \text{six number summary} &: \text{summary}(x) & (8) \\ \bar{\mu} &= \frac{\sum x_i}{N} = \text{mean}(x) & (9) \\ \bar{x} &= \frac{\sum x_i}{n} = \text{mean}(x) & (10) \\ \tilde{x} &= P_{50} = \text{median}(x) & (11)\end{aligned}$$

$$\begin{aligned}\sigma &= \sqrt{\frac{\sum (x_i - \mu)^2}{N}} & (12) \\ s &= \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} = \text{sd}(x) & (13) \\ CV &= \frac{\sigma}{\mu} = \frac{s}{\bar{x}} & (14)\end{aligned}$$

### 2.2 RELATIVE STANDING

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

Percentiles:

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

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!)

$$\begin{aligned}P(A) &= \frac{x}{n} & (17) \\ P(\bar{A}) &= 1 - P(A) & (18) \\ P(A \text{ or } B) &= P(A) + P(B) - P(A \text{ and } B) & (19) \\ P(A \text{ or } B) &= P(A) + P(B) \quad \text{if } A, B \text{ mut. excl.} & (20) \\ P(A \text{ and } B) &= P(A) \cdot P(B|A) & (21) \\ P(A \text{ and } B) &= P(A) \cdot P(B) \quad \text{if } A, B \text{ independent} & (22) \\ n! &= n(n-1) \cdots 1 = \text{factorial}(n) & (23) \\ n!_k &= \frac{n!}{(n-k)!} \quad \text{Perm. no elem. alike} & (24) \\ n!_1 n!_2 \cdots n!_k &= \frac{n!}{n_1! n_2! \cdots n_k!} \quad \text{Perm. } n_1 \text{ alike, } \dots & (25) \\ nC_k &= \frac{n!}{(n-k)!k!} = \text{choose}(n, k) & (26)\end{aligned}$$

## 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.

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

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

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

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

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

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

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

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

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

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

### 5.1 UNIFORM DISTRIBUTION

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

### 5.2 NORMAL DISTRIBUTION

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

### 5.3 t-DISTRIBUTION

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

### 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

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

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

## 7 Estimation

### 7.1 CONFIDENCE INTERVALS

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

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

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

$$df = n - 1$$

$$\begin{aligned}\text{variance: } \frac{(n-1)s^2}{\chi^2_{\alpha/2}} < \sigma^2 < \frac{(n-1)s^2}{\chi^2_{1-\alpha/2}} & (62) \\ df = n - 1\end{aligned}$$

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

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

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

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

$$df = n - 1$$

### 7.2 CI CRITICAL VALUES (TWO SIDED)

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

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

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

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

### 7.3 REQUIRED SAMPLE SIZE

$$\begin{aligned}\text{proportion: } n &= \hat{p}\hat{q} \left(\frac{z_{\alpha/2}}{E}\right)^2, & (70) \\ (\hat{p} = \hat{q} = 0.5 \text{ if unknown})\end{aligned}$$

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

## 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  $\alpha$ .

### 8.1 1-SAMPLE PROPORTION

$H_0: p = p_0$

**prop.test(x, n, p=p<sub>0</sub>, alternative="two.sided")**

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

### 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= $\mu_0$ , 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= $\sigma$ , sig.level= $\alpha$ , power=1 -  $\beta$ , 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(x<sub>1</sub>, x<sub>2</sub>)** and **n=c(n<sub>1</sub>, n<sub>2</sub>)**

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

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

Required Sample size:

**power.prop.test(p1=p<sub>1</sub>, p2=p<sub>2</sub>, power=1 -  $\beta$ , sig.level= $\alpha$ , 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= $\sigma$ , sig.level= $\alpha$ , power=1 -  $\beta$ , 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= $\sigma$ , sig.level= $\alpha$ , power=1 -  $\beta$ , 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  $2 \times 2$  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.

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

### 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=aoov(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
- 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 **MyTable** in 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
```