Introduction, Sections 1.1 and 1.3, The R Statistical Package

Note made by: Dr. Timothy Hanson Instructor: Peijie Hou

Department of Statistics, University of South Carolina

Stat 205: Elementary Statistics for the Biological and Life Sciences

Course and instructor information

- Stat 205: Elementary Statistics for the Biological and Life Sciences, Section 002, WMBB Nursing 127. Tuesday and Thursday 8:30am–9:45am
- Instructor: Peijie Hou (PJ)
- Office: LeConte 209E
- Email: houp@email.sc.edu
- Office hours: T/W/Th 10am 11:30am or by appointment
- Instructor assistant (grader): TBA
- **Textbook:** Samuels, M. L., and Witmer, J. A. & Schaffner, A. (2010). *Statistics for the Life Sciences,* 4th Ed
- Course website: http://www.stat.sc.edu/~houp/ Stat205/stat205_spring2015.html

Homework and exams

- About 8–9 homeworks you will turn in for credit (45% of your grade). Each will be one page, one-sided. Most will be a statistical analysis in R, with pertinent output included and a short write-up. I will give you an example, and your first homework, next week.
- Homework problems from each section will be assigned but not collected or graded. These will form the basis of exam questions.
- Three exams (mid February, late March and Thursday, April 30 at 9am), each covering about a third of the course. Not cumulative (except for a short answer part of the final). Each exam is worth 15% of your grade. (45% total)
- Attendance is 10%. 2.5 hours a week is a small investment for the wealth of knowledge you will gain.
- Strong, positive correlation between attendance and your final grade.

Grading scales; class rules; expectation

- A: 90-100% B+: 85-89.9% B: 80-84.9% and so on.
- Scientific calculator and access to the Internet is required.
- NO cell phones! NO talking! NO helping each other in the exam!
- My expectation for you:
 - *READ* the sections of the text to be covered prior to the class section.
 - ATTEND class regularly and on time.
 - BRING lecture notes with you.
 - ATTEMPT to do all assigned homework.
 - ASK questions when you dont understand!

Topics we'll cover

- Graphical displays and summary statistics
- Probability, random variables (normal and binomial)
- Confidence intervals for μ and p
- Two-sample testing and CI
- 2 × 2 tables: relative risk & odds ratios
- Analysis of variance
- Linear regression
- Logistic regression, survival analysis, ROC curves
- Use of the statistical package R to analyze real data

Tentative lecture schedule is on syllabus.

- Clinical trials/drug development compare existing treatments with new methods to cure disease.
- Agriculture enhance crop yields, improve pest resistance.
- Ecology study how ecosystems develop/respond to environmental impacts.
- Lab studies learn more about biological tissue/cellular activity.

1.1 Statistics and the life sciences

- Statistics is the science of
 - collecting,
 - summarizing,
 - analyzing, and
 - interpreting

data.

- Goal: to understand the underlying biological phenomena that generate the data.
- Statistics separates signal from noise.
- Are there *associations* or *relationships* among variables in the data?

Example 1.1.1: vaccine for Anthrax

Table 1.1.1 Response of sheep to anthrax						
	Treatment					
Response	Vaccinated	Not vaccinated				
Died of anthrax Survived	0 24	24 0				
Total Percent survival	24 100%	24 0%				

Is there variability in the data?

Example 1.1.2: liver tumors in mice

Table 1.1.2 Incidence of liver tumors in mice						
	Treatment					
Response	E. coli	Germ free				
Liver tumors	8	19				
No liver tumors	5	30				
Total	13	49				
Percent with liver tumors	62%	39%				

- Is there an association between germ environment (germ-free vs. *E. coli*) and whether liver tumors develop?
- Is the association perfect?
- Statistics can help answer *whether* there's a difference and further *quantify* the effect of germ exposure (Chapter 10).

Example 1.1.4: MAO and schizophrenia

Table 1.1.4 MAO activity in schizophrenic patients							
Diagnosis			MAO activity				
I:	6.8	4.1	7.3	14.2	18.8		
Chronic undifferentiated	9.9	7.4	11.9	5.2	7.8		
schizophrenic	7.8	8.7	12.7	14.5	10.7		
(18 patients)	8.4	9.7	10.6				
II:	7.8	4.4	11.4	3.1	4.3		
Undifferentiated with	10.1	1.5	7.4	5.2	10.0		
paranoid features	3.7	5.5	8.5	7.7	6.8		
(16 patients)	3.1						
III:	6.4	10.8	1.1	2.9	4.5		
Paranoid schizophrenic (8 patients)	5.8	9.4	6.8				

- Monoamine oxidase (MAO) enzyme thought to regulate behavior.
- Blood from n = 42 schizophrenia patients collected, stratified by diagnosis (I, II, III).
- Is there an association between MAO and diagnosis?

Example 1.1.4: MAO and schizophrenia



Figure 1.1.2 MAO activity in schizophrenic patients

- What happens to MAO as severity of diagnosis increases? Is the relationship perfect?
- These are side-by-side dotplots, described in Sec. 2.2. Formal approach in Chapter 11.

- Data can come from observational studies, planned experiments, clinical trials, etc.
- Data are *random*. Formally, a piece of data is a random variable (Chapter 3).
- The underlying mathematics that drives the methodology in this course relies on assuming data are a **random sample** from their population.
- A **random sample** is one in which each subject has the same probability of being measured, and subjects are chosen independently of each other.
- This provides a representative set of observations from the population, the data Y₁, Y₂,..., Y_n.

Random sampling

- The **population** is *all* the subjects/animals/specimens/etc. of interest.
- Since we can't measure the entire population (usually) we take a small sample of size *n* and use the data collected to *infer* about the population.



R computing & graphics package

- R is a powerful, free statistical computing and graphics package.
- Popular with many researchers due to contributed packages: R functions to do specialized, advanced, & often complex statistical analyses.
- R can also do many important, routine calculations, analyses, and provide common graphical displays used in this course.
- Installed in several of the computing labs across campus, e.g. Sloan 108 & 109, Gambrell 003.
- You can download it and install it from CRAN: http://cran.r-project.org/

The Comprehensive R Archive Network



Here is where you download R.

Installing R

- From http://cran.r-project.org/, under Download and Install R click on your platform (Linux, MacOS X, or Windows).
- **for Windows** click on base and on the next page click on Download R 3.1.2 for Windows (this is the latest release as of August 2012).
- Click Save File and when it's done downloading run the executable by clicking on it alternatively you can choose to Run Program directly after downloading from the web.
- The installation program will ask you a series of questions; choose the defaults. (e.g. English language, the suggested installation folder, the checked selected components to install, not to customize startup options, shortcut in the Start Menu, and additional tasks).
- When it's done, click on the new R desktop icon. Click on the console. This is where you will type commands to R.

The R interface



Initially, there is only the console window open. If you make plots, other windows will open too.

Note that the # sign is a "comment" – R ignores anything after #.

```
# generate some random normal data
data=rnorm(100)
# look at a histogram and a boxplot
hist(data)
boxplot(data)
# compute the sample mean, median, variance, standard deviation
mean(data)
median(data)
var(data)
sd(data)
# if you have a question about a command, preface it with ?
?hist
```

```
# read data from web, take 1st & 2nd columns as MAO and group indicators, plot
stuff=read.table("http://www.stat.sc.edu/~houp/Stat205/mao.txt",header=FALSE)
stuff
MAO=stuff[,1]
MAO
group=stuff[,2]
group
plot(group,MAO)
# you can also read data from a file on your computer (text, Excel, etc.)
```

MAO data: output

```
> stuff=read.table("http://www.stat.sc.edu/~houp/Stat205/mao.txt",header=FALSE)
> stuff
V1 V2
1
   6.8
        1
2
  4.1
       1
3
  7.3 1
4
  14.2
       1
5
  18.8
       1
  9.9
6
       1
7
  7.4
       1
8
  11.9
       1
9
  5.2
       1
10 7.8
       1
11 7.8
       1
12 8.7
       1
13 12.7
       1
14 14.5 1
15 10.7
       1
16 8.4
       1
17 9.7
       1
18 10.6
        1
19 7.8
       2
20 4.4
       2
21 11.4
       2
22 3.1
       2
23 4.3 2
24 10.1 2
25 1.5 2
```

MAO data: output continued

26 7.4 2 27 5.2 2 28 10.0 2 29 3.7 2 30 5.5 2 31 8.5 2 32 7.7 2 33 6.8 2 34 3.1 2 35 6.4 3 36 10.8 3 37 1.1 3 38 2.9 3 39 4.5 3 40 5.8 3 41 9.4 3 42 6.8 3 > MAO=stuff[,1] > MAO [1] 6.8 4.1 7.3 14.2 18.8 9.9 7.4 11.9 5.2 7.8 7.8 8.7 12.7 14.5 10.7 9.7 8.4 [18] 10.6 7.8 4.4 11.4 3.1 4.3 10.1 1.5 7.4 5.2 10.0 3.7 5.5 8.5 7.7 6.8 3.1 [35] 6.4 10.8 1.1 2.9 4.5 5.8 9.4 6.8 > group=stuff[,2] > group [1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 > plot(group,MAO)

Plot of MAO data from R



You can right click on an R plot to save it to the clipboard as a metafile or bitmap. These can be saved into Microsoft applications such as Word. You can also leftclick on the plot then under Save choose Save as and save the plot, e.g. PDF.

Plot of MAO data from R

IR RGui		
Elle History Resize Windows		
R Console	R Graphics: Device 2 (ACTIVE)	
<pre>99 3.7 2 90 5.5 2 10 5.5 2 11 0.5 2 13 0.5 2 13 0.5 2 14 0.1 2 14 0.1 2 14 0.1 2 14 0.1 2 15 0.4 3 15 0.4</pre>	9 - - 0	° ° 8 8 ° ° ° • • • • • • • • • • • • • • • • •

R window after cutting and pasting the commands a few slides ago.

- R will allow you to do all analyses covered in this course, and beyond.
- There are some tutorials, both installed in R and on the web. See course website. This can get you started.
- For homework, I'll give you a skeleton set of commands to get the basic job done with no frills.
- R's error messages can be cryptic and therefore R is not as "user friendly" as some other packages such as Minitab.
- However it is free; now being used by hundreds of thousands of people.