

Data Analysis & Graphics Using R – Solutions to Exercises (May 1, 2010)

Exercise 1

An experimenter intends to arrange experimental plots in four blocks. In each block there are seven plots, one for each of seven treatments. Use the function `sample()` to find four random permutations of the numbers 1 to 7 that will be used, one set in each block, to make the assignments of treatments to plots.

```
> for(i in 1:4)print(sample(1:7))

[1] 7 4 1 5 2 3 6
[1] 7 1 4 5 3 2 6
[1] 6 1 2 7 4 3 5
[1] 6 7 4 3 2 5 1

> ## Store results in the columns of a matrix
> ## The following is mildly cryptic
> sapply(1:4, function(x)sample(1:7))

      [,1] [,2] [,3] [,4]
[1,]    2    7    1    5
[2,]    4    1    3    4
[3,]    1    4    5    1
[4,]    6    3    7    3
[5,]    3    5    6    6
[6,]    5    6    2    2
[7,]    7    2    4    7
```

Exercise 2

Use `y <- rnorm(100)` to generate a random sample of 100 numbers from a normal distribution. Calculate the mean and standard deviation of `y`. Now put the calculation in a loop and repeat 25 times. Store the 25 means in a vector named `av`. Calculate the standard deviation of the values in `av`.

```
> av <- numeric(25)
> sdev <- numeric(25)
> for(i in 1:25){
+ y <- rnorm(100)
+ av[i] <- mean(y)
+ sdev[i] <- sd(y)
+ }
> sd(av)

[1] 0.1197
```

Exercise 3

Create a function that does the calculations of exercise 2.

```

> avfun <- function(m=50, n=25){
+   for(i in 1:25){
+     y <- rnorm(50)
+     av[i] <- mean(y)
+   }
+   sd(av)
+ }

```

It is insightful to run the function several times, and see how the value that is returned varies.

Exercise 7

The function `pexp(x, rate=r)` can be used to compute the probability that an exponential variable is less than `x`. Suppose the time between accidents at an intersection can be modeled by an exponential distribution with a rate of .05 per day. Find the probability that the next accident will occur during the next 3 weeks.

We require the probability that the time to the next accident is less than or equal to 21 days.

```
> pexp(21, .05)
```

```
[1] 0.65
```

Note that the rate is both the waiting time from an arbitrary time to the next accident, and the “interarrival” time between accidents. The expected time to the next accident is unaffected by whether or not an accident has just occurred.

Exercise 8

Use the function `rexp()` to simulate 100 exponential random numbers with rate .2. Obtain a density plot for the observations. Find the sample mean of the observations. Compare with the the population mean. (The mean for an exponential population is $1/\text{rate}$.)

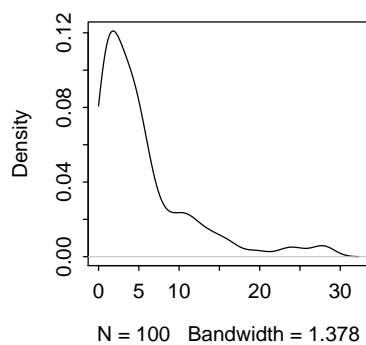


Figure 1: Density plot, for 100 random values from an exponential distribution with `rate = 0.2`

```

> ## Code
> z <- rexp(100, .2)
> plot(density(z, from=0), main="")

```

Notice the use of the argument `from=0`, to prevent `density()` from giving a positive density estimate to negative values.

Compare `mean(z) = 5.42` with $1/0.2 = 5$.

Exercise 10

The following data represent the total number of aberrant crypt foci (abnormal growths in the colon) observed in 7 rats that had been administered a single dose of the carcinogen azoxymethane and sacrificed after six weeks:

```
87 53 72 90 78 85 83
```

Enter these data and compute their sample mean and variance. Is the Poisson model appropriate for these data. To investigate how the sample variance and sample mean differ under the Poisson assumption, repeat the following simulation experiment several times:

```
x <- rpois(7, 78.3)
mean(x); var(x)
```

```
> y <- c(87, 53, 72, 90, 78, 85, 83)
> c(mean=mean(y), variance=var(y))
```

```
mean variance
78.29 159.90
```

Then try

```
> x <- rpois(7, 78.3)
> c(mean=mean(x), variance=var(x))
```

```
mean variance
75.57 56.29
```

It is unusual to get as big a difference between the mean and the variance as that observed for these data, making it doubtful that these data are from a Poisson distribution.

*Exercise 11**

A Markov chain is a data sequence which has a special kind of dependence. For example, a fair coin is tossed repetitively by a player who begins with \$2. If ‘heads’ appear, the player receives one dollar; otherwise, she pays one dollar. The game stops when the player has either \$0 or \$5. The amount of money that the player has before any coin flip can be recorded – this is a Markov chain. A possible sequence of plays is as follows:

Player’s fortune:	2	1	2	3	4	3	2	3	2	3	2	1	0
Coin Toss result:	T	H	H	H	T	T	H	T	H	T	T	T	T

Note that all we need to know in order to determine the player’s fortune at any time is the fortune at the previous time as well as the coin flip result at the current time. The probability of an increase in the fortune is .5 and the probability of a decrease in the fortune is .5. Such transition probabilities are usually summarized in a transition matrix:

$$P = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ .5 & 0 & .5 & 0 & 0 & 0 \\ 0 & .5 & 0 & .5 & 0 & 0 \\ 0 & 0 & .5 & 0 & .5 & 0 \\ 0 & 0 & 0 & .5 & 0 & .5 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

The (i, j) entry of this matrix refers to the probability of making a change from the value i to the value j . Here, the possible values of i and j are $0, 1, 2, \dots, 5$. According to the matrix, there is a probability of 0 of making a transition from \$2 to \$4 in one play, since the $(2,4)$ element is 0; the probability of moving from \$2 to \$1 in one transition is 0.5, since the $(2,1)$ element is 0.5.

The following function can be used to simulate N values of a Markov chain sequence, with transition matrix P :

```
Markov <- function (N=100, initial.value=1, P)
{
  X <- numeric(N)
  X[1] <- initial.value + 1
  n <- nrow(P)
  for (i in 2:N){
    X[i] <- sample(1:n, size=1, prob=P[X[i-1],])
  }
  X - 1
}
```

- (a) Simulate 15 values of the coin flip game, starting with an initial value of \$2.
- Simulate 100 values of the Markov chain which has the following transition matrix. Save the result to a vector and use `ts.plot()` to plot the sequence.

$$P = \begin{bmatrix} 0.10 & 0.90 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0.50 & 0.00 & 0.50 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.50 & 0.00 & 0.50 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.50 & 0.00 & 0.50 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.50 & 0.00 & 0.50 \\ 0.00 & 0.00 & 0.00 & 0.00 & 1.00 & 0.00 \end{bmatrix}$$

- (b) Now simulate 1000 values from the above Markov chain, and calculate the proportion of times the chain visits each of the states. It can be shown, using linear algebra, that in the long run, this Markov chain will visit the states according to the following *stationary* distribution

0	1	2	3	4	5
0.1098901	0.1978022	0.1978022	0.1978022	0.1978022	0.0989011

There is a result called the *ergodic* theorem which allows us to estimate this distribution by simulating the Markov chain for a long enough time. Compare your calculated proportions with above theoretical proportions. Repeat the experiment

Here is code that may be used for these calculations.

```
(a) > P <- matrix(c(1, rep(0,5), rep(c(.5,0,.5, rep(0,4))),4), 0,1),
+               byrow=TRUE,nrow=6)
> Markov(15, 2, P)

(b) > Pb <- matrix(c(0.10,0.90,0.00,0.00,0.00,0.00,
+                   0.50,0.00,0.50,0.00,0.00,0.00,
+                   0.00,0.50,0.00,0.50,0.00,0.00,
+                   0.00,0.00,0.50,0.00,0.50,0.00,
+                   0.00,0.00,0.00,0.50,0.00,0.50,
+                   0.00,0.00,0.00,0.00,1.00,0.00),
+               byrow=TRUE, nrow=6)
> xb <- Markov(100, 1, Pb)
> xb
> ts.plot(xb)

(c) > xc <- Markov(1000, 1, Pb)
> table(xb)/1000 # one of several ways to calculate the proportions
> xc
> xc2 <- Markov(10000, 1, Pb)
> table(xc2)/10000

(d) > Pd <- matrix(c(0.50, 0.50, 0, 0, 0, 0,
+                   0.50, 0.45, 0.05, 0, 0, 0,
+                   0, 0.01, 0, 0.90, 0.09, 0,
+                   0, 0, 0.01, 0.40, 0.59, 0,
+                   0, 0, 0, 0.50, 0, 0.50,
+                   0, 0, 0, 0, 0.50, 0.50),
+               nrow=6, byrow=TRUE)
```

The following function may be helpful, in examining results.

```
> `plotmarkov` <-
+ function(n=10000, start=1, window=100, transition=Pd){
+   xc2 <- Markov(n, start, transition)
+   z4 <- as.integer(xc2==4)
+   z5 <- as.integer(xc2==5)
+   mav4 <- rollmean(z4,window)
+   mav5 <- rollmean(z5,window)
+   df <- data.frame(av4=mav4, av5=mav5, x=rep(1:1000, length=length(mav4)),
+                   gp=(0:(length(mav4)-1))%/%1000)
+   print(xyplot(av4+av5 ~ x | gp, data=df, layout=c(1,10), type="l",
+             par.strip.text=list(cex=0.65)))
+ }
> ## Use thus
> library(zoo) # Use rollmean() [moving average] from zoo
> library(lattice)
> plotmarkov(start=1)
> plotmarkov(start=4)
```