

Lab: Week 4 – Solutions

36-350 – Statistical Computing

Week 4 – Spring 2021

Apply: Base R

You are given the following 8×8 matrix:

```
set.seed(1001)
mat = matrix(rnorm(64),nrow=8)
mat[5,8] = mat[6,7] = mat[4,2] = NA
```

Question 1

Notes 4A (6,10,12)

Compute the mean for each row and for each column using both `apply()` and either `rowMeans()` or `colMeans()`. (So there should be four function calls overall.) Deal with the NAs by passing (an) additional argument(s) to these functions, when possible.

```
apply(mat,MARGIN=1,FUN=mean,na.rm=TRUE)
```

```
## [1]  0.29520657  0.51323956 -0.46939506 -0.46438383 -0.19427723  0.02530351
## [7]  0.06268717  0.16192061
```

```
rowMeans(mat,na.rm=TRUE)
```

```
## [1]  0.29520657  0.51323956 -0.46939506 -0.46438383 -0.19427723  0.02530351
## [7]  0.06268717  0.16192061
```

```
apply(mat,MARGIN=2,FUN=mean,na.rm=TRUE)
```

```
## [1] -0.1141279 -0.1176691  0.3097320 -0.3976337  0.2625528 -0.1828815
0.0382058
## [8] 0.2301255
```

```
colMeans(mat,na.rm=TRUE)
```

```
## [1] -0.1141279 -0.1176691  0.3097320 -0.3976337  0.2625528 -0.1828815
0.0382058
## [8] 0.2301255
```

Question 2

Function writing review

How does the `Income` variable in R's `state.x77` matrix correlate with other variables? Write a function called `cor_var()` that takes two inputs: `v1`, a numeric vector; and `v2`, another numeric vector whose default value is `state.x77[, "Income"]`. Its output should be the correlation of `v1` and `v2`, computed

via the `cor()` function. Check that `cor_var(v1=state.x77[, "Life Exp"])` gives you 0.3402553, and `cor_var(v1=state.x77[, "Income"])` gives you 1.

```
cor_var = function(v1,v2=state.x77[, "Income"]) { cor(v1,v2) }
cor_var(v1=state.x77[, "Life Exp"])
```

```
## [1] 0.3402553
```

```
cor_var(v1=state.x77[, "Income"])
```

```
## [1] 1
```

Question 3

Notes 4A (6-7,9-10)

Using `apply()` and the function `cor_var()` that you defined in the last question, calculate the correlation between each one of the 8 variables in the `state.x77` matrix and the `Population` variable. Display these correlations.

```
apply(state.x77,MARGIN=2,FUN=cor_var,v2=state.x77[, "Population"])
```

```
## Population      Income Illiteracy   Life Exp      Murder    HS Grad
## 1.000000000 0.20822756 0.10762237 -0.06805195 0.34364275 -0.09848975
##      Frost      Area
## -0.33215245 0.02254384
```

Question 4

Notes 4A (6,10)

Using `apply()` and the base R `stats` package function `cor()`, display the Spearman correlation between each one of the eight variables in the `state.x77` matrix and the `Frost` variable. (Note that `Spearman` is not the default value for the `method` argument to the `cor()` function.)

```
apply(state.x77,MARGIN=2,FUN=cor,y=state.x77[, "Frost"],method="kendall")
```

```
## Population      Income Illiteracy   Life Exp      Murder    HS Grad
## -0.30454484 0.13098703 -0.53610593 0.22213145 -0.41362334 0.29262334
##      Frost      Area
## 1.00000000 0.05403215
```

Variations on Apply: Base R

Question 5

Notes 4B (6)

Create a data frame called `state.df` from the matrix `state.x77` and the factors `state.region` and `state.division`. Be sure to name the two new columns appropriately. Using `state.df` and `tapply()`, compute the average population in each of the four defined regions of the U.S. Display the name of the region has the largest average population (and only that name). Then compute the average population in each of the nine defined divisions of the U.S., and display the name of the division has the largest average population (and only that name). Hint: the names may be displayed using a combination of `names()` and `which.max()`.

```
state.df = data.frame(state.x77,Region=state.region,Division=state.division)
(population.by.reg = tapply(state.df$Population,INDEX=state.df$Region,FUN=mean))
```

```
##      Northeast      South North Central      West
```

```
##          5495.111          4208.125          4803.000          2915.308
(population.by.div = tapply(state.df$Population,INDEX=state.df$Division,FUN=mean))

##          New England      Middle Atlantic      South Atlantic East South Central
##          2031.167          12423.000          4118.250          3379.000
## West South Central East North Central West North Central          Mountain
##          5217.000          8189.000          2384.429          1203.125
##          Pacific
##          5654.800

names(population.by.reg)[which.max(population.by.reg)]

## [1] "Northeast"

names(population.by.div)[which.max(population.by.div)]

## [1] "Middle Atlantic"
```

Question 6

Notes 4A (5) and Notes 4B (3-4)

Split the rows of the data frame `state.df` by state divisions, and call the resulting list `state.df.by.div`. Then use `lapply()` to display just the first two rows of each data frame in the list `state.df.by.div`.

```
state.df.by.div = split(state.df,f=state.df$Division)
lapply(state.df.by.div,head,2)

## $`New England`
##          Population Income Illiteracy Life.Exp Murder HS.Grad Frost Area
## Connecticut      3100      5348          1.1      72.48      3.1      56.0      139 4862
## Maine              1058      3694          0.7      70.39      2.7      54.7      161 30920
##          Region      Division
## Connecticut Northeast New England
## Maine          Northeast New England
##
## $`Middle Atlantic`
##          Population Income Illiteracy Life.Exp Murder HS.Grad Frost Area
## New Jersey      7333      5237          1.1      70.93      5.2      52.5      115 7521
## New York        18076      4903          1.4      70.55     10.9      52.7      82 47831
##          Region      Division
## New Jersey Northeast Middle Atlantic
## New York    Northeast Middle Atlantic
##
## $`South Atlantic`
##          Population Income Illiteracy Life.Exp Murder HS.Grad Frost Area
## Delaware          579      4809          0.9      70.06      6.2      54.6      103 1982
## Florida           8277      4815          1.3      70.66     10.7      52.6      11 54090
##          Region      Division
## Delaware South South Atlantic
## Florida   South South Atlantic
##
## $`East South Central`
##          Population Income Illiteracy Life.Exp Murder HS.Grad Frost Area
## Alabama          3615      3624          2.1      69.05     15.1      41.3      20 50708
## Kentucky          3387      3712          1.6      70.10     10.6      38.5      95 39650
```

```

##           Region           Division
## Alabama    South East South Central
## Kentucky   South East South Central
##
## $`West South Central`
##           Population Income Illiteracy Life.Exp Murder HS.Grad Frost Area
## Arkansas      2110   3378         1.9   70.66   10.1   39.9    65 51945
## Louisiana      3806   3545         2.8   68.76   13.2   42.2    12 44930
##           Region           Division
## Arkansas      South West South Central
## Louisiana      South West South Central
##
## $`East North Central`
##           Population Income Illiteracy Life.Exp Murder HS.Grad Frost Area
## Illinois      11197   5107         0.9   70.14   10.3   52.6   127 55748
## Indiana        5313   4458         0.7   70.88    7.1   52.9   122 36097
##           Region           Division
## Illinois North Central East North Central
## Indiana  North Central East North Central
##
## $`West North Central`
##           Population Income Illiteracy Life.Exp Murder HS.Grad Frost Area
## Iowa          2861   4628         0.5   72.56    2.3   59.0   140 55941
## Kansas         2280   4669         0.6   72.58    4.5   59.9   114 81787
##           Region           Division
## Iowa  North Central West North Central
## Kansas North Central West North Central
##
## $Mountain
##           Population Income Illiteracy Life.Exp Murder HS.Grad Frost Area
## Arizona      2212   4530         1.8   70.55    7.8   58.1    15 113417
## Colorado      2541   4884         0.7   72.06    6.8   63.9   166 103766
##           Region Division
## Arizona      West Mountain
## Colorado      West Mountain
##
## $Pacific
##           Population Income Illiteracy Life.Exp Murder HS.Grad Frost Area
## Alaska         365   6315         1.5   69.31   11.3   66.7   152 566432
## California     21198   5114         1.1   71.71   10.3   62.6    20 156361
##           Region Division
## Alaska         West Pacific
## California      West Pacific

```

Below, we read in a data table showing the fastest women's 100-meter sprint times.

```
sprint.df = read.table("http://www.stat.cmu.edu/~mfarag/350/women_100m_with_header.dat",header=TRUE,stringsAsFactors=FALSE)
```

```
## [1] "data.frame"
```

```
head(sprint.df)
```

```
## Rank Time Wind First.Name Last.Name Country Birthdate Race
```

```
## 1  1 10.49 0.0  Florence Griffith-Joyner      USA 21.12.59 1q1
## 2  2 10.61 1.2  Florence Griffith-Joyner      USA 21.12.59 1
## 3  3 10.62 1.0  Florence Griffith-Joyner      USA 21.12.59 1q3
## 4  4 10.64 1.2  Carmelita                      Jeter  USA 24.11.79 1
## 5  5 10.65 1.1  Marion                      Jones  USA 12.10.75 1
## 6  6 10.67 -0.1 Carmelita                      Jeter  USA 24.11.79 1
##      Location      Date
## 1 Indianapolis 16.07.1988
## 2 Indianapolis 17.07.1988
## 3      Seoul 24.09.1988
## 4      Shanghai 20.09.2009
## 5 Johannesburg 12.09.1998
## 6 Thessaloniki 13.09.2009
```

Question 7

Review of string processing

Extract the last four digits of each entry of the `Date` column. (Hint: you will have to use `as.character()` to convert `sprint.df$Date` from a factor variable to strings.) Create a new data frame called `new.sprint.df` that combines `sprint.df` and a new column called `Year` that contains your extracted four-digit years. Display the first five rows and all 11 columns of `new.sprint.df`. Display the class of the newly created `Year` column.

```
d = sprint.df$Date
year.str = substr(d,nchar(d)-3,nchar(d))
new.sprint.df = data.frame(sprint.df,Year=year.str)
head(new.sprint.df,5)
```

```
##      Rank  Time Wind First.Name      Last.Name Country Birthdate Race
## 1  1 10.49 0.0  Florence Griffith-Joyner      USA 21.12.59 1q1
## 2  2 10.61 1.2  Florence Griffith-Joyner      USA 21.12.59 1
## 3  3 10.62 1.0  Florence Griffith-Joyner      USA 21.12.59 1q3
## 4  4 10.64 1.2  Carmelita                      Jeter  USA 24.11.79 1
## 5  5 10.65 1.1  Marion                      Jones  USA 12.10.75 1
##      Location      Date Year
## 1 Indianapolis 16.07.1988 1988
## 2 Indianapolis 17.07.1988 1988
## 3      Seoul 24.09.1988 1988
## 4      Shanghai 20.09.2009 2009
## 5 Johannesburg 12.09.1998 1998
```

```
class(new.sprint.df$Year)
```

```
## [1] "character"
```

Question 8

Notes 4B (6)

Using `tapply()` and the newly created `Year` column, compute the median 100-meter sprint time in each year of the data frame. Call the resulting vector `med.time.by.year`. Create a table of median times. Which median time appears the most, and how many times does it appear? When is the last year that that particular median time appeared in the data?

```
med.time.by.year = tapply(new.sprint.df$Time,INDEX=new.sprint.df$Year,FUN=median)
table(med.time.by.year)
```

```
## med.time.by.year
## 10.98 10.99 10.995      11 11.005 11.01 11.015 11.02 11.03 11.035 11.04
##      1      1      1      6      1      5      2      4      7      1      7
## 11.045 11.05 11.055 11.06 11.07 11.08
##      1      2      1      2      2      1
```

```
med.time.by.year
```

```
## 1968 1972 1973 1976 1977 1978 1979 1980 1981 1982 1983
## 11.080 11.070 11.070 11.055 11.030 11.050 11.040 11.060 11.040 11.010 11.035
## 1984 1985 1986 1987 1988 1989 1990 1991 1992 1993 1994
## 10.990 11.010 11.030 11.040 11.000 11.040 11.050 10.995 10.980 11.000 11.015
## 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005
## 11.040 11.000 11.030 11.010 11.020 11.030 11.020 11.020 11.045 11.030 11.030
## 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016
## 11.060 11.040 11.020 11.010 11.015 11.000 11.000 11.000 11.030 11.010 11.005
## 2017
## 11.040
```

11.04 seconds appears seven times in the data. This median time last occurred in 2017.

Below, we read in a data table related to the political economy of strikes.

```
strikes.df = read.csv("http://www.stat.cmu.edu/~mfarag/350/strikes.csv")
class(strikes.df)
```

```
## [1] "data.frame"
```

```
head(strikes.df)
```

```
## country year strike.volume unemployment inflation left.parliament
## 1 Australia 1951      296          1.3      19.8          43.0
## 2 Australia 1952      397          2.2      17.2          43.0
## 3 Australia 1953      360          2.5       4.3          43.0
## 4 Australia 1954        3          1.7       0.7          47.0
## 5 Australia 1955      326          1.4       2.0          38.5
## 6 Australia 1956      352          1.8       6.3          38.5
## centralization density
## 1      0.3748588      NA
## 2      0.3751829      NA
## 3      0.3745076      NA
## 4      0.3710170      NA
## 5      0.3752675      NA
## 6      0.3716072      NA
```

```
dim(strikes.df) # Note that since 18 × 35 = 630 > 625, some years missing from some countries
```

```
## [1] 625  8
```

Question 9

Notes 4A (5) and Notes 4B (5)

Split `strikes.df` by country, using the `split()` function. Call the resulting list `strikes.by.country`.

Using `strikes.by.country` and `sapply()`, compute the average centralization metric (a quantity related to unionization) for each country over the range of years in the file. Display the names of the countries that had the highest and lowest average centralization metric (and only the names of those countries).

```
strikes.by.country = split(strikes.df,strikes.df$country)
sapply.out = sapply(strikes.by.country,function(df)mean(df$centralization))
cat("Maximum was in ",names(sapply.out)[which.max(sapply.out)],".\n",sep="")
```

```
## Maximum was in Austria.
```

```
cat("Minimum was in ",names(sapply.out)[which.min(sapply.out)],".\n",sep="")
```

```
## Minimum was in Canada.
```

Question 10

Notes 4B (5)

Using `strikes.by.country` and `sapply()`, compute a summary of the long-term centralization metric for each country. Study the output—do its dimensions make sense to you?

```
sapply(strikes.by.country, function(df)summary(df$centralization))
```

```
##      Australia  Austria  Belgium      Canada  Denmark  Finland
## Min.   0.3701921 0.9951362 0.7451018 4.985230e-06 0.4951243 0.7453985
## 1st Qu. 0.3723613 0.9963630 0.7480245 8.232258e-04 0.4971313 0.7486803
## Median 0.3745076 0.9977592 0.7489699 2.206919e-03 0.5003940 0.7501793
## Mean   0.3746440 0.9976705 0.7494852 2.244134e-03 0.4999586 0.7503741
## 3rd Qu. 0.3763172 0.9988332 0.7514769 3.468929e-03 0.5022077 0.7521806
## Max.   0.3798597 0.9997884 0.7544044 4.849537e-03 0.5048790 0.7549842
##
##      France  Germany  Ireland      Italy      Japan Netherlands
## Min.   0.0002446096 0.2453393 0.4951136 0.2454353 0.1205130 0.7454194
## 1st Qu. 0.0013202927 0.2477310 0.4974278 0.2490072 0.1233528 0.7474436
## Median 0.0028737475 0.2493486 0.4994846 0.2507560 0.1247869 0.7491107
## Mean   0.0027299088 0.2499682 0.4997119 0.2506995 0.1246753 0.7496027
## 3rd Qu. 0.0042529214 0.2524444 0.5022122 0.2527474 0.1261252 0.7520595
## Max.   0.0049236913 0.2548710 0.5048117 0.2547880 0.1297671 0.7540260
##
## New.Zealand  Norway  Sweden Switzerland      UK      USA
## Min.   0.3706028 0.8700540 0.8701569 0.4956250 0.3701746 0.000109027
## 1st Qu. 0.3730609 0.8730289 0.8723843 0.4976971 0.3738972 0.001355673
## Median 0.3761876 0.8750384 0.8756796 0.4993706 0.3756106 0.002406464
## Mean   0.3759404 0.8753418 0.8752538 0.4999900 0.3759468 0.002390639
## 3rd Qu. 0.3786986 0.8780262 0.8778525 0.5024074 0.3785299 0.003252352
## Max.   0.3798821 0.8799584 0.8794025 0.5048787 0.3797725 0.004975356
```

When the output of the function is always of the same length, then `sapply()`'s output will be simplified to a matrix. Summary has six outputs, so `sapply()`'s output has six rows.

Question 11

Notes 4B (5)

Using `strikes.by.country` and just *one* call to `sapply()`, compute the average unemployment rate, average inflation rate, and average strike volume for each country. The output should be a matrix of dimension 3 x 18. Also, within that call, give the output matrix appropriate row names.

```
sapply(strikes.by.country,function(df) {
  c("Unemployment Average"=mean(df$unemployment),
    "Inflation Average"=mean(df$inflation),
    "Strike Average"=mean(df$strike.volume))
})
```

```
##           Australia   Austria   Belgium   Canada   Denmark
## Unemployment Average 3.505714 2.540000 3.646667 6.042857 5.711429
## Inflation Average    6.594286 5.102857 4.150000 4.797143 6.582857
## Strike Average      378.600000 25.600000 244.000000 749.542857 194.828571
##           Finland   France   Germany   Ireland   Italy
## Unemployment Average 2.571429 3.182857 3.117143 7.771429 6.725714
## Inflation Average    7.317143 6.948571 3.294286 8.151429 8.005714
## Strike Average      448.542857 185.400000 43.828571 547.428571 997.685714
##           Japan Netherlands New.Zealand   Norway   Sweden
## Unemployment Average 1.602857 3.691429 1.002857 1.428571 2.137143
## Inflation Average    5.820000 4.814286 7.691429 6.320000 6.434286
## Strike Average      165.828571 26.114286 259.257143 75.114286 73.485714
##           Switzerland   UK   USA
## Unemployment Average 0.3285714 3.451429 5.542857
## Inflation Average    3.4171429 7.105714 4.428571
## Strike Average      3.6571429 322.714286 448.228571
```

Question 12

Notes 4B (5)

Using `strikes.df`, `split()`, and `sapply()`, compute the average unemployment rate for each country, before and during 1970, and after 1970. Display the output; it should be a numeric vector of length 36. One way to perform the splitting is to define a new column called `pre1970` that indicates that a year column is less than or equal to 1970. Then use both `country` and `pre1970` to do the splitting. If you are not sure how to use both factor variables at once, look at the documentation for `split()`, specifically its argument `f`.

```
strikes.df$pre1970 = strikes.df$year<=1970
(unemp.vec = sapply(split(strikes.df,list(strikes.df$country,strikes.df$pre1970)),
  function(df)mean(df$unemployment)))
```

```
##   Australia.FALSE   Austria.FALSE   Belgium.FALSE   Canada.FALSE
##           5.5066667           2.1000000           4.7700000           8.0000000
##   Denmark.FALSE   Finland.FALSE   France.FALSE   Germany.FALSE
##           6.2800000           4.3266667           5.6800000           4.1933333
##   Ireland.FALSE   Italy.FALSE   Japan.FALSE Netherlands.FALSE
##           9.2200000           7.3200000           2.0000000           6.7800000
## New.Zealand.FALSE   Norway.FALSE   Sweden.FALSE Switzerland.FALSE
##           2.0066667           1.9933333           2.4000000           0.3866667
##           UK.FALSE   USA.FALSE   Australia.TRUE   Austria.TRUE
##           6.1333333           6.9466667           2.0050000           2.8700000
##   Belgium.TRUE   Canada.TRUE   Denmark.TRUE   Finland.TRUE
##           3.0850000           4.5750000           5.2850000           1.2550000
##   France.TRUE   Germany.TRUE   Ireland.TRUE   Italy.TRUE
##           1.3100000           2.3100000           6.6850000           6.2800000
##   Japan.TRUE   Netherlands.TRUE   New.Zealand.TRUE   Norway.TRUE
##           1.3050000           1.3750000           0.2500000           1.0050000
##   Sweden.TRUE   Switzerland.TRUE   UK.TRUE   USA.TRUE
##           1.9400000           0.2850000           1.4400000           4.4900000
```


Question 13

Review of matrices

Using the result from above, display the difference in the average unemployment rate before and after 1970 for each country. (To be clear: subtract the pre-1970 results from the post-1970 results.) Which country had the biggest increase in average unemployment from before to after? The biggest decrease? (Hint: use the output from Q12 to populate a matrix, with pre-1970 results in one column and post-1970 results in another.)

```
unemp.mat = matrix(unemp.vec, ncol=2)
unemp.diff = unemp.mat[,1]-unemp.mat[,2]
names(unemp.diff) = unique(strikes.df$country)
unemp.diff
```

```
##   Australia   Austria   Belgium   Canada   Denmark   Finland
##   3.5016667 -0.7700000  1.6850000  3.4250000  0.9950000  3.0716667
##      France   Germany   Ireland   Italy   Japan Netherlands
##   4.3700000  1.8833333  2.5350000  1.0400000  0.6950000  5.4050000
## New.Zealand   Norway   Sweden Switzerland   UK   USA
##   1.7566667  0.9883333  0.4600000  0.1016667  4.6933333  2.4566667
```

```
names(unemp.diff)[which.max(unemp.diff)]
```

```
## [1] "Netherlands"
```

```
names(unemp.diff)[which.min(unemp.diff)]
```

```
## [1] "Austria"
```

Question 14

Pipes + Notes 4D (6)

How does the `Frost` variable in R's `state.x77` matrix correlate with other variables? Cast `state.x77` to a data frame, and, using pipes, generate the correlation matrix for `Frost` and `Life.Exp`. (Note that the act of casting changed the name of the life expectancy column from `Life Exp` to `Life.Exp`.) The off-diagonal elements of the matrix should be 0.262068.

```
suppressWarnings(suppressMessages(library(dplyr)))
data.frame(state.x77) %>% select("Life.Exp", "Frost") %>% cor()
```

```
##           Life.Exp   Frost
## Life.Exp 1.000000 0.262068
## Frost    0.262068 1.000000
```

Question 15

Pipes + Notes 4D (9)

Take the `state.df` data frame defined in Q5 above and mutate it so as to create a new column: `GradLit`. This column should have, for each row in the data frame, the percentage of high school graduates divided by the percentage of literate (note: *literate*, not *illiterate*) individuals, times 100. Then pipe the output so as to compute the median value of `GradLit`. (There is a bit of weirdness here: due to environmental issues, your call to `median()` will not work unless it is placed within curly braces. You are only surrounding `median()` with curly braces...not the entire pipe stream!) Your final value should be 53.59844.

```
state.df %>% mutate(GradLit=100*HS.Grad/(100-Illiteracy)) %>% {median(.$GradLit)}
```

```
## [1] 53.59844
```

Question 16

Pipes + Notes 4D (5,6,8)

Take the `state.df` data frame and (1) select all states in the South region, and (2) display the result ordered by the decreasing product of income and life expectancy. In the end, display just the state name and the computed product. There is a quirk here: selecting rows can lead to the loss of row names. (This means that here, you will have a final result but not know which states they correspond to.) To preserve the identity of the states, pipe `state.df` to the function `rownames_to_column("give column name here, like State")`, then do the rest of your piping.

```
suppressWarnings(suppressMessages(library(tidyverse)))
state.df %>% rownames_to_column("State") %>% filter(Region=="South") %>%
mutate(ILE=Income*Life.Exp) %>% select(State,ILE) %>% arrange(desc(ILE))
```

```
##           State      ILE
## 1      Maryland 372095.8
## 2        Florida 340227.9
## 3      Delaware 336918.5
## 4      Virginia 329446.1
## 5         Texas 296929.2
## 6      Oklahoma 284465.9
## 7        Georgia 280397.1
## 8 North Carolina 268188.8
## 9      Tennessee 267890.3
## 10      Kentucky 260211.2
## 11 West Virginia 251309.2
## 12        Alabama 250237.2
## 13 South Carolina 247034.6
## 14      Louisiana 243754.2
## 15      Arkansas 238689.5
## 16 Mississippi 210942.8
```