

Python Notes

Steven Wang

March 30, 2016

```
1 #The parts of the Marin videos that have been noted:
2 #0
3 #1,2,3,4,5,6,7,8,9
4 #2.3,
5 #3.3,3.4,
6 #,4.9
7 #5.1,5.2,5.3
8
9
10 R is case sensitive
11
12 x = 11
13 y <- 11
14
15 ls() # show everything stored in workspace
16 rm(y) # delete y
17
18 x.1 = 14 # variable name can have period and number, but it could not
    start with number
19
20 xx = "marin" # assign characters to variable
21 yy = "1" # if include numbers in quotations, R will treat them as
    characters
22
23 sqrt(y)
24 y^(1/2)
25 log(y) # ln(y)
26 exp(y) # anti-log
27
28 log2(y) # log based on 2
29 abs(-14) # absolute value
30
31 # if you input an incomplete command, R would return with a "+" sign
32
33 c # concatenate
34 x1 = c(1,3,5,7,9)
```

```

35
36 gender = c("male","female")
37
38 2 7 # integer values
39 seq(from=1, to=7,by=1)
40 seq(from=1, to=7,by=1/3)
41 rep(1,times=10)
42
43 rep(c("m","f"),times=5)
44
45 x = 1:5
46 x + 10
47 x * 10
48
49 # if two vectors are of the same length, we may add/subtract/mult/div
    corresponding elements // element-wise
50
51 y[3] #extract the 3rd element of y
52 y[-3] #extract all elements except the 3rd element of y
53 y[1:3] # extract 1st to 3rd elements
54 y[c(1,5)] # extract 1st and 5th
55 y[-c(1,5)] # extract all except 1st and 5th
56 y[y<6] # extract the elements that are less than 6
57
58 mat=matrix(c(1,2,3,4,5,6,7,8,9),nrow=3,byrow=TRUE) # "nrow": the number
    of rows; "byrow=TRUE": elements will be entered in a row-wise //"TRUE"
    has to be capital letters
59
60 mat[1,2] # extract element from row 1, column 2
61 mat[c(1,3),2] # extract elements from row 1 and 3, column 2
62 mat[2,] # extract row 2, all columns // leaving row or column blank to
    extract all columns or rows
63
64 mat*10 # multiply element-wise
65
66 #IMPORTING DATA AND WORKING WITH DATA
67 help(read.table) #help(command you want to know more about)
68 ?read.table #show help
69 Data1 <- read.table(file="/DirectoryPath/data.txt", header=TRUE, sep="\t"
    ) # "header=TRUE" tells R that 1st row is header; "sep=\t": separate by
    table(since the data in the example is table delimited). or "sep
    =",",", "sep=",""
70 read.csv(file, header = TRUE, sep = ",", quote = "/", dec = ".", fill =
    TRUE, comment.char = "", ...)
71
72 Data2 = read.table(file.choose(), header=TRUE, sep = "\t") #press enter
    and R will let you choose your data source file

```

```

73
74 #if using R-studio, GUI, just click some buttons
75 #in Europe, it's common to use comma representing a decimal point
76
77 dim(mat) # show dimensions of the data, rows and columns
78
79 head(Data1) # show the 1st 6 rows of the data
80 tail(Data1) # show the last 6 rows of the data
81
82 Data1[c(5,6,7,8,9),] #square brackets: subset
83 Data1[5:9,]
84 Data[-(4:722)   ]
85
86 #subsetting data
87 names(LungCapData) #show the names of the variables
88 mean(LungCapData$Age) #extract variable "Age" from LungCapData
89 LungCapData$Age
90
91 attaching data #pros:able to call variables by there names without "$"
    cons:put data in R's memory
92
93
94 mean(Age) #this would not work
95
96 attach(LungCapData)
97 detach(LungCapData)
98
99 class(Age) # return the type of the variable: integer, numeric, factor(/
    categorical)
100
101 levels(smoke) # show factors of a factor variable, like("yes" "no")
102
103 summary(LungCapData)
104
105 x = (c(0,1,1,1,0,0,0,0,0))
106 class(x)
107 x = as.factor(x)
108 class(x)
109
110 length(Age) # show how many obeservations are there under Age
111 Age[11:14]
112
113 mean(Age[Gender=="female"])
114
115 FemData = LungCapData[Gender=="female",] # create a subset of data
    containing only female and include all columns
116

```

```

117 MaleOver15 = LungCapData[Gender=="male" & Age>15,]
118
119 #Logic Statements and some other
120 temp = Age>15
121 temp[1:5] #return FALSE TRUE TRUE FALSE FALSE
122
123 temp2 = as.numeric(Age>15)
124 temp2[1:5] #return 0 1 1 0 0
125
126 FemSmoke <- Gender=="female" & Smoke=="yes"
127
128 MoreData <- cbind(LungCapData, FemSmoke) #cbind: add the new data in each
      row
129
130 rm(list=ls()) # remove all the thing in the workspace
131
132 #setting up working directory
133 getwd() #show working directory
134
135 setwd("/home/coupe")
136 setwd("~")
137
138 projectWD <- "/home/coupe/"
139 setwd(projectWD)
140
141 save.image("Marin.Rdata") #save this session .Rdata
142 load("Marin.Rdata")
143 load(file.choose())
144
145 #using scripts .R
146 #select the commands and Run it
147
148
149 #histograms
150 hist(LungCap)
151 hist(LungCap, freq=FLASE)
152 hist(LungCap, freq=F)
153
154 hist(LungCap, prob=T)
155
156
157 hist(LungCap, prob=T, ylim=c(0, 0.2))
158 hist(LungCap, prob=T, ylim=c(0, 0.2), breaks=7)
159 hist(LungCap, prob=T, ylim=c(0, 0.2), breaks=seq(from=0, to=16, by=2))
160 hist(LungCap, prob=T, ylim=c(0, 0.2), breaks=seq(from=0, to=16, by=2),
      main="Boxplot of Lung Capacity", xlab="Lung Capacity", las=1) #las=1:
      rotate the y axis

```

```

161 lines(density(lungCap))
162 lines(density(lungCap), col=2, lwd=3) #"col=2": set color as 2//2 is red;
    "lwd=3": the width of the line
163
164 #3.3 Normal Distribution, Z Scores, and Normal Probabilities in R
165 #X~N(75,5^2)
166 help(pnorm)
167 pnorm(q=70, mean=75, sd=5, lower.tail=T) #by lower.tail, we mean less
    than 70
168 pnorm(q=70, mean=75, sd=5, lower.tail=T) #In R, by default, lower.tail=T)
169 pnorm(q=85, mean=75, sd=5, lower.tail=F)
170
171 pnorm(q=1, mean=0, sd=1, lower.tail=T) #calculate the probability of Z,
    the standard normal
172
173 qnorm() #calculate quantile or percentagetile
174 qnorm(p=0.25, mean=75, sd=5, lower.tail=T) #find Q1
175
176 dnorm() #probability function
177 x <- seq(from=55, to=95, by=0.25)
178 dens <- dnorm(x, mean=75, sd=5)
179
180 plot(x, dens) #"x" work as variable of "dens"
181 plot(x, dens, type="l") #change the plot from dots to a line
182
183 abline(v=75) #vertical line
184
185 rnorm() #draw a random sample from a normally distributed population
186 rand <- rnorm(n=40, mean=75, sd=5) #with 40 observations
187
188 #3.4 t Distribution and t Scores in R
189 #mean=0, sd=1, tf=25 //25 degrees of freedom
190 help(pt)
191 #P(t > 2.3)
192 pt(q=2.3, df=25, lower.tail=F) # one-sided p-value(the tail area which is
    greater than 2.3 //higher tail)
193 pt(q=2.3, df=25, lower.tail=F) + pt(q=-2.3, df=25, lower.tail=T) #two-
    sided p-value
194 pt(q=2.3, df=25, lower.tail=F)*2 #two-sided p-value
195
196 #find t-value for 95% confidence
197 #value of t with 2.5% in each tail
198 qt(p=0.025, df=25, lower.tail=T)
199
200 help(pf) #F probability
201 help(pexp) #exponential probability
202

```

```

203 #4.9 Correlation and Covariance in R
204 #LungCap and Age
205 help(cor.test)
206
207 (Age, LungCap, main="Scatterplot", las=1)
208
209 cor(Age, LungCap, method="pearson") #pearson is the default; the order of
    "Age" and "LungCap" does not matter)
210 cor(Age, LungCap, method="kendall")
211 cor.test(Age, LungCap, method="pearson")
212 cor.test(Age, LungCap, method="spearman") #R will return a warning
213 cor.test(Age, LungCap, method="spearman". exact=F)
214
215 cor.test(Age, LungCap, method="pearson", alt="greater", conf.level=0.99)
    #"alt":by default, the alternative is a two-sided test
216
217 cov(Age, LungCap) #coveriance
218
219 pairs(LungCapData) #produce all possible pair-wise plots
220 pairs(LungCapData[,1:3])
221
222 cor(LungCapData[,1:3]) #matrix of correlation
223 cor(LungCapData[,1:3], method="spearman")
224
225 cov(LungCapData[,1:3]) #matrix of covariance
226
227 #5.1 Linear Regression in R
228 #LungCap is the outcome or dependent (Y) variable
229 plot(Age, LungCap, main="Scatterplot")
230 cor(Age, LungCap)
231
232 help(lm)
233 ?lm
234
235 mod <- lm(LungCap ~ Age) #1st variable is Y variableA, 2nd is X variable
236 summary(mod) #we have got understand this
237
238 attributes(mod) #return what's stored in mod
239 mod$coefficients
240 mod$coef
241 coef(mod)
242
243 abline(mod)
244
245 confint(mod) #confidence interval
246 confint(mod, level=0.99) # change confidence interval
247

```

```

248 anova(mod) #this ANOVA table corresponds to the f-test presented in the
      last row of the linear regression summary
249
250 #5.2 Checking Linear Regression Assumptions in R
251 plot(Age, LungCap)
252 mod <- lm(LungCap ~ Age)
253
254 #intercept slop
255 #standard deviation of residual errors is called Residual Standard Error
256
257 #R has a set of built in regression diagnostic plots to check regression
      through residuals, r of the 4 assumptions
258 plot(mod) # actually shows 4 plots
259 par(mfrow=c(2,2)) #if you want to show all plots on the screen at the
      same time
260 plot(mod)
261
262 #How non-constant variance will show up in a residual plot
263
264 #5.3 Multiple Linear Regression in R
265 help(lm)
266 model1 <- lm(LungCap ~ Age + Height)
267 summary(model1)
268
269 cor(Age, Height, method="pearson")
270 #return shows high collinearity, not good
271 confint(model1, conf.level=0.95) # set confidence interval
272
273 model2 <- lm(LungCap ~ Age + Height + Smoke + Gender + Caesarean)
274 summary(model2)
275
276 #
      #####
277 ##not from marin
278 #
      #####
279
280 edit(file = "test.R")
281 file.edit('foo.R')
282
283 20150203
284 as.POSIXITlt("YYYY-MM-DD")->A
285 as.Date(as.character("19990101"), "%Y%m%d")
286 # "%M": minute, "%m": month
287

```

```

288 %d   day as a number (0-31)
289 %a   abbreviated weekday
290 %A   unabbreviated weekday
291 %m   month (00-12)
292 %b   abbreviated month
293 %B   unabbreviated month
294 %y   2-digit year
295 %Y   4-digit year
296
297
298 my_list<-list()
299 my_list<-list(vec1,vec2,mat1)
300 mylist[[1]]
301 #call the first element in the list
302
303 mylist<-list(number=vec1, people=vec2,m1=mat1)
304 mylist$number
305 #call the first element in the list
306
307 while~
308 {
309     if~
310         next
311     ~
312 }
313
314 while~
315 {
316     if~
317         break
318     ~
319 }
320
321 mat[logical mat]
322 mat[mat<1]
323 a=sort(a)
324
325 sum(data)
326 rowSums(data)
327 colSums(data)
328 as.matrix(data)
329
330 #20150328
331
332 Normal Probability Plot(QQ plot)
333 -The normal probability plot is a graphical technique for assessing
    whether or not a data set is approximately normally distributed.

```



```

334 -The data are plotted against a theoretical normal distribution in such a
      way that the points should form an approximate straight line.
335 -Departures from this straight line indicate departures from normality.
336 -The relevant R functions are qqnorm() and qqline().
337
338 qqnorm(x)
339 qqline(x)
340
341 Time Series
342 The ts() function will convert a numeric vector into an R time series
      object. The format is ts(vector, start=, end=, frequency=) where start
      and end are the times of the first and last observation and frequency
      is the number of observations per unit time (1=annual, 4=quarterly, 12=
      monthly, etc.).
343 a=ts(b)
344 plot(a)
345
346 diff(x)
347 #takes differences between current value and its previous one.
348
349 A * B #Element-wise multiplication
350 A % * % B #Matrix multiplication(no space between operators)
351 solve(A, b) Returns vector x in the equation b = Ax (i.e., A-1b)
352 solve(A) Inverse of A where A is a square matrix.
353
354 solve(A, b) #Returns vector x in the equation b = Ax (i.e., A-1b)
355 solve(A) #Inverse of A where A is a square matrix.
356
357 lm(y ~ x_1 + x_2 + x_3)

```

R