

Prog. Brief Reference Sheet

Marco Sousa

JULY 2019

Contents

1	Introduction (future me, you better not suck bro)	2
2	RStudio	2
2.1	Elementary Vector Operations	2
2.2	Elementary Matrix Operations	4
2.3	Sequence Operations	5
2.4	Data Frames	6
2.5	Logical Operators and NA	6
2.6	Files and Linear Fitting	6
2.7	Conditionals (If,else)	7
2.8	Loops (For,While,Repeat,Apply)	8
2.9	Functions	9
2.9.1	Number of positive items in a list	9
2.9.2	Finding non-duplicate items in a list	10
2.9.3	Frequency of letters	10
2.9.4	List of primes within n	10
2.9.5	Second max value of a list	11
2.10	Probability Distributions	11
2.10.1	Uniform distribution	11
2.10.2	Normal distribution	12
2.10.3	Binomial distribution	12
2.10.4	Poisson distribution	12
2.10.5	t distribution	12
2.10.6	Sampling data from a list	12
2.11	Graphs and Descriptive Stats	12
2.12	Hypothesis Testing	16
2.13	Logistic Regression	16
2.14	dplyr	17
2.14.1	Some comments	17
2.14.2	select()	18
2.14.3	filter()	18
2.14.4	Pipe Operator	18
2.14.5	arrange()	18
2.14.6	mutate()	19
2.14.7	summarise()	19
2.14.8	groupby()	19
2.14.9	Filtering out NA rows	19
2.15	R Comments	20
3	Python	21
4	Java	22
5	Notes to self	22

1 Introduction (future me, you better not suck bro)

The purpose of this document is simply to add *quick and brief* notes/suggestions regarding programming.

USE THE TABLE OF CONTENTS PAGE. Click them to skip to the section.

You can technically copy paste, but I'd be careful, sometimes it doesn't ASCII correctly.

Not necessarily a tutorial, but can be very informative.

Also note that my R code looks like file .R code; not console line code (with > and +).

Special thanks to the UMassD professors for their instruction.

Special thanks to Prof Yan Donghui who much of this is taken from, and also for his teaching.

This will be updated as I have time x.x

Thank you for your time, and enjoy.

2 RStudio

A good intro source for R-PRogramming:

<https://sites.google.com/site/rprogrammingumassd/home>

Another source for more R Programming:

<https://www.datacamp.com/>

2.1 Elementary Vector Operations

Vectors are [1:n]

Creating a vector:

```
1 mydata <- c(2.9, 'string', TRUE, value, value,...)
```

List of numbers:

```
1 numberList <- 25:30
2 # Note: Will contain 25,26,27,28,29,30 (BOTH start and endpoint)
```

Access a value in VECTOR:

```
1 mydata[1], or mydata[1:n]
```

Can test entire vector using logic

```
1 mydata > 3
2 #Result: [1] FALSE TRUE TRUE TRUE TRUE FALSE FALSE FALSE FALSE
```

Adding names:

```
1 names(mydata) <- c('name1', 'name2',...);
```

Changing a particular name:

```
1 names(mydata)[1] <- 'newname'
```

letters is a pre-defined R vector of the 26 lowercase letters

```
1 letters [1:5]
2 #Result: [1] "a" "b" "c" "d" "e"
```

Finding a vector value associated with a name!

```
1 mydata['name2']
```

Important note about VECTOR logic!

```
1 mydata > 3 ##returns TRUE FALSE vector
2 mydata[mydata > 3] ##returns VALUES which are TRUE
```

MODE tells you the type of values in a vector

```
1 mode(mydata)
```

DIM returns the dimensions of the data structure

```
1 dim(mydata)
2 #Result: 2 5 (if it's a 2 by 5, for example)
```

DIM can also make a vector into a data frame structure

```
1 dim(mydata) <- c(m,n) (makes it an m by n structure)
2 dim(mydata) <- NULL (makes structure into vector)
```

ADDING/INSERTING a vector value is complicated...

Appending is discouraged because it copies the vector and adds the last part repeatedly in a for loop. Consider Python/Java for list appending in for loops.

DELETING a vector value

```
1 mydata2[-c(1,3)]; (deletes at location 1 and 3)
```

Making a list of vectors (list is a data structure)

```

1 mydatalist <- list(mydata,mydata1,mydata2,...)
2 ## RESULT: [[1]] [1] 1 3 6, [[2]] [1] "red" "blue" "green" "black ",[[3]]...

```

Accessing a list

```

1 mydatalist [[1]][2]

```

LENGTH can be found of vectors

```

1 length(mydata)

```

2.2 Elementary Matrix Operations

Matrix function:

```

1 matrix(c(value,value ,...) ,nrow,ncolumn,byrow)
2 A<-matrix(c(1,2,4,7,5,9),2,3, byrow=T); #example

```

Class tells object type, so matrix or vector...etc.

```

1 class(A)

```

typeof tells MODE of the ELEMENTS of a matrix

```

1 typeof(A)

```

MATRIX dimensions

```

1 dim(A),nrow(A), nrow(A) #fairly obvious

```

Changing column name or ROW name

```

1 colnames(A) <- c("COL1", "COL2",...)
2 rownames(A) <- c("ROW1", "ROW2",...)

```

Transpose a matrix:

```

1 A <- t(A)

```

Deleting a column or columns:

```

1 ##Deleting a single column
2 Data$genome <- NULL
3
4 ##Deleting multiple columns
5 Data[1:2] <- list(NULL) #If you do not include the list, it will not work on the frame
6 #or
7 Data <- Data[,-(2:3)] # vector
8 Data <- Data[,-(2:3),drop=FALSE] # still a data.frame

```

Matrix multiplication:

```

1 ##Element-wise multiplication:
2 A*C (sizes 2x3 * 2x3)
3 ##Linear Algebraic multiplication:
4 A %*% B

```

Binding rows/columns:

```

1 ##rbind STACKS rows on top of one another of same dimension
2 rbind(A,C)
3 #Result: puts row of A on top of B in single matrix
4
5 ##cbind STACKS columns similarly

```

Return/Replace Diagonal values along matrix

```

1 ##Return diagonal values
2 diag(F)
3 ##Replace diagonal values
4 diag(F)<-c(1,1);

```

2.3 Sequence Operations

Create a sequence

```

1 s4<-seq(length=51,from=-1,by=0.2)

```

Remember: Can return value at index n with vector[n]

Repeat a number x, n times:

```

1 rep(x,n)
2 #Fill a matrix:
3 B <- matrix(rep(1,6),3,2)

```

A sequence can be repeated:

```

1 This repeats x 2 times and stores into y:
2 x<-1:4;
3 y<-rep(x,2)
4 #Result:(1234)(1234)
5 #This repeats each x i times.
6 i<-rep(2,4)
7 z<-rep(x,i)
8 #Result:(11)(22)(33)(44)
9 #Note the difference from above: rep(x,2) versus rep(x,rep(2,4))
10 w<-rep(x,x)
11 #Result: 1223334444
12 x<-rep(rep(1:4,rep(2,4)),3)
13 #Looks complicated by remember the TWO parts of the rep function
14 #Result: (11223344)(11223344)(11223344) – Notice how it repeats the inner 3 times

```

2.4 Data Frames

A data frame is a matrix-like structure whose columns have the same length but may be of differing data types (e.g., numeric, logical, factor and character and so on).

Creating a data frame:

```

1 dataFrame <- data.frame(x = 1:2, y = 1:10, z=rep('A', 10))
2 #x,y,z are column names. It will fill to the largest column's dimension (this means columns must align, so x
  =1:3 would not work for example.

```

2.5 Logical Operators and NA

Logical Operators:

```

1 <, <=, >, >=, ==, !=
2 #Element-wise AND: &, Element-wise OR: |
3 #Logical AND: &&, Logical OR: ||
4 ##The logical operators only consider the first element of a vector

```

When an element or value is “not available,” the location is reserved as NA.

In general, any operation on an NA is NA.

The function `is.na(x)`, where `x` is a vector, gives a logical vector of the same length as `x` with values set to be true if the corresponding element in `x` is NA.

The function `is.na()` is also good for finding locations in which a value is NA.

2.6 Files and Linear Fitting

Read in files, such as .dat, .txt, .csv (Comma separated files), SPSS, and so on

```

1 tmp<-read.table("car.txt", header=TRUE, sep=",")

```

Note that excel files are xlsx, and should be saved as .csv first

```
1 dataFile<-read.csv("storeSales.csv",na.strings="N/A");
```

Other files can be imported from top right (Import Dataset)

Take a column and assign it to a vector:

```
1 x<-tmp[,1]
2 y<-tmp[,2]
```

Sorting vectors; index.return keeps track of the index for each value

```
1 xs<-sort(x, index.return =TRUE);
```

lm collects data to fit linear models

```
1 mylm<-lm(y ~ x);
```

The plot function:

```
1 plot(x,y,xlab="Age", ylab="Price",cex.lab=1.5,cex.axis=1.5, bty="l",pch=20, font.axis=2, font.lab=2);
```

Create the fitted line from lm

```
1 points(xs$x, fitted (mylm)[xs$ix],type="l",lwd=2)
2 #Note that xs$x and xs$ix select x from xs, and ix from xs
```

2.7 Conditionals (If,else)

```
1 if (test Expression) {
2   statement
3 }
```

```
1 if (x<=2) {print("less than 3")}
2 else if (x>=4) {print("greater than 3")}
3 else print("equal to 3")
```

IMPORTANT NOTE: the } MUST exist on the next line

Since R reads line by line, if it is on the same line, it will COMPLETE the statement and the next

line will error.

2.8 Loops (For,While,Repeat,Apply)

The for loop operates similarly to pythonic loops w/ a sequence
FOR Loop:

```
1 x <- seq(1, 100, by=1);
2 xsquared = NULL
3 for (n in 1:100)
4 {
5   xsquared[n] = x[n]^2
6 }
```

A minor note about the for loop:

```
1 ##The -1 will not subtract JUST the variable; add parenthesis to do so.
2 for (i in 1:variable -1) {
3   #do things
4 }
```

WHILE Loop:

```
1 x<-0;
2 while (x < 10) {
3   x <- x + 4;
4   print (x);
5   if (x == 8) {
6     break;
7   }
8 }
```

REPEAT Statement

```
1 sum <- 1;
2 repeat
3 {
4   sum <- sum + 2;
5   print(sum);
6   if (sum > 11)
7     break;
8 }
```

apply(x, MARGIN, FUN, ARGS) for 2-dimensional data

```
1 x <- cbind(x1 = 3, x2 = c(4:1, 2:5));
2 dimnames(x)[[1]] <- letters[1:8];
3 apply(x, 2, mean);
4 col.sums <- apply(x, 2, sum);
5 row.sums <- apply(x, 1, sum);
```

Apply a function to each cell of a ragged array
tapply(vector, factor, FUN) Apply a funct

```
1 x<-c(1,2,3,4,5,6,7,8);
2 names(x)<-c("a","a","b","a","b","a","b","a");
3 tapply(x,factor(names(x)),mean);
```

lapply() and apply()

```
1 x <- list(a = 1:10, b = (-3:3)^2, lv = c(TRUE,FALSE,FALSE,TRUE));
2 # compute the list mean for each list element
3 lapply(x, mean);
4 sapply(x,mean);
```

2.9 Functions

General FUNCTION:

```
1 myfunction<-function(argument1, argument2, ..., argumentn )
2 {
3     R statements; ##Body of a function
4     return(output);
5 }
```

Some notes about functions:

You can make a function in the CONSOLE. It can be edited with fix(func)

However, it is better to save functions in .R files!

Also note that RStudio can run Python/C++ scripts

5 R function examples are provided below:

2.9.1 Number of positive items in a list

```
1 nPositives<-function(x)
2 {
3     n<-length(x);
4     counts<-0;
5     for(i in 1:n)
6     {
7         if(x[i]>0) {
8             counts<-counts+1;
9         }
10    }
11    return(counts);
12 }
```

2.9.2 Finding non-duplicate items in a list

```
1 nonDup<-function(x)
2 {
3   y<-NULL;
4   while(length(x)>=1)
5   {
6     tmp<-(1:length(x))[x==x[1]];
7     if (length(tmp)==1) y<-c(y,x[1]);
8     x<-x[-tmp];
9   }
10  return(y);
11 }
12 >x<-c(1,2,3,2,5,3,6);
13 >nonDup(x)
```

2.9.3 Frequency of letters

```
1 letterFreq<-function(x)
2 {
3   n<-length(x);
4   fTable<-matrix(0,26,1);
5   for(i in 1:n)
6   {
7     idx<-(1:26)[letters==tolower(x[i])];
8     fTable[idx] <- fTable[idx] + 1;
9   }
10  return(data.frame(letters,fTable));
11 }
12 >x<-c("D","a","t","a","s","c","i","e","n","c","e","i","s","f","u","n");
13 >letterFreq(x)
```

Note:idx[-(1:26)[letters==tolower(x[i])];

The above BRACKET RETURNS the number index in which the case is TRUE.

2.9.4 List of primes within n

```
1 is.prime<-function(n)
2 {
3   if (n==2) return(1);
4   if (n==1 || n%%2==0) {
5     return(0);
6   }
7   m<-floor(sqrt(n))+1;
8   for(i in 2:m) { if (n %% i ==0) return(0);
9   }
10  return(1);
11 }
12
13 print.primes<-function(n, h)
14 {
```

```

15 | z<-tapply(1:n,as.factor(1:n),is.prime);
16 | primes<-(1:n)[z==1];
17 | plot(primes, col="blue", type="b", pch=5,
18 |       xlab="Index",ylab="");
19 | text(1:length(primes),primes+h,labels=primes);
20 | return(length(primes));
21 | }
22 | >print.primes(100,2);

```

Note the above uses tapply, which applies a function to each term in a designated sequence. Also note n is how far out in primes and h simply adjusts labels.

2.9.5 Second max value of a list

```

1 | secMax<-function(x)
2 | {
3 |   infty<-10^20; #define this as a very small number;
4 |   n<-length(x);
5 |   if (n <2) return(-infy);
6 |   mysec<- -infy;
7 |   mymax<-x[1];
8 |   for (i in 2:n)
9 |   {
10 |     if (x[i] > mymax)
11 |     {
12 |       mysec<-mymax;
13 |       mymax<-x[i];
14 |     }
15 |     else
16 |     {
17 |       if (x[i] > mysec)
18 |       {
19 |         mysec<-x[i];
20 |       }
21 |     }
22 |   }
23 |   return(mysec);
24 | }

```

2.10 Probability Distributions

2.10.1 Uniform distribution

In interval [a,b]

```

1 | runif(n, min = a, max = b)

```

2.10.2 Normal distribution

With mean `u` and standard deviation `sigma`

```
1 rnorm(n, mean = u, sd = sigma)
```

2.10.3 Binomial distribution

With size `n` and success probability `prob`

```
1 rbinom(n, size, prob)
```

2.10.4 Poisson distribution

With parameter `lambda`

```
1 rpois(n, lambda)
```

2.10.5 t distribution

With `df` degrees of freedom, and non-central parameter (by default 0)

```
1 rt(n, df, ncp)
```

2.10.6 Sampling data from a list

```
1 sample(x, size, replace = FALSE, prob = NULL)
```

`x`: Either a vector of one or more elements from which to choose, or a positive integer.

`n`: A positive number, the number of items to choose from. See ‘Details.’

`size`: A non-negative integer giving the number of items to choose.

`replace`: Should sampling be with replacement?

`prob`: A vector of probability weights for obtaining the elements of the vector being sampled.

2.11 Graphs and Descriptive Stats

```
1 ##Histogram:
2 hist(rnorm(100,0,1),breaks=8);
3 ##Boxplot:
4 x<-read.table("newbedfordTempMonth.Data");
5 names(x)<-c("Jan", "Feb", "Mar", "Apr", "May", "Jun", "Jul", "Aug", "Sep", "Oct", "Nov", "Dec");
```

```

6 | boxplot(x[,1], col="orange", cex.lab=1.5, cex.axis=1.5, font=2);
7 | boxplot(x, col="orange", cex.lab=1.5, cex.axis=1.5, font=2);
8 |
9 | ##Stem-and-leaf plot:
10 | x<-c (2.21,2.24,2.28,2.33,2.37,2.52,2.79,2.81,2.82,2.82,2.83,2.83) ;
11 | stem(x, scale=2);
12 |
13 | ##QQ-plot
14 | qqnorm(rnorm(100,0,1));
15 | abline(0,1);
16 |
17 | qqnorm(runif(100,0,1));
18 | abline(0,1);
19 |
20 | ##Uncomment when running for the first time
21 | ##install.packages("moments");
22 | library(moments);
23 |
24 | x<-read.table("eruption.Data");
25 | kurtosis(x);
26 | skewness(x);

```

```

1 | #Data Visualization with some examples:
2 |
3 | numbers <- c(3574,1328,6548,1316,1053,626);
4 | colors <- c("blue", "green", "orange", "cyan", "cadetblue4", "darkmagenta");
5 | cats <- c("Connecticut", "Maine", "Massachusetts", "New Hampshire", "Rhode Island", "Vermont");
6 |
7 | ##Barplot
8 | barplot(numbers, names.arg=cats, col=colors, cex.lab=1.5, cex.axis=1.5, cex.names=1.5, font=2, font.main=3);
9 |
10 | ##Pie Chart
11 | pie(numbers, labels=cats, col=colors, font=2, font.main=3);
12 |
13 | ##Pareto Chart
14 | install.packages("qcc");
15 | library(qcc);
16 | names(numbers) <-cats;
17 | pareto.chart(numbers, ylab= "Frequency", main= , col=colors, font=2, cex.lab=1.5, cex.axis=1.5);
18 |
19 | ##Box Plot (Gives 5 Number Summary, Tukey)
20 | x<- read.table(...);
21 | boxplot(x, col= orange , cex.lab=1.5, cex.axis=1.5, font=2);
22 |
23 | ##Dotplot
24 | stripchart(x, method= stack , pch=19, col= blue , offset=0.5, ,cex.lab=1.5,cex.axis=1.5);
25 |
26 | ##Histogram
27 | hist(x, main= , xlab= Length (seconds) , ylab= Count , breaks=20, freq=TRUE, cex.lab=1.5, cex.
axis=1.5);
28 |
29 | ##Stem and Leaf Plot
30 | stem(x,scale=2);
31 |
32 | ##Grouped-Bar Chart
33 | #What the Data looks like:

```

```

34 Soybean Spect ImgSeg Heart Wine WDBC Robot Madelon
35 CF 84.43 68.02 48.24 68.26 79.19 88.70 55.12 41.20
36 RP 71.83 61.11 47.69 60.54 70.79 85.41 55.19 35.50
37 BC2 72.34 56.28 49.91 59.10 70.22 85.38 50.21 35.37
38 EA 76.59 56.55 51.31 59.26 70.22 85.41 50.32 37.19
39
40 cfr<-read.table("CFrt.Data", header=TRUE);
41 row.names<-c("CF","RP","BC2","EA");
42 colors<-c("blue","green","pink","orange");
43 barplot(as.matrix(cfr),main="", ylab = "Rate", cex.lab=0.8, cex.main=1.4, xpd=FALSE, beside=TRUE, col=
    colors,ylim=c(30,90));
44 legend("topright", legend=c("CF","RP","BC2","EA"), fill=colors);
45
46 ##Grouped-Bar Chart – UC Enrollment example:
47 NYR<-15; NRACE<-9; x<-matrix(0,NYR,NRACE+2);
48 x[1,]<-c(1996,44796,6477,5883,20251,1510,2807,71008,6992,6994,166718);
49 x[2,]<-c(1997,46237,6350,6211,20072,1491,3142,71658,7730,6971,169862);
50 x[3,]<-c(1998,47066,6023,6499,19764,1399,3293,71009,8196,10394,173643);
51 x[4,]<-c(1999,48462,5722,6782,19739,1265,3469,72606,8200,12165,178410);
52 x[5,]<-c(2000,49735,5654,7043,20083,1144,3618,73714,9139,13225,183355);
53 x[6,]<-c(2001,52339,5572,7455,21386,1140,3938,75647,10168,14258,191903);
54 x[7,]<-c(2002,55852,5807,7973,22976,1204,4290,77833,11165,14197,201297);
55 x[8,]<-c(2003,57951,6103,8232,24771,1232,4733,79145,11708,14516,208391);
56 x[9,]<-c(2004,58762,6049,8120,24992,1208,5006,77655,11816,14301,207909);
57 x[10,]<-c(2005,60297,6083,8214,25518,1188,5182,77374,11456,13768,209080);
58 x[11,]<-c(2006,63064,6281,8217,26907,1196,5424,78453,11469,13287,214298);
59 x[12,]<-c(2007,65534,6647,8437,28613,1258,5700,78957,12182,12706,220034);
60 x[13,]<-c(2008,68031,7137,8550,30981,1264,5811,79080,12996,12190,226040);
61 x[14,]<-c(2009,69236,7508,8653,33398,1344,5994,79134,13809,11452,230528);
62 x[15,]<-c(2010,69871,7990,8636,35993,1539,17544,77715,15176,0,234464);
63
64 #Combine Other and No response (i.e., column 7 and 10==>column8)
65 #Combine Native and Fillipino (i.e., column 4 and 6==>column4)
66 y<-matrix(0,NYR,NRACE+2);
67 for(i in 1:5) {y[,i]<-x[,i];}
68 y[,4]<-y[,4]+x[,6];
69 y[,6]<-x[,8];
70 y[,7]<-x[,7]+x[,9]+x[,10];
71 y[,8]<-x[,11];
72 z<-y[,2];y[,2]<-y[,4];y[,4]<-z; ##Swap Asian and Native
73 for(i in 1:NYR) { y[i,]<-y[i,]*100/y[i,8];}
74 z<-t(y[,2:7]);
75 colors<-c("blue","red","gold","grey","cyan","brown");
76 ethnics<-c("Other/Unknown","White","Hispanic","Asian",
77
78 b<-barplot(z, main="UC enrollment 1996–2010", col=colors, las=2, cex.names=0.5, axes=FALSE, legend.text=
    FALSE, xlim=c(0,28),width=1);
79 legend("bottomright", bty="n", cex=0.9, legend=ethnics, fill=rev(colors));
80
81 ##Pyramid Plot – US Population example
82
83 library(plotrix);
84 #The Pop. Data
85 m1950<-c(8236164,6714555,5660399,5311342,5606293,5972078,5624723,5517544,5070269,4526366,
86 4128648,3630046,3037838,2424561,1628829,1001798,504958,236828);
87 m2000<-c(9810733,10523277,10520197,10391004,9687814,9798760,10321769,11318696,11129102,9889506,
88 8607724,6508729,5136627,4400362,3902912,3044456,1834897,1226998);
89 f1950<-c(7927407,6485130,5458869,5305256,5875535,6270182,

```

```

90 5892284,5728842,5133704,4544099,4143540,3605074,3021637,2578375,1783120,1150609,620386,340073);
91 f2000<-c(9365065,10026228,10007875,9828886,9276187,9582576,
92 10188619,11387968,11312761,10202898,8977824,6960508,5668820,5133183,4954529,4371357,3110470,3012589);
93
94 agelabels<-c("0-4", "5-9", "10-14", "15-19", "20-24", "25-29",
95 "30-34", "35-39", "40-44", "45-49", "50-54", "55-59",
96 "60-64", "65-69", "70-74", "75-79", "80-84", "85+");
97 GPS<-18;
98 ma1950<-matrix(0,1,GPS); ma2000<-matrix(0,1,GPS);
99 fa1950<-matrix(0,1,GPS); fa2000<-matrix(0,1,GPS);
100 pop1950<-sum(m1950+f1950); pop2000<-sum(m2000+f2000);
101 ma1950p<-m1950/pop1950; fa1950p<-f1950/pop1950;
102 ma2000p<-m2000/pop2000; fa2000p<-f2000/pop2000;
103 males<-matrix(0, GPS,2);
104 males[,1]<-t(ma1950p); males[,2]<-t(ma2000p);
105 females<-matrix(0,GPS,2)
106 females[,1]<-t(fa1950p); females[,2]<-t(fa2000p);
107
108 myPyramid<-pyramid.plot(males*100, females*100, labels=agelabels, lxc=c("blue", "green"), rxc=c("blue", "
109 green"), laxlab=seq(0,10,by=2), raxlab=seq(0,10,by=2), top.labels=c("Males", "Age", "Females"), gap=3);
110 mtext("US population by age and gender",3,2,cex=1.5);
111 legend(par("usr")[1], GPS,c("1950", "2000"),fill=c("blue", "green"));
112
113 #To restore the margins and background
114 par(mar=myPyramid,bg="transparent");
115
116 ##Bubble Chart – Crime Example
117
118 crime<-read.csv("http://datasets.flowingdata.com/crimeRatesByState2005.tsv", header=TRUE, sep="\t");
119 radius<-sqrt( crime$population/pi);
120 symbols(crime$murder, crime$burglary, circles=radius,inches=0.35, fg="white", bg="red", xlab="Murder Rate",
121 ylab="Burglary Rate");
122 text(crime$murder, crime$burglary, crime$state, cex=0.5);
123
124 ##Bubble Chart – US Population by States
125 library( maps);
126 library( ggplot2);
127 #Get US map
128 usa <- map_data("state");
129 myData = data.frame(name=c("FL", "CO", "CA", "MA", "NY"),
130 lat=c(29,39,38,42,43),
131 long=c(-82,-105,-120,-71,-75),
132 pop=c(19.3,5.2,38.0,6.6,19.6) )
133
134 #Not exactly sure how this works out...
135 ggplot() + geom_path(data = usa, aes(x=long, y=lat, group = group)) + geom_point(data = myData, aes(x=long
136 , y=lat, size=pop), color = "red");
137
138 ##More data to analyze:
139 Age Male (1950) Male (2000) Female (1950) Female (2000)
140 0-14 20611118 30854207 19871406 29399168
141 15-29 16889713 29877578 17450973 28687649
142 30-44 16212536 32769567 16754830 32889348
143 45-64 15322898 30142586 15314350 31810050
144 65-84 5560146 13182627 6132490 17569539
145 85+ 236828 1226998 340073 3012589

```



```
145 #Can make a scatterplot, or scatterplot by proportion, or by pyramid!
```

2.12 Hypothesis Testing

```
1 ##t-test
2 x<-c(49, 50, 47.5, 49.6, 48.6, 49.0, 48.6, 49, 49.4, 50.2);
3 t.test(x, alternative="two.sided", mu=50);
4
5 ##Wilcox test
6 x <- c(1.83, 0.50, 1.62, 2.48, 1.68, 1.88, 1.55, 3.06, 1.30);
7 y <- c(1.67, 0.647, 1.35, 2.05, 1.06, 1.29, 1.06, 3.14, 1.09);
8 wilcox.test(x, y, paired = TRUE, alternative = "greater");
9
10 ##Kolmogorov-Smirnov test
11 ks.test(rnorm(100,0,1), runif(100,-4,4));
12
13 ##Shapiro-Wilk test for normality
14 shapiro.test(rnorm(100,0,1));
15 shapiro.test(runif(100,-4,4));
```

2.13 Logistic Regression

A tutorial: <https://www.datacamp.com/community/tutorials/logistic-regression-R>

An alternative tutorial:

http://rstudio-pubs-static.s3.amazonaws.com/14336_8e6cdabdc9434030b194dd9ba7b22ea2.html

```
1 ## download for testing Smarket dataset
2 install.packages("ISLR");
3 library(ISLR)
4 ## Viewing the distribution of the data?
5 par(mfrow=c(1,8))
6 for(i in 1:8) {
7   hist(Smarket[,i], main=names(Smarket)[i])
8 }
9
10 ## To show correlations! ***
11 install.packages("corrplot")
12 library(corrplot)
13 correlations <- cor(Smarket[,1:8])
14 corrplot(correlations, method="circle")
15
16 ## pairs() will create a scatterplot matrix
17 pairs(Smarket, col=Smarket$Direction)
18
19
20 ## LOGISTIC MODEL! ***
21 glm.fit <- glm(Direction ~ Lag1 + Lag2 + Lag3 + Lag4 + Lag5 + Volume, data = Smarket, family = binomial)
22 summary(glm.fit)
23
```

```

24 ## Logistic prediction function! ***
25 glm.probs = predict(glm.fit, type = "response")
26 glm.probs[1:n] ## Views selected probabilities in vector
27
28 ##making predictions from each probability
29 glm.pred <- ifelse(glm.probs > 0.5, "Up", "Down")
30
31 ##Creating a table of correct/wrong predictions.
32 ##Smarket[,9] accesses the 9th column (the binomial column; Direction)
33 table(glm.pred, Smarket[,9])
34
35 ##mean value to see how well logistic fit is; [,9] is Direction.
36 mean(glm.pred == Smarket[,9])
37
38 ## train allows you to select rows to test with.
39 train = Year<2005
40 glm.fit <- glm(.. subset = train)
41 glm.probs <- predict(...newdata = Smarket[!train,])
42 Direction.2005 = Smarket$Direction[!train]
43 table(glm.pred, Direction.2005)
44 mean(glm.pred == Direction.2005)

```

There are necessary conditions for Logistic regression and this site MAY outline them: <https://www.statisticshowto.datasciencecentral.com/assumptions-conditions-for-regression/>
Some additional notes:

```

1 ## The par function allows you to edit how many figures can be plotted, and the margin.
2 ## par(mar) tells you the current plot margins
3 par("mar")
4 ## dev.off() will reset par in case there is an issue w/ dimensions
5 ## Smarket[,n] returns the WHOLE column for index n ***

```

2.14 dplyr

A Tutorial: https://genomicsclass.github.io/book/pages/dplyr_tutorial.html

```

1 install.packages("dplyr")
2 library(dplyr)

```

2.14.1 Some comments

dplyr is for cleaning and ordering data. There are other ways, but I like this one.

```

1 #Know that you don't have to use the following to save to a new dataframe
2 { . ->> newDataFrame }
3 #Instead, you can also just do it normally:
4 newDataFrame <- myData %>% select( blah blah blah)

```

Since this is an old version, my examples might instead save in the former, where the latter is much nicer.

2.14.2 select()

```
1 ## Select chooses COLUMNS with arguments
2 ## select(datafile, columnName1,columnName2,...)
3 sleepData <- select(msleep, name, sleep_total)
4 head(sleepData)
5 ## can select with other arguments
6 head(select(msleep, -name))
7 head(select(msleep, name:order))
8 head(select(msleep, starts_with("sl")))
9
10 ##Some common arguments:
11 ##ends_with() = Select columns that end with a character string
12 ##contains() = Select columns that contain a character string
13 ##matches() = Select columns that match a regular expression
14 ##one_of() = Select columns names that are from a group of names
```

2.14.3 filter()

```
1 ## filter selects ROWS with argument
2 ## filter(data, argument to select row)
3 filter(msleep, sleep_total >= 16)
4 filter(msleep, order %in% c("Perissodactyla", "Primates"))
5
6 ##See boolean operators (>, <, >=, <=, !=, %in%)
```

2.14.4 Pipe Operator

```
1 head(select(msleep, name, sleep_total))
2
3 msleep %>%
4   select(name, sleep_total) %>%
5   head
6
7 ## Are effectively the same; does multiple operations orderly left to right pipeing output to next as input.
8 ## The below suggests one way to save data in between piping
9 {. ->> newDataFrame }
```

2.14.5 arrange()

```
1 ## To arrange (or ORDER) ROWS by a particular column
2 head(arrange(msleep, order))
3 msleep %>% arrange(order) %>% head
4
5 msleep %>%
6   select(name, order, sleep_total) %>%
7   arrange(order, sleep_total) %>%
```

```

8 |     filter (sleep_total >= 16)
9 |
10 | ##This orders rows BIGGEST to SMALLEST (because of desc) for Length then Width.
11 | arrange( iris , desc(Sepal.Length),desc(Sepal.Width))

```

2.14.6 mutate()

Adds new COLUMNS to the data frame

```

1 | ## In this occasion we are STORING rem_proportion as name of new column BASED on sleep_rem / sleep_total.
2 | ## We also saved it as a new data frame with the last line!
3 | msleep %>%
4 |   mutate(rem_proportion = sleep_rem / sleep_total) %>%
5 |   { . ->> newmsleep }

```

2.14.7 summarise()

```

1 | ##summarise() provides summary statistics.
2 | msleep %>%
3 |   summarise(avg_sleep = mean(sleep_total),
4 |             min_sleep = min(sleep_total),
5 |             max_sleep = max(sleep_total),
6 |             total = n())
7 |
8 | # Note the n() function counts the number of observations

```

2.14.8 groupby()

```

1 | ##group_by() splits the data frame by a variable, applies a function to the individual data frames, and then
  |   combine the output.
2 | ## The following example outputs a frame of summary statistics for each taxonomic order.
3 | msleep %>%
4 |   group_by(order) %>%
5 |   summarise(avg_sleep = mean(sleep_total),
6 |             min_sleep = min(sleep_total),
7 |             max_sleep = max(sleep_total),
8 |             total = n())

```

2.14.9 Filtering out NA rows

```

1 | ## The below filters out na values for column a.
2 | filter (!is.na(column.a))
3 |
4 | ## Here's an example in use:

```

```
5 | msleep %>%  
6 |   filter (!is.na(sleep_rem)) %>%  
7 |   summarise(  
8 |     avg_sleeprem = mean(sleep_rem),  
9 |     totalrem = n())
```

2.15 R Comments

Comments done with .

Ctrl + Shift + C will comment out a highlighted block.

An introduction to docstrings for R functions can be found here: [/urlhttps://cran.r-project.org/web/packages/docstring/](https://cran.r-project.org/web/packages/docstring/)

3 Python

TBA

Wish I had more time to fill this out; love Python ;-;

4 Java

yikes

5 Notes to self

I should look at my old data and relearn ANOVA analysis and Tukey Analysis and correlations and such...

sum stuff I gott test

```
1 testPlot <- ggplot() + geom_path(data = usa, aes(x=long, y=lat, group = group)) + xlab("Longitude") + ylab("Latitude") + ggtitle("2010 Population Bubble Chart")
2
3
4
5 testPlot + geom_point(data = newtestSelect, aes(x=Longitude, y=Latitude, size=Pop2010), color = "red") +
6   geom_text(data=newtestSelect, aes(x=Longitude,y=Latitude), label=newtestSelect$Abbreviation);
7
8 ##I also can slice out 2 and 11 to make a better map:
9
10
11 testPlot + geom_point(data = newtestSelect2, aes(x=Longitude, y=Latitude, size=Pop2010), color = "red") +
12   geom_text(data=newtestSelect2, aes(x=Longitude,y=Latitude), label=newtestSelect2$Abbreviation, nudge_x = 0, nudge_y = 0, fontface = "bold");
```

```
1 #Pairs example:
2 data(iris);
3 pairs(iris[1:4], main = "", pch = 21,
4       bg = c("red", "green3", "blue")[unclass(iris$Species)]);
5
6 #Chernoff Faces example:
7 ##Uncomment if running for the first time
8 ##install.packages("aplpack");
9 library(aplpack);
10 ##May need to fix Tcl/Tk
11 ##May need to install XQuartz/X11
12 data(iris);
13 ##postscript("irisFaces.eps")
14 faces(iris[1:4],
15
16       labels=c(rep("ST",50),rep("VS",50), rep("VG",50)));
17
18 ##dev.off();
19
20 #Chernoff Faces ex. 2:
21 USArrests<-read.table("USArrests.Data", header=TRUE);
22 faces(USArrests);
23
24
25 #Radial plot example R:
26 ##Uncomment if used for the first time
27 ##install.packages("plotrix");
```

```

28 library(plotrix);
29 radmat<-matrix(sample(1:4,40,replace=TRUE),nrow=4);
30 radial.plot(radmat,rp.type="l",
31 radial.pos=seq(0,20*pi/11.1,length.out=10),
32 label.pos=seq(0,20*pi/11.1,length.out=10),start=pi/2,
33 clockwise=TRUE, labels=2001:2010,radial.lim=c(0.2,4),
34 cex=1.5, font=2, lty=1, lwd=2);
35 legend(-5,4,c("Black", "Red", "Green", "Blue"),col=1:4,lty=1);

```

```

1  ##CLUSTERING TECHNIQUES
2
3  ##Kmeans clustering
4  kmeans(x, centers, iter.max = 10, nstart = 1, algorithm = c("Hartigan-Wong", "Lloyd", "Forgy", "MacQueen")
5    , trace=FALSE)
6  # centers = either # clusters, or a set of initial (distinct) cluster centers
7  # iter.max = maximum # iterations allowed
8  # nstart = # random starts for K-means if centers is a number
9  # algorithm determines the choice of initial cluster centers
10
11 #Example: K-means clustering of 2-F Gaussian Data
12 require(graphics);
13 x <- rbind(matrix(rnorm(100, sd = 1), ncol = 2), matrix(rnorm(100, mean = 1, sd = 1), ncol = 2));
14
15 colnames(x) <- c("x", "y");
16 cl <- kmeans(x, 2);
17 colors<-rep("blue",nrow(x)); colors[cl$cluster==2]<-"red";
18 pchs<-rep(17,nrow(x)); pchs[cl$cluster==2]<-19;
19 plot(x, col = colors, pch=pchs);
20 txts<-c(rep("1",50), rep("2",50));
21 text(x, labels=txts);
22
23 #Example: K-means clustering on Iris data (pairwise)
24 data(iris);
25 mykms<-kmeans(iris[,1:4],centers=3);
26 pairs(iris[,1:4], main = "", pch = 21, bg = c("red", "green3", "blue")[mykms$cluster]);
27
28 # Agglomerative Clustering
29
30 library(cluster);
31 agnes(x, diss = inherits(x, "dist"), metric = "euclidean", stand = FALSE, method = "average", par.method,
32   keep.diss = n < 100, keep.data = !diss);
33 # metric = euclidean or manhattan
34 # method = complete , single or average .
35
36 #Example: Agglo clustering on iris
37
38
39 Duplicated, tail, and, rev

```