

Tools for Analytics

- Elaborate tools with nifty visualizations; expensive licensing fees: Ex: Tableau, Tom Sawyer
- Software that you can buy for data analytics: Brilig, small, affordable but short-lived
- Open sources tools: Gephi, sporadic support
- Open source, freeware with excellent community involvement: R system
- Some desirable characteristics of the tools: simple, quick to apply, intuitive, useful, flat learning curve
- A demo to prove this point: data → actions /decisions

R Language

- R is a software package for statistical computing.
- R is an interpreted language
- It is open source with high level of contribution from the community
- "R is very good at plotting graphics, analyzing data, and fitting statistical models using data that fits in the computer's memory."
- "It's not as good at storing data in complicated structures, efficiently querying data, or working with data that doesn't fit in the computer's memory."

Why R?

- There are many packages available for statistical analysis such as SAS and SPSS but there are expensive (user license based) and are proprietary.
- R is open source and it can pretty much do what SAS can do but free.
- R is considered one of the best statistical tools in the world.
- For R people can submit their own packages/libraries, using the latest cutting edge techniques.
- To date R has got almost 15,000 packages in the CRAN (Comprehensive R Archive Network – The site which maintains the R project) repository.
- R is great for exploratory data analysis (EDA): for understanding the nature of your data before you launch serious analytics.
- Many tutorial vignettes are available for you to learn.

R Packages

- An R package is a set of related functions
- To use a package you need to load into R
- R offers a large number of packages for various vertical and horizontal domains:
- Horizontal: display graphics, statistical packages, machine learning
- Verticals: wide variety of industries: analyzing microarray data, modeling credit risks, social sciences, automobile data (none so far on sensor data from automobiles!)

Library

- Library \rightarrow Package \rightarrow Class \rightarrow
- R considers every item as a class/object
- Thousands of Online libraries
- 150000 packages
- CRAN: Comprehensive R Archive Network
- Look at all the packages available in CRAN <u>http://cran.r-project.org/</u>

• R-Forge is another source for people to collaborate on R projects

Approach to learning R

- R Basics, fundamentals
- The R language
- Working with data
- Statistics with R language



R Basics

- Obtaining the R package
- Installing it
- Install and use packages
- Quick overview and tutorial



A quick demo of R's capabilities

 See p.98 onwards till p.102 of simpleR: Using R for introductory statistics By J. Verzani
 <u>http://cran.r-project.org/doc/contrib/Verzani-SimpleR.pdf</u>

R in a nutshell, by Joseph Adler, O'reilly, 2010 Chapter 3 Basics, Ch.4 packages, (search for this online) Look for these resources online...and try these.
See Rhandout pdf linked to today's lecture

See Rhandout.pdf linked to today's lecture

Packages

- A package is a collection of functions and data files bundled together.
- In order to use the components of a package it needs to be installed in the local library of the R environment.
- Loading packages
- Custom packages
- Building packages



More R

- R syntax
- R Control structures
- R Objects
- R formulas



R Syntax

- R language is composed of series of expression resulting in a value
- Examples of expression include assignment statements, conditional statements, and arithmetic expressions
 - > a<- 42
 - > b <- a% 5
 - > if (b == 0) " a divisible evenly by 5" else " not evenly divisible by 5"
- [1] " not evenly divisible by 5"
- Variables in R are called symbols



CSE4/587



Special Values

- > v <- c(1,2,3)
- > v
- [1] 1 2 3
- > length(v) <- 4
- > v
- [1] 1 2 3 NA
- NA : not defined or not available
- Very Large and very small numbers:
- > 2 ^ 1024
- [1] Inf
- > 2 ^ 1024
- [1] -Inf



Curly braces

15

Curly braces are used to group a set of operations in the body of a function:
f <- function() {x <- 1; y <- 2; x + y}
f()
[1] 3



Control Structures

16

- > i <-4
- > repeat {if (i > 25) break else {print(i); i <- i + 5;}}</p>

2/7/2018

- [1] 4
- [1] 9
- [1] 14
- [1] 19
- [1] 24



Demo: Exam Grade: Traditional reporting 1

Q1	Q2	Q3	Q4	Q5	Total
16.7	13.9	9.6	18.5	13.7	72.4
20.0	16.0	9.0	19.0	17.0	76.0
20.0	20.0	15.0	25.0	20.0	90.0
Q1	Q2	Q3	Q4	Q5	Total
16.0	14.2	9.6	19.4	14.0	73.2
80.1%	71.1%	64.0%	77.4%	70.2%	73.2%
Q1	Q2	Q3	Q4	Q5	Total
17.3	13.6	9.7	17.6	13.3	71.5
86.7%	67.8%	64.6%	70.3%	66.7%	71.5%

Question 1..5, total, mean, median, mode; mean ver1, mean ver2





Distribution of exam1 points

CSE4/587



CSE4/587





R-code



```
exam1<-data2$midterm
hist(exam1, col=rainbow(8))
boxplot(data2, col=rainbow(6))
```

```
boxplot(data2,col=c("orange","green","blue","grey","yellow", "sienna"))
fn<-boxplot(data2,col=c("orange","green","blue","grey","yellow", "pink"))$stats</pre>
```

```
text(5.55, fn[1,6], paste("Minimum =", fn[1,6]), adj=0, cex=.7)
text(5.55, fn[2,6], paste("LQuartile =", fn[2,6]), adj=0, cex=.7)
text(5.0, fn[3,6], paste("Median =", fn[3,6]), adj=0, cex=.7)
text(5.55, fn[4,6], paste("UQuartile =", fn[4,6]), adj=0, cex=.7)
text(5.55, fn[5,6], paste("Maximum =", fn[5,6]), adj=0, cex=.7)
```

```
grid(nx=NA, ny=NULL)
```



One last example

- survey.results <- factor(c("Disagree", "Neutral", "Strongly Disagree", "Neutral", "Agree", "Strongly Agree", "Disagree", "Strongly Agree", "Neutral", "Strongly Disagree", "Neutral", "Agree"), levels=c("Strongly Disagree", "Disagree", "Neutral", "Agree", "Strongly Agree"), ordered=TRUE)
- survey.results
- R will automatically compute the numbers in each category!
- There are many more functions and operations available in R that are related to data.

Lets explore RStudio

~/EDASigmoid - RStudio

0

CSE4/587

_

6

File Edit Code View Plots Session Build Debug Tools Help

👰 🗸 🥶 🖌 🔒 💷 🖉 🖉 🚱 🖓

data1 × mtcars ×											Environment History			-		
										32 observations of 11 variables 🛛 💣 🕞 📰 Import Dataset 🗕 🎸 Clear 🎯		📃 List 🗸				
	row.names	mpg	cyl	disp	hp	drat	wt	qsec	VS	am	gear	carb	Global Environment -	Q		
1	Mazda RX4	21.0	6	160.0	110	3.90	2.620	16.46	0	1	4	4	Data			7
2	Mazda RX4 Wag	21.0	6	160.0	110	3.90	2.875	17.02	0	1	4	4	• data1 25 obs. of 3 variables			
3	Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4	1	Omtcars 32 obs. of 11 variables			
4	Hornet 4 Drive	21.4	6	258.0	110	3.08	3.215	19.44	1	0	3	1	Values			
5	Hornet Sportabout	18.7	8	360.0	175	3.15	3.440	17.02	0	0	3	2	≡ oglm.out List of 30			
6	Valiant	18.1	6	225.0	105	2.76	3.460	20.22	1	0	3	1	coefficients : Named num [1:2] -21.23 1.63			
7	Duster 360	14.3	8	360.0	245	3.21	3.570	15.84	0	0	3	4	attr(*, "names")= chr [1:2] "(Intercept)" "Age"	attr(*, "names")= chr [1:2] "(Intercept)" "Age" residuals : Named num [1:25] -1.002 -1.01 -1.019 -0.413 -0.482		
8	Merc 240D	24.4	4	146.7	62	3.69	3.190	20.00	1	0	4	2	residuals : Named num [1:25] -1.002 -1.01 -1.019 -0.413 -0.482			
9	Merc 230	22.8	4	140.8	95	3.92	3.150	22.90	1	0	4	2	attr(*, ''names) = chr [1:25] 11 '' 2'' '' '' 4''	041.22		
10	Merc 280	19.2	6	167.6	123	3.92	3.440	18.30	1	0	4	4	Titted.values : Named num [1:25] 0.00203 0.01051 0.0187 0.02780 0	04132	·	
11	Merc 280C	17.8	6	167.6	123	3.92	3.440	18.90	1	0	4	4				_
12	Merc 450SE	16.4	8	275.8	180	3.07	4.070	17.40	0	0	3	3	rites Piots Packages nelp viewer	(0)		
13	Merc 450SL	17.3	8	275.8	180	3.07	3.730	17.60	0	0	3	3	Install Packages Q Check for Updates G	Q		
14	Merc 4505LC	15.2	8	275.8	180	3.07	3.780	18.00	0	0	3	3	boot Bootstrap Functions (originally by Angelo Canty for S)	1.3-7	/	8
15	Lincoln Continental	10.4	8	4/2.0	205	2.93	5.250	17.98	0	0	3	4	Class Functions for Classification	7.3-5	j	8
10	Chrysler Imperial	14.7	•	400.0	215	2.00	5.424	17.02	0	0	2	4	🔲 📃 <u>cluster</u> Cluster Analysis Extended Rousseeuw et al.	1.14.5	.3	8
10	Cinyster Imperior	22.4	0	70 7	250	1 00	2 200	19.47	1	1	3	1	Code Analysis Tools for R	0.2-8	3	()
19	Honda Civic	30.4	4	75.7	52	4.93	1,615	18.52	1	1	4	2	🗐 📄 colorspace Color Space Manipulation	1.2-4	4	0
Const. (DAGeneral/ -										Compiler The R Compiler Package	2.15.	.2	8			
Console ~/EDAsigmoid/ @										datasets The R Datasets Parkane	215	2	0			
R is a collaborative project with many contributors.						5.						20.0	-	0		
'citation()' on how to cite R or R packages in publications.							in publ	licatio	ons.				2.0-0	, 	0	
							1100	help	~ ~				digest Create cryptographic nash digests of K objects	0.6.4		8
Type 'demo()' for some demos, 'help()' for on-line help, or 'help.start()' for an HTML browser interface to help. Type 'q()' to quit R. [Workspace loaded from ~/EDASigmoid/.RData]								lp.	OL.			doBy - Groupwise summary statistics, general linear contrasts, population means (squares-means), and other utilities	ast- 4.5-9)	8	
													foreiqn Read Data Stored by Minitab, S, SAS, SPSS, Stata, Systat, dBase,	0.8-5	j 1	8
													gqmap A package for spatial visualization with Google Maps and OpenStreetMap	2.3		8
> getOption("defaultPackages")													gqplot2 An implementation of the Grammar of Graphics	0.9.3	.1	0
[1] "datasets", "utils", "grDevices" "graphics", "stats", "methods"								"stat	ts"		metho	graphics The R Graphics Package	2.15.	.2	8	
> [1	(.packages()) l "stats" "or:	anhics		arDevi	ices"	"uti	ls"	"data	aset		metho	ds"	se"	2.15.	.2	0
> View(data1)						Gall				grid The Grid Graphics Package	215	2	8			
> View(datal)												atable Arrange grobs in tables	012	-	0	
5	View(mtcars)												Analyze yous in tables.	2.22		0
>													Functions for kernel smoothing for wand & Jones (1995)	2.23-	0	

▲ 🇞 🛊 📭 🚰 🌵 😻 9:14 AM 6/4/2014



2/7/2018

CSE4/587

Let do a EDA of cars data

- Look at the tutorial in handout#1
- R is good for Exploratory Data Analytics
- It is really good for most statistical computing you will you in your domain.
- You can repeat the same on Jupyter.
- We will also look at real "click" data from Nytimes from Oneill and Schutt's text. See the data in today's notes.