BIG DATA Solution of the second sec



TEXAS A&M UNIVERSITY
Superfund Research Center

Welcome !

superfund.tamu.edu

This Session will Begin at 2:00 pm Eastern US Time "Manipulating and Displaying Big(ish) Data in R "

Allison Dickey– North Carolina State University Fred Wright– North Carolina State University Burcu Beykal– University of Connecticut



BIG DATA Superfund.tamu.edu IN ENVIRONMENTAL SCIENCE AND TOXICOLOGY



TEXAS A&M UNIVERSITY
Superfund Research Center

- All participants are muted to enable the speakers to present without interruption.
- Please rename yourself and designate Full Name and Affiliation.
- Last-minute installation issues? Use private chat to any of the speakers when they're not speaking.
- Use the reaction icon at the bottom of your screen to raise your hand.
 However, due to time constraints, most questions will be answered in chat window.
 - This meeting will be recorded and posted on the @tamusuperfund website <u>https://superfund.tamu.edu/big-data-series-2021/</u> in the coming weeks.



BIG DATA SIENCE AND TOXICOLOGY



BIG DATA (Construction) In the second second







Part 3: More Advanced Visualization





What is this session about?

- This session is NOT a comprehensive introduction to R
- We present a number of topics that have been useful to the speakers for Superfund-relevant data.
- In other words, what do we and collaborators do most often?
- We do not assume familiarity with R, and you might need to review material later!





Do you want to just drive it, or understand it? (a bit of both)



BIG DATA

Æ

IN ENVIRONMENTAL SCIENCE AND TOXICOLOGY

VS.







Why are there so many... programs about R?





RStudio

WALT DISNEY PICTURES/PHOTOFEST





R



RMarkdown









Why are there so many... programs about R?

Actual language





RStudio

WALT DISNEY PICTURES/PHOTOFEST



R Open

R

Another version of R, can be faster on some computers for some code



-RMarkdown hanced reproducibility, save as htm



RShiny

Web tools, sliders, radio buttons, etc.



R: a pre-introduction (Part O) Fred Wright, NC State University





A schematic view This will make more sense after today!



A schematic view of R (from E. Paridis, R for Beginners)





- *R* stores variables, data, functions, results, etc, in the form of *objects* which have a *name*.
- The user can do actions on these objects with *operators* (arithmetic, logical, and comparison) and *functions* (which are themselves objects)
- Let's illustrate by starting with the command line in a new session. You'll see ">" starting each command, to distinguish it from output. Later we may drop this

BIG DATA &

Simple Arithmetic Operations











If the object already exists, its previous value is erased (the modification affects only the objects in the active memory, not the data on the disk).





Object names : applies to variables & functions

Object names cannot start with digits or periods, but they can have periods (.), underscores (_) and digits within them.

Names are case sensitive.

```
> d <- 5
Error: unexpected input in " "
>
> d <- 5
Error: unexpected input in " "
> d <- 5
> d
[1] 5
> 1f <- 5
Error: unexpected symbol in "1f"
> f1 <- 5
> f1
[1] 5
> var.name <- 10
> var.name
[1] 10
```

BIG DATA So Superfund.com/.edu N ENVIRONMENTAL SCIENCE AND TOXICOLOGY		
> n == 20 +	Logical "equals" comparison.	
[1] FALSE		
> n == 15		
[1] TRUE		
> n < 10		
[1] FALSE		
> n > 10		
[1] TRUE		
> n <= 20 +	"Less than or equal to"	
[1] TRUE		
> n == n		
[1] TRUE		
> a <- 15	Assignment statement	
> a < −15 ←	$\Delta query (le 'a' less than 152)$	I Start I
[1] FALSE		
> n == a		
[1] TRUE		AT AT A





Primitive data types : character, numeric, logical

CHARACTER :

> a <- "scooby doo"

> a

- [1] "scooby doo"
- > is.character(a) ←
- [1] TRUE
- > typeof(a)
- [1] "character"

functions that start with "is" often can answer a question about an object

BIG DATA SIPERIAL SCIENCE AND TOXICOLOGY

NUMERIC :

- > a <- 10
- > is.numeric(a)
- [1] TRUE
- > is.character(a)
- [1] FALSE
- > a <- 10/3
- > a
- [1] 3.333333
- > is.numeric(a)
- [1] TRUE
- > a^100
- [1] 1.940325e+52
- > a^1000
- [1] Inf •
- > a^(-1000)
- [1] 0 •

Overflow. Value too large.

Underflow. Value too small.







Functions

Functions need an argument, provided in parentheses
> a <- 10
> sqrt(a)
[1] 3.162278
> is.numeric(a)
[1] TRUE





R object types (not exhaustive)

• Vector:

- a one-dimensional "array." All elements of the vector must be of the same data type--numerical, character, etc.
- Matrix:
 - a two-dimensional array with an arbitrary number of rows and columns. Again all elements of the matrix must be of the same data type.

• Data frame:

 Also used to store tabular data. Each column can be a different data type. Columns can be thought of as different "variables," and column names can be used in R built-in functions.

BIG DATA SIENCE AND TOXICOLOGY

Here is a matrix I made



TEXAS A&M UNIVERSITY
Superfund Research Center

> m

	[,1]	[,2]	[,3]	[,4]	[,5]	[,6]	[,7]	[,8]	[,9]	[,10]	[,11]
[1,]	-5	-4	-3	-2	-1	0	1	2	3	4	5
[2,]	-10	-8	-6	-4	-2	0	2	4	6	8	10
[3,]	-15	-12	-9	-6	-3	0	3	6	9	12	15
[4,]	-20	-16	-12	-8	-4	0	4	8	12	16	20
[5,]	-25	-20	-15	-10	-5	0	5	10	15	20	25
[6,]	-30	-24	-18	-12	-6	0	6	12	18	24	30
[7,]	-35	-28	-21	-14	-7	0	7	14	21	28	35
[8,]	-40	-32	-24	-16	-8	0	8	16	24	32	40
[9,]	-45	-36	-27	-18	-9	0	9	18	27	36	45
[10,]	-50	-40	-30	-20	-10	0	10	20	30	40	50
[11,]	-55	-44	-33	-22	-11	0	11	22	33	44	55
$> \dim(m)$											

> dim(m) [1] 11 <u>11</u>

Single Element

> m[,2] -

> $m[3,2] \leftarrow [row, column]$ [1] -12 Column Selection [all rows, column 2]

Note : element selection used square brackets [].

[1] -4 -8 -12 -16 -20 -24 -28 -32 -36 -40 -44





> smaller.m <- m[2:3,4:8]+ > smaller.m [,1] [,2] [,3] [,4] [,5] [1,] -4 -2 02 4 [2,] -6 -3 0 3 6 > smaller.m^2 -[,1] [,2] [,3] [,4] [,5] [1,] 16 4 0 16 4 [2,] 36 9 0 9 36 Subset of matrix is a matrix

Many operations operate element-wise





A simple data frame

- 3 9 1.9
- 4 6 4.5
- > mydf[,2]
- [1] 1.2 3.4 1.9 4.5
- > mydf\$y
- [1] 1.2 3.4 1.9 4.5



Next up: Part 1 - The actual Intro! (shown as .html and R code)

BIG DATA 🍌

IN ENVIRONMENTAL SCIENCE AND TOXICOLOGY



22



Part 2: Data visualization and basic analysis

Fred Wright, NC State University





Review A: Got an R question?

- In RStudio, use the Help menu
- In the console window (whether Rstudio or basic R), help()
- Or ? before a command
- If you're still not sure, try googling the question, which often leads to a keyword
- More advanced code and solutions appear on sites such as stackoverflow (again, google-able)





Review B: Be careful not to lose data!

- R can "hang" if you try to plot something huge, or start a very data intensive procedure, etc.
- Sometimes you can "escape." Literally!
- R objects are stored in memory
- You might need to kill an R session.
- The best defense against data loss? Scripts that can reproduce everything.



Plotting in "base" R

- Simple code for simple plots
- More advanced plots require some understanding of R syntax, but can be fancy
- Want fancy without lots of coding yourself? Stay tuned for next section (but still pay attention now!)



Reading in the data (review from previous)

this code picks up where Part 1 left off

count_table <- read.csv("GSE62902.csv", row.names = 1) # read data
alignment_subset <- count_table[39181:39185,] # figure out which are the pesky last rows
count_table <- count_table[-(39181:39185),] # keep only the non-pesky rows
dosage <- count_table[1,] # save dosage as its own object
count_table <- count_table[-1,] # now remove dosage, because it's not a sequencing count</pre>



Normalizing data (not publication-quality)

install.packages("preprocessCore") # need only do once for R version
library(preprocessCore) # need only do once for R session

dosage[dosage==0]<-0.003 # I'll eventually take logs, so can't have zero log10dose<-log10(as.numeric(dosage)) # I like the log scale for dose X<-as.matrix(count_table)+0.5 # to make sure it's a matrix, +.5 to zero-protect Xnorm1<-log10(t(t(X)/colSums(X))) # tricky, divide by column sums and take logs Xnorm2<-normalize.quantiles(Xnorm1) # quantile normalization, "final" expression data #Xnorm2<-data.frame(Xnorm2) # data frame has some nice properties colnames(Xnorm2)<-colnames(X) # but now we should add the column names back rownames(Xnorm2)<-rownames(X) # and we should add the row names back</pre>



Histograms

hist(Xnorm2[,1]) # histogram of the first column hist (Xnorm2\$GSM1535917) # won't work becuase Xnorm2 is not a data frame hist(Xnorm2, col="yellow") # color! par(mfrow=c(2,2)) # sets "parameters" for plotting, in this case 2*2 rows/columns plotting hist(Xnorm2[1,],col=1) # first row=first gene hist(Xnorm2[2,],col=2) # why did it do that? hist(Xnorm2[3,],col=3) hist(Xnorm2[4,],col=4)



Boxplots

boxplot(Xnorm1) # a boxplot
boxplot(Xnorm2) # reflects alignment
dev.off() # turn off graphics device





Summary stats

summary(Xnorm2[,1]) # summarizes basic stats
summary(Xnorm2[1,])
print(mean(Xnorm2)) # the mean. try sd,, median,
etc.





Scatterplots

plot(Xnorm2[,2],Xnorm2[,30]) # mouse 2 vs. mouse 30
abline(0,1,col=2) # col=2 means "red"
identify(Xnorm2[,2],Xnorm2[,30]) # click near points
and then to console, hit escape





Scatterplots cont.

plot(log10dose,Xnorm2[4305,]) # gene 4305 vs. log10dose
plot(log10dose,Xnorm2[4305,],pch=16,col="blue")
plot(log10dose,Xnorm2[4305,],pch=16,col="blue",cex=1.5)
plot(log10dose,Xnorm2[4305,],pch=16,col="blue",cex=1.5,
ylab="normalized expression",main="My plot!")





Regression

- mylm<-lm(Xnorm2[4305,]~log10dose) # fit linear model
 summary(mylm) # shows what's in mylm
 abline(mylm) # add overlay line</pre>
- log10dose.2<-log10dose^2 # squared term</pre>
- mylm<-lm(Xnorm2[4305,]~log10dose+log10dose.2) # fit again, now with quadratic term
- points(log10dose,mylm\$fitted.values,col=2,pch=3) # overlay
 fitted values
- points(log10dose,mylm\$fitted.values,col=2,pch=3,lwd=3) #
 overlay fitted values, thicker





t-tests

t.test(Xnorm2[4305,dosage==0.003],Xnorm2[4305,dosage==30]) # ttest of two groups at extremes

?t.test



Analysis of variance

- myaov<-aov(Xnorm2[4305,]~dosage) # analysis of variance.
 what's wrong?</pre>
- myaov<-aov(Xnorm2[4305,]~as.factor(dosage)) # okay now
 it works</pre>
- summary(myaov) # results
Principal components

- mypca<-prcomp(t(Xnorm2)) #PCA, we use t() to turn data on side
- summary(mypca)
- str(mypca)
- #https://www.datacamp.com/community/tutorials/pca-analysisr
- pairs(mypca\$x[,1:3]) # pairs command does multiple pairwise
 scatterplots
- pairs(mypca\$x[,1:3],col=as.factor(dosage),pch=16) # color
 by dose
- plot(log10dose,mypca\$x[,1],pch=1) # wow! The first PC is
 largely driven by dose



Who says base R can't make fancy plots? (this is utterly gratuitous)

#gratuitous plot

```
a <- (-440:200) / 200
b <- c(-300:300) / 200
b <- complex(length(b), 0, b)
m <- outer(a, b, FUN = "+")
magnitude.m <- sqrt(Re(m)^2 + Im(m)^2)
z <- matrix(0,length(a), length(b))
for (i in (1:20)) {
    z <- z^2 + m</pre>
```

```
}
```

```
magnitude.z <- sqrt(Re(z)^2 + Im(z)^2)
magnitude.z[is.na(magnitude.z)] <-magnitude.m[is.na(magnitude.z)]^500
image(a, Im(b), -log(magnitude.z) * (magnitude.z > 10), col =
c(heat.colors(100),1))
```





Next up: More advanced visualization!



http://publicdomainaudiovideo.blogspot.com/

MORE ADVANCED VISUALIZATION

Burcu Beykal, PhD



.



.....



Advanced Visualization with the ggplot2 Library

- ggplot2 is based on the concept called grammar of graphics
- Main idea: any graph can be made from the same 3 components:
 - 1. Dataset
 - 2. Coordinate system
 - 3. Geoms visual marks to represent data points





Advanced Visualization with the ggplot2 Library

- ggplot2 is based on the concept called grammar of graphics
- Main idea: any graph can be made from the same 3 components:
 - 1. Dataset
 - 2. Coordinate system
 - 3. Geoms visual marks to represent data points



Final Plot







Getting started with ggplot2



•Generic syntax (help you get things plotted):

ggplot(data=<DATA>) + <GEOM_FUNCTION>(mapping = aes(<MAPPINGS>))

•Optional layers (help you beautify plots and improve visualization):

+ optional layers (<COORDINATE FUNCTION>, <FACET_FUNCTION>, <SCALE_FUNCTION>,

<THEME_FUNCTION>)



•Once we have a clean dataset:











•Adding elements like the dose information:













Improving visualization – changing shape





Boxplot Example

first gene is selected for the plot

•Getting started with boxplot: Analyzing 1 Gene across doses







Improve boxplot visualization: flip axis







Improve boxplot visualization: Change x-axis label & remove legend





Boxplot Example



There are still many ways improving this plot. Changing the background color, font styles, and sizes are some examples.





Improve boxplot visualization: Customizing labels and other texts in ggplot





Boxplot Example





- •Using the "theme" command is a great way to customize plots.
- •There are many complete themes to choose from and they are fully customizable

ggplot(df,aes(x=CLASS.DOSE, y=ENSMUSG0000000001, fill=CLASS.DOSE)) +
 geom_boxplot() + coord_flip() + labs(x = "Dose") +
 theme bw() + theme(legend.position="none")

theme black and white is one of the complete themes that can be used. Note how this is placed before the "theme()" command. If this is placed after, your customizations will be overwritten by the theme_bw()









Try new themes from the library

theme_solarized()





TEXAS A&M UNIVERSITY

RESEARCH CENTER





Adding Individual Observations on the Boxplot

Improve boxplot visualization: Overlaying samples on to the plot



geom_jitter will add the samples used to construct the boxplot. You can specify preferred colors to improve visualization.



Adding Individual Observations on the Boxplot





Changing Colors in ggplot

Improve boxplot visualization: Adding color elements





Changing Colors in ggplot





Choosing Colors in R

• Colors in R are defined using hexadecimals.

Colors in R document

- list of color names and options to chose from
- <u>http://www.stat.columbia.edu/~tzheng/files/Rcolor.pdf</u>

R Color Cheatsheet document

- <u>https://www.nceas.ucsb.edu/sites/default/files/2020-04/colorPaletteCheatsheet.pdf</u>
- Explains how colors are defined in R
- Calling colors using RColorBrewer library

• For color advice:

- <u>https://colorbrewer2.org/#type=sequential&scheme=Greens&n=3</u>
- Great for maps and other plot types



Final Boxplot







•Create density plot of measurements across doses:





Density Plots







•Let's change transparency to visualize better:













•Faceting plots with respect to dosage:





Faceting Plots





Faceting Plots








•Creating maps in R using ggplot2:





•Remove the fill color and change border color:







•Getting states on the map:







•Getting states on the map:











•Isolating a state for plotting:







•Changing the colors of counties:







•Changing the colors of counties:





+



•Add text to map:

find the centroid of the regions to place the text

```
centroid <- aggregate(cbind(long,lat) ~ subregion, data=ca_county, FUN=mean)
head(centroid)</pre>
```

ggplot(ca_df, aes(x = long, y = lat, group = group)) +
 geom_polygon(data = ca_county, aes(fill = subregion, alpha=0.6),
 geom_text(data = centroid, aes(x=long, y=lat, label=subregion),

color = "black") + inherit.aes = FALSE, fontface = "bold

```
coord_fixed(1.3) + theme_bw() + theme(legend.position="none")
```

place the text at the centroids and add the subregion labels



The figure looks busy. Let's highlight only a few counties to improve the visualization





•Select a subset of counties for labeling:

find the centroid of the regions to place the text

select_counties <- c("los angeles", "san diego", "monterey", "santa barbara")
select_ca <- centroid[centroid\$subregion %in% select_counties,]</pre>

```
ggplot(ca_df, aes(x = long, y = lat, group = group)) + geom_polygon(data = ca_county, aes(fill = subregion,
alpha=0.6), color = "black") +
        geom_text(data = select_ca, aes(x=long, y=lat, label=subregion), inherit.aes = FALSE, fontface = "bold") +
```

```
coord_fixed(1.3) + theme_bw() + theme(legend.position="none")
```

only add the selected text to the plot



Let's push the labels outward for better 40.0 visualization r ^{37.5} <mark>ع</mark>ل nonteres 35.0 santa barbara los angeles san diego 32.5 -120.0 -122.5 -117.5 -115.0 long





•Select a subset of counties for labeling:









•Load the Canada region data:

library(autoimage) data(canada) library(broom) canada_df <- tidy(canada)

Load the environmental data:

env_data <- read.csv("~/Desktop/workshop/complex_exposure_raw_data.csv")</pre>

•Plot:

```
ggplot(canada_df, aes(x=long, y = lat, group = group)) + geom_polygon(fill="gray70", color="black") +
    geom_point(data = env_data, aes(x = Long, y = Lat), inherit.aes = FALSE, size = 2, shape = 23,
fill = "darkred") +
    coord_fixed(1.3) + theme_bw()
```

add sampling locations from the environmental data as points to the map. The shape of the point is selected to be diamonds.







•Redefining axis limits:

```
ggplot(canada_df, aes(x=long, y = lat, group = group)) + geom_polygon(fill="gray70",color="black") +
geom_point(data = env_data, aes(x = Long, y = Lat), size = 2, shape = 23, fill = "darkred") +
theme_bw() +
coord_map(xlim=c(-122, -100), ylim=c(49, 62))
```

define x and y limits to zoom in





Let's remove the map fill color and plot lead levels for each sample.









The legend shows the lead levels in the samples



•Change the gradient colors for improved communication:

adjust the gradient color with low levels shown in blue (cold) and high lead levels with red (hot)





You can explore different metals and chemicals in the dataset using the same set of codes



•Show species on the map:

env_data\$Species <- factor(env_data\$Species)</pre>

make species factor

color samples with respect to the studies species



Canadian Oil Sands Region Chemical exposures among biota





Last Step: Saving Plots





Last Step: Saving Plots

•Using ggsave command to save plots:





Last Step: Saving Plots

•Where does R save the plots? Check and change the working directory



If you don't change the working directory, R will save your plots to the main directory.



Interactive Plots - Plotly



Interactive density plot:

save the ggplot command to a variable

panel.spacing = unit(0.1, "lines")) +
scale fill brewer(palette="Reds")

ggplotly(myplot)



call the variable in ggplotly to view the interactive plot



Interactive Plots - Plotly

Interactive scatter plot:

```
labs(col='Dose') +
```

```
theme_bw() +
theme(text = element_text(size=14))
```

ggplotly(my_scatter)

call the variable in ggplotly to view the interactive plot



Interactive Plots - Plotly

Interactive boxplot plot:

ggplotly(myboxplot)

call the variable in ggplotly to view the interactive plot

Saving an Interactive Plot as an HTML File



TEXAS A&M UNIVERSITY

RESEARCH CENTER



Cheatsheet for ggplot2

Link to the cheatsheet

https://www.maths.usyd.edu.au/u/UG/SM/STAT3022/r/current/Misc/data-visualization-2.1.pdf





Stackoverflow

ggplot heatmap labels		× 🕴 Q	
🔍 All 🖾 Images 🔗 Shopping 🗉 News 🕞 Video	os : More	Tools	5
About 45,400 results (0.55 seconds)			
https://stackoverflow.com > questions > how-to-get-labe			
How to get labels in my ggplot heatmap? - S Apr 22, 2015 · 1 answer Layers in the ggplot are plotted in order as they are written. To c	tack Overflow	stack overflow	About Products For Teams Q Search
the geom_text() call should be place after heatmap with values (ggplot2) - Stack Overflow	Sep 9, 2019	Home	How to get labels in my ggplot heatmap? Asked 6 years, 5 months ago Active 4 years, 10 months ago Viewed 3k times
Modify axis and label on a heatmap - Stack OverflowAprGGplot heatmap has 2 labels on each tile - Stack OverflowNovHow to add new y axis label on heatmap at specific locationFebMore results from stackoverflow.comNov	Apr 5, 2016 Nov 27, 2016 Feb 24, 2016	Collectives	I want to make a heatmap with the actual values as labels in the map, but the labels don't show. Here's my data: df <- structure(list(a = c(0.39, 0.26, 0.39, 0.45, 0.41, 0.42, 0.42, 0.34, 0.36,
		FIND A JOB Jobs Companies	<pre> Here's my code: library(ggplot2) library(grid) ggplot(df_m_aes(Scale_Item_fill=abs((orrelation))) + </pre>
		TEAMS Stack Overflow for Teams – Collaborate and share knowledge with a private group.	<pre>geom_title() + theme_bw(base_size=10) + theme(axis.text.x = element_text(angle = 90),</pre>



Stackoverflow

lome	How to get labels in my ggplot heatm	ap?	Scroll down to see the answer
UBLIC	Asked 6 years, 5 months ago Active 4 years, 10 months ago Viewe	d 3k times	
Questions	I want to make a heatmap with the actual values as labe	els in the map, but the labels don't show.	
Tags			
Users			
OLLECTIVES 0	df <- structure(list(a = c(0.39, 0.26, 0.39, 0	.45, 0.41, 0.42, 0.42, 0.34, 0.36, 0	
Explore Collectives			
IND A JOB	Here's my code:	1 Answer	Active Oldest V
Jobs	library(ggplot2) library(grid)		
Companies	<pre>gaplot(df.m. aes(Scale. Item. fill=abs(Correlat</pre>	Layers in the ggplot are plotted in	n order as they are written. To display text above the heatm
	<pre>geom_text(aes(label = round(Correlation, 2)), geom_tile() +</pre>	size=2.5) + the geom_text() call should be pla	ace after geom_tile() call.
EAMS			
EAMS Stack Overflow for Teams – Collaborate	<pre>theme_bw(base_size=10) + theme(axis.text.x = element_text(angle = 90),</pre>	Ŭ	
EAMS Stack Overflow for Teams – Collaborate and share knowledge with a private group.	<pre>theme_bw(base_size=10) + theme(axis.text.x = element_text(angle = 90),</pre>	geom_tile() +	
EAMS Stack Overflow for Teams – Collaborate and share knowledge with a private group.	<pre>theme_bw(base_size=10) + theme(axis.text.x = element_text(angle = 90),</pre>) + geom_tile() + geom_text(aes(label = round	(Correlation, 2)), size=2.5) +
EAMS Stack Overflow for Teams – Collaborate and share knowledge with a private group.	<pre>theme_bw(base_size=10) + theme(axis.text.x = element_text(angle = 90),</pre>) + geom_tile() + geom_text(aes(label = round	(Correlation, 2)), size=2.5) +
EAMS Stack Overflow for Teams – Collaborate and share knowledge with a private group.	<pre>theme_bw(base_size=10) + theme(axis.text.x = element_text(angle = 90),</pre>) + geom_tile() + geom_text(aes(label = round	(Correlation, 2)), size=2.5) +





Questions??




Thank you for joining us !

The next session is on November 3, 2021 2:00 – 4:00 pm Eastern US Time "PLACING TOXICOLOGY DATA IN THE CONTEXT OF EXPOSURE"

Caroline Ring—US Environmental Protection Agency