

Introduction to



Jamie Ford, NISD
Paulina Cano, CI:NOW

What is R?

- ◎ One of the most widely used data analysis software, used by statisticians, analysts, data scientists, etc.
- ◎ Powerful statistical programming language with unique data visualizations
- ◎ More than 14,000 libraries approved on CRAN (plus others on GitHub, etc.)
- ◎ R has more than 2 million users worldwide and is growing rapidly
- ◎ R can be downloaded online for free along with **Rstudio**

How does R compare to other statistical software?

SPSS

STATA

sas

R

	SPSS	STATA	sas	R
Ease of learning	✓✓	✓	✓	✗
Good user interface	✓✓	✓	✓	✗
Programming Capabilities	✗	✓	✓	✓✓
Support from company	✓	✓	✓	✗
Price	✗	✗	✗	✓
Advanced Visualization capabilities	✗	✗	✗	✓✓
Handle complex models	✗	✓✓	✓✓	✓✓
Handle large sets of data	✓	✓✓	✓✓	✓✓



Rstudio

1) Script

The screenshot displays the RStudio interface with four main panels:

- Script Editor:** Contains R code for loading data, plotting, and fitting lines. The code includes comments and function calls like `plot(google$degree, google$data_viz)` and `abline(lm(google$data_viz ~ google$degree, col="red"))`.
- Console:** Shows the output of the R commands, including variable types and values, such as `$facebook : num 1.93 -0.52 -1.18 2.21 -1.28 -1.33 -0.14 -0.34 -2.2 0.1 ...`.
- Workspace:** Lists the objects in the current environment, including `google` (51 obs. of 9 variables), `regmodel.1`, and `reg1`.
- Plots:** Displays a scatter plot titled "Interest in Data Visualization Searches by Percent of Population with College Degrees". The plot shows a positive correlation between the percentage of the population with college degrees (x-axis, 15-45) and searches for data visualization (y-axis, -1 to 3). Two regression lines are overlaid: a red line for the linear model and a blue line for the lowess smoother.

3) Workspace

4) Results/Plots

2) Console

The logo features a large, white, stylized letter 'R' with a thick, rounded top bar. To the right of the 'R', the word 'LIVE' is written in a clean, white, sans-serif font.

R LIVE

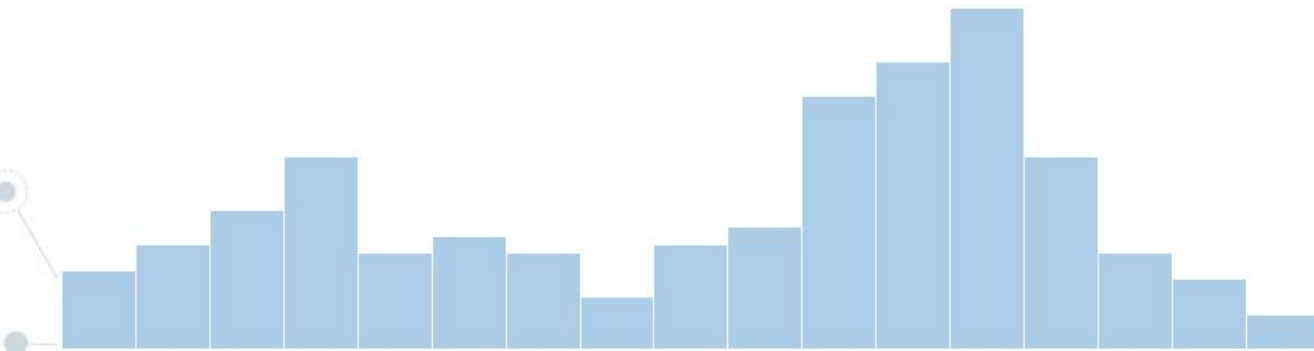
Is R right for you?

Advantages

- ⊙ Open-Source
- ⊙ Community support
- ⊙ Automation
- ⊙ Flexibility
- ⊙ Dynamic output

Disadvantages

- ⊙ Steep learning curve
- ⊙ Programming and capacity limitations when compared to Python or similar
- ⊙ Some libraries may not be updated
- ⊙ Not standardized



Using R to Work with Census Data

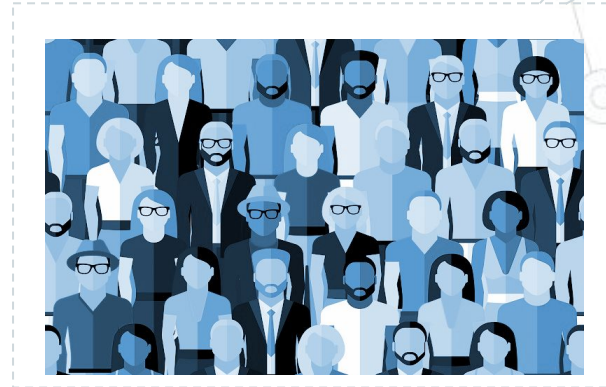
R allows you to download census data directly.

Steps:

1. Request a free Census Bureau API key
https://api.census.gov/data/key_signup.html
2. Download a few packages: **tigris** (shapefiles), **tidycensus** (Census and ACS data with feature geometries) and **sf**, (simple features is use to represent geographic vector data).
3. Load variables of interest.

```
nc_pop <-  
  get_acs(geography = "county",  
          variables = "B01003_001",  
          state = "NC",  
          geometry = TRUE)
```

4. You are now ready to interact with the data



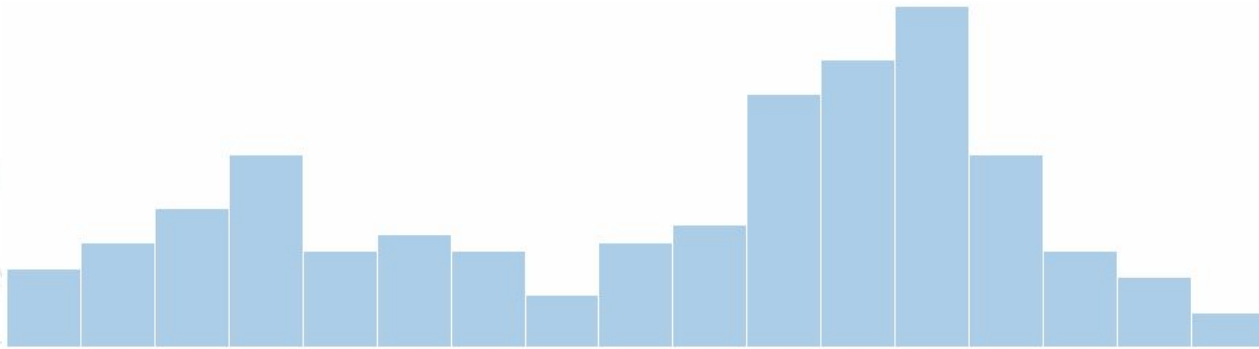
United States™
Census
Bureau

Continuation of Census Data and R

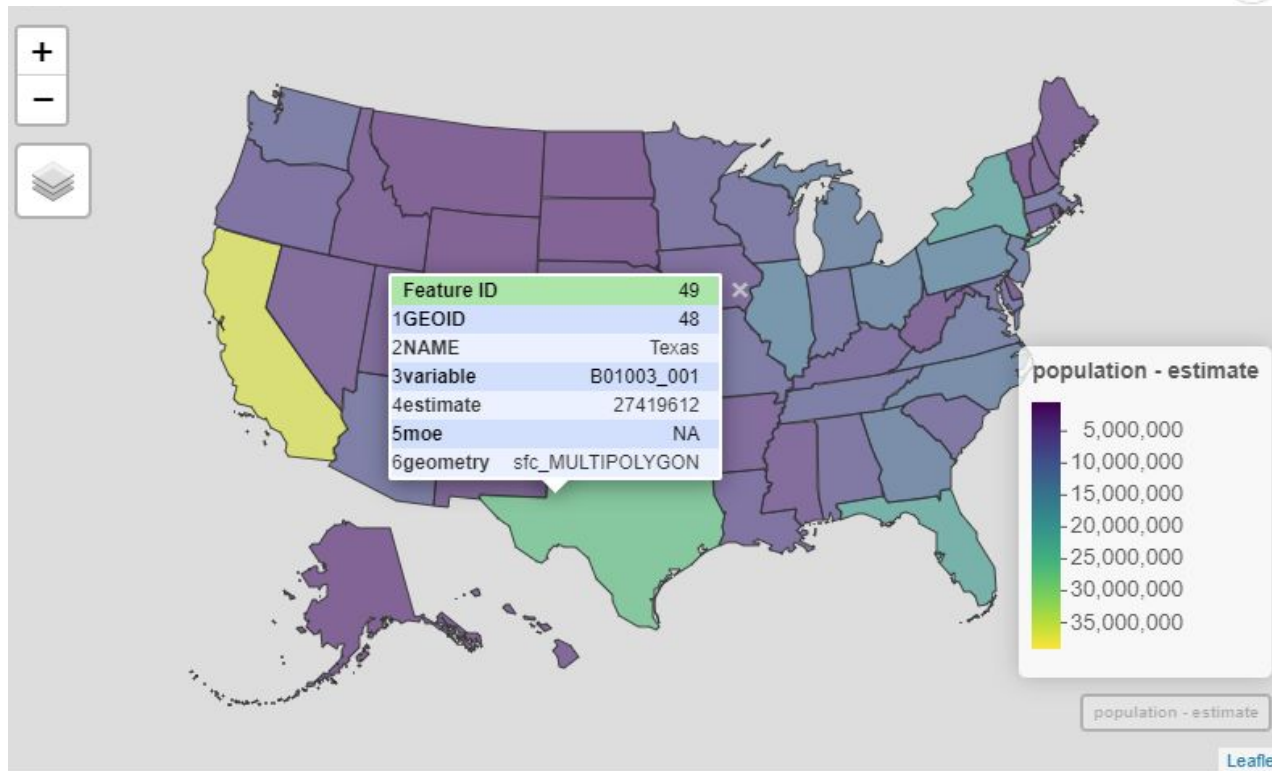
If we install the **leaflet** and **mapview** packages, we can visualize the data:

```
population <- get_acs(geography = "state",  
  variables = "B01003_001",  
  geometry = TRUE,  
  shift_geo = TRUE)
```

```
mapviewOptions(legend.pos = "bottomright")  
mapviewOptions(leafletWidth = 800)  
#mapviewOptions()  
#mapviewOptions(default = TRUE)  
mapview(population, zcol = "estimate", native.crs = TRUE, crs = 5070)
```



Example Output



https://map-rfun.library.duke.edu/02_choropleth.html

Data Manipulation



Derive new variables



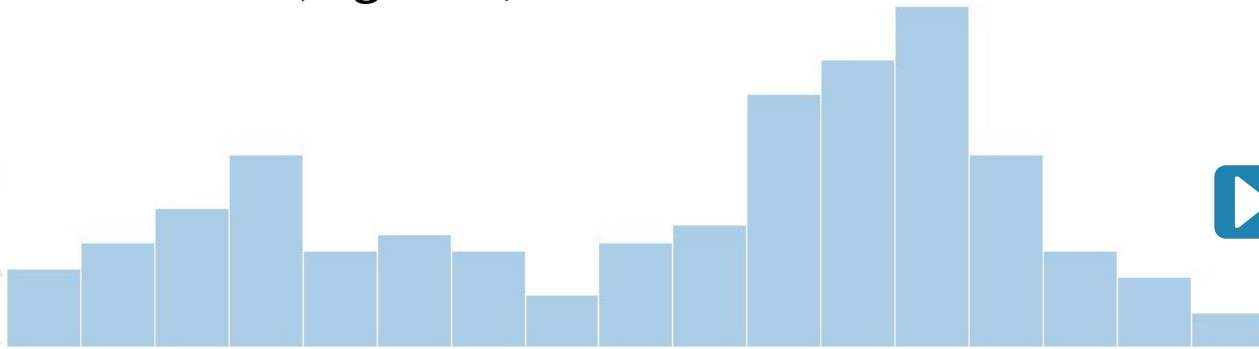
Join multiple data sets of data together



Create summaries of your dataset



Pull information directly from websites and/or public data sets (e.g. ACS)





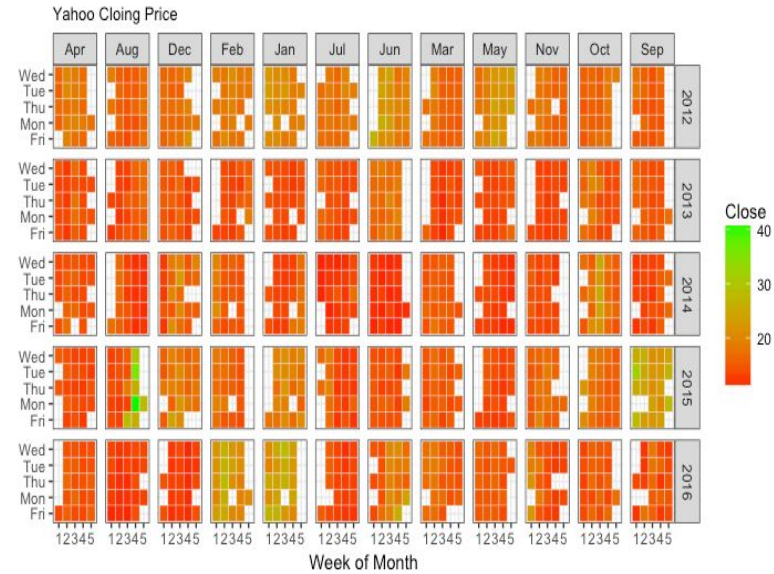
R LIVE

Data Visualization

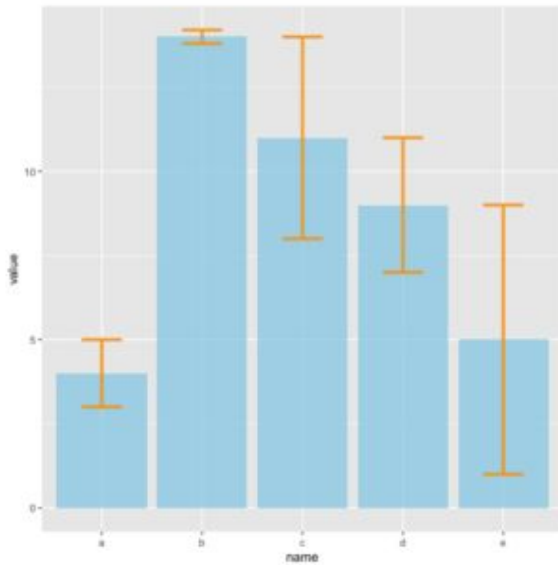
R has several packages that enable visualizing data:

- ⦿ BaseR
- ⦿ Ggplot2
- ⦿ Leaflet (interactive)
- ⦿ Plotly (interactive)
- ⦿ Other specialized (various models, EDA, GIS, network, etc.)

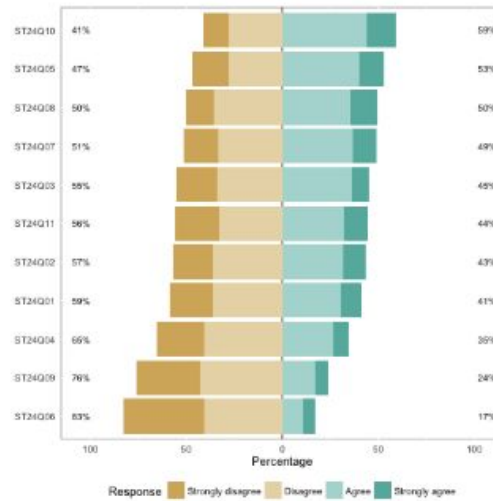
Time-Series Calendar Heatmap



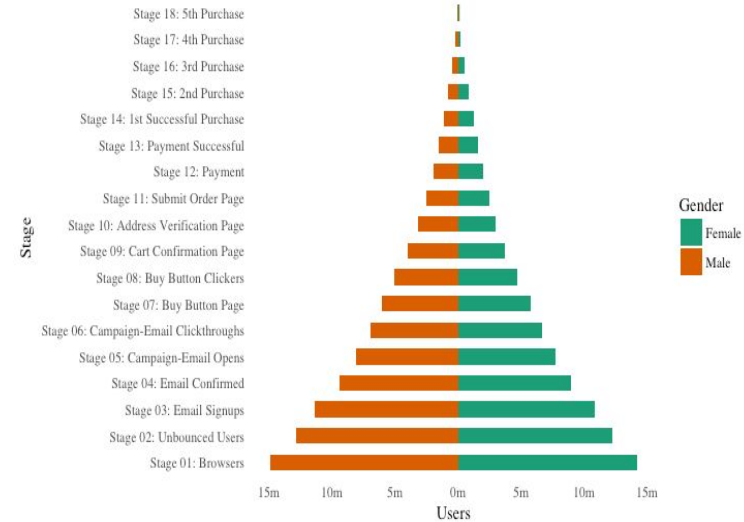
Bar Plots



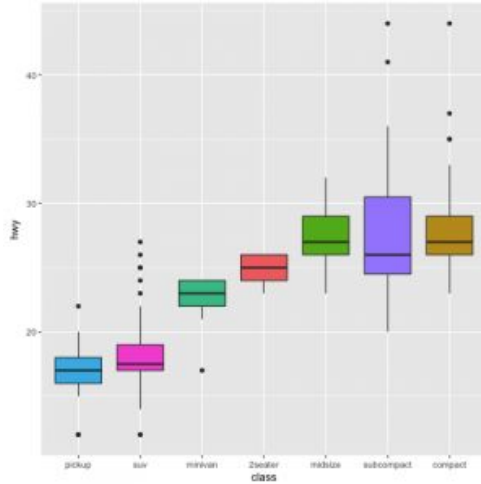
#4 Error bars on barplot



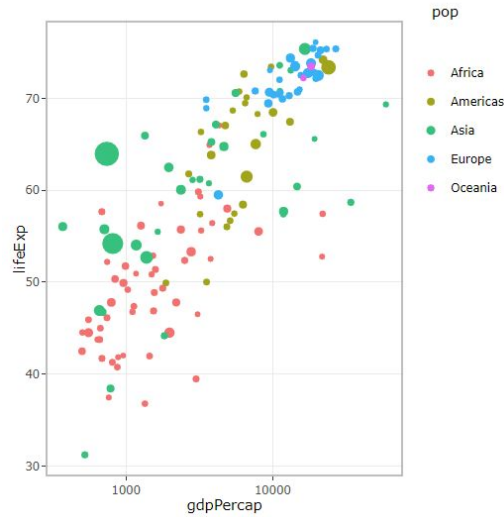
Email Campaign Funnel



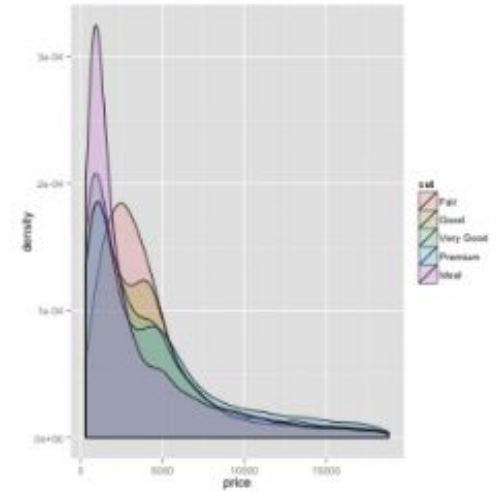
Data Visualization



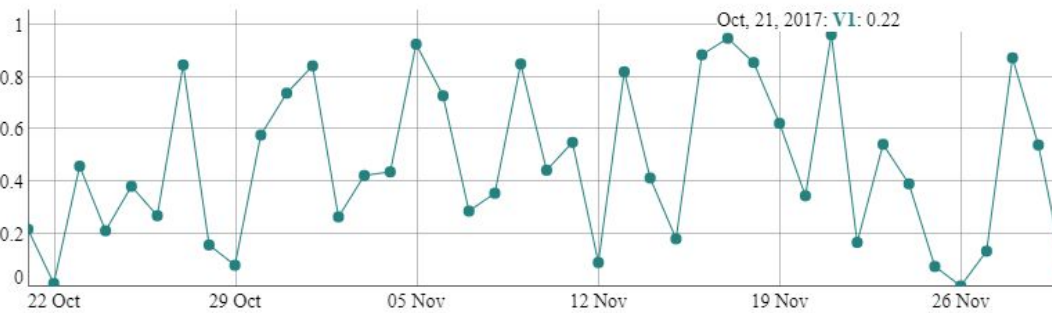
Boxplots



Bubble Graphs



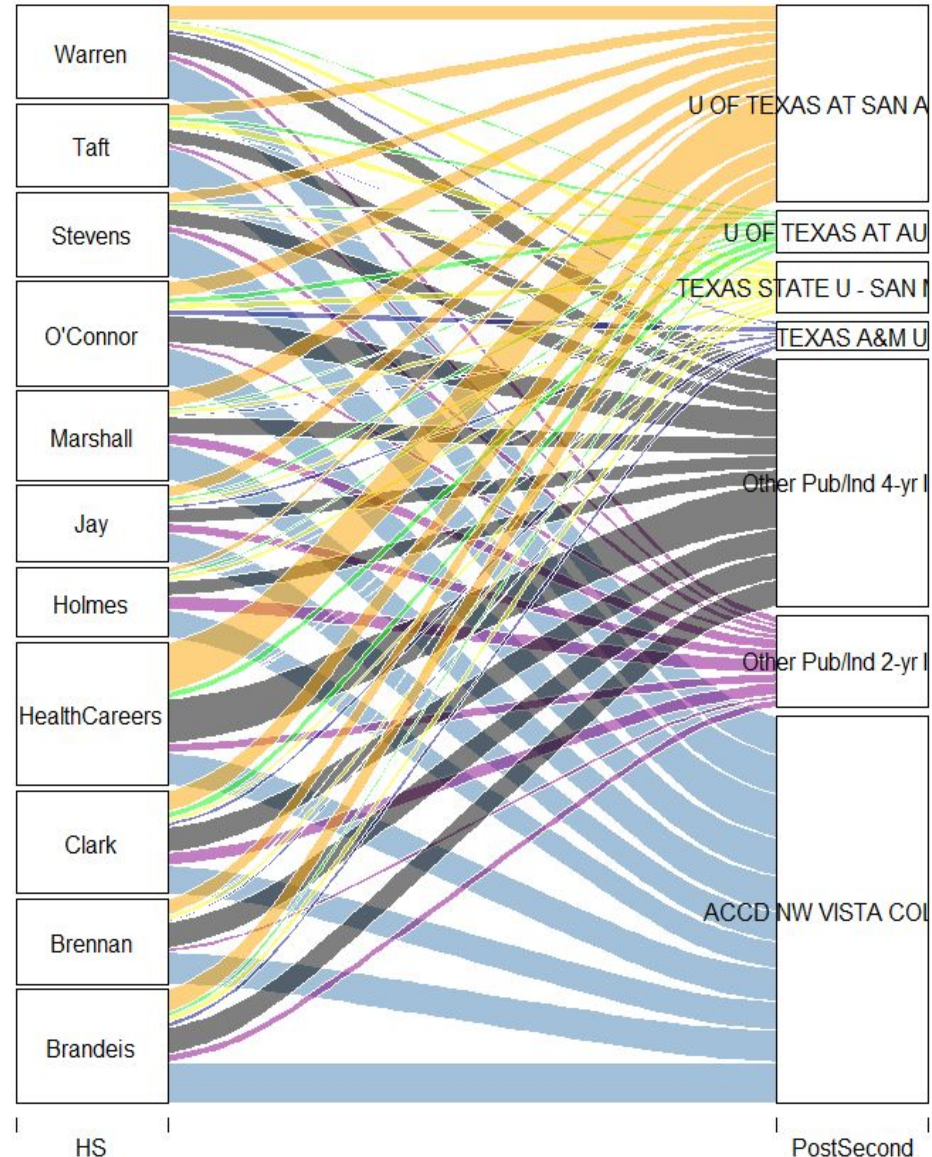
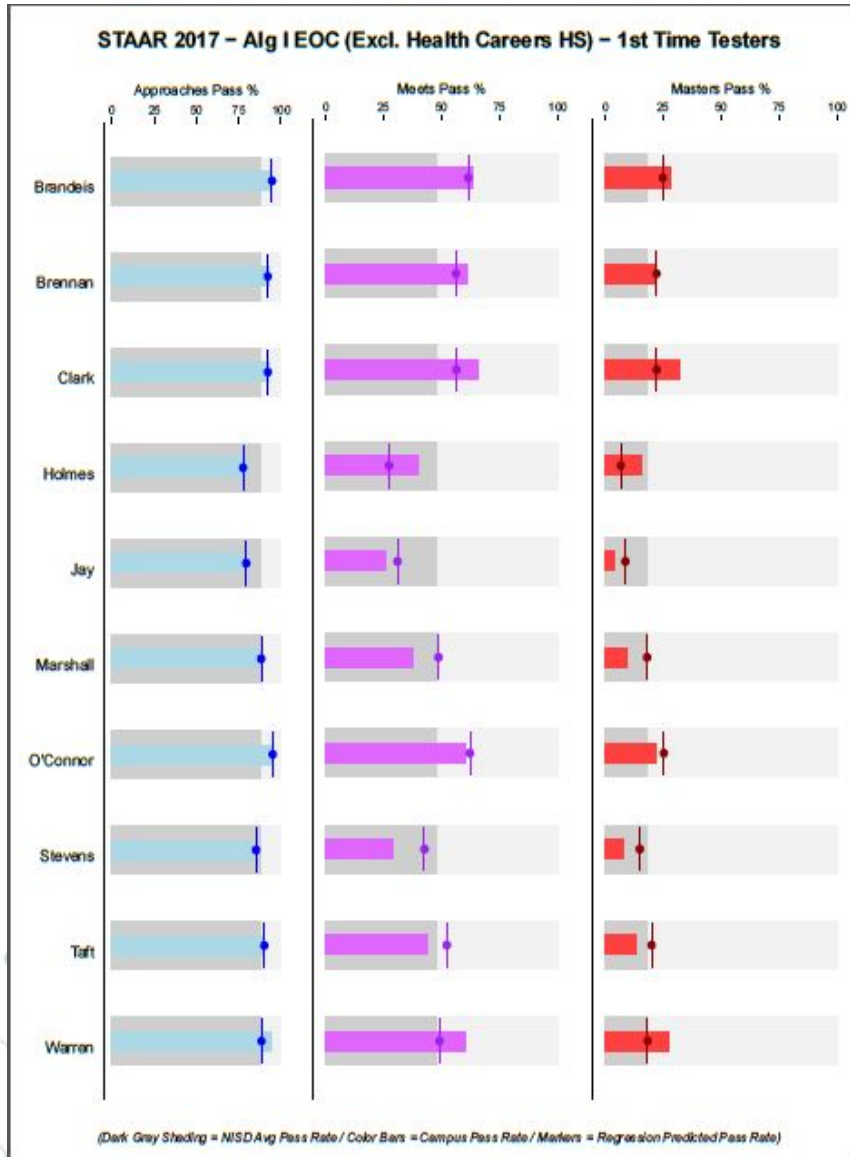
Density Graphs



Time Series Graph

EXAMPLES OF PROJECTS

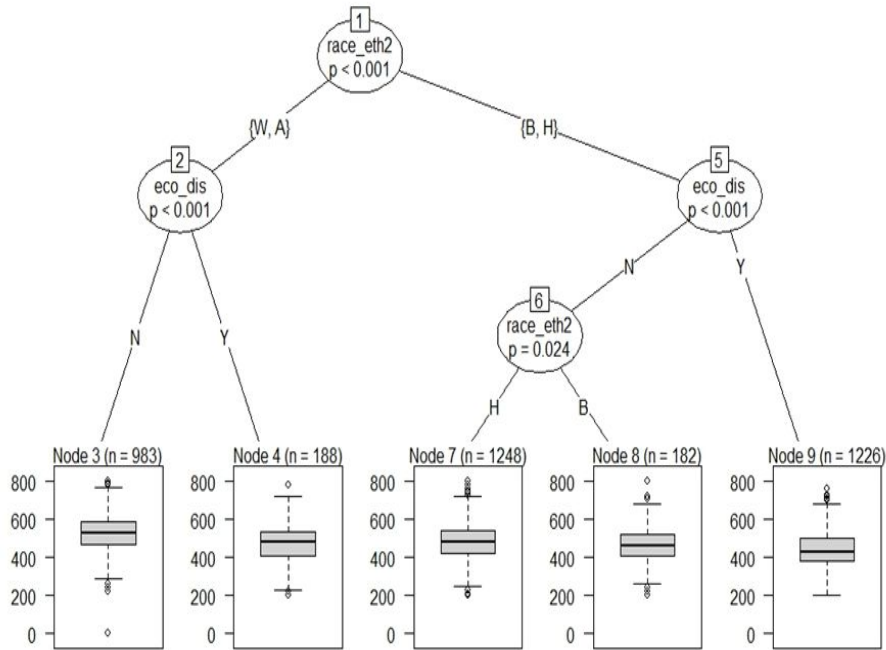
Visualizations of STAAR Results & College Enrollment Flows



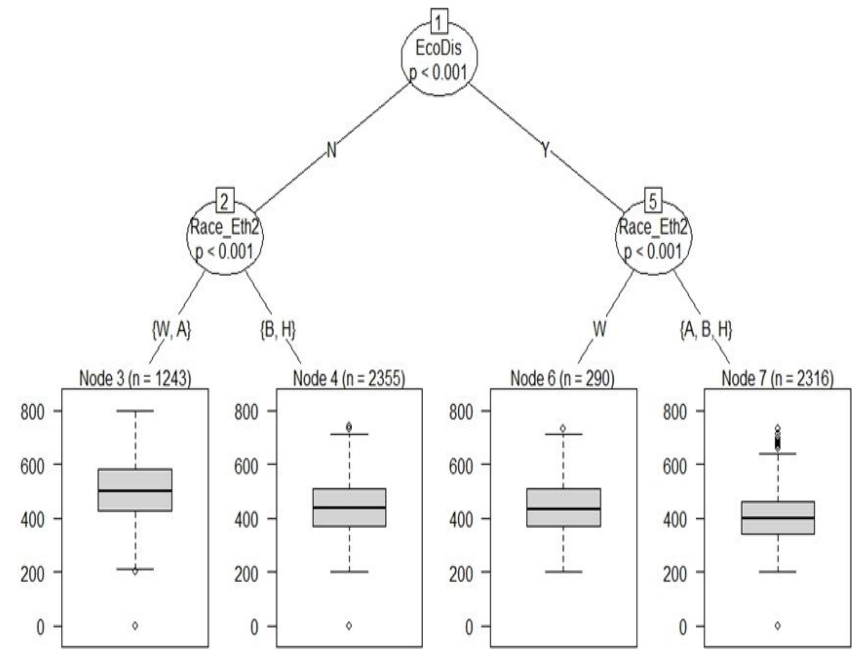
EXAMPLES OF PROJECTS

Decision Trees using CTREE

SAT SY1314 READING



SAT SPR 2015 READING



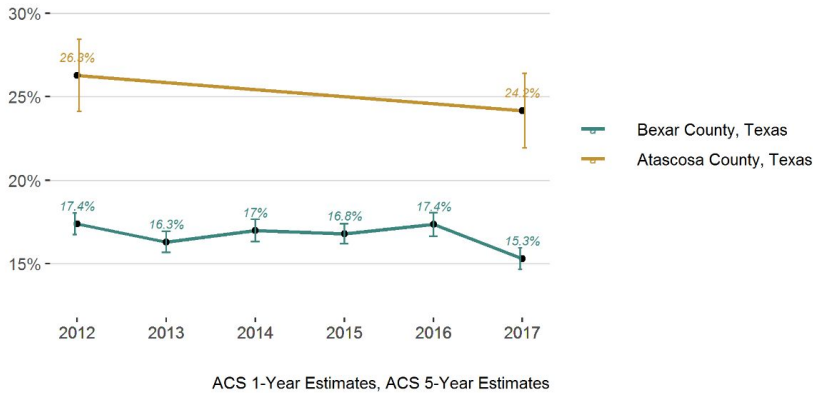
EXAMPLES OF PROJECTS

Automation and customization of over 200 Trendlines

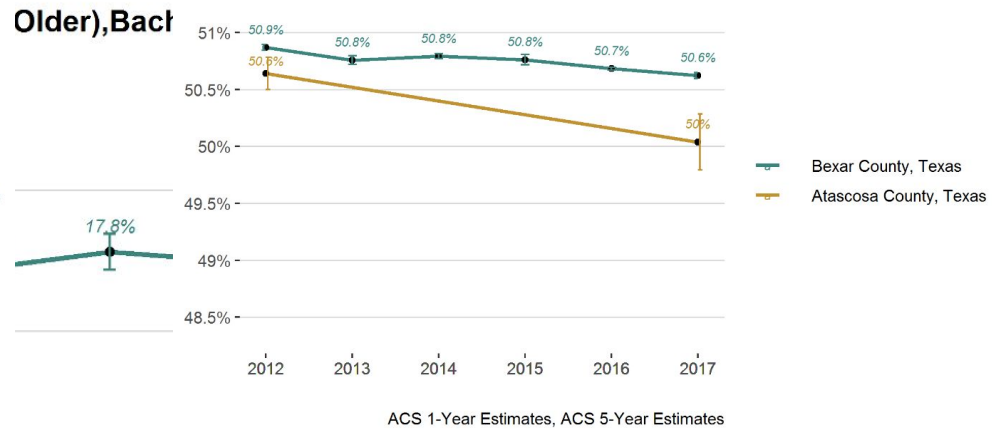
Geoid	Title	Subtitle	Source	Year	Estimate	Margin of Error (Moe)	Min Moe	Max Moe
Atascosa County	Educational Attainment (25 and Older), Bachelors Degree	Highest Degree Obtained	ACS 1-Year Estimates, ACS 5-Year Estimates	2012	8.15	1.38	6.77	9.53
Atascosa County	Educational Attainment (25 and Older), Bachelors Degree	Highest Degree Obtained	ACS 1-Year Estimates, ACS 5-Year Estimates	2017	9.61	1.64	7.97	11.25
Bexar County	Educational Attainment (25 and Older), Bachelors Degree	Highest Degree Obtained	ACS 1-Year Estimates, ACS 5-Year Estimates	2012	16.5	0.59	15.91	17.10
Bexar County	Educational Attainment (25 and Older), Bachelors Degree	Highest Degree Obtained	ACS 1-Year Estimates, ACS 5-Year Estimates	2013	17	0.58	16.42	17.58

EXAMPLE OF PROJECT: AUTOMATION AND CUSTOMIZATION OF TRENDLINES

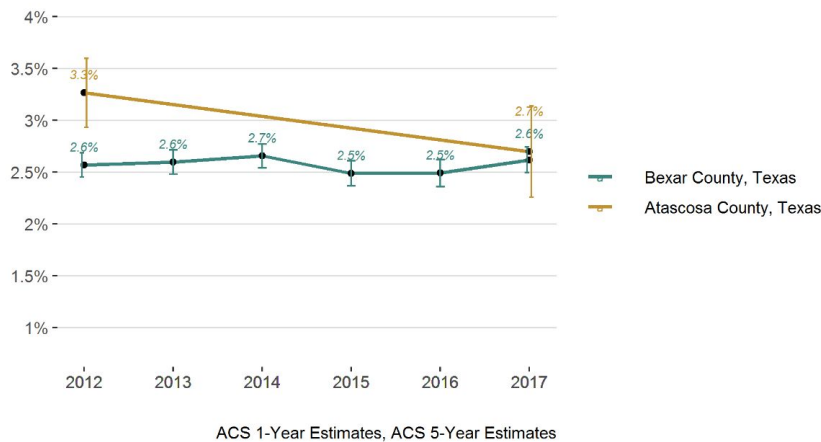
Educational Attainment (25 and Older), Less than High School
Highest Degree Obtained



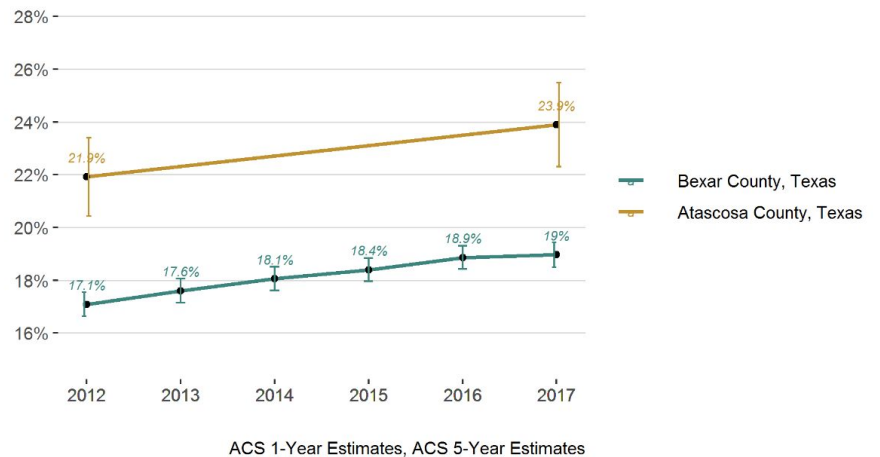
Gender, Female
Percent of Total Population



Disability Status by Age, 75 years or older
Percent of Insured Civilian Non-Institutionalized Population



Age Dependency, Senior Dependency (65 and Older)
Dependant Populations per Adults 18 to 64



```

1- {r}
2 trendpercent$geoid <- factor(trendpercent$GEOID, levels = c("Bexar County, Texas","Atascosa County, Texas"))
3 # create graphing function
4- trend.graph <- function(trendpercent, na.rm = TRUE, ...){
5
6 # create list of title in data to loop over
7 title_list <- unique(trendpercent$CHARTTITLE ) } Create an Index
8
9 # create for loop to produce ggplot2 graphs
10- for (i in seq_along(title_list)) {
11
12 # create plot for each title in trend and adds values for x and y plotting
13 plot <-
14 ggplot(subset(trendpercent, trendpercent$CHARTTITLE ==title_list[i]),
15 aes(x=YEAR, y=ESTIMATE, ymin=minmoe, ymax=maxmoe, group=(GEOID), color=GEOID)) + theme_hc() +
16
17 geom_line(aes(colour=GEOID),size = 1.0) + geom_point(color="black")+
18
19 #labels in bars, paste % symbol, adjust color and position
20 geom_text_repel(aes(label = paste0(round(ESTIMATE,1),"%"), size=2.5 , fontface='italic', force=1,
21 segment.color = 'transparent', min.segment.length = 0, segment.size = 0, point.padding = .3,
22 nudge_x = 0, nudge_y = .1, direction = "y")+
23
24 #add the % symbol to the values in axis and adjust limits to be automated and adds the function "expand".
25 scale_y_continuous(labels=function(ESTIMATE) paste0(ESTIMATE,"%"), expand = expand_scale(mult = c(.1,.3),
26 add = c(1.5, 0)),breaks=pretty_breaks(n=5),limits=c(min(trendpercent$minmoe[trendpercent$CHARTTITLE ==title_list[i]]),
27 ,max(trendpercent$maxmoe[trendpercent$CHARTTITLE ==title_list[i]])))+
28
29 scale_x_continuous(breaks=c(trendpercent$YEAR))+
30
31 #specify colors for lines
32 scale_color_manual(values=c("Bexar County, Texas"="#3c857f", "Atascosa County, Texas"="#c09231"))+
33
34 #modify legend
35 theme(legend.position="right", legend.title = element_blank(),legend.spacing.x = unit(.6, 'cm'),
36
37 #remove panel, set background color and reorder items in legend
38 panel.grid.major = element_blank(), panel.background = element_rect(fill = "white")) +
39
40 #error bars
41 geom_errorbar(aes(ymin=minmoe, ymax=maxmoe), width=.1, position = position_dodge(0.1)) +
42
43 #automates title and specified x and y labels as blank
44 ggtitle(title_list[i])+labs( x="", y="", caption = trendpercent$SOURCEID [trendpercent$CHARTTITLE ==title_list[i]],
45 subtitle = trendpercent$CHARTSUBTITLE[trendpercent$CHARTTITLE ==title_list[i]]) +
46 theme(plot.title=element_text (hjust = 0, face='bold', size=12, margin=margin(0,0,10,0))) +
47 theme(plot.subtitle = element_text(hjust=0, size=10,margin=margin(0,0,20,0)))+
48 theme(plot.margin = unit(c(.5,.9,.5,.5), "cm"))#top, right, bottom, left
49
50 #Export as images
51 ggsave(paste0(title_list[i],".png"))
52
53 # print plots to screen
54 print(plot)
55 }
56 # run graphing function
57 trend.graph(trendpercent, ESTIMATE)
58

```

Y Axis
X Axis
Group
Min Moe
Max Moe

Labels

Limits
and
breaks

Color of
lines

Legend

Background

Error bars

Title,
subtitle,
source

Save as an
image

Resources



News & Tutorials

R-bloggers

Blogs related to R and its applications

<https://www.r-bloggers.com/>

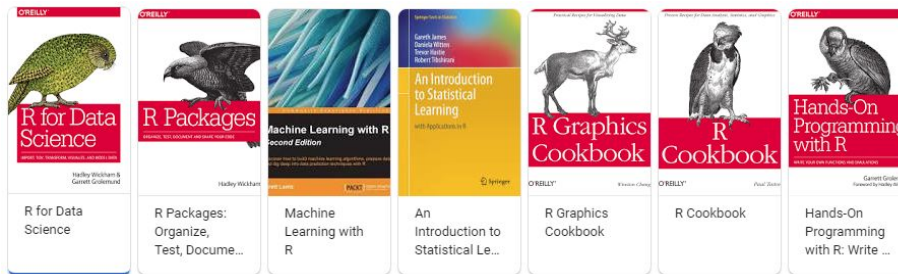
R Graph Gallery

Examples of visualizations with code samples

<https://www.r-graph-gallery.com/>



Books



Troubleshooting

Rdocumentation

Manuals and information for packages

<https://www.rdocumentation.org/>

Stackoverflow

Developers share knowledge

<https://stackoverflow.com/>



Quick Tips

Rtips

List of common tasks performed in R

<http://pj.freefaculty.org/R/Rtips.html>.

ImpatientR

Introduction to basic functions

<https://www.burns-stat.com/documents/tutorials/why-use-the-r-language/>

YaRrr! The Pirates Guide to R

Intro to basic analytical tools in R, from basic coding and analyses, to data wrangling, plotting, and statistical inference.

<https://bookdown.org/ndphillips/YaRrr/>



Thanks!

Questions?

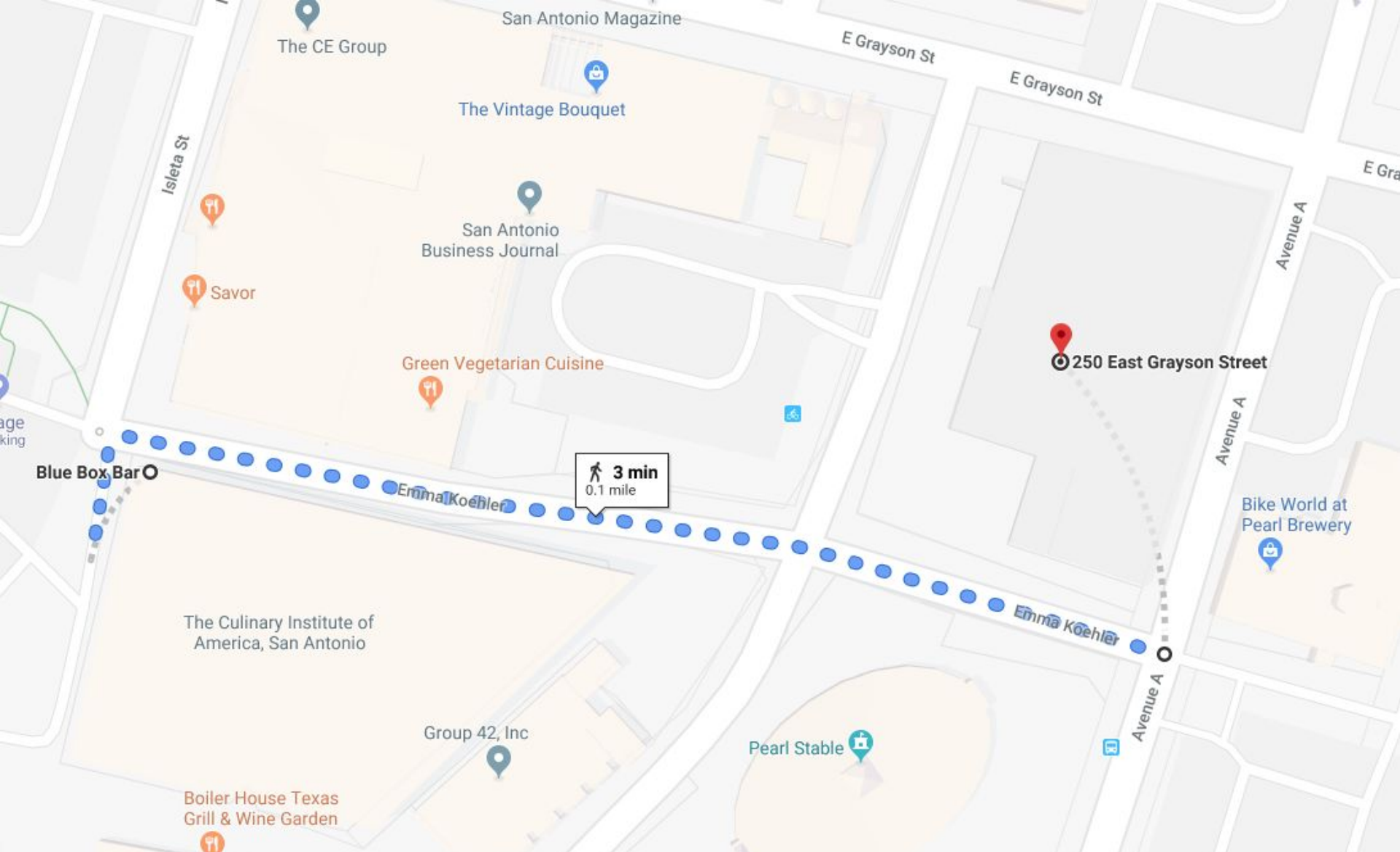
Paulina Cano:

paulina.canomccutcheon@uth.tmc.edu

Jamie Ford:

jamie.ford@nisd.net





Data Drinks @ 4:30 PM

EST. 2012

BLUE BOX

BAR