

Visualizing Data

Summer Training Academy for Research Success

Michael F. Seese

Department of Political Science
University of California San Diego

Summer 2021

Outline

The Visual Presentation of Data

Base Graphics

The Grammar of Graphics

Replication code is available on GitHub

 https://github.com/mfseese/STARS_Summer2021

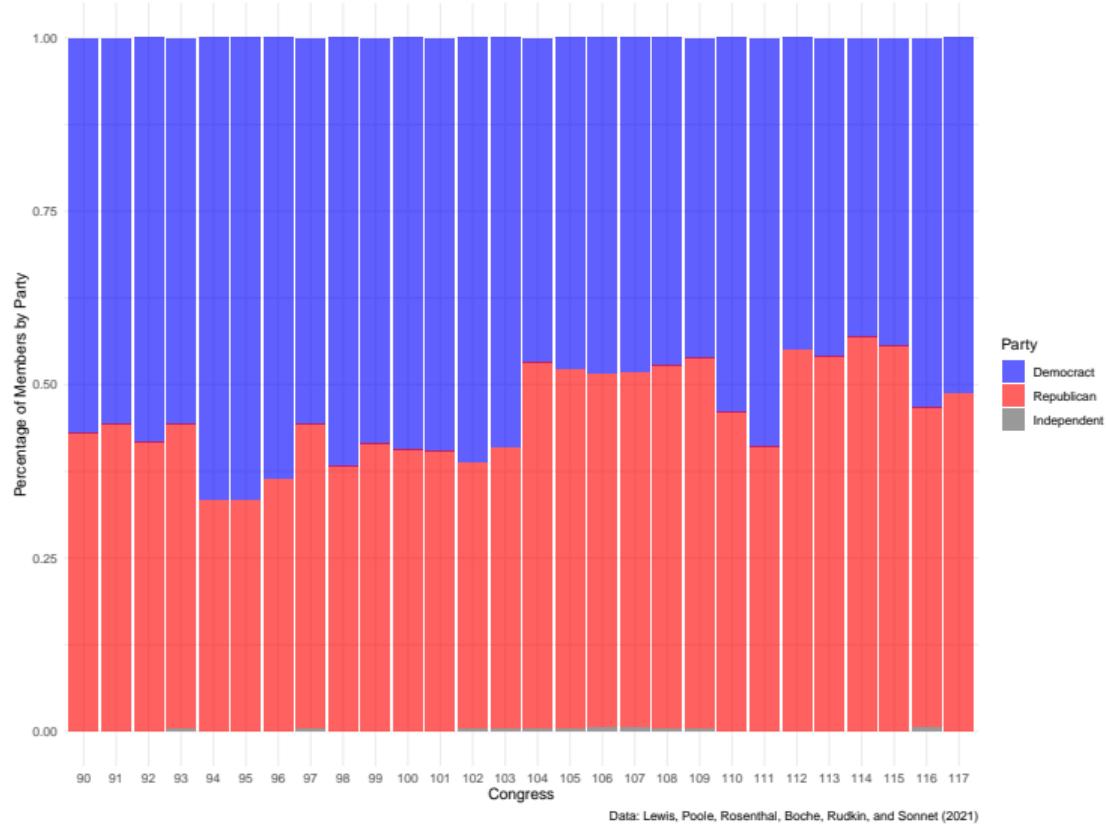
The Visual Presentation of Data

Why Should We Visualize Data?

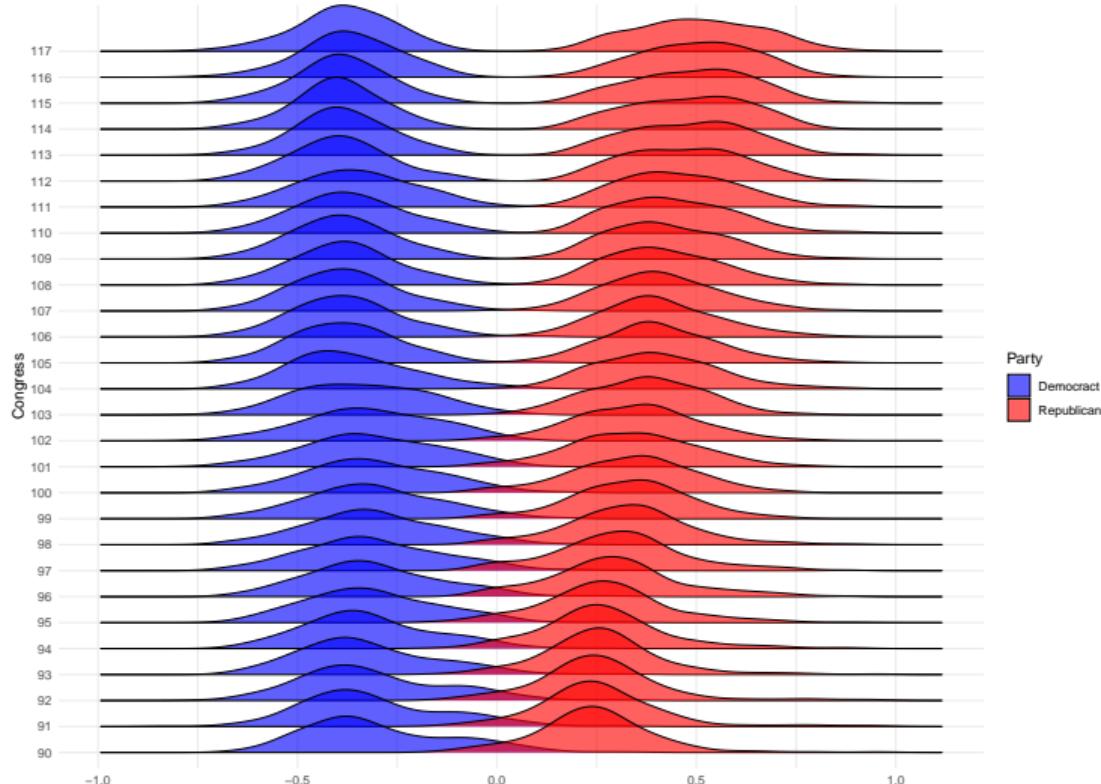
1. To make sense of the data
 - 1.1 Identify broad patterns
 - 1.2 Check our assumptions and intuition
2. To tell stories
 - 2.1 Underscore the **puzzle**
 - 2.2 Present the results of our analysis

	A	B	C	D	E	F	G	H	I	K	L	M	N	O	P	Q	R	S	T	U	V	
1	member_id	upper_state_gov	lower_state_gov	district_code	state_id	state_name	party	ideology	last_meeting	name	longitude	latitude	icon	dead	senateate	dist1	nominate_log	dist2	nominate_log	dist3	nominate_log	dist4
49851	117 Senate	48398	30	O MA	180	KENYON, Robert K.	DEM	1954	1954	-72.295	42.39	0.000	0	0	0.295	0.000	0.000	0.000	0.000	-0.402	0.217	
49852	117 Senate	48399	30	O MA	180	KENNEDY, Ted	DEM	1954	1954	-71.045	42.39	0.000	0	0	0.295	0.000	0.000	0.000	0.000	-0.402	0.217	
49853	117 Senate	29934	46	O MA	280	KENYON, Roger F.	DEM	1954	1954	-7.576	8.933	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.371	
49854	117 Senate	41303	46	O MA	280	KING-SMITH, Cindy	DEM	1954	1954	0.39	0.281	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.434	
49855	117 Senate	41304	46	O MA	280	KIRK, Jim	DEM	1954	1954	0.39	0.281	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.395	
49856	117 Senate	41305	34	O MA	280	KIRKLEY, John David	DEM	1954	1954	0.37	-0.357	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.438	
49857	117 Senate	25138	46	O MA	280	KIRKLEY, Steve	DEM	1954	1954	0.36	-0.381	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.318	
49858	117 Senate	48703	64	O MA	180	KISTER, Jim	DEM	1954	1954	-2.225	61.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.382	
49859	117 Senate	48704	64	O MA	180	KLEIN, Jeff	DEM	1954	1954	0.36	0.295	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.398	
49860	117 Senate	41306	35	O MA	280	KLINE, Benjamin Eric	DEM	1954	1954	0.36	0.295	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.369	
49861	117 Senate	21743	65	O MA	180	KLINE, Jacqueline Sheryl	DEM	1954	1954	-2.395	62.283	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.369	
49862	117 Senate	41307	46	O MA	280	KLINE, Jennifer Marie	DEM	1954	1954	0.36	0.295	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.398	
49863	117 Senate	48996	4	O MA	180	KNAUF, Jerome	DEM	1954	1954	0.37	-0.143	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.318	
49864	117 Senate	41302	46	O MA	280	KNAPP, Margaret [Maggie]	DEM	1954	1954	-2.315	62.306	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.315	
49865	117 Senate	41308	22	O MA	180	KNOX, James	DEM	1954	1954	0.36	0.295	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.377	
49866	117 Senate	29810	66	O MA	180	KOONIK, Gary Anthony	DEM	1954	1954	0.36	-0.261	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.418	
49867	117 Senate	29810	66	O MA	180	KOONIK, Martin	DEM	1954	1954	-2.306	62.021	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.384	
49868	117 Senate	29810	66	O MA	180	KOONIK, William	DEM	1954	1954	-2.367	62.025	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.383	
49869	117 Senate	29810	66	O MA	180	KOONIK, William	DEM	1954	1954	-2.306	62.021	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.388	
49870	117 Senate	28735	3	O MA	180	KOOLBECK, Kristen	DEM	1954	1954	-0.474	-0.405	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.375	
49871	117 Senate	29846	47	O MA	280	KOHL, Richard M.	DEM	1954	1954	0.441	-0.096	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.382	
49872	117 Senate	21350	36	O ND	280	KOHLER, Robert Edward [Thom]	DEM	1954	1954	-0.474	-0.177	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.377	
49873	117 Senate	41337	36	O ND	280	KOMAROFF, Kevin	DEM	1954	1954	0.39	0.383	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.418	
49874	117 Senate	41337	36	O ND	280	KOHREN, John	DEM	1954	1954	0.34	0.233	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.315	
49875	117 Senate	29848	24	O OH	180	KOHLBERG, Robert Jones [Bob]	DEM	1954	1954	0.36	0.206	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.377	
49876	117 Senate	29848	24	O OH	180	KOHLBERG, Steven	DEM	1954	1954	0.36	-0.111	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.377	
49877	117 Senate	15424	55	O OK	280	KOHRT, James Mountain	DEM	1954	1954	0.554	0.045	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.358	
49878	117 Senate	21184	55	O OK	280	KOHLBERG, James	DEM	1954	1954	0.36	0.156	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.361	
49879	117 Senate	21184	55	O OK	280	KOHLBERG, James	DEM	1954	1954	0.36	0.156	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.361	
49880	117 Senate	41308	72	O OR	180	KOHLBERG, James	DEM	1954	1954	0.35	0.470	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.399	
49881	117 Senate	29810	55	O OR	180	KOHLBERG, James	DEM	1954	1954	0.36	0.470	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.399	
49882	117 Senate	29810	55	O OR	180	KOHLBERG, James	DEM	1954	1954	0.36	0.470	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.399	
49883	117 Senate	29810	55	O OR	180	KOHLBERG, James	DEM	1954	1954	0.36	0.470	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.399	
49884	117 Senate	29810	55	O OR	180	KOHLBERG, James	DEM	1954	1954	0.36	0.470	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.399	
49885	117 Senate	29810	55	O OR	180	KOHLBERG, James	DEM	1954	1954	0.36	0.470	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.399	
49886	117 Senate	29810	55	O OR	180	KOHLBERG, James	DEM	1954	1954	0.36	0.470	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.399	0.399	
49887	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49888	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49889	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49890	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49891	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49892	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49893	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49894	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49895	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49896	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49897	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49898	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49899	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49900	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49901	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49902	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49903	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49904	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49905	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49906	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49907	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49908	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.295	0.000	0.000	0.000	0.000	0.482	0.031	
49909	117 Senate	29814	37	O SD	280	KOHLBERG, James	DEM	1954	1954	0.425	0.144	0.000	0	0	0.29							

Party Composition of the U.S. House, 1967 – Present

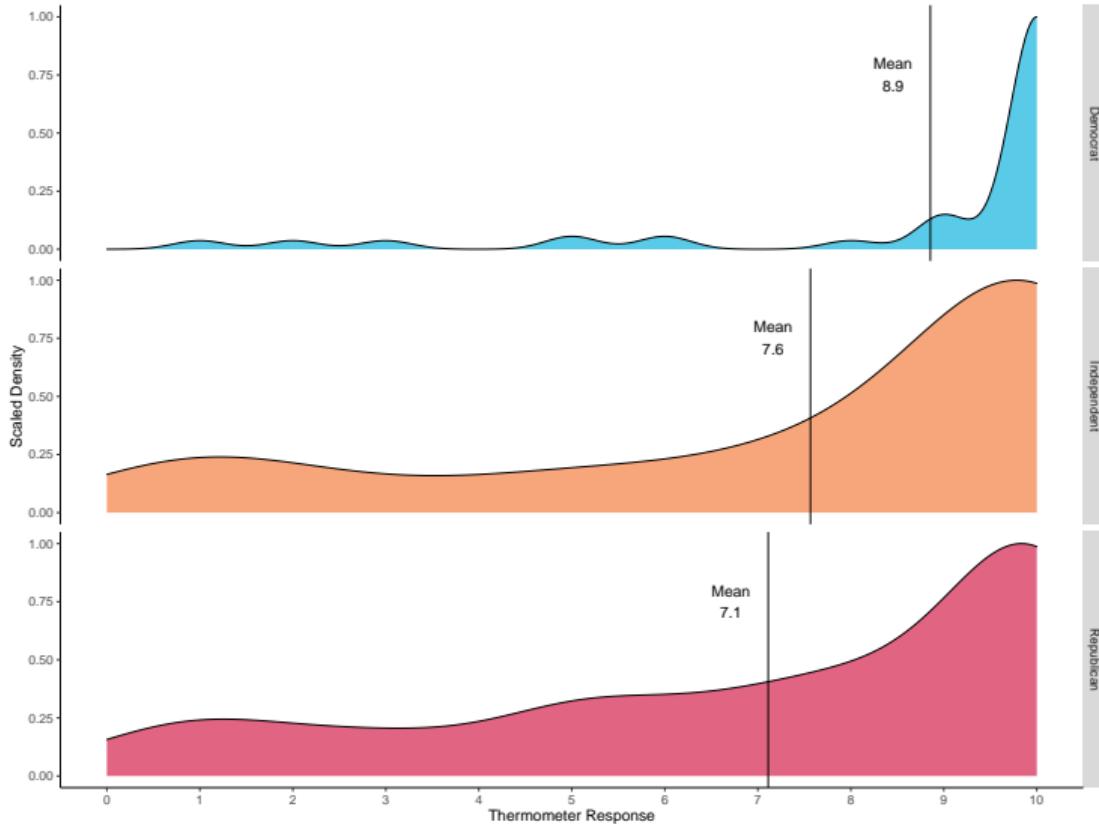


The Growing Ideological Divide in the U.S. House

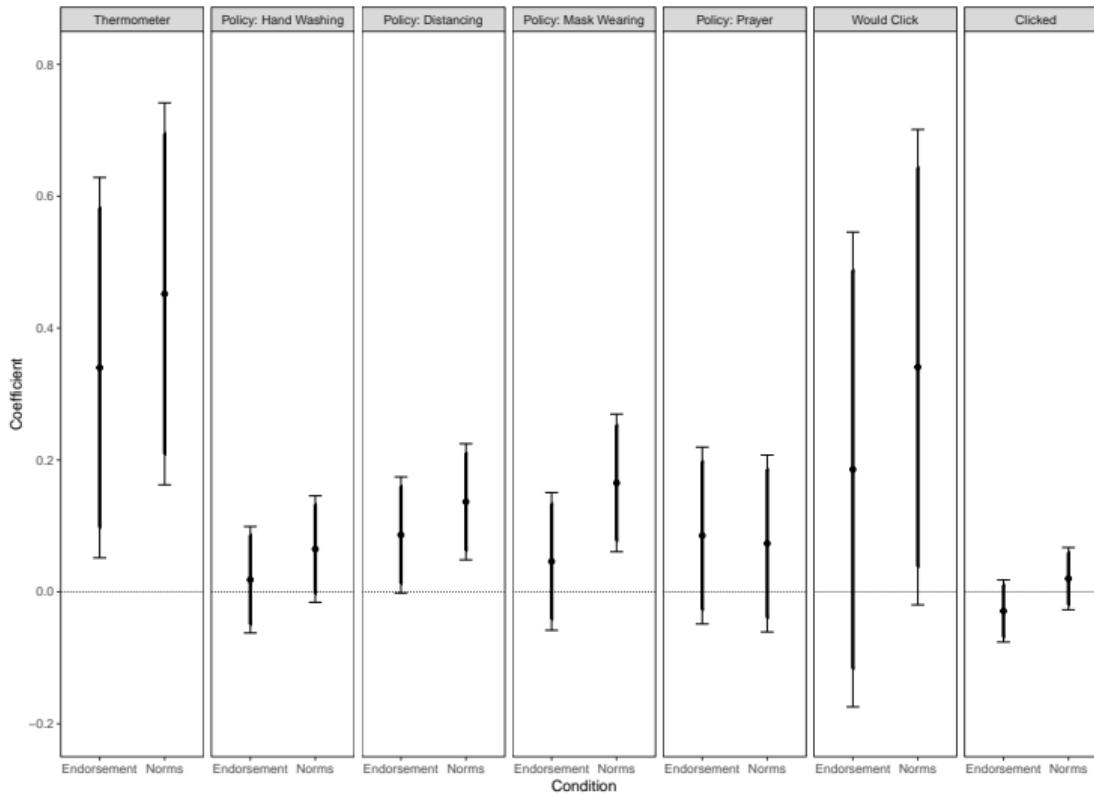


Data: Lewis, Poole, Rosenthal, Boche, Rudkin, and Sonnet (2021)

Mask Wearing Thermometer Response by Party: Control Group



Coefficient Plot: OLS Model with Controls



Shown with 90% and 95% confidence intervals.
Note: Model for "Clicked" outcome uses subset data (only petition-eligible respondents).

Worth a Thousand Words Data Points

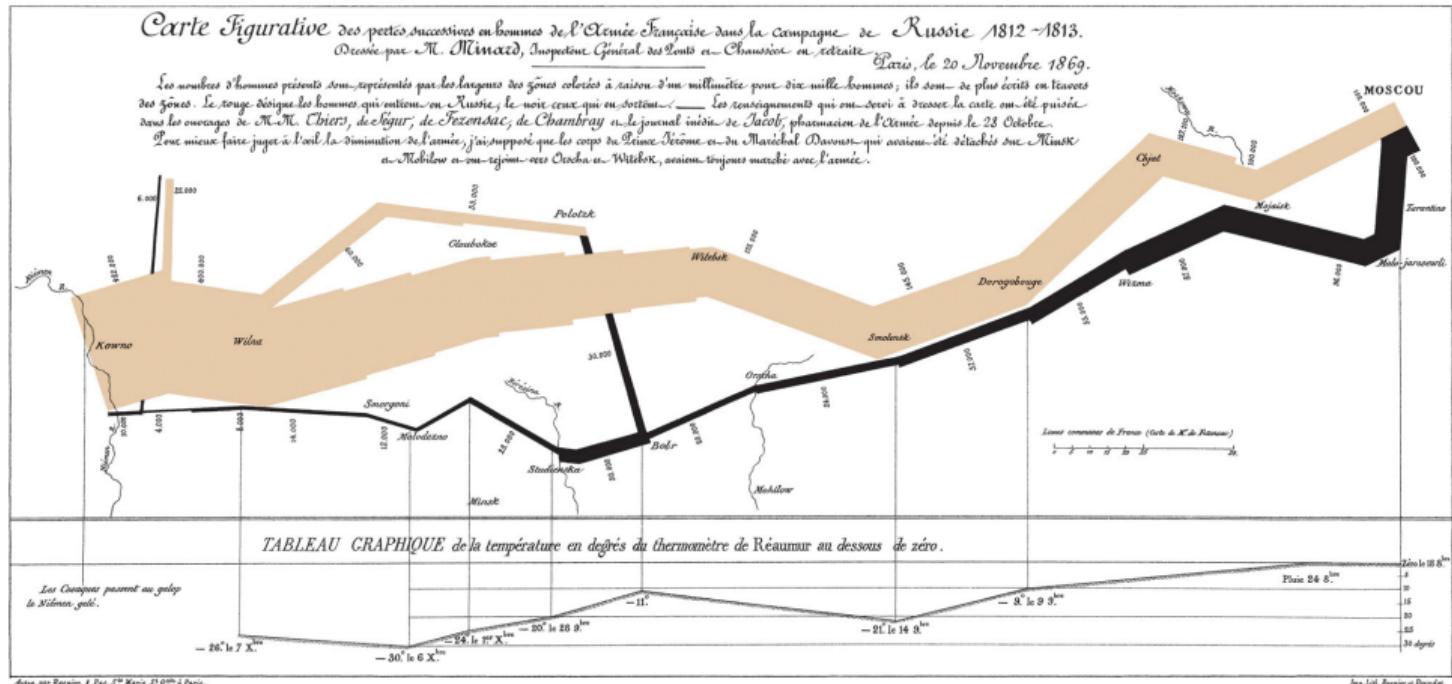


Figure: Minard's *Carte Figurative*, 1869

The Process

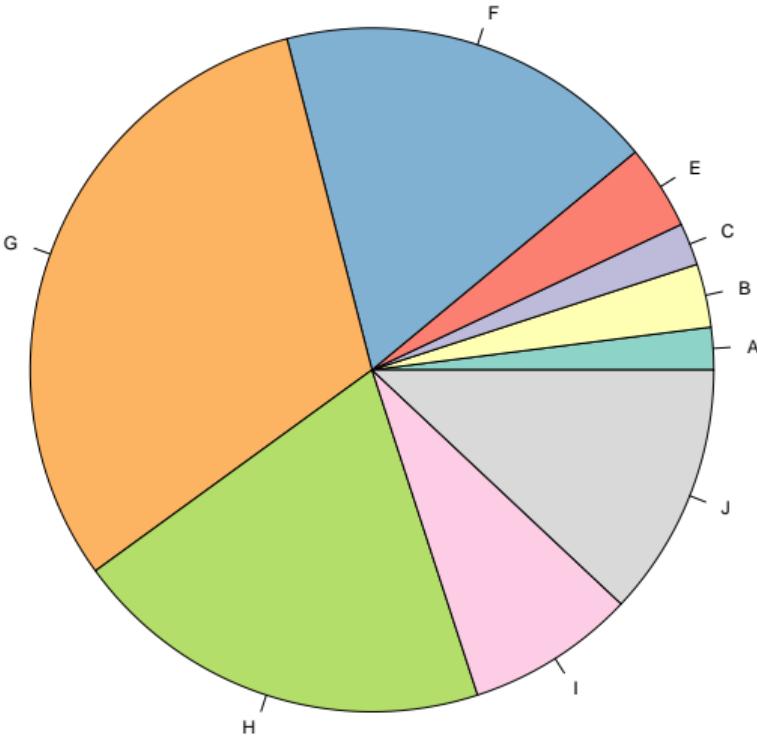
1. Select the appropriate plot type
2. Generate the plot with **proper** formatting
3. Ask yourself (and other people): “Is this legible? Is this intuitive?”
4. Revise with aesthetics in mind
5. Clean your code to ensure replicability

Plot Types

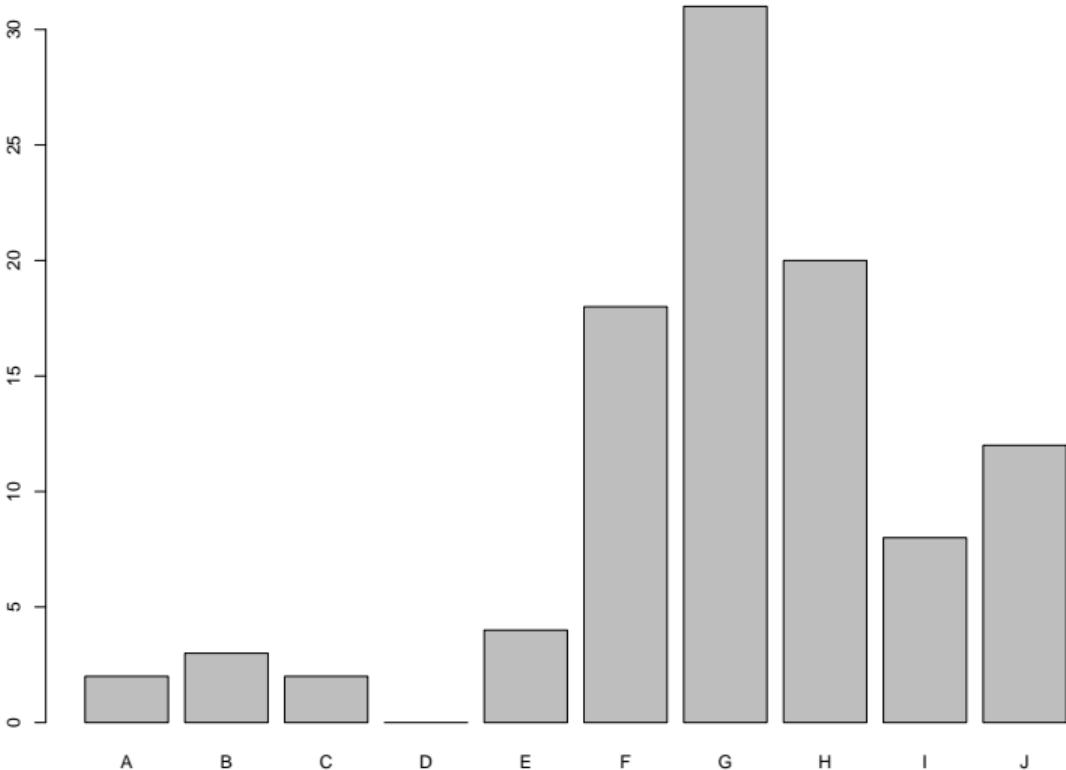
Table: Selecting the Appropriate Plot

1 Variable	≥ 2 Variables			
Distributions	Comparison	Relationship	Results	Avoid
Histogram	Bar*	Scatter	Coefficient Plot	Pie
Density Plot	Dot Plot*	Bubble	Margins Plot	Donut
Box Plot	Line	Fit Plots	Diagnostic	Stacked Area
Violin Plot	Time Series			
Thematic Maps	Range Plots			

<i>x</i>	<i>y</i>
A	2
B	3
C	2
D	0
E	4
F	18
G	31
H	20
I	8
J	12

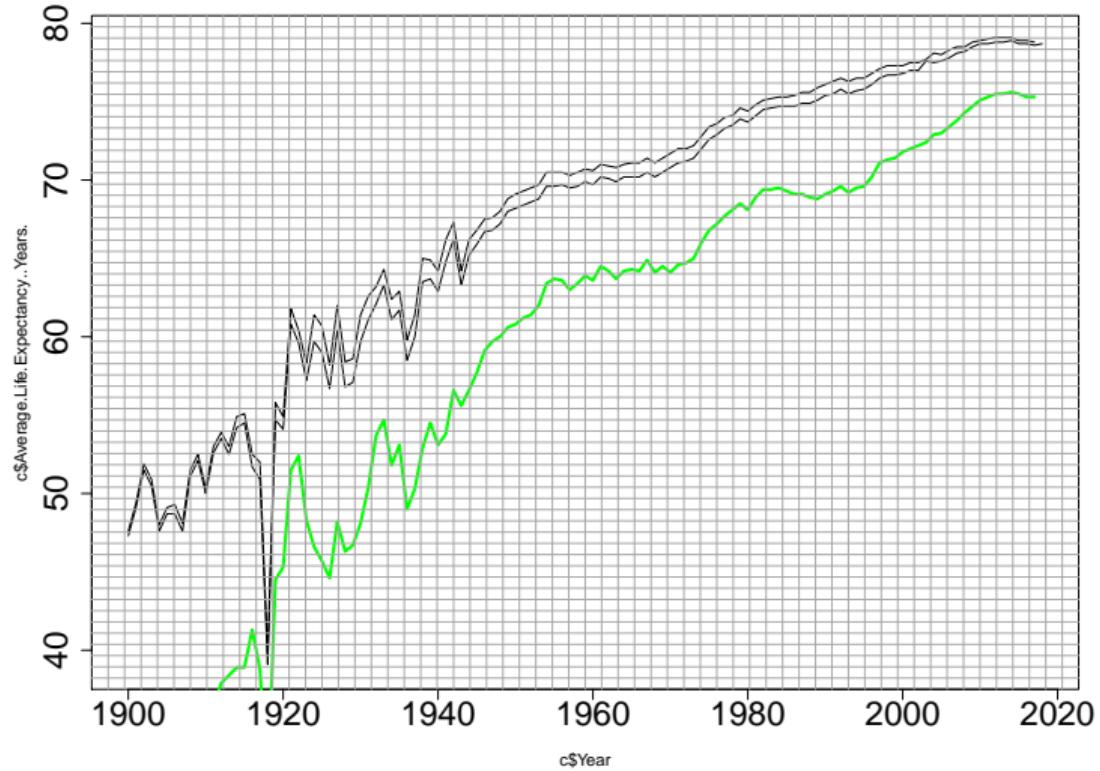


x	y
A	2
B	3
C	2
D	0
E	4
F	18
G	31
H	20
I	8
J	12

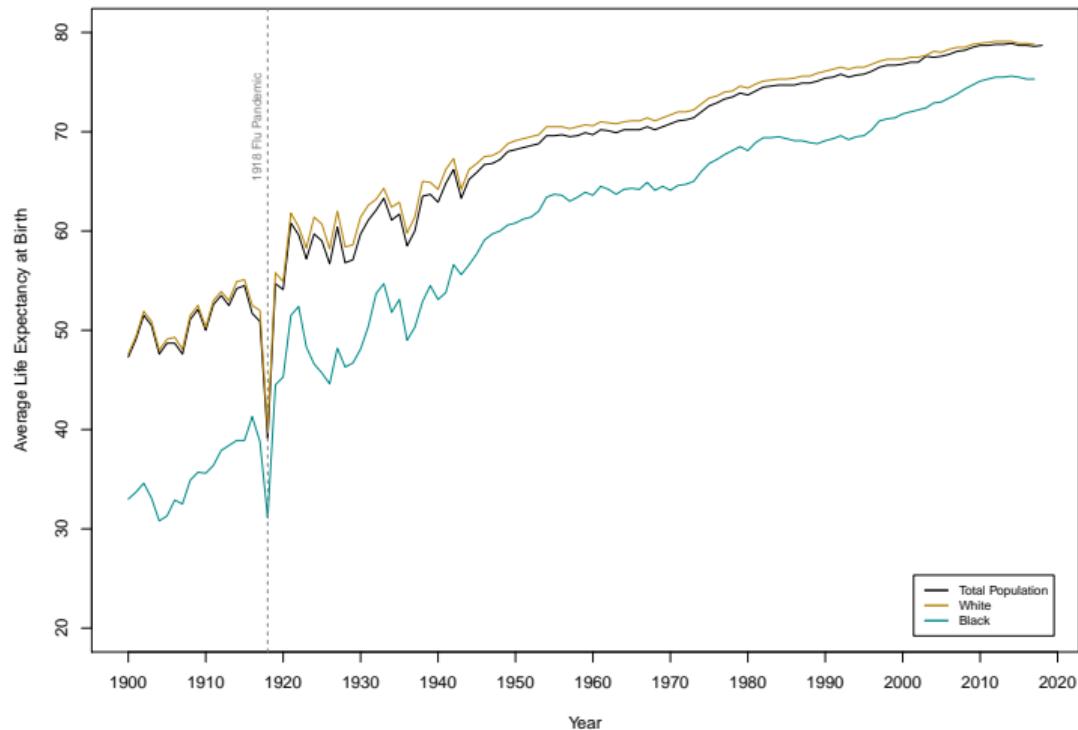


Proper Formatting

- ▶ Include
 - ▶ A (descriptive) title
 - ▶ A note or caption with data sources and supplementary information
 - ▶ Marker labels where appropriate
 - ▶ A complete legend
 - ▶ Confidence intervals when possible
- ▶ Display
 - ▶ All axes and their labels
 - ▶ The full range of relevant variation
- ▶ Select appropriate color palette
- ▶ Always relabel variables for display
- ▶ Ensure text is legible (size, direction, etc.)



U.S. Life Expectancy by Race and Year, 1900 – 2018

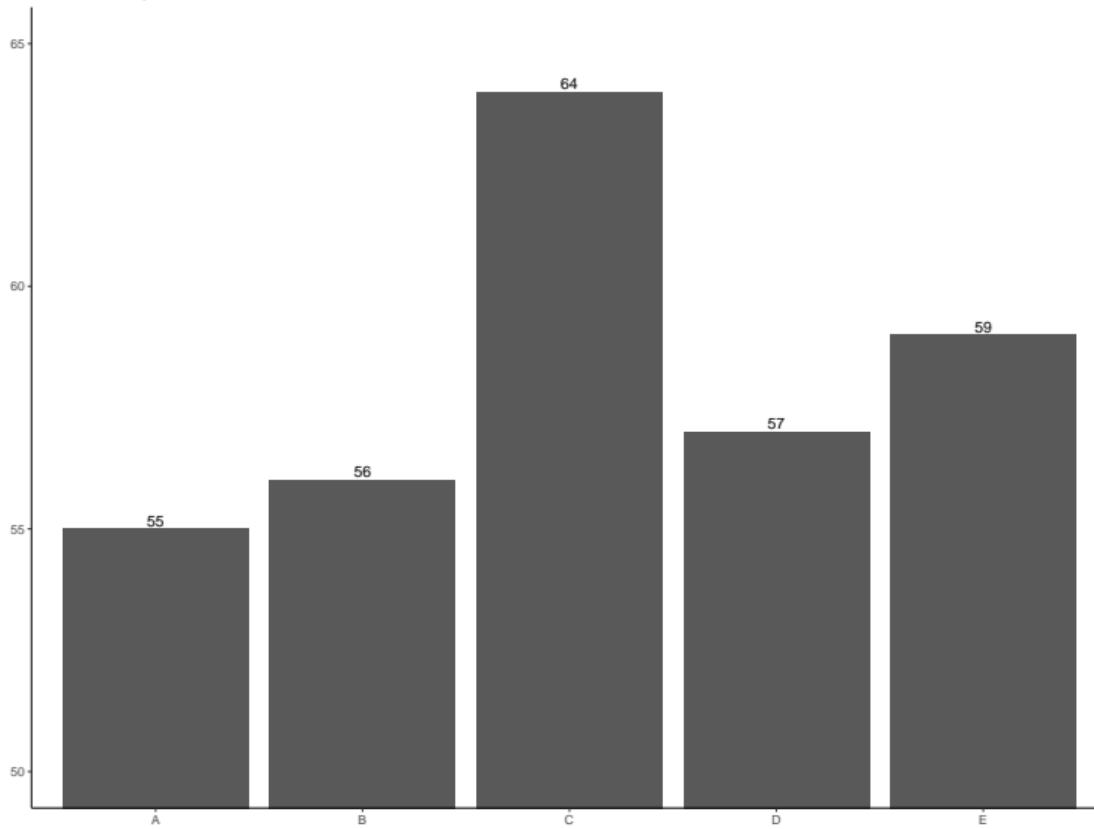


Data: U.S. Centers for Disease Control and Prevention, National Center for Health Statistics

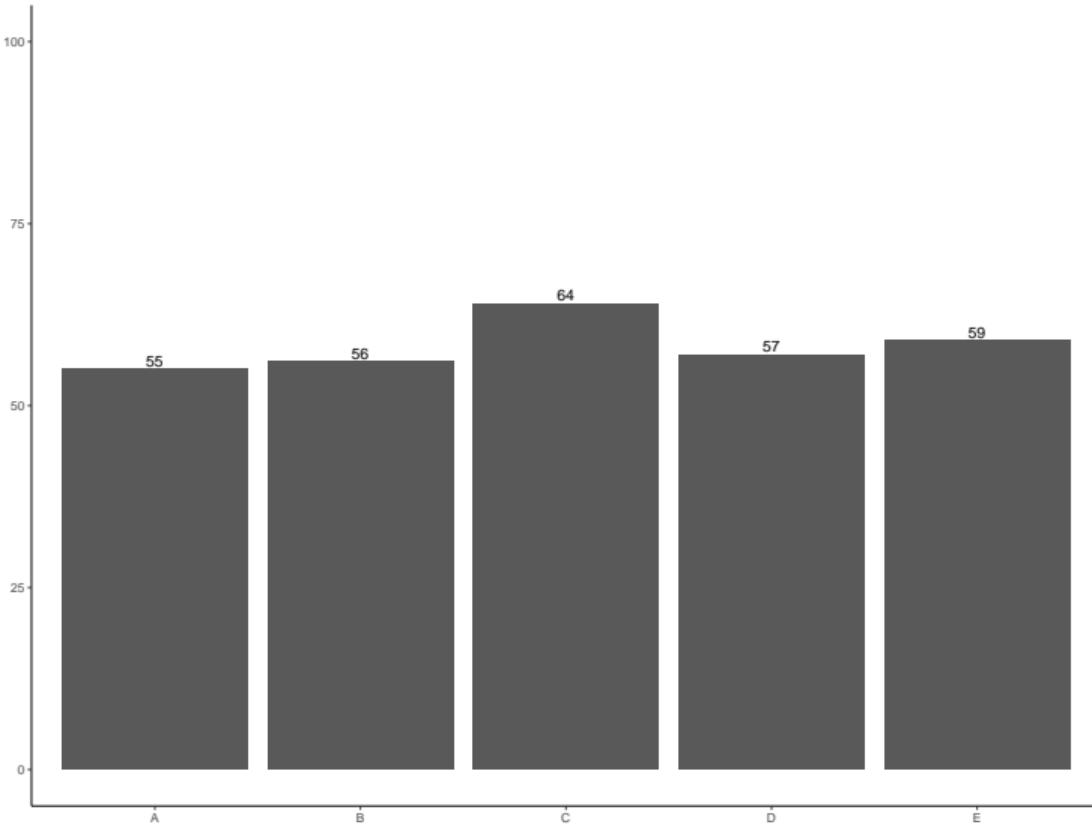
Legibility & Aesthetics

- ▶ Design for the “interocular”
- ▶ Avoid deceptive or confusing presentation of the data
- ▶ Maximize the data-to-ink ratio
- ▶ Optimize ticks and gridlines
- ▶ Ensure plot is as accessible as possible (e.g., colorblind palette, etc.)
- ▶ Simplify wherever possible

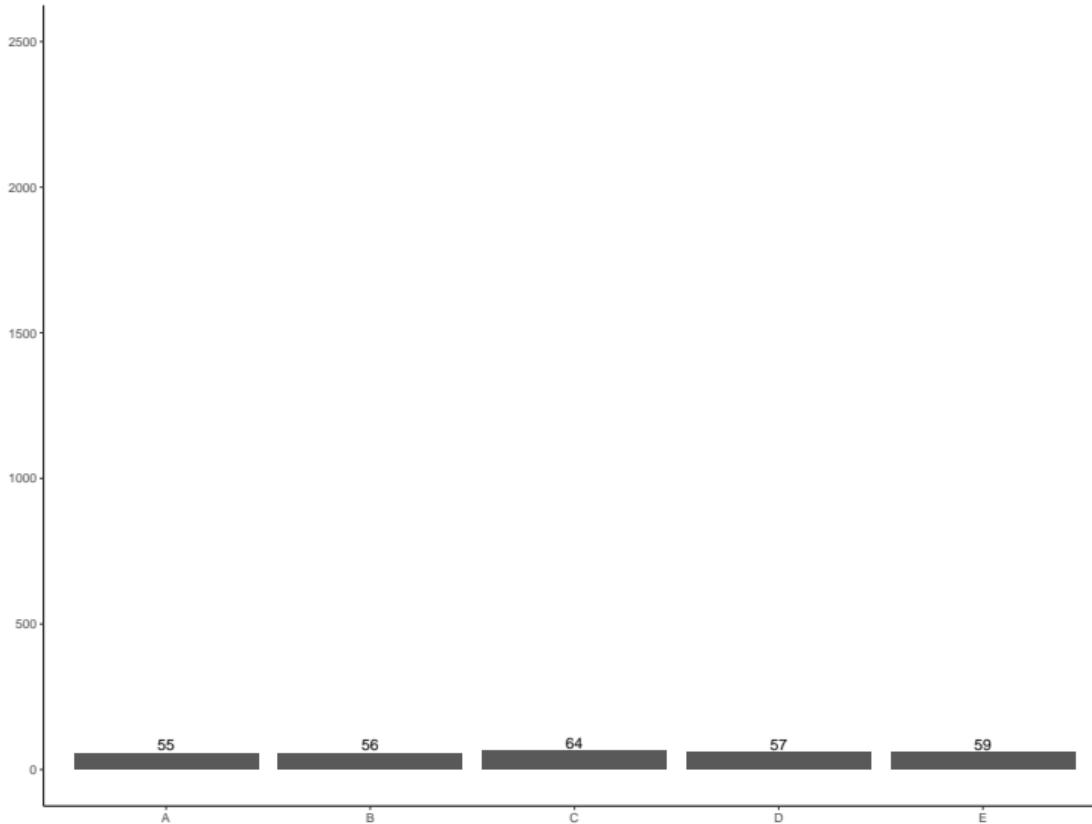
This chart emphasizes the extreme value of C



This chart implies more uniformity

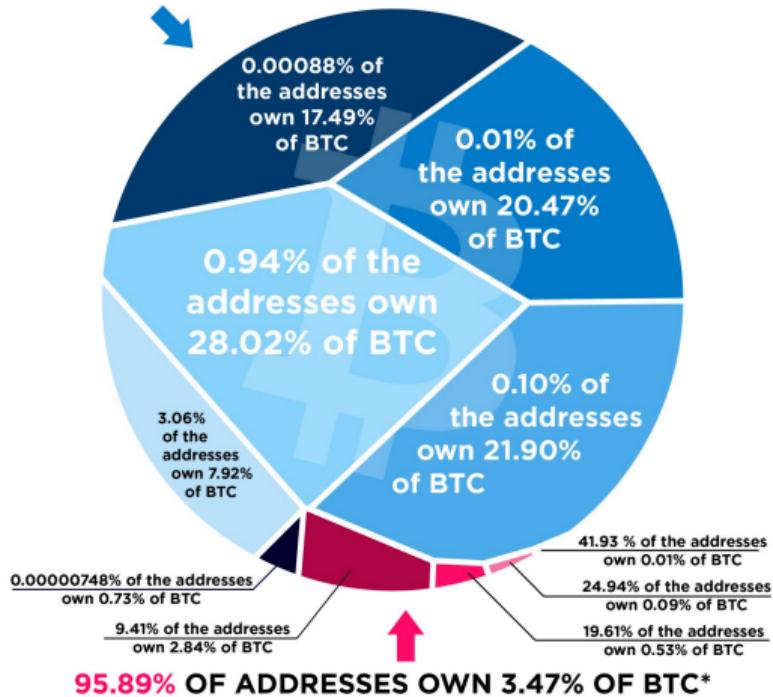


This chart gives the perception of equivalence



The Bitcoin Wealth Distribution

4.11% OF ADDRESSES OWN 96.53% OF BTC*

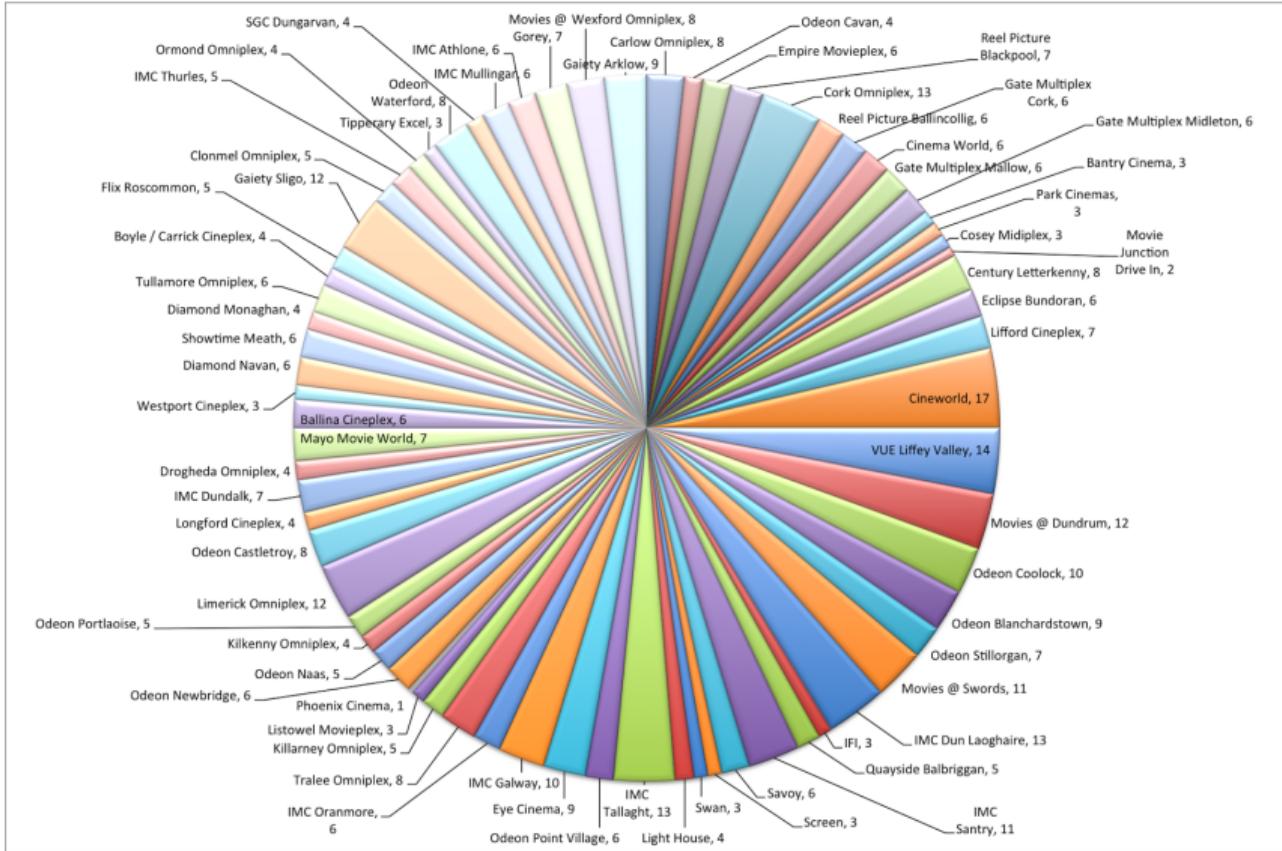


* Data as of September 12th, 2017

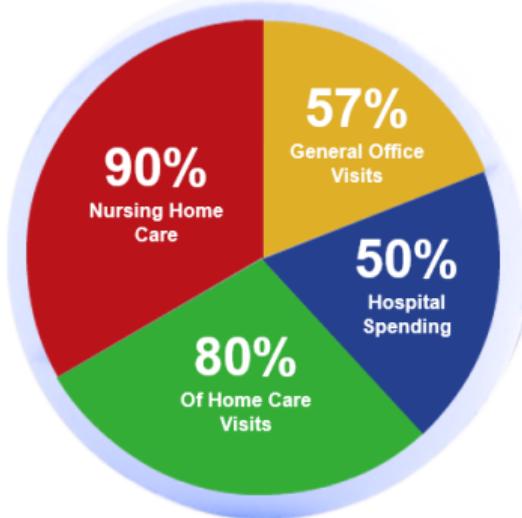
Article and Sources:

<https://howmuch.net/articles/bitcoin-wealth-distribution>

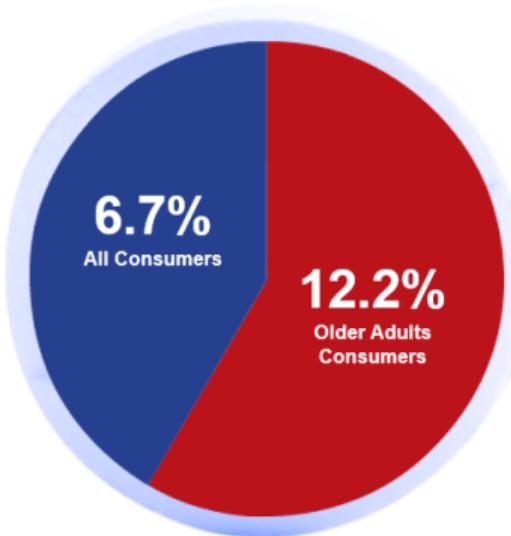
<https://bitcoinprivacy.net/>



The Numbers on Older Adults



The Majority



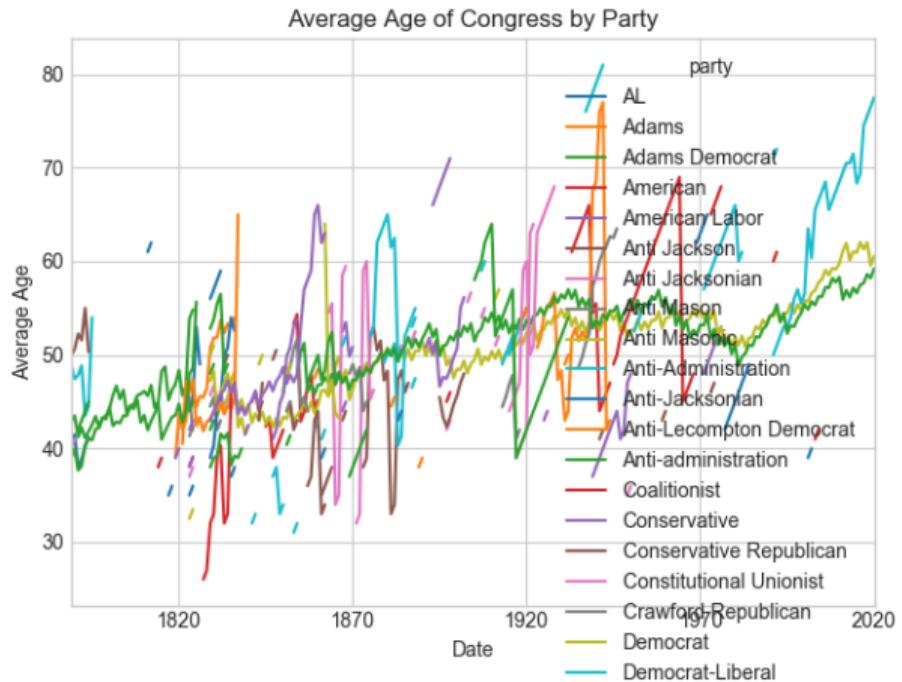
Healthcare Spending

CORONAVIRUS

FLORIDA CONFIRMED CASES

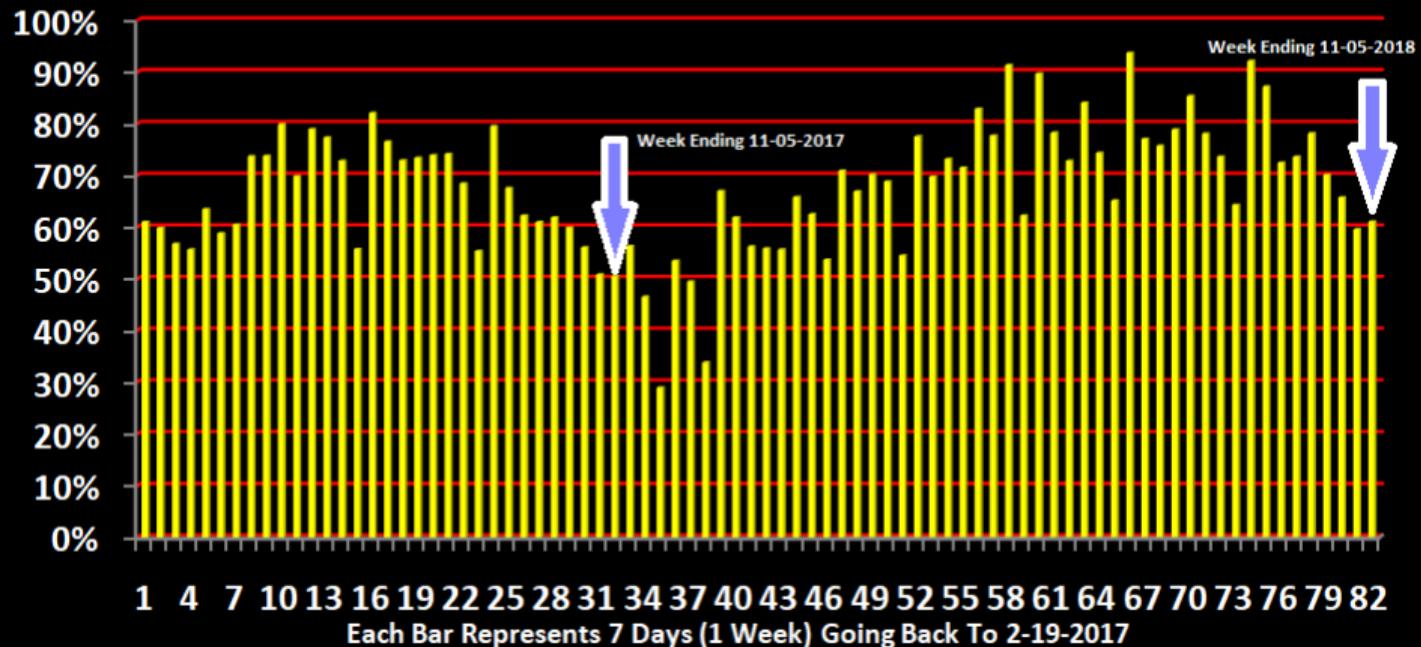


SOURCE: FLORIDA DEPT. OF HEALTH

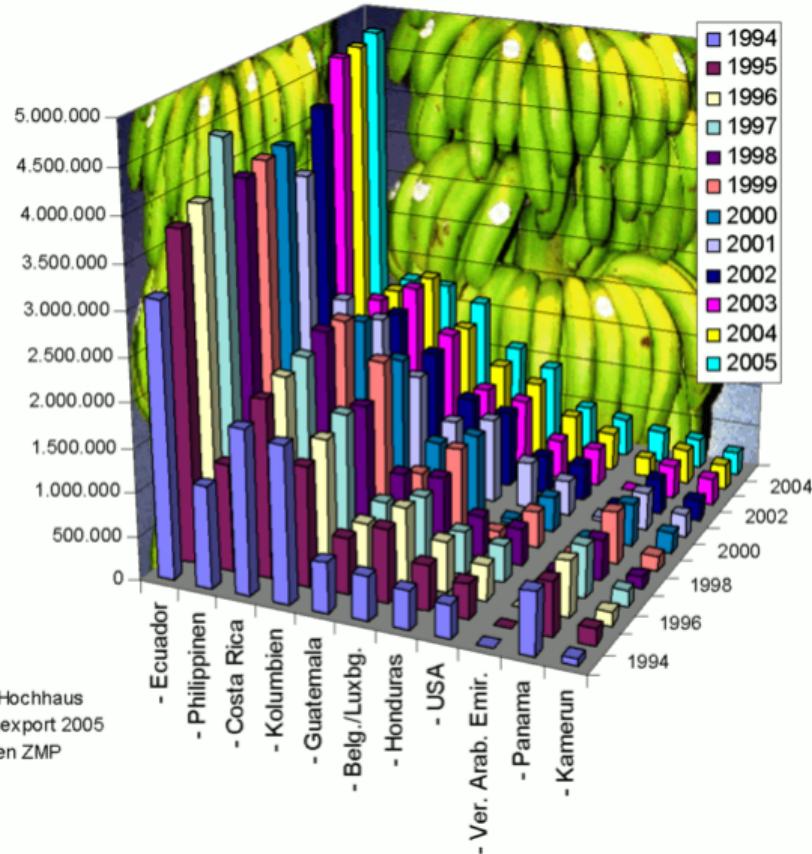


List/Pend

SimpsonRealEstateNW.Com



Export von Bananen in Tonnen von 1994-2005



Dr. Hochhaus
Banlexport 2005
Daten ZMP

Legibility & Aesthetics

Think through the presentation medium:

- ▶ Video and image compression tend to eat at fine details; best to go a bit chonky
- ▶ Pure red, green, and blue all tend to bleed, but because of how math works red is most noticeable
- ▶ Slides, digital papers, and paper papers all require different font and color selections
- ▶ In general, save as .pdf or other vector format (\LaTeX handles .pdf well)

Replicability

- ▶ Keep a lab notebook
 - ▶ Write down **everything**
 - ▶ Bookmark useful guides
- ▶ Commit to a filing convention and stick to it
- ▶ Use a repository such as GitHub for versioning
- ▶ Save an archive copy of your data and never touch it again
 - ▶ Work from duplicate in your WD
 - ▶ Try to avoid overwriting files
- ▶ **Everything** gets done in code
- ▶ Extensively comment your code
- ▶ Save often
 - ▶ Make sure autosave is on
 - ▶ Back up to the cloud or a secondary drive

Base Graphics

Base Graphics

- ▶ Play with some data compiled by IPI-WZB
Data available on GitHub and at this link: <https://wzb-ipi.github.io/corona/>
- ▶ Includes country-level Covid data and various political correlates

Base Graphics

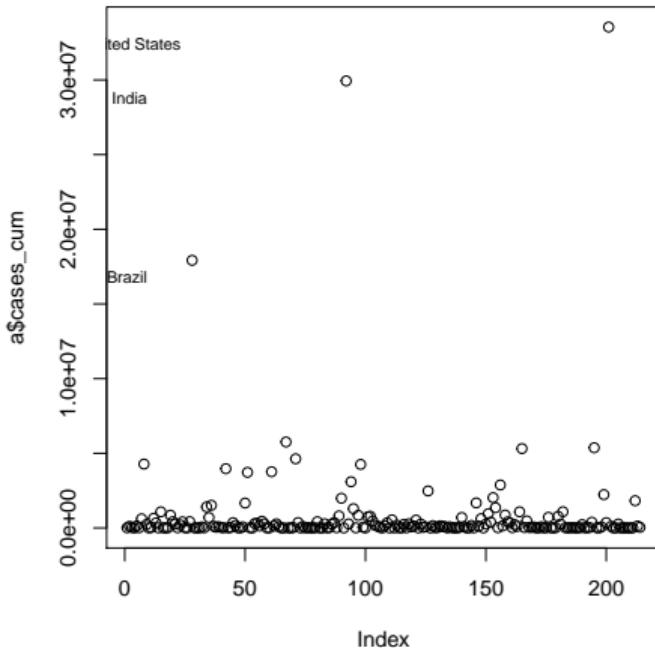
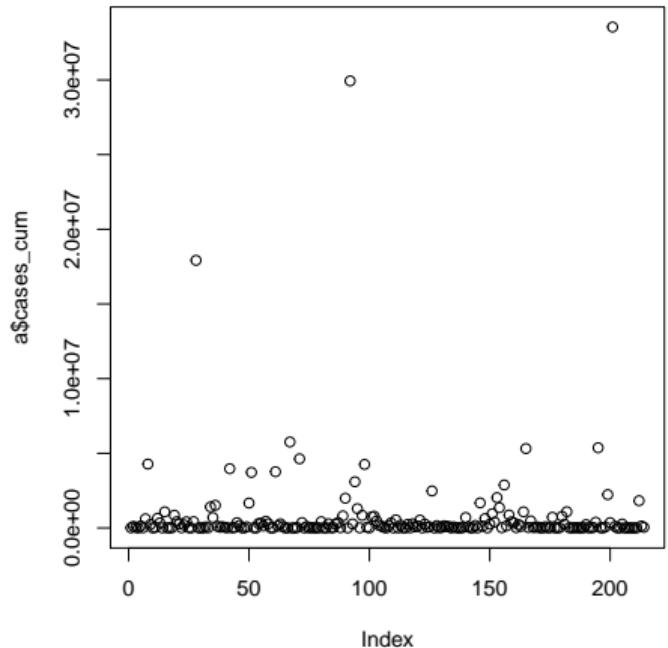
General Format: `plot(x, y, ...)`

```
1 model <- lm(y ~ x1, data = df)
2
3 pdf("plot.pdf", width = 11, height = 8.5)
4
5 par(mfrow = c(1, 2))
6
7 plot(df$x, df$y, type = "p")
8 abline(model, col = "red")
9
10 plot(df$x2, df$y,
11       xlab = "x label",
12       main = "title")
13
14 dev.off()
```

Base Graphics

```
1 a <- read.csv("wzb_covid_june2021.csv", header = TRUE)
2
3 # What happens if we just use the command "plot"?
4 plot(x = a$cases_cum)
5
6 # There are some wild outliers. What countries are these?
7 plot(x = a$cases_cum)
8 text(a$cases_cum[a$cases_cum > 15000000],
9      labels = a$X[a$cases_cum > 15000000],
10      cex = 0.75, pos = 1)
```

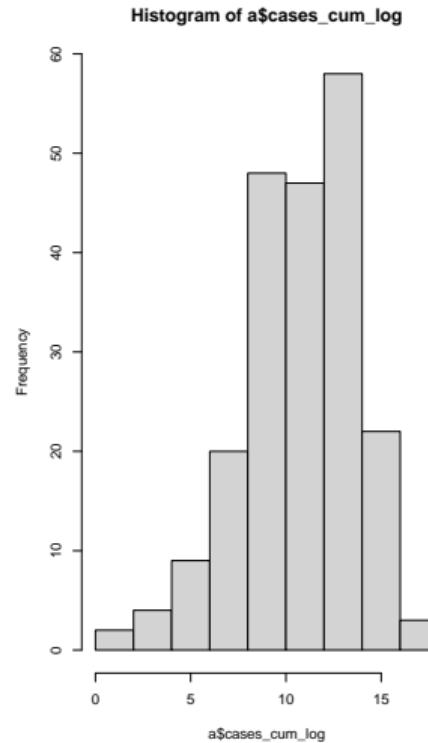
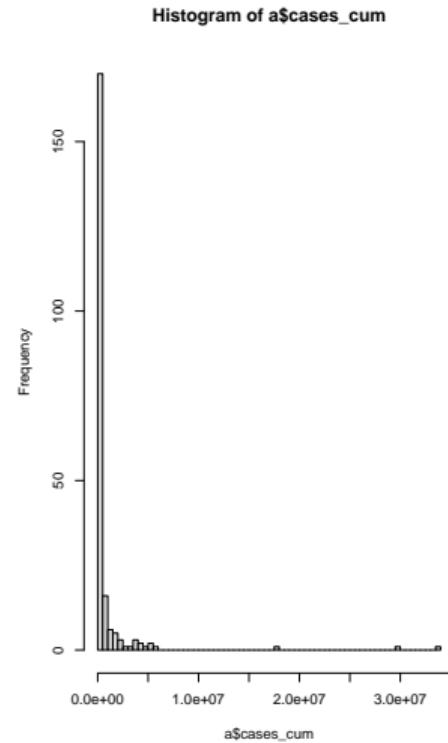
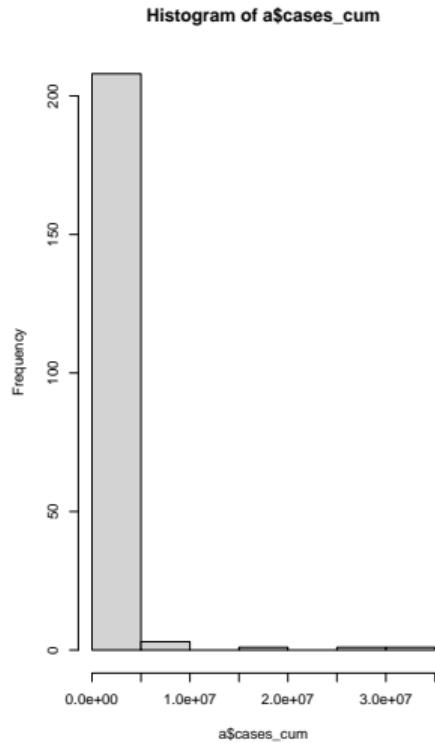
Base Graphics



Base Graphics: Histograms

```
1 # How about a histogram
2 hist(a$cases_cum)
3
4 # Let's add a few more breaks
5 hist(a$cases_cum, breaks = 50)
6
7 # I can't read the exponential notation, so lets try log:
8 a$cases_cum_log <- log(a$cases_cum)
9 hist(a$cases_cum_log)
```

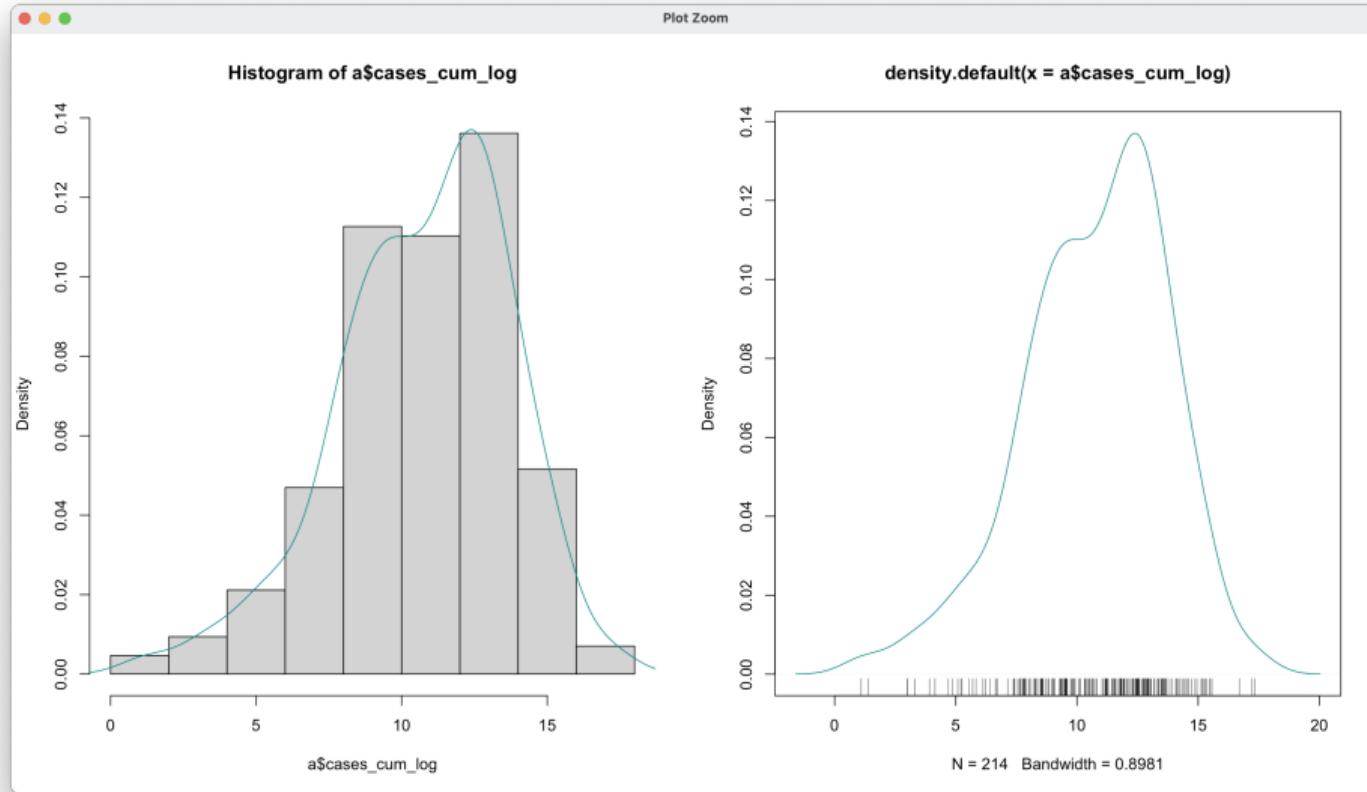
Base Graphics: Histogram



Base Graphics: Histograms, Density, Rugs

```
1 # Let's add the density now
2 d_ccl <- density(a$cases_cum_log)
3
4 # Plot matrix using mfrow / mfcol (rows, columns)
5 par(mfrow = c(1, 2))
6
7 ## Plot 1
8 hist(a$cases_cum_log, freq = FALSE)
9 lines(d_ccl, col = "darkcyan")
10
11 ## Plot 2
12 plot(d_ccl, col = "darkcyan")
13 rug(a$cases_cum_log)
14
15 ## Clear
16 dev.off()
```

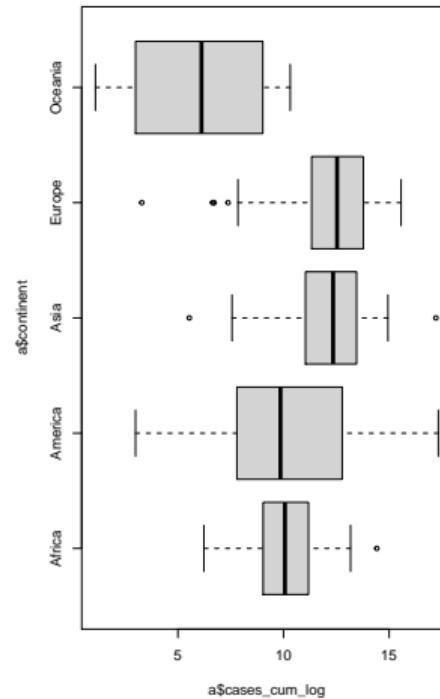
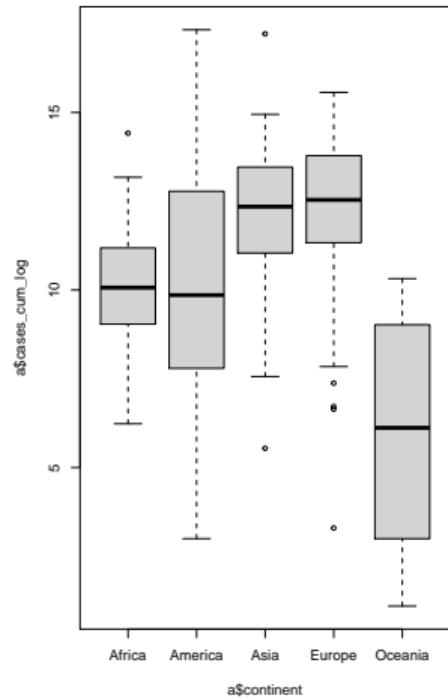
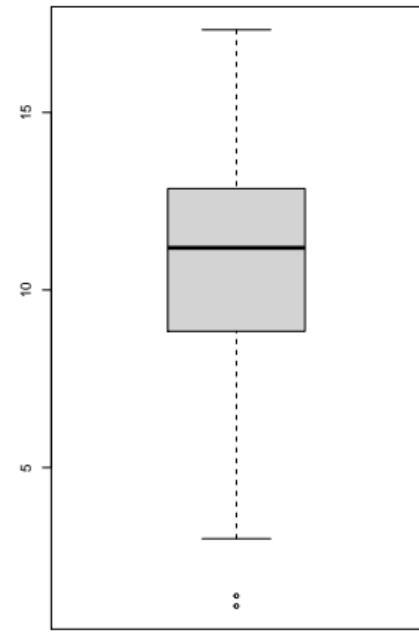
Base Graphics: Histograms, Density, Rugs



Base Graphics: Box Plots

```
1 boxplot(a$cases_cum_log)
2
3 boxplot(a$cases_cum_log ~ a$continent)
4
5 boxplot(a$cases_cum_log ~ a$continent, horizontal = TRUE)
```

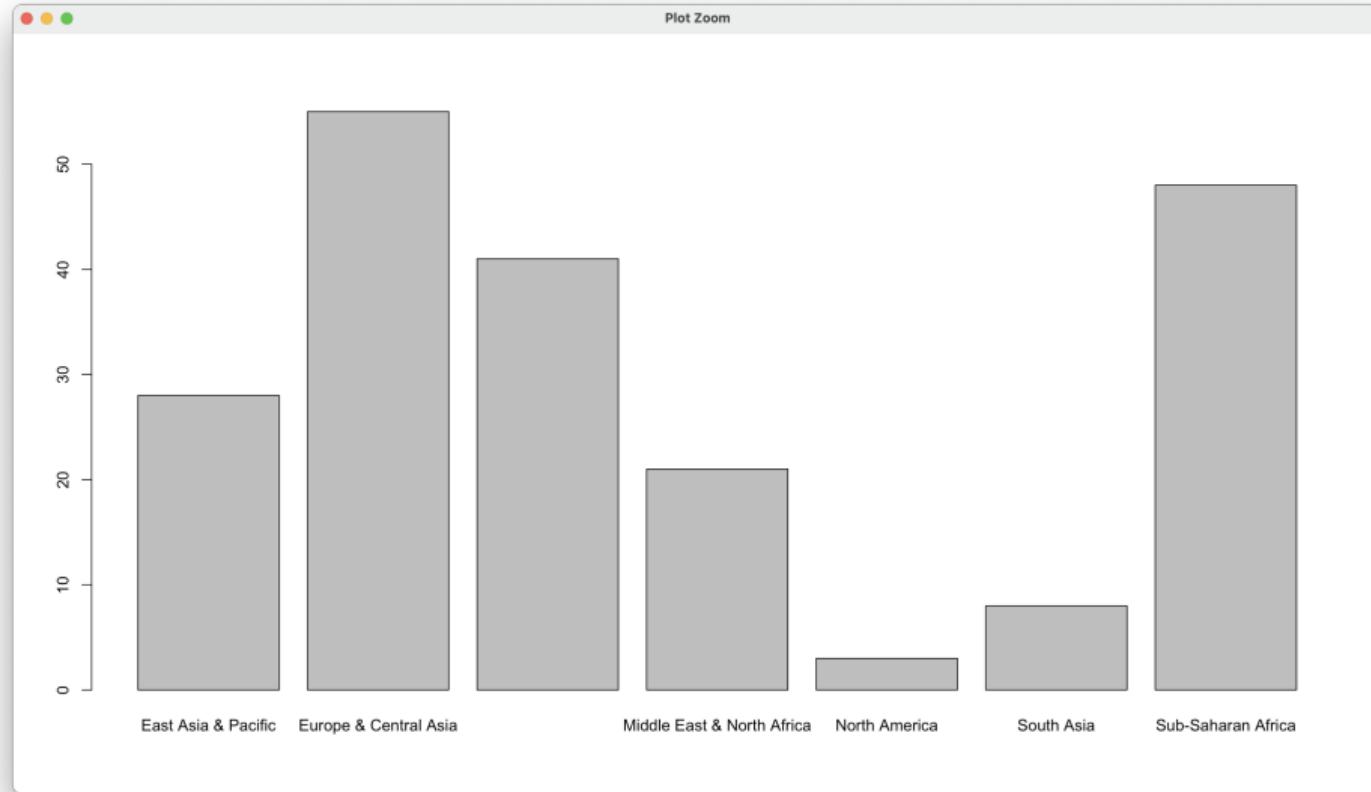
Base Graphics: Box Plots



Base Graphics: Bar Plots

```
1 # How many countries are there in each region?  
2 b <- table(a$region)  
3  
4 barplot(b[c(-1)])
```

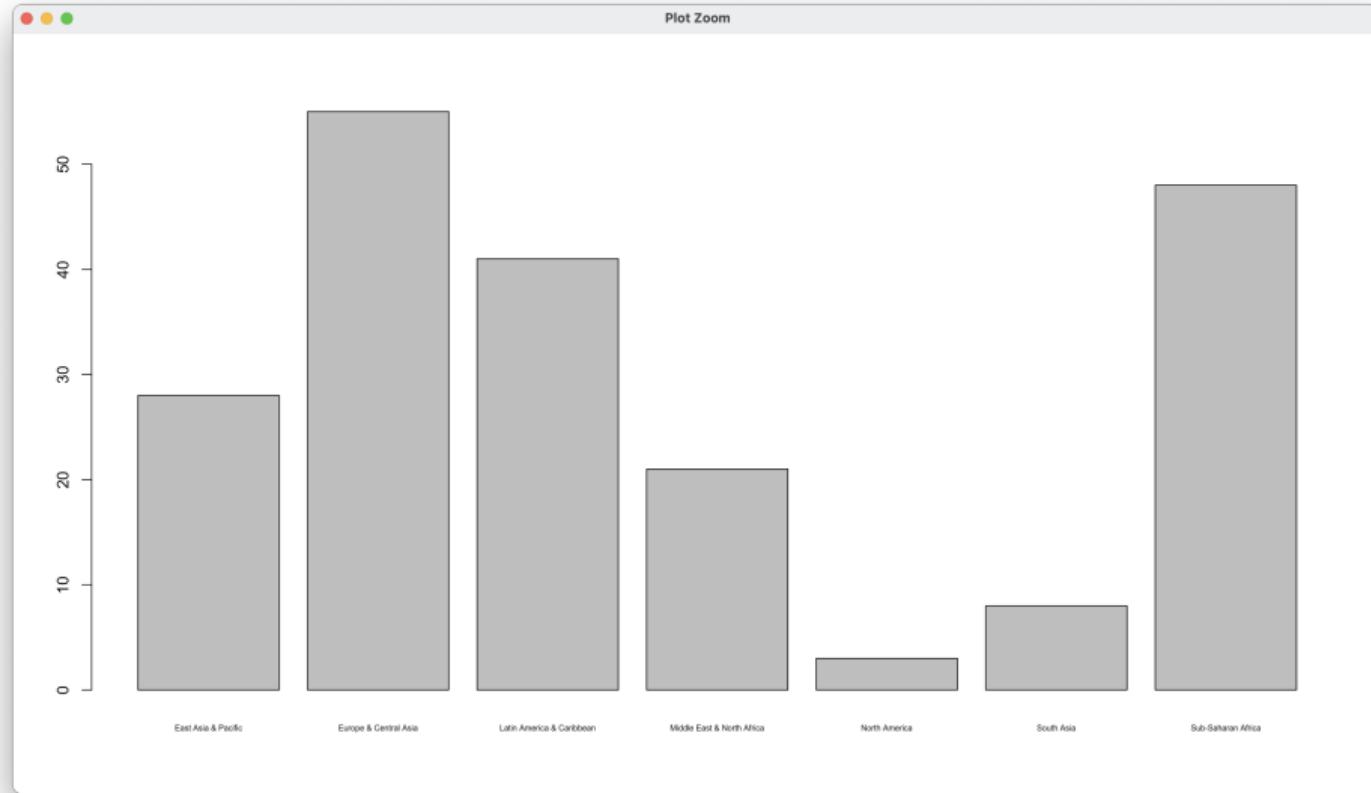
Base Graphics: Bar Plots



Base Graphics: Bar Plots

```
1 # wtf... R isn't plotting all the names? Because they don't fit...
2 barplot(b[c(-1)], cex.names = 0.5)
```

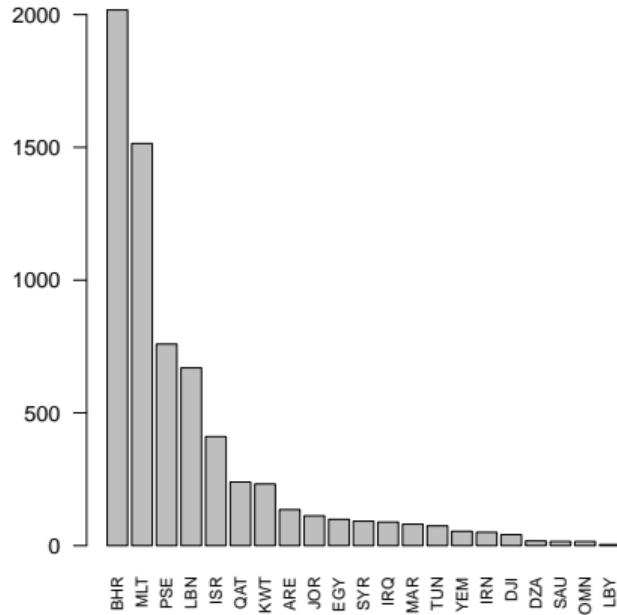
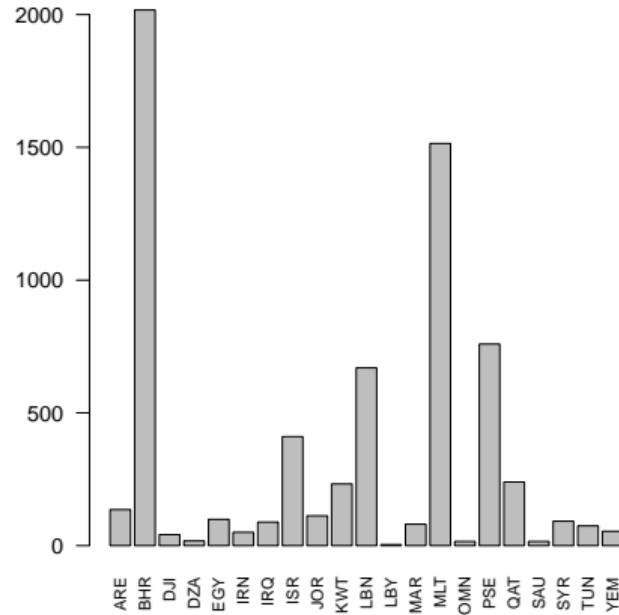
Base Graphics: Bar Plots



Base Graphics: Bar Plots

```
1 # Try something else: Filter down to just MENA countries
2 c <- a %>% filter(region == "Middle East & North Africa")
3
4 # New bar plot
5 barplot(c$pop_density, names.arg = c$geoid2, las = 2, cex.names = 0.75)
6
7 # Can we sort these bars?
8 d <- c[order(c$pop_density, decreasing = TRUE),]
9 barplot(d$pop_density, names.arg = d$geoid2, las = 2, cex.names = 0.75)
```

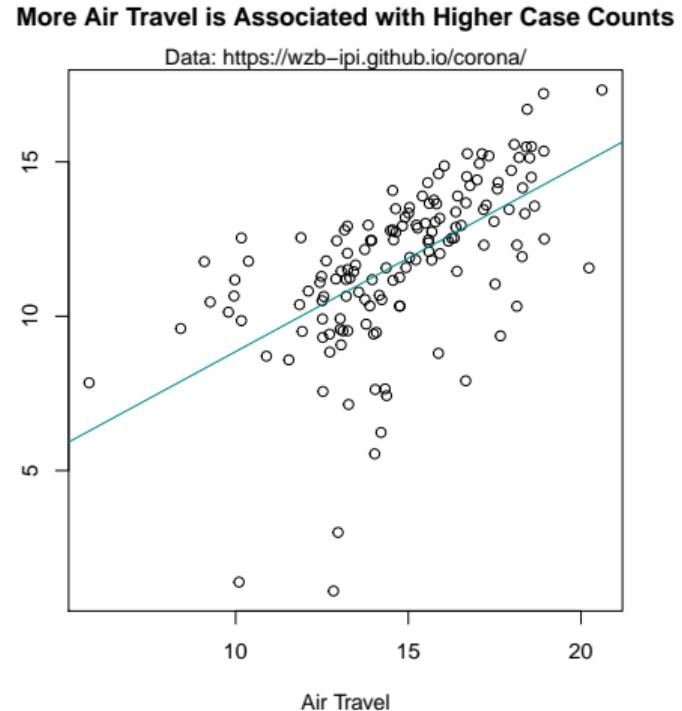
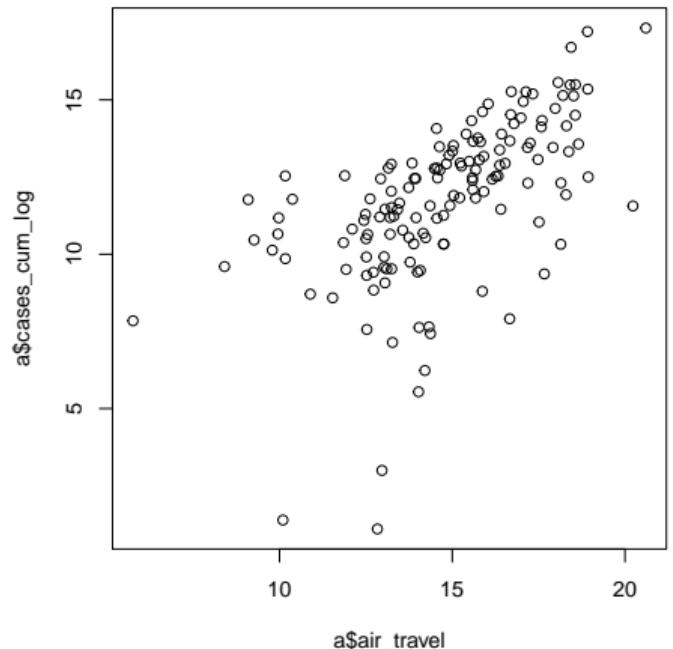
Base Graphics: Bar Plots



Base Graphics: Scatter Plots

```
1 plot(a$air_travel, a$cases_cum_log)
2
3 # Wow, that was easy. But let's make it nicer:
4 plot(a$air_travel, a$cases_cum_log,
5       xlab = "Air Travel",
6       ylab = "Cummulative Covid Cases (Log)")
7 abline(lm(a$cases_cum_log ~ a$air_travel), col = "darkcyan")
8 mtext("Data: https://wzb-ipi.github.io/corona/", side = 3)
9 title(main = "More Air Travel is Associated with Higher Case Counts")
```

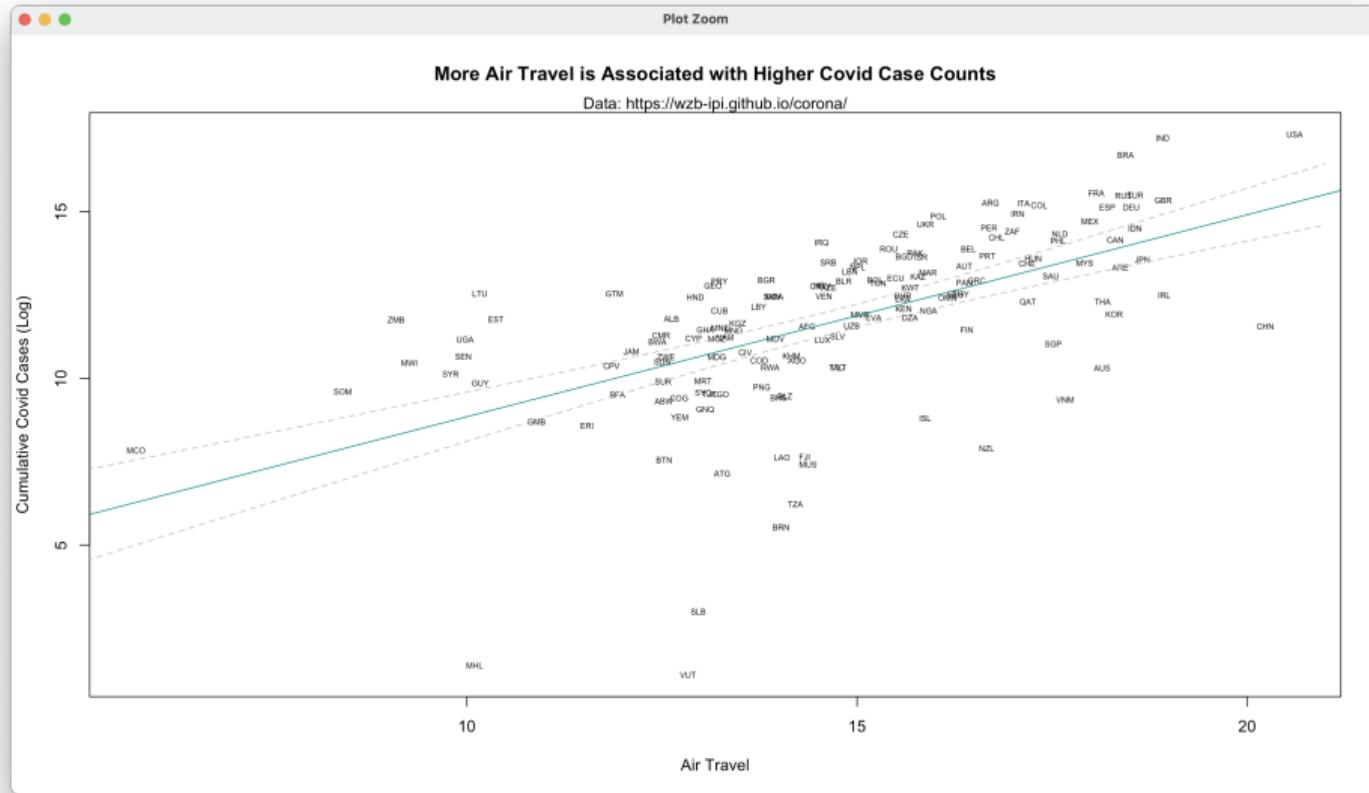
Base Graphics: Scatter Plots



Base Graphics: Scatter Plots

```
1 # But what about a confidence interval? We need to run the model and predict:
2 model <- lm(cases_cum_log ~ air_travel, data = a)
3 xvalues <- data.frame(air_travel = seq(5 , 21, length.out = 100))
4 predictions <- predict(model, newdata = xvalues, interval = "confidence")
5
6 plot(a$air_travel, a$cases_cum_log,
7       col = "white",
8       xlab = "Air Travel",
9       ylab = "Cumulative Covid Cases (Log)")
10 text(a$air_travel, a$cases_cum_log, labels = a$geoid2, cex = 0.5)
11 abline(lm(a$cases_cum_log ~ a$air_travel), col = "darkcyan")
12 lines(xvalues[,1], predictions[,2], col = "gray", lty = 2)
13 lines(xvalues[,1], predictions[,3], col = "gray", lty = 2)
14 mtext("Data: <a href='https://wzb-ipi.github.io/corona/">https://wzb-ipi.github.io/corona/", side = 3)
15 title(main = "More Air Travel is Associated with Higher Covid Case Counts")
```

Base Graphics: Scatter Plots



The Grammar of Graphics

ggplot2

General Format: ggplot (df, aes (x, y)) + geom_point ()

```
1 library(ggplot2)
2
3 plot <- ggplot(df, aes(x = x, y = y)) +
4   geom_point() +
5   geom_smooth(method = "lm") +
6   coord_cartesian(xlim = c(0, 100), ylim = c(0, 100)) +
7   labs(title = "Title",
8       subtitle = "Subtitle",
9       y = "Y Title",
10      x = "X Title",
11      caption = "Caption") +
12   theme_minimal()
13 plot
14 ggsave("plot.pdf", plot = plot, width = 11, height = 8.5)
```

ggplot2

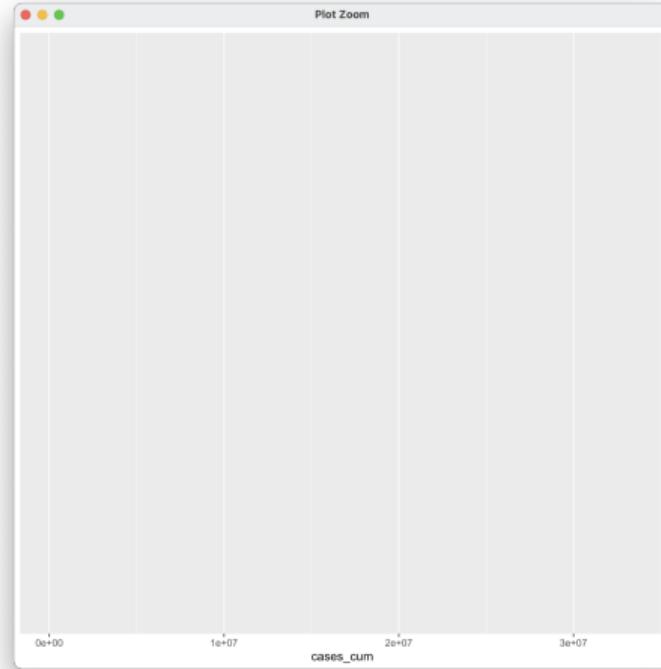
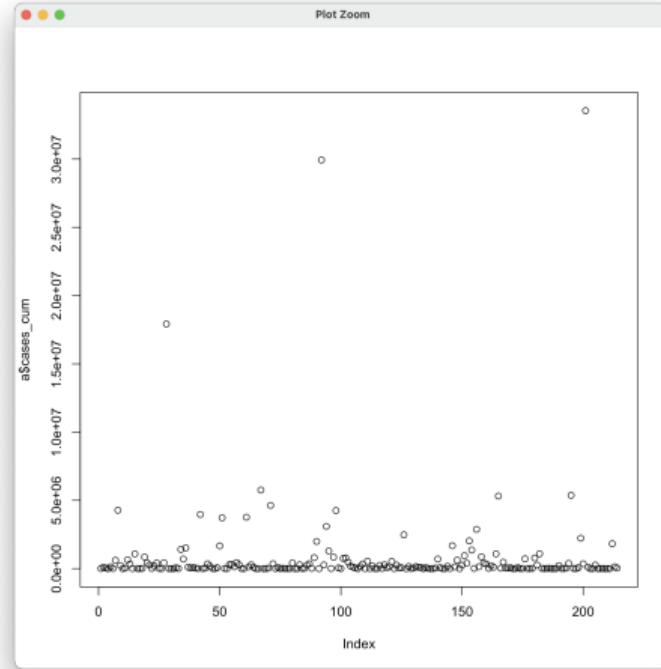
- ▶ Remember in base graphics we could generate a plot with:

```
plot(x = a$cases_cum)
```

- ▶ What happens if we just type:

```
ggplot(a, aes(x = cases_cum))
```

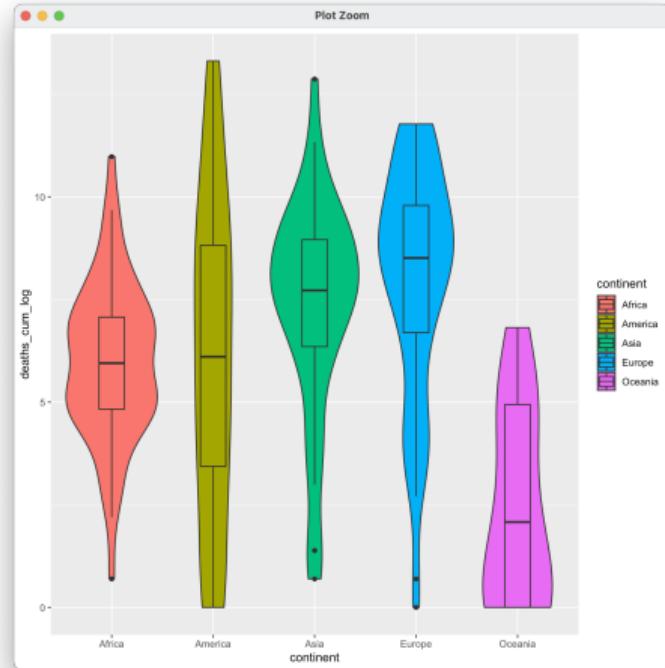
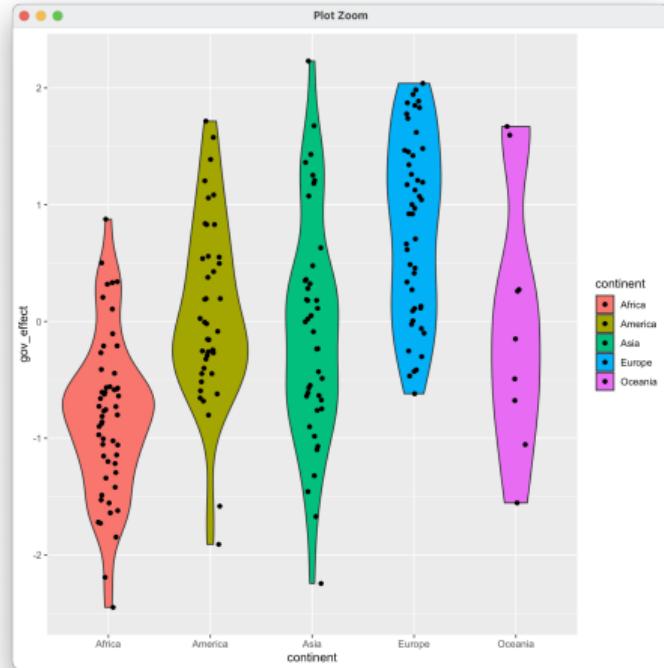
ggplot2



ggplot2: Violin Plots, Box Plots

```
1 gg1 <- ggplot(a, aes(x = continent, y = gov_effect, fill = continent)) +
2   geom_violin() +
3   geom_jitter(position=position_jitter(0.1))
4 gg1
5
6 gg2 <- ggplot(a, aes(x = continent, y = deaths_cum_log, fill = continent)) +
7   geom_violin() +
8   geom_boxplot(width = 0.25)
9 gg2
```

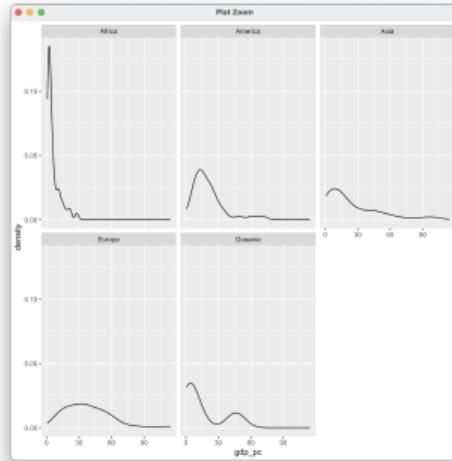
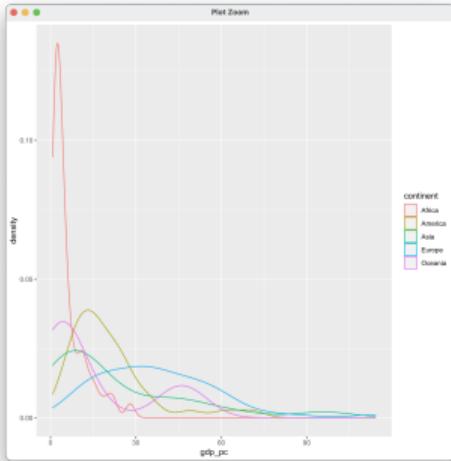
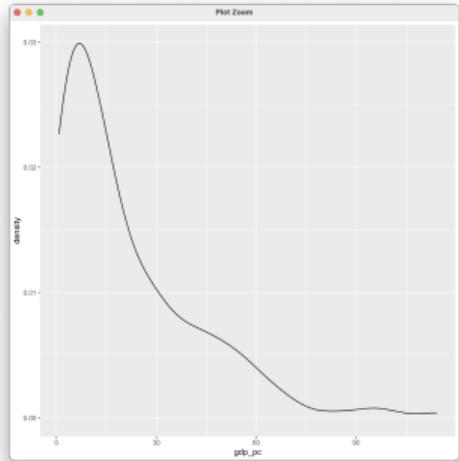
ggplot2: Violin Plots, Box Plots



ggplot2: Density Plots

```
1 gg3 <- ggplot(a, aes(x = gdp_pc)) +  
2   geom_density()  
3 gg3  
4  
5 gg4 <- ggplot(a, aes(x = gdp_pc)) +  
6   geom_density(aes(color = continent))  
7 gg4  
8  
9 gg5 <- ggplot(a, aes(x = gdp_pc)) +  
10  geom_density() +  
11  facet_wrap(~ continent)  
12 gg5
```

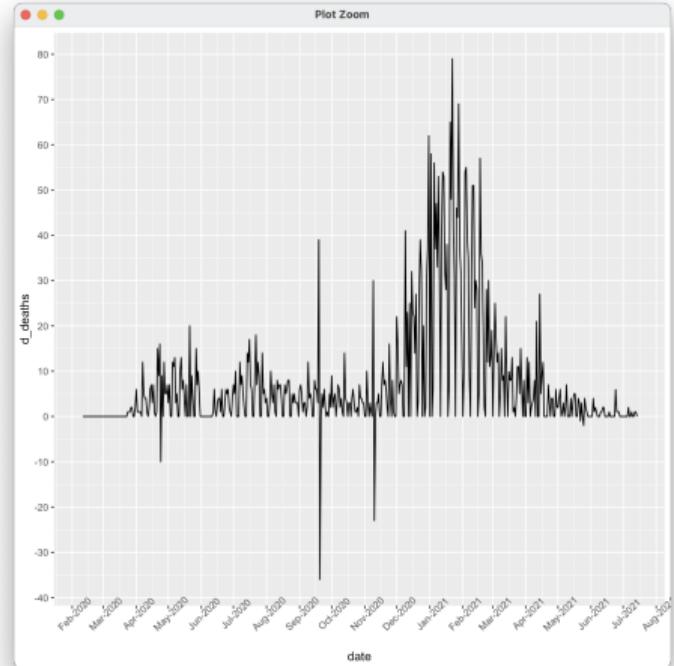
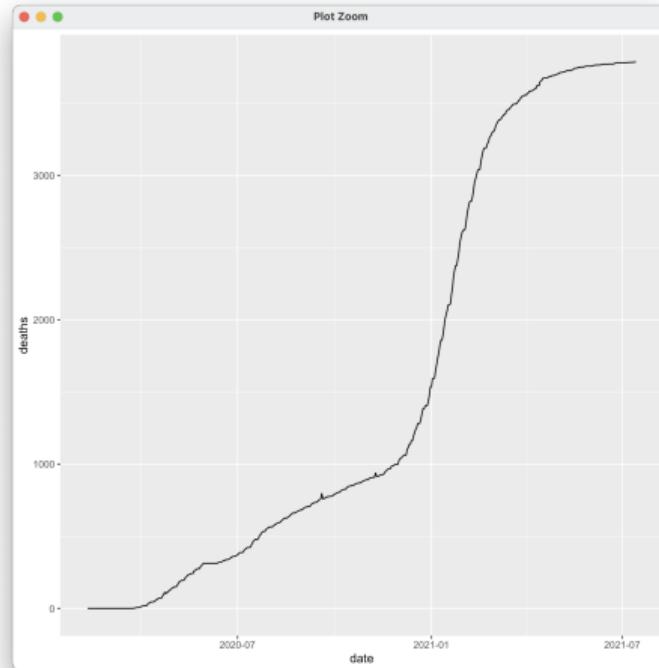
ggplot2:



ggplot2: Line / Time Series

```
1 # Some other Covid data
2 c <- covid19(c("US"), level = 3)
3 d <- c %>% filter(administrative_area_level_3 == "San Diego")
4
5 gg6 <- ggplot(d, aes(x = date, y = confirmed)) +
6   geom_line()
7 gg6
8
9 gg7 <- ggplot(d) +
10   geom_line(aes(x = date, y = deaths))
11 gg7
12
13 # What if we wanted daily deaths?
14 d <- d %>% mutate(d_deaths = deaths - lag(deaths))
15
16 gg8 <- ggplot(d) +
17   geom_line(aes(x = date, y = d_deaths)) +
18   scale_x_date(date_labels = "%b-%Y", breaks = breaks_pretty(20)) +
19   scale_y_continuous(breaks = breaks_pretty(10)) +
20   theme(axis.text.x = element_text(angle = 45))
21 gg8
```

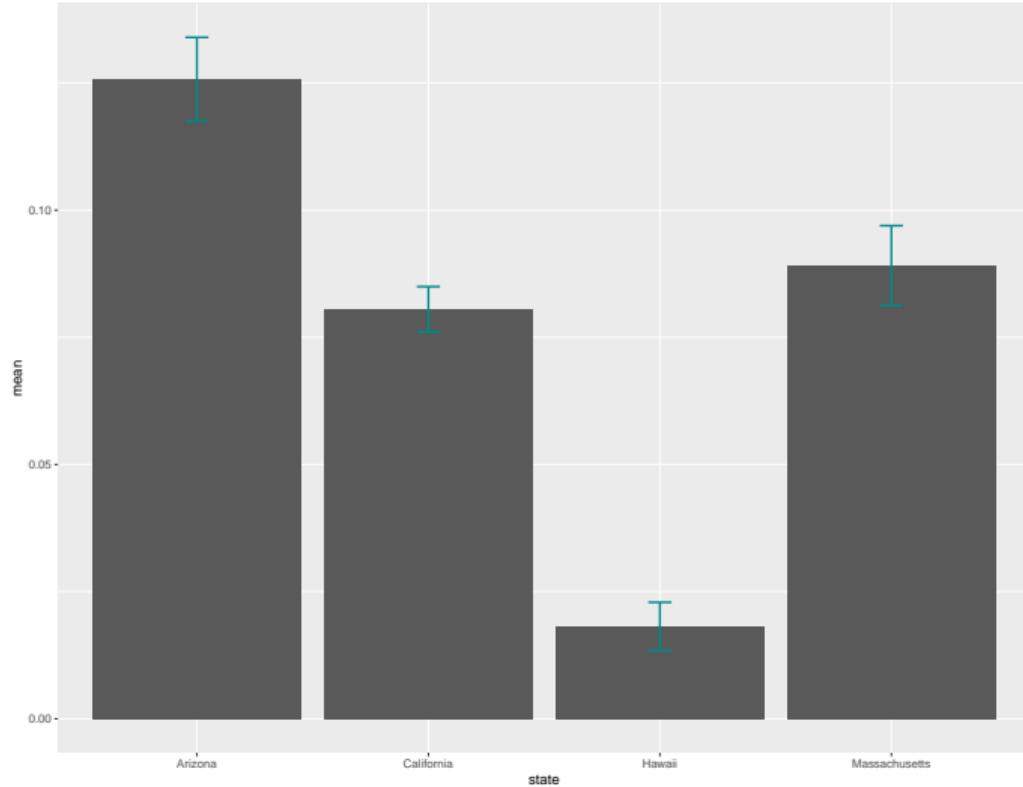
ggplot2: Line / Time Series



ggplot2: Bar Plots + CI

```
1 # Let's look at 1 July in the 4 states that I've lived in for more than a year:
2 e <- c %>% filter(administrative_area_level_2 == "Hawaii" |
3                     administrative_area_level_2 == "California" |
4                     administrative_area_level_2 == "Massachusetts" |
5                     administrative_area_level_2 == "Arizona")
6 e <- e %>% filter(date == "2021-07-01")
7
8 # Let's normalize the case counts by the population:
9 h <- e %>% mutate(mconfirmed = confirmed / population) %>%
10    data.frame() %>%
11    select(state = administrative_area_level_2, mconfirmed) %>%
12    group_by(state) %>%
13    summarize(n = n(),
14              mean = mean(mconfirmed),
15              sd = sd(mconfirmed)) %>%
16    mutate(se = sd/sqrt(n))
17
18 # Now we can plot:
19 gg11 <- ggplot(h) +
20   geom_bar(aes(x = state, y = mean), stat = "identity") +
21   geom_errorbar(aes(x = state, ymin = mean-se, ymax = mean+se),
22                 width=0.1, colour="darkcyan", size=0.75)
23 gg11
```

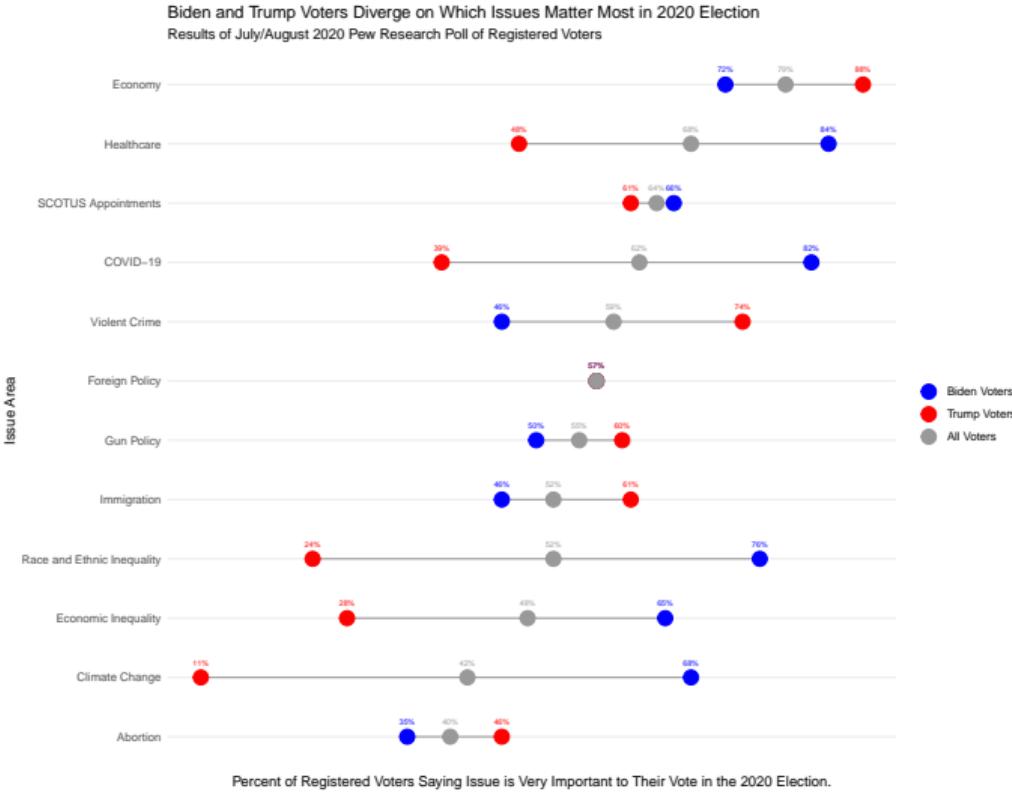
ggplot2: Bar Plots + CI



ggplot2: Dot Plots

```
1 # Using some Pew data:  
2 b <- read.csv("pew_data.csv")  
3  
4 # We have a problem though! Our data is wide, and we need it to be long...  
5 c <- melt(b, id = "issues")  
6  
7 gg14 <- ggplot(c, aes(x = value, y = reorder(issues, value))) +  
8   geom_line(aes(group = issues), color = "#999999") +  
9   geom_point(aes(color = variable), size = 5) +  
10  scale_color_manual(values = c("#0000FF", "#FF0000", "#999999"),  
11    labels = c("Biden Voters", "Trump Voters", "All Voters")) +  
12  geom_text(aes(label = sprintf("%1.0f%%", value), color = variable),  
13    size = 2, nudge_y = 0.25) +  
14  scale_x_continuous(breaks = NULL) +  
15  labs(title = "Biden and Trump Voters Diverge on Which Issues Matter Most in 2020 Election",  
16    subtitle = "Results of July/August 2020 Pew Research Poll of Registered Voters",  
17    x = "Percent of Registered Voters Saying Issue is Very Important to Their Vote in the 2020 Election.",  
18    y = "Issue Area") +  
19  theme_minimal() +  
20  theme(legend.title = element_blank())  
21 gg14
```

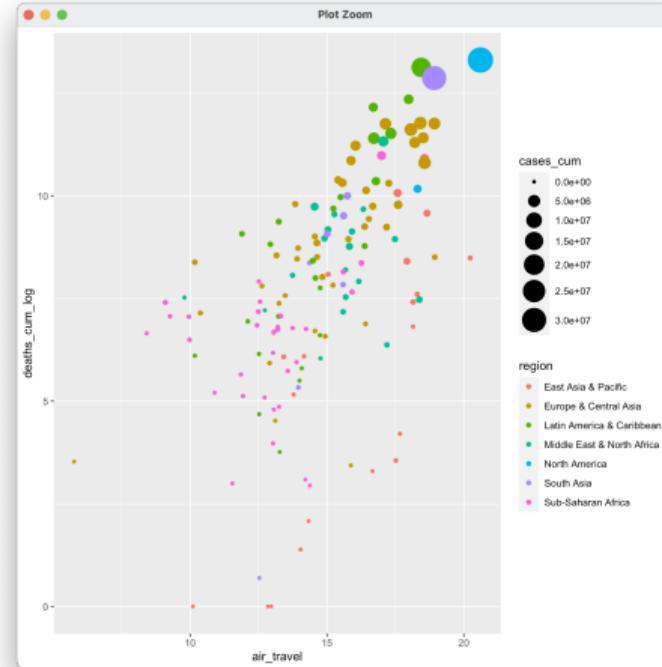
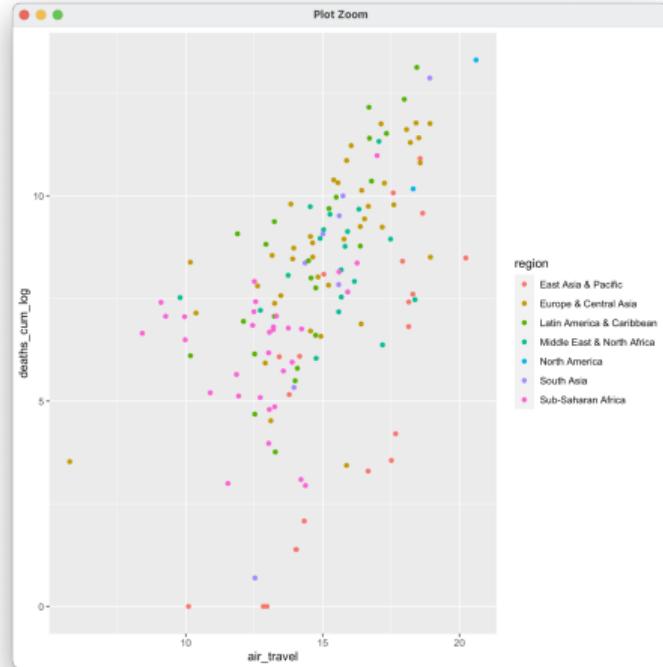
ggplot2: Dot Plots



ggplot2: Scatter Plots

```
1 b <- a %>% filter(region != "")  
2  
3 gg15 <- ggplot(b, aes(x = air_travel, y = deaths_cum_log, color = region)) +  
4   geom_point()  
5 gg15  
6  
7 # Let's try to scale by case count  
8 gg16 <- ggplot(b, aes(x = air_travel, y = deaths_cum_log,  
9                     size = cases_cum,  
10                    color = region)) +  
11   scale_size_continuous(range = c(1, 10), breaks = pretty_breaks()) +  
12   geom_point()  
13 gg16
```

ggplot2: Scatter Plots

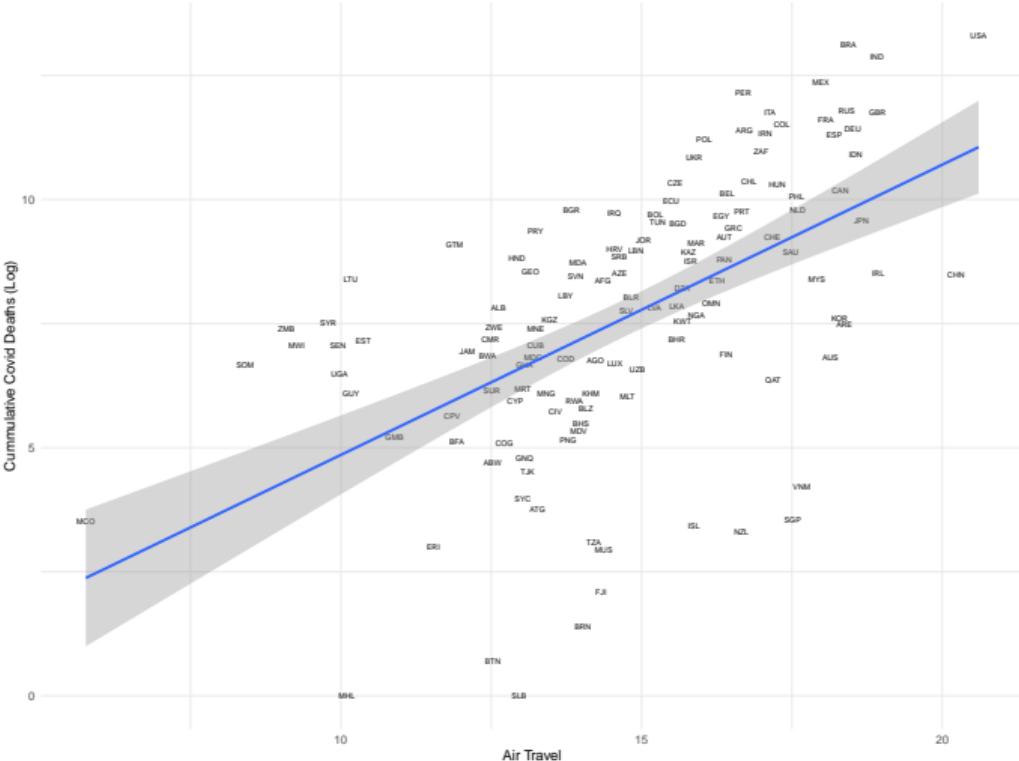


ggplot2: Scatter Plots

```
1 # Or better yet, can you make it look like the plot we made with base graphics?
2 gg20 <- ggplot(b, aes(x = air_travel, y = deaths_cum_log)) +
3   geom_point() +
4   geom_smooth(method = lm) +
5   labs(title = "More Air Travel is Associated with Higher Covid Case Counts",
6        x = "Air Travel",
7        y = "Cummulative Covid Deaths (Log)") +
8   theme_minimal()
9 gg20
10
11 # Almost there, but I want ISO codes instead of points
12 gg20 <- ggplot(b, aes(x = air_travel, y = deaths_cum_log, label = geoid2)) +
13   geom_text(size = 2, check_overlap = TRUE) +
14   geom_smooth(method = lm,) +
15   labs(title = "More Air Travel is Associated with Higher Covid Case Counts",
16        x = "Air Travel",
17        y = "Cummulative Covid Deaths (Log)") +
18   theme_minimal()
19 gg20
```

ggplot2: Scatter Plots

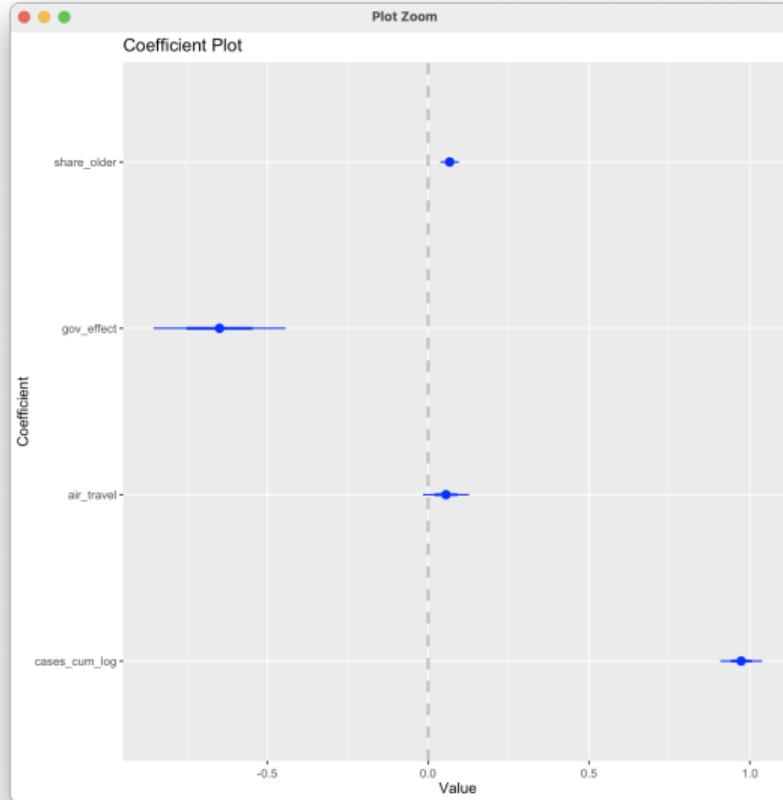
More Air Travel is Associated with Higher Covid Case Counts



ggplot2: Coefficient Plots

```
1 m <- lm(deaths_cum_log ~ cases_cum_log + air_travel + gov_effect + share_older,  
2           data = b)  
3 m  
4  
5 coefplot(m, intercept = FALSE)  
6 # Wow, that was easy
```

ggplot2: Coefficient Plots



ggplot2: Coefficient Plots

Wrangle the data to make a coefficient plot by hand:

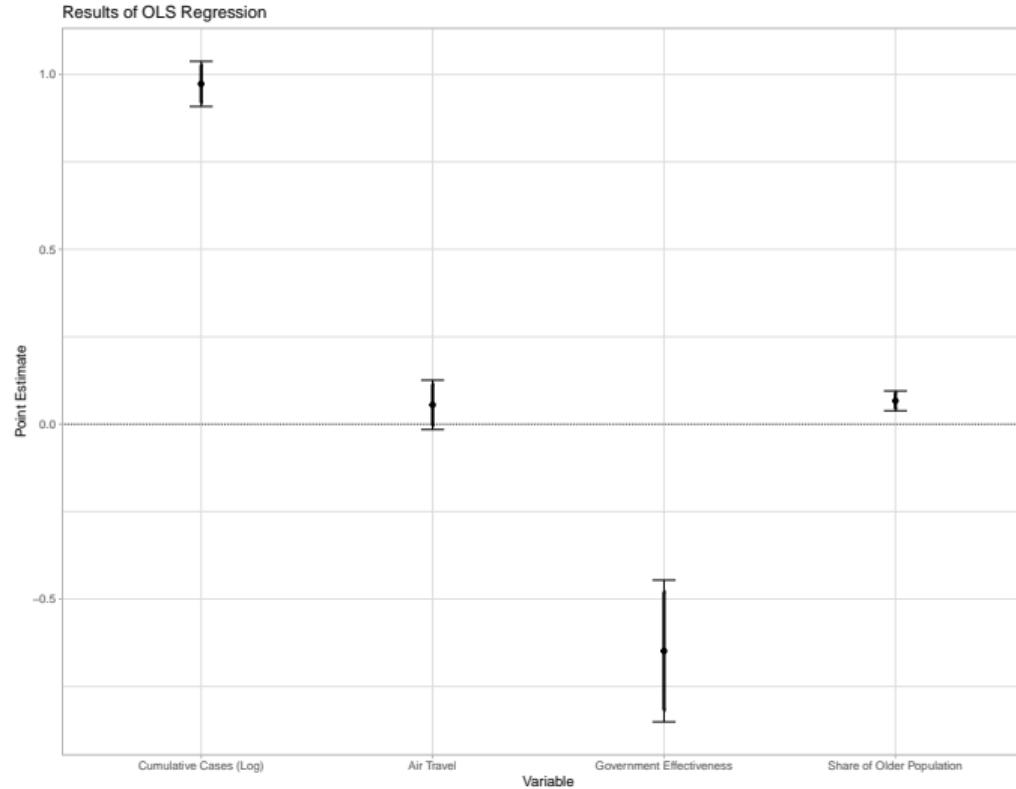
```
1 # Can we make it harder? Let's try to replicate this manually in ggplot
2 z1 <- tidy(m, conf.int = TRUE, conf.level = .95)
3 z1 <- z1 %>% dplyr::rename(low95 = conf.low, high95 = conf.high)
4
5 z2 <- tidy(m, conf.int = TRUE, conf.level = .9)
6 z2 <- z2 %>% dplyr::rename(low90 = conf.low, high90 = conf.high)
7
8 z3 <- z1 %>% full_join(z2)
9 z3$term <- dplyr::recode(z3$term,
10                         "(Intercept)" = "Intercept",
11                         "cases_cum_log" = "Cumulative Cases (Log)",
12                         "air_travel" = "Air Travel",
13                         "gov_effect" = "Government Effectiveness",
14                         "share_older" = "Share of Older Population")
15
16 z3$term <- factor(z3$term,
17                     levels = c("Intercept",
18                               "Cumulative Cases (Log)",
19                               "Air Travel",
20                               "Government Effectiveness",
21                               "Share of Older Population"),
22                     ordered = TRUE)
```

ggplot2: Coefficient Plots

Now we can actually plot it:

```
1 gg21 <- ggplot(z3 %>% filter(term != "Intercept")) +
2   geom_hline(yintercept=0, lty="11", colour="grey30") +
3   geom_errorbar(aes(term, ymin = low95, ymax = high95), width=0.1) +
4   geom_errorbar(aes(term, ymin = low90, ymax = high90), lwd = 1.15, width=0) +
5   geom_point(aes(term, estimate)) +
6   labs(title = "Results of OLS Regression",
7        y = "Point Estimate",
8        x = "Variable") +
9   theme_light()
10 gg21
```

ggplot2: Coefficient Plots



Final Tips

- ▶ Sketch out your plot on paper before you begin
- ▶ Build plots iteratively, one element at a time
 - ▶ ggplot is particularly conducive to this method
 - ▶ Allows you to see what works and what doesn't
- ▶ Try to write generic and reusable code
- ▶ Remember that coding is nothing more than learning what to ask Google so that it gives you a useful response
 - ▶ You are almost certainly not the first person to ever encounter a particular coding issue
 - ▶ Don't forget about YouTube; if it's not on YouTube it doesn't exist
- ▶ Remember that R isn't the only statistical programming package out there

Feel Free to Reach Out!

Mike Seese

 mseese@ucsd.edu

 [@mikeseese](https://twitter.com/mikeseese)

 <https://github.com/mfseese>