

Categorical Variables (Chapter 4)

Prof. Joyce Robbins

Numeric data

```
## 'data.frame': 15 obs. of 20 variables:  
## $ a1 : num 18.6 37.6 71.6 94.2 100.2 ...  
## $ a2 : num 17 38.2 67.8 106.8 64.2 ...  
## $ a3 : num 19 36.2 90.4 110.9 83.4 ...  
## $ a4 : num 6 48.6 77 115.5 94.1 ...  
## $ a5 : num 15.8 43.6 81.6 133 87.6 ...  
## $ a6 : num 0 22.8 36.6 111.2 54.8 ...  
## $ a7 : num 6.2 31 62 101.5 66.8 ...  
## $ a8 : num 5 30.2 31.1 89.7 53.5 ...  
## $ a9 : num 7.2 27 65 124.1 104.9 ...  
## $ a10: num 0 25.8 60.8 69.5 81.9 ...  
## $ a11: num 8 19.4 60.2 102.7 56.5 ...  
## $ a12: num 15 38 71.4 106.9 67.4 ...  
## $ a13: num 2.8 35.8 66.6 121.5 67.7 ...  
## $ a14: num 4.4 35.4 48 120.7 41 ...  
## $ a15: num 6.6 34.8 52 100.6 78 ...  
## $ a16: num 4 28.6 34.1 101.5 40.1 ...  
## $ a17: num 2.4 41.2 30 116.4 11.2 ...  
## $ a18: num 9.6 24.4 54 103.9 67.4 ...  
## $ a19: num 0 33.8 47.6 111.7 79.7 ...  
## $ a20: num 2.2 31.2 57.6 127.7 65.5 ...
```

Categorical data

```
## Classes 'tbl_df', 'tbl' and 'data.frame': 1373 obs. of 12 variables:  
## $ respondent_id : num 3308895255 3308891308 3308891135 3308879091 3308871671 ...  
## $ knowledge     : Ord.factor w/ 4 levels "Novice"<"Intermediate"<...: 2 1 2 1 1 3 1 3 1 1 ...  
## $ interest      : Ord.factor w/ 4 levels "Not at all"<"Not much"<...: 3 3 4 2 2 4 3 4 2 3 ...  
## $ gender        : chr "Male" "Male" "Male" "Male" ...  
## $ age           : Factor w/ 4 levels "18-29","30-44",...: 1 1 2 3 2 2 3 3 2 NA ...  
## $ household_income: Factor w/ 5 levels "$0 - $24,999",...: 4 4 3 1 2 3 NA 1 3 NA ...  
## $ education     : Ord.factor w/ 5 levels "Less than high school degree"<...: 1 3 5 1 2 5 2 3 3  
NA ...  
## $ location       : chr "West South Central" "West South Central" "Pacific" "New England" ...  
## $ algeria         : chr "N/A" "N/A" "3" "N/A" ...  
## $ argentina      : chr "3" "N/A" "4" "3" ...  
## $ australia      : chr "5" "3" "N/A" "N/A" ...  
## $ belgium         : chr "4" "3" "3" "3" ...
```

Warnings

- words are hard to work with!
- not a lot of options (esp. for 1 dimension): bar plot, Cleveland dot plot
- data cleaning takes more time
- main choices: *which* categories to plot, *order* of categories

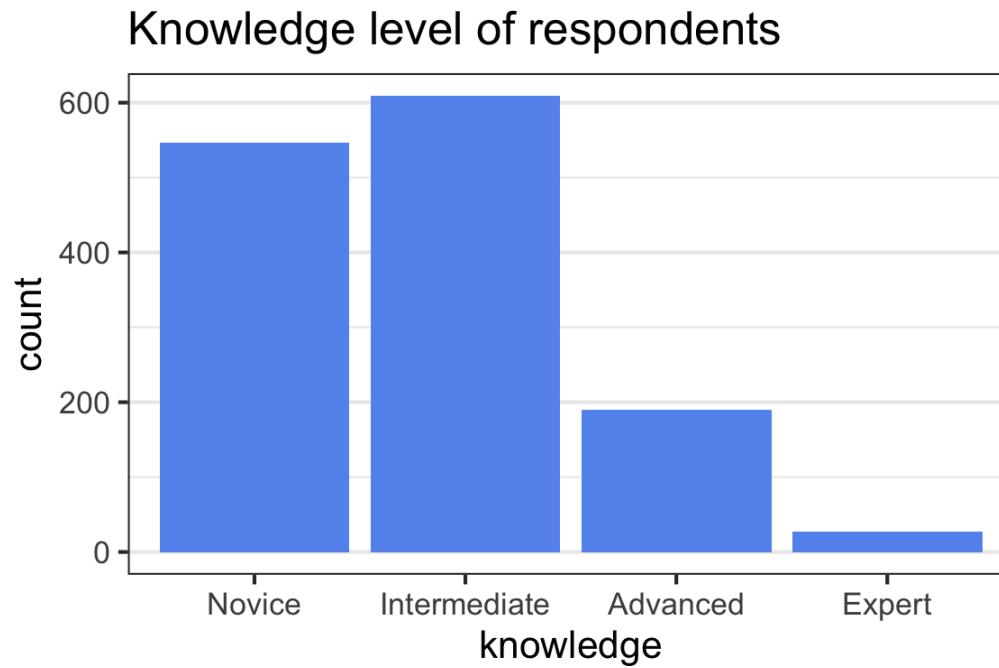
Types of data

- nominal – no fixed category order
- ordinal – fixed category order
- (“real”) discrete, small # of possibilities
- Not always clearcut: nominal vs. ordinal, ordinal vs. discrete, and...
- Sometimes numbers = nominal, not discrete

Ordinal data

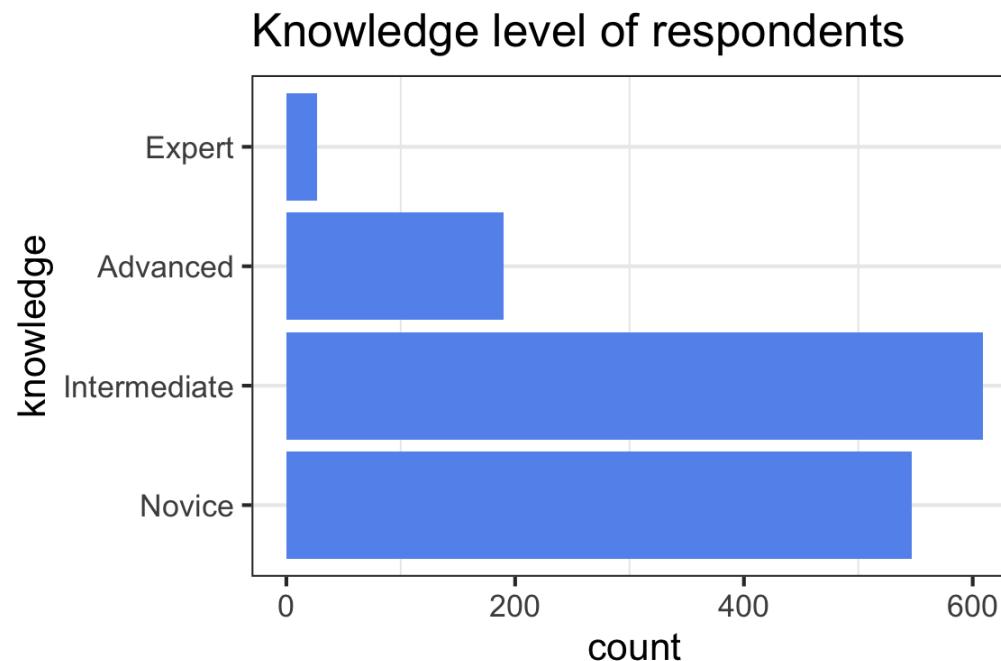
Sort in logical order of the categories (left to right)

ORDERING IS VERY IMPORTANT FOR ALL GRAPHS



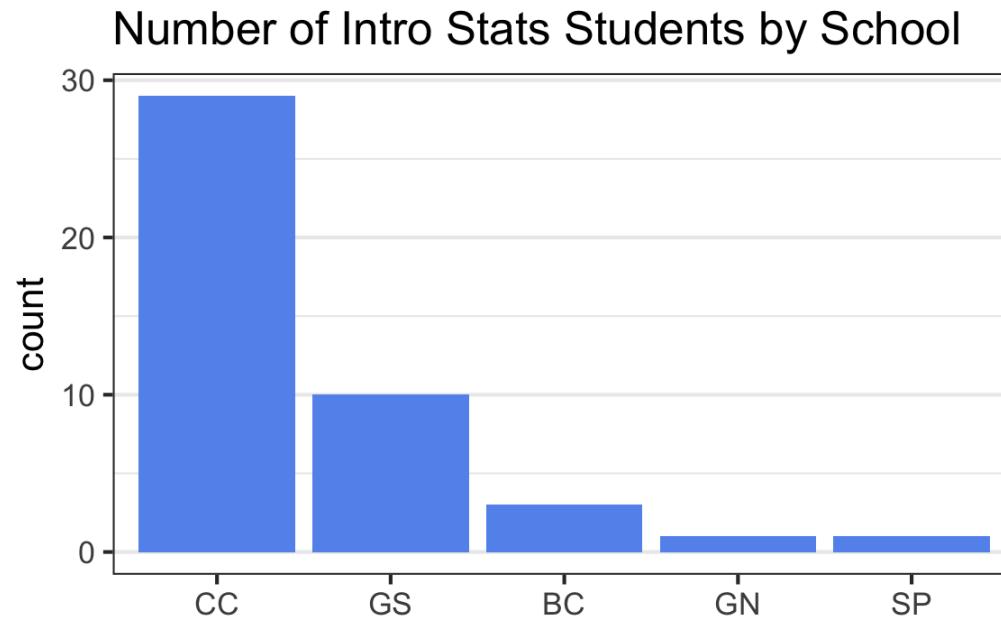
Ordinal data

Sort in logical order of the categories (starting at bottom OR top)



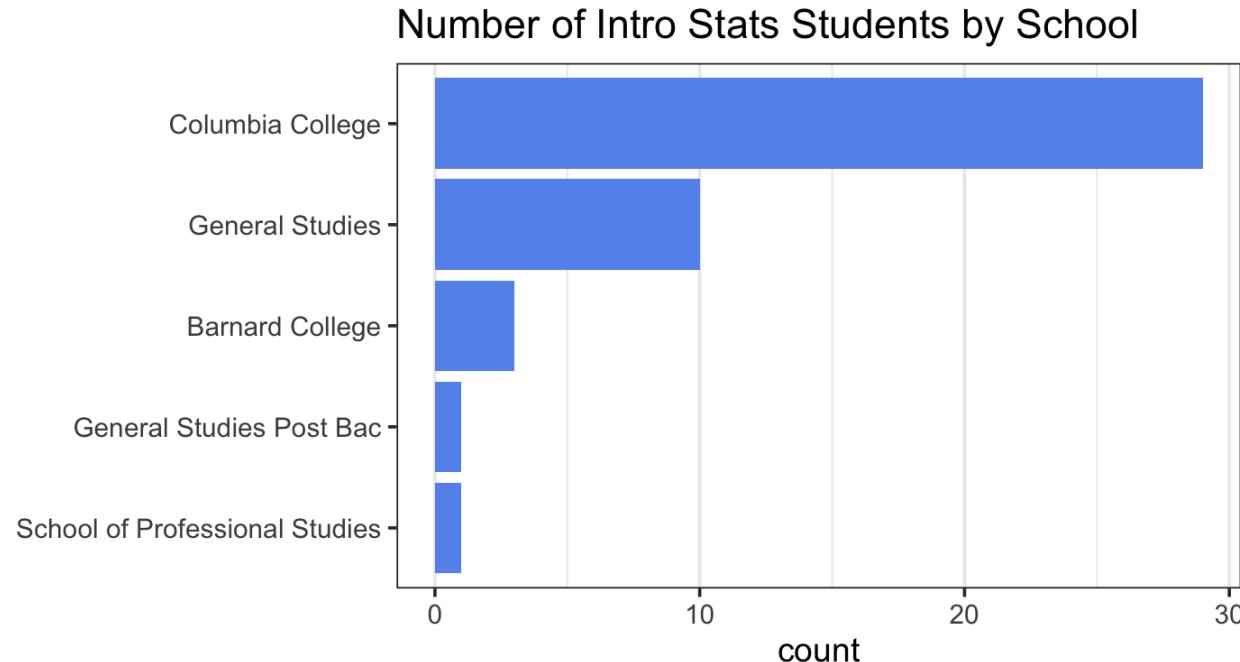
Nominal data

Sort from highest to lowest count (left to right, or top to bottom)



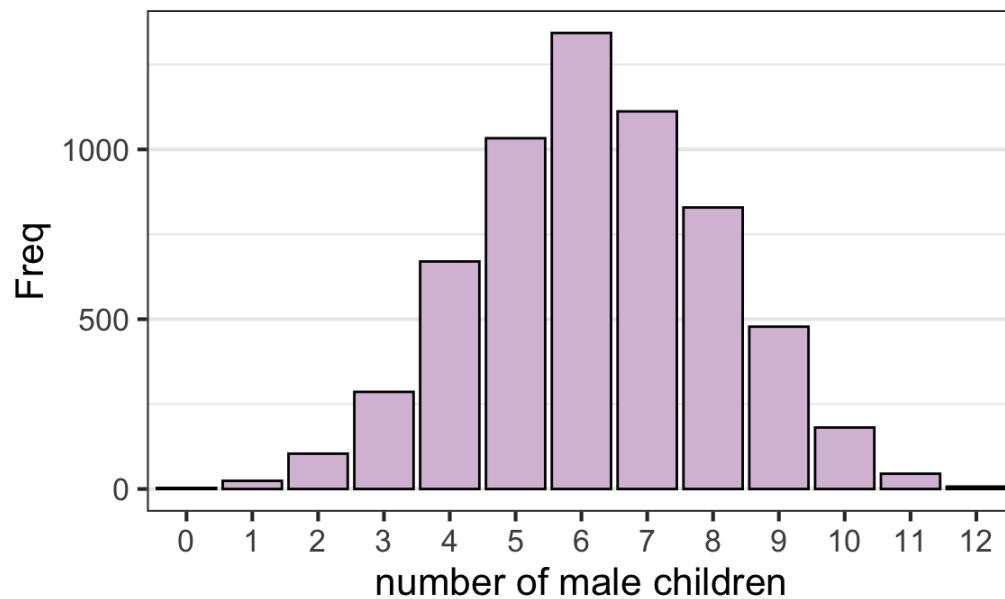
Nominal data

... or top to bottom

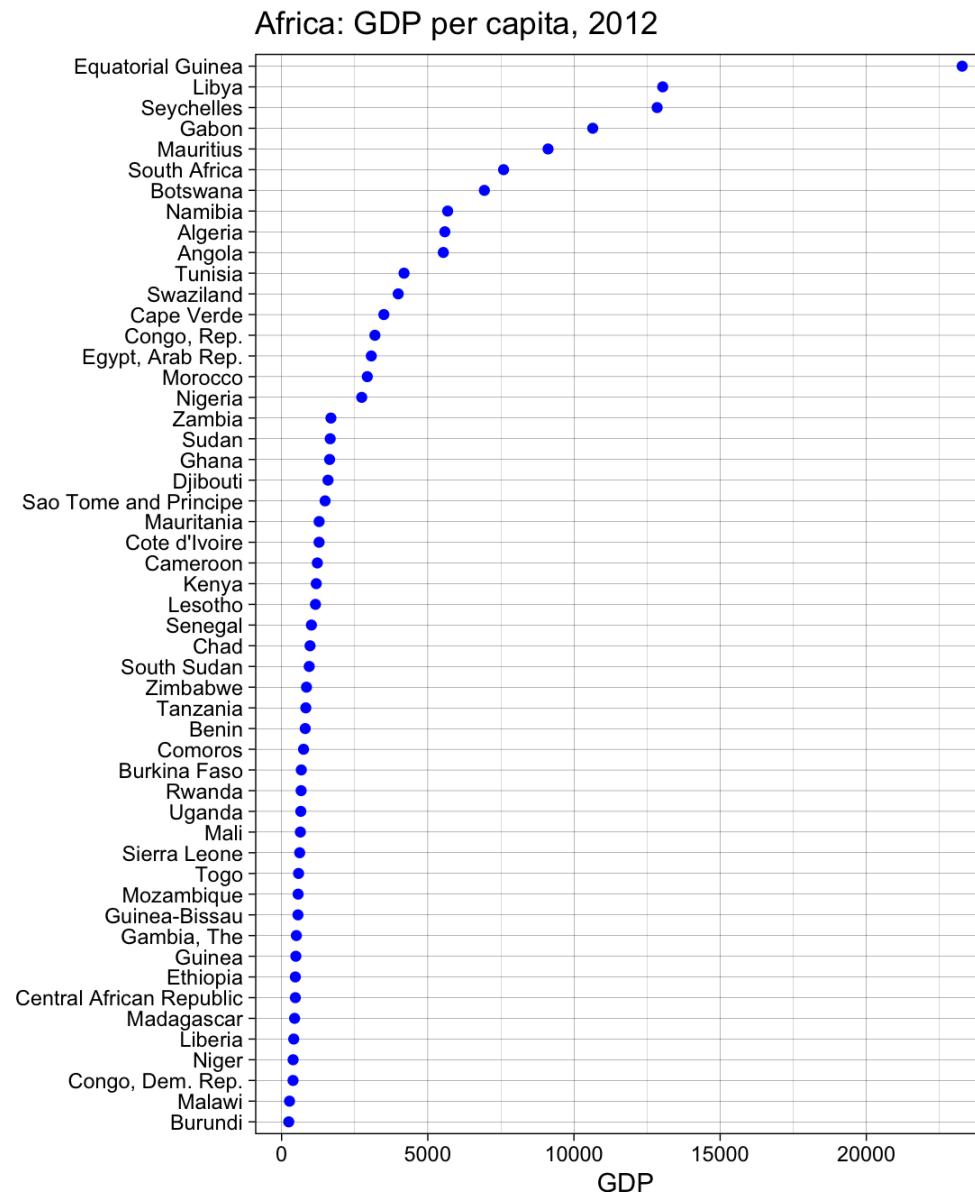


Discrete data

19c Saxony: # of males in families with 12 c

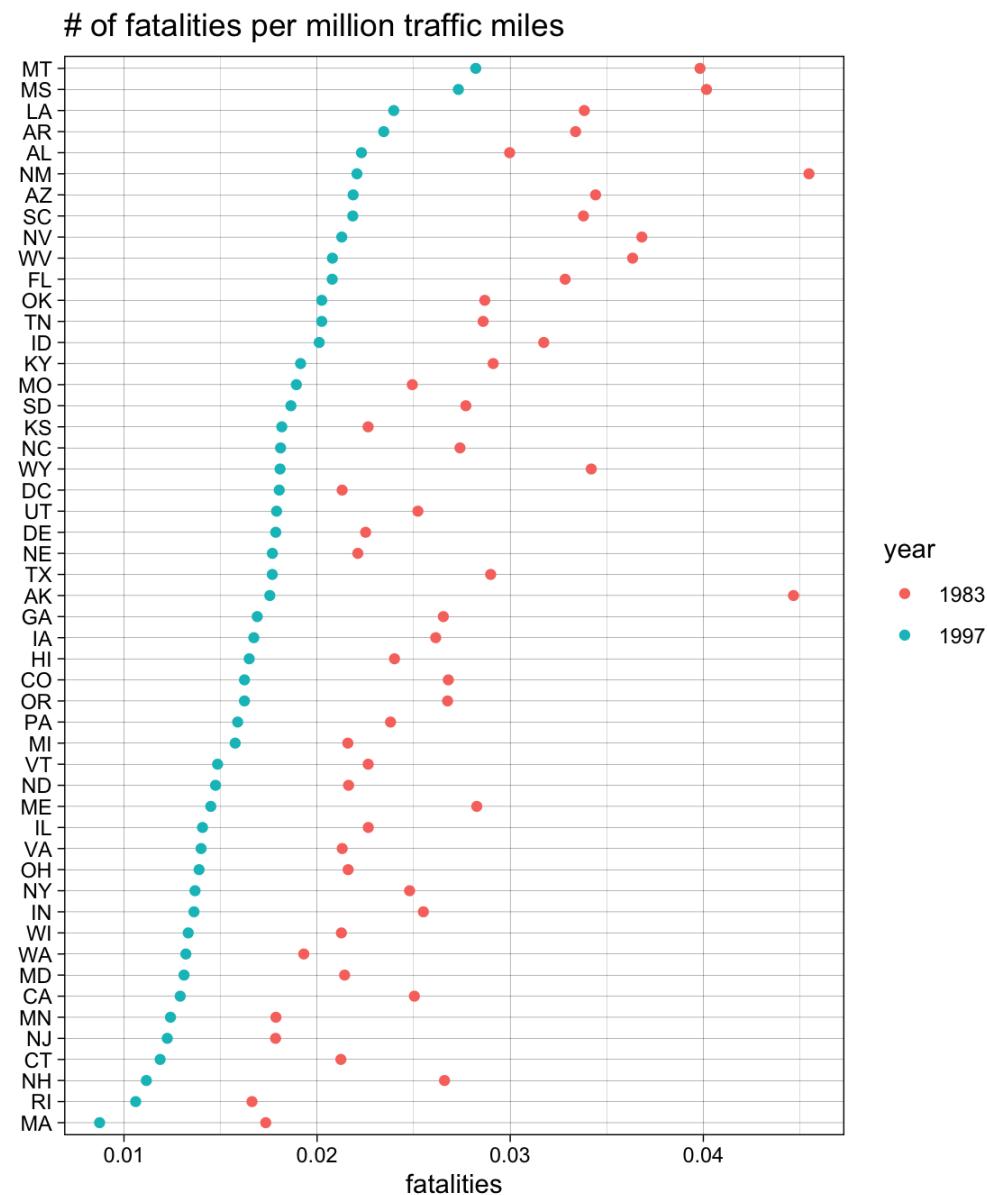


Cleveland dot plot



Cleveland dot plot with multiple dots

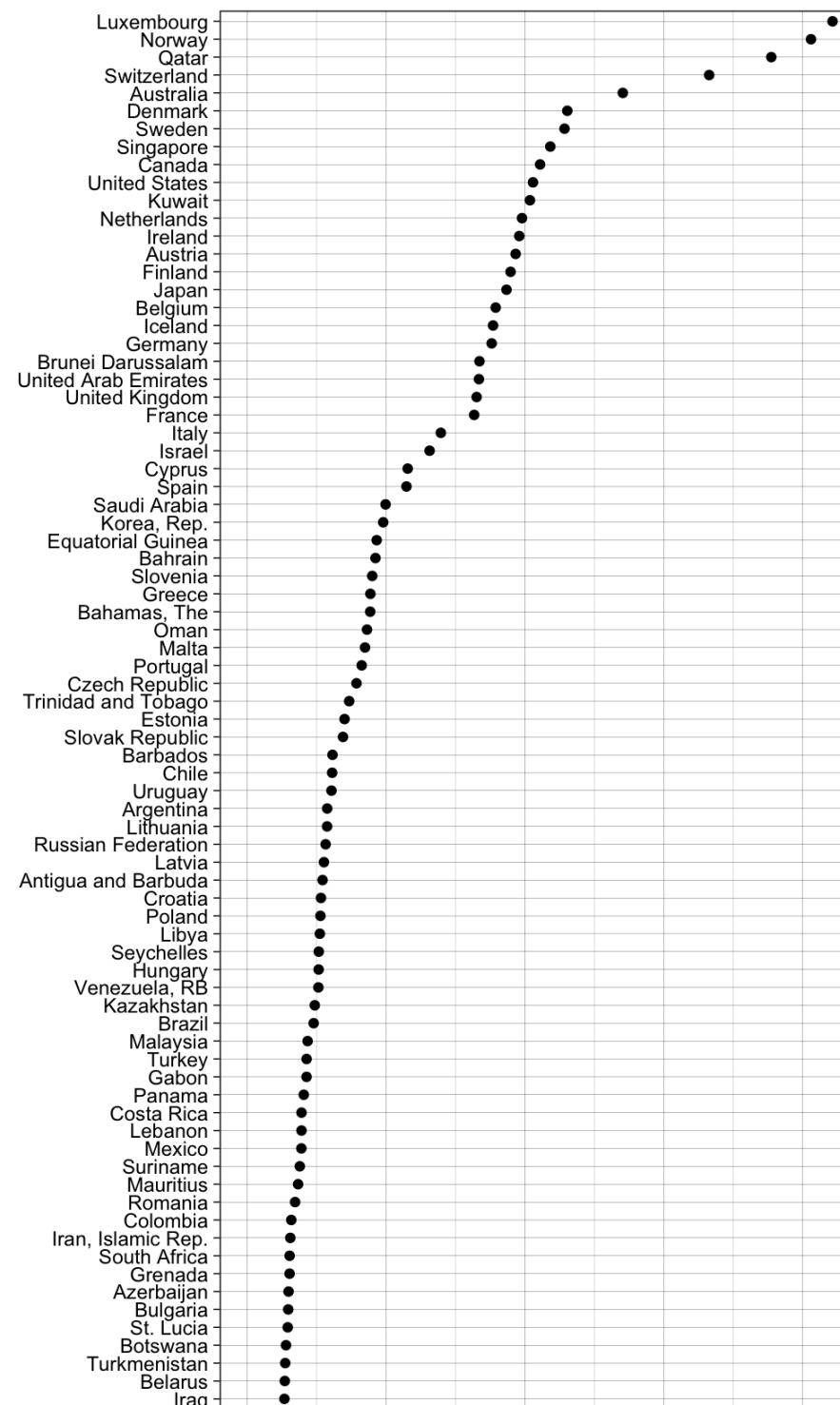
Sorted by 1997 fatality rate

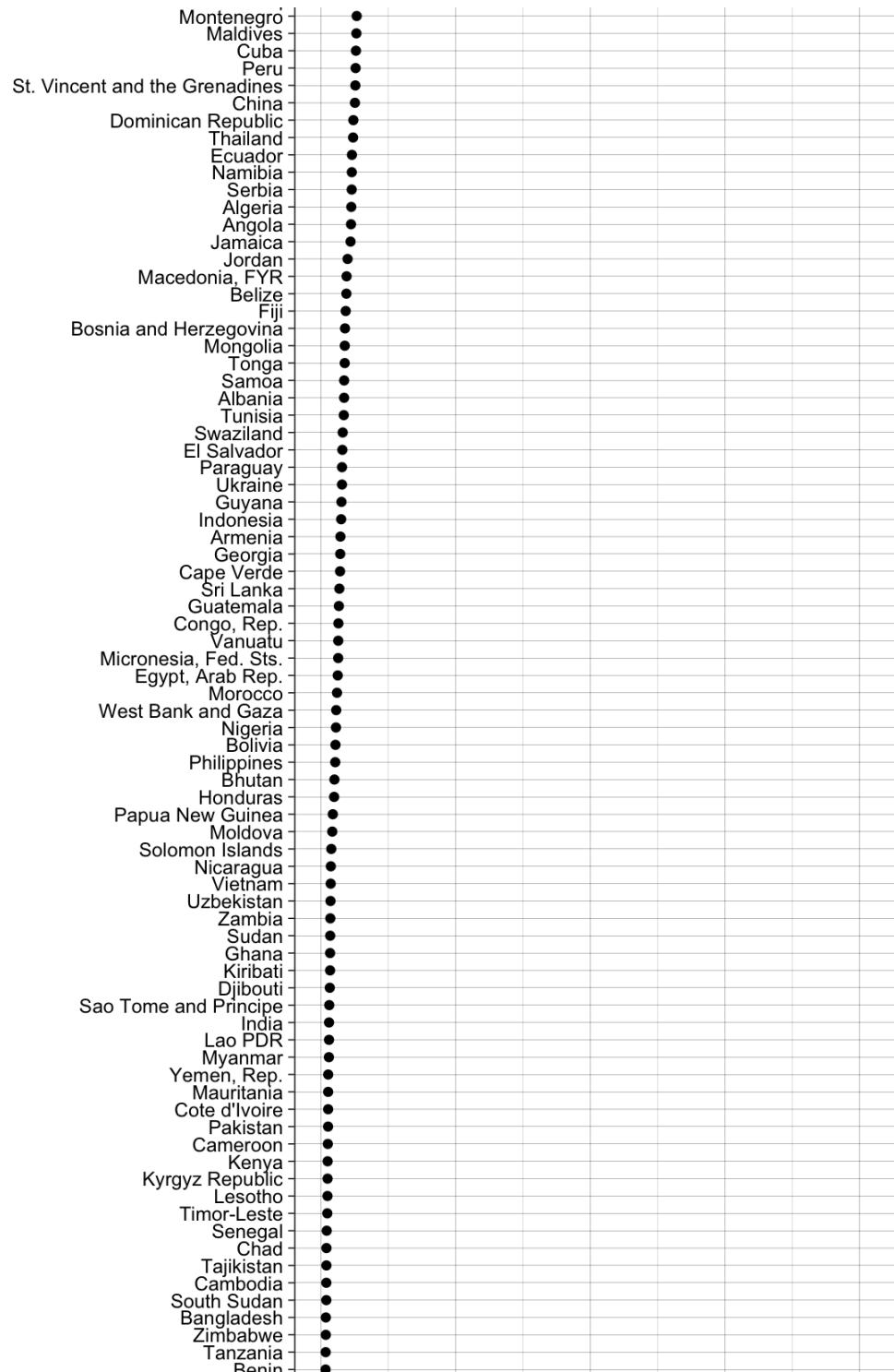


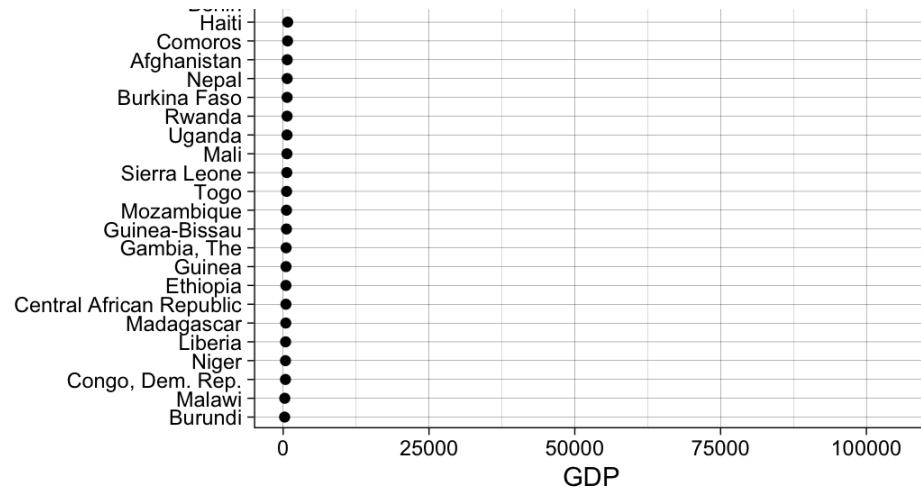
Large number of categories

Scroll

In chunk options: {r fig.height=20}



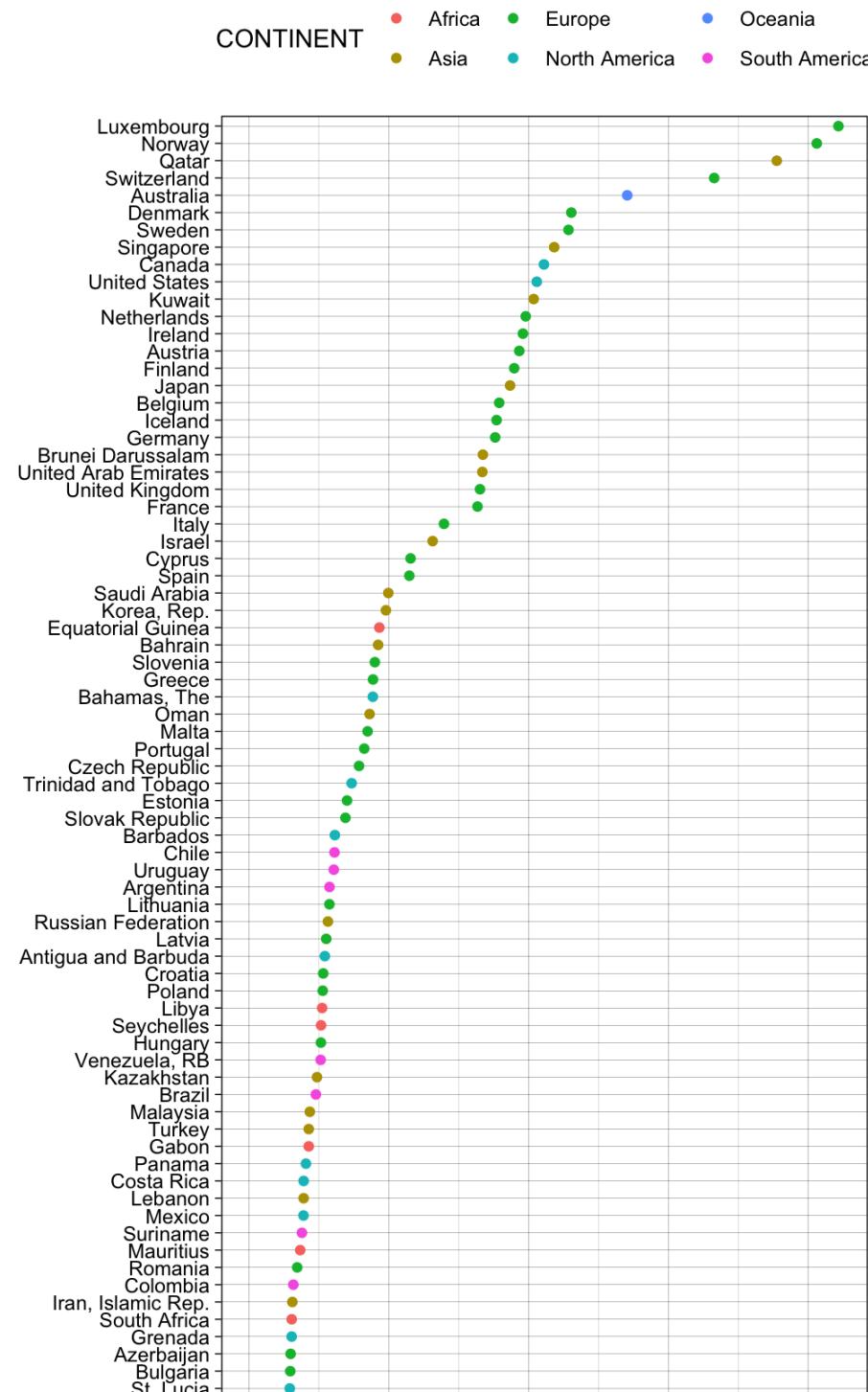


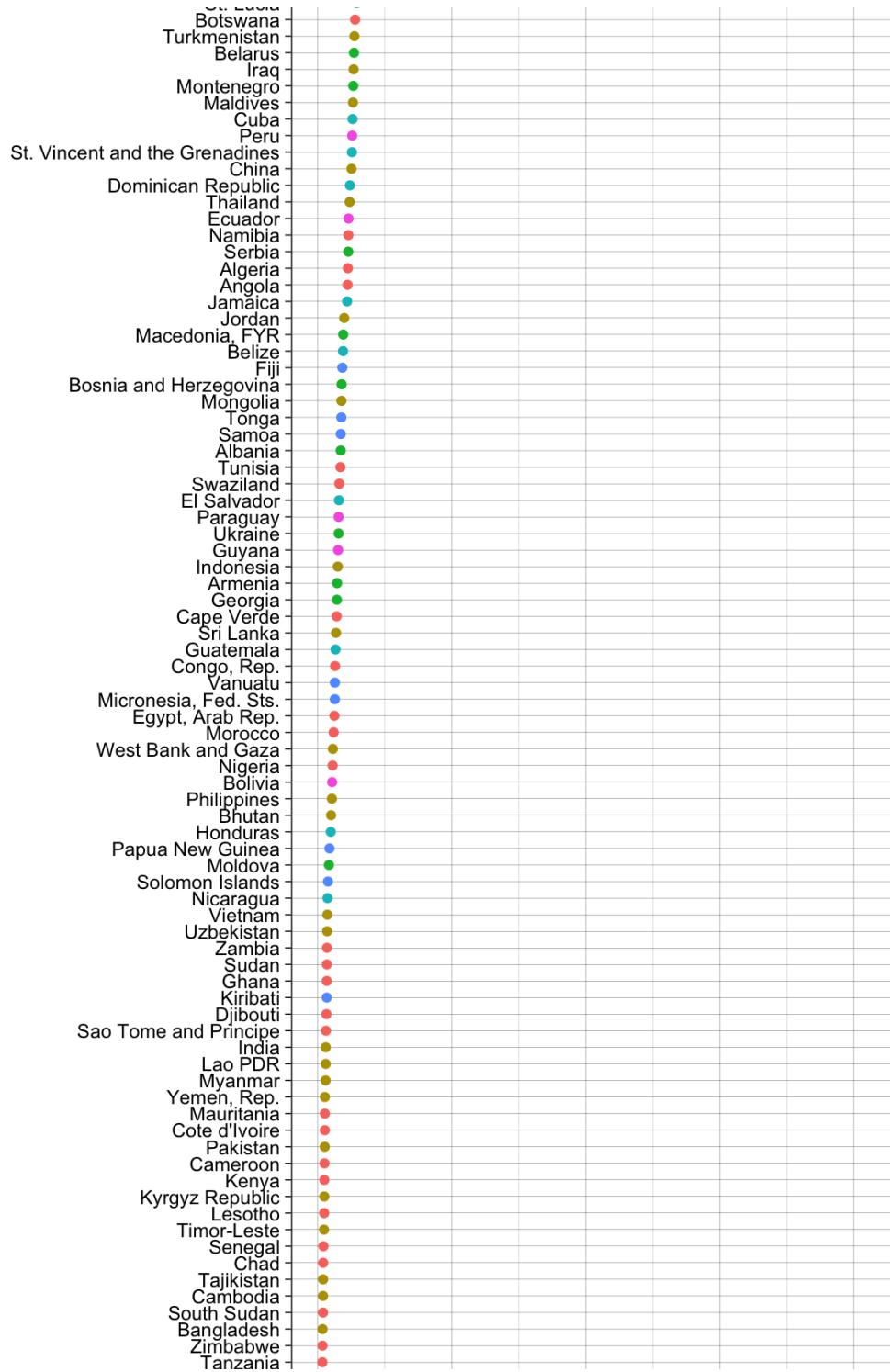


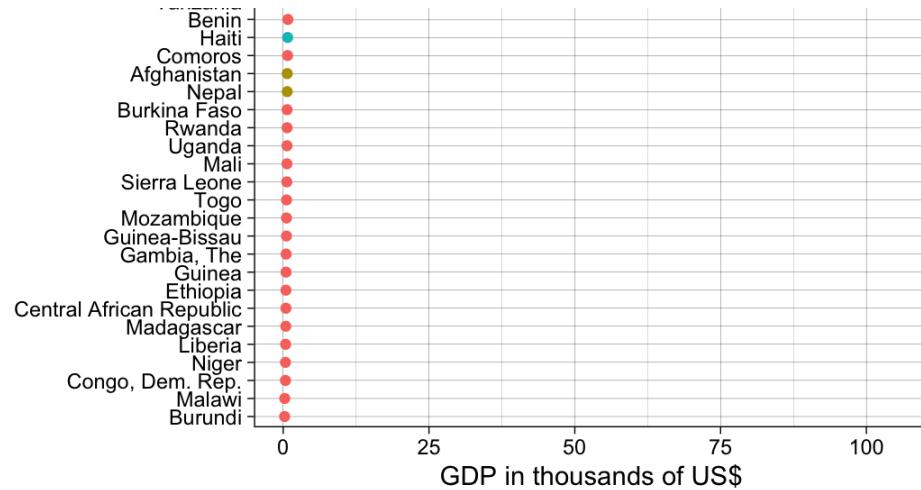
Large number of categories (with color)

Scroll

In chunk options: {r fig.height=20}

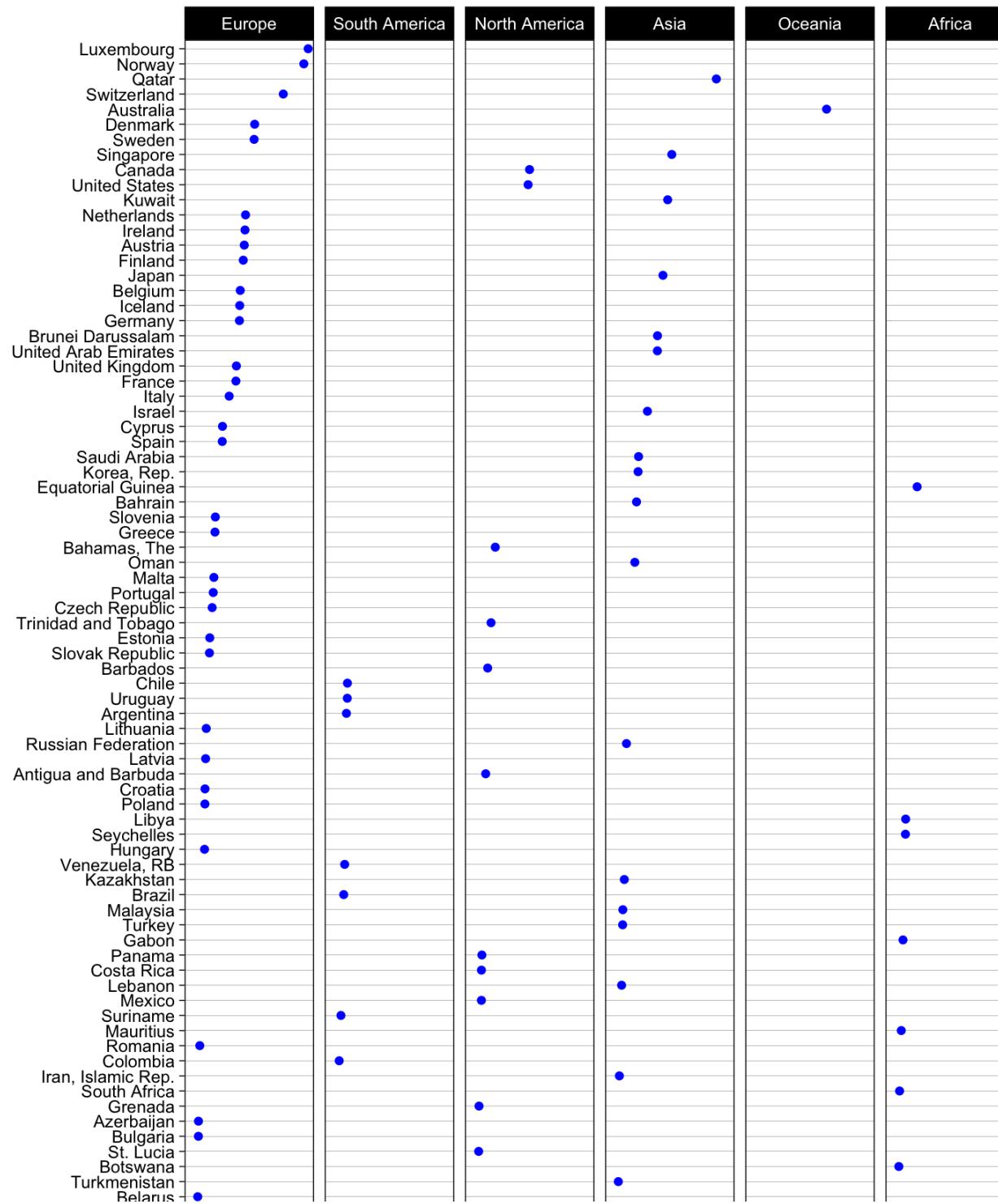




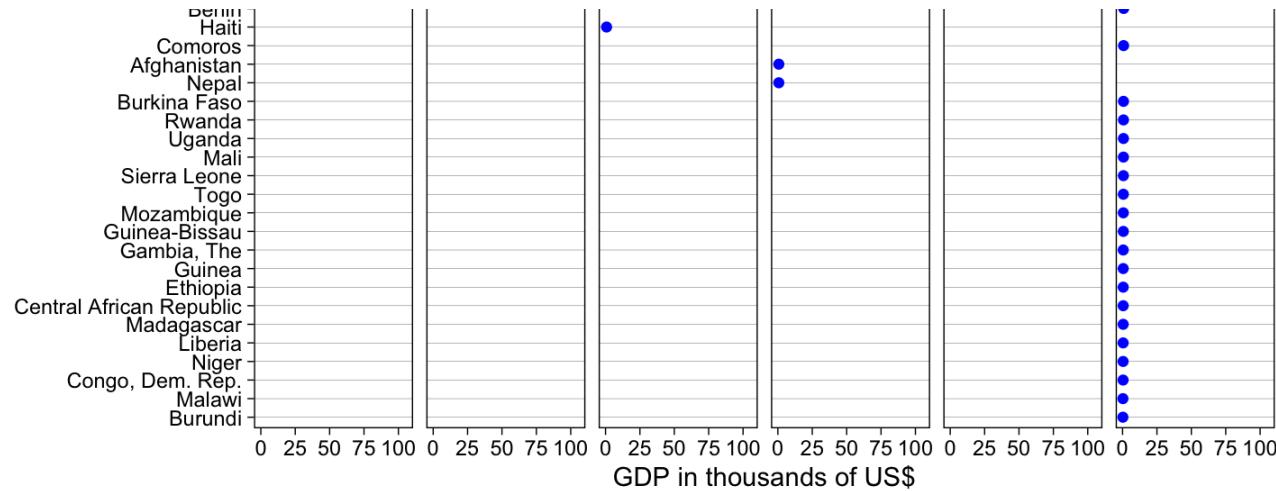


Cleveland dot plot with facets

What's wrong?

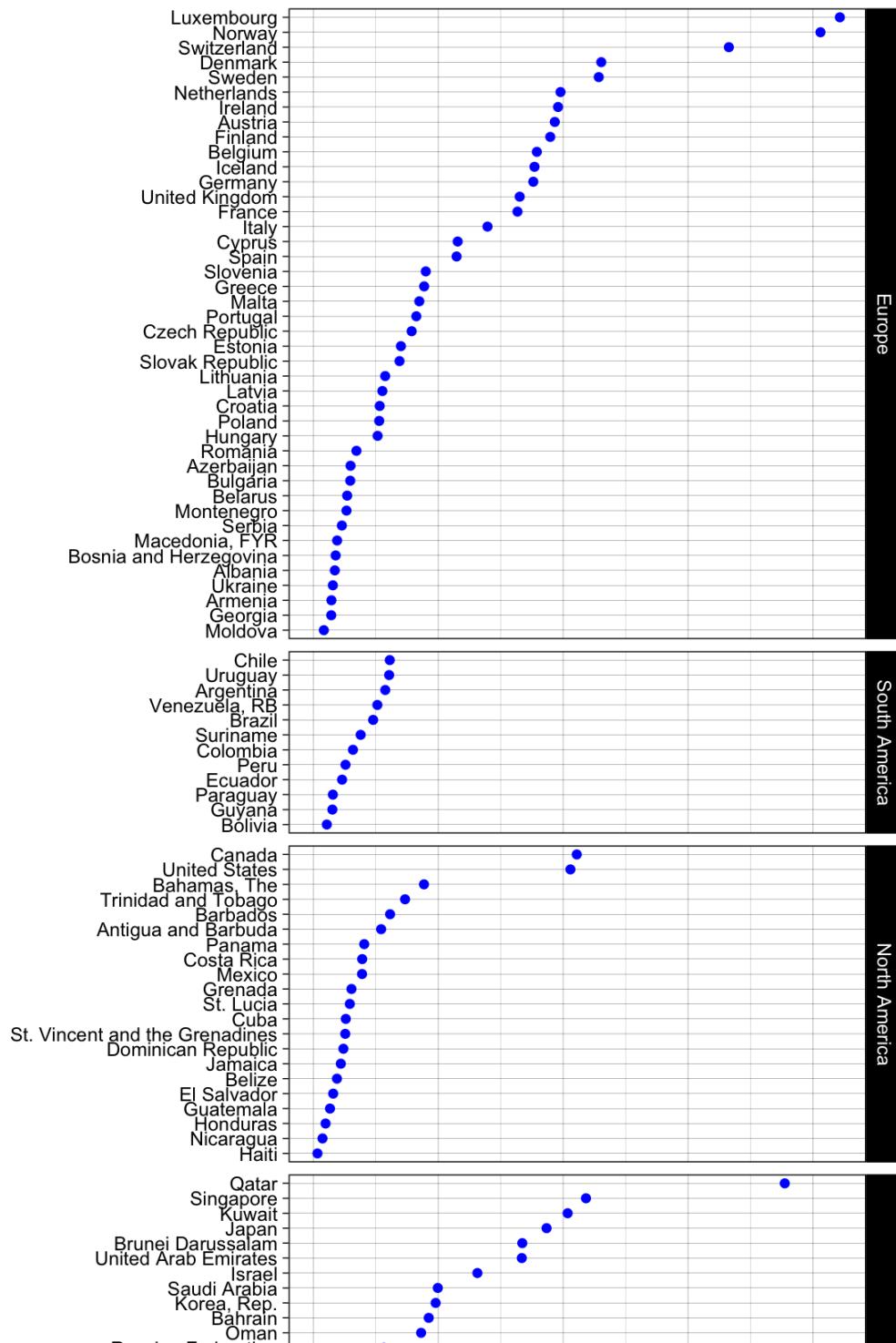


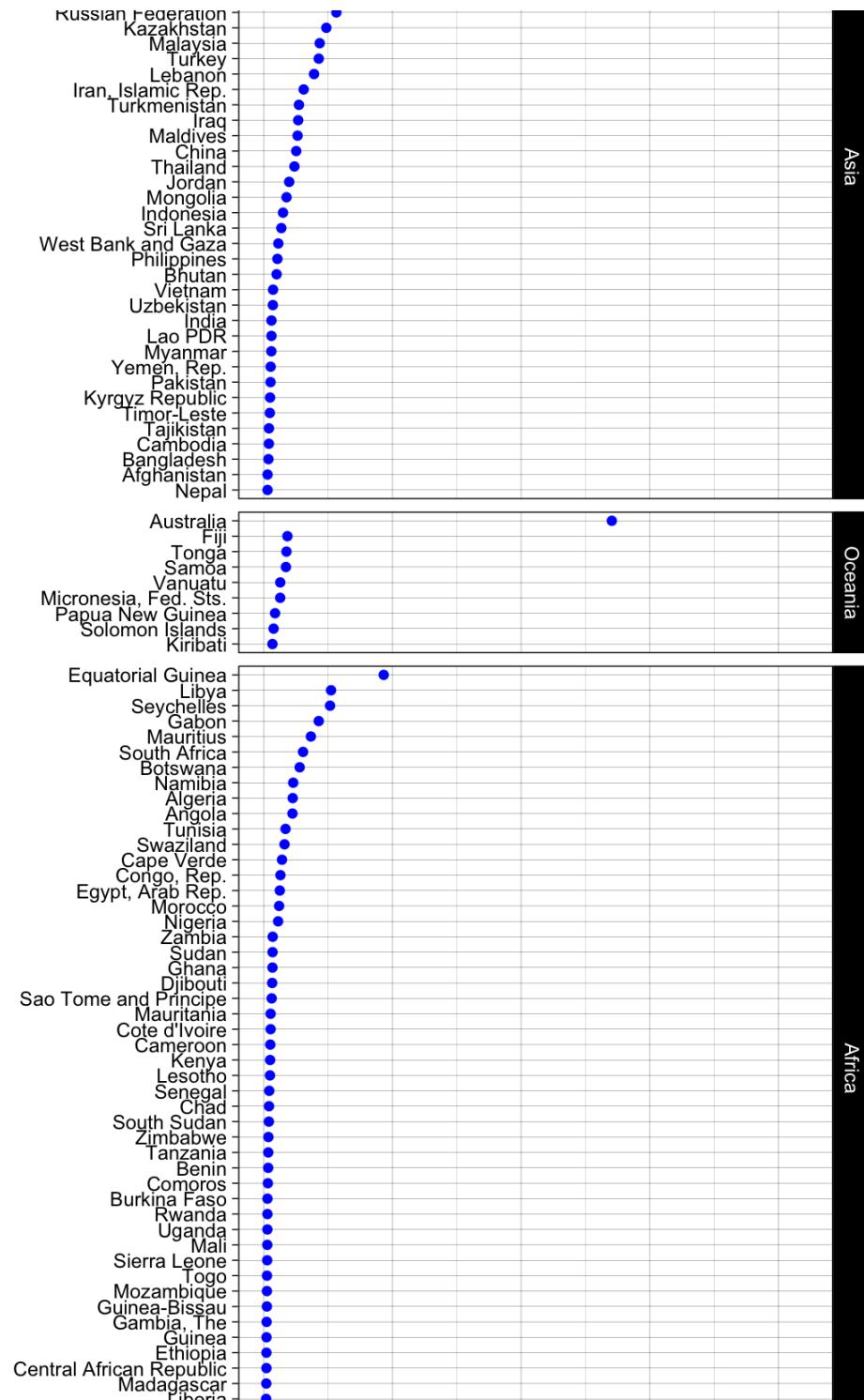


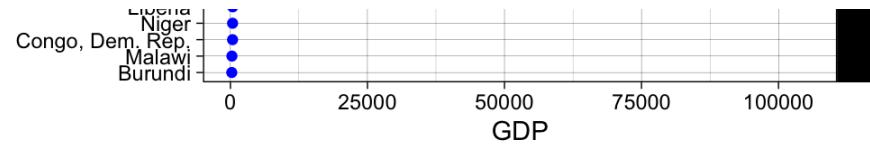


Cleveland dot plot with facets

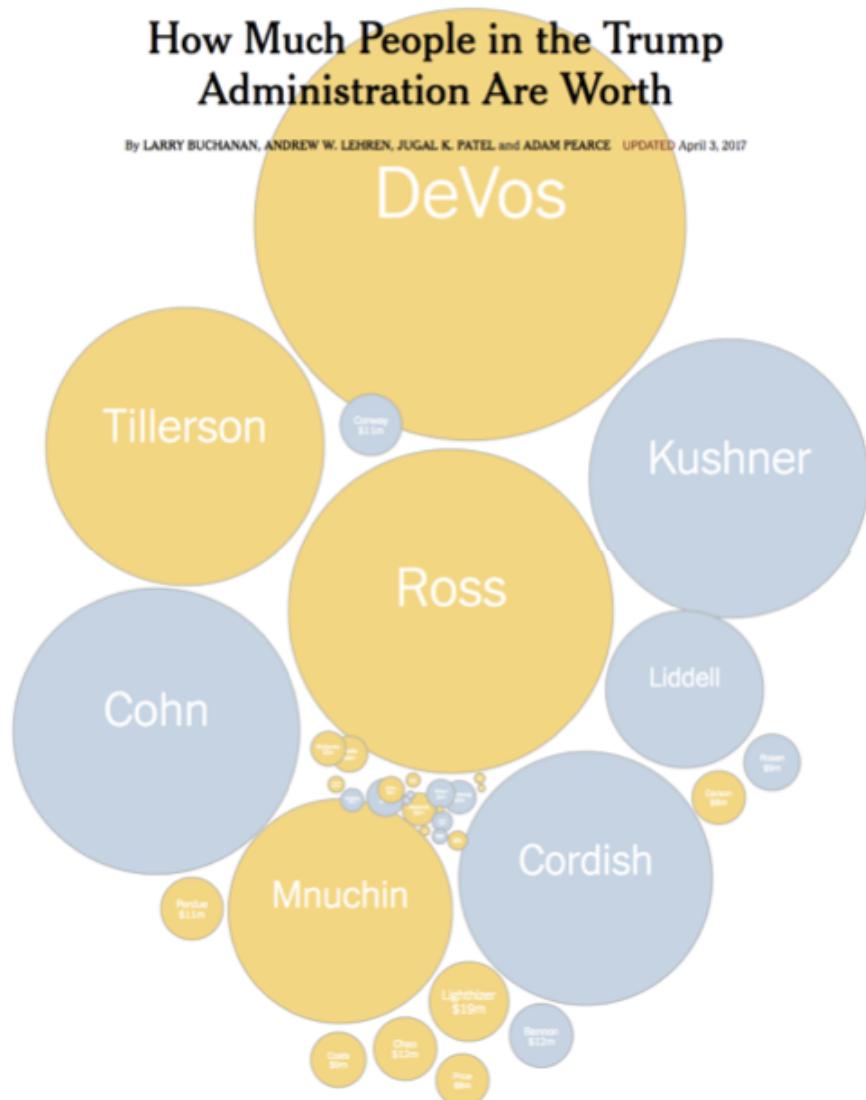
```
scales = "free_y", space = "free_y"
```



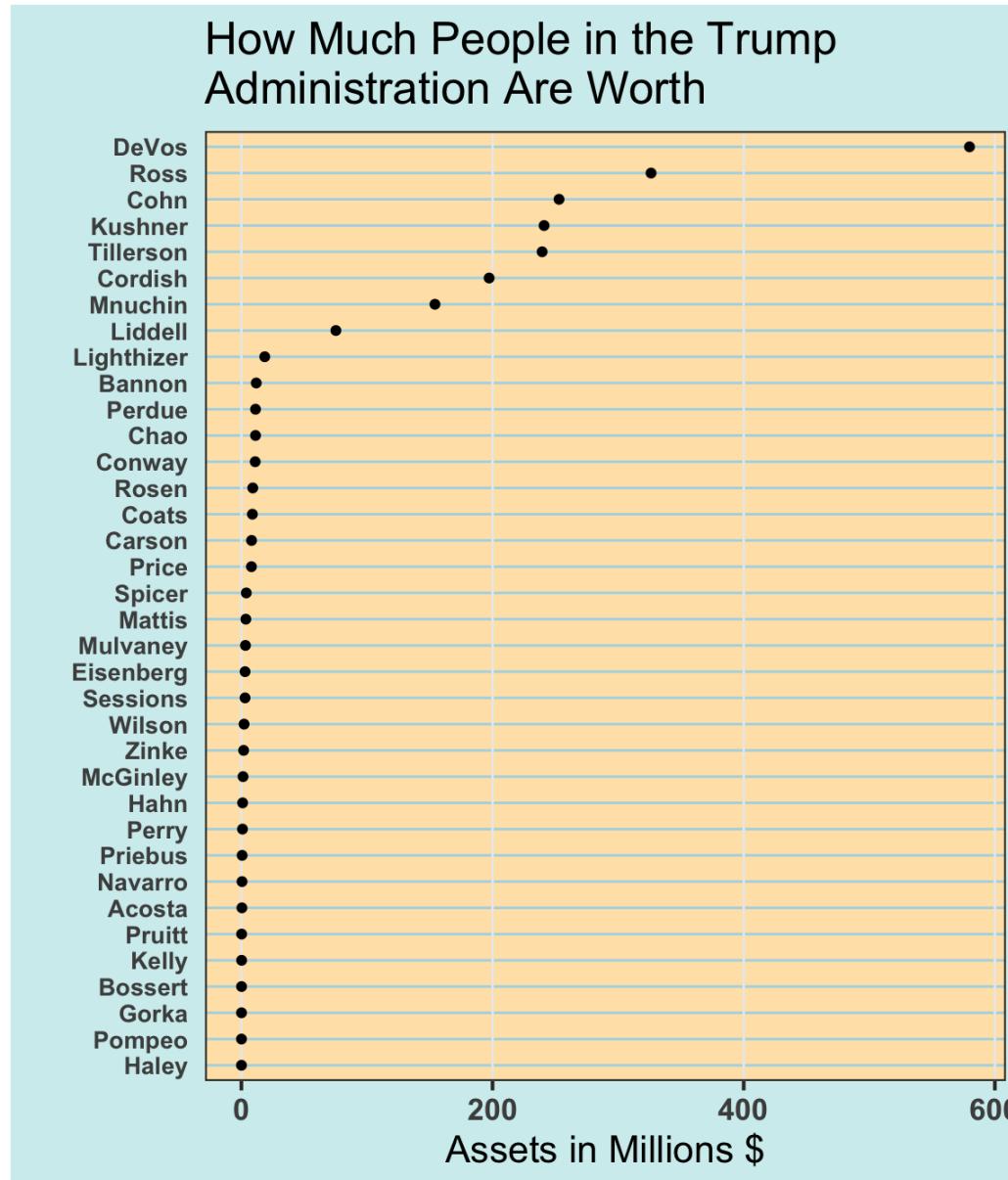




Trump Administration Assets

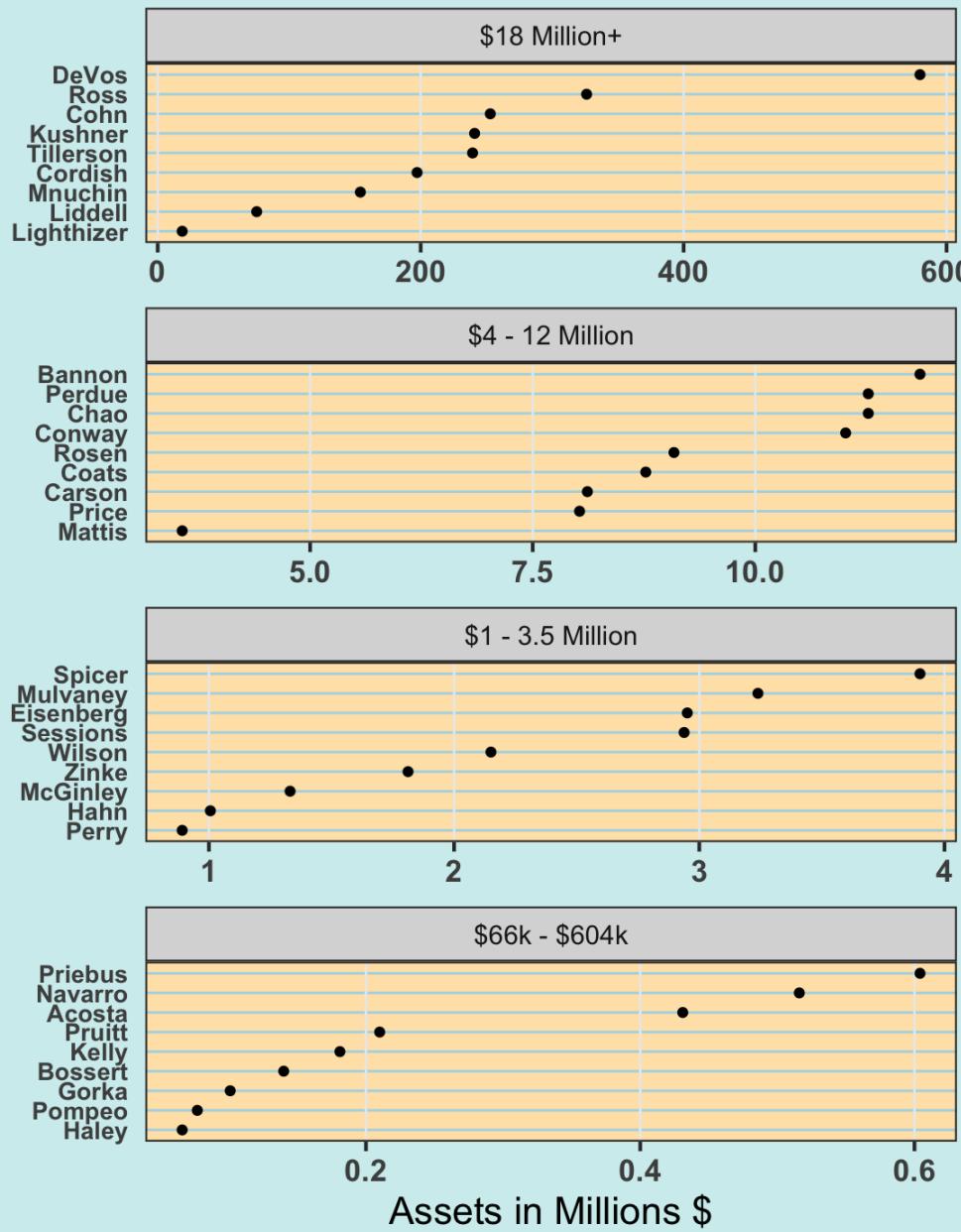


Redraw as Cleveland Dot Plot

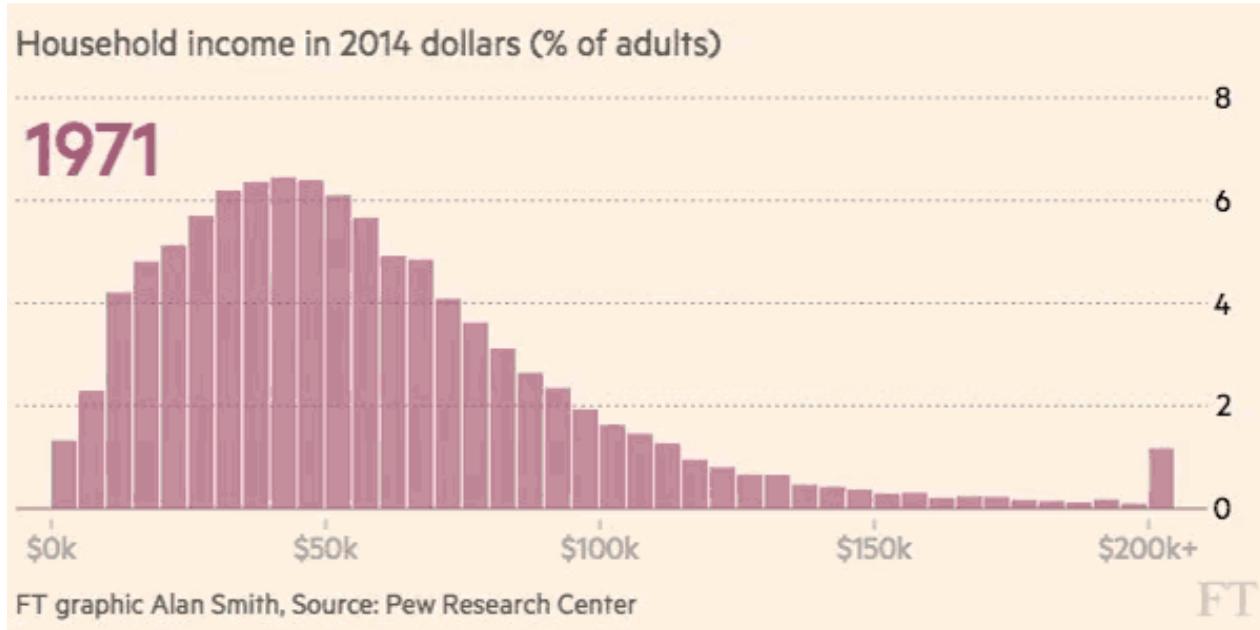


Cleveland Dot Plot with Panels

How Much People in the Trump Administration Are Worth



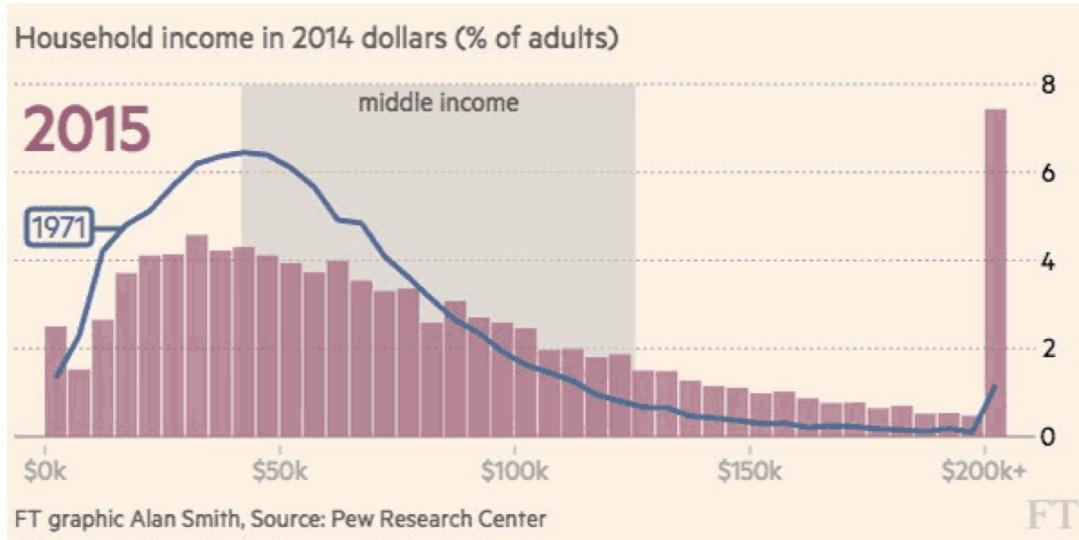
“Other” “or more” categories: “topcoding”



Source: “America’s explosion of income inequality, in one amazing animated chart”

<http://www.latimes.com/business/hiltzik/la-fi-hiltzik-ft-graphic-20160320-snap-htmlstory.html>

“Other” “or more” categories: “topcoding”



Source: “America’s explosion of income inequality, in one amazing animated chart”

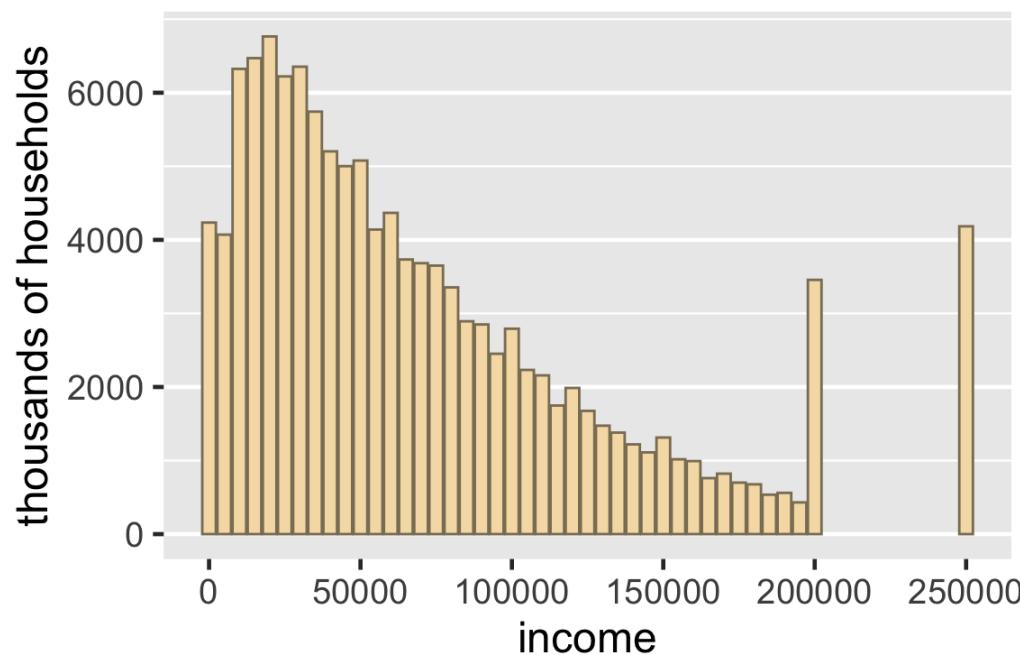
<http://www.latimes.com/business/hiltzik/la-fi-hiltzik-ft-graphic-20160320-snap-htmlstory.html>

“Other” “or more” categories

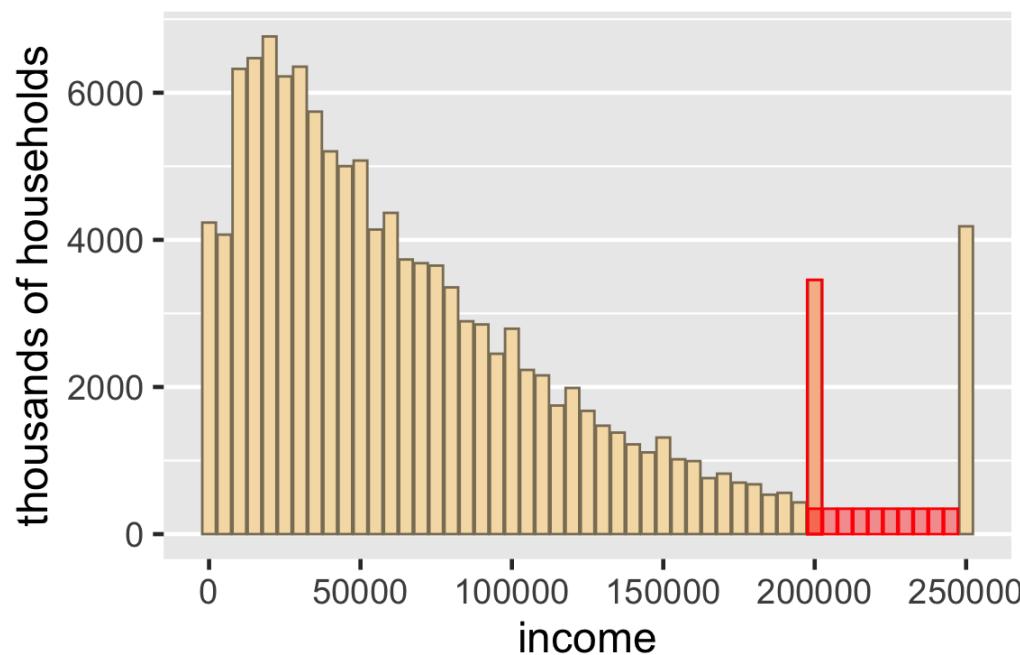
\$82,500 to \$84,999	\$85,000 to \$87,499	\$87,500 to \$89,999	\$90,000 to \$92,499	\$92,500 to \$94,999	\$95,000 to \$97,499	\$97,500 to \$99,999	\$100,000 and over	Va (D)
1,102	1,683	892	2,065	894	1,306	770	22,426	
\$82,500 to \$84,999	\$85,000 to \$87,499	\$87,500 to \$89,999	\$90,000 to \$92,499	\$92,500 to \$94,999	\$95,000 to \$97,499	\$97,500 to \$99,999	\$100,000 and over	Va (D)
973	1,520	775	1,880	824	1,172	711	20,773	
323	550	265	634	309	381	238	7,479	
650	970	509	1,246	515	790	473	13,295	
129	163	117	185	70	134	59	1,653	

Source: https://www2.census.gov/programs-surveys/cps/tables/pinc-01/2017/pinc01_1_1_1.xls

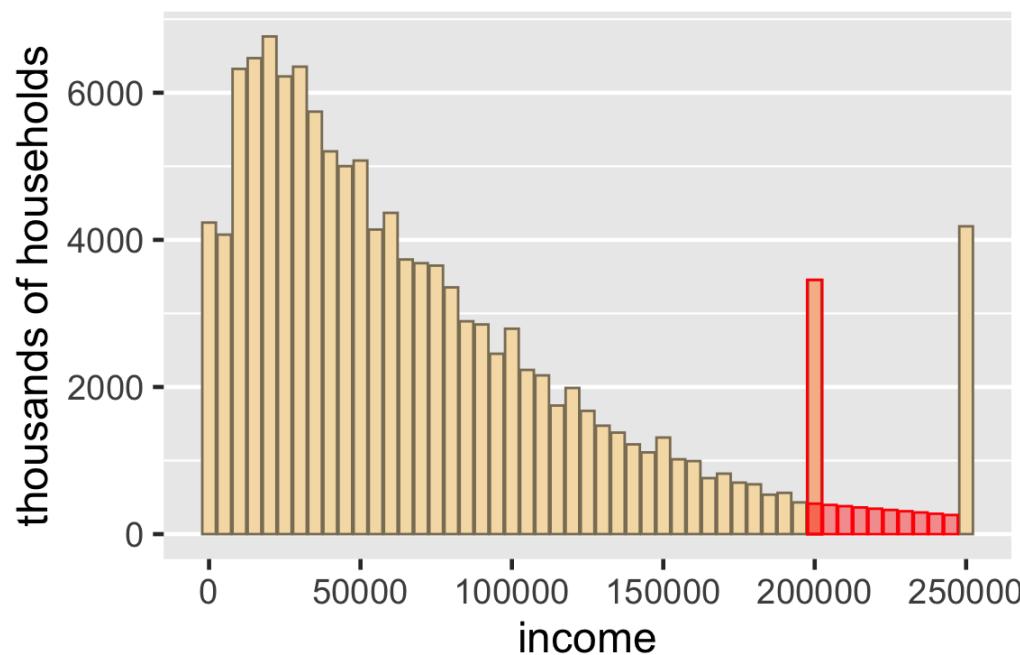
Household Income in 2015



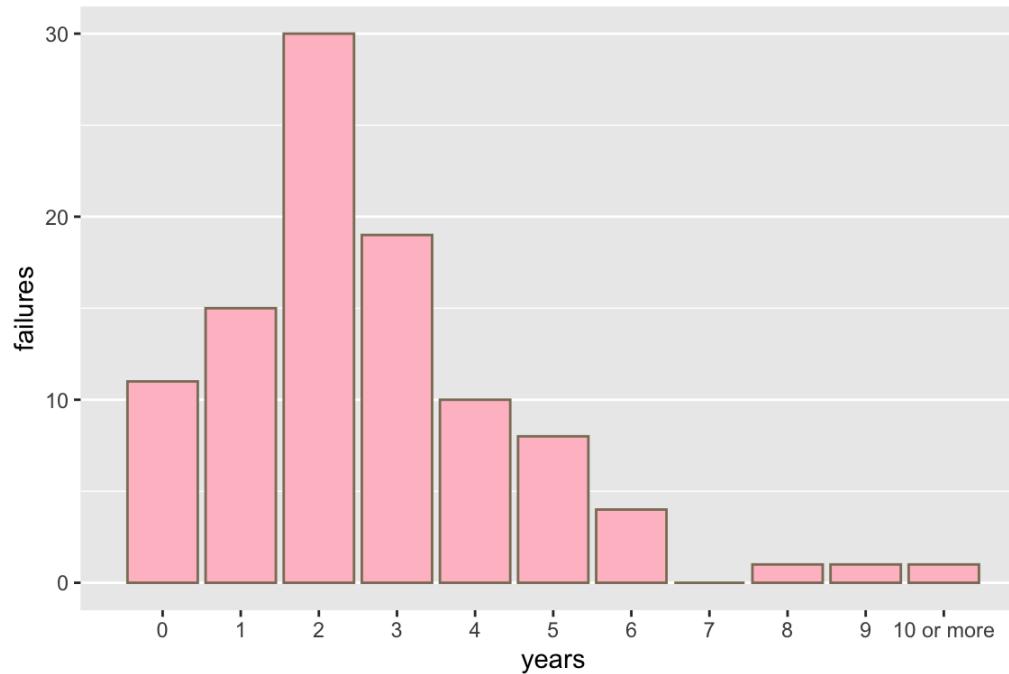
Household Income in 2015



Household Income in 2015



Reasonable use of “or more”



Data cleaning / transforming

<http://toddwschneider.com/posts/the-simpsons-by-the-data/>