

Categorical Variables (Chapter 4)

Prof. Joyce Robbins

Categorical Data

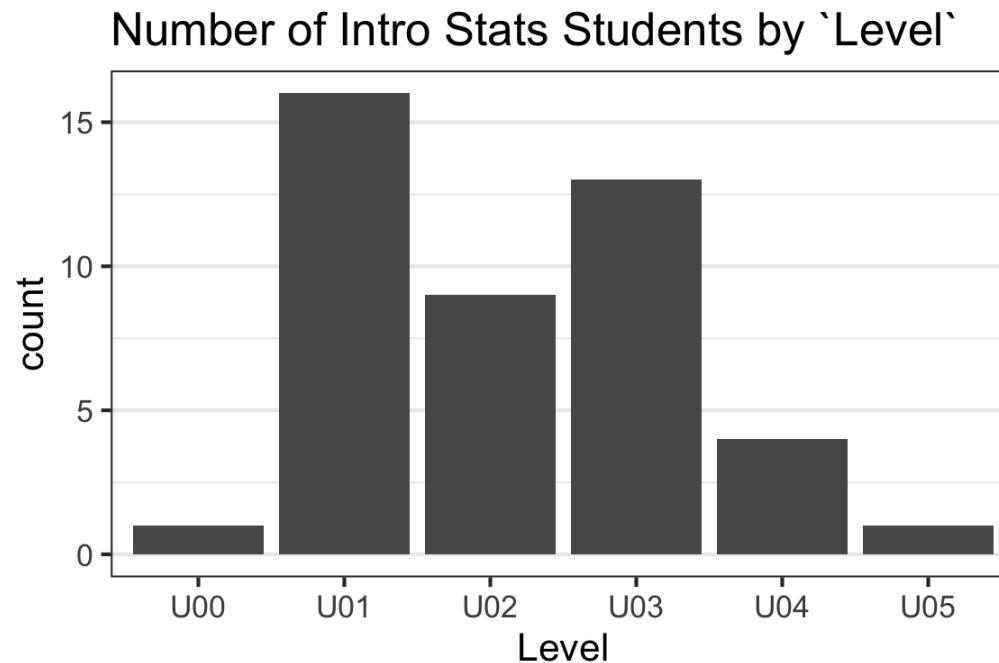
- hard to work with
- not a lot of options (esp. for 1 dimension)
- choice about which categories to display
- choice of the order of categories
- data cleaning takes more time

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

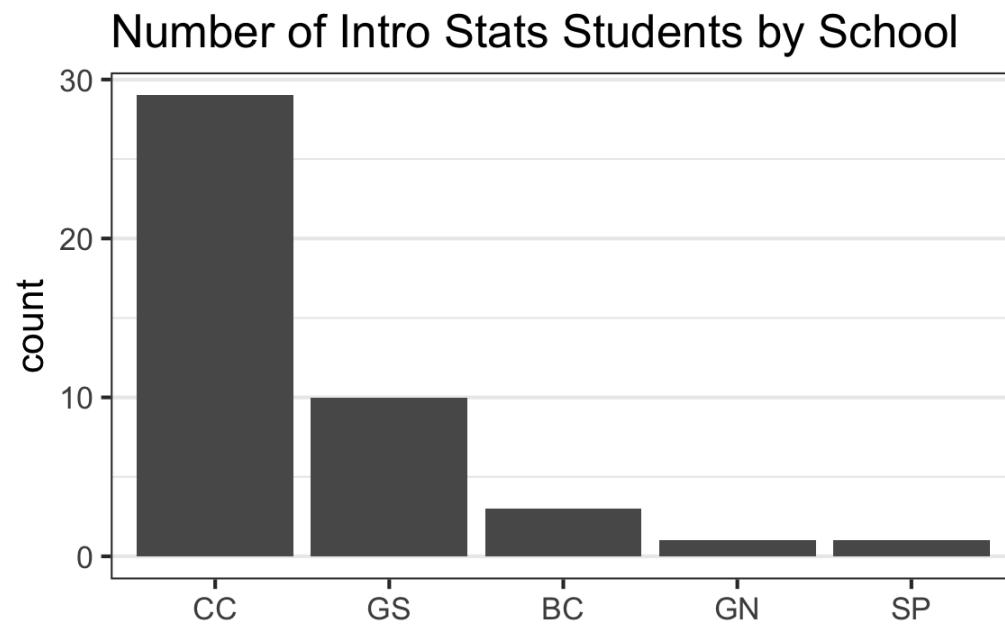
Ordinal data

Sort in logical order of the categories



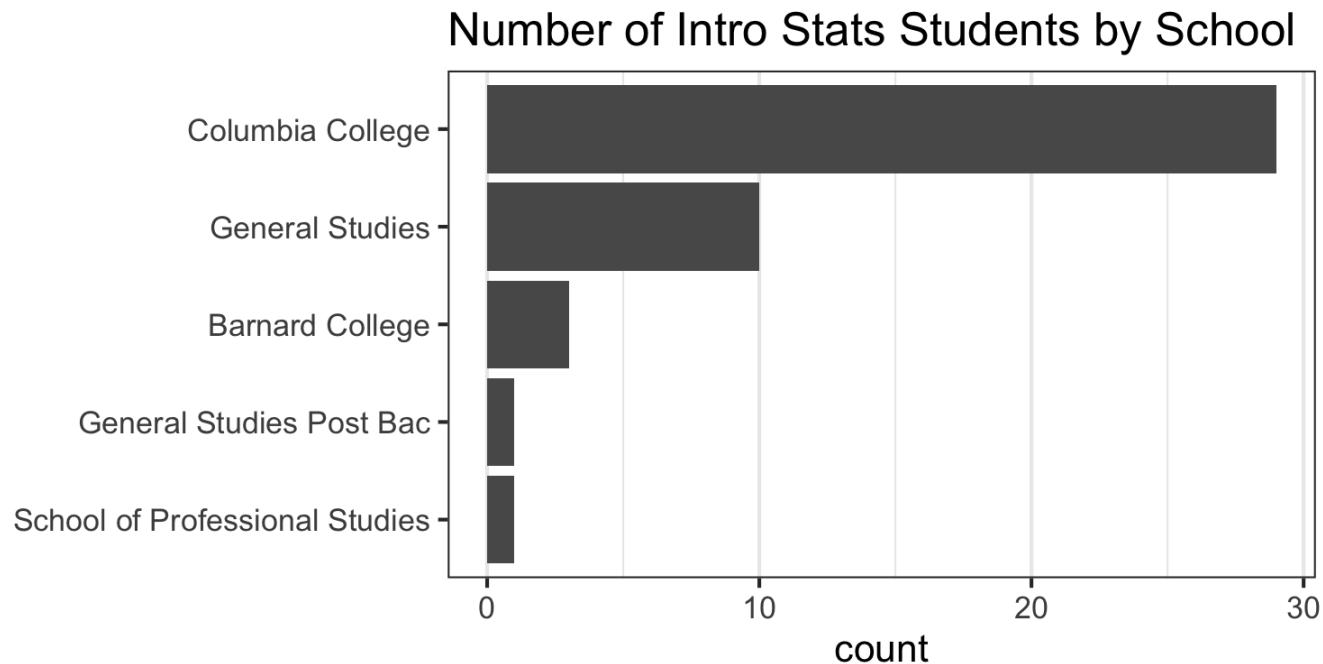
Nominal data

Sort from highest to lowest count (left to right)



Nominal data

... or top to bottom



Unbinned data

`forcats` package (part of tidyverse)

`geom_bar()` takes x only

- Vertical bars:

```
ggplot(df, aes(fct_infreq(x))) + geom_bar()
```

- Horizontal bars:

```
ggplot(df, aes(fct_rev(fct_infreq(x)))) + geom_bar() + coord_flip()
```

Binned data

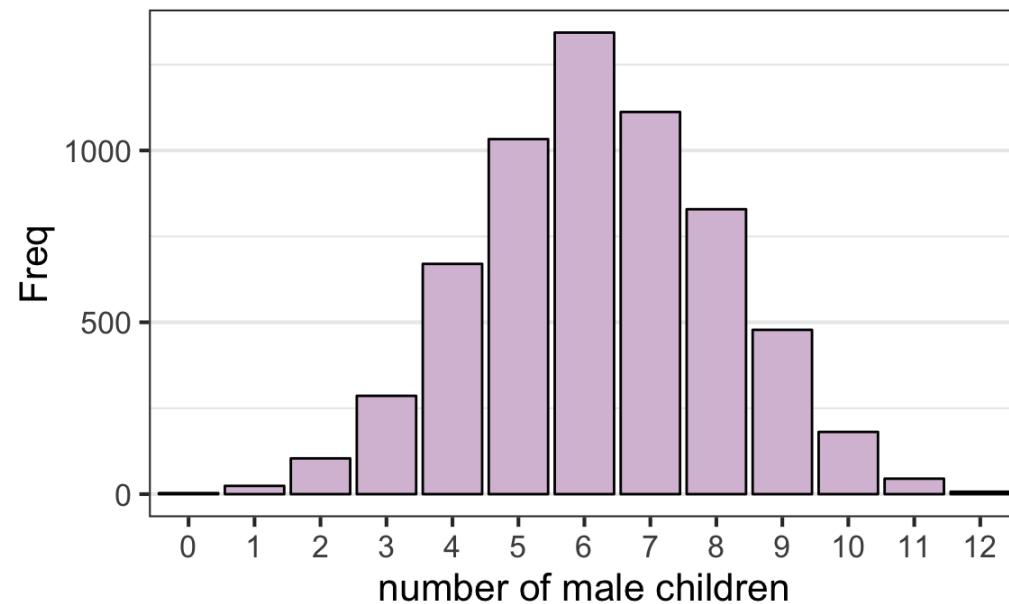
`geom_col()` - takes x & y

`...aes(fct_reorder(x, y), y)`

<https://github.com/jtr13/codehelp/blob/master/R/reorder.md>

Discrete data

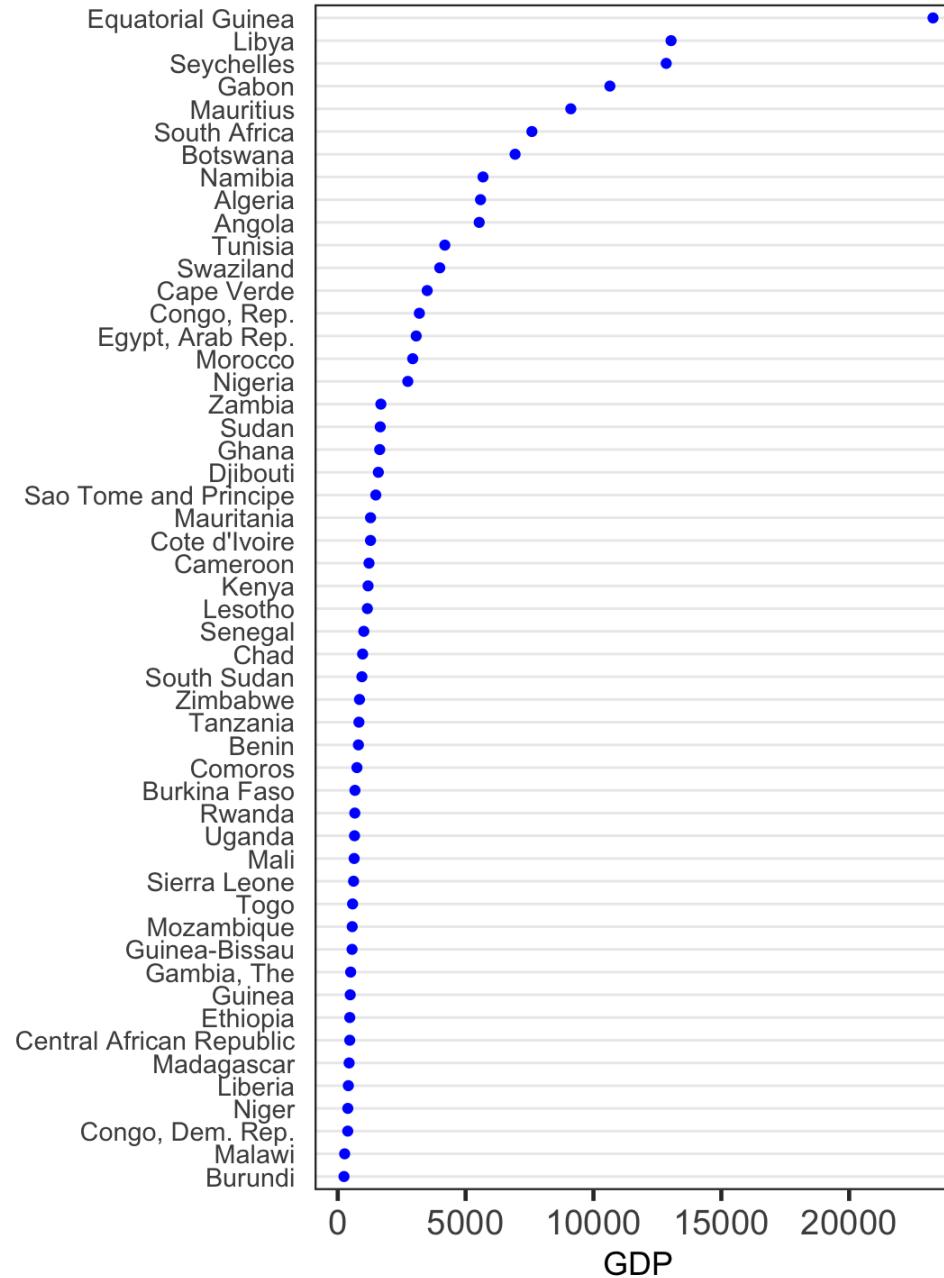
19c Saxony: # of males in families with 12 c



Large number of categories

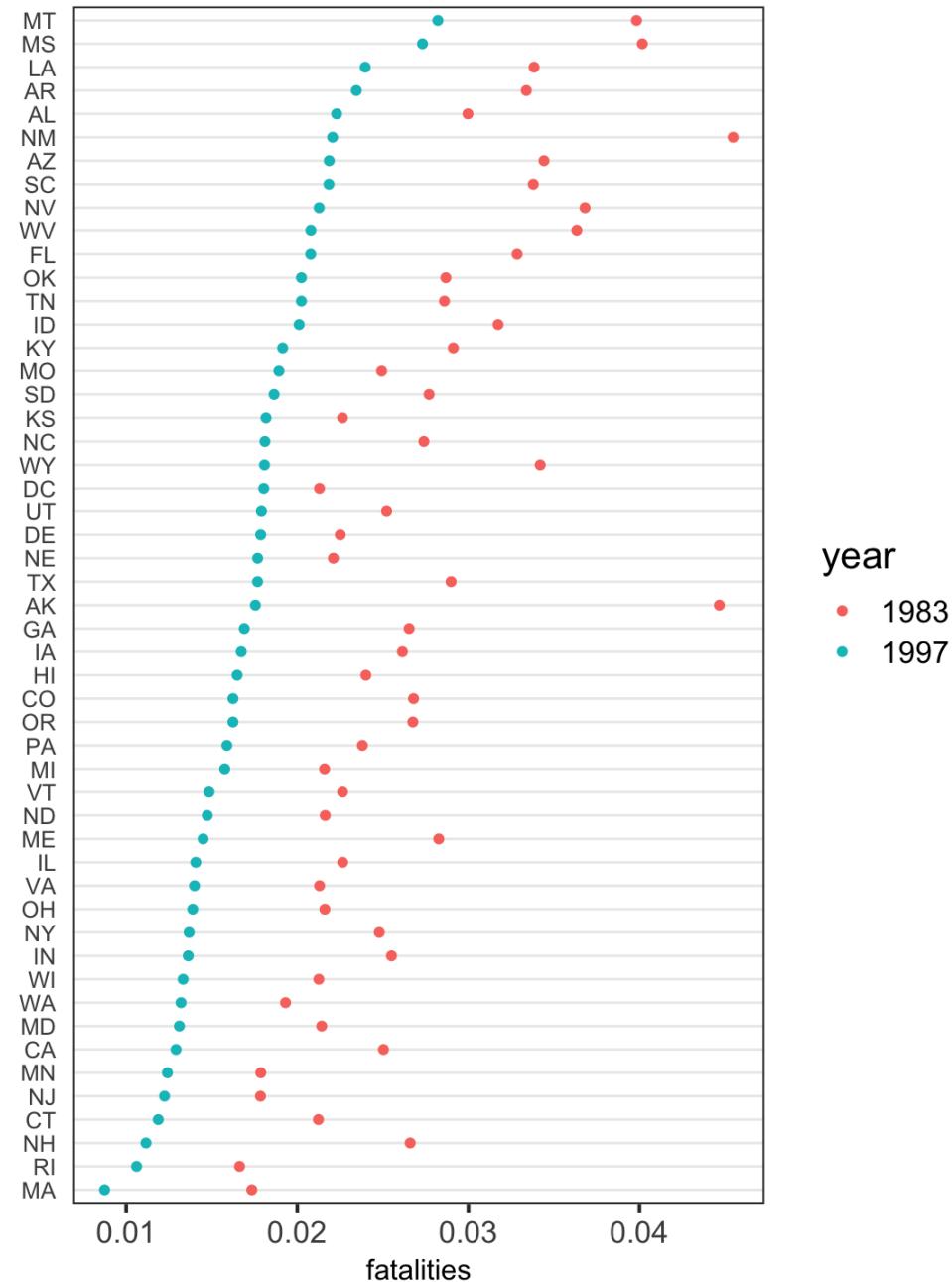
Cleveland dot plot

Africa: GDP per capita, 2012



Cleveland dot plot with multiple dots

of fatalities per million traffic miles

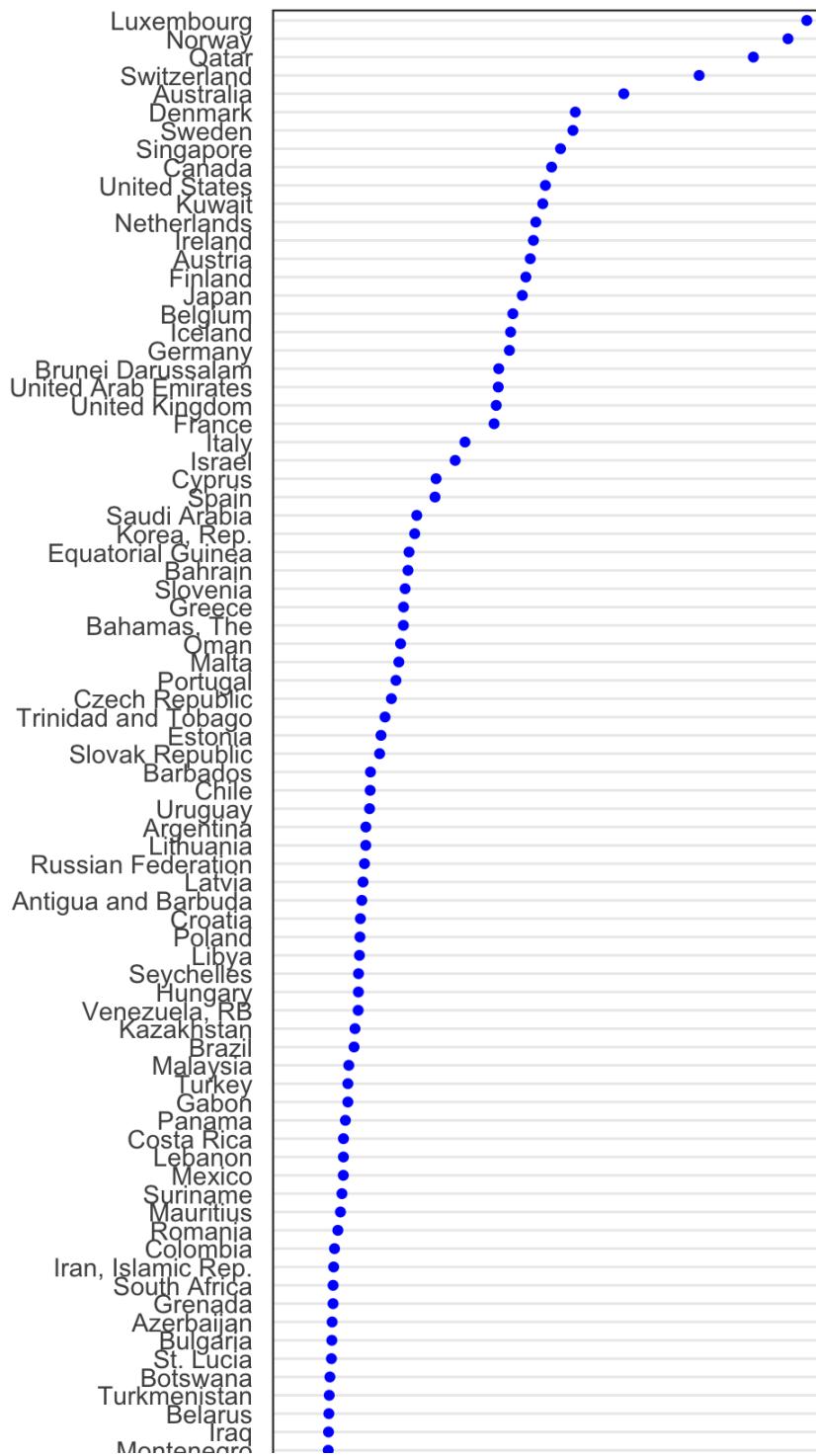


fct_reorder2 -> double sort: year, then fatalities

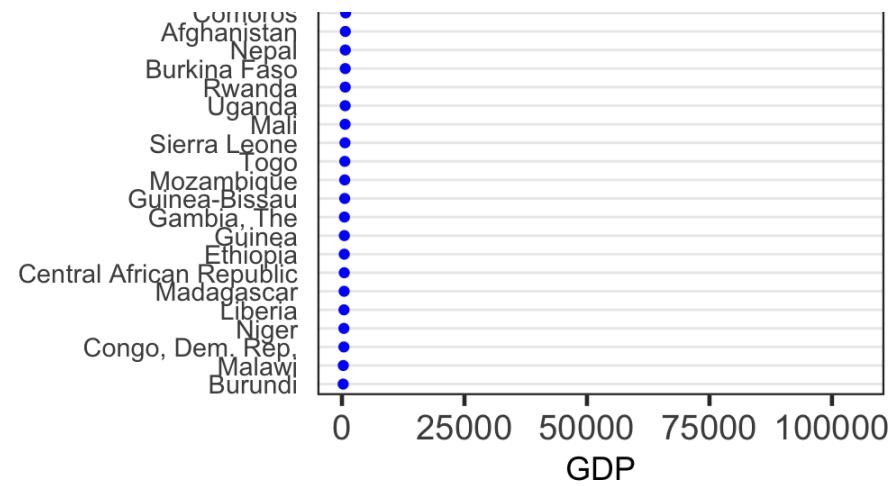
Large number of categories

Scroll

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

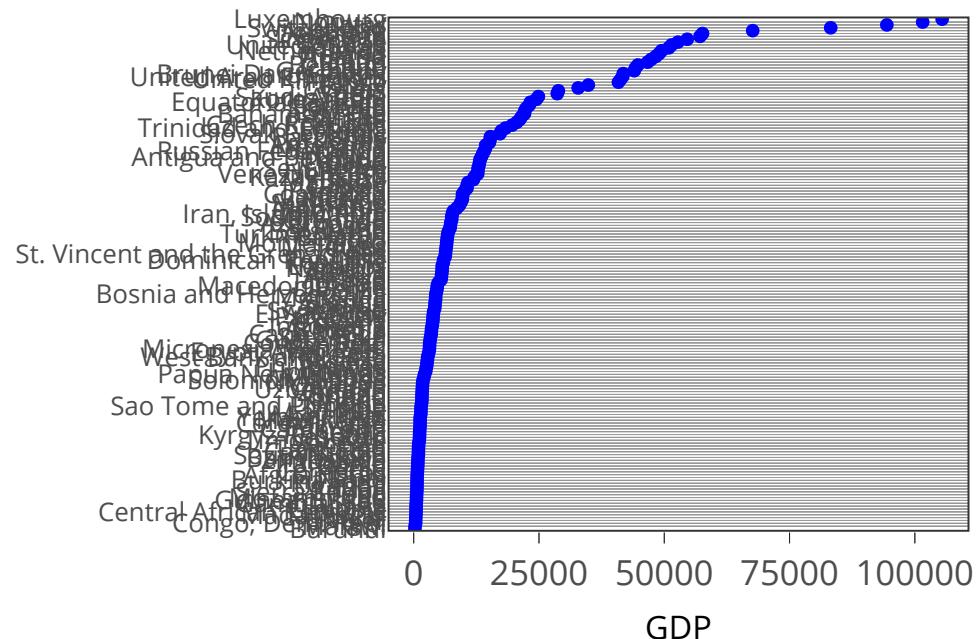


Maldives
Cuba
Peru
St. Vincent and the Grenadines
China
Dominican Republic
Thailand
Ecuador
Namibia
Serbia
Algeria
Angola
Jamaica
Jordan
Macedonia
FYR
Belize
Fiji
Bosnia and Herzegovina
Mongolia
Tonga
Samoa
Albania
Tunisia
Swaziland
El Salvador
Paraguay
Ukraine
Guyana
Indonesia
Armenia
Georgia
Cape Verde
Sri Lanka
Guatemala
Congo, Rep.
Vanuatu
Micronesia, Fed. Sts.
Egypt, Arab Rep.
Morocco
West Bank and Gaza
Nigeria
Bolivia
Philippines
Bhutan
Honduras
Papua New Guinea
Moldova
Solomon Islands
Nicaragua
Vietnam
Uzbekistan
Zambia
Sudan
Ghana
Kiribati
Djibouti
Sao Tome and Principe
India
Lao PDR
Myanmar
Yemen, Rep.
Mauritania
Cote d'Ivoire
Pakistan
Cameroon
Kenya
Kyrgyz Republic
Lesotho
Timor-Leste
Senegal
Chad
Tajikistan
Cambodia
South Sudan
Bangladesh
Zimbabwe
Tanzania
Benin
Comoros
Haiti



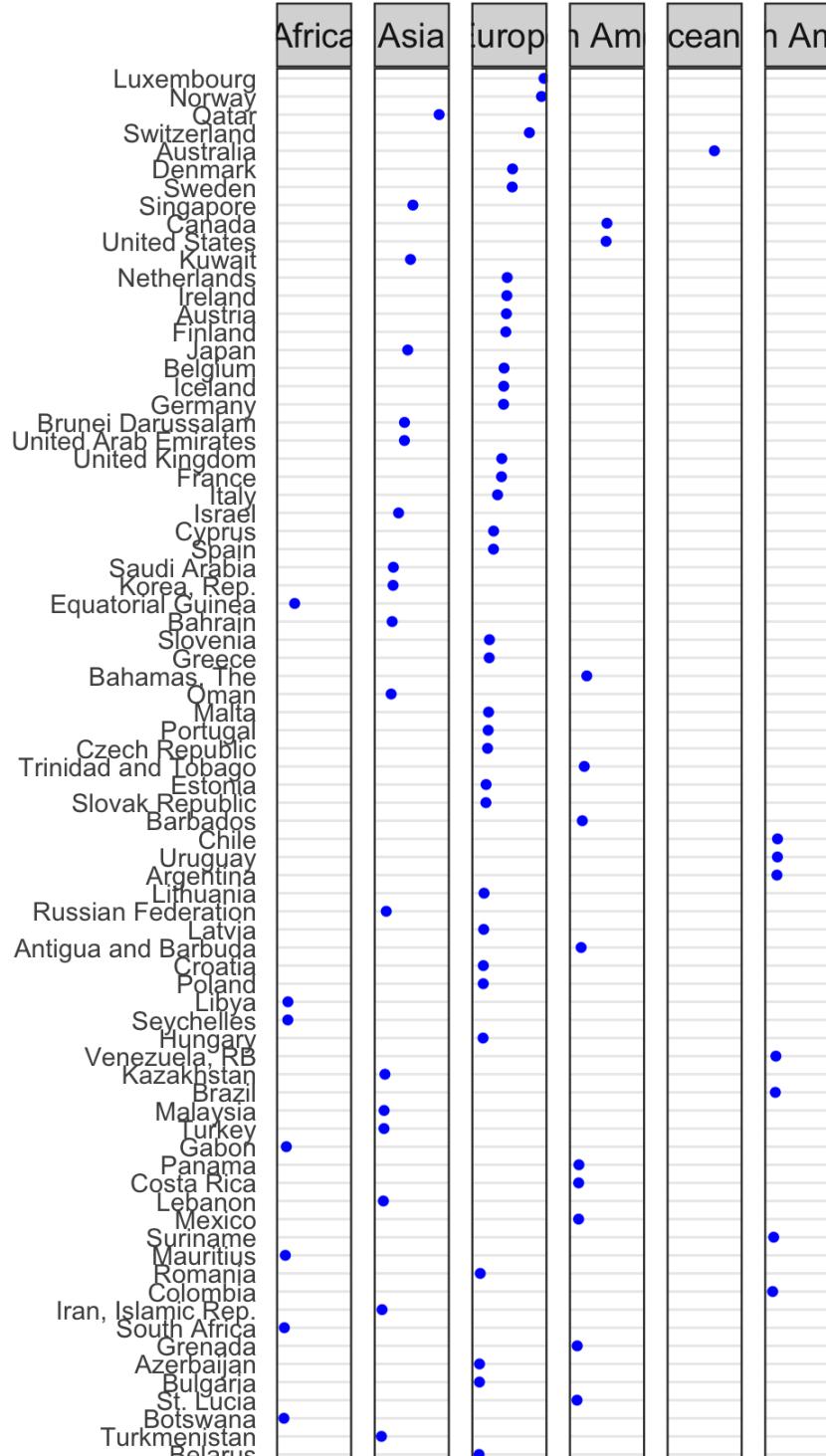
Interactive Cleveland Dot Plot

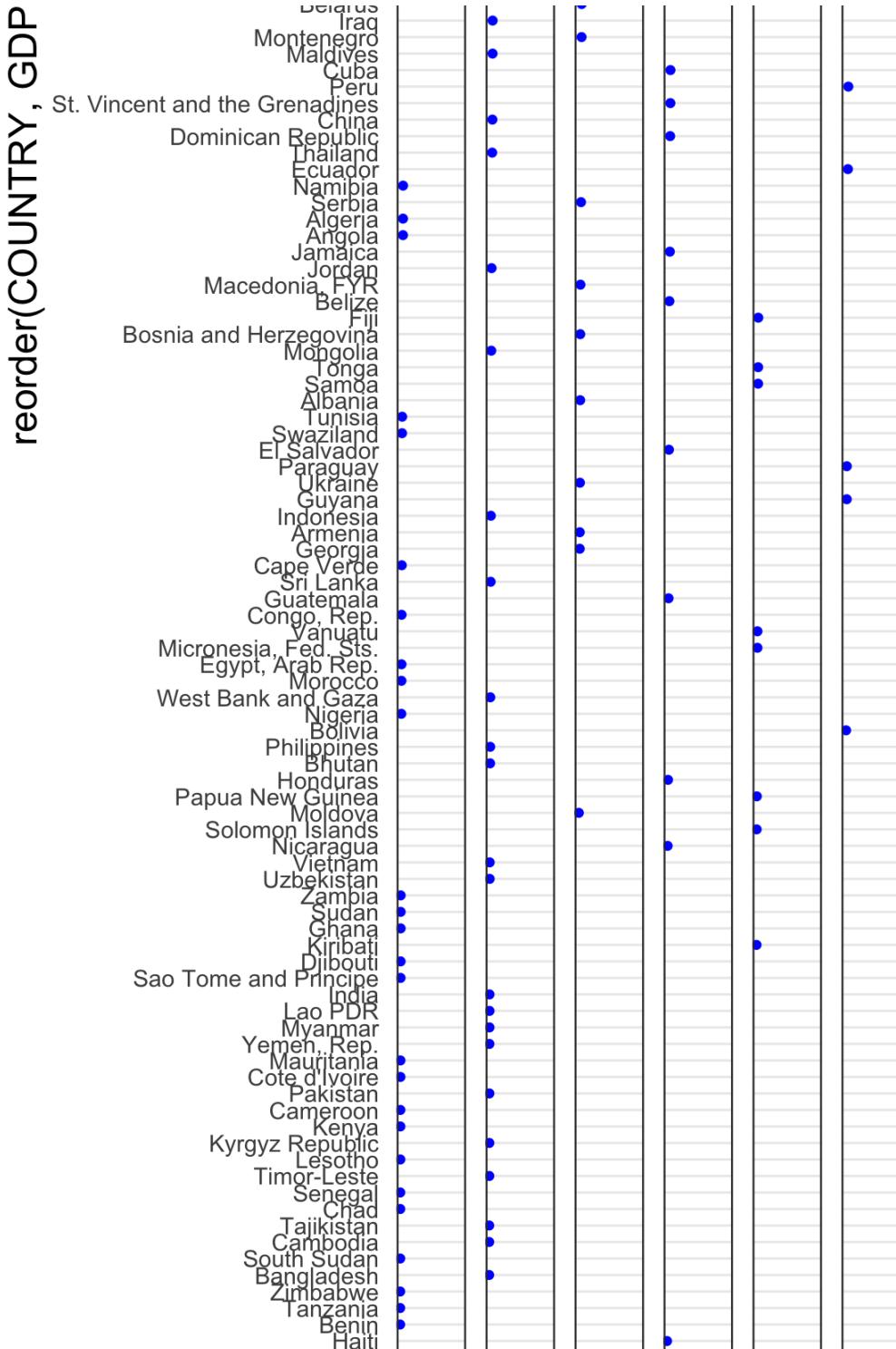
```
library(plotly)
ggplotly(g)
```

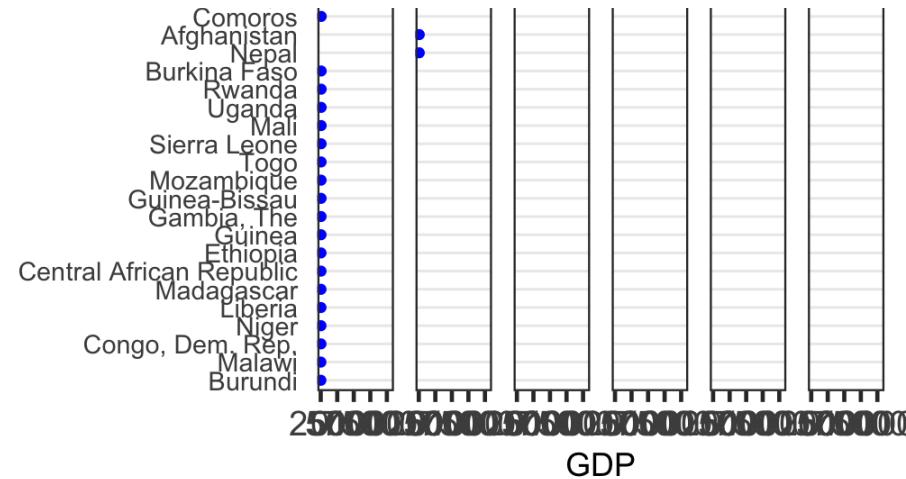


Cleveland dot plot with facets

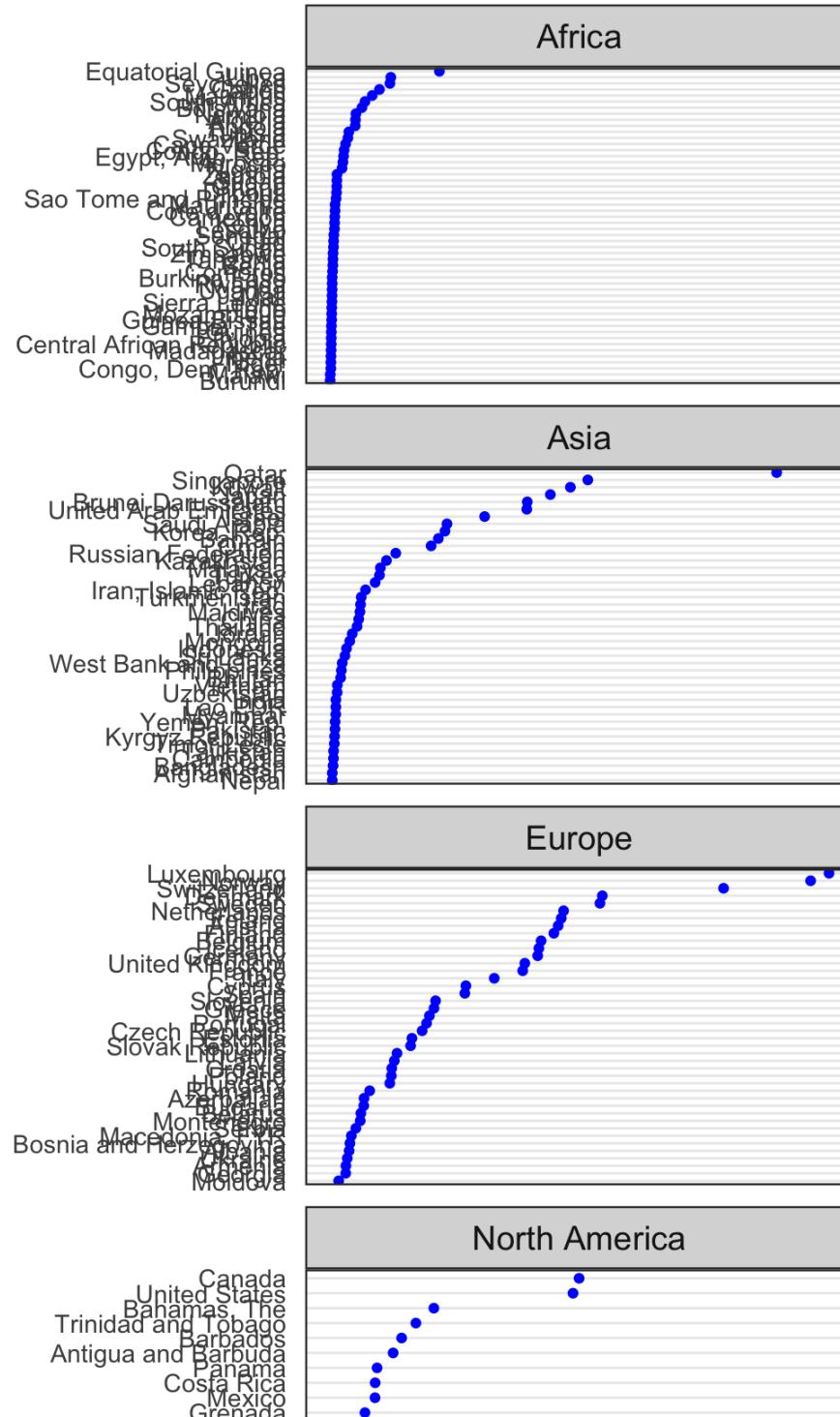
What's wrong?

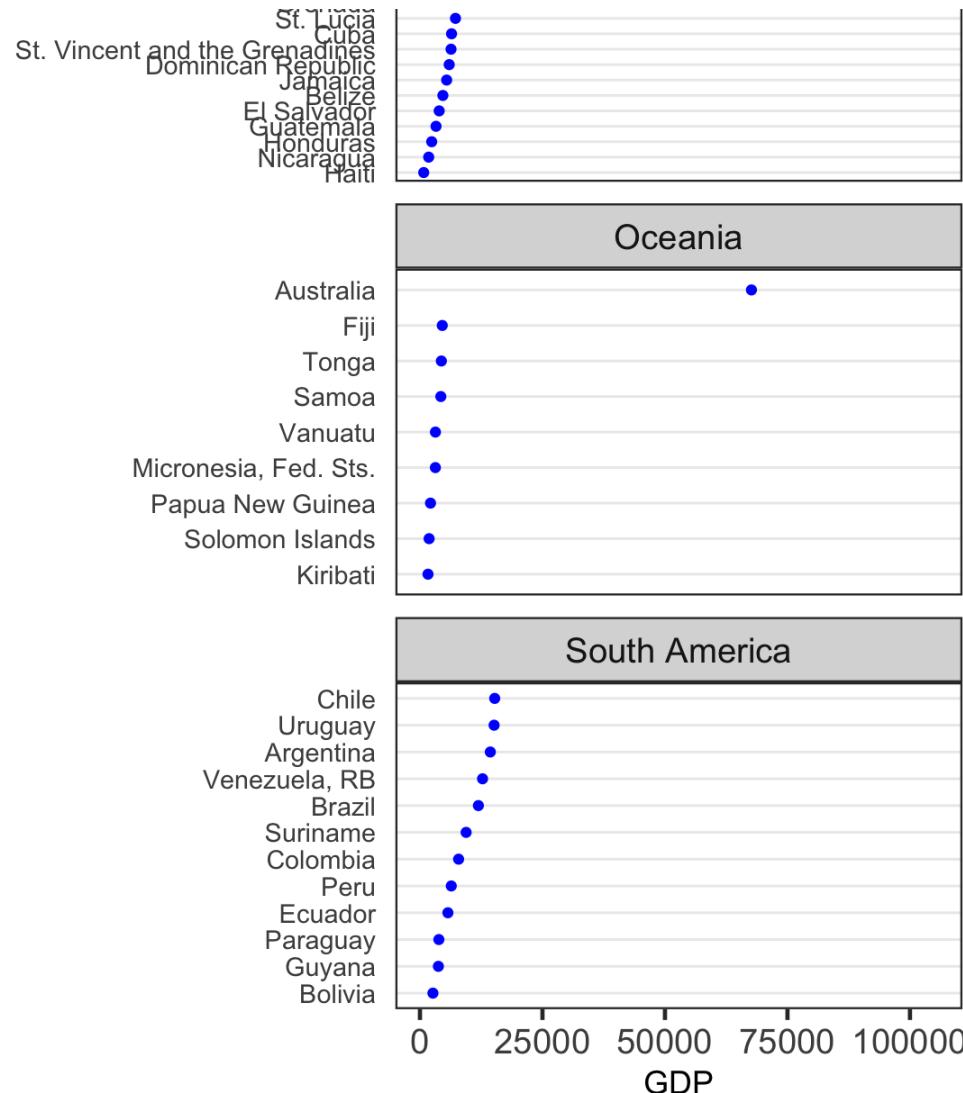




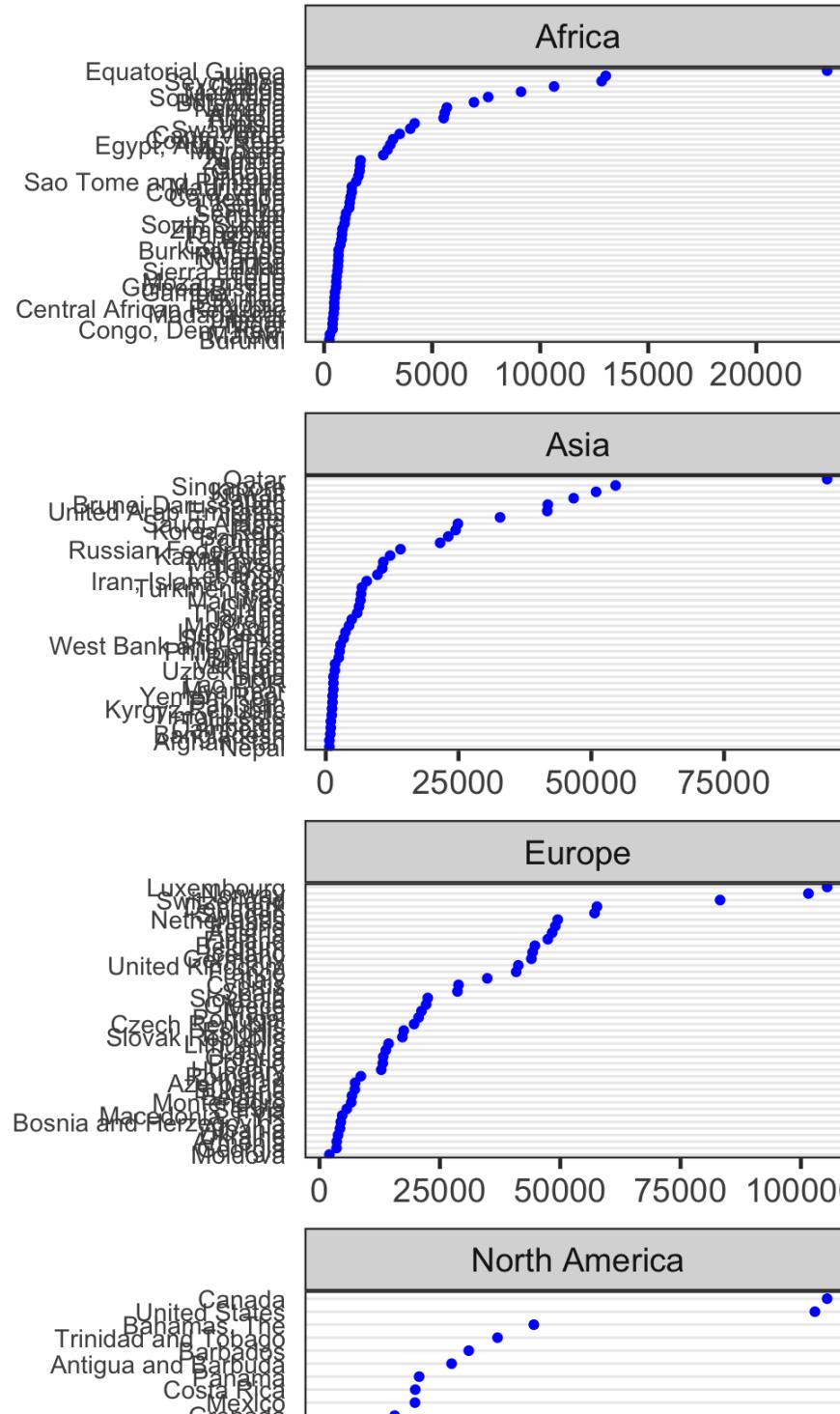


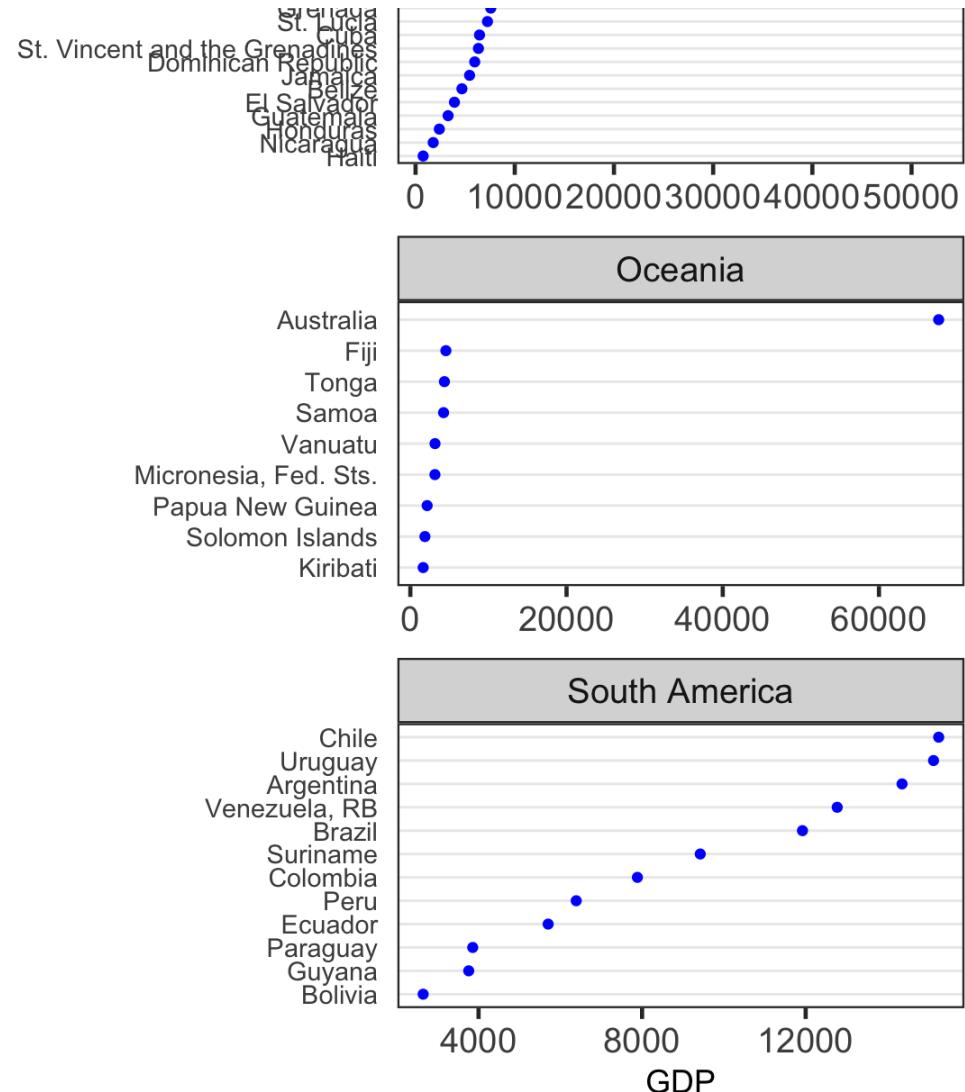
Cleveland Dot Plot with Facets scales = "free_y"



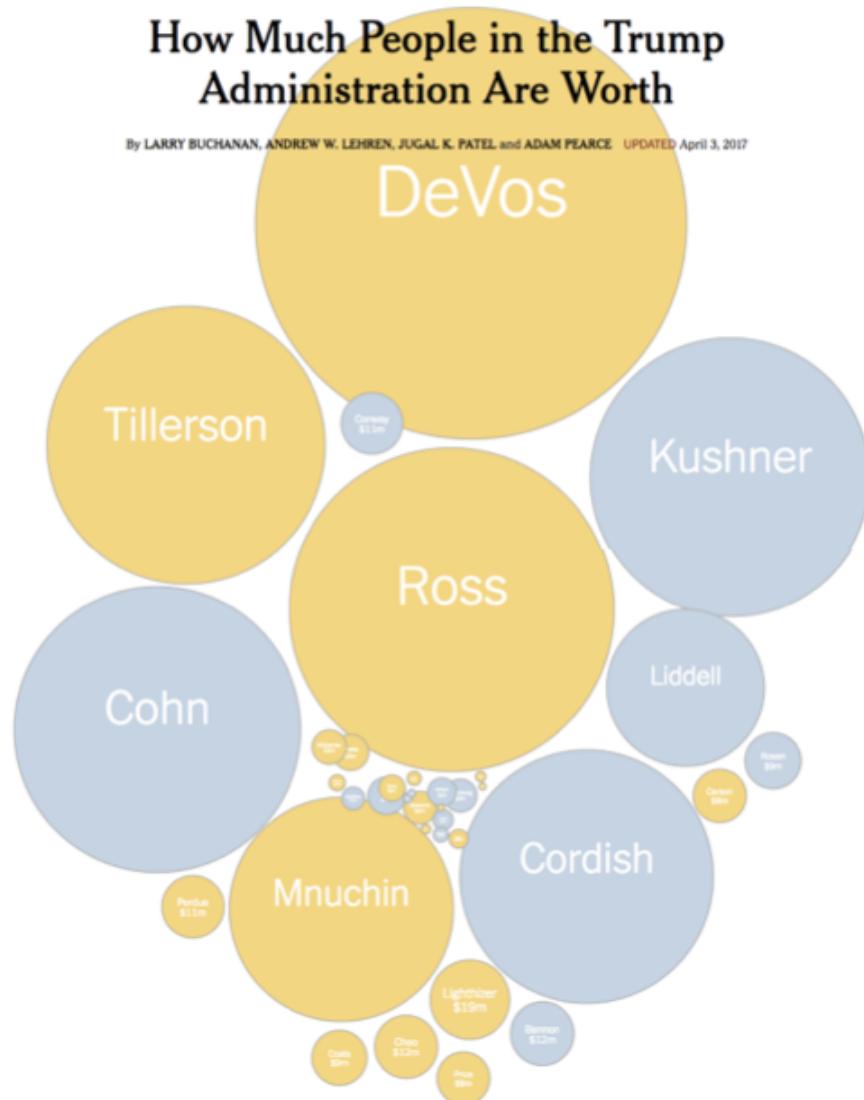


Cleveland Dot Plot with Facets scales = "free"



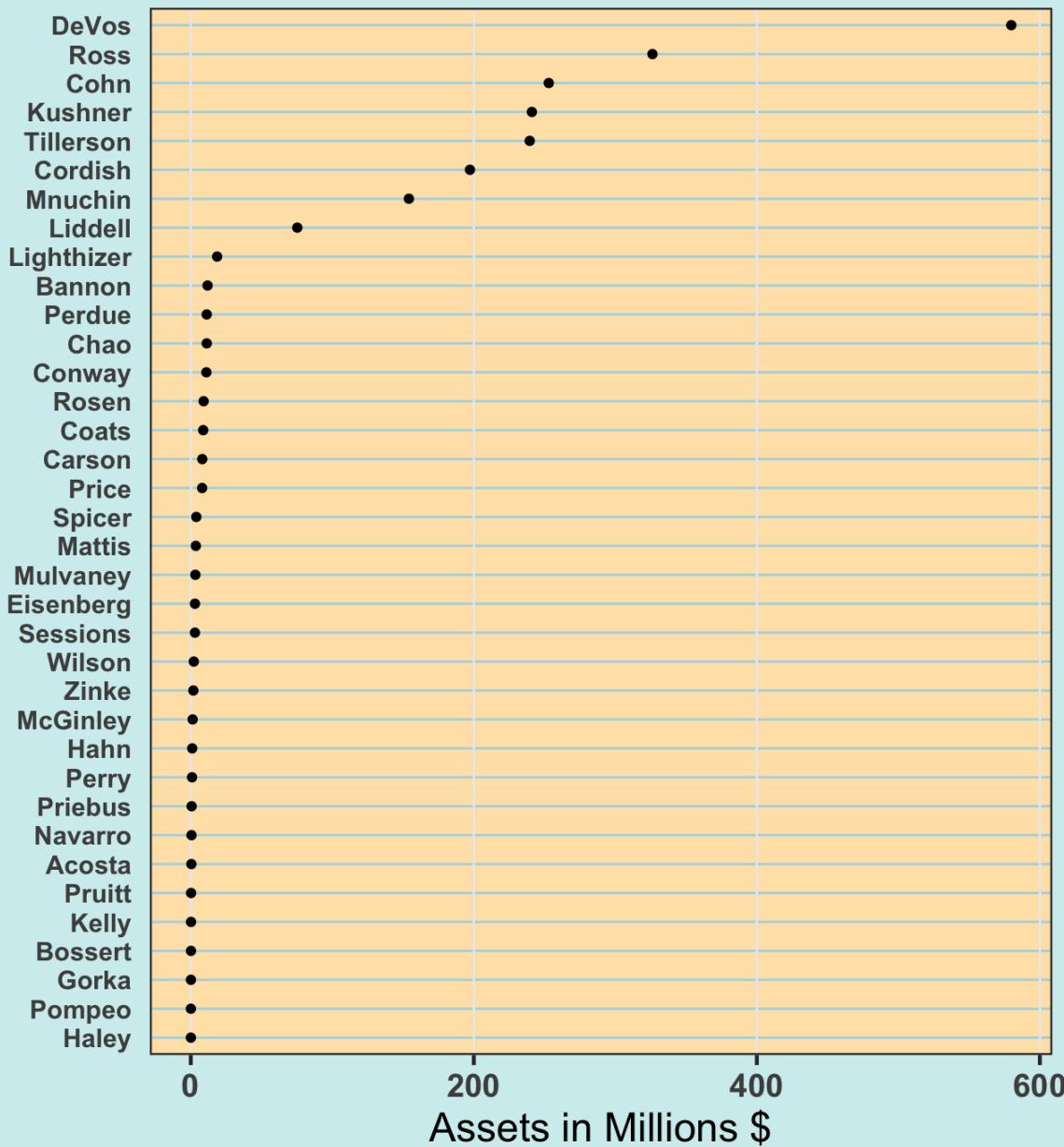


Trump Administration Assets

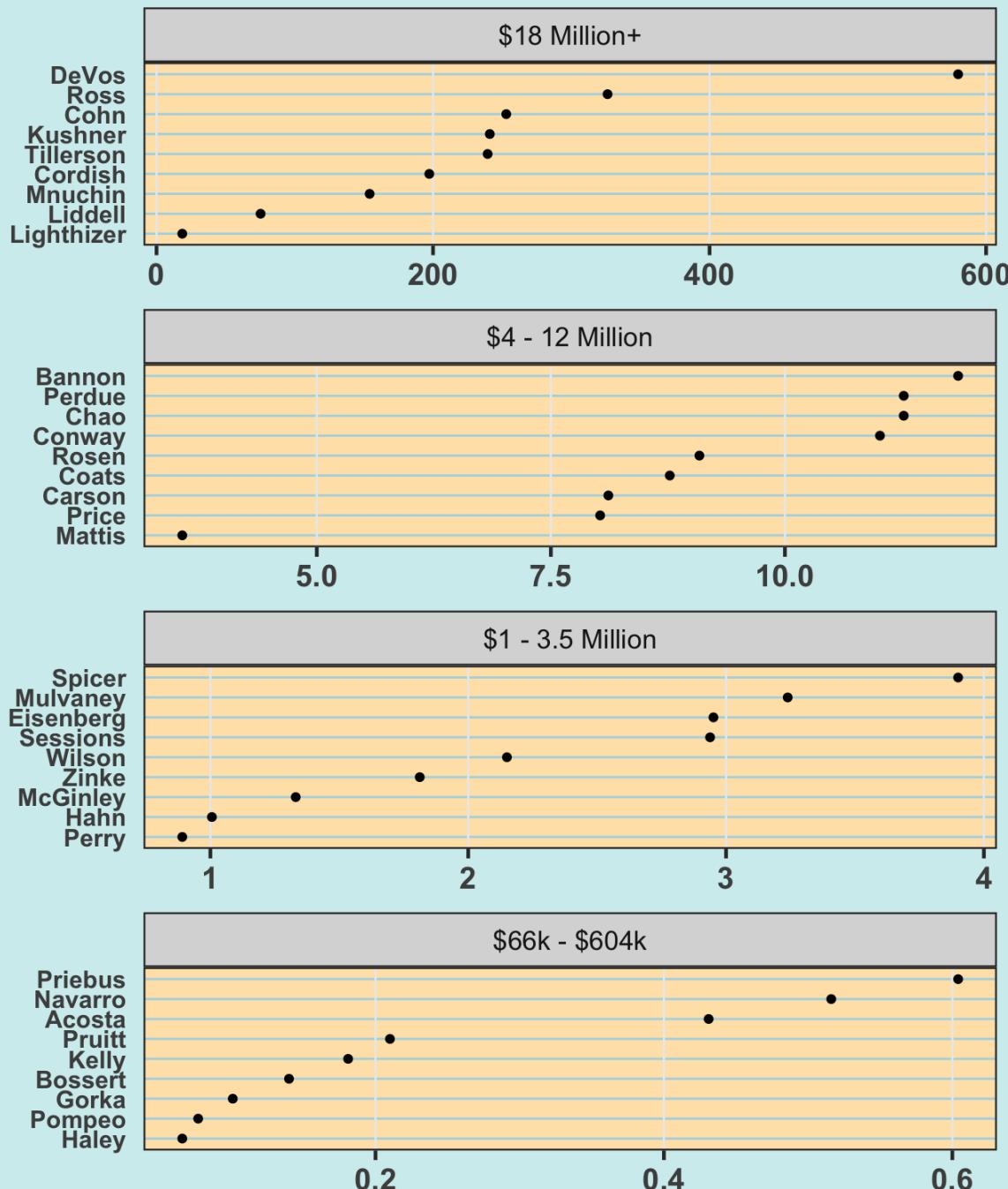


Redraw as Cleveland Dot Plot

How Much People in the Trump Administration Are Worth



How Much People in the Trump Administration Are Worth



Assets in Millions \$

Recoding factor levels

Keep a trail of breadcrumbs

```
x <- factor(c("G234", "G452", "G136"))
levels(x)
```

```
## [1] "G136" "G234" "G452"
```

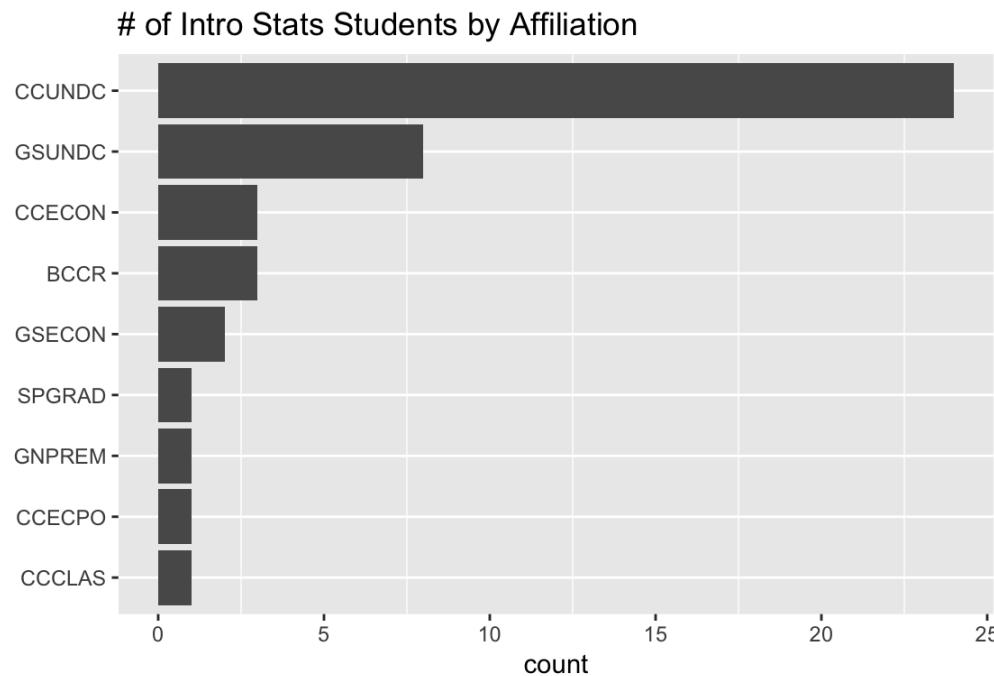
```
levels(x) <- c("Physics", "Math", "Chemistry")
x
```

```
## [1] Math      Chemistry Physics
## Levels: Physics Math Chemistry
```

```
# clear connections between old and new names, old version is preserved
x <- factor(c("G234", "G452", "G136"))
y <- fct_recode(x, Physics = "G234", Math = "G452", Chemistry = "G136")
y
```

```
## [1] Physics  Math      Chemistry
## Levels: Chemistry Physics Math
```

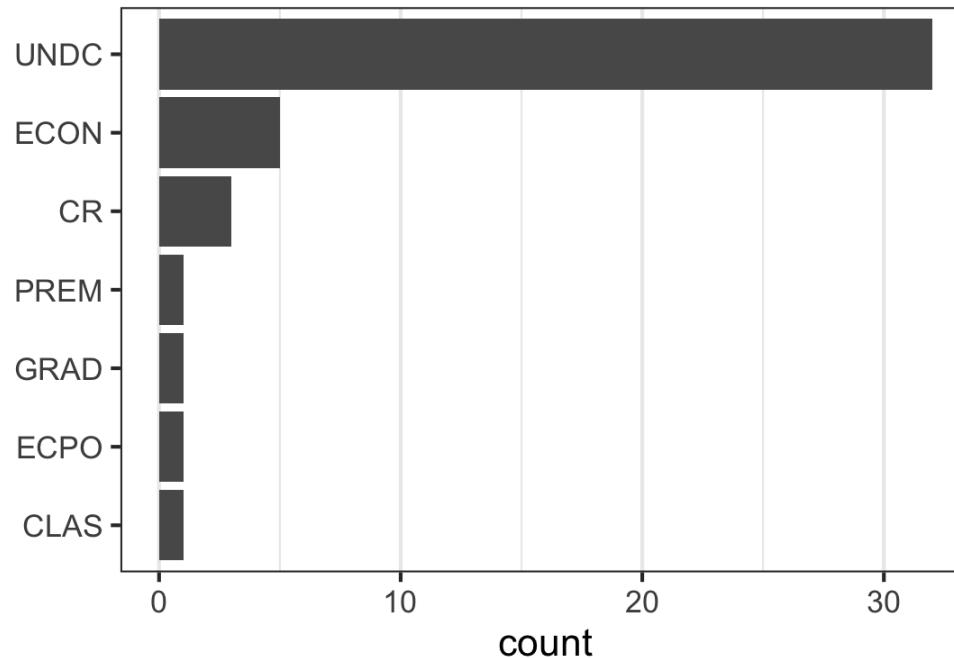
Cleaning data



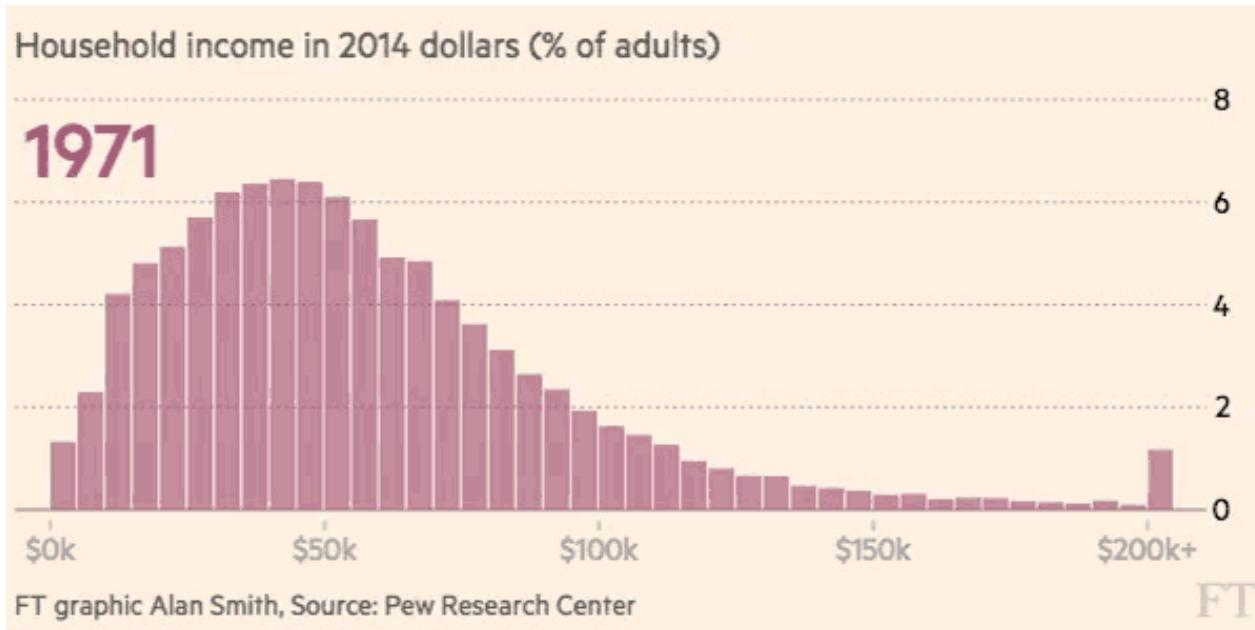
What's the problem?

Major only

Remove school from affiliation



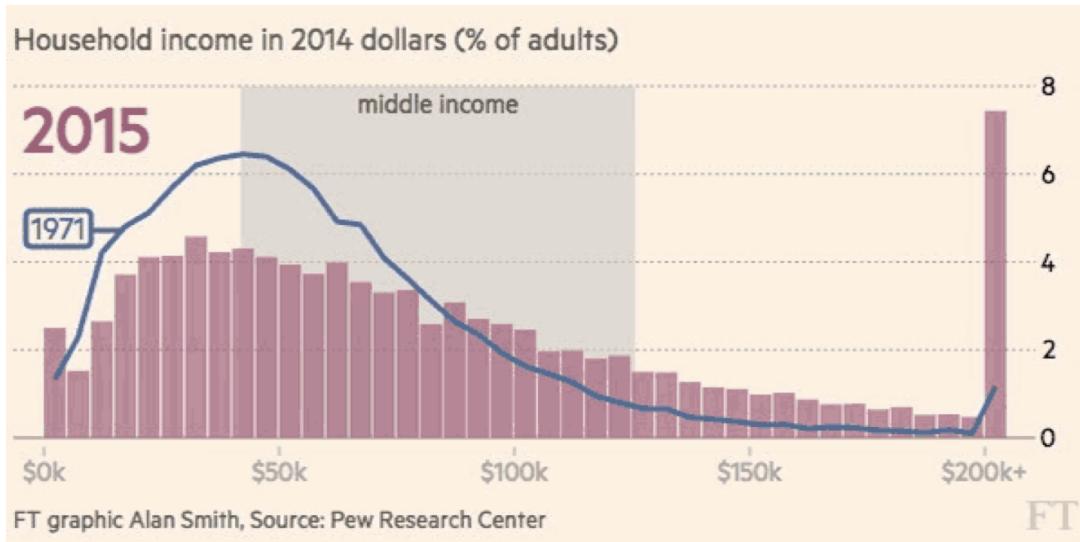
“Other” “or more” categories: “topcoding”



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”

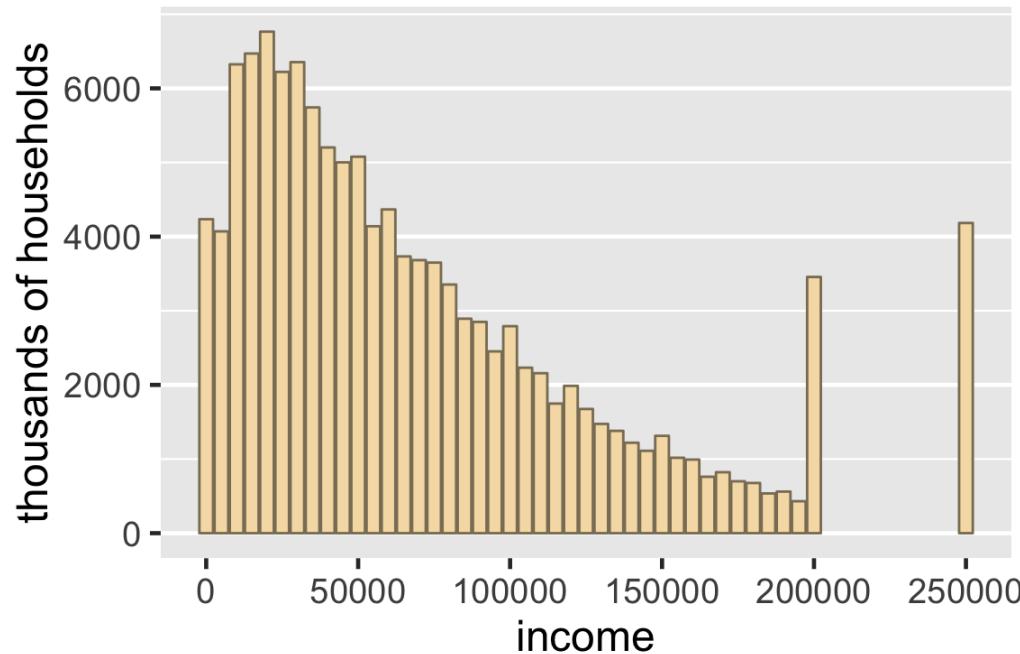
<http://www.latimes.com/business/hiltzik/la-fi-hiltzik-ft-graphic-20160320-snap-htmlstory.html>

“Other” “or more” categories

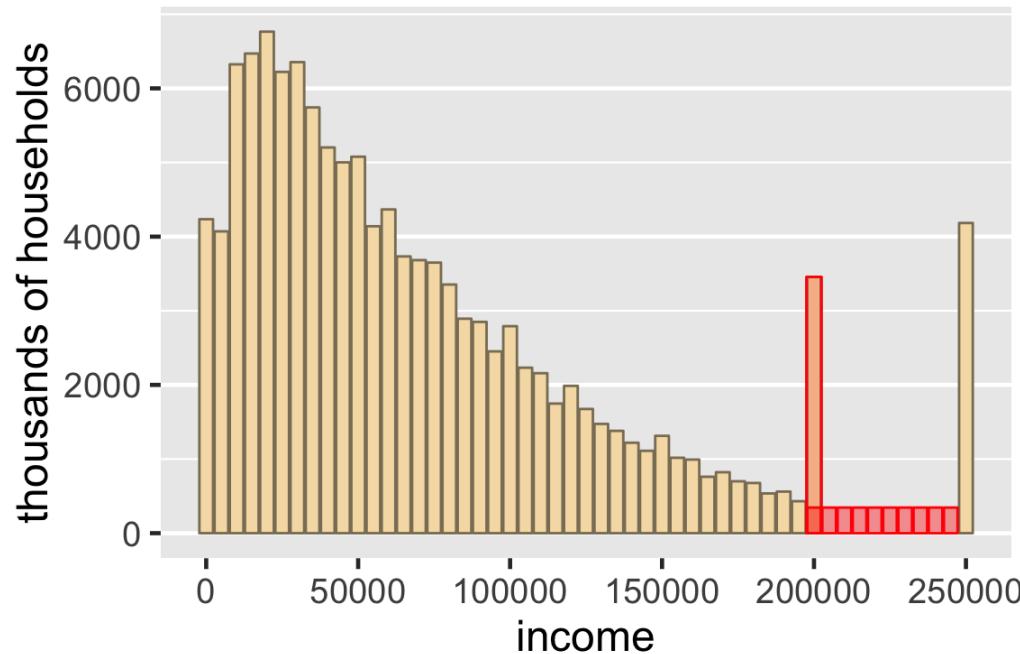
\$82,500 to \$84,999	\$85,000 to \$87,499	\$87,500 to \$89,999	\$90,000 to \$92,499	\$92,500 to \$94,999	\$95,000 to \$97,499	\$97,500 to \$99,999	\$100,000 and over	Va (D)
1,102	1,683	892	2,065	894	1,306	770	22,426	
\$82,500 to \$84,999	\$85,000 to \$87,499	\$87,500 to \$89,999	\$90,000 to \$92,499	\$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.xls

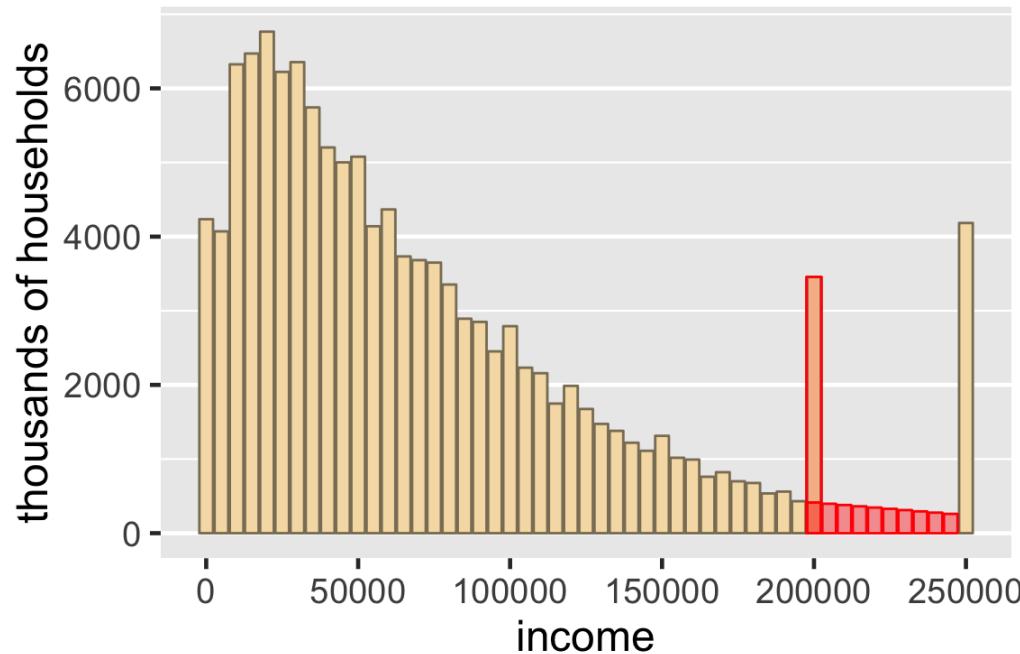
Household Income in 2015



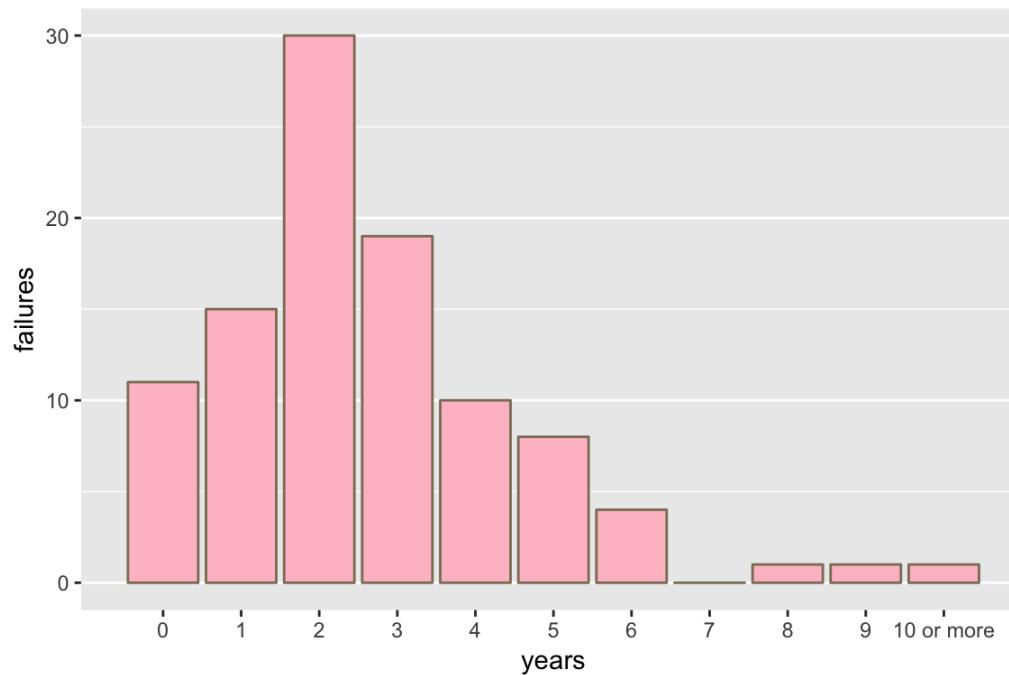
Household Income in 2015



Household Income in 2015



Reasonable use of “or more”



Data cleaning / transforming

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