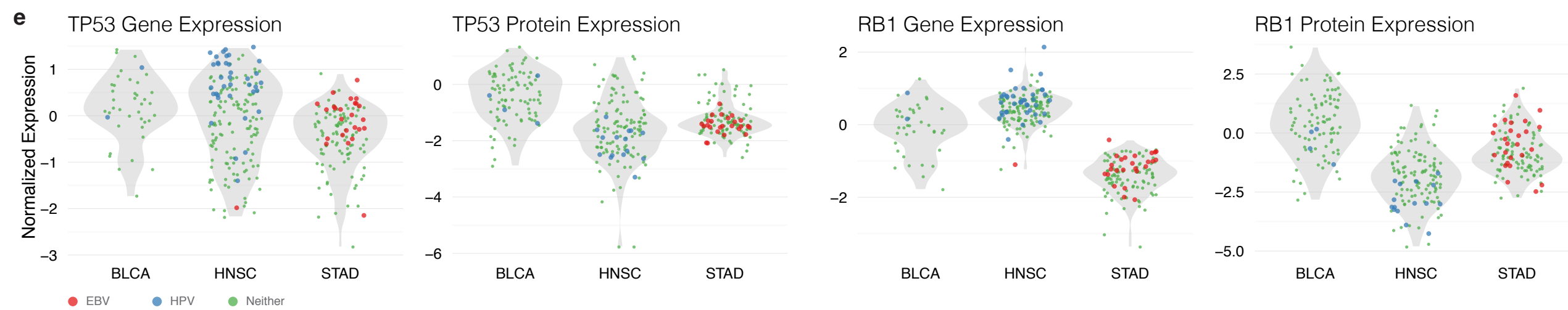
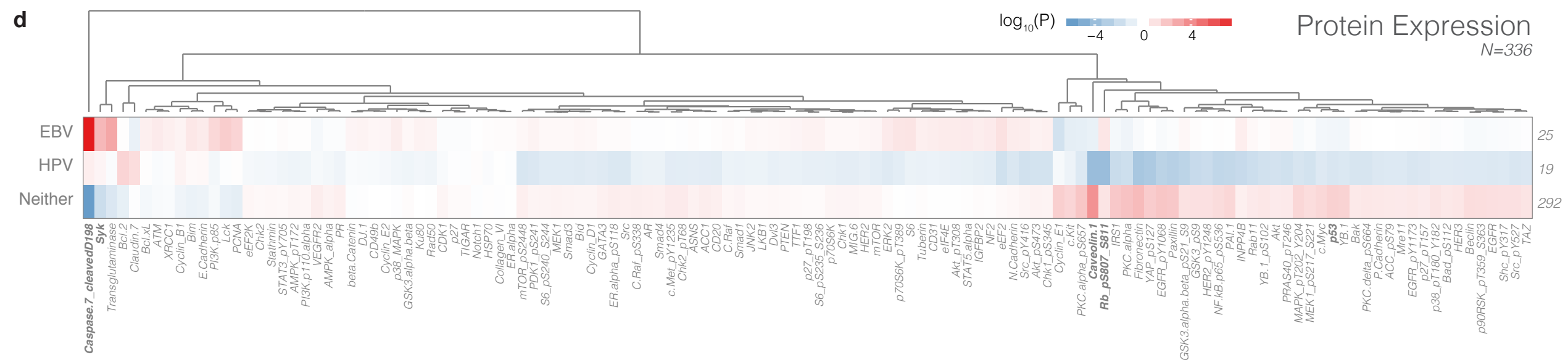
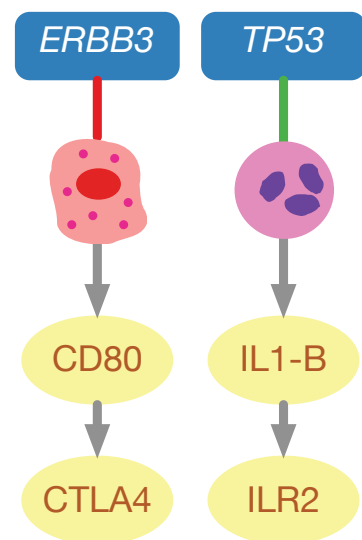


# Data Visualization with ggplot2

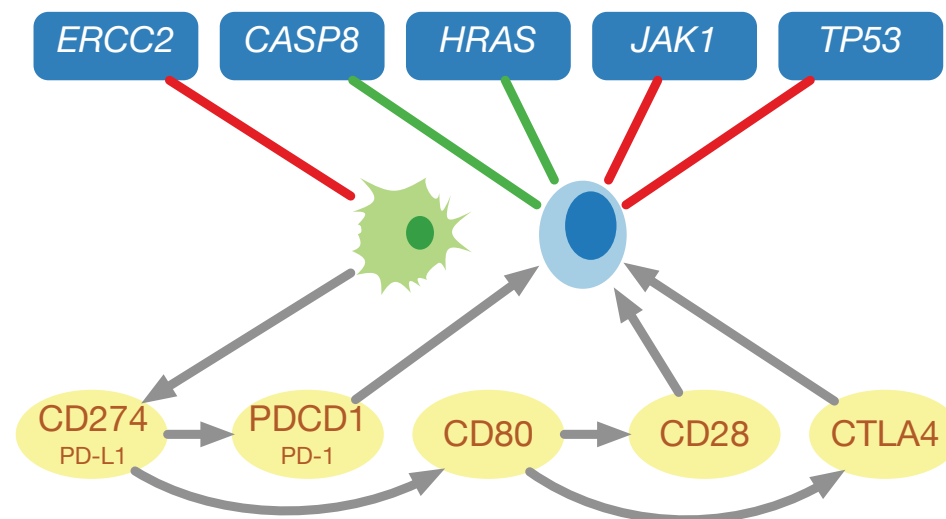
Matt Wyczalkowski  
February 8, 2018



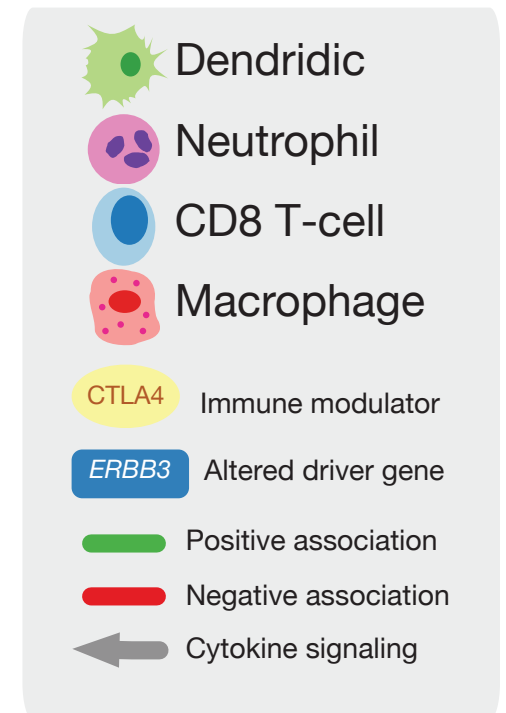
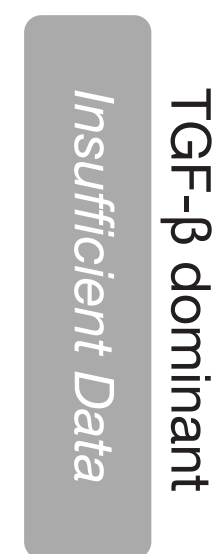
C1: Wound healing



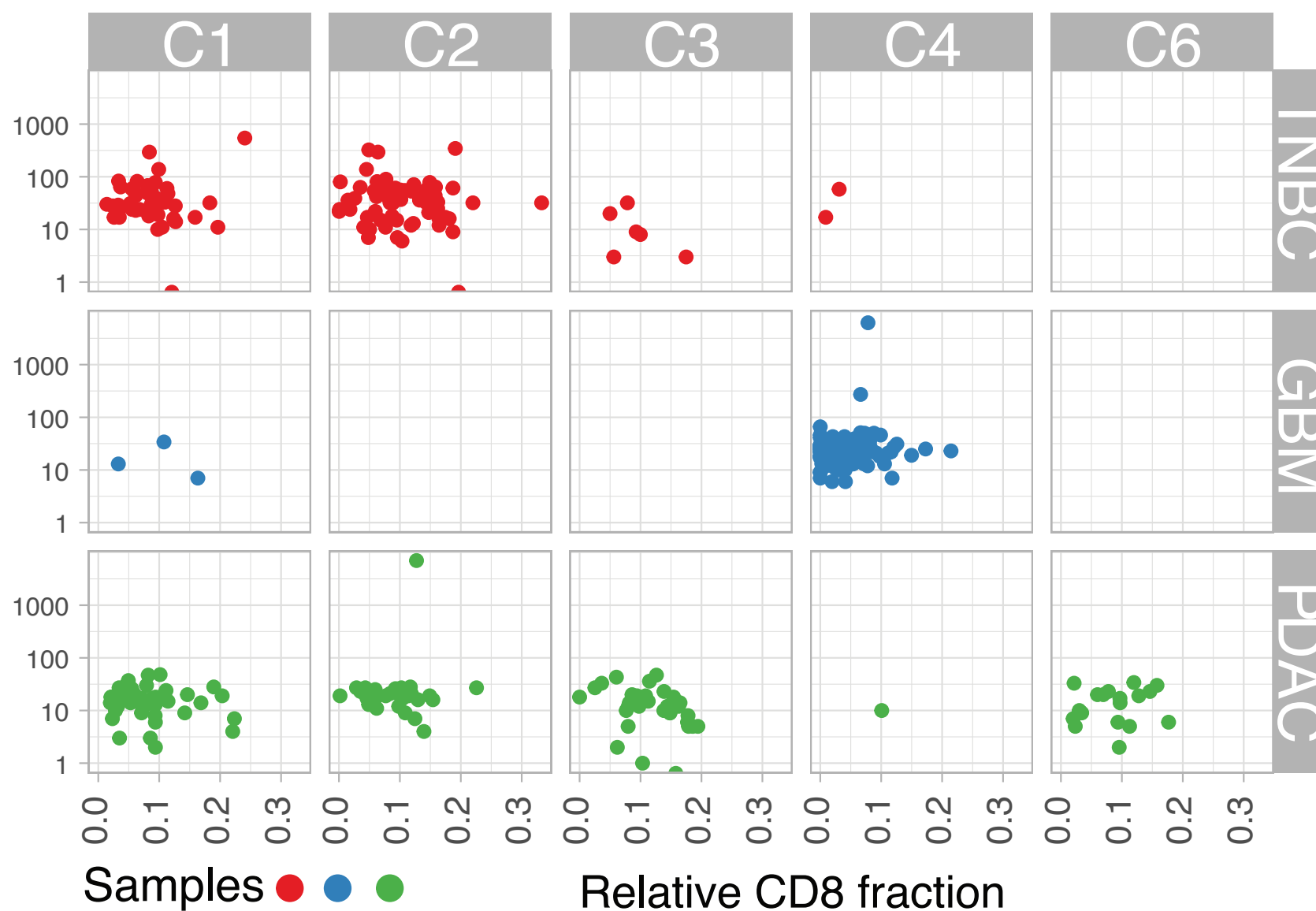
C2: IFN-γ dominant



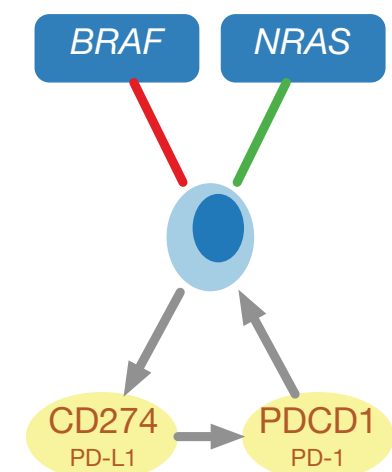
C6



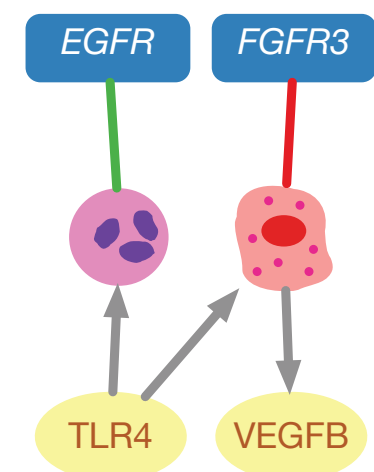
Number of immunogenic mutations

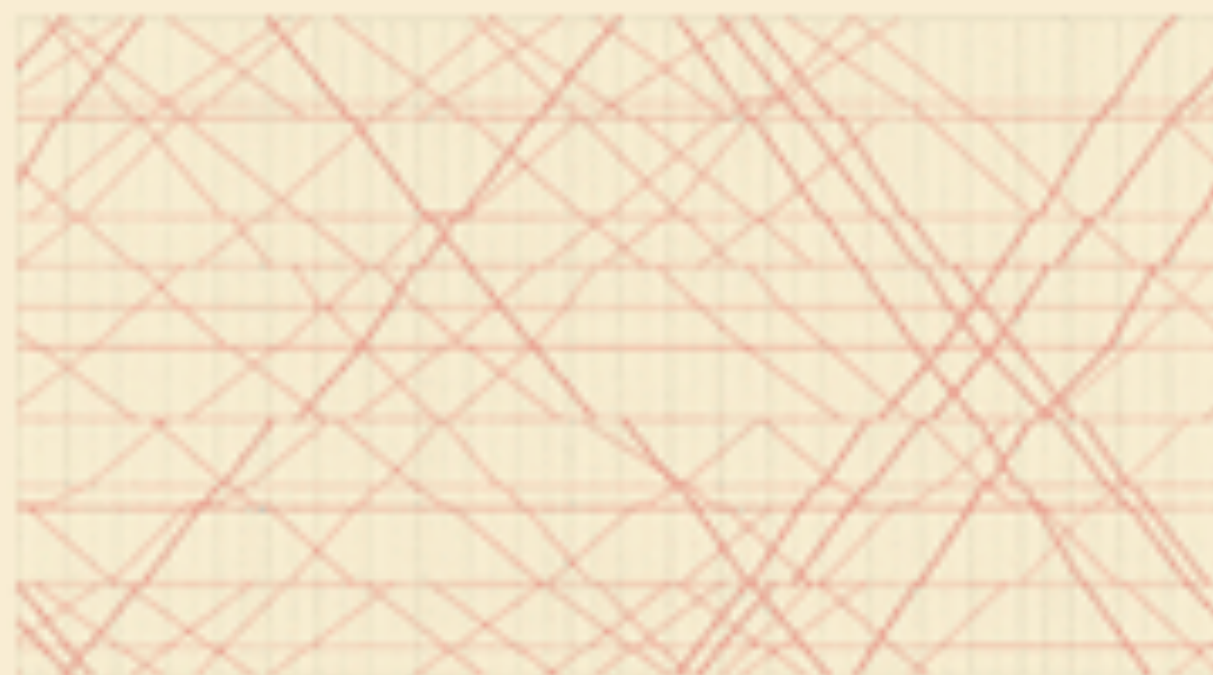


C3: Inflammatory



C4: Lymphocyte depleted





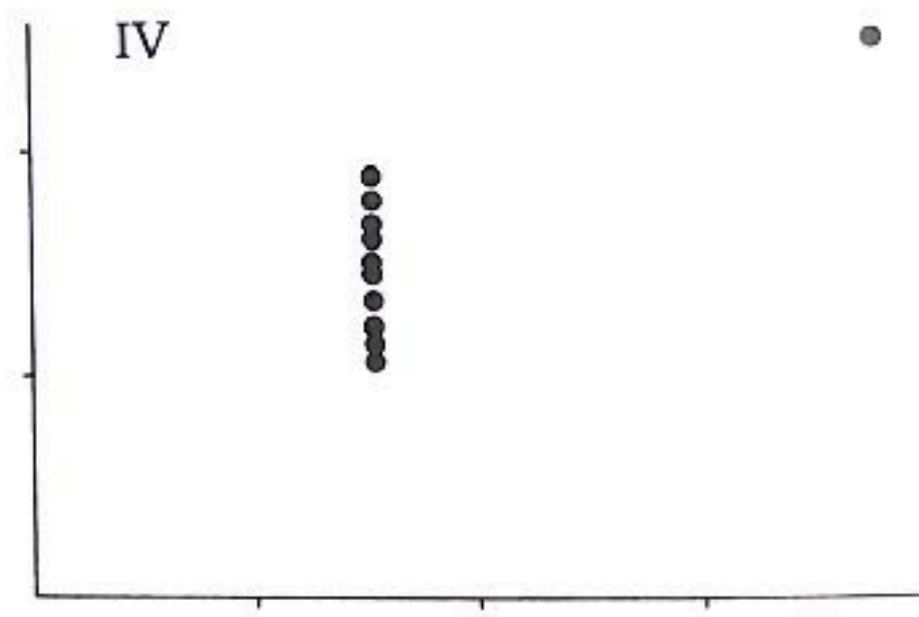
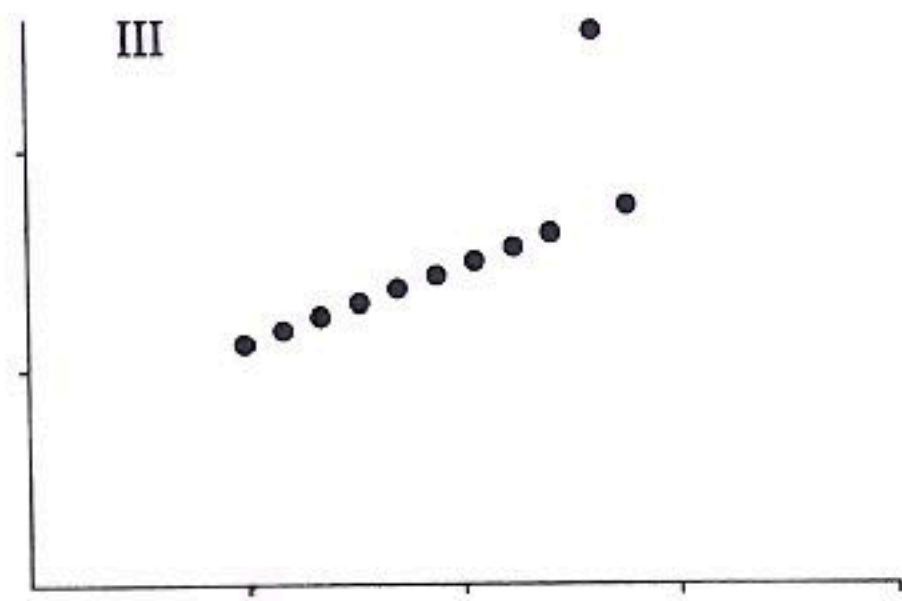
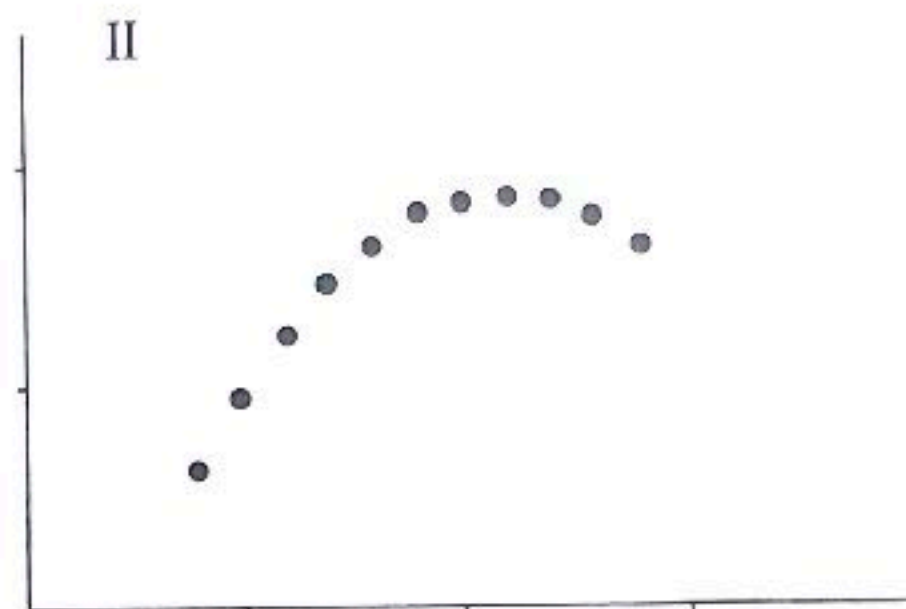
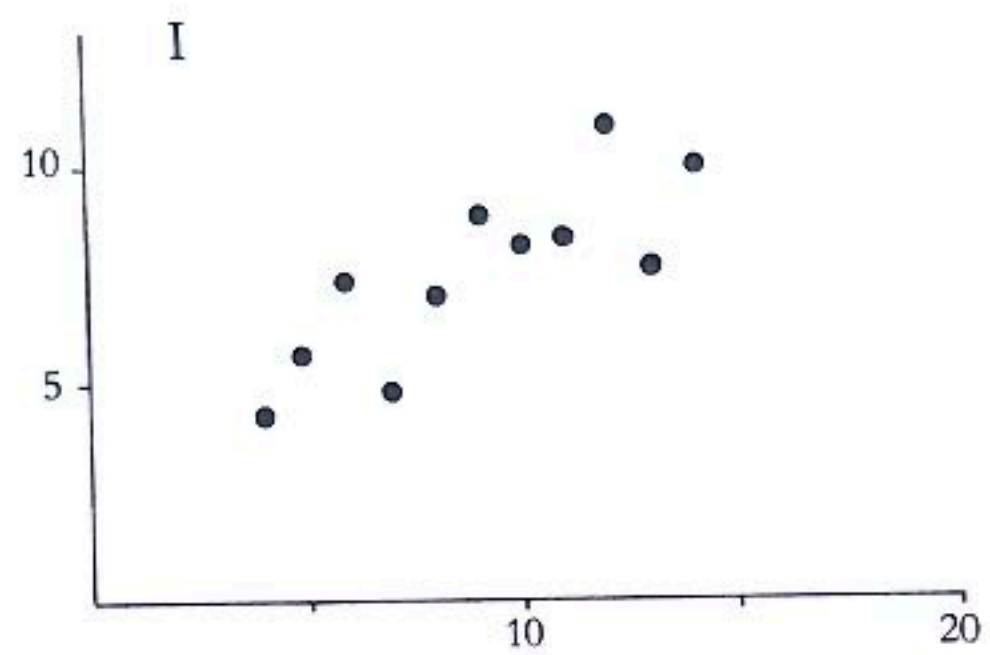
**SECOND EDITION**

The Visual Display  
of Quantitative Information

EDWARD R. TUFTE

I		II		III		IV	
X	Y	X	Y	X	Y	X	Y
10.0	8.04	10.0	9.14	10.0	7.46	8.0	6.58
8.0	6.95	8.0	8.14	8.0	6.77	8.0	5.76
13.0	7.58	13.0	8.74	13.0	12.74	8.0	7.71
9.0	8.81	9.0	8.77	9.0	7.11	8.0	8.84
11.0	8.33	11.0	9.26	11.0	7.81	8.0	8.47
14.0	9.96	14.0	8.10	14.0	8.84	8.0	7.04
6.0	7.24	6.0	6.13	6.0	6.08	8.0	5.25
4.0	4.26	4.0	3.10	4.0	5.39	19.0	12.50
12.0	10.84	12.0	9.13	12.0	8.15	8.0	5.56
7.0	4.82	7.0	7.26	7.0	6.42	8.0	7.91
5.0	5.68	5.0	4.74	5.0	5.73	8.0	6.89

$N = 11$   
 mean of X's = 9.0  
 mean of Y's = 7.5  
 equation of regression line:  $Y = 3 + 0.5X$   
 standard error of estimate of slope = 0.118  
 $t = 4.24$   
 sum of squares  $X - \bar{X} = 110.0$   
 regression sum of squares = 27.50  
 residual sum of squares of Y = 13.75  
 correlation coefficient = .82  
 $r^2 = .67$

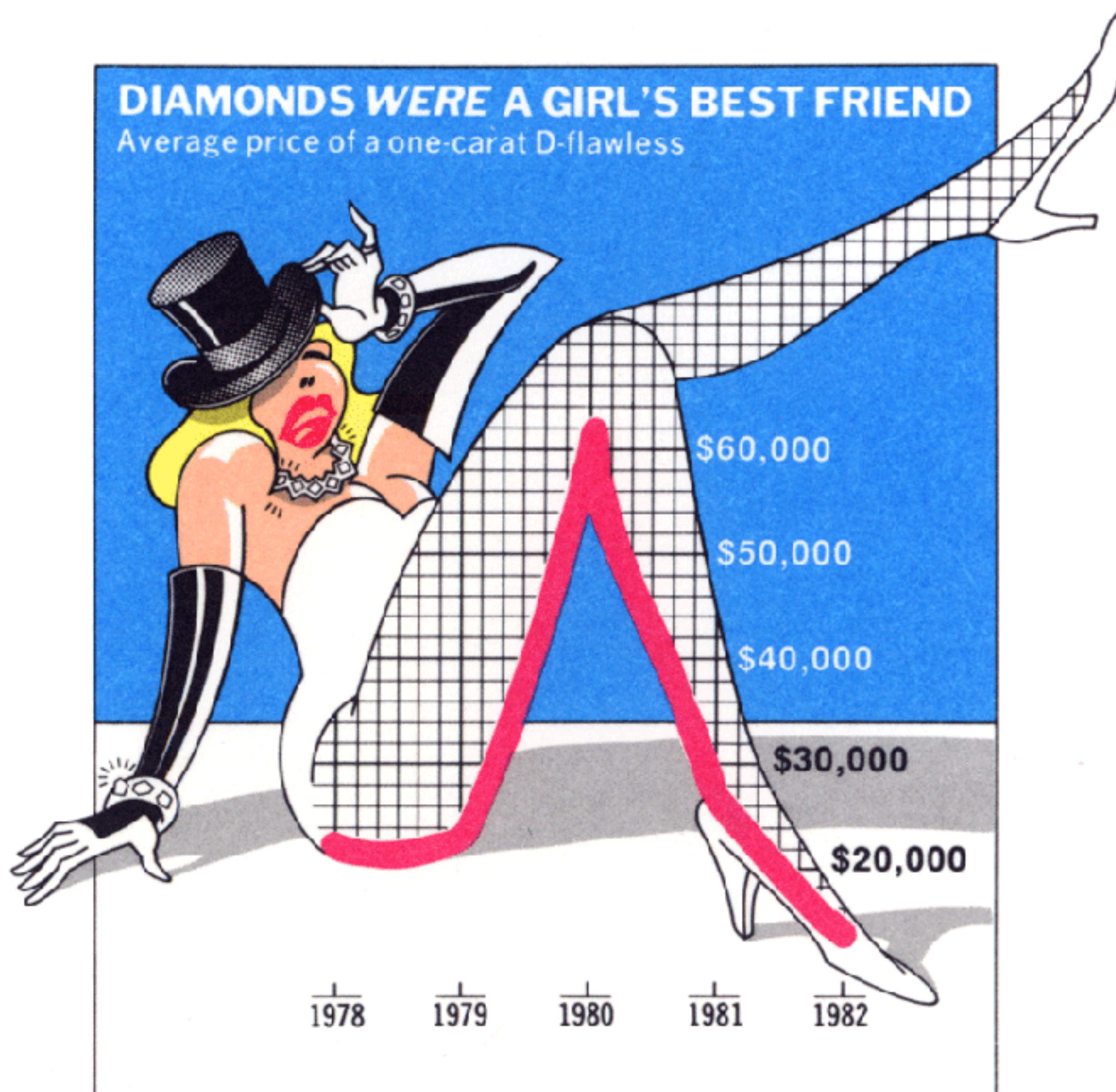


# Graphics Reveal Data

- Graphical displays should...
  - Show the data
  - Avoid distortion
  - Induce viewer to think about substance
  - Make large datasets coherent
  - Encourage eye to compare different pieces of data

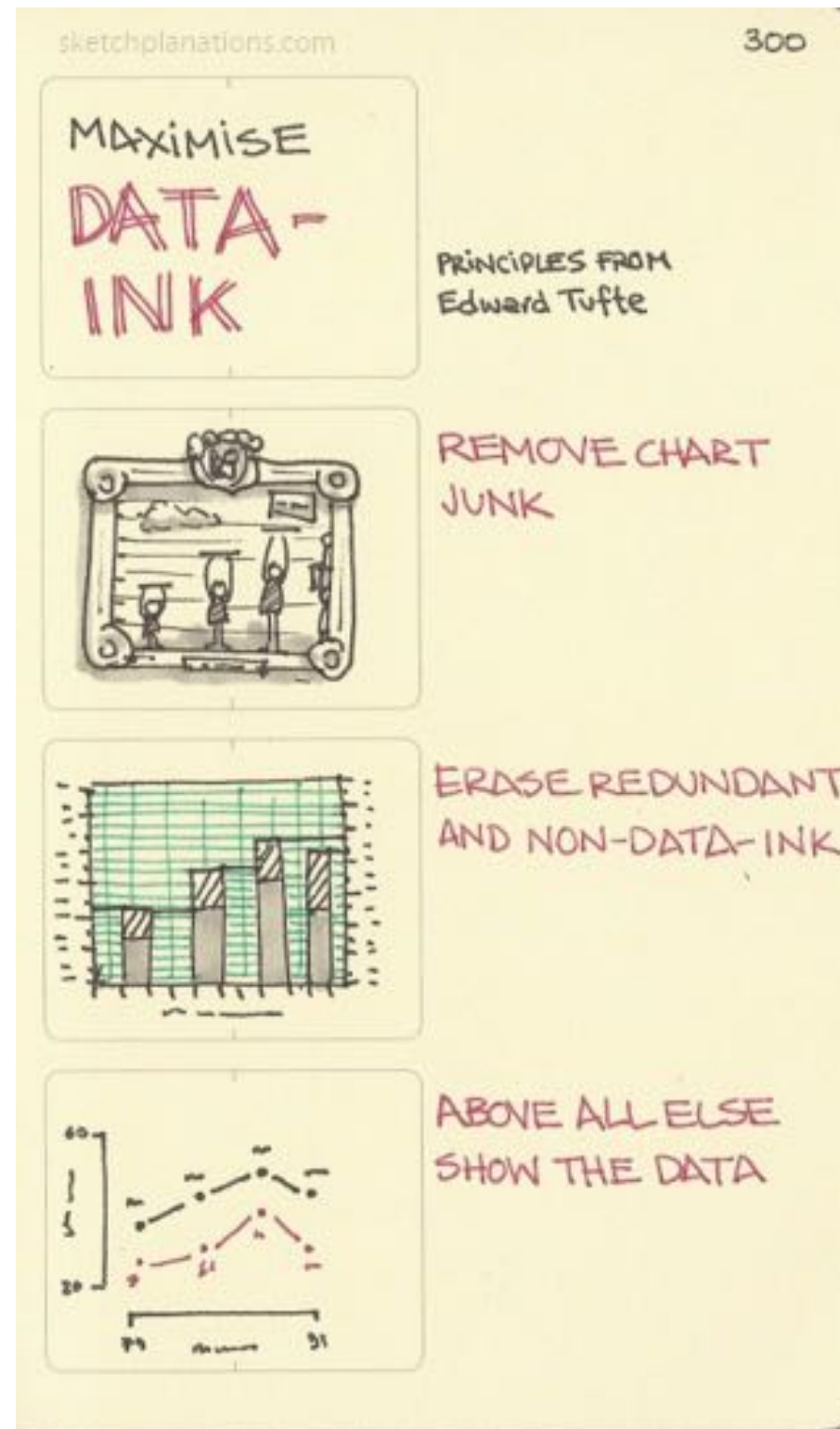


# Chartjunk





# Maximize Data-ink Ratio



# Grammar of Graphics

- A rational way of thinking about and making data graphics
- All plots are composed of:
  - Data
  - Layers / Geoms
  - Scales
  - Coord
  - Facet
  - Theme
- ggplot2 is an R-language implementation of Grammar of Graphics principles

# Data

- The data you want to visualize, and *aesthetic mappings* describing how variables in data map to attributes you can see
- Much of the work with ggplot2 is creating the data frame - the rest is easy!

# Layers and Geoms

- Layers are made of geometric elements (geoms)
  - `geom_point`
  - `geom_bar`
  - `geom_violin`
- Layers also have statistical transformations (e.g. binning or counting observations)

# Scales

- Scales map values in data to values of aesthetics
  - Color
  - Size
  - Shape
- Define contents of legend

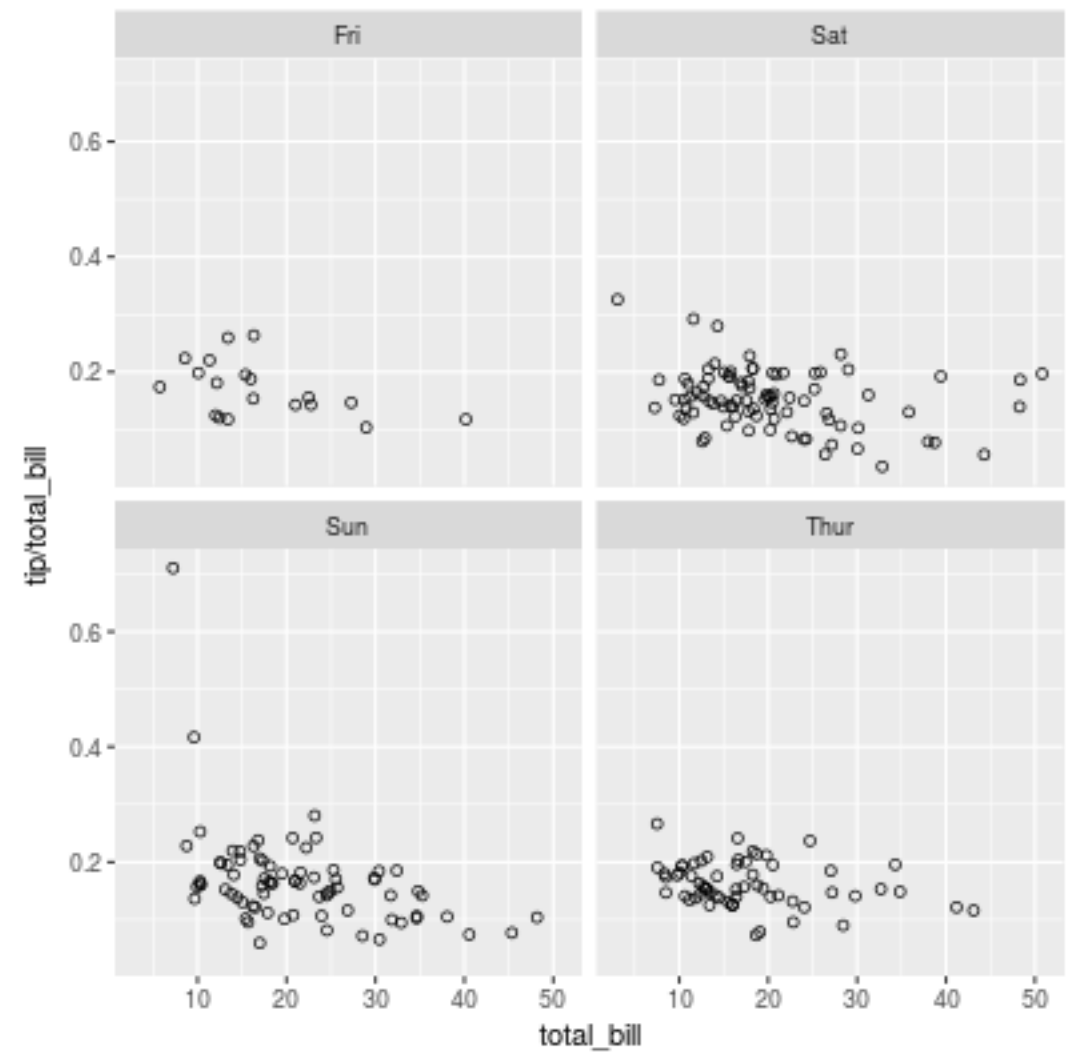


# Coord

- Coordinate system describes how data coordinates are mapped to the plane of the graphic.
  - Cartesian most common
  - Make pie charts by changing to polar coordinates

# Facets

- Breaks up data into subsets of small multiples
- Create an array of plots



# Themes

- Controls details of display
  - Size of fonts
  - Grid lines
  - Background color

# Example

```
p <- ggplot(data=BRFSS)
```

```
p <- p + geom_point(aes(x=height, y=weight),  
color="blue")
```

```
p
```

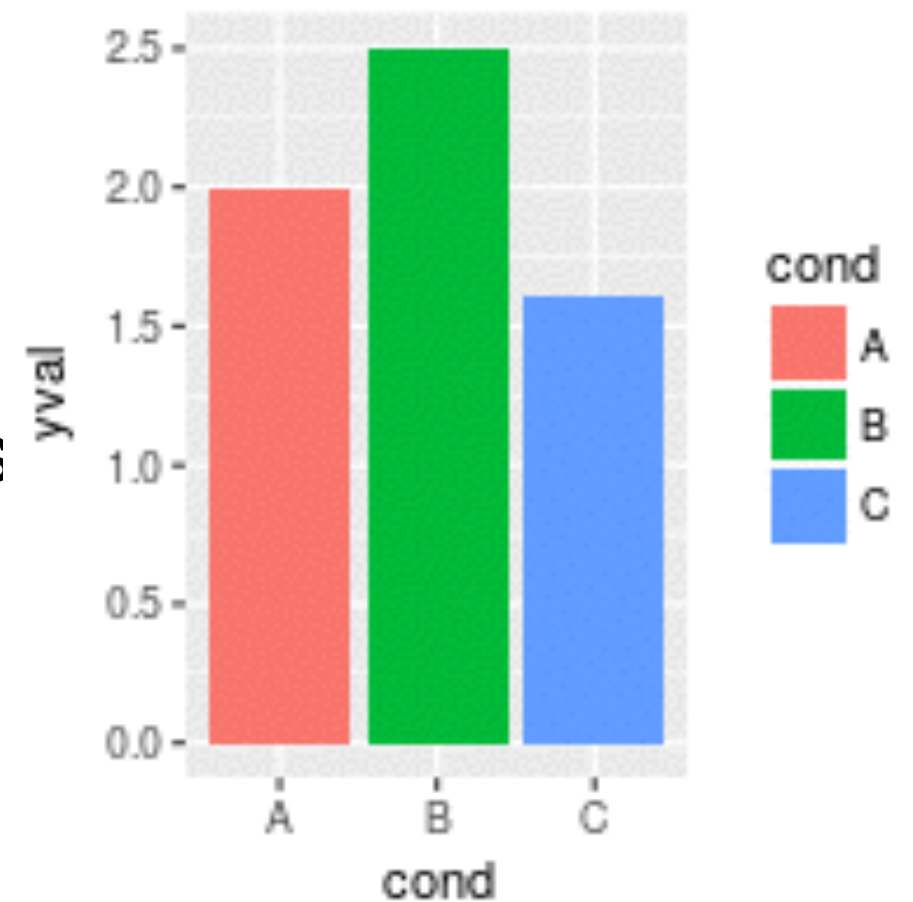
# Data: Categorical vs. Quantitative

- Basic distinction in types of data
- Categorical
  - Take on discrete values, e.g., names
- Quantitative (or continuous)
  - Take on continuous values, e.g., height
- A figure will show relationship between two or more variables
  - Type of figure to use is influenced by type of data



# Colors

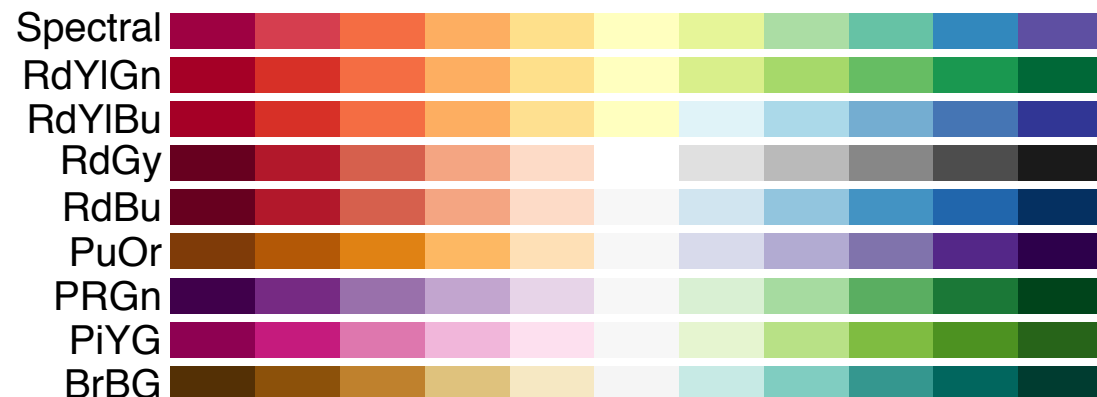
- ggplot2 has default color scheme
  - Relatively nice, but has drawbacks
- Controlled by setting scales



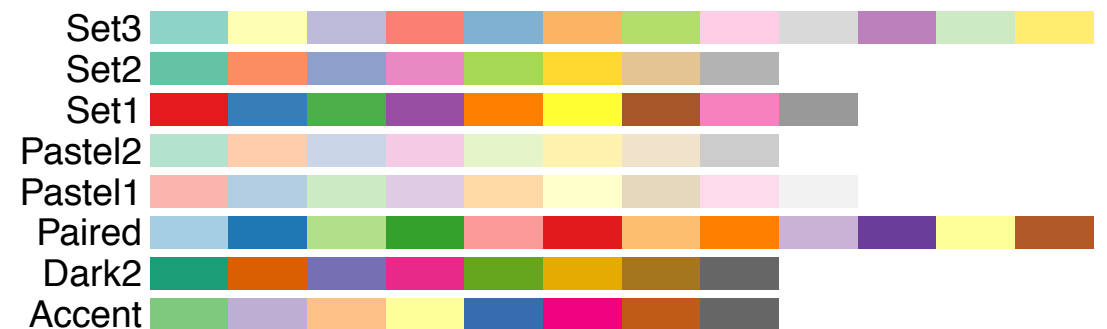
# RColorBrewer

- Attractive set of alternative color palettes
- Color-blind friendly, looks distinguishable in monochrome
- My go-to colors for illustrations and figures!

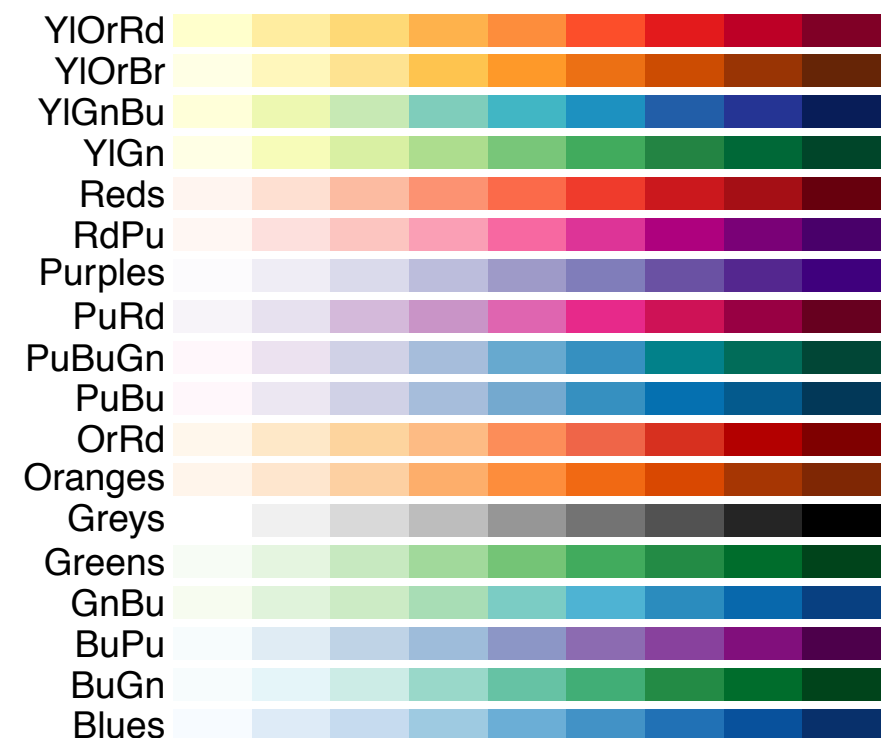
## Diverging



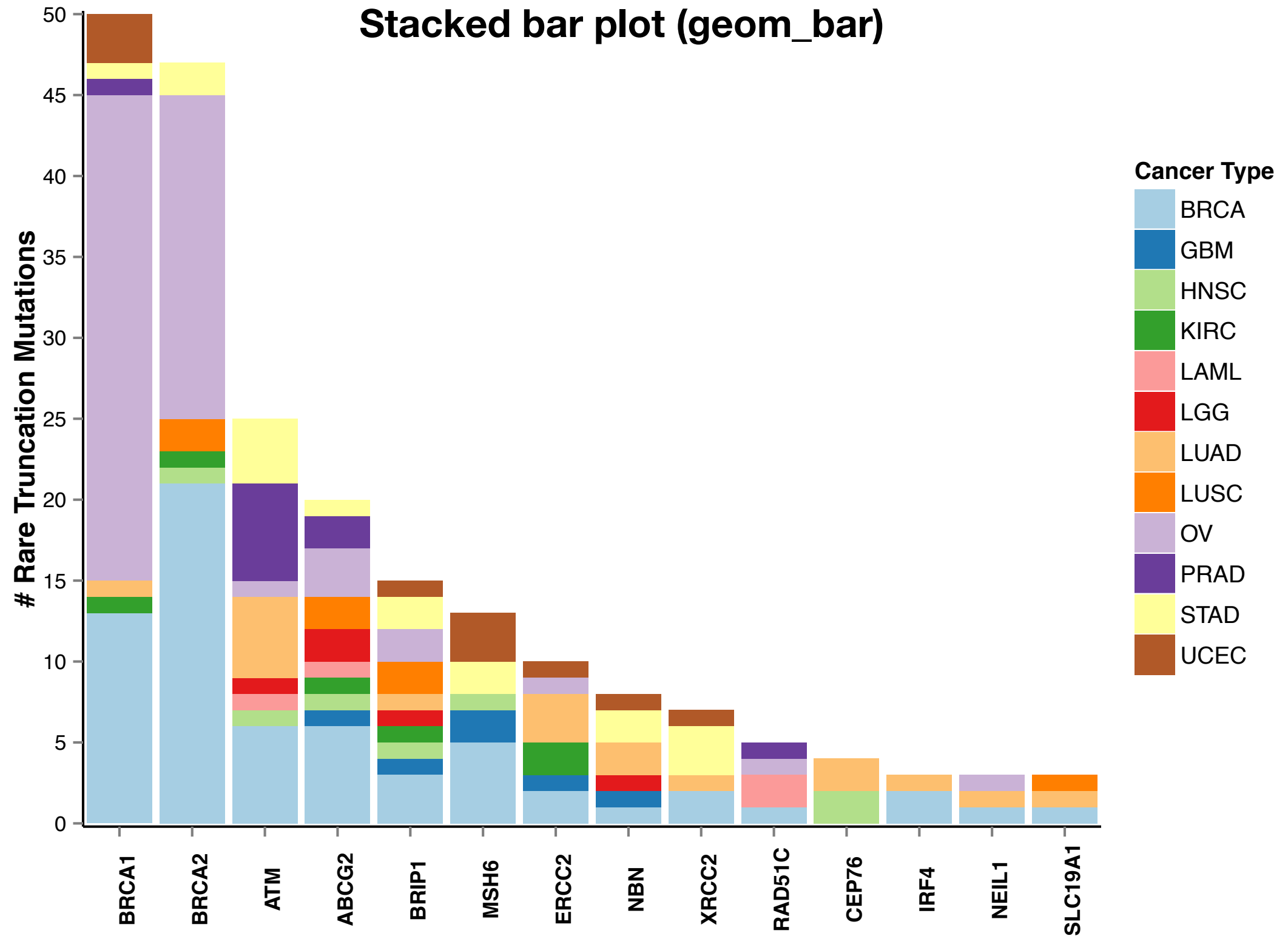
## Qualitative



## Sequential

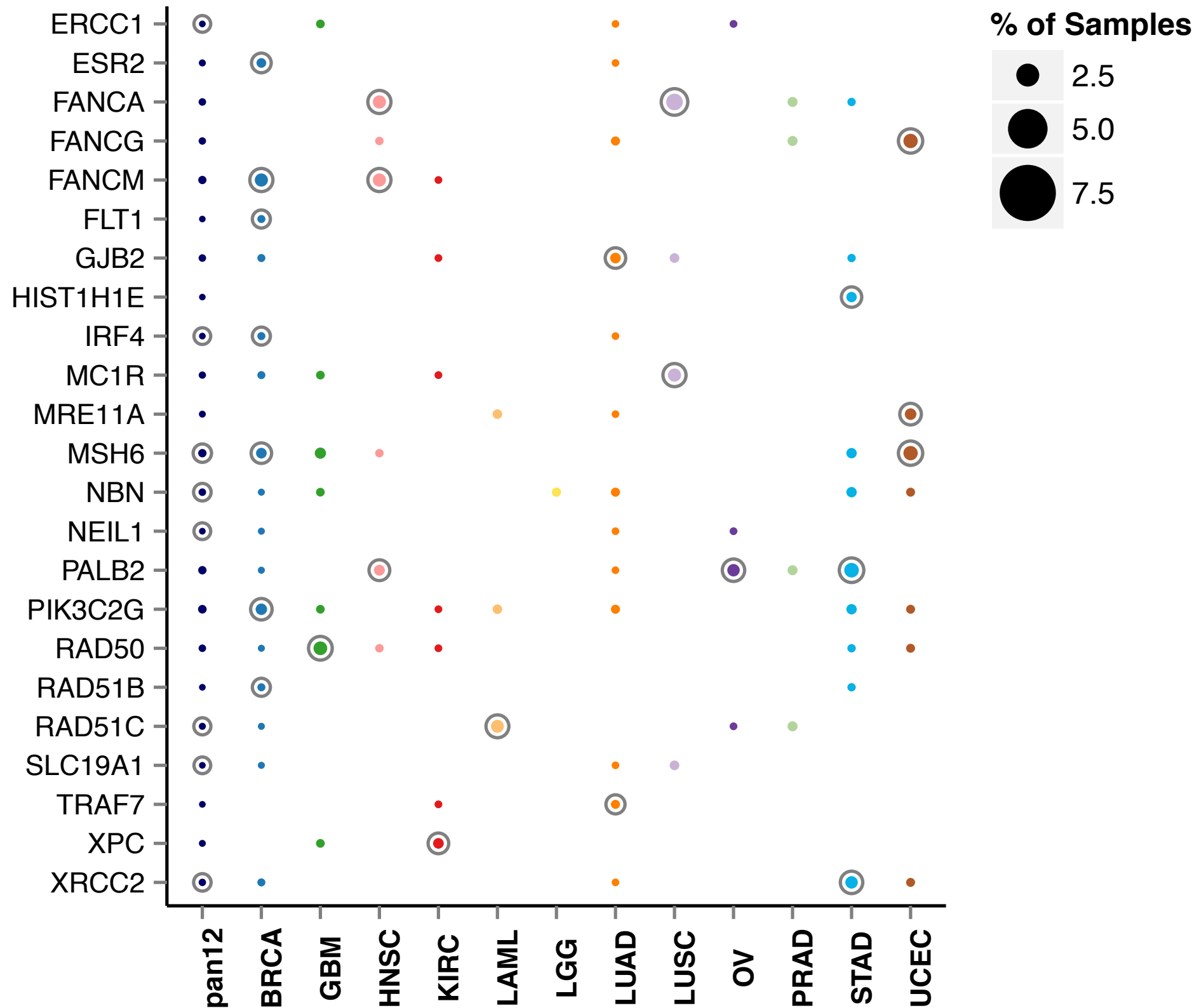


# Case Study: Before



# Case Study: After

Bubble plot (geom\_point)



# BRFSS Dataset

- 2016 BRFSS Survey Data and Documentation
  - Behavioral Risk Factor Surveillance System
  - [https://www.cdc.gov/brfss/annual\\_data/annual\\_2016.html](https://www.cdc.gov/brfss/annual_data/annual_2016.html)
- 486303 observations of 275 variables
  - 1.0 Gb in size
- Will subsample and select specific columns



# Subsampled Dataset:

## BRFSS.48K.csv

- 48,630 observations across 10 variables
  - state
  - employed
  - income
  - seatbelt
  - diabetes
  - height
  - weight
  - bmi
  - age
  - sex

# Getting more help

- ggplot2\_book.pdf
  - For sale on amazon
  - Available for free <https://github.com/hadley/ggplot2-book.git>
  - Included in class notes
- R Cookbook: <http://www.cookbook-r.com/Graphs/>
- Google