

Data Visualization

Spring 2019

Data Management

~ /> -

Spring 2019

<http://data880.co>

<http://socviz880.co>

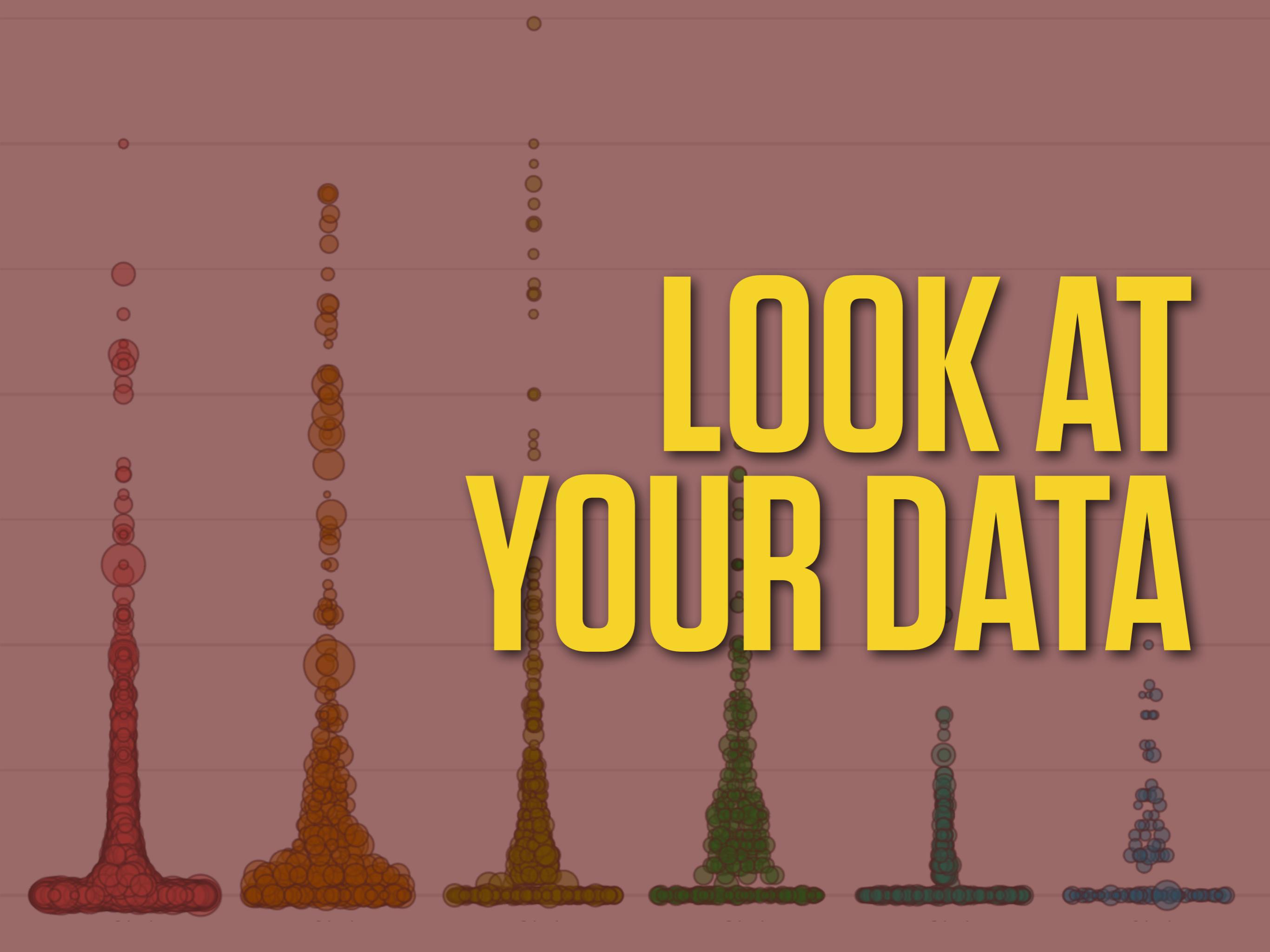
```
my_packages <- c("tidyverse", "fs", "devtools")
install.packages(my_packages)
```

```
devtools::install_github("kjhealy/socviz")
```

```
library(socviz)
```

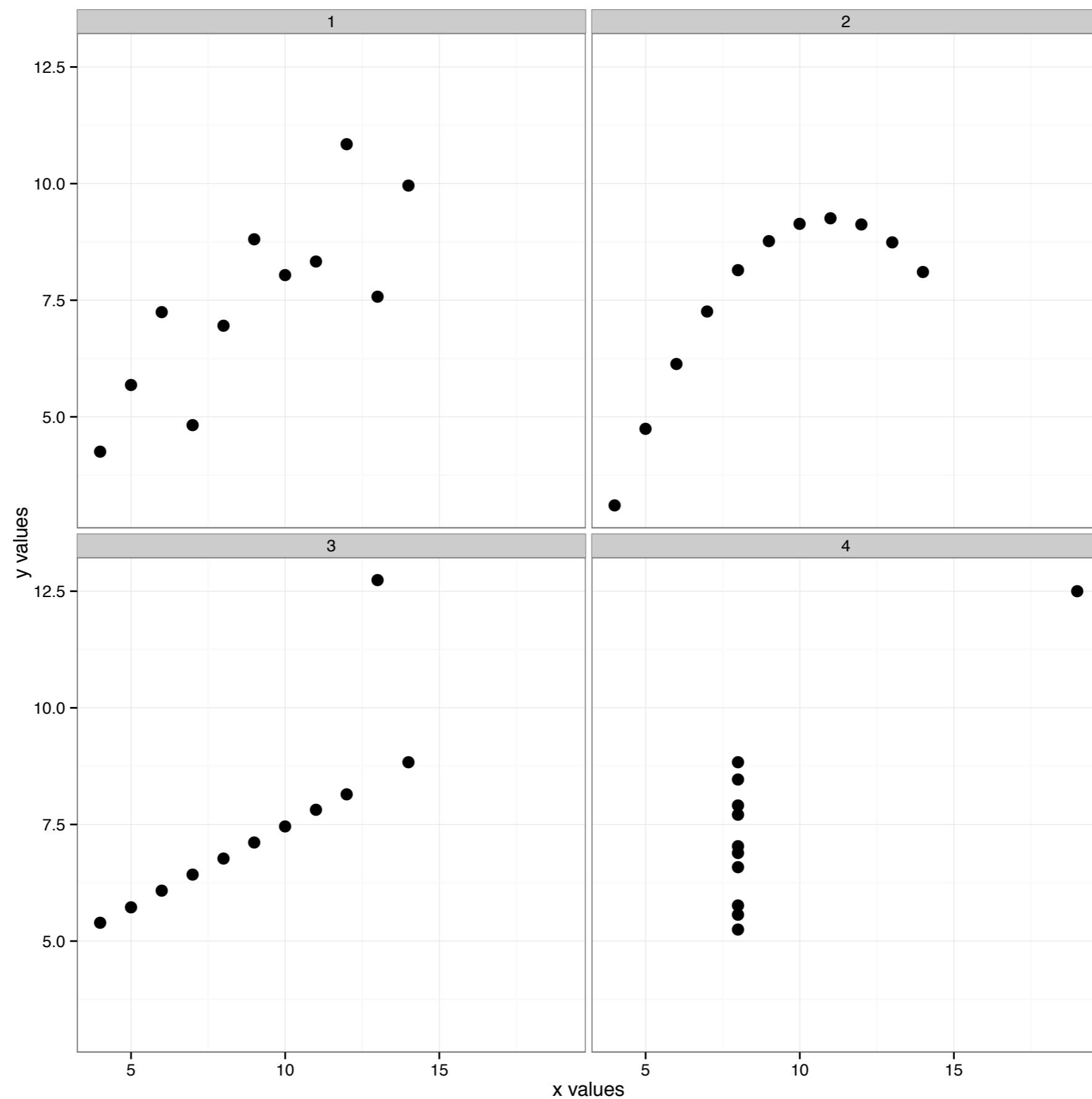
```
setup_course_notes()
```

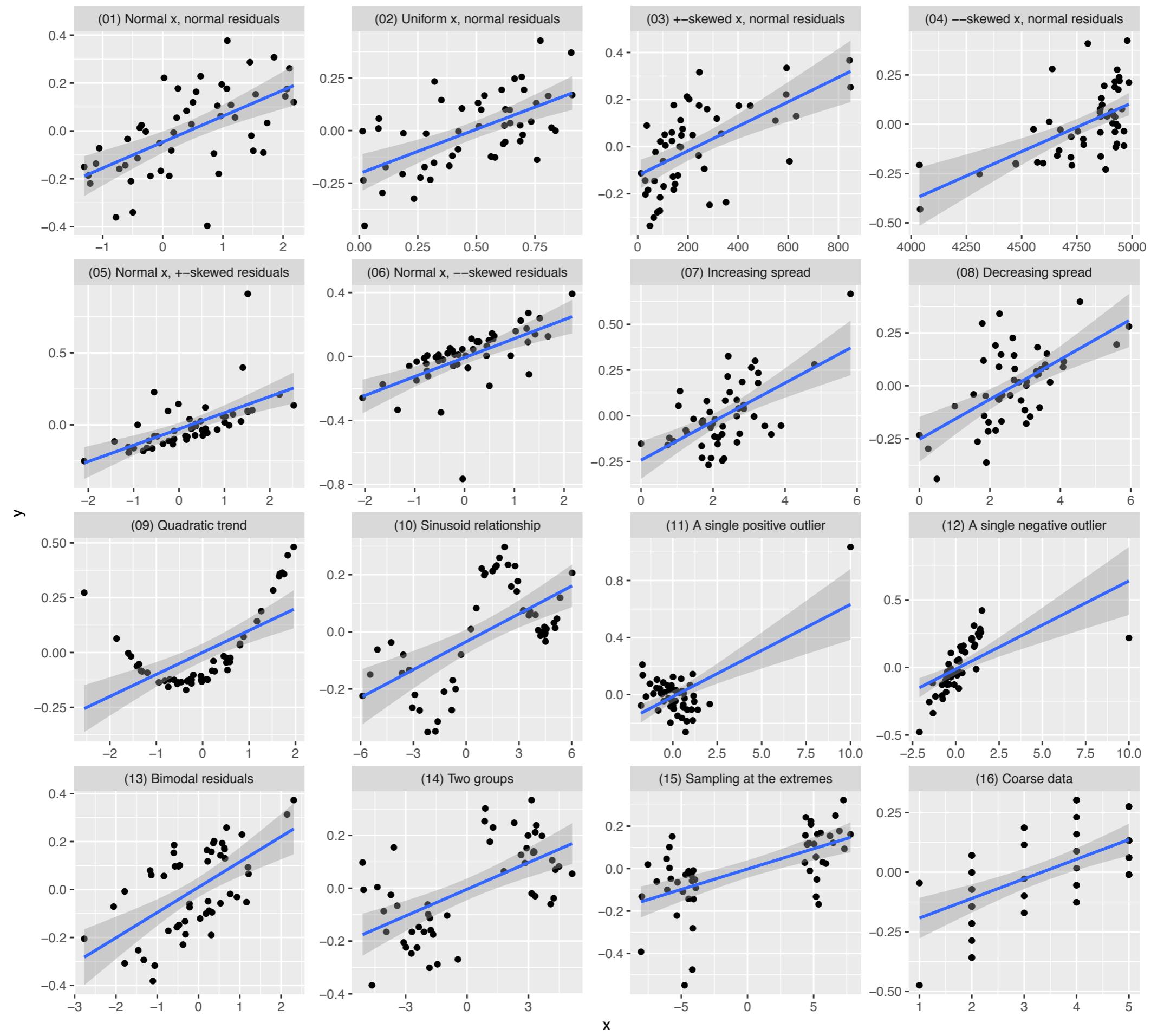
The course packet gets
extracted to `~/Desktop/`

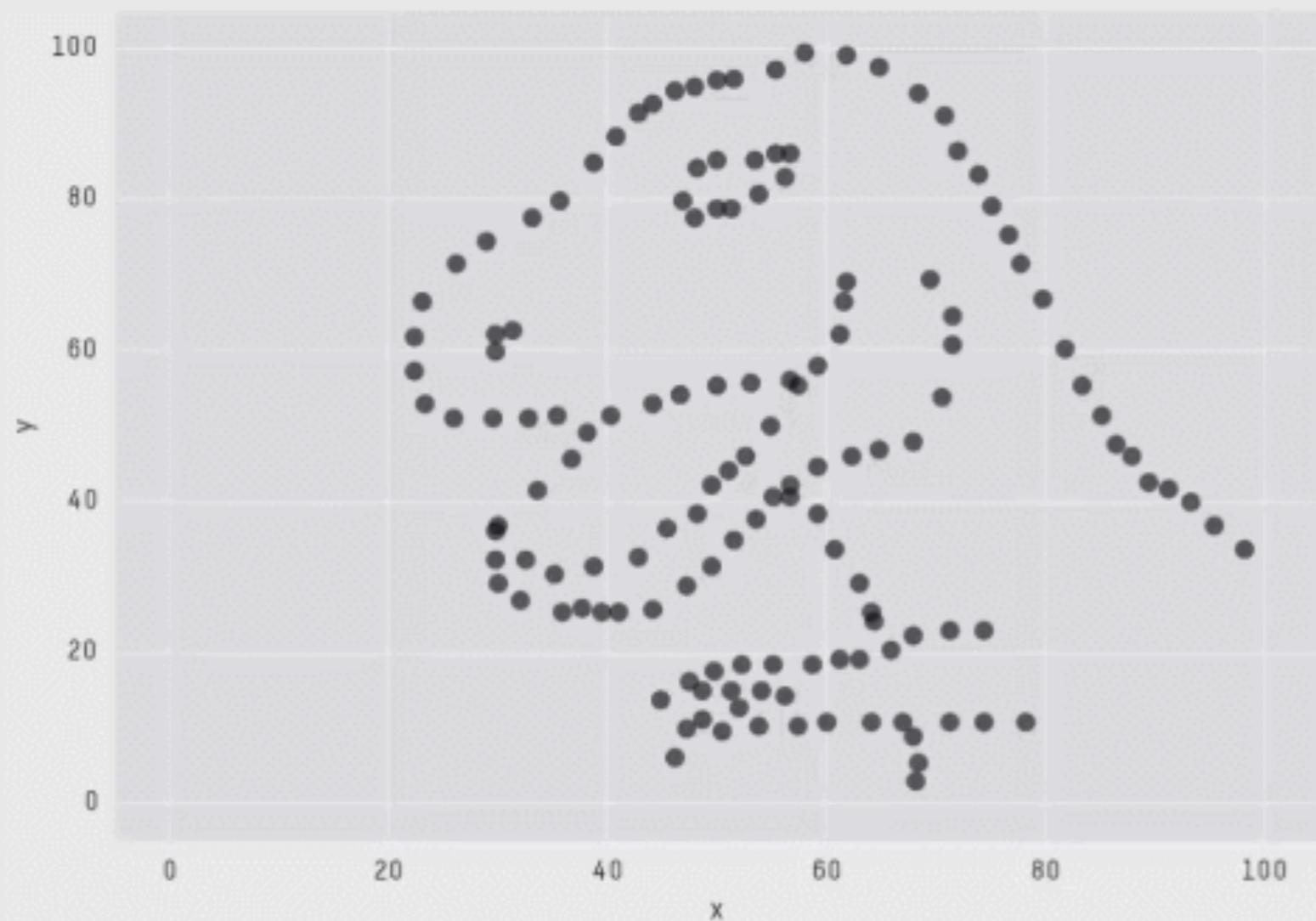


LOOK AT YOUR DATA

SEEING THINGS

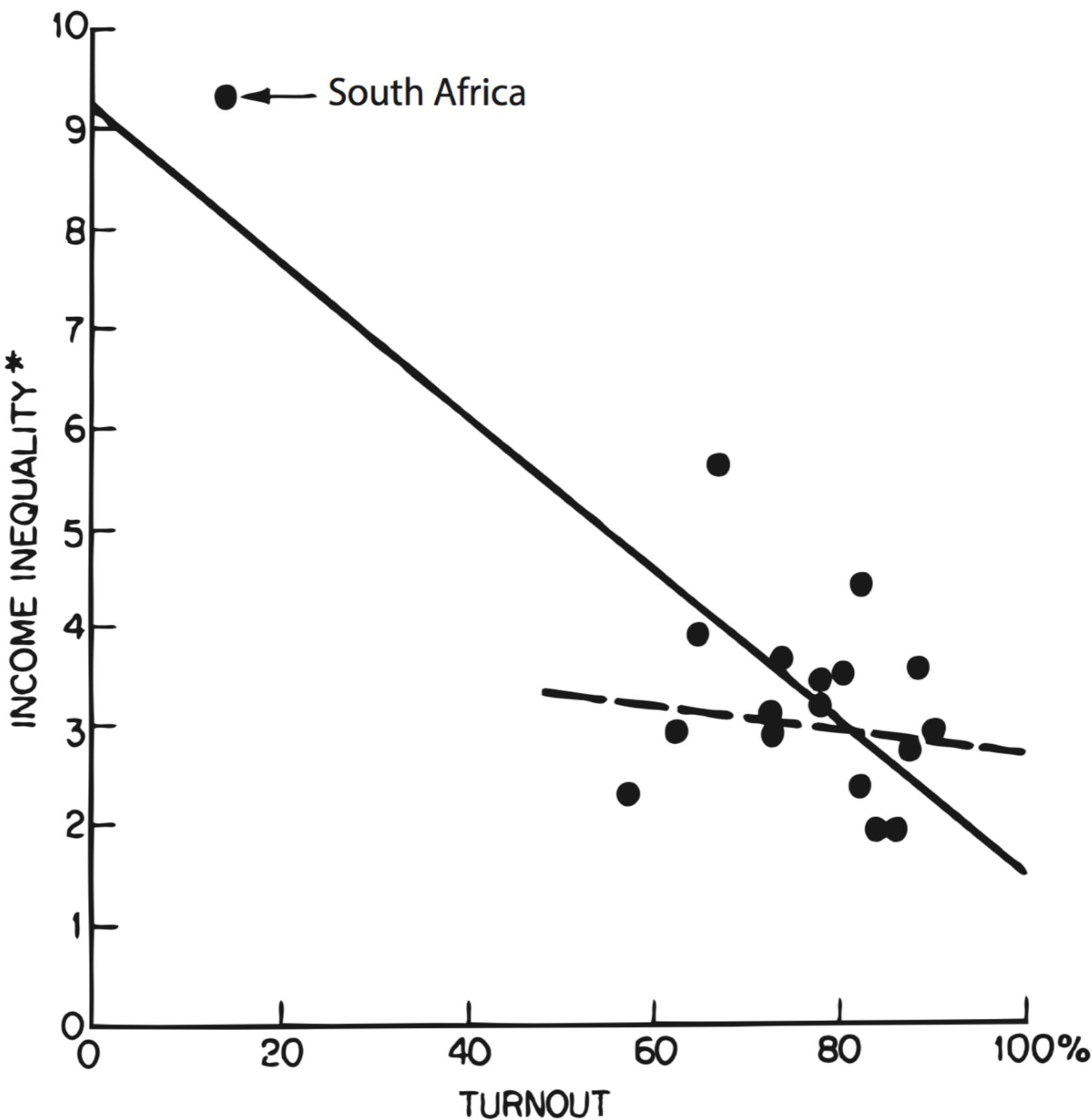






X Mean: 54.2659224
Y Mean: 47.8313999
X SD : 16.7649829
Y SD : 26.9342120
Corr. : -0.0642526

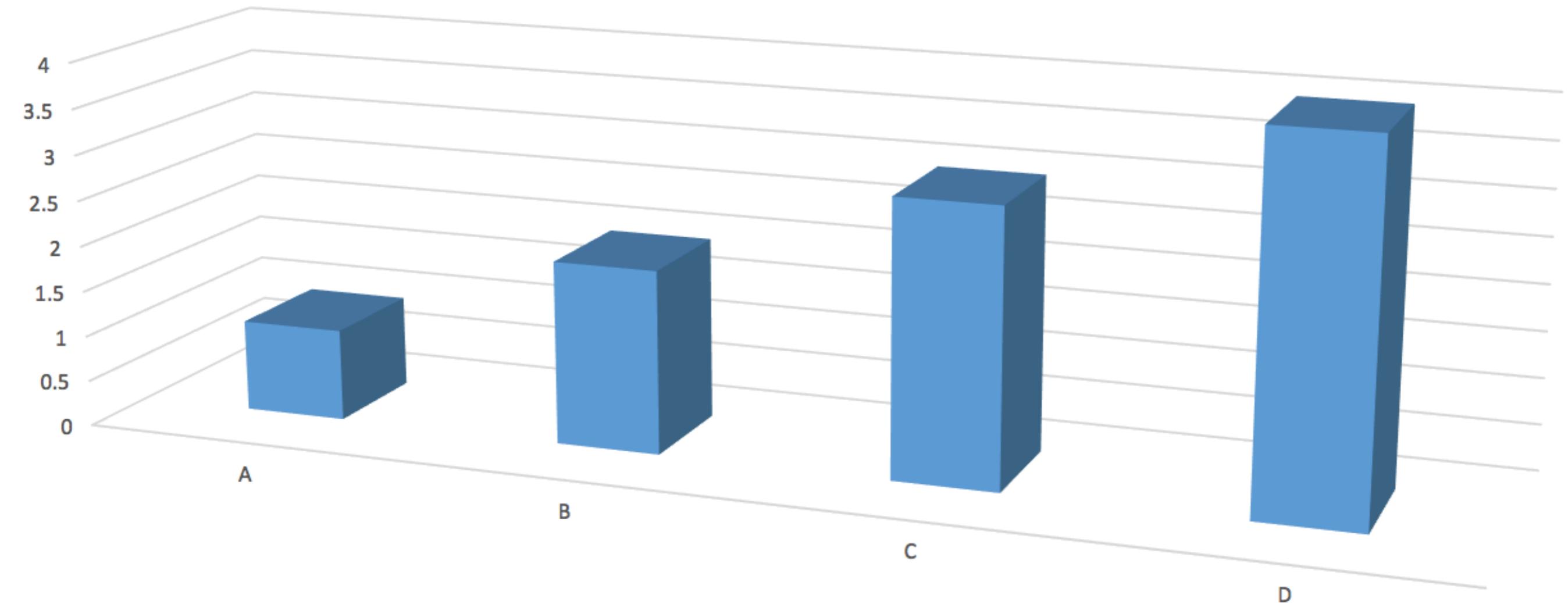
Cairo; Matejka & Fitzmaurice

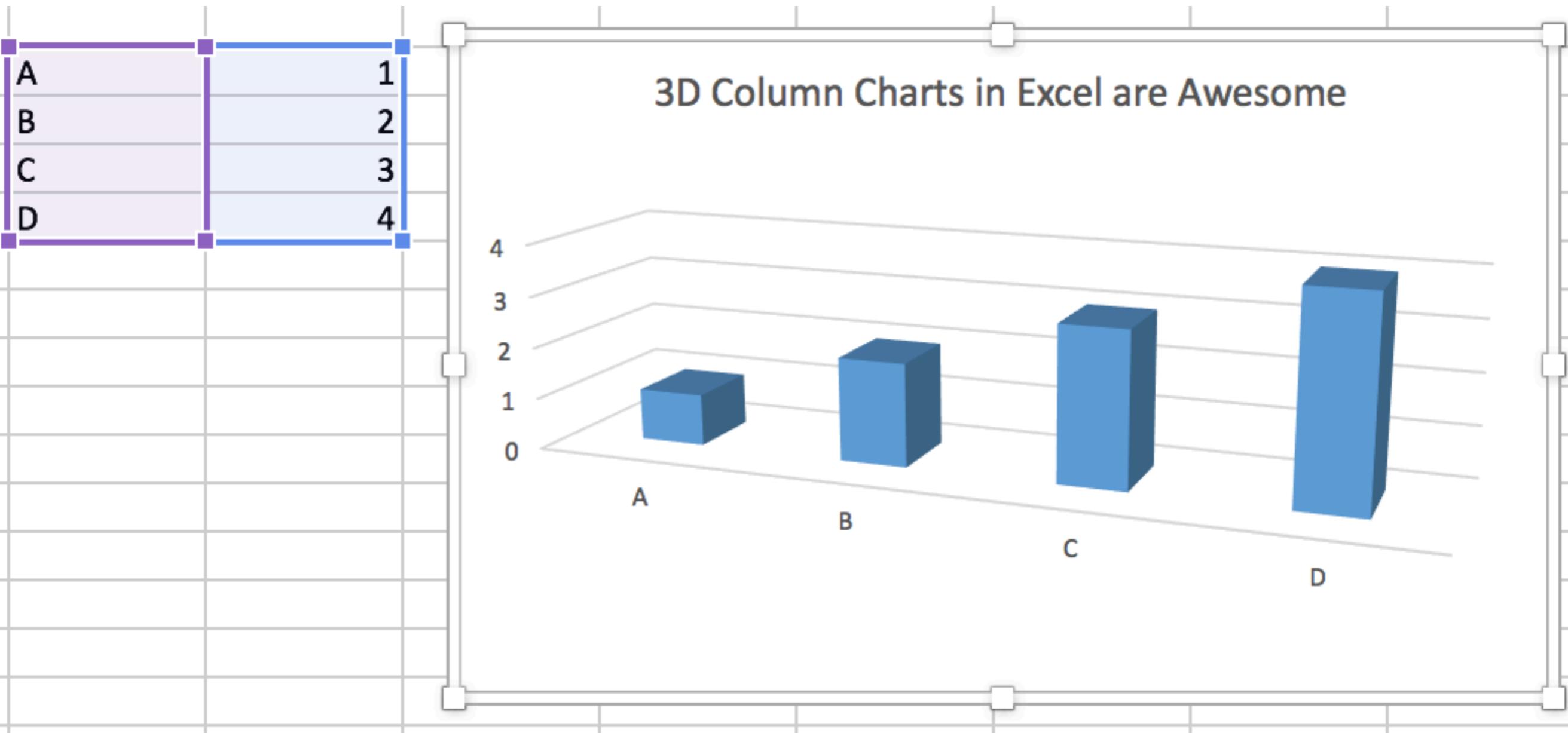


Key. — Bivariate slope including South Africa ($N = 18$)
— Bivariate slope excluding South Africa ($N = 17$)

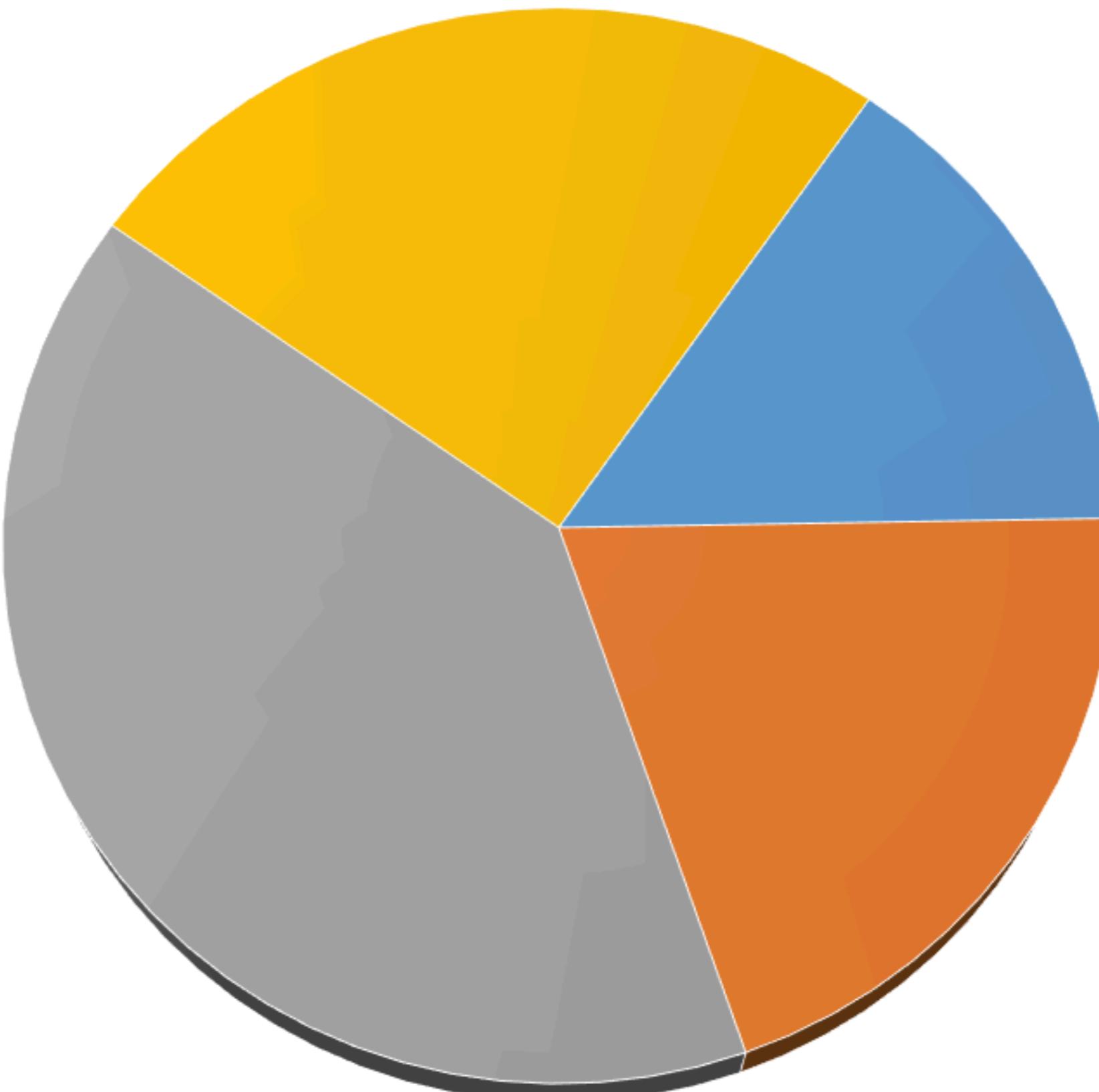
NOT
SEEING THINGS

3D Column Charts in Excel are Awesome





The Infamous Pie Chart

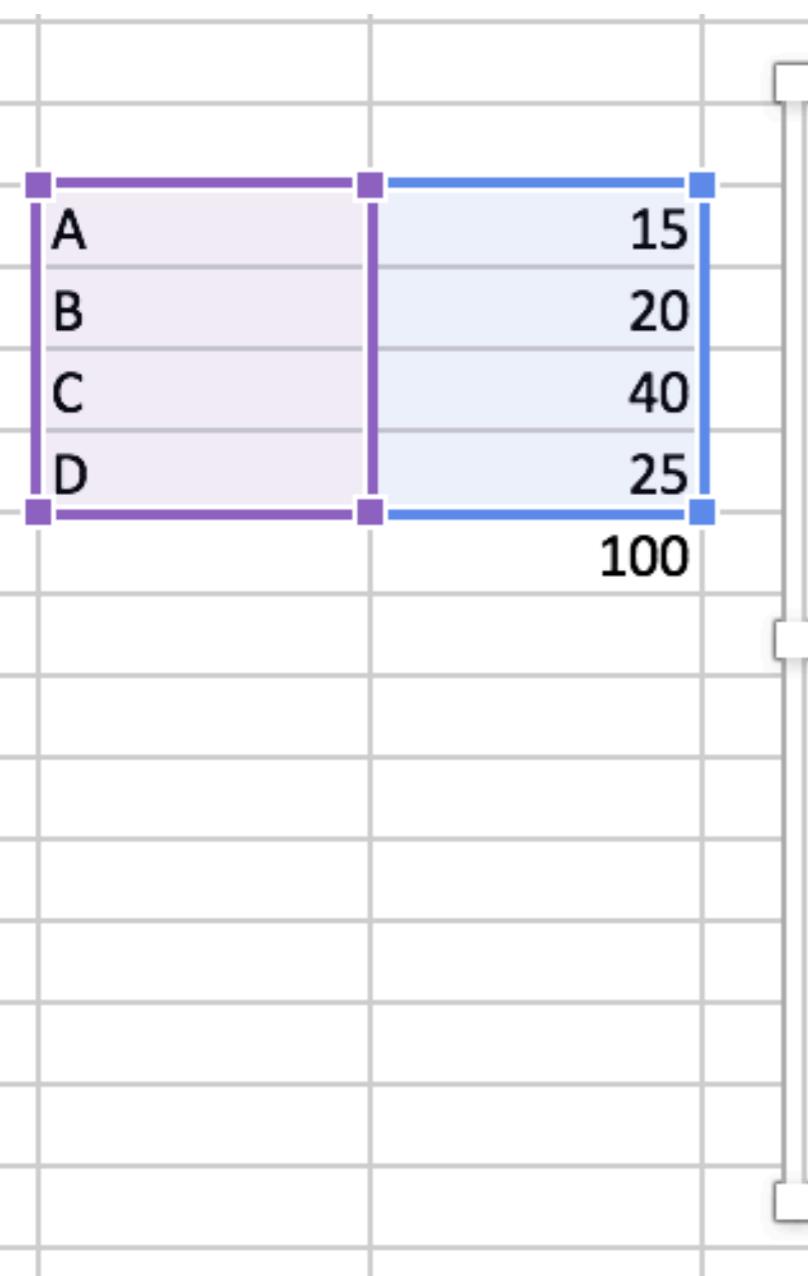


■ A ■ B ■ C ■ D

The Infamous Pie Chart



■ A ■ B ■ C ■ D



BAD TASTE

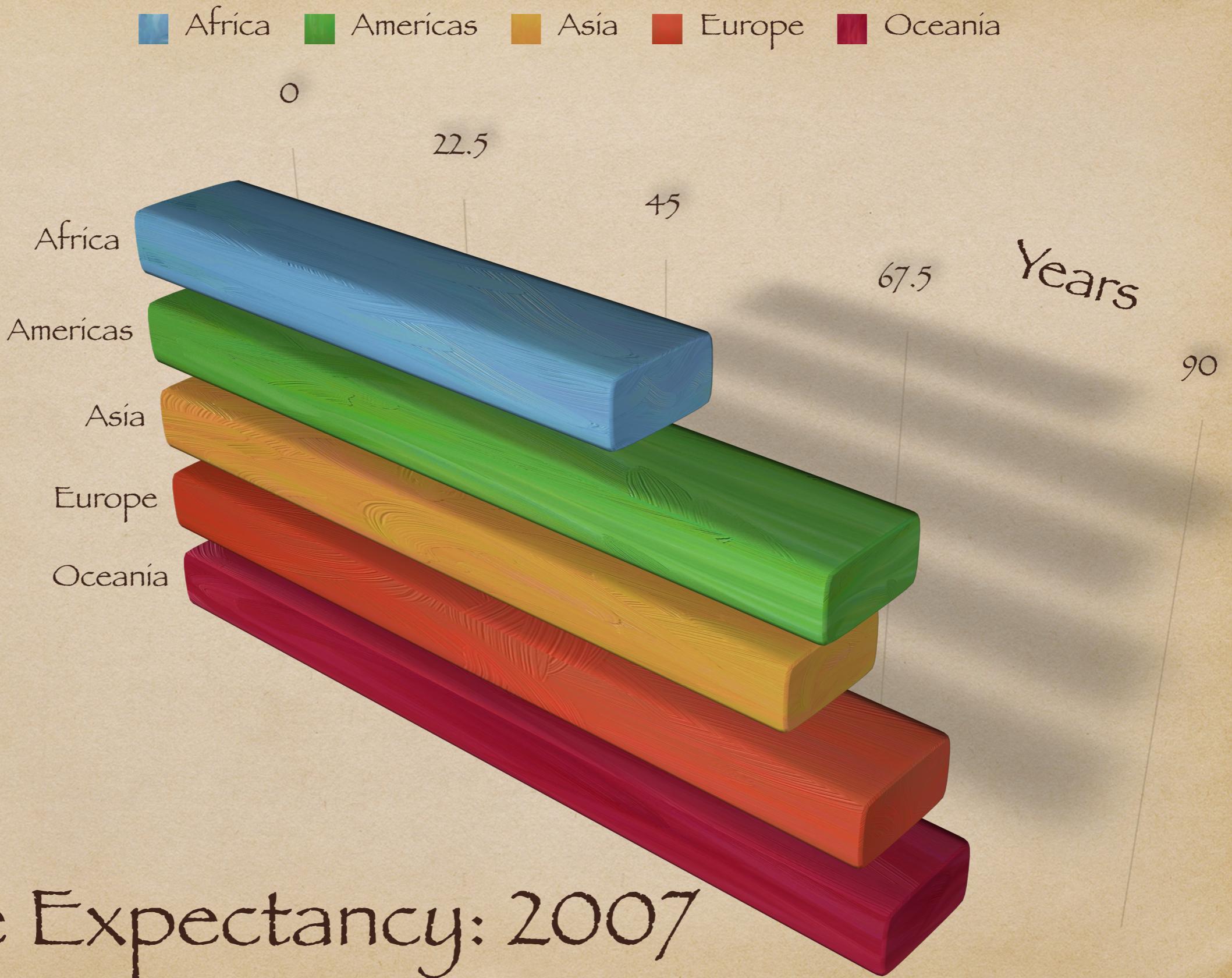
BAD DATA

BAD PERCEPTION

TASTE

**SIMPLIFY,
SIMPLIFY?**

CONTINENT



Life Expectancy: 2007

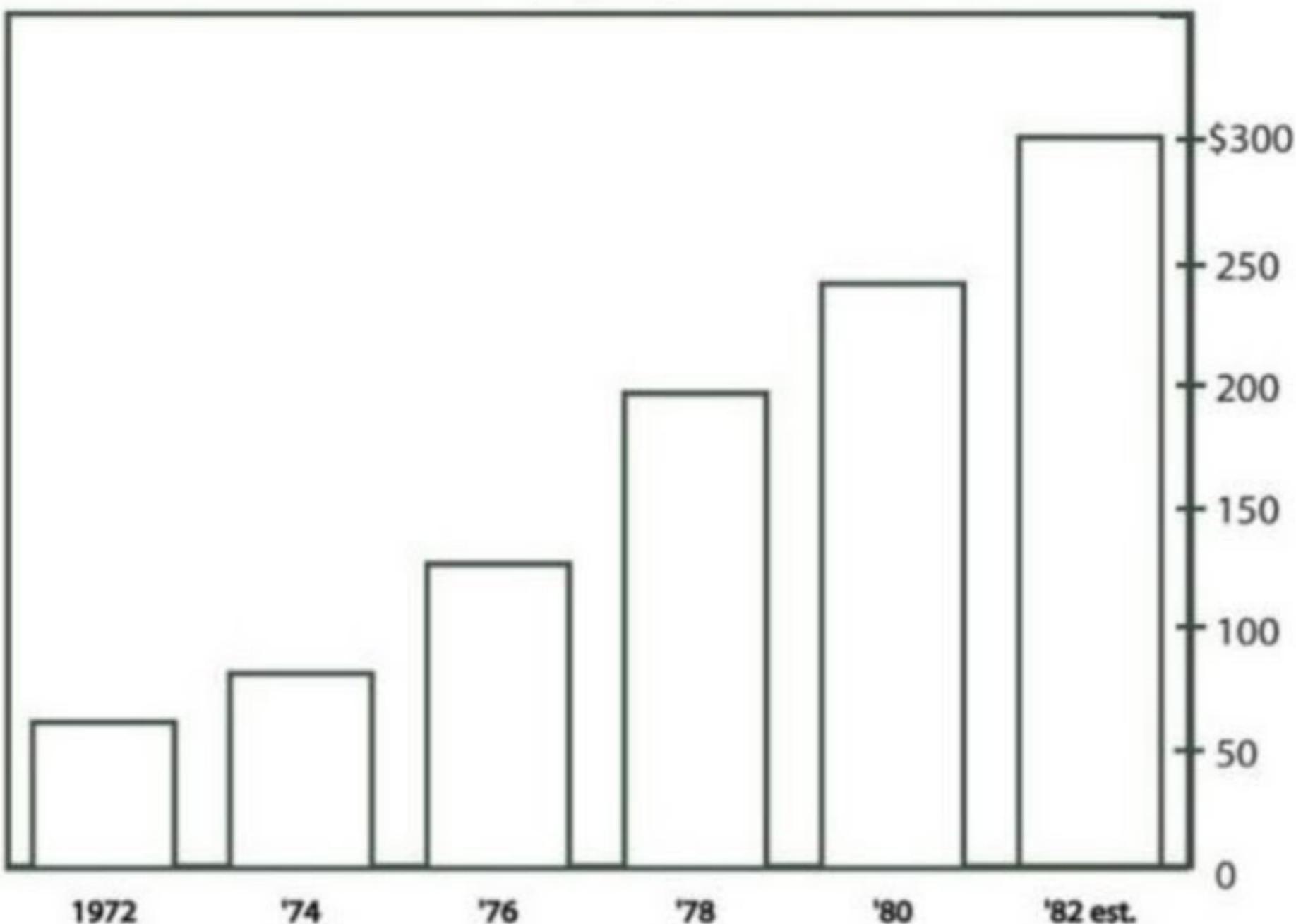
MONSTROUS COSTS

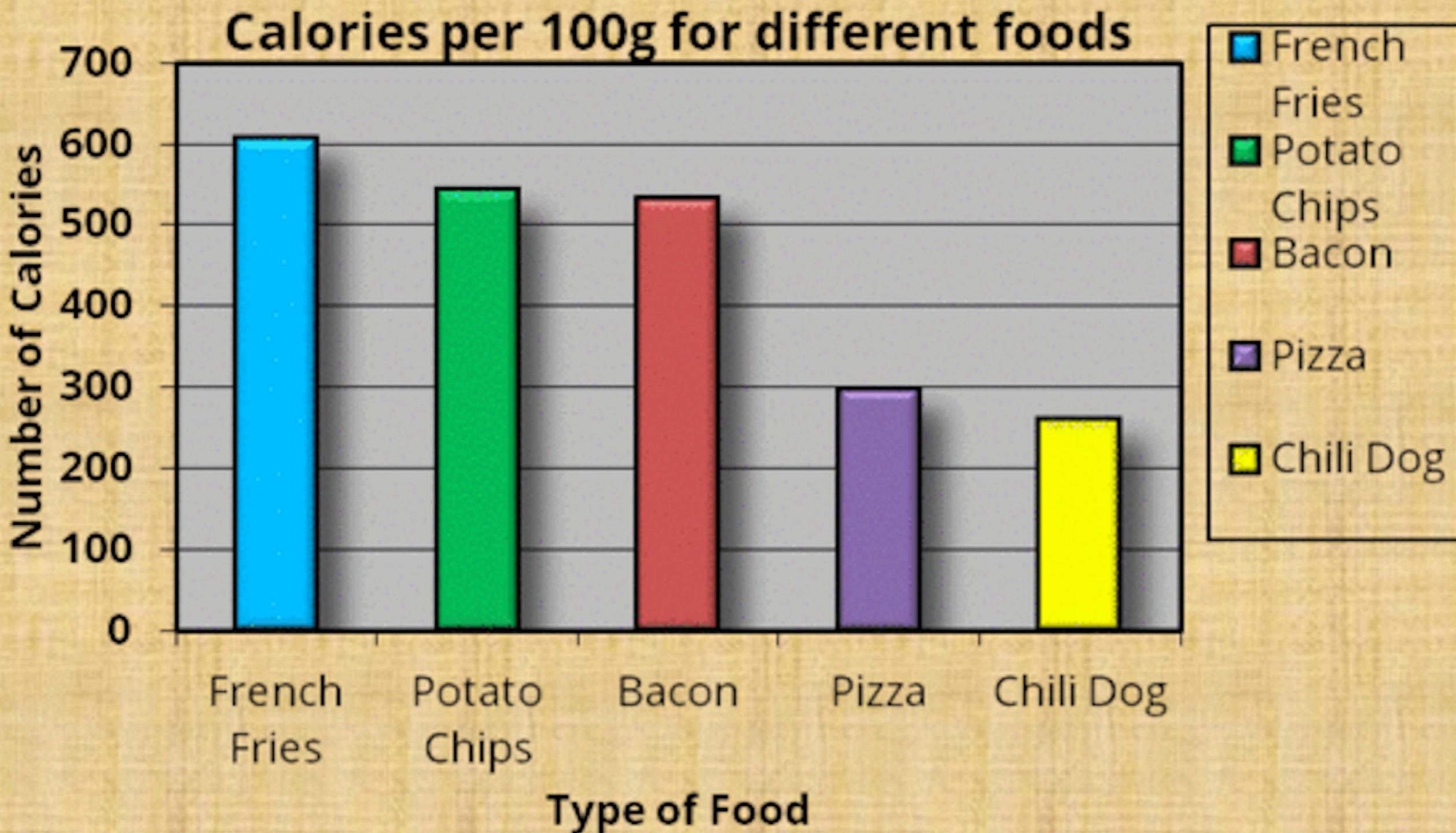
Total House and Senate campaign expenditures,
in millions



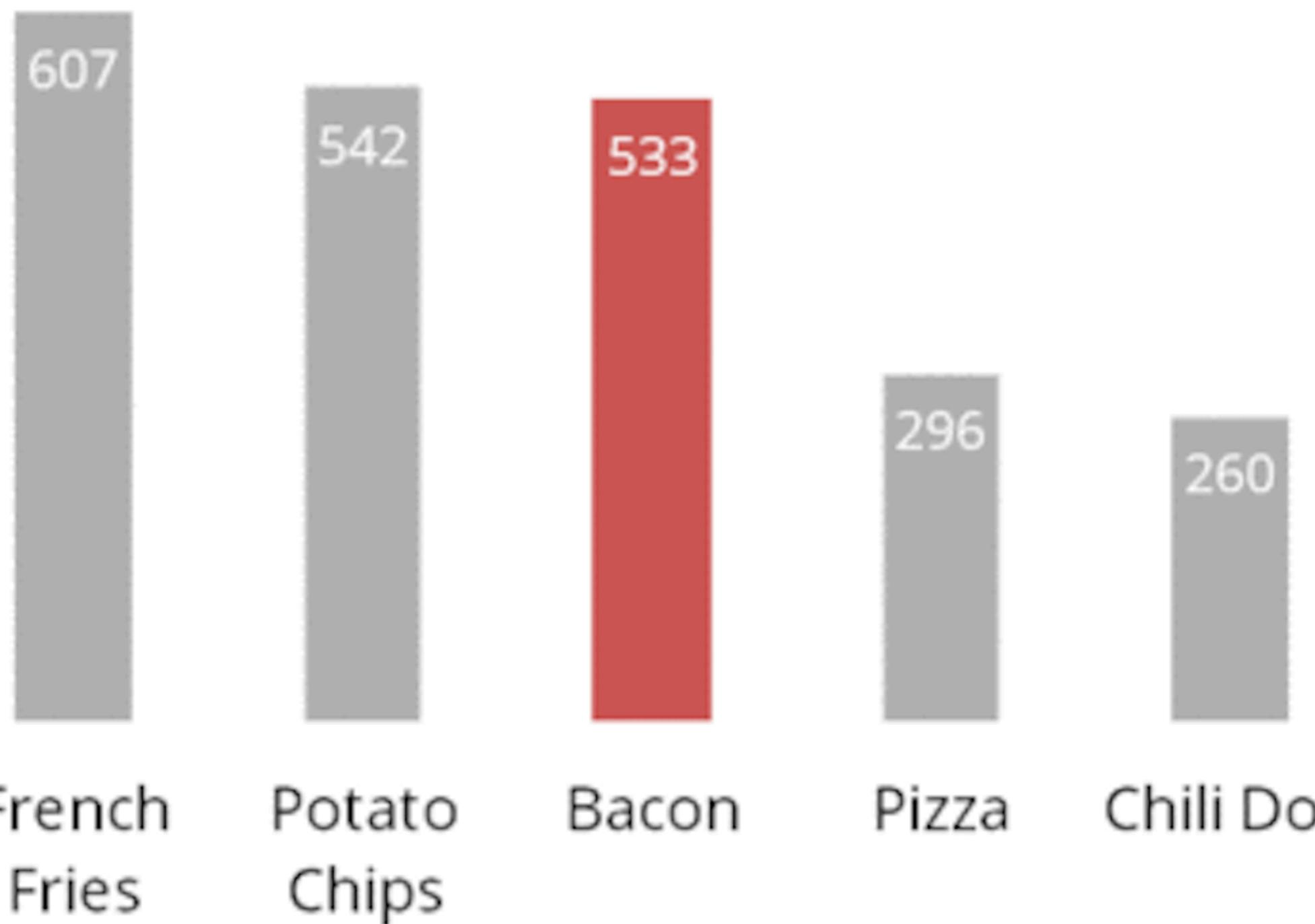
MONSTROUS COSTS

Total House and Senate campaign expenditures, in millions





Calories per 100g



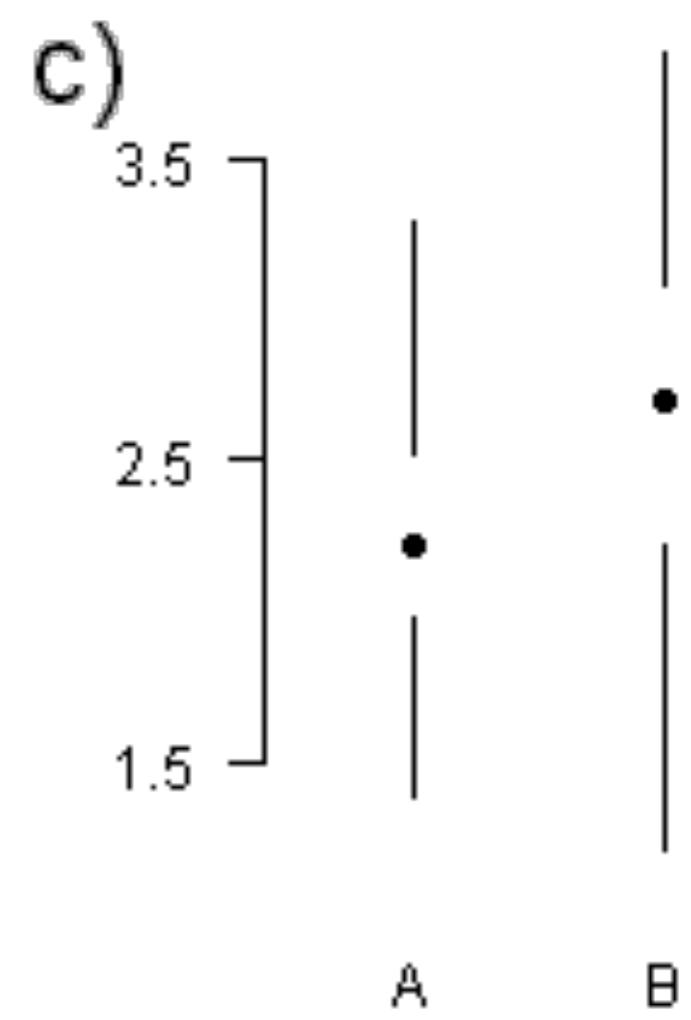
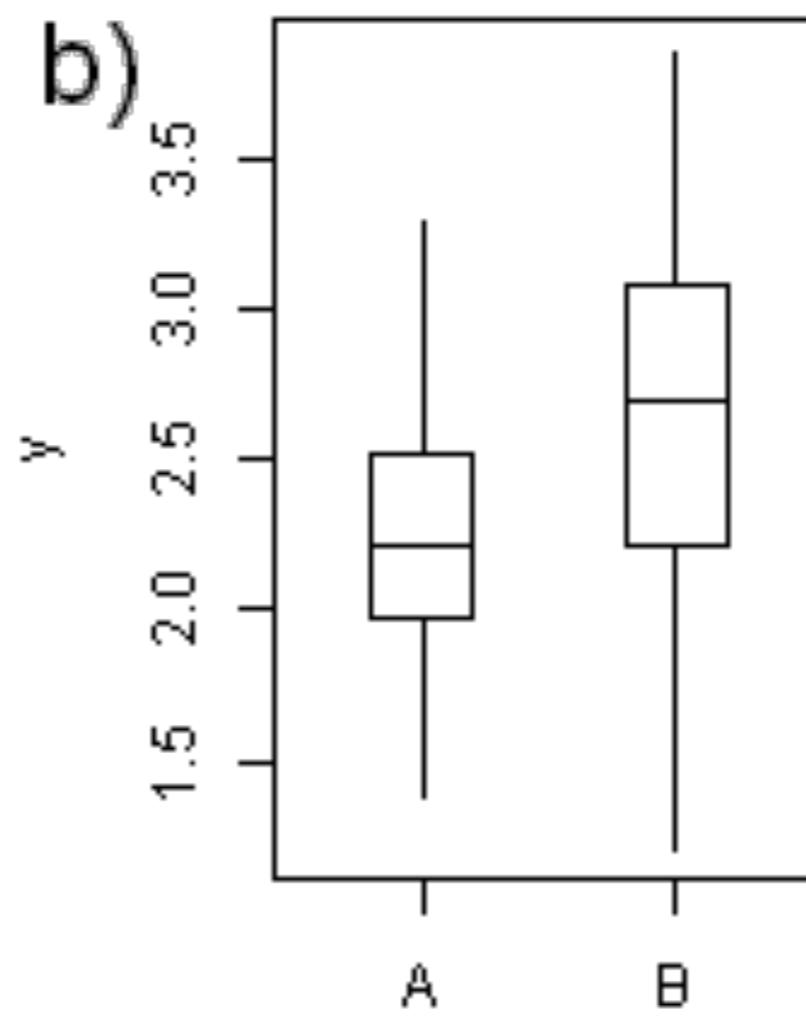
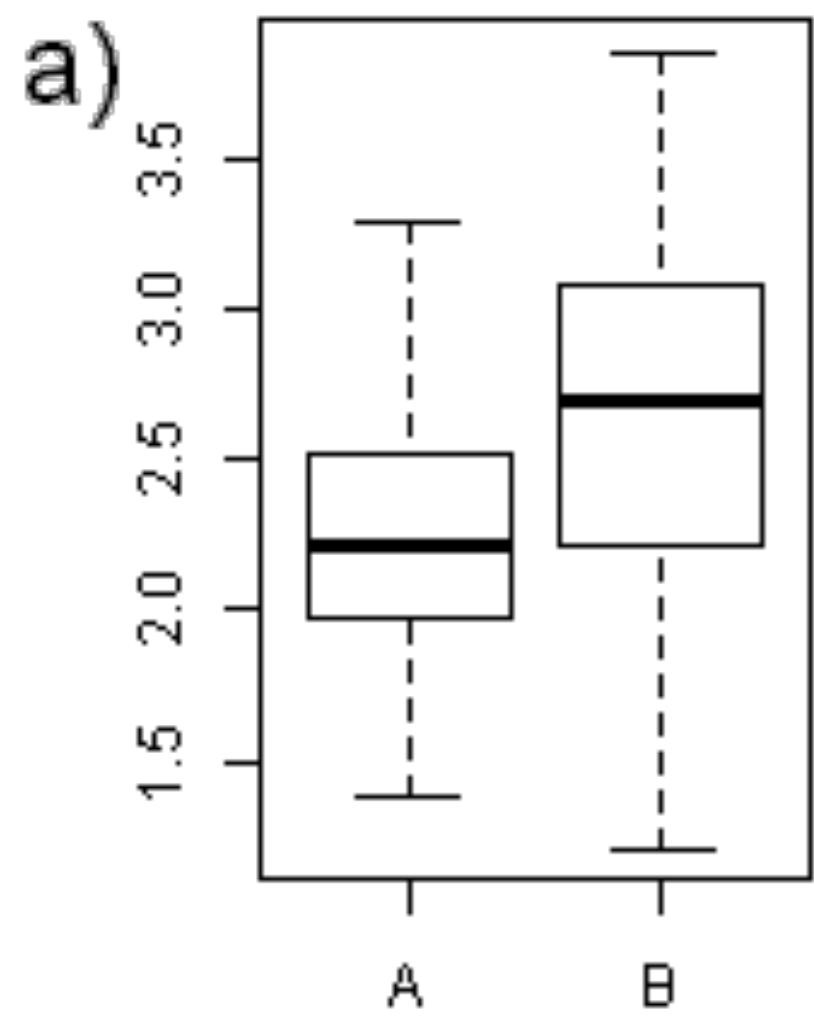
French
Fries

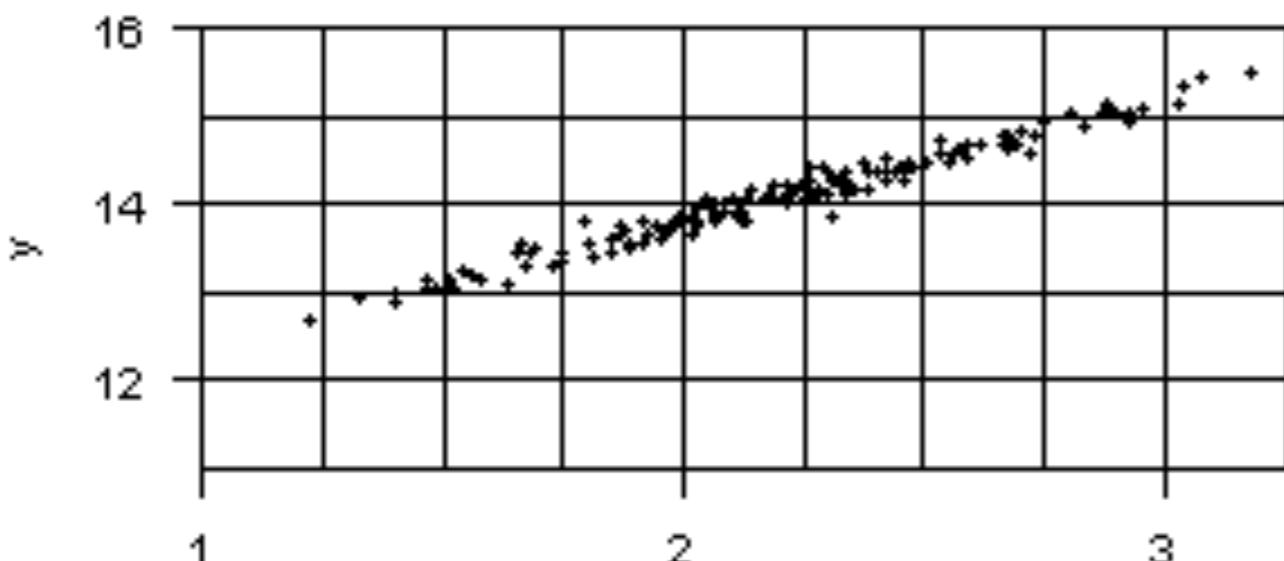
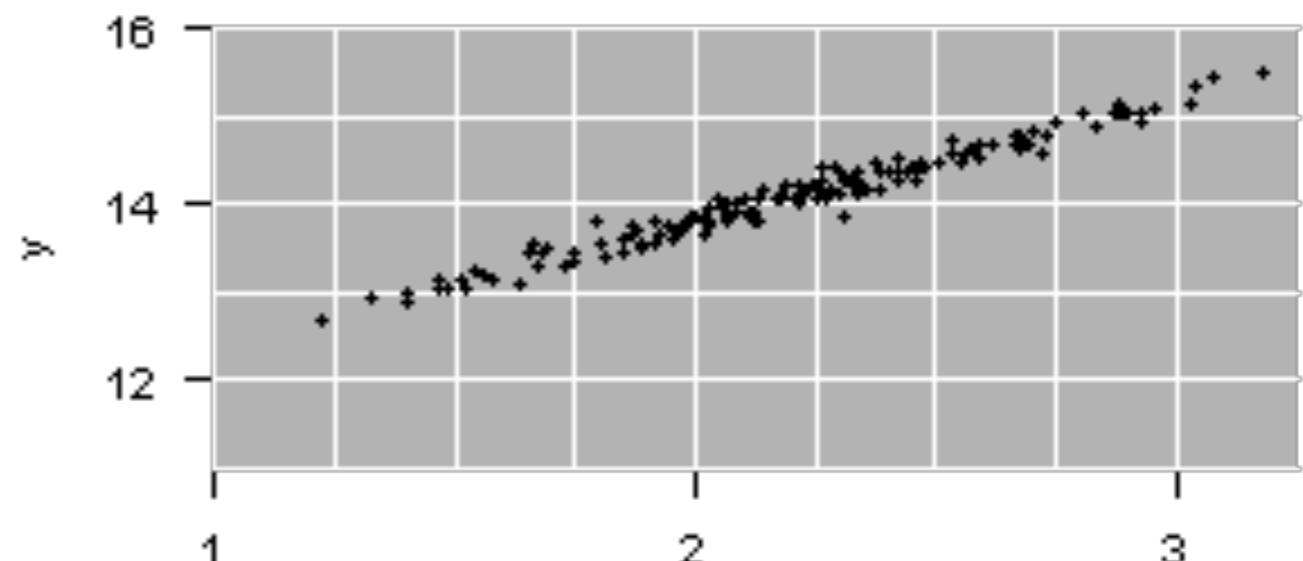
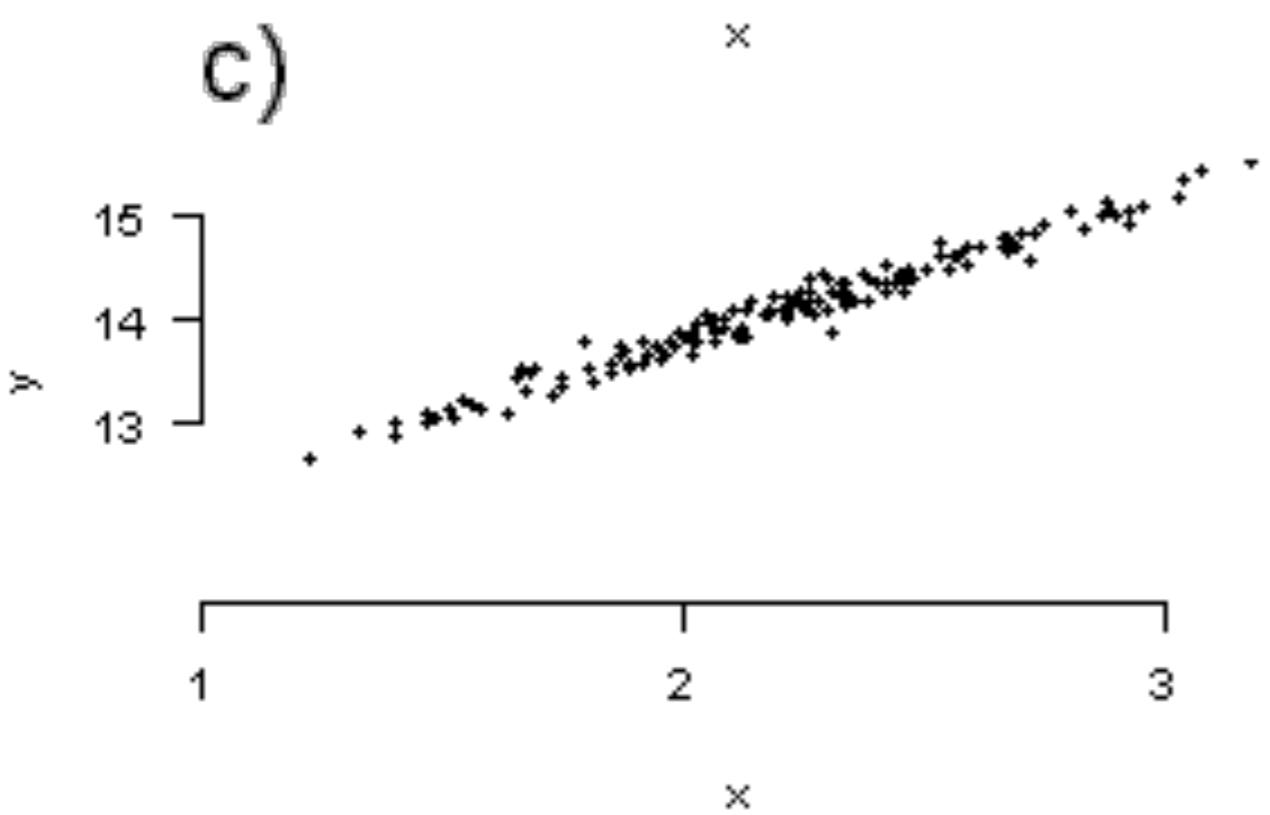
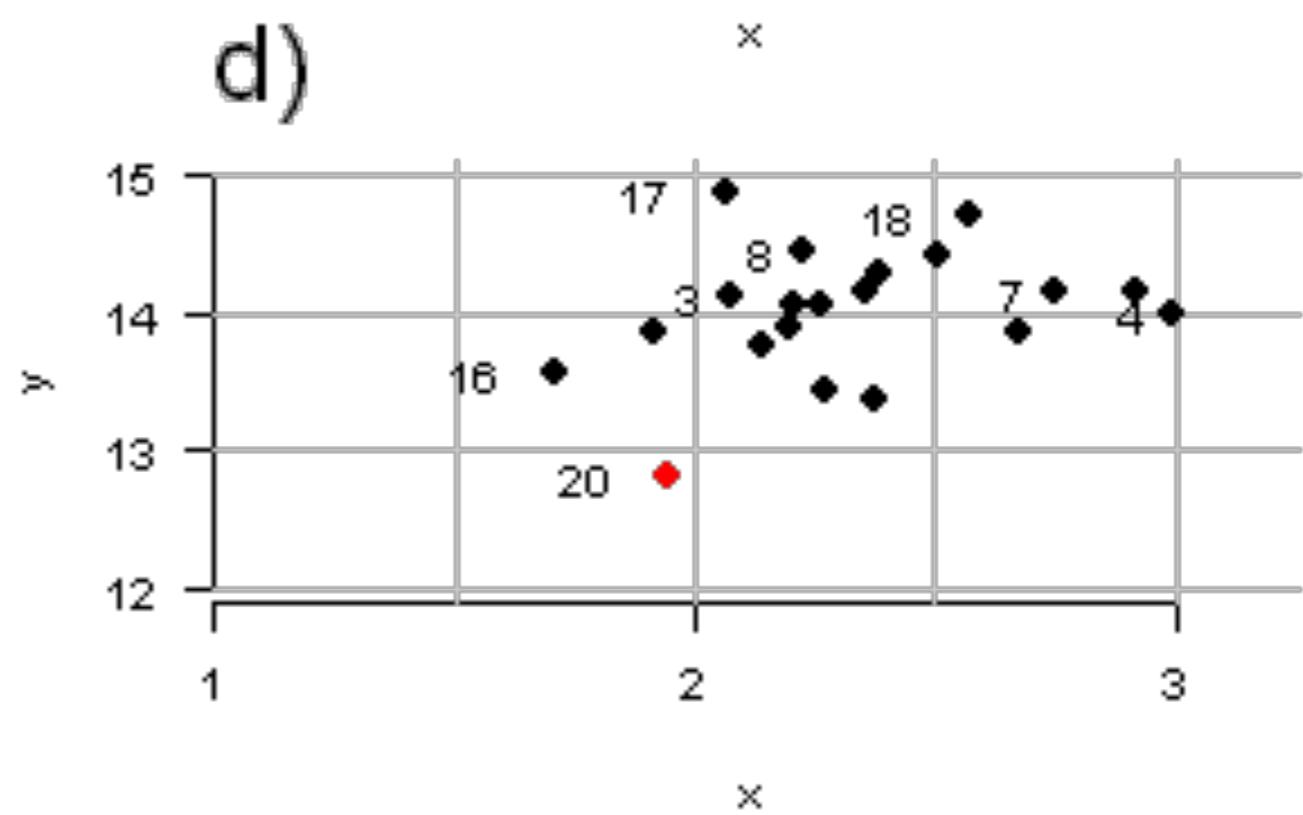
Potato
Chips

Bacon

Pizza

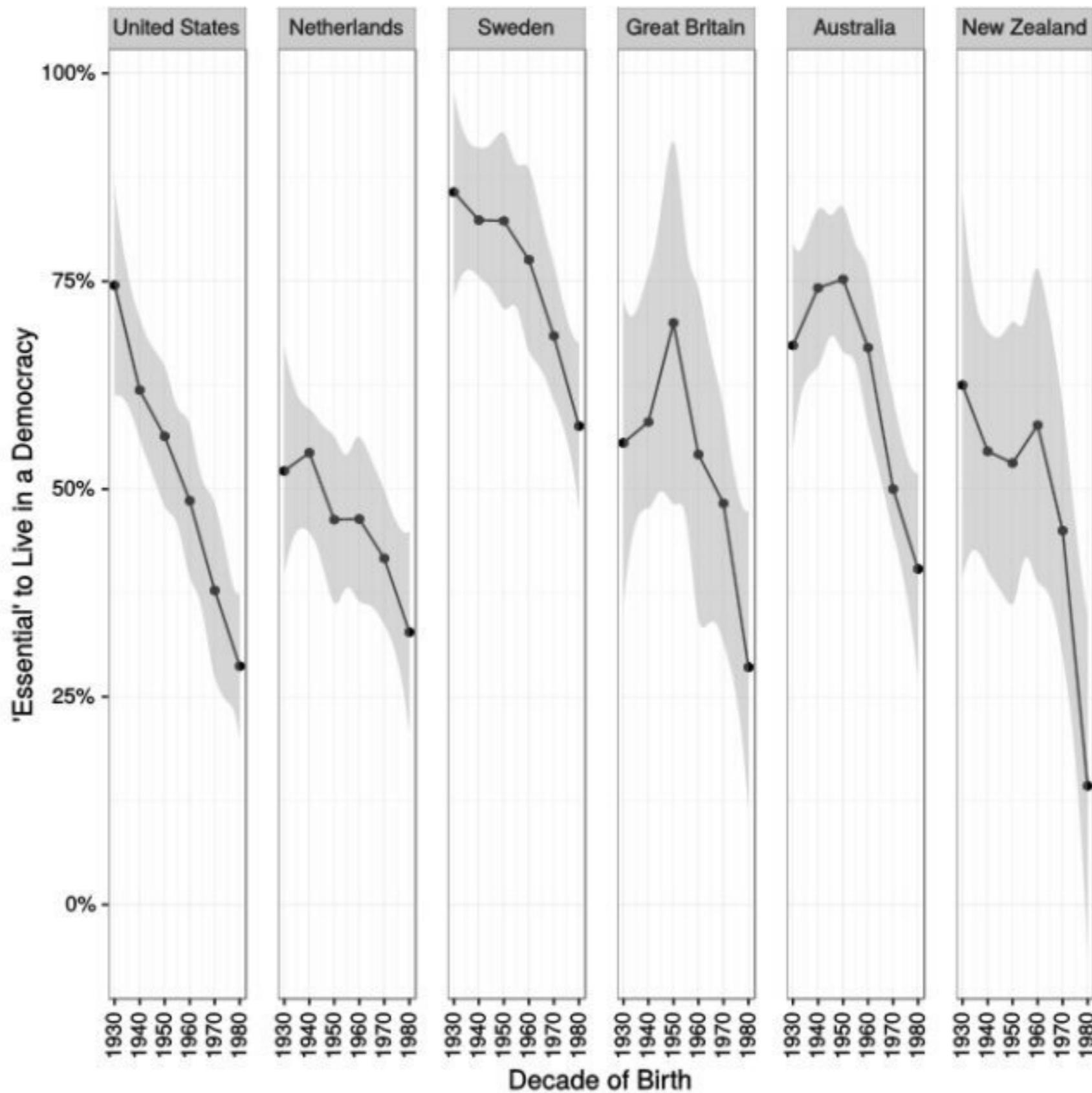
Chili Dog



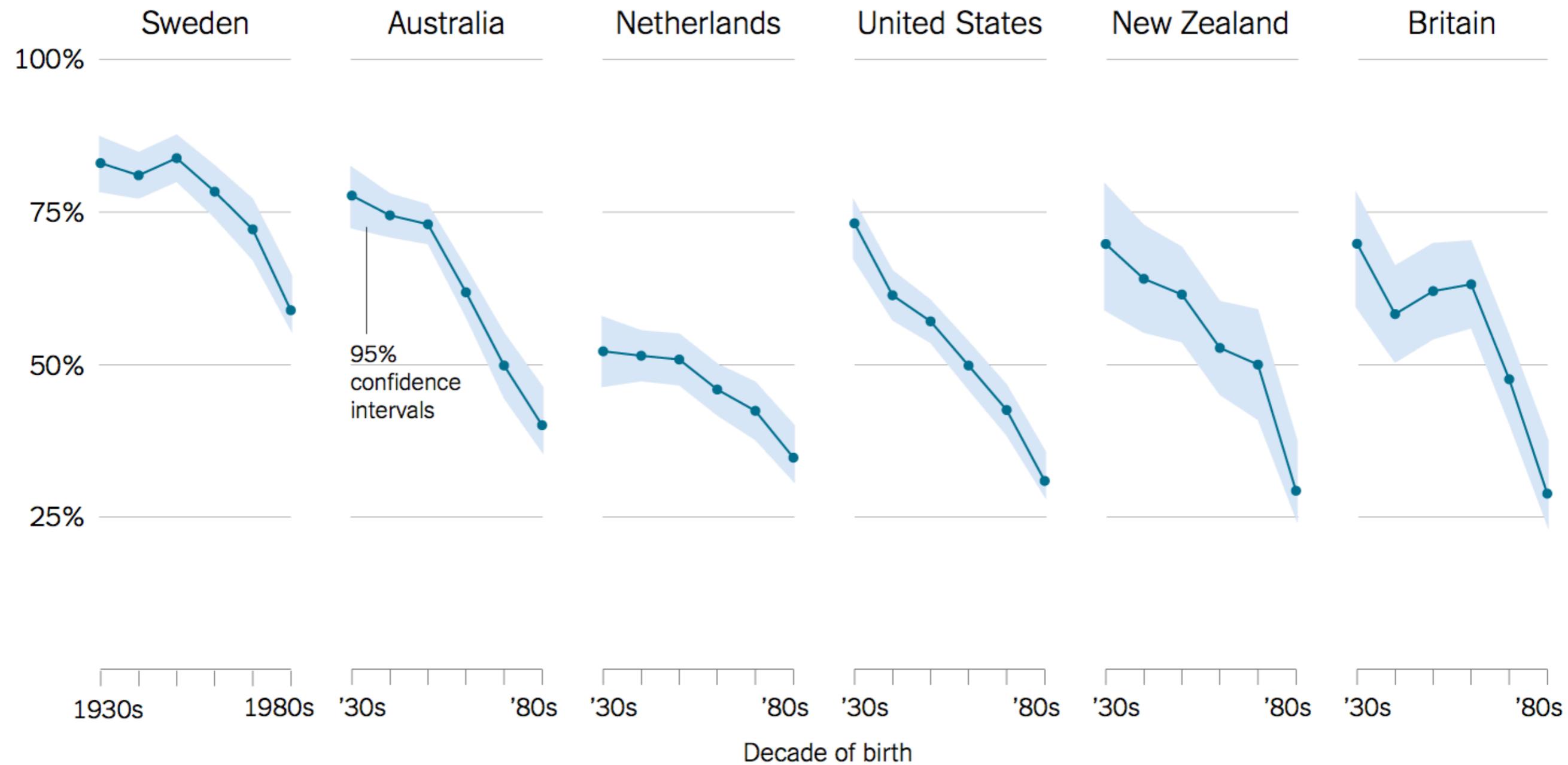
a)**b)****c)****d)**

DATA

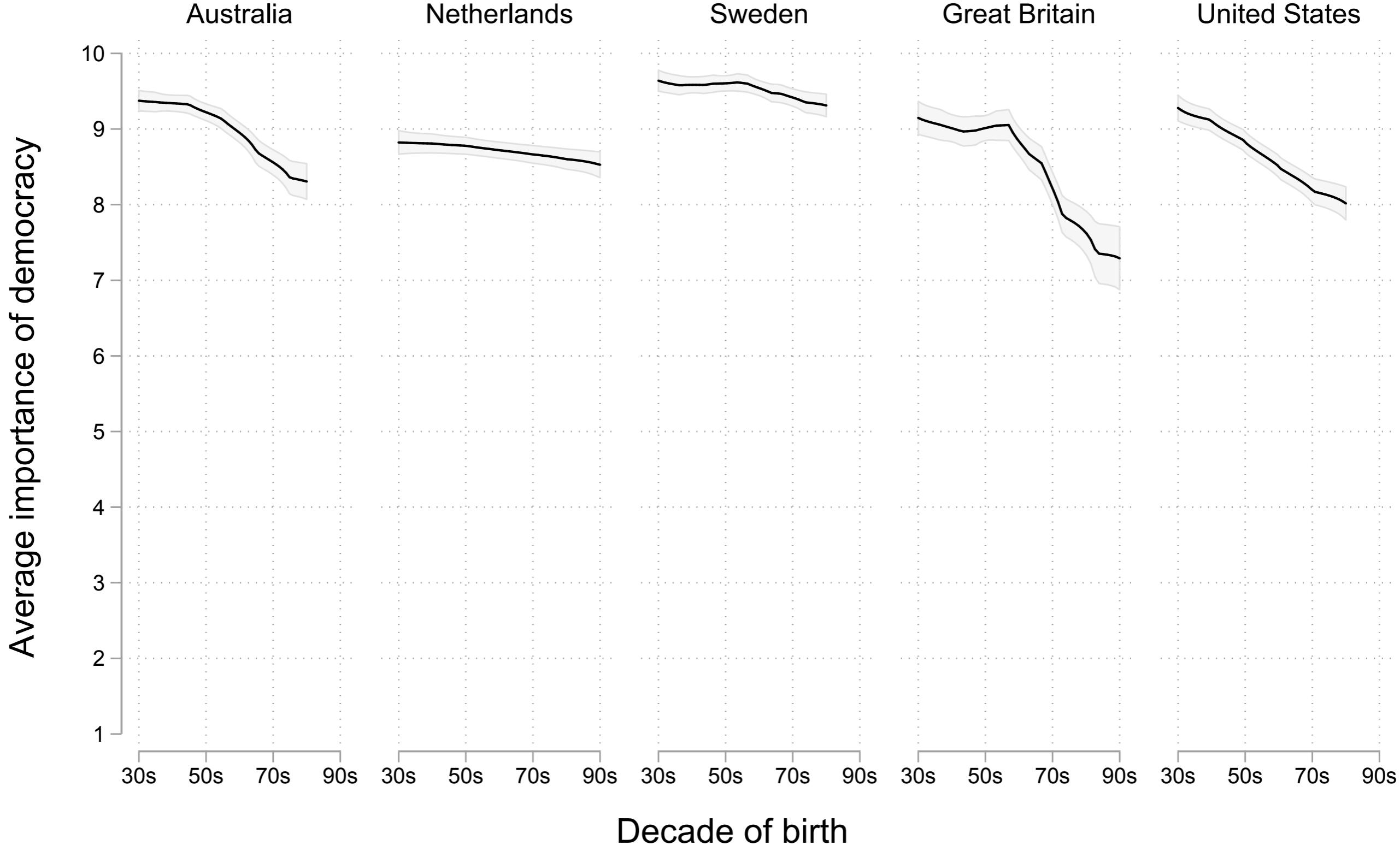
**GOOD DESIGN
BAD DATA**



Percentage of people who say it is “essential” to live in a democracy



Source: Yascha Mounk and Roberto Stefan Foa, “The Signs of Democratic Deconsolidation,” Journal of Democracy | By The New York Times

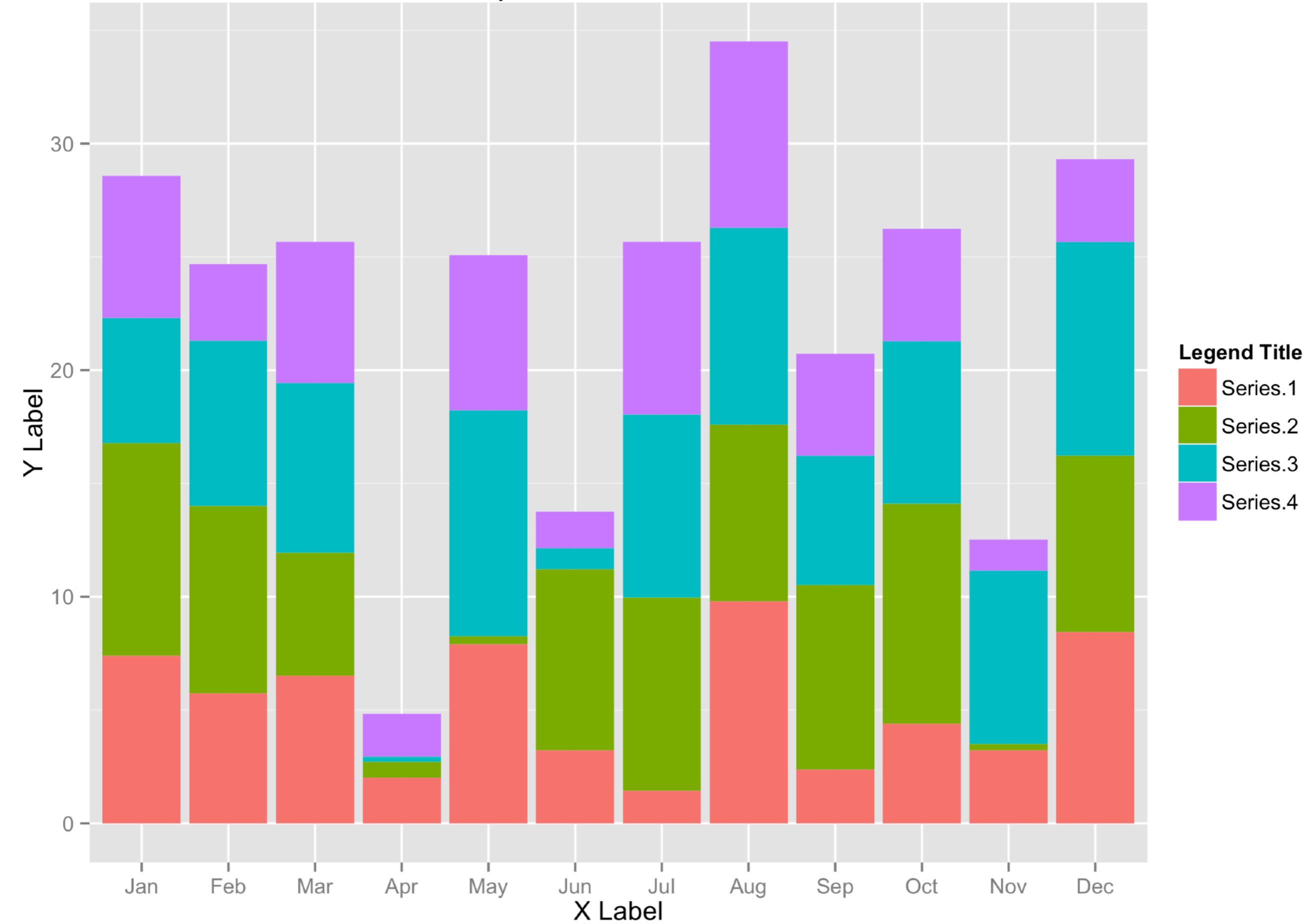


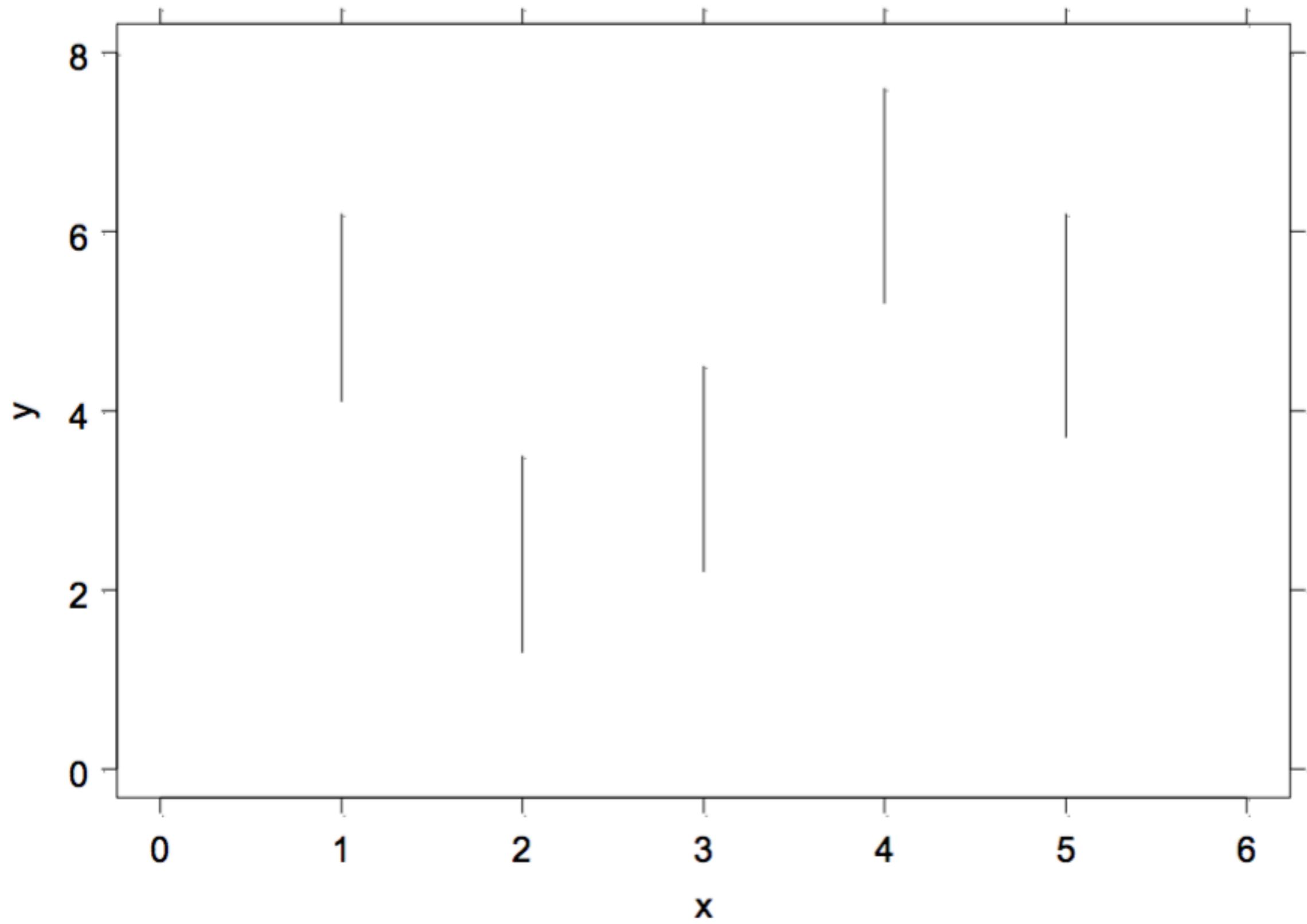
Graph by Erik Voeten, based on WVS 5

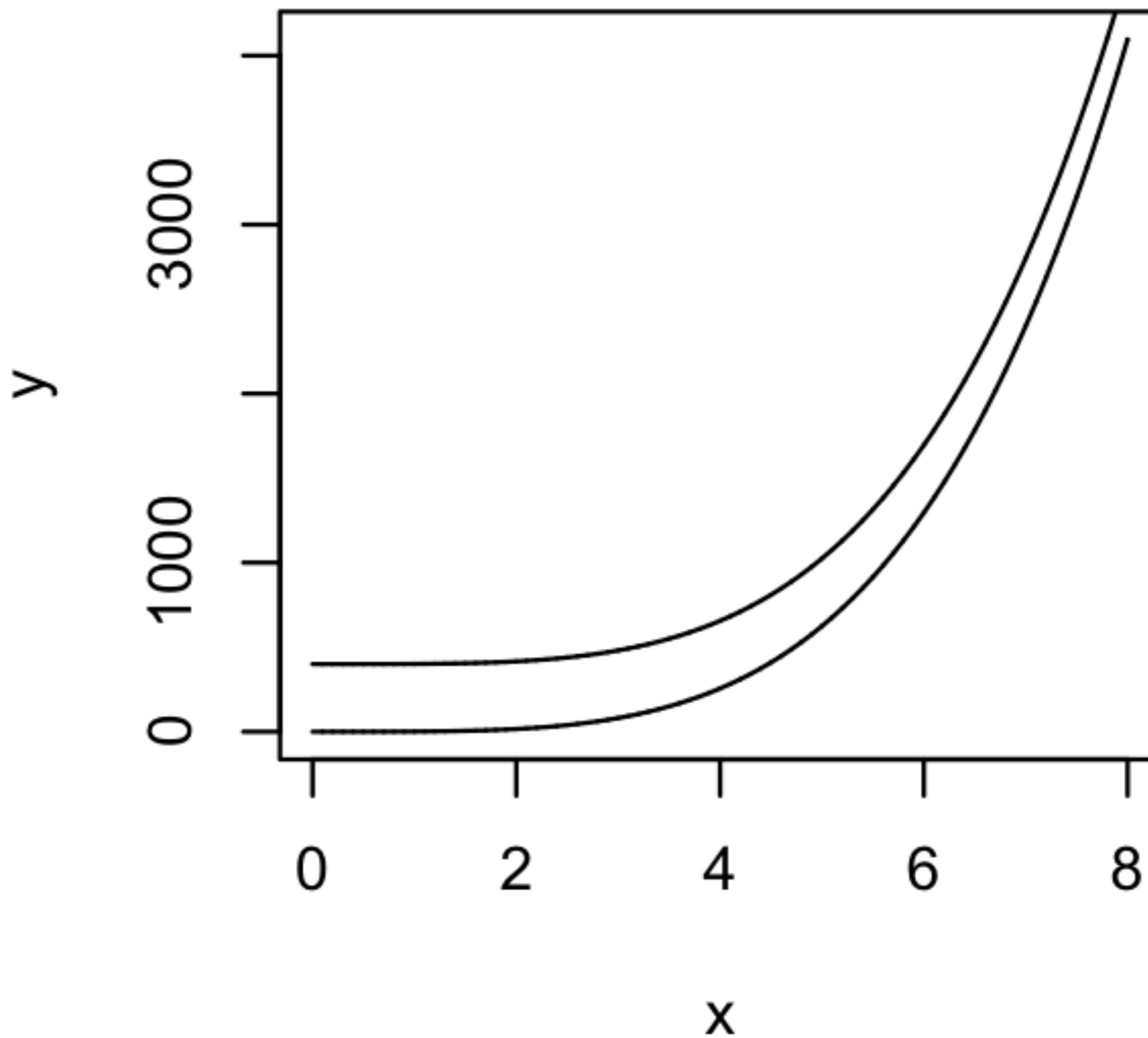
PERCEPTION

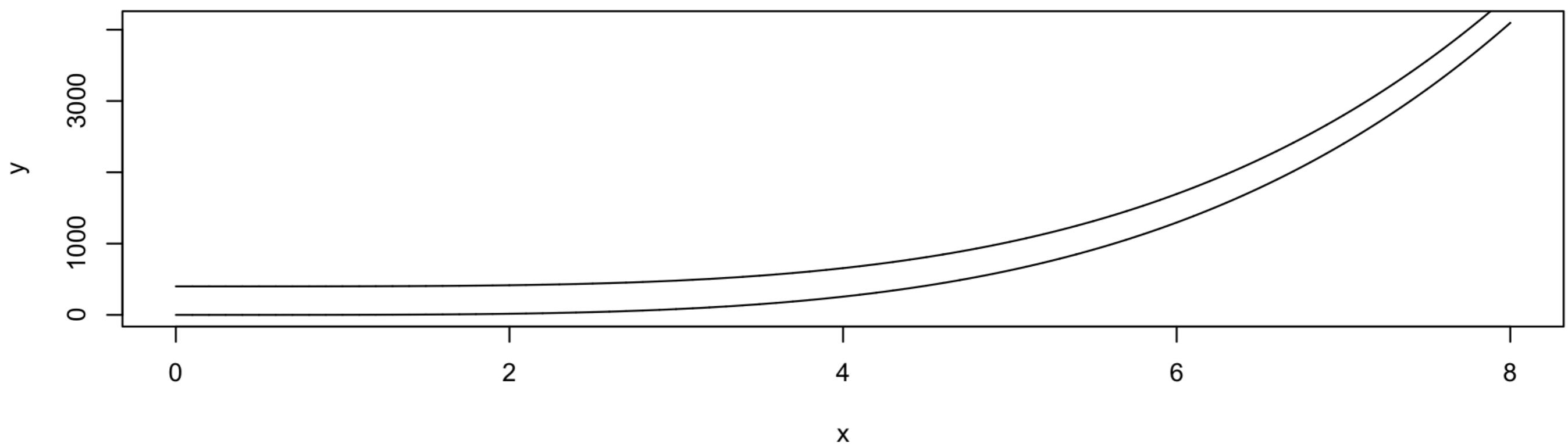
GOOD DESIGN
IS MORE THAN
GOOD TASTE

An Example Stacked Column Chart









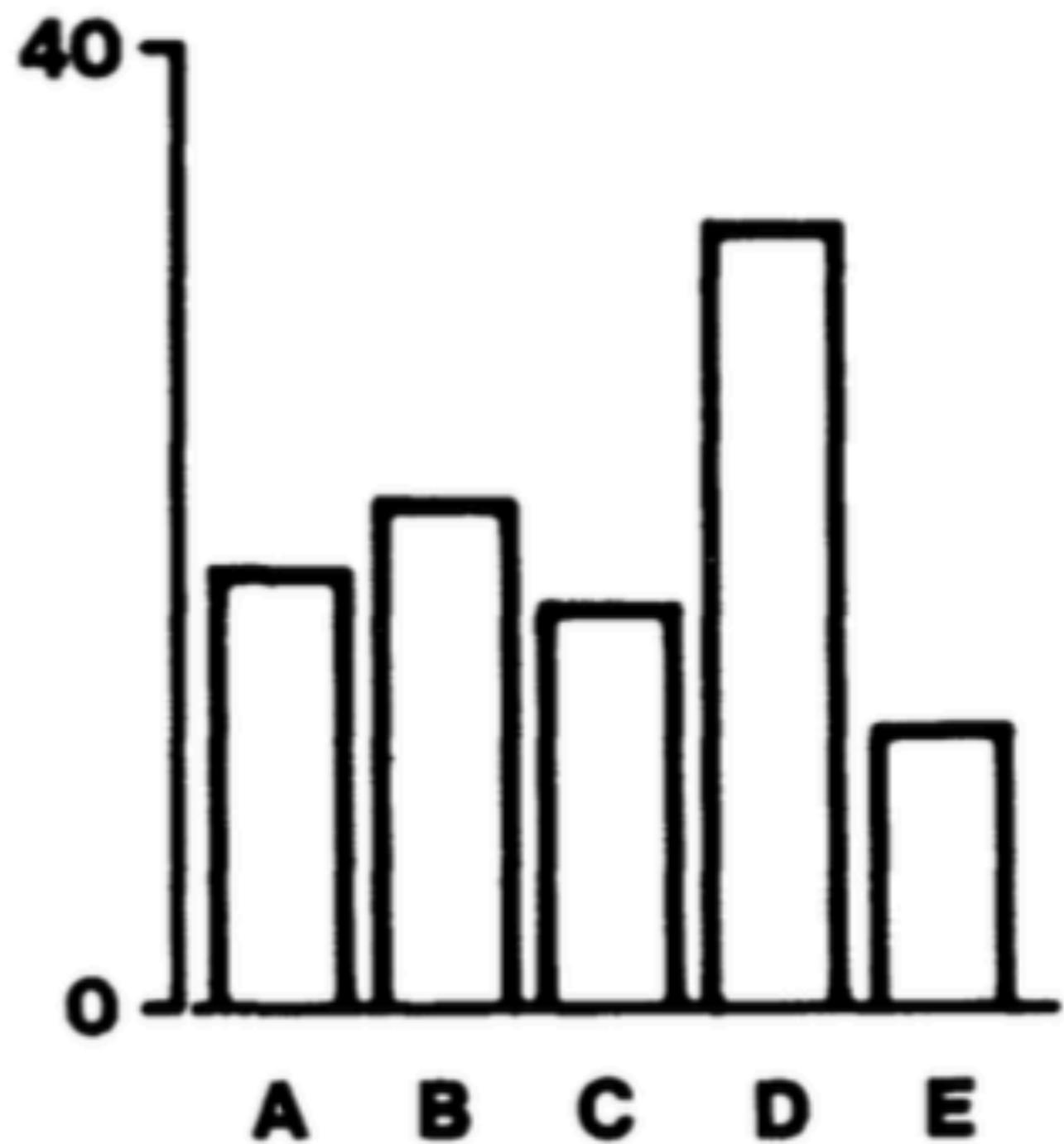
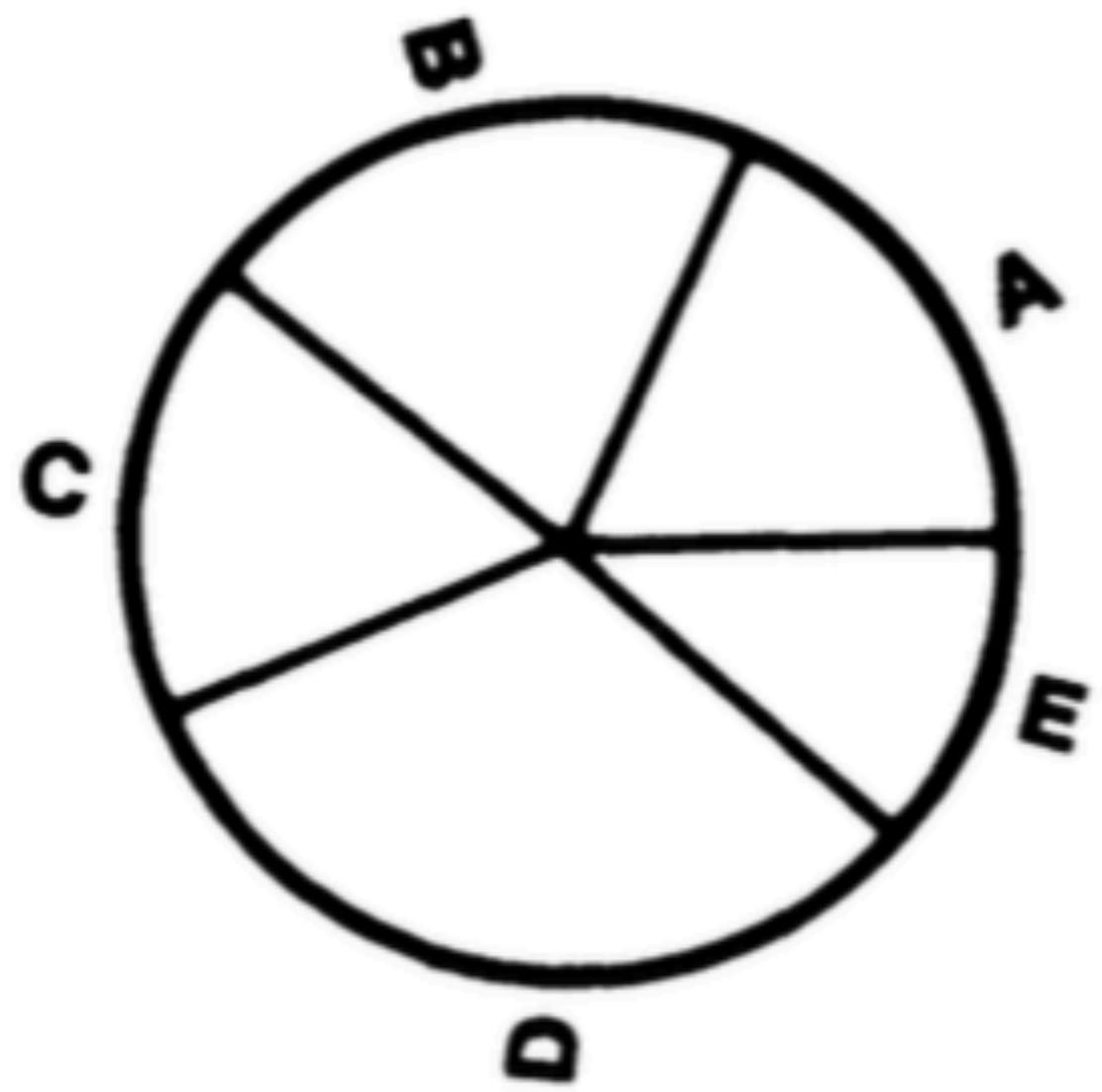
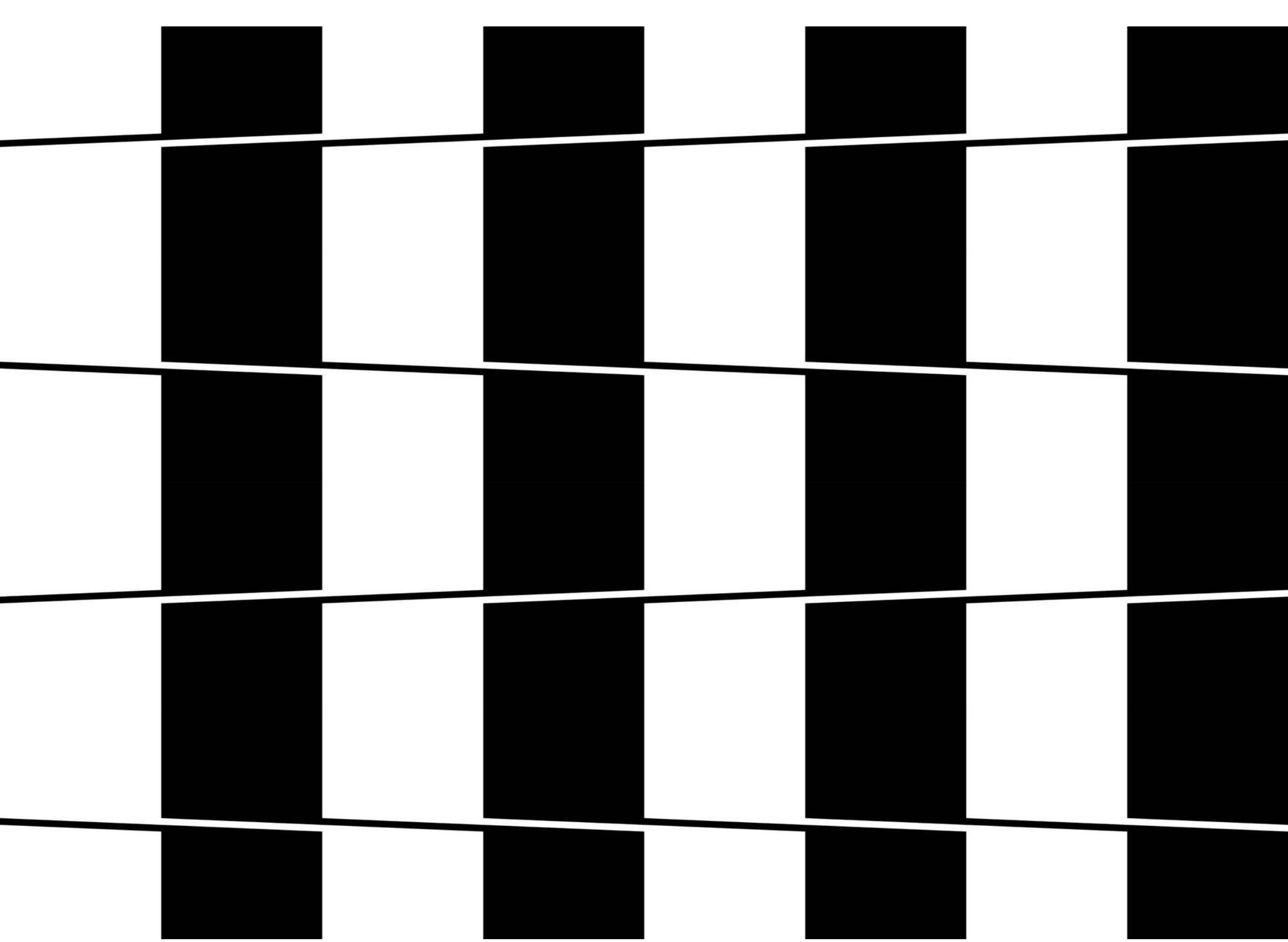


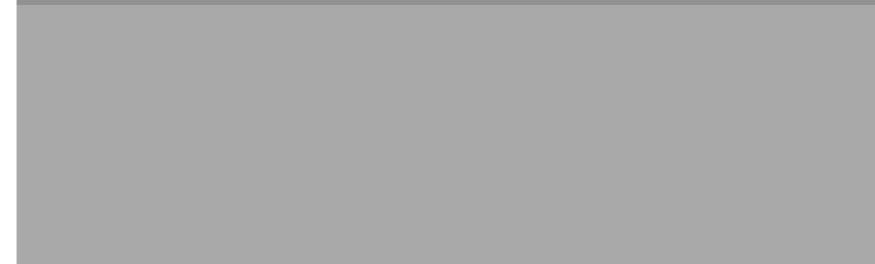
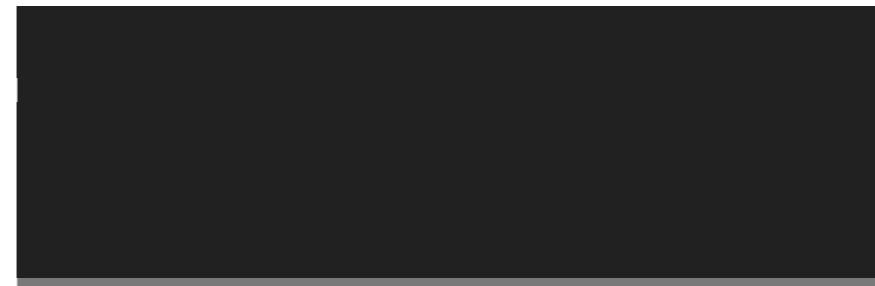
Figure 3. Graphs from position-angle experiment.

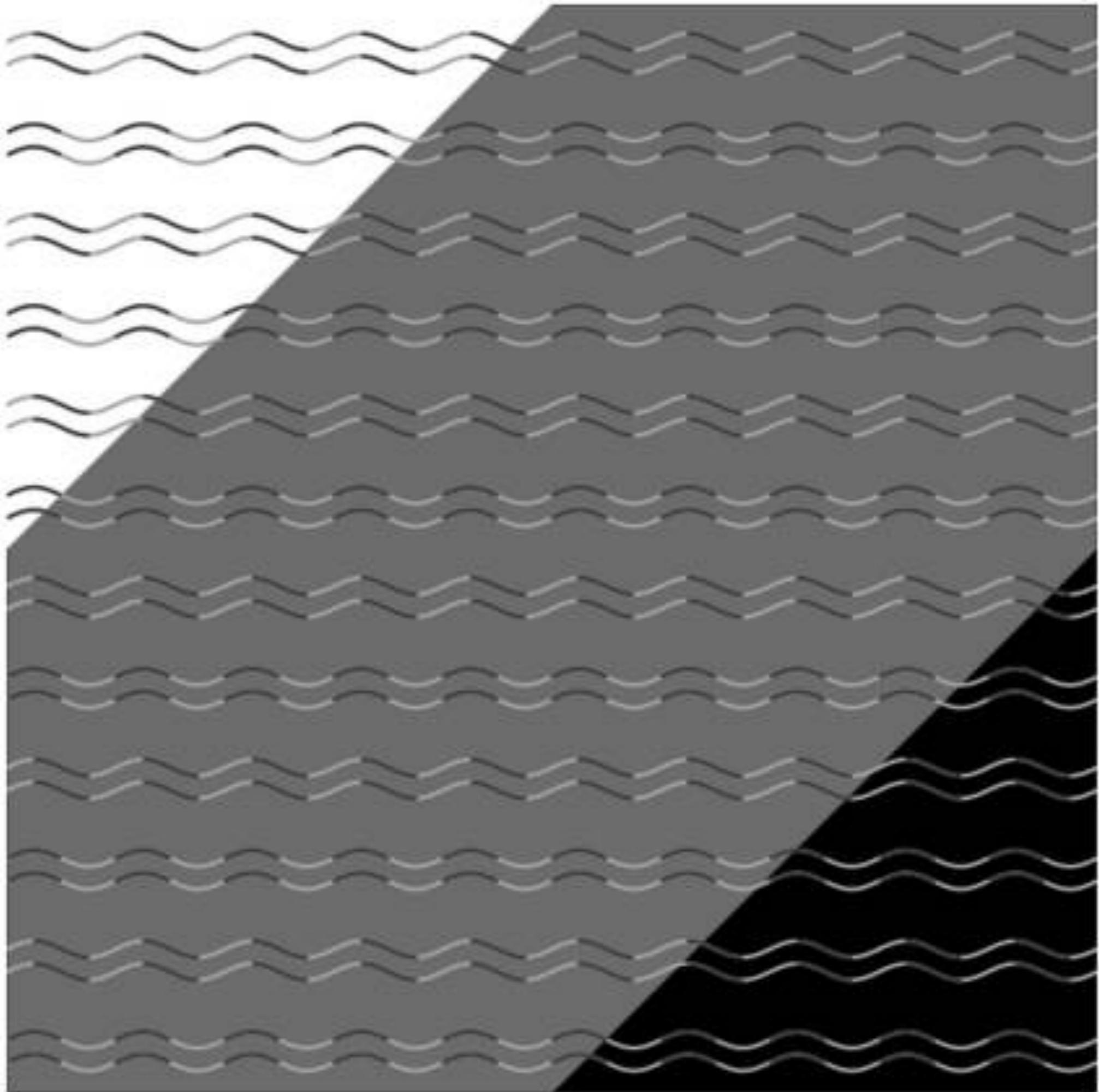


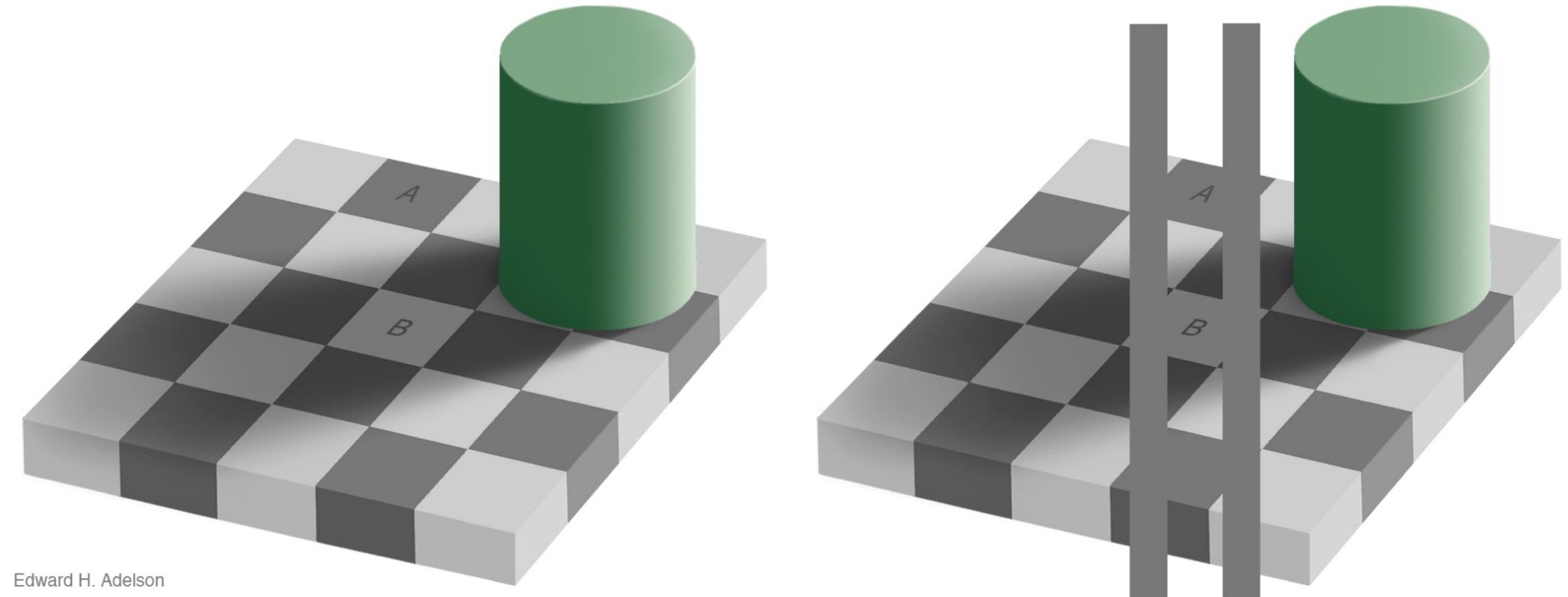
fine







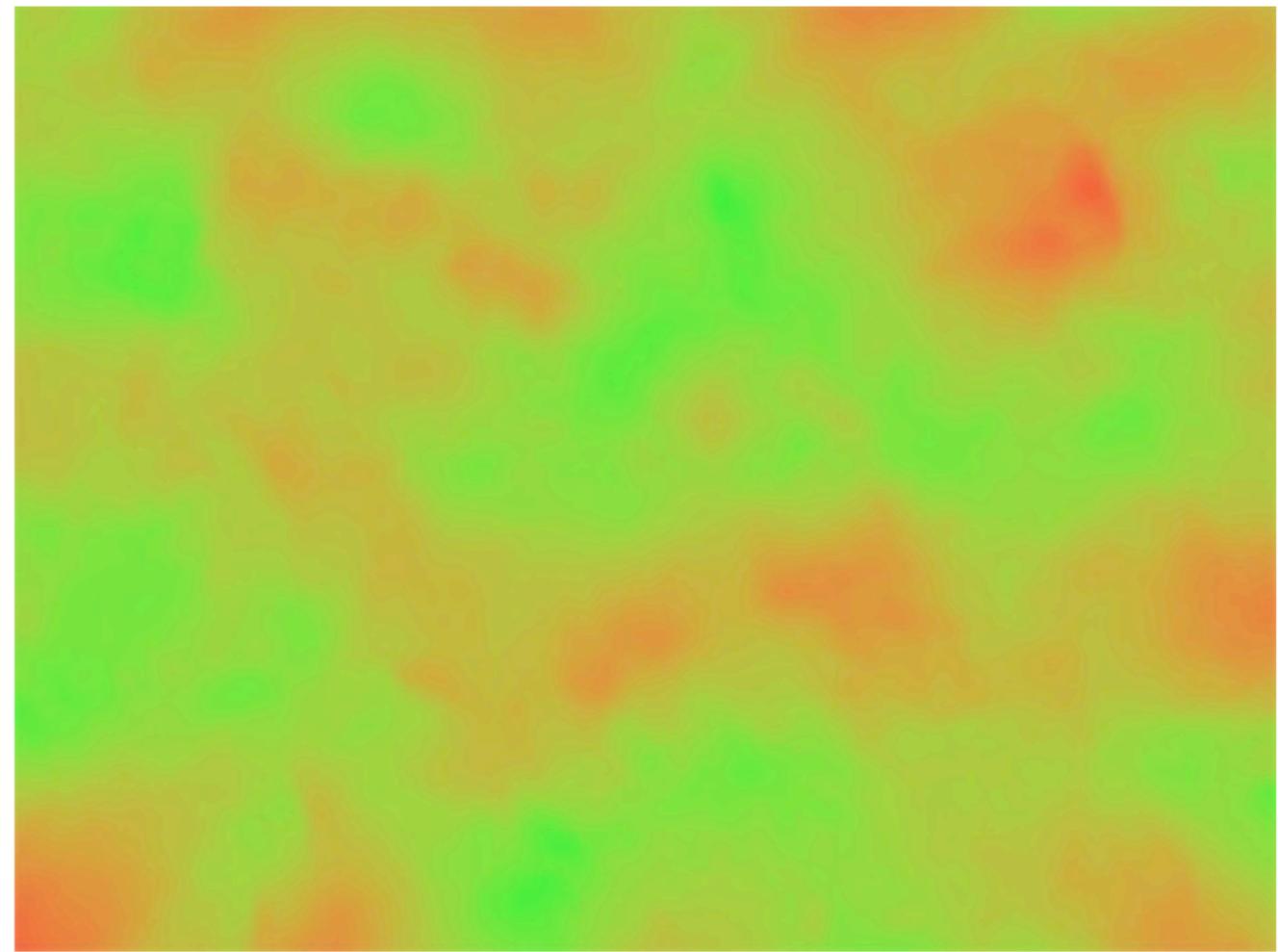
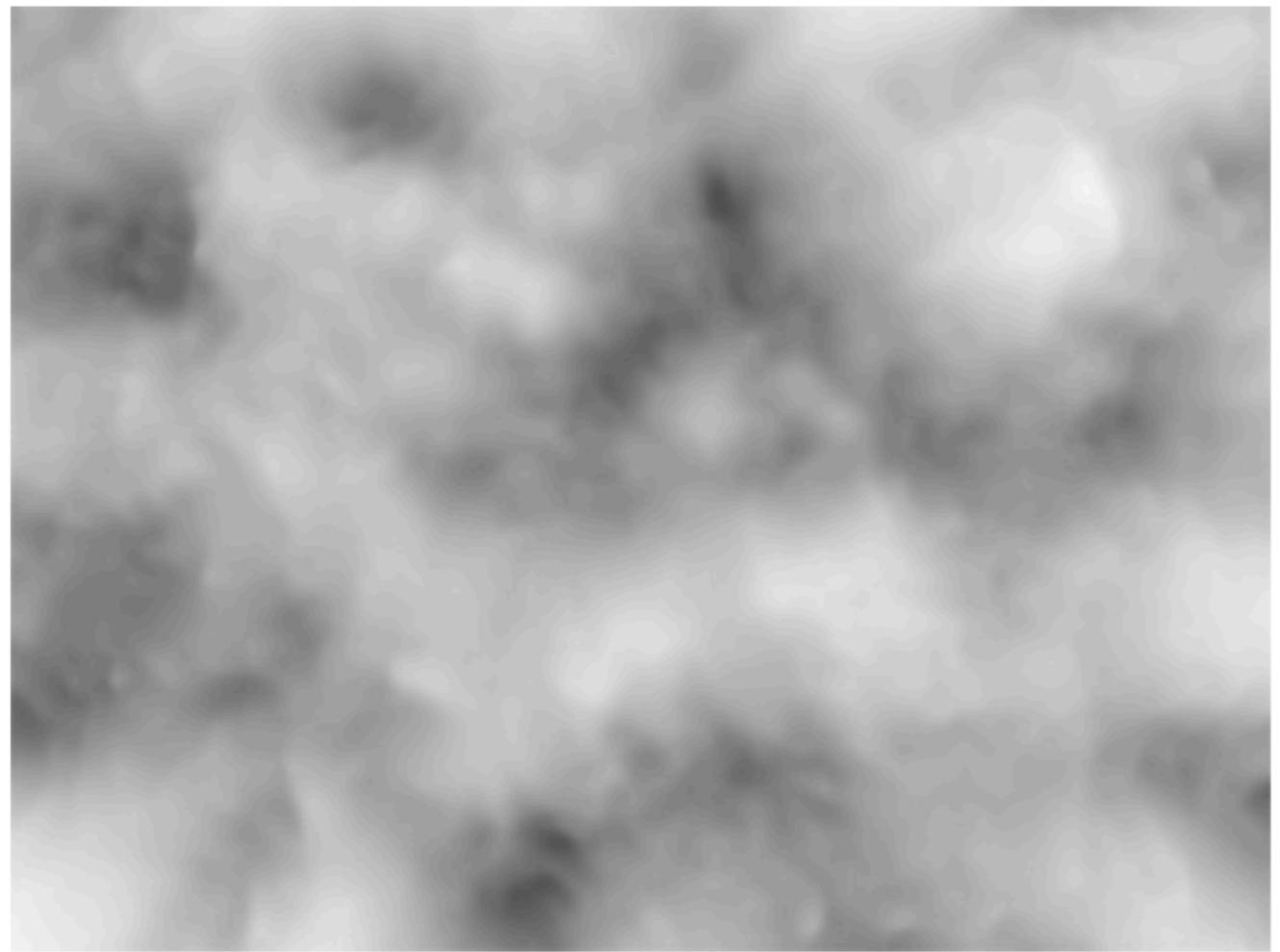




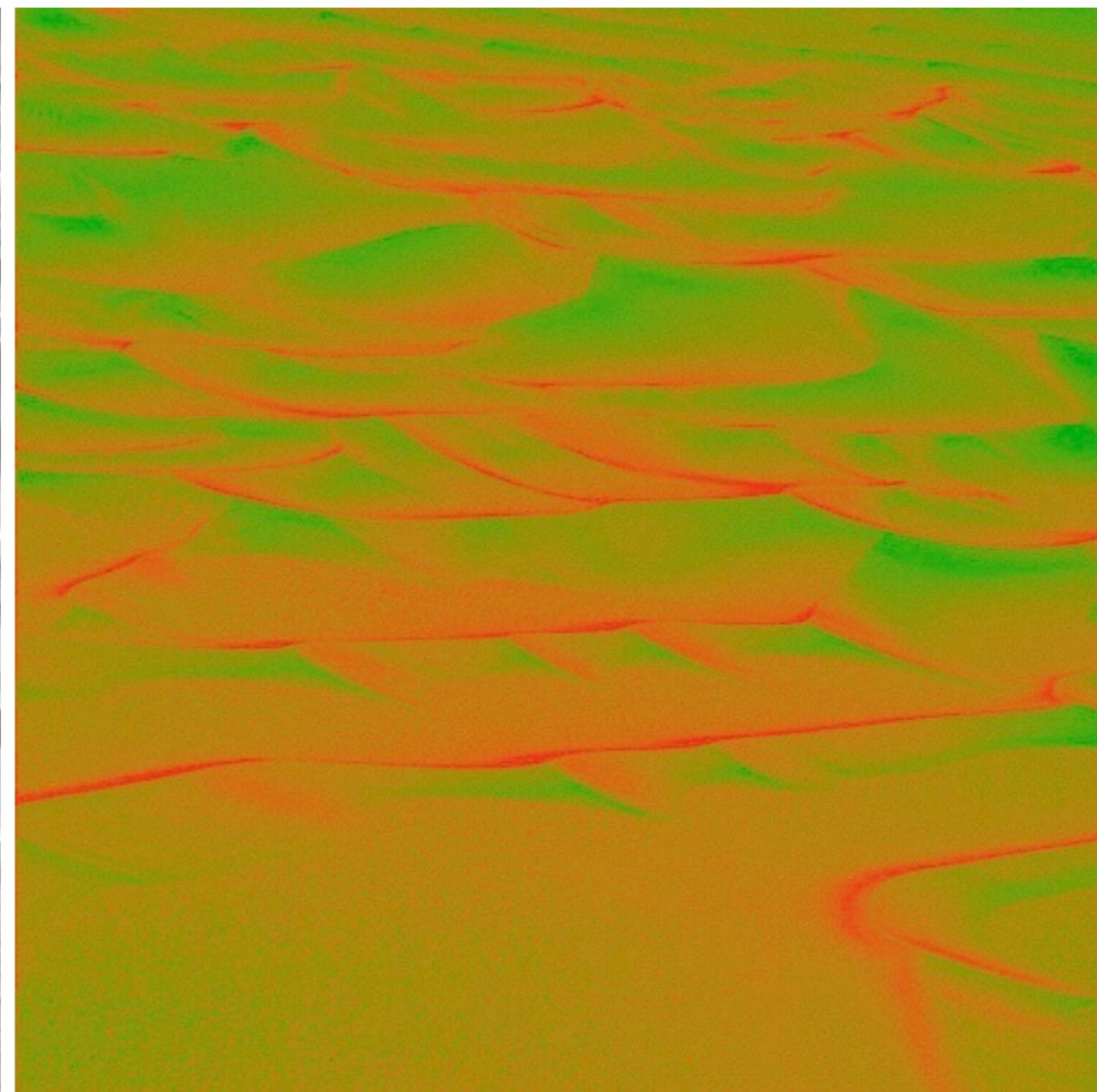
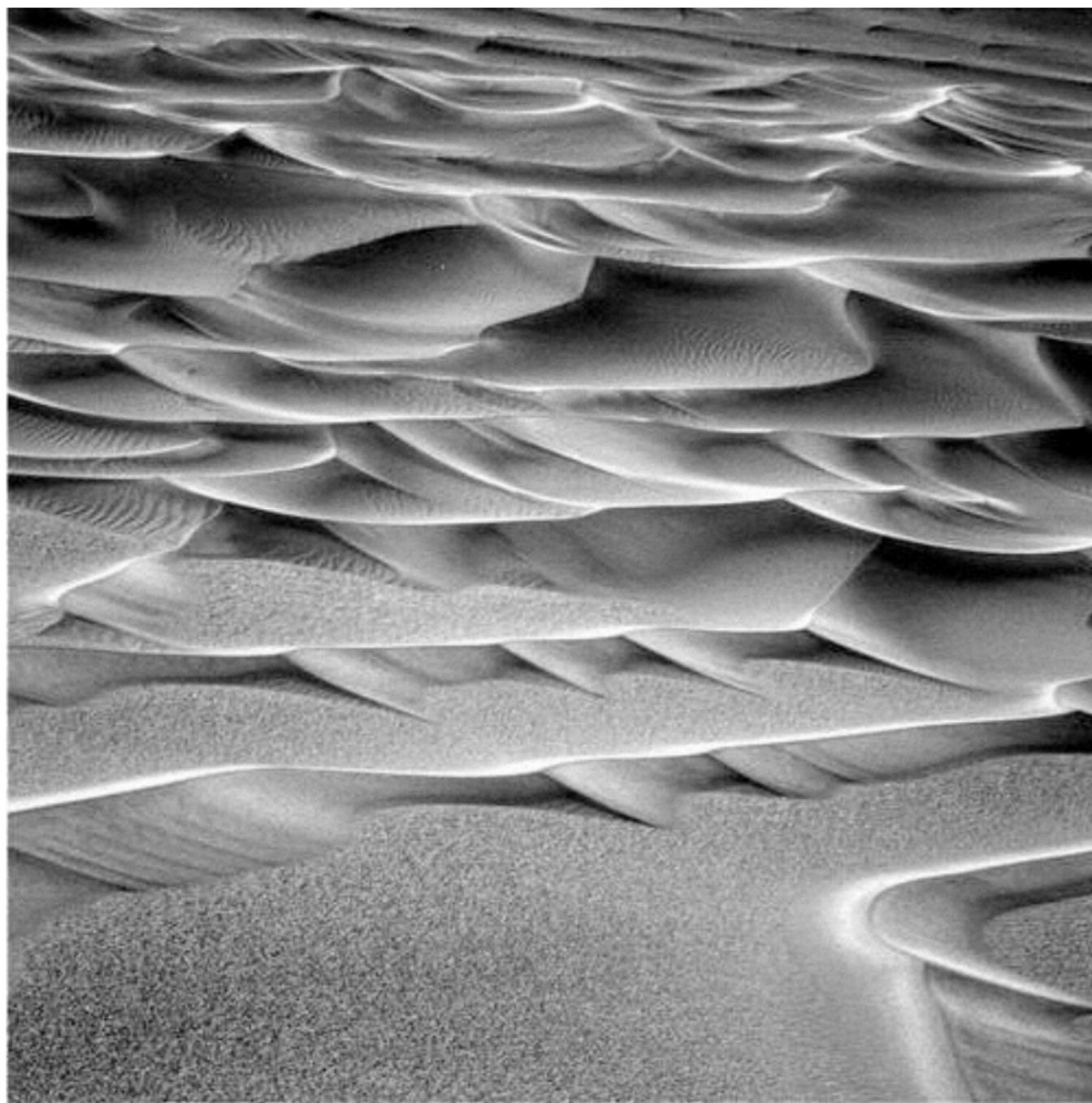
Edward H. Adelson

A

B

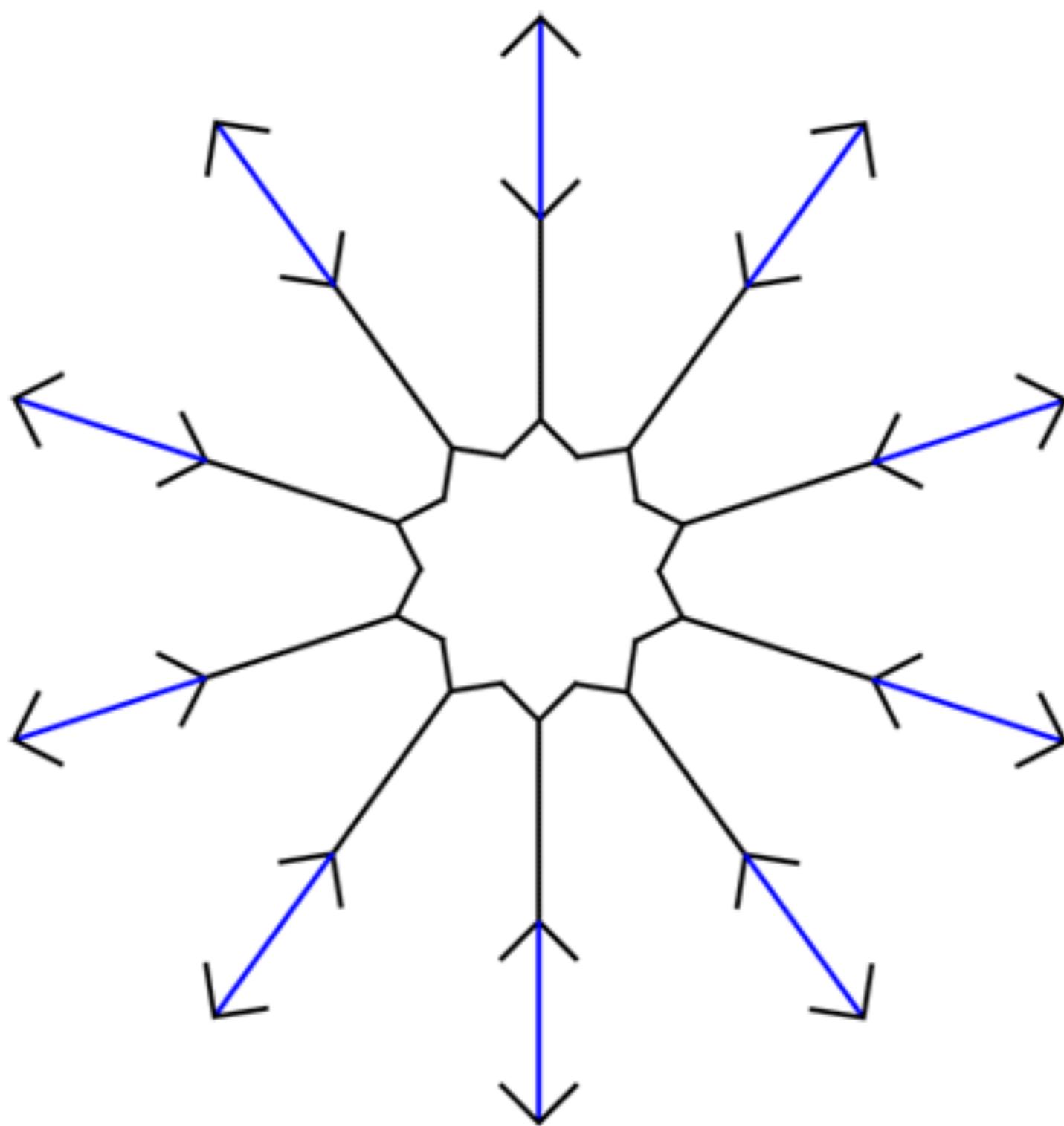


Ware (2008)



Ware (2008)

Sarcone's Dynamic Müller-Lyer Illusion



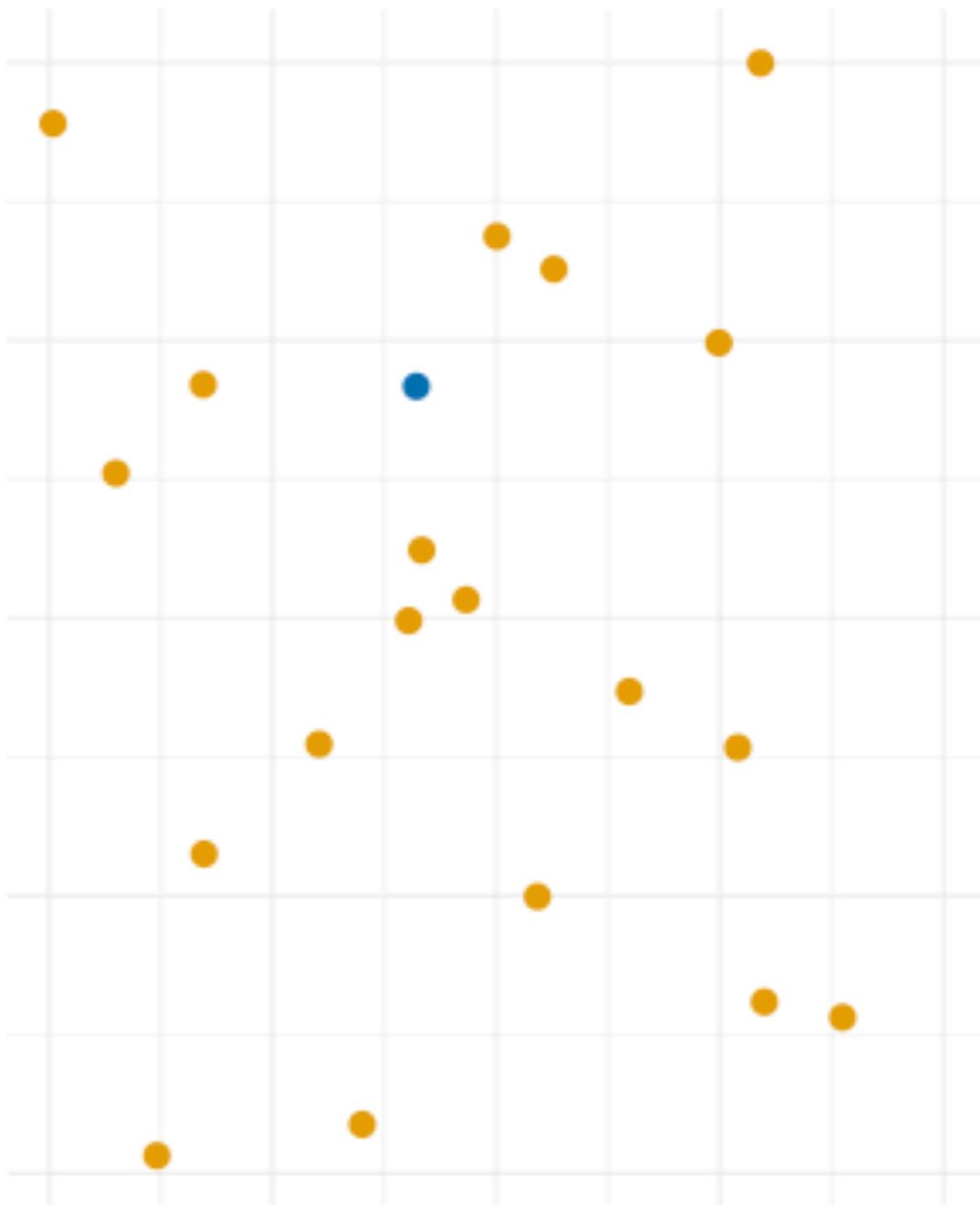
+



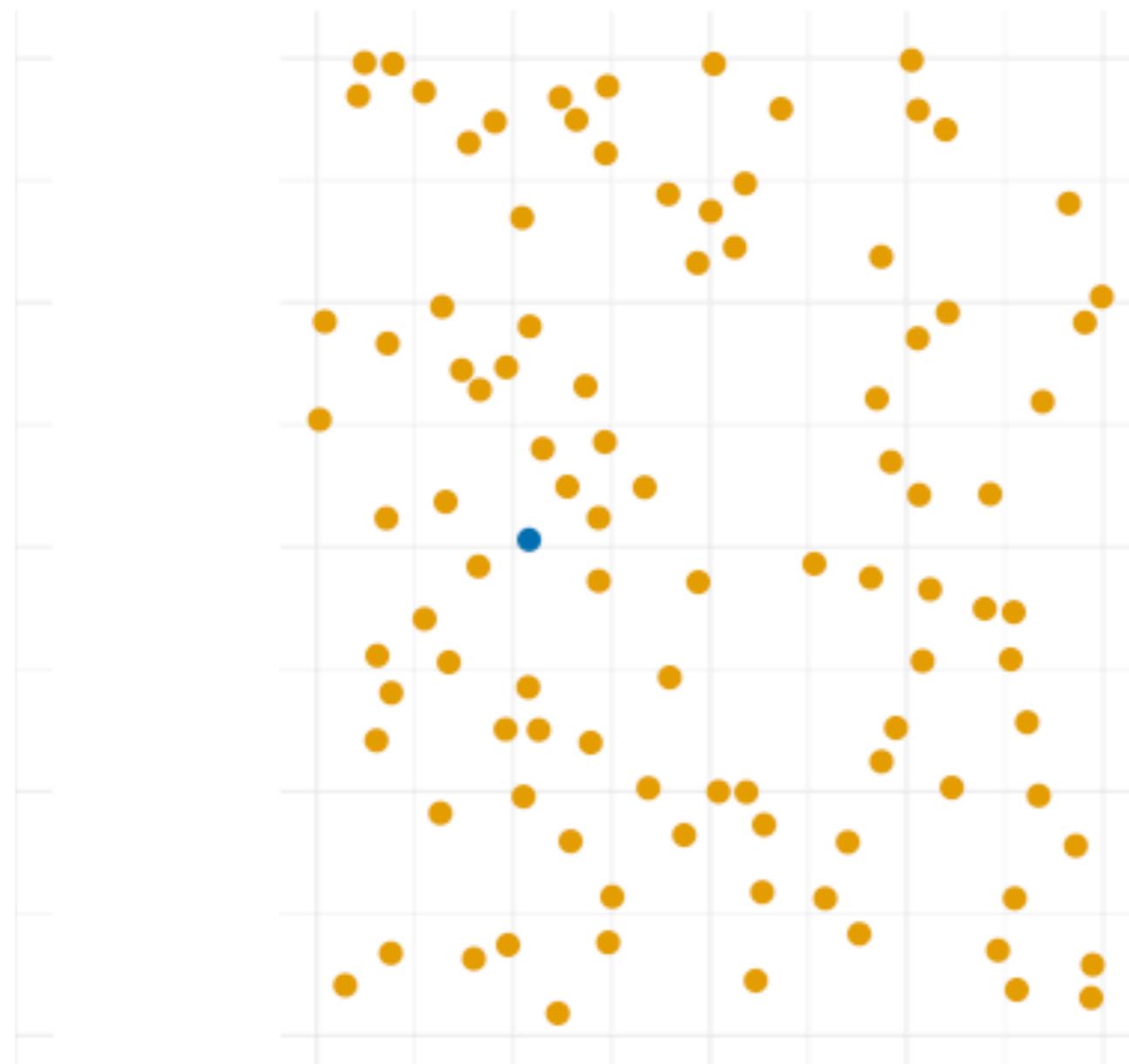


PRE-ATTENTIVE PROCESSING

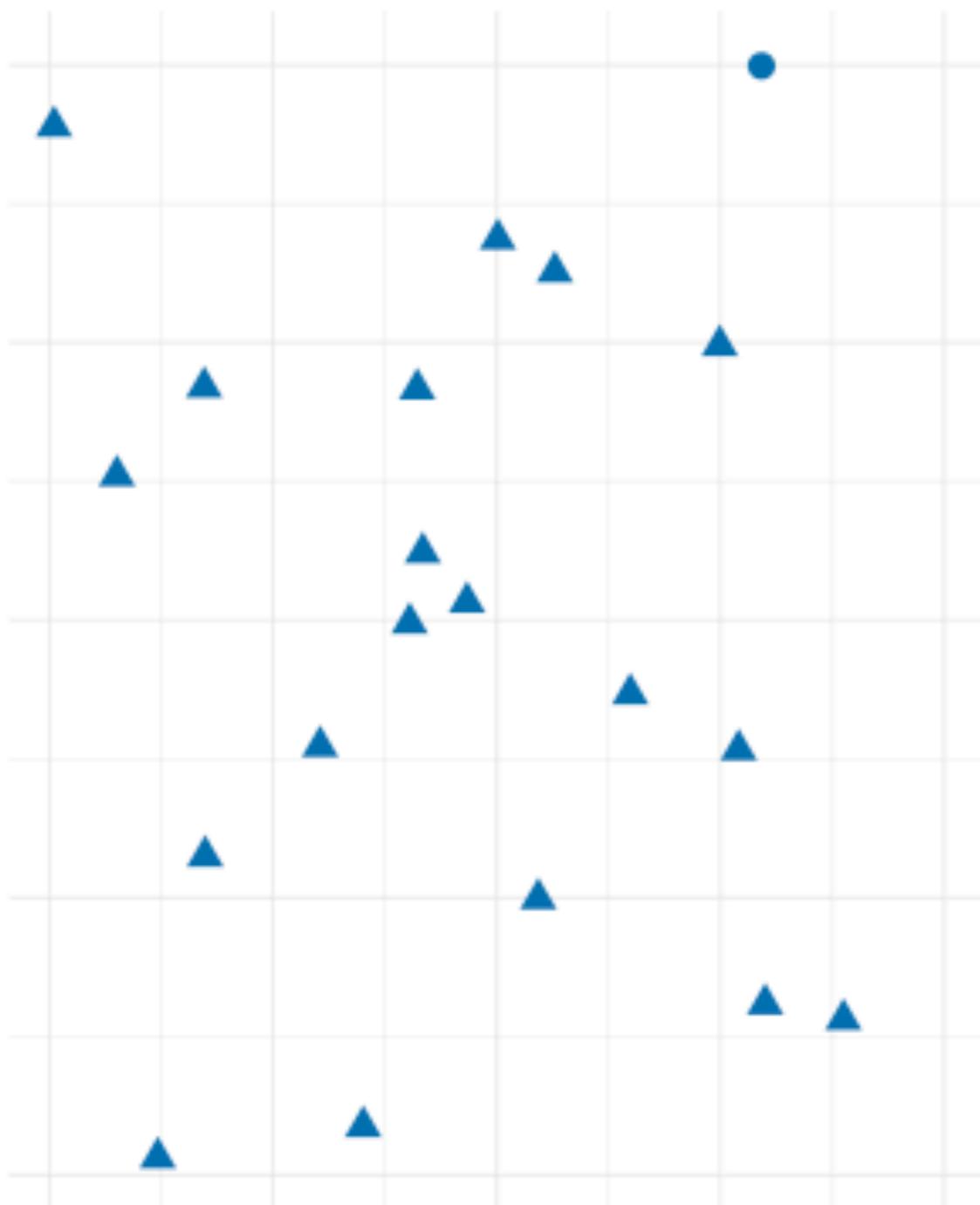
Color Only, N=20



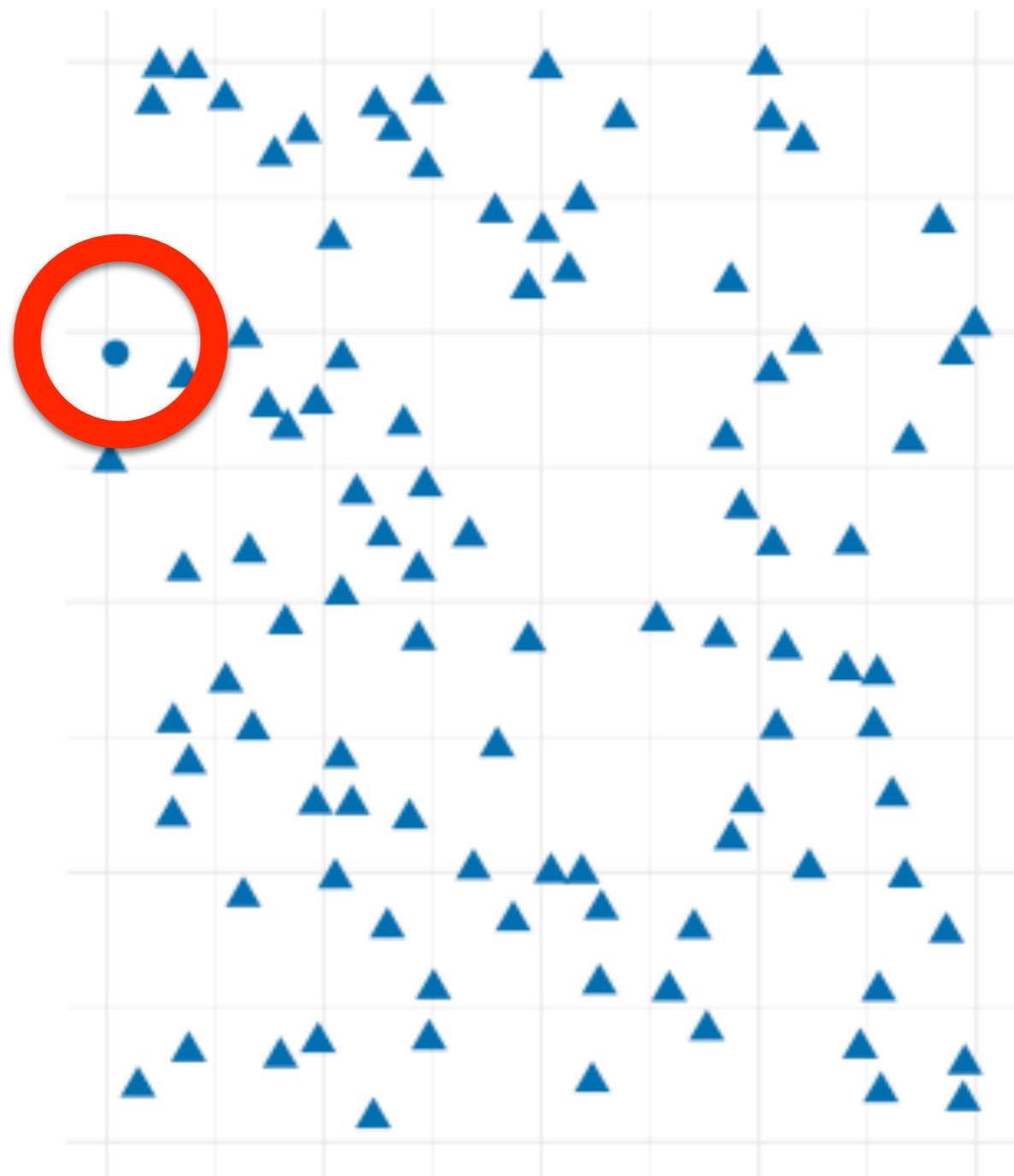
Color Only, N=100



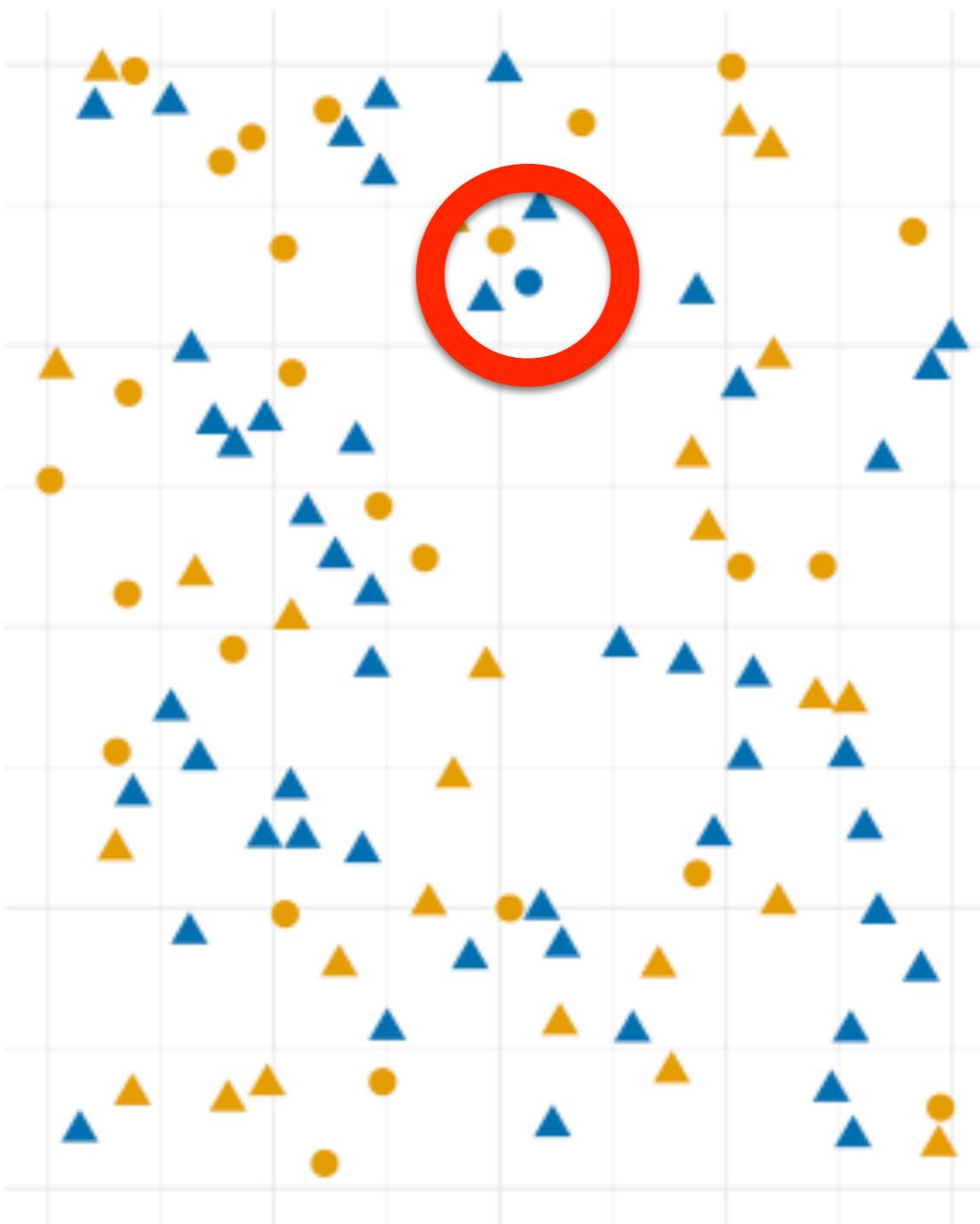
Shape Only, N=20

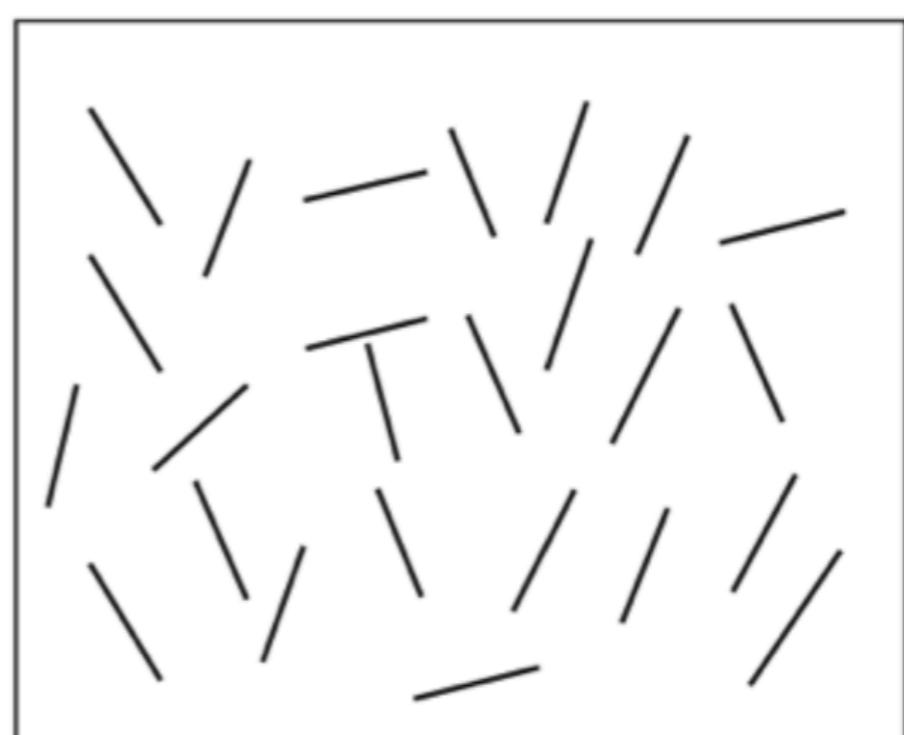
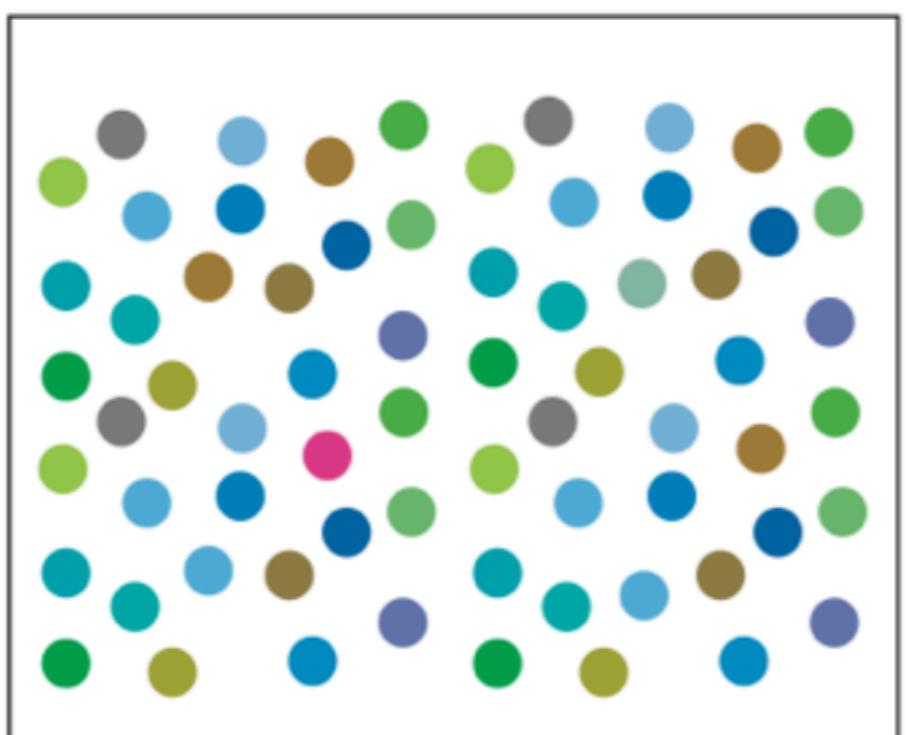
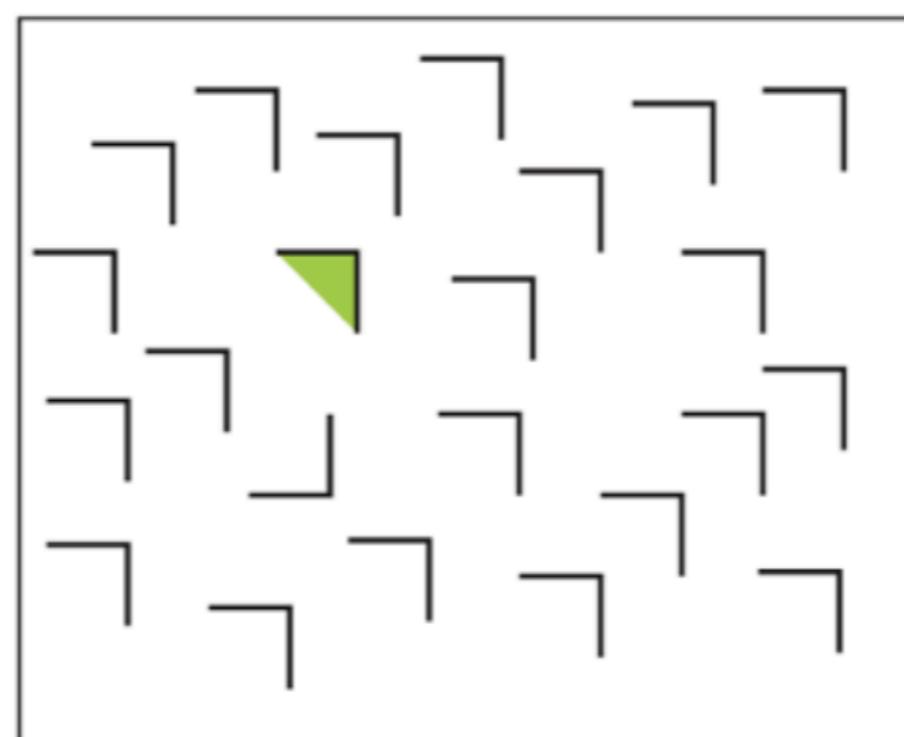
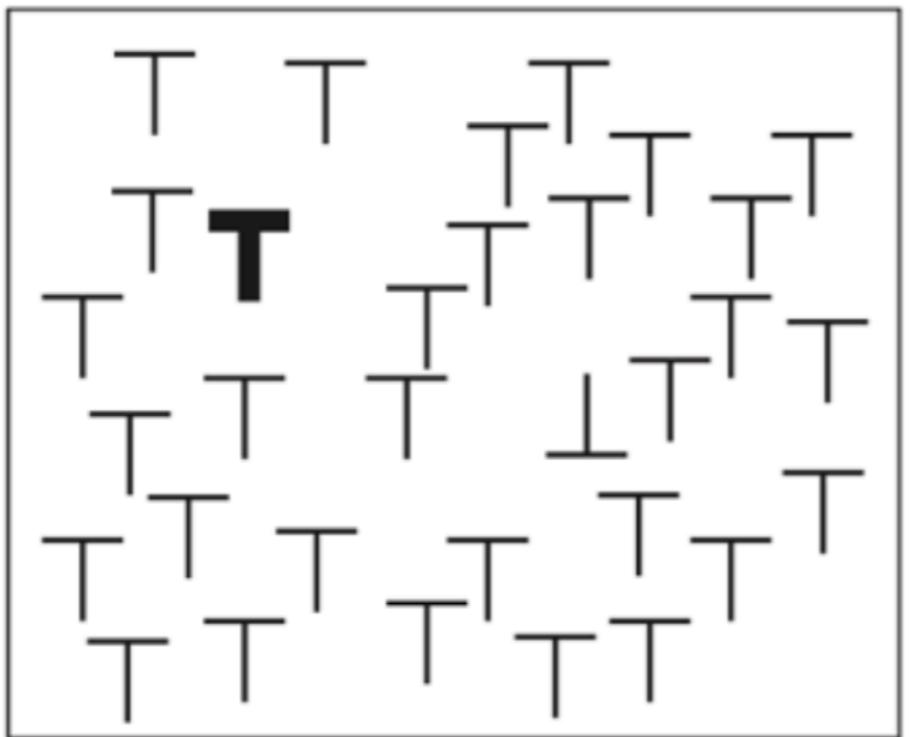


Shape Only, N=100

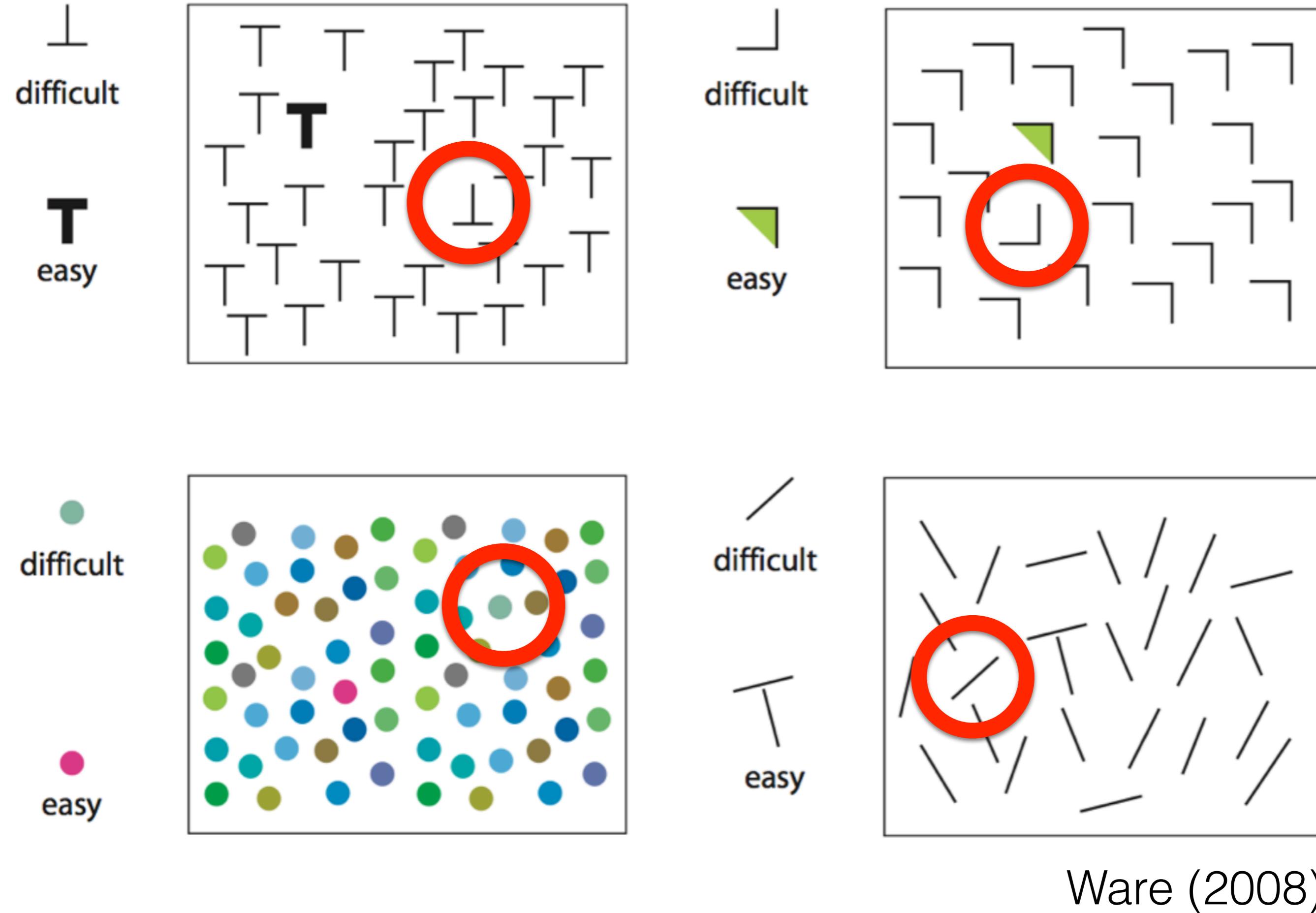


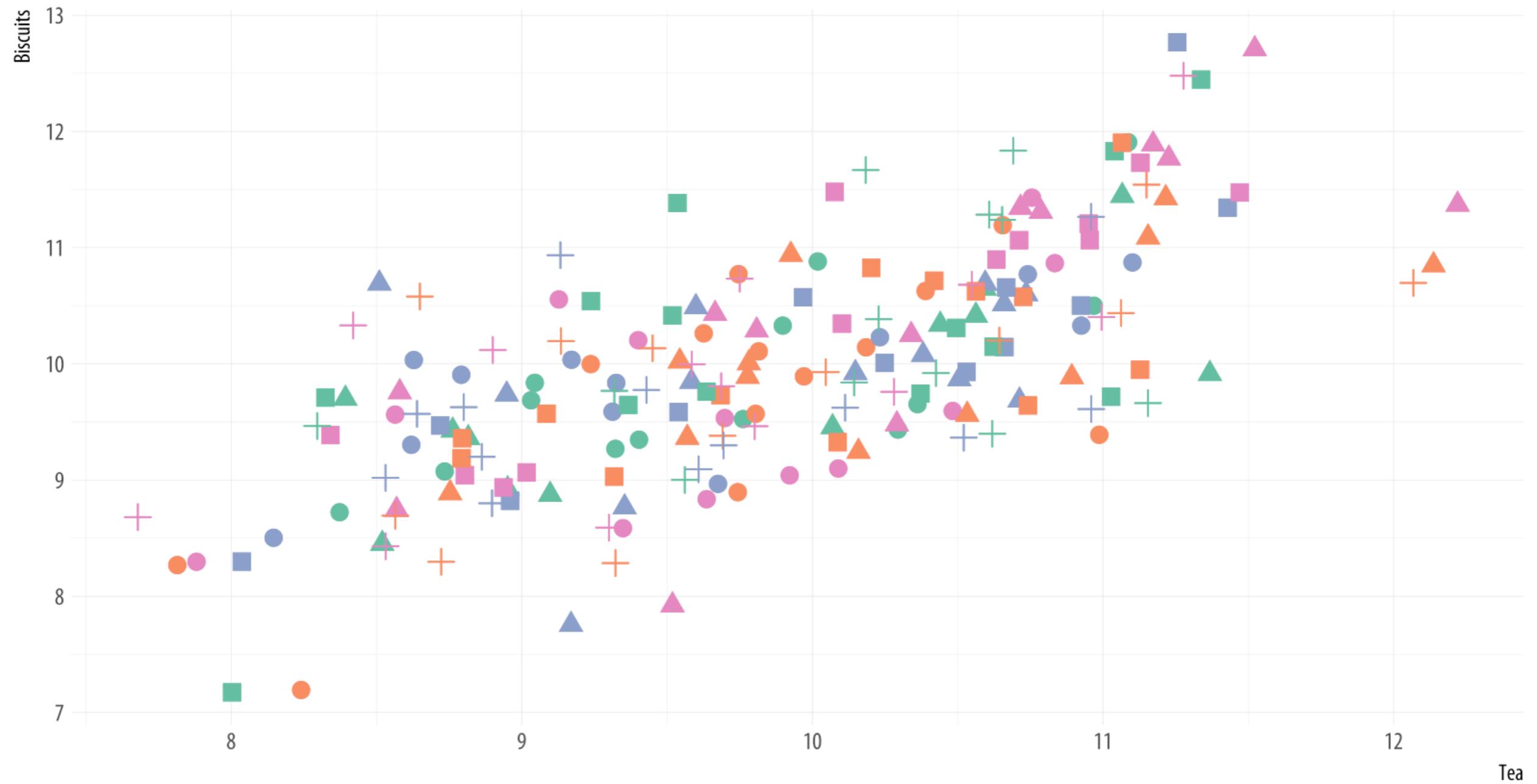
Color & Shape, N=100

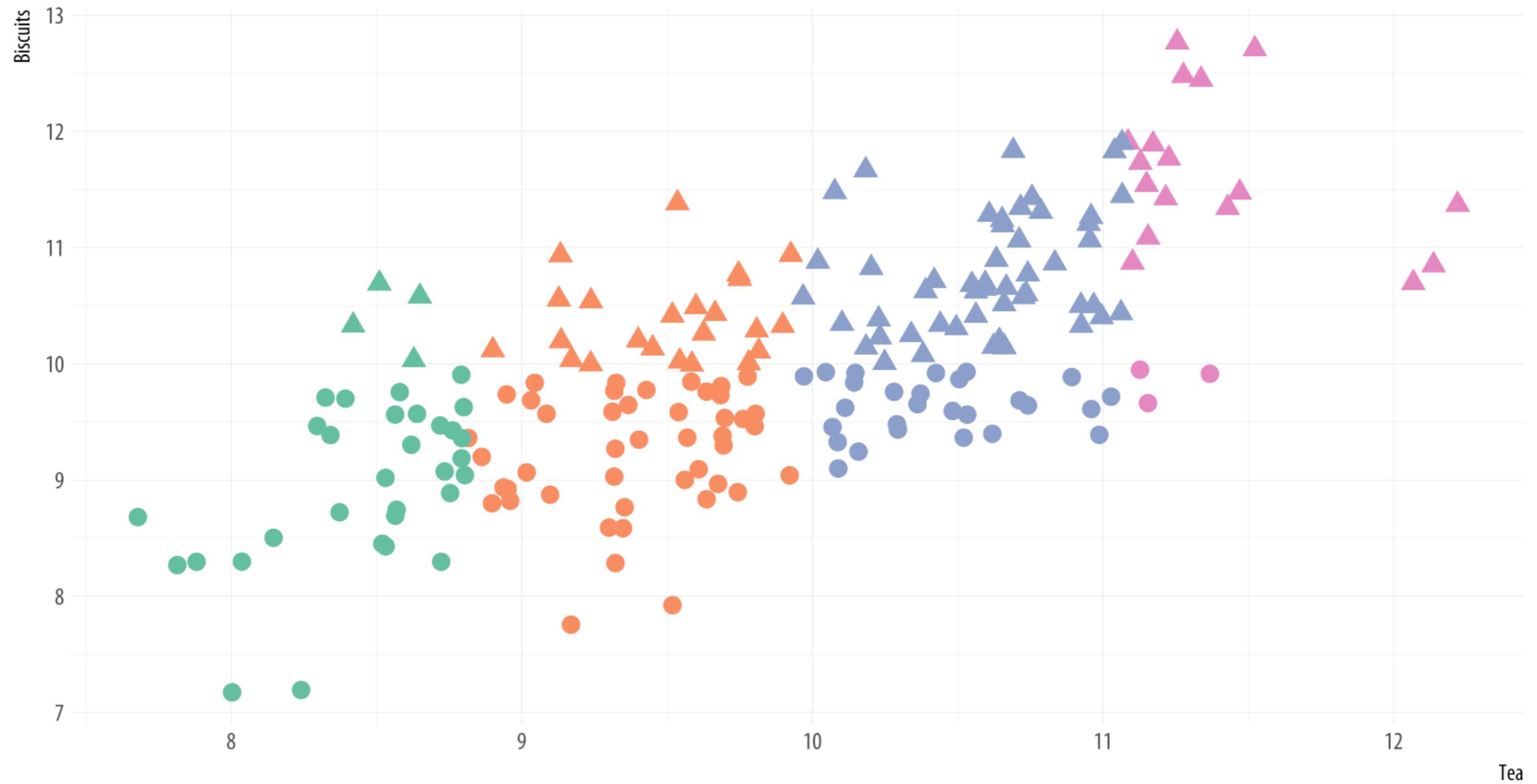




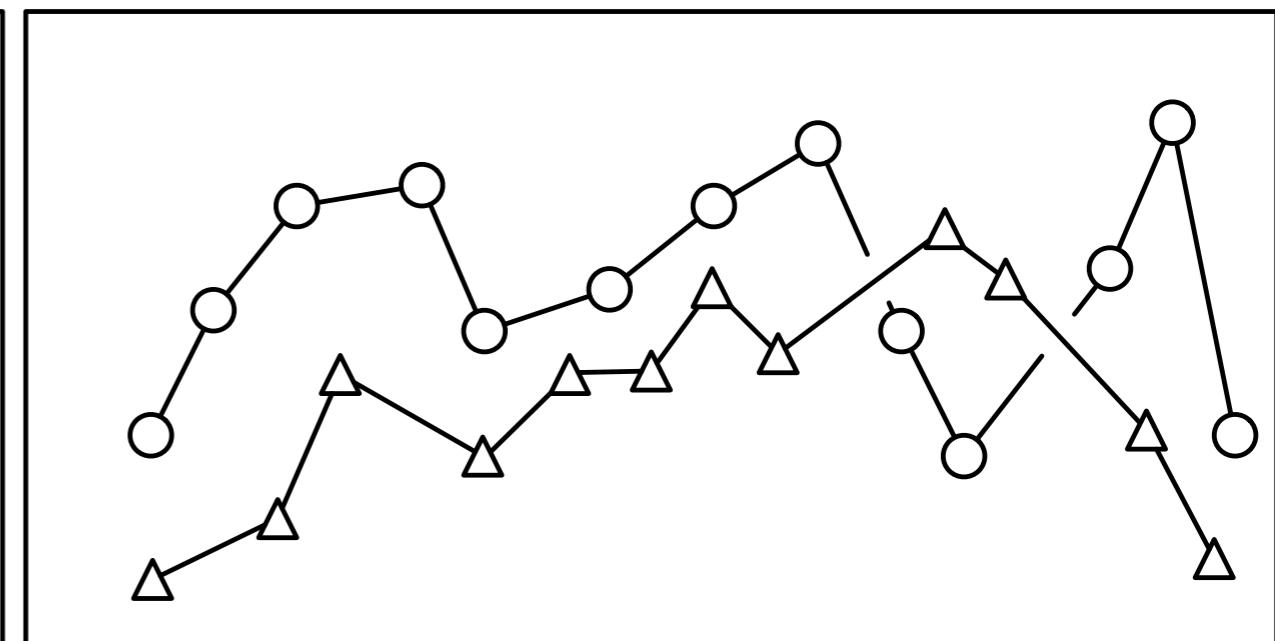
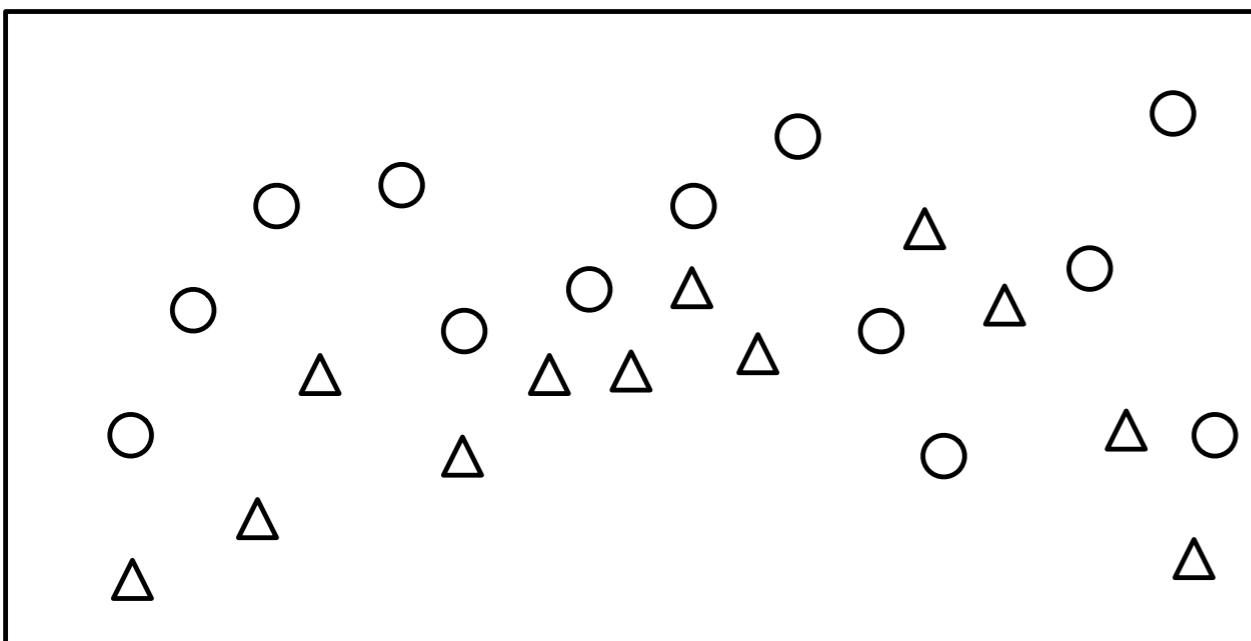
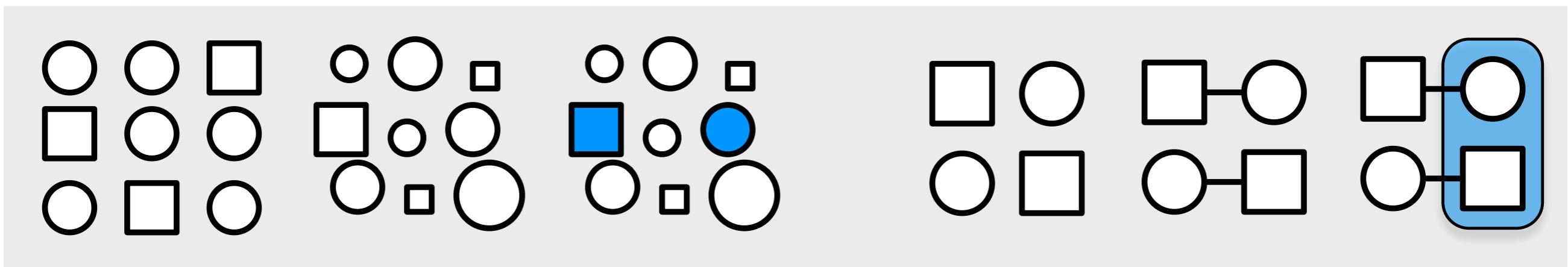
Ware (2008)



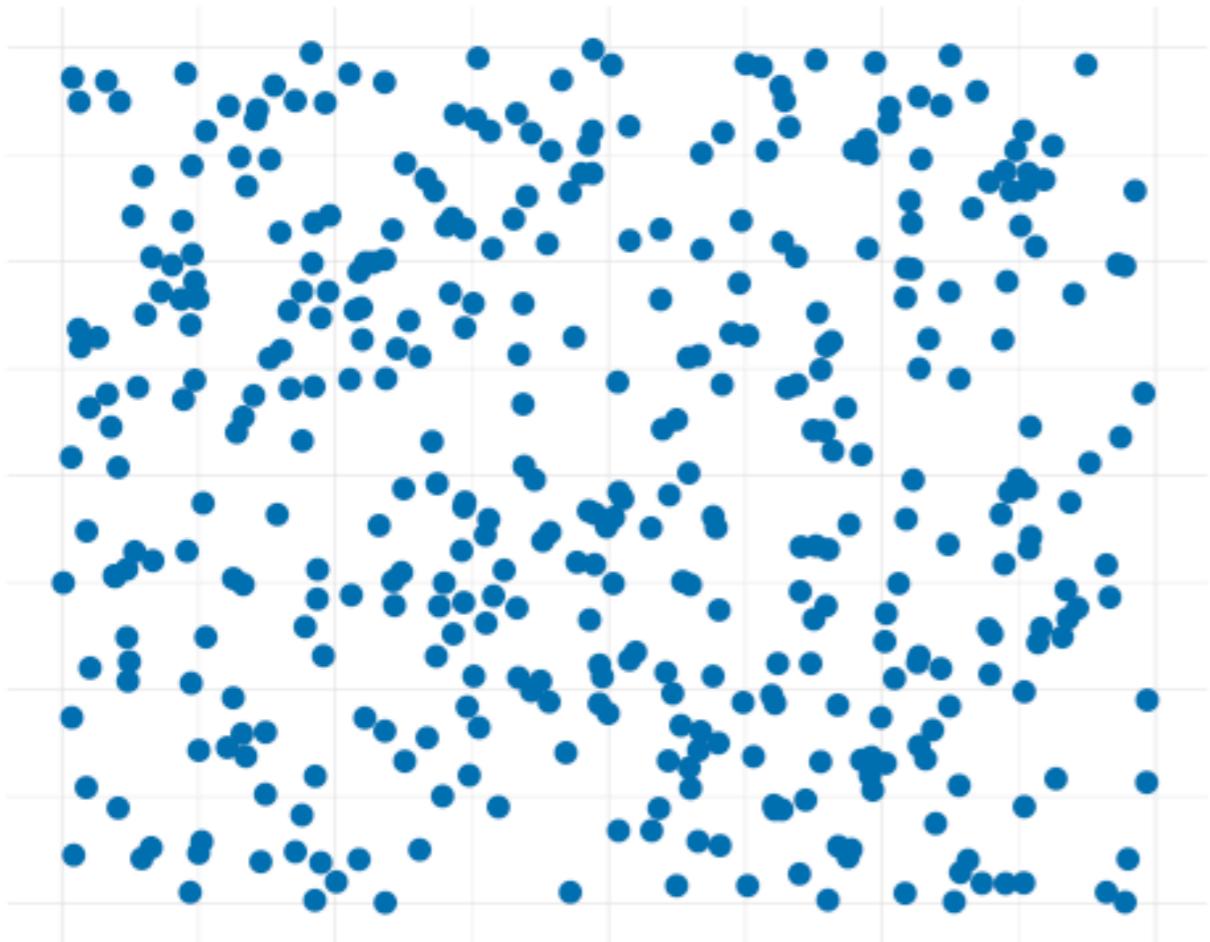




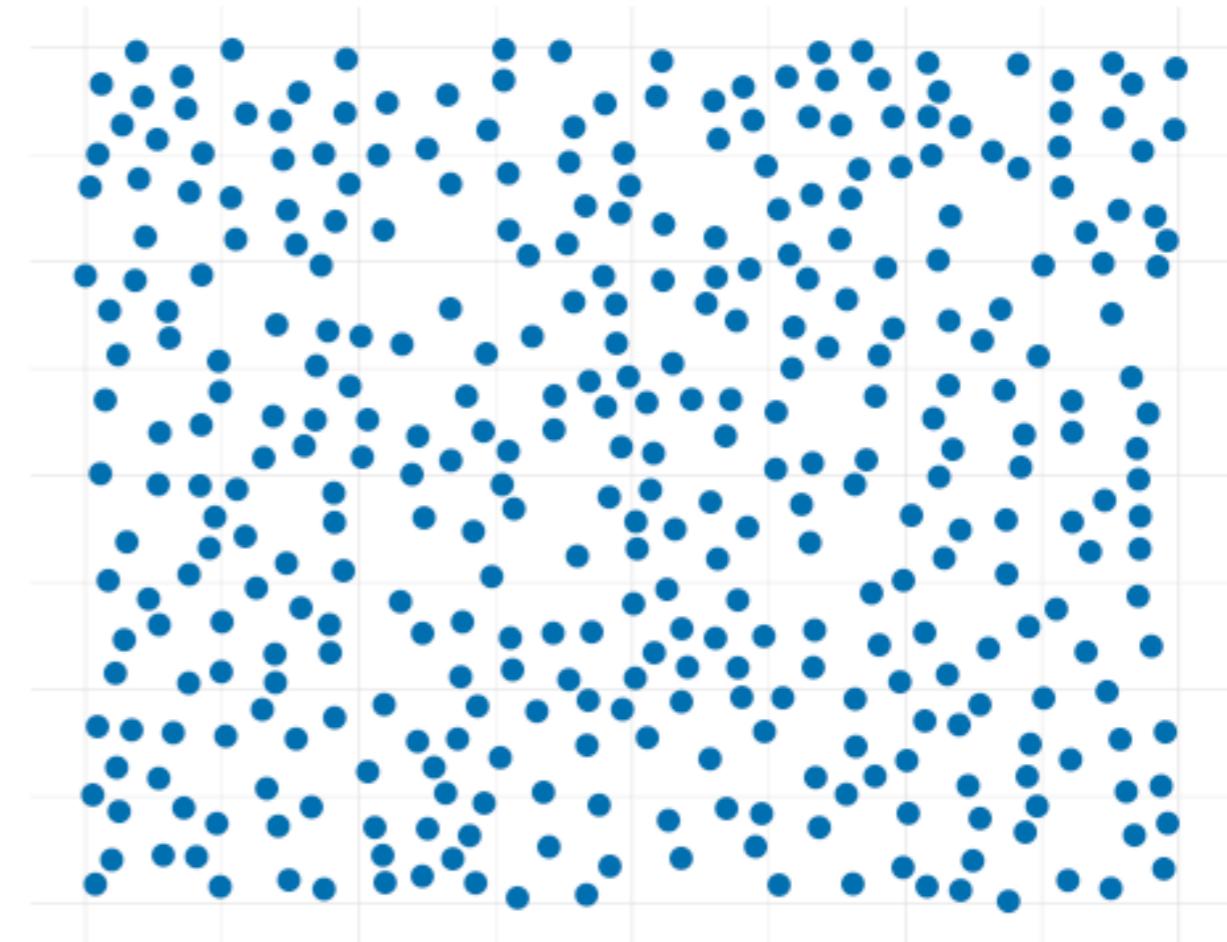
GESTALT INFERENCES



Poisson



Matérn



VISUAL TASKS IN DECODING GRAPHS

**GRAPH ELEMENTS
ENCODE QUANTITIES**

**VIEWERS MUST
THEN DECODE THEM**

QUALITIES & QUANTITIES

nominal
ordinal
interval
ratio

Stevens (1946)

names
grades
ranks
counted fractions
amounts
balances

Mosteller & Tukey (1977)

VISUAL ELEMENTS

position

volume

color

length

shape

motion

angle

gradient

area

pattern

QUALITIES & QUANTITIES

I

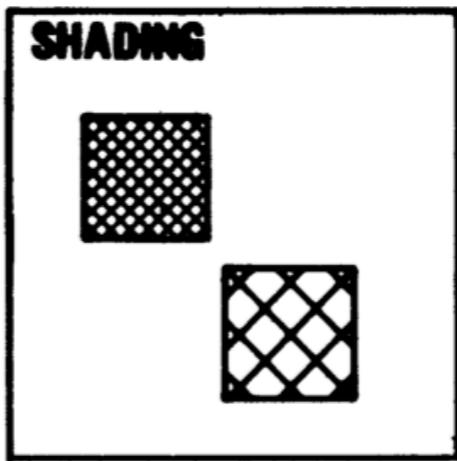
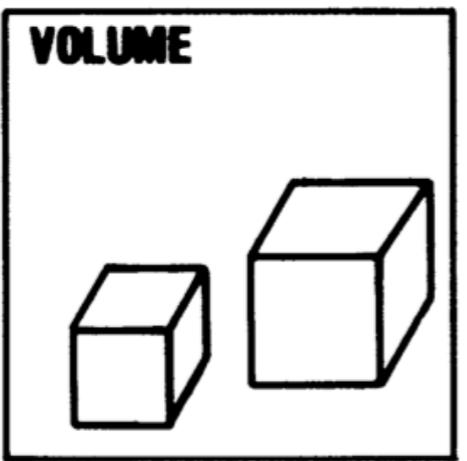
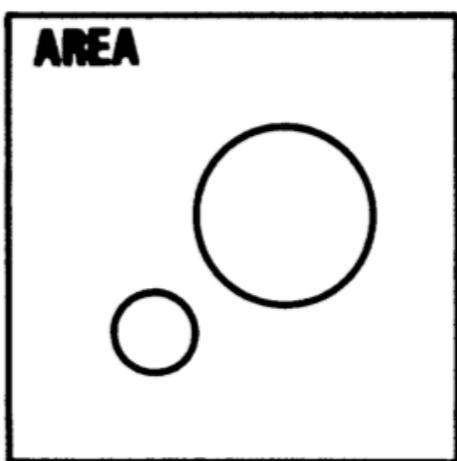
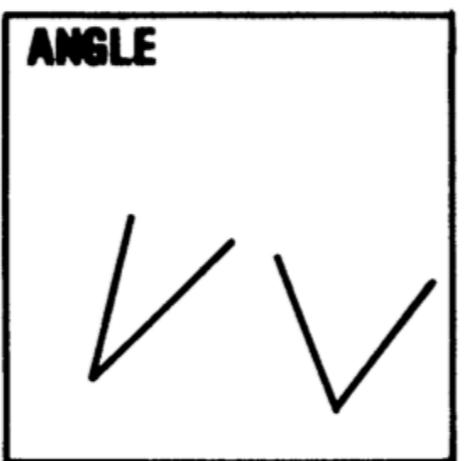
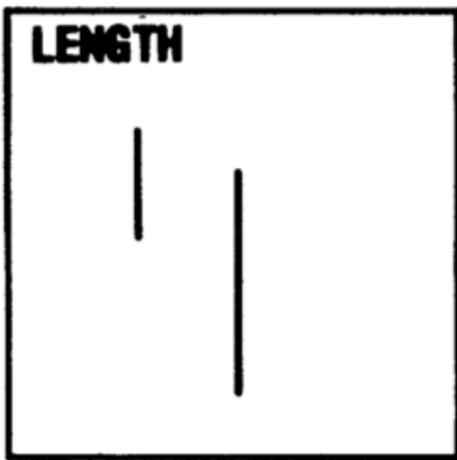
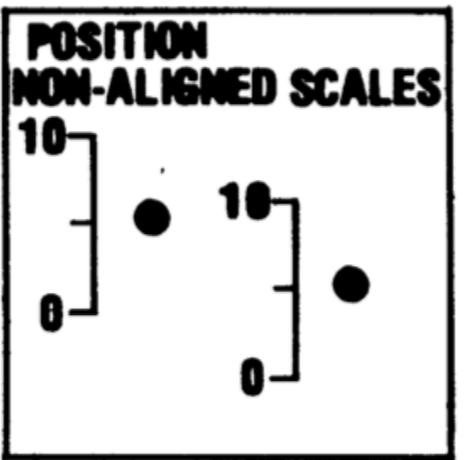
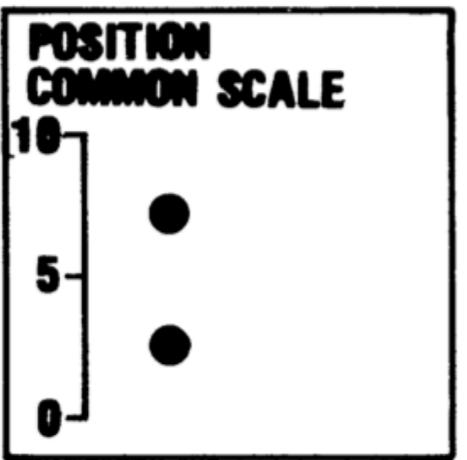
Mapping



Decoding

I

VISUAL ELEMENTS



COLOR SATURATION

Figure 1. Elementary perceptual tasks.

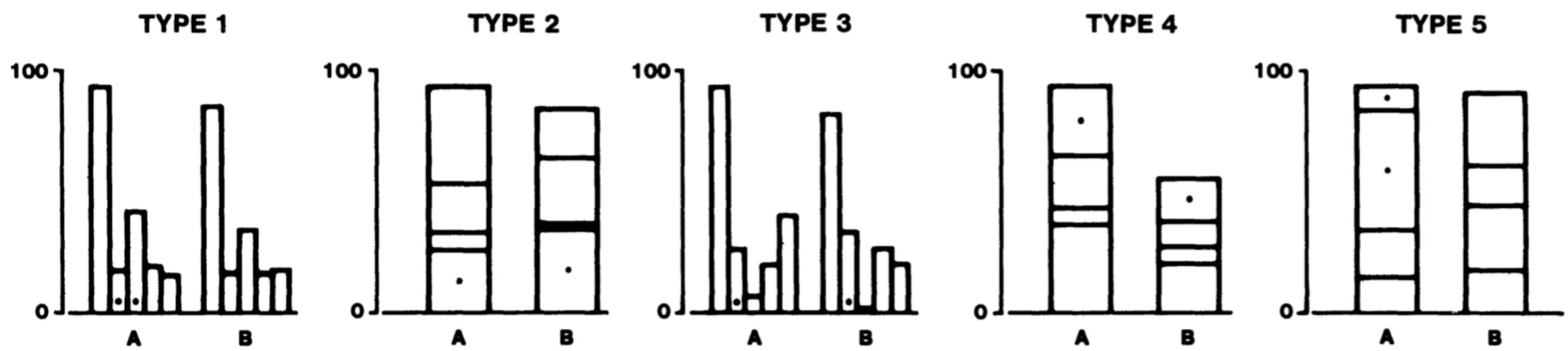
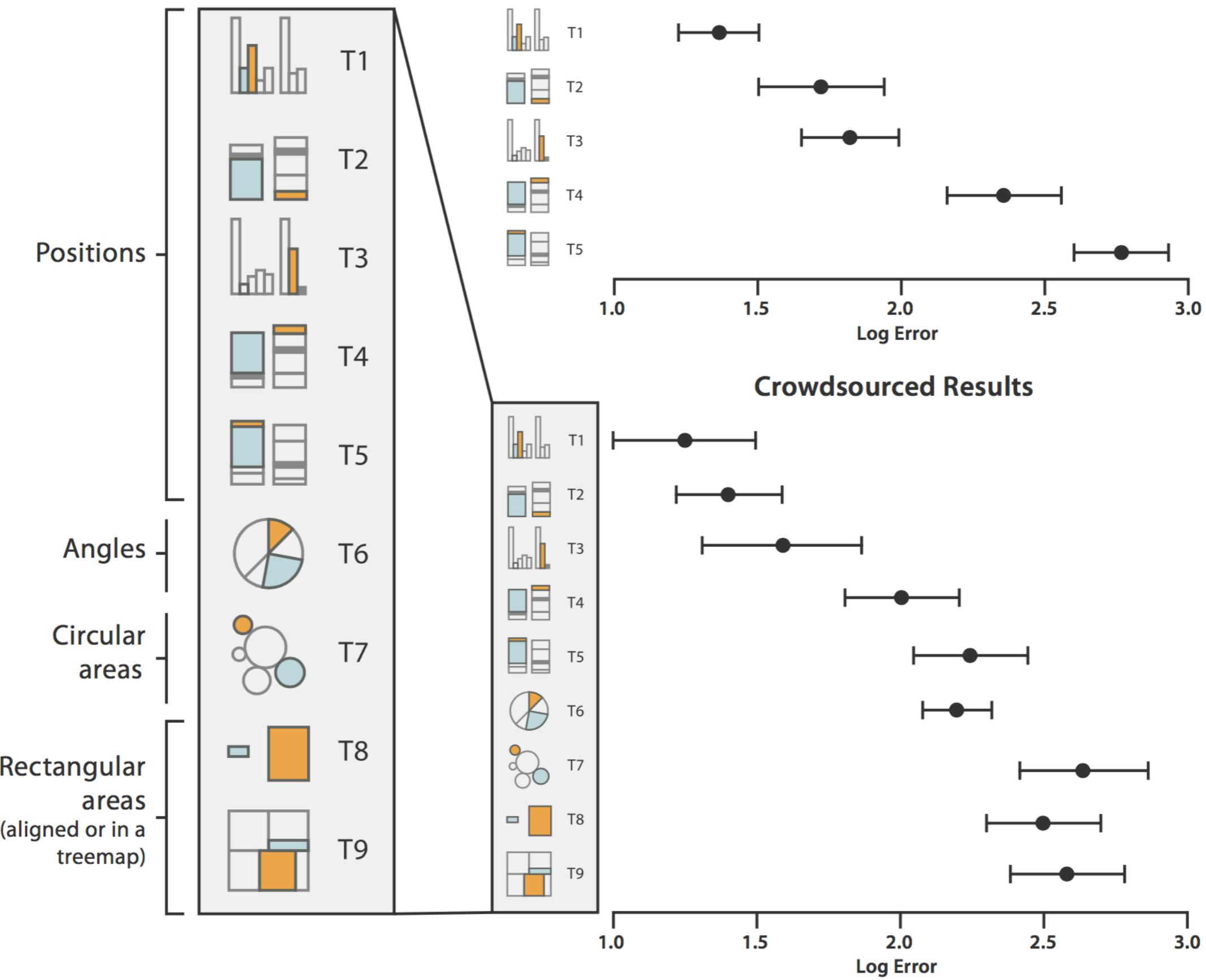


Figure 4. Graphs from position-length experiment.

Cleveland & McGill's Results



See: Heer & Bostock (2010)

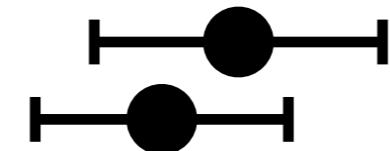
**ENCODINGS
or MAPPINGS
FOR DATA**

Position on a common scale



Better

Position on unaligned scale



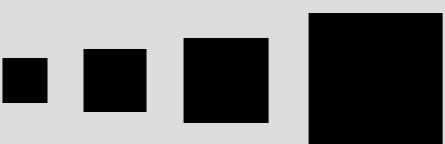
Length (1D as size)



Tilt or Angle



Area (2D as size)

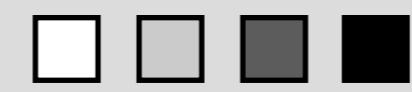


Effectiveness

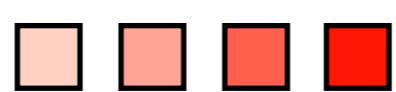
Depth (3D as Position)



Color luminance [brightness]



Color saturation [intensity]



Curvature

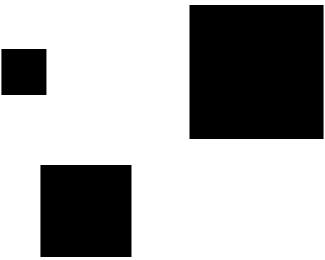


Worse

Volume (3D as size)

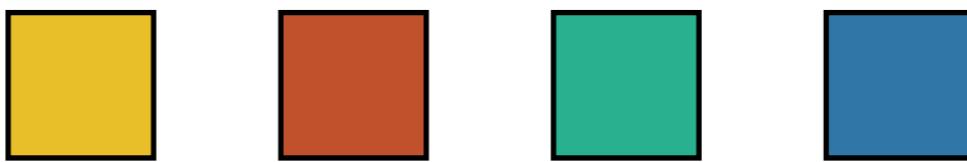


Spatial Region



Better

Color [hue]



Effectiveness

Motion



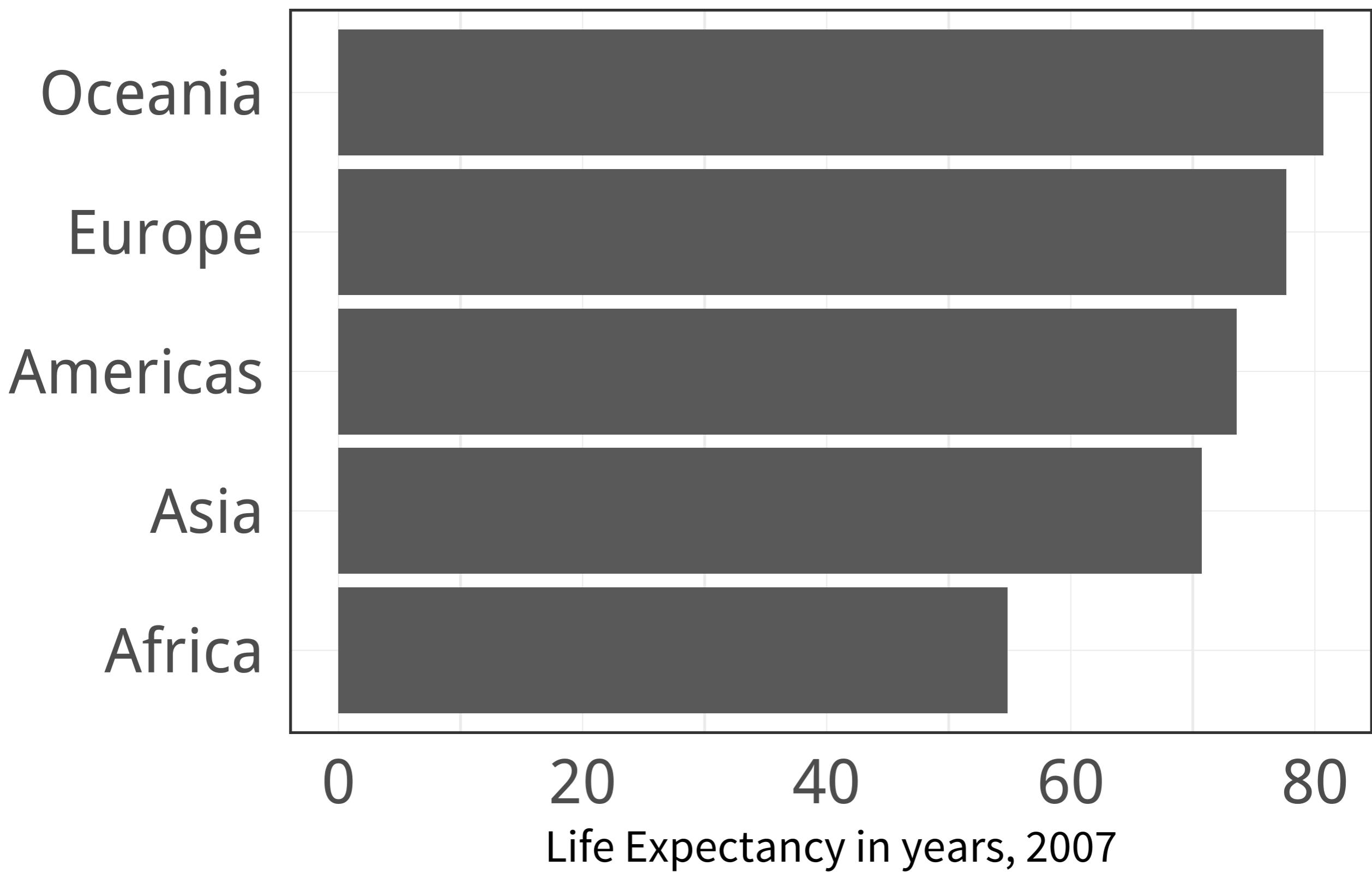
Shape



Worse

For more see: Munzer (2014)

**HONESTY AND
GOOD JUDGMENT**



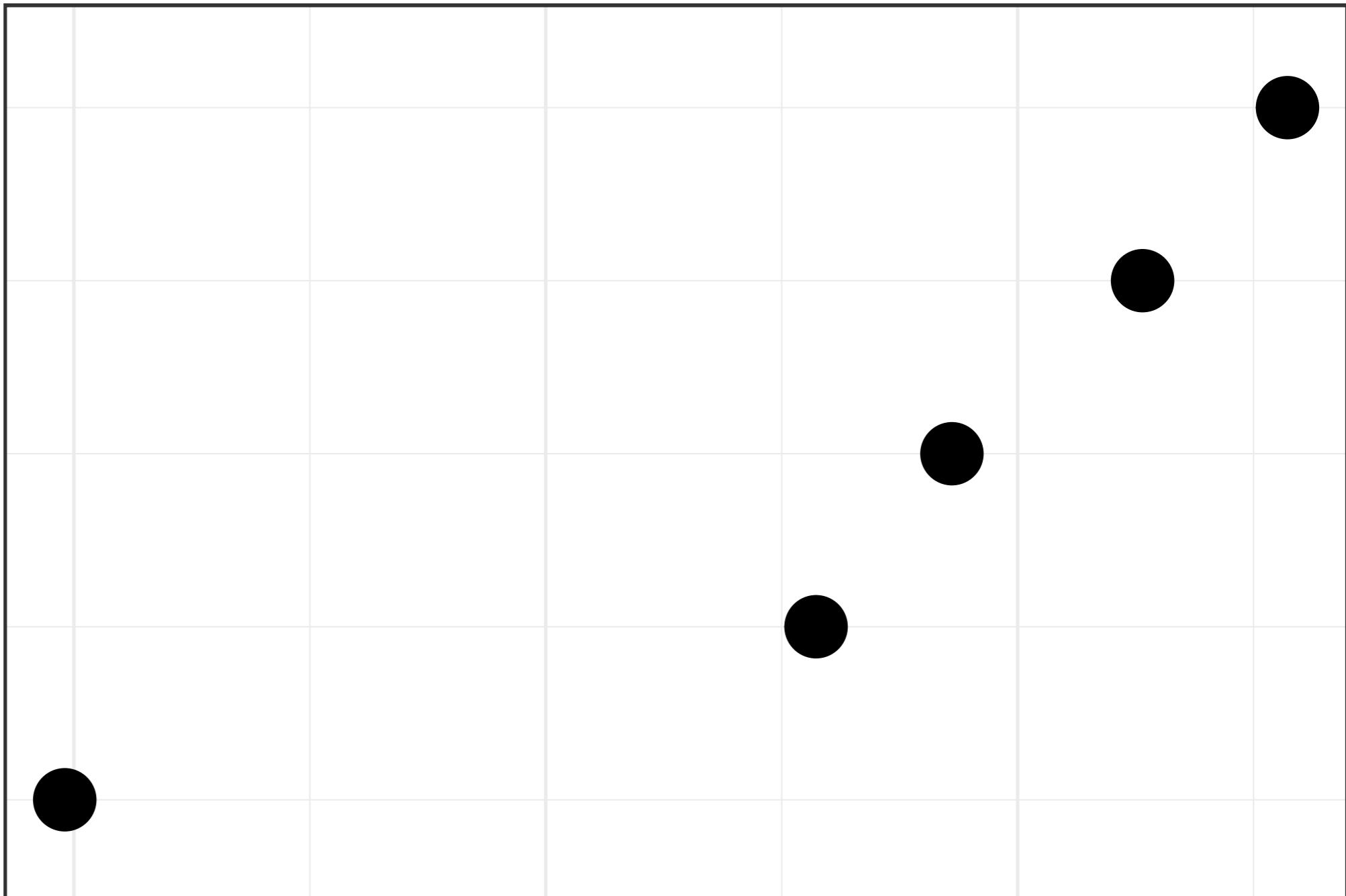
Oceania

Europe

Americas

Asia

Africa



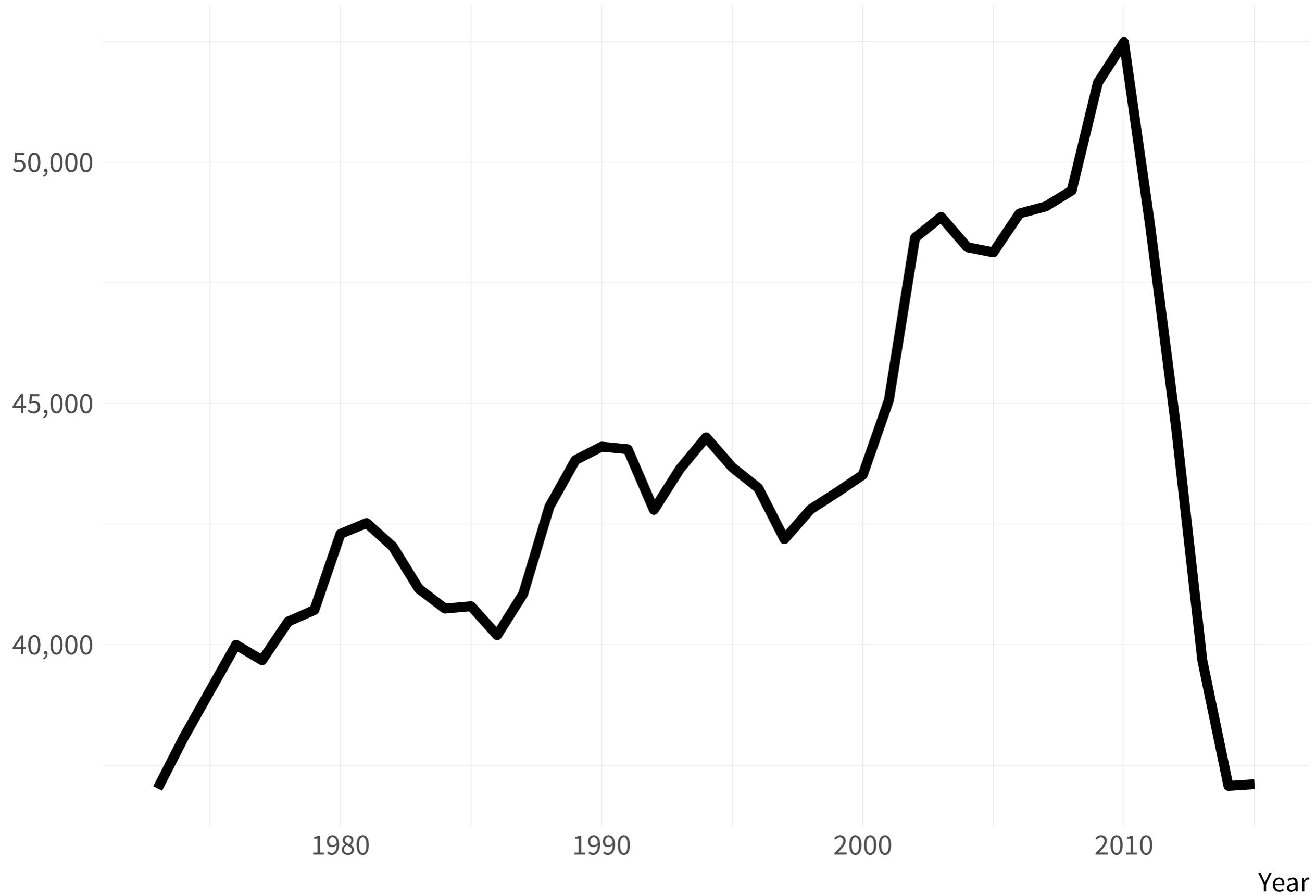
60

70

80

Life Expectancy in years, 2007

First Year Enrollment



First Year Enrollment

40,000

20,000

0

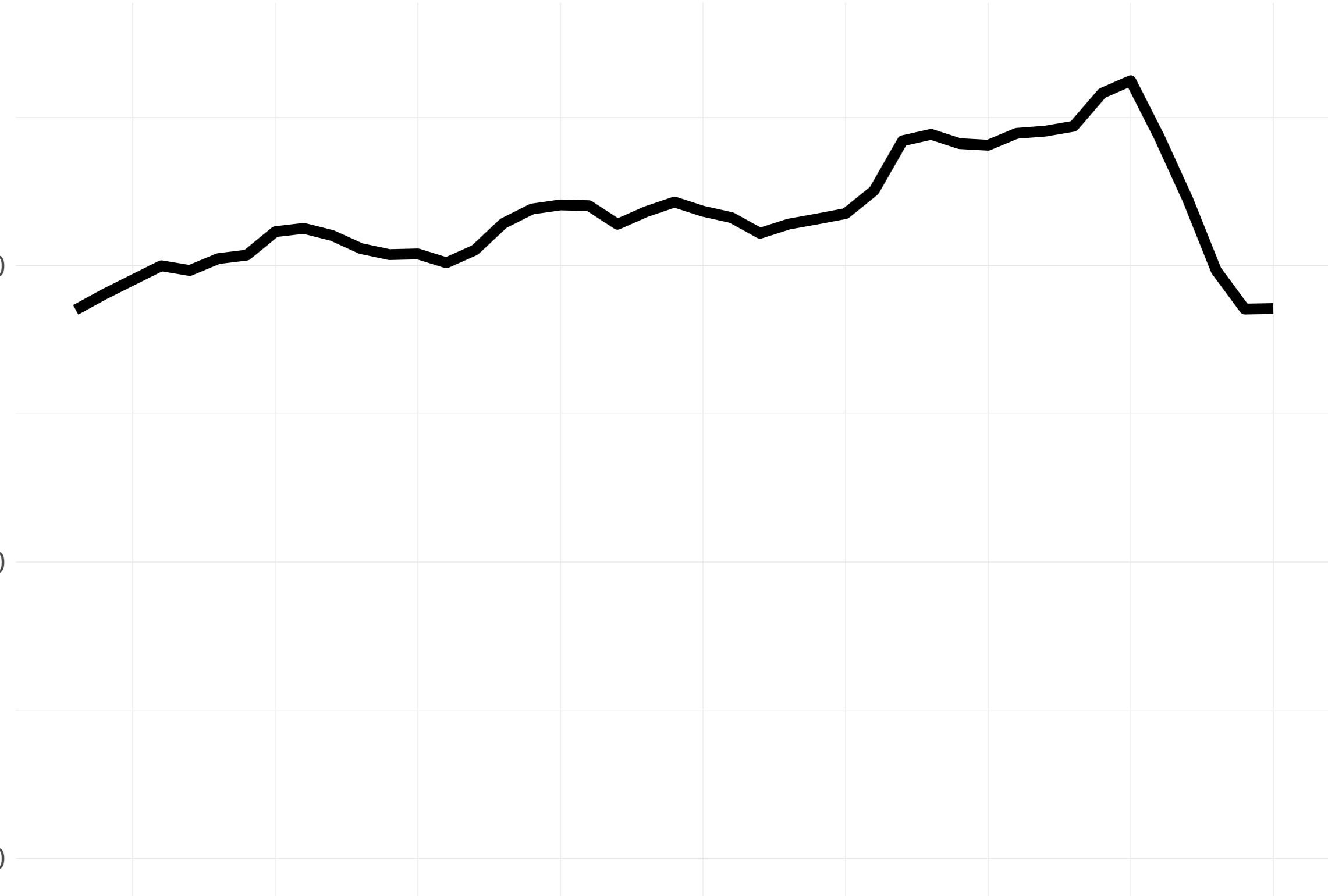
1980

1990

2000

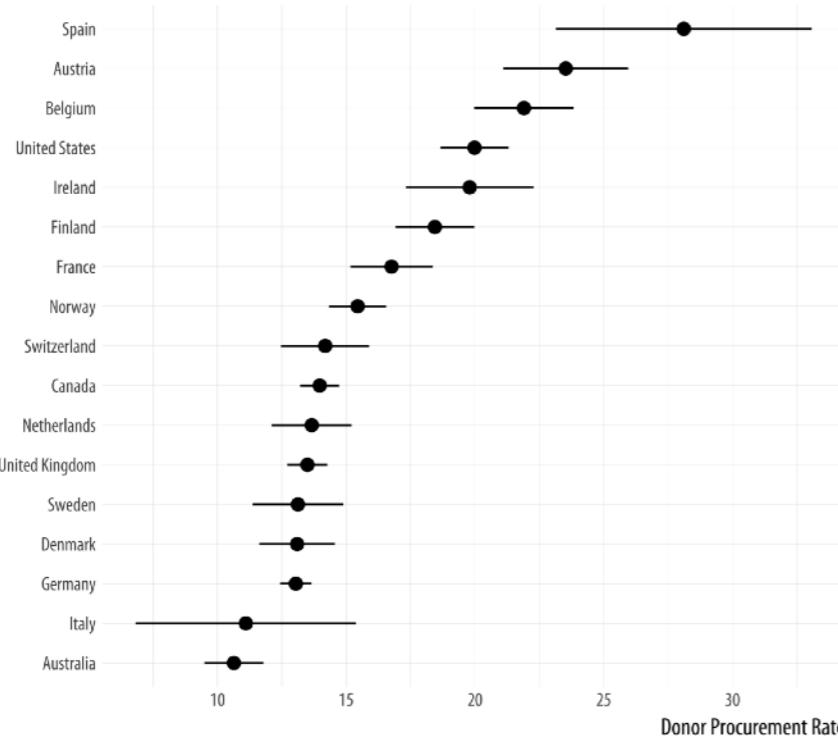
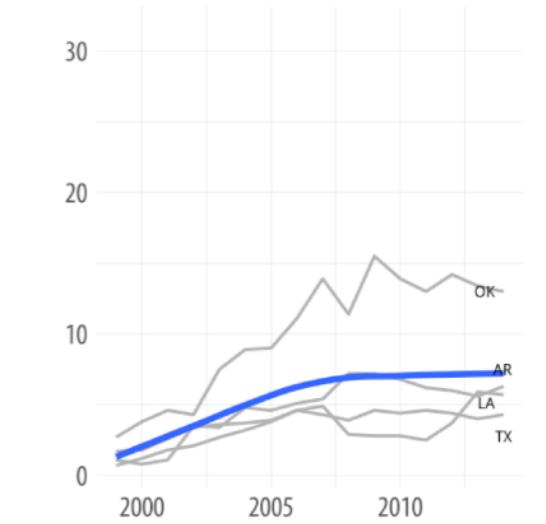
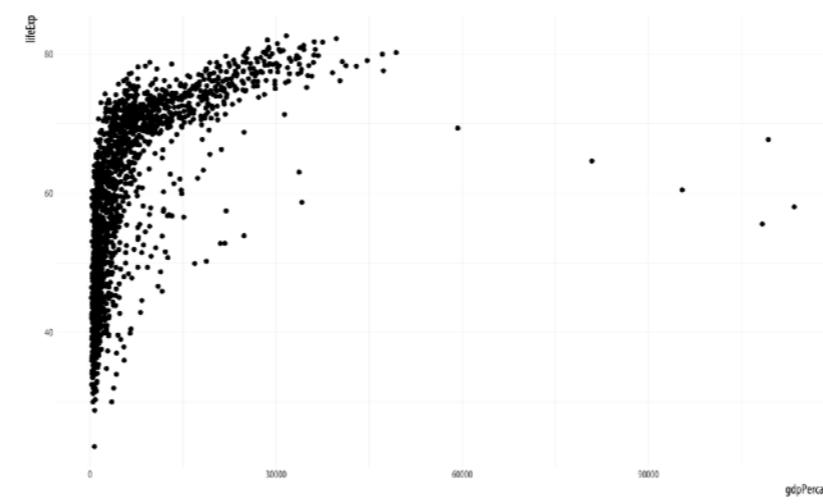
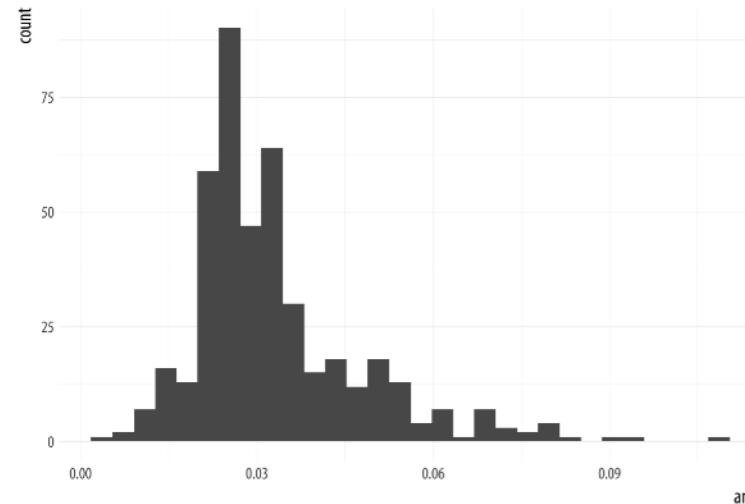
2010

Year

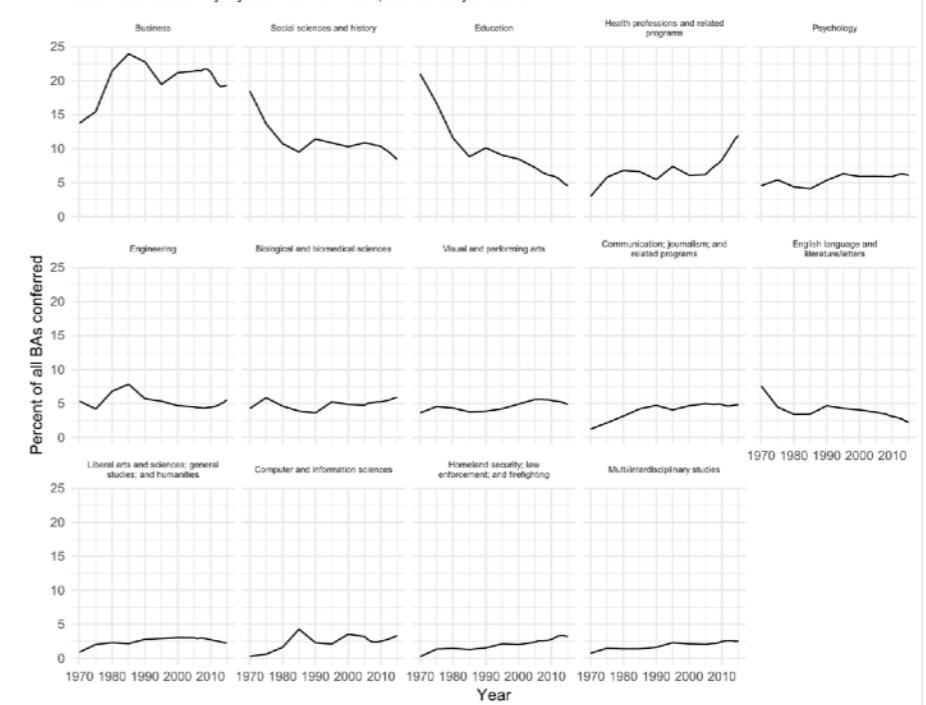


VISUALIZING DATA IN SOCIAL SCIENCE

Workhorses

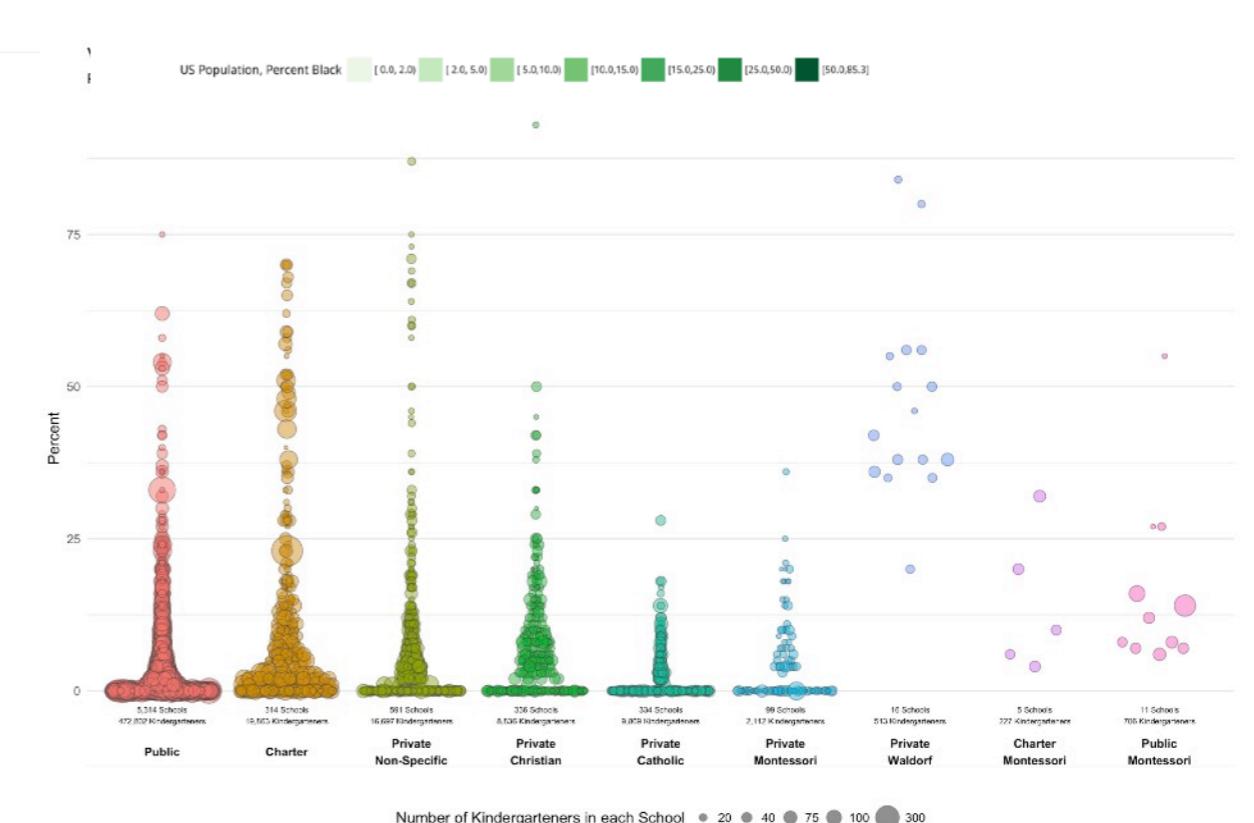
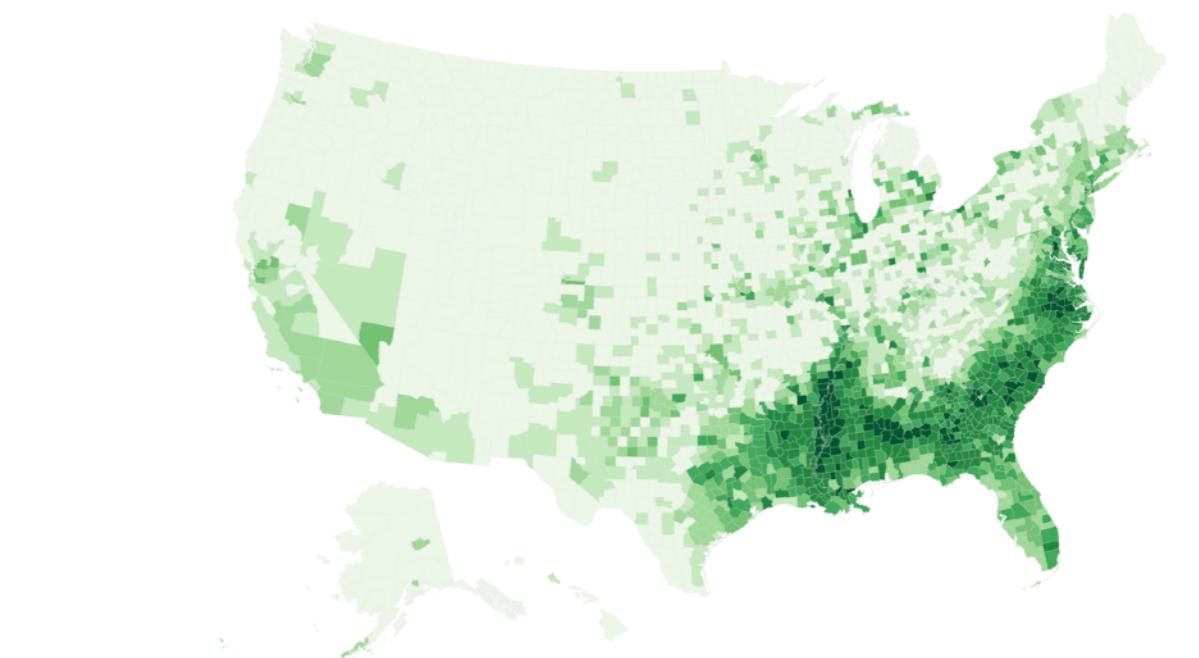
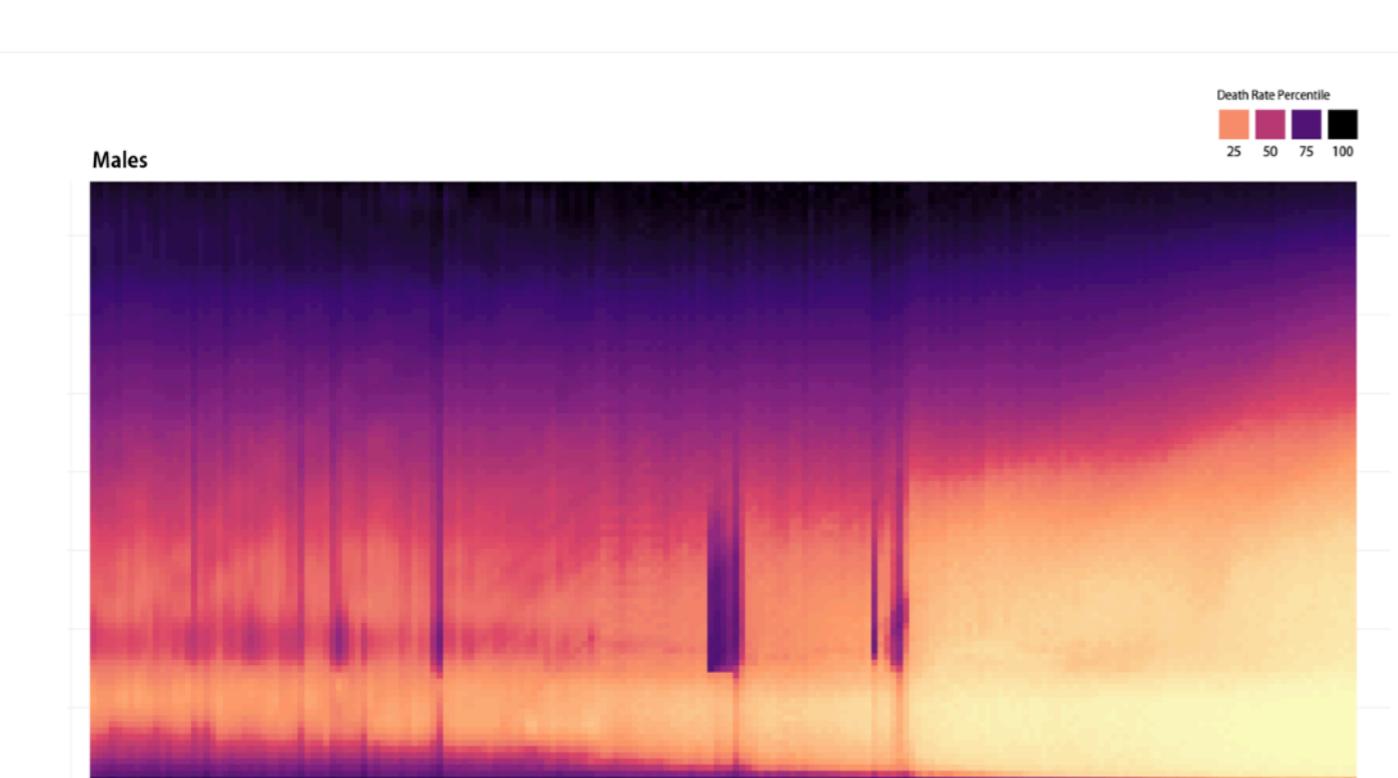
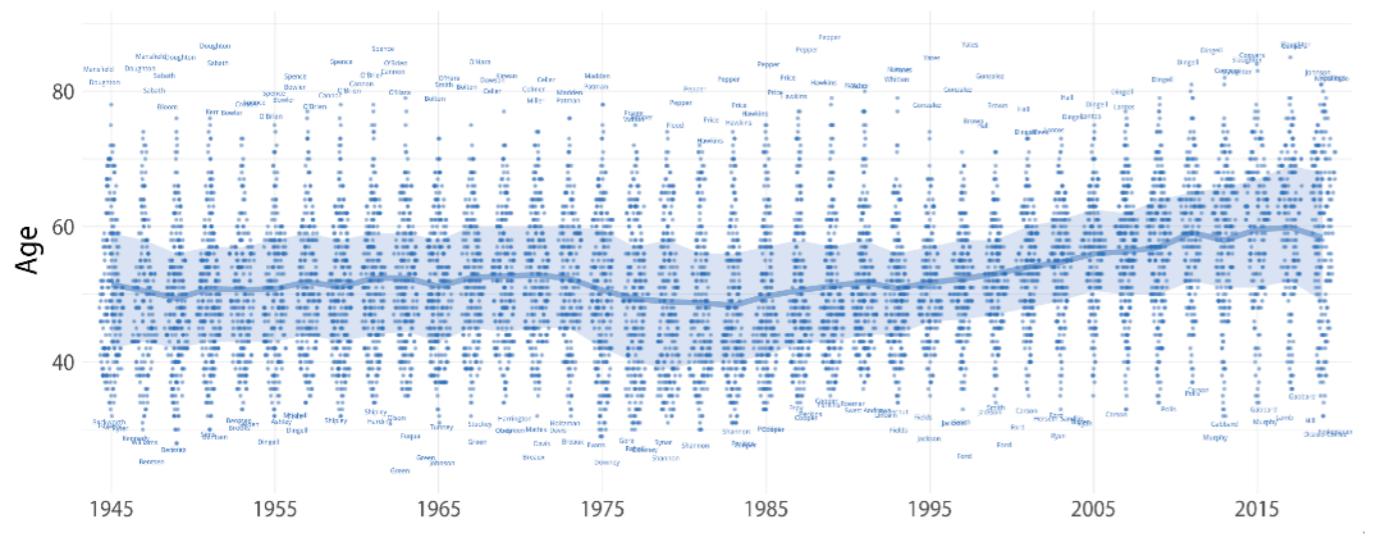


US Trends in Bachelor's Degrees Conferred, 1970-2015,
for Areas averaging more than 2% of all degrees
Observations are every 5 years from 1970-1995, and annually thereafter

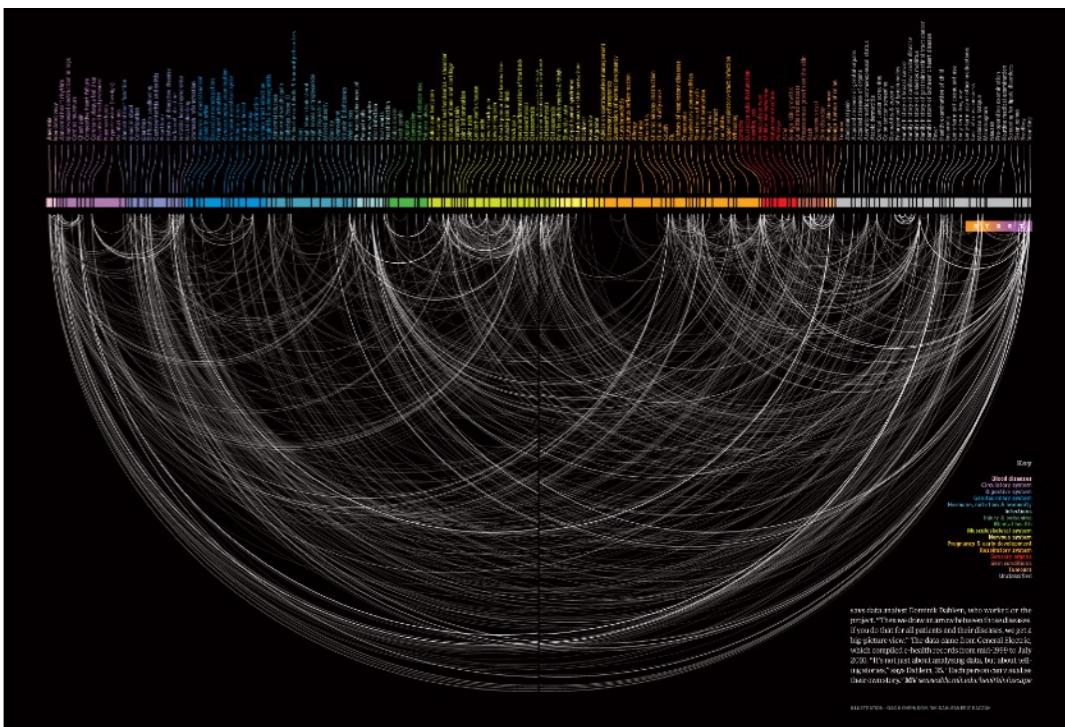
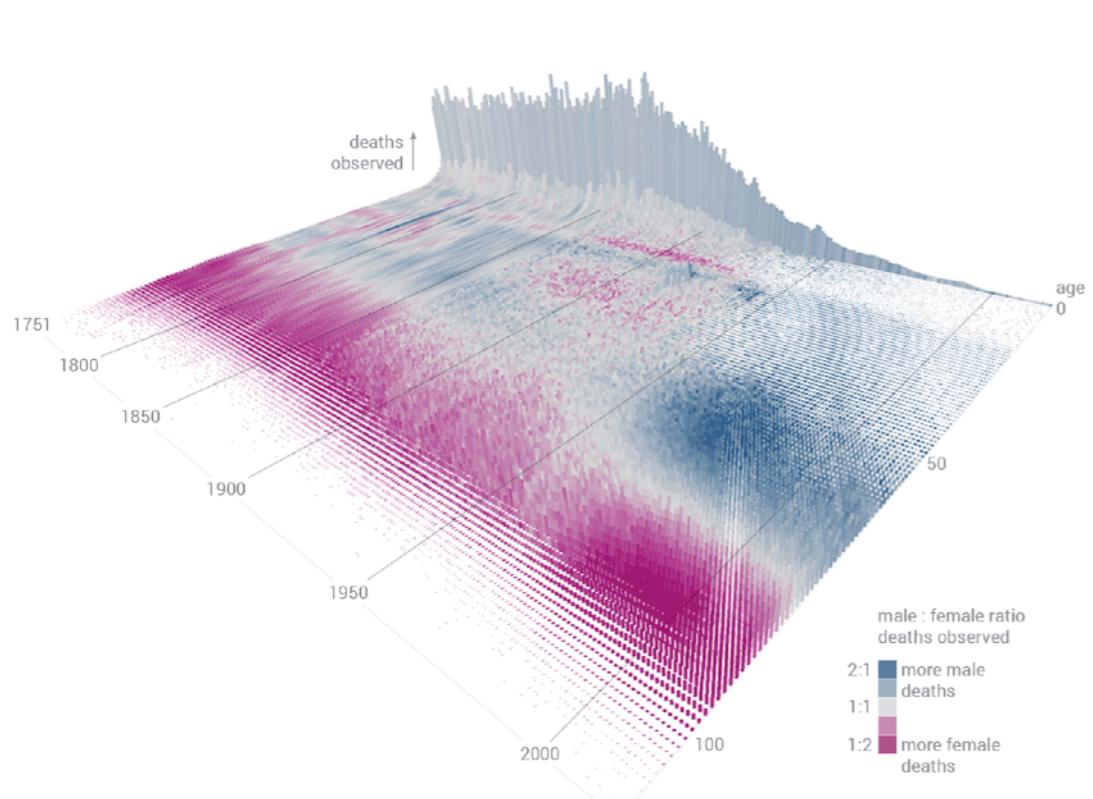
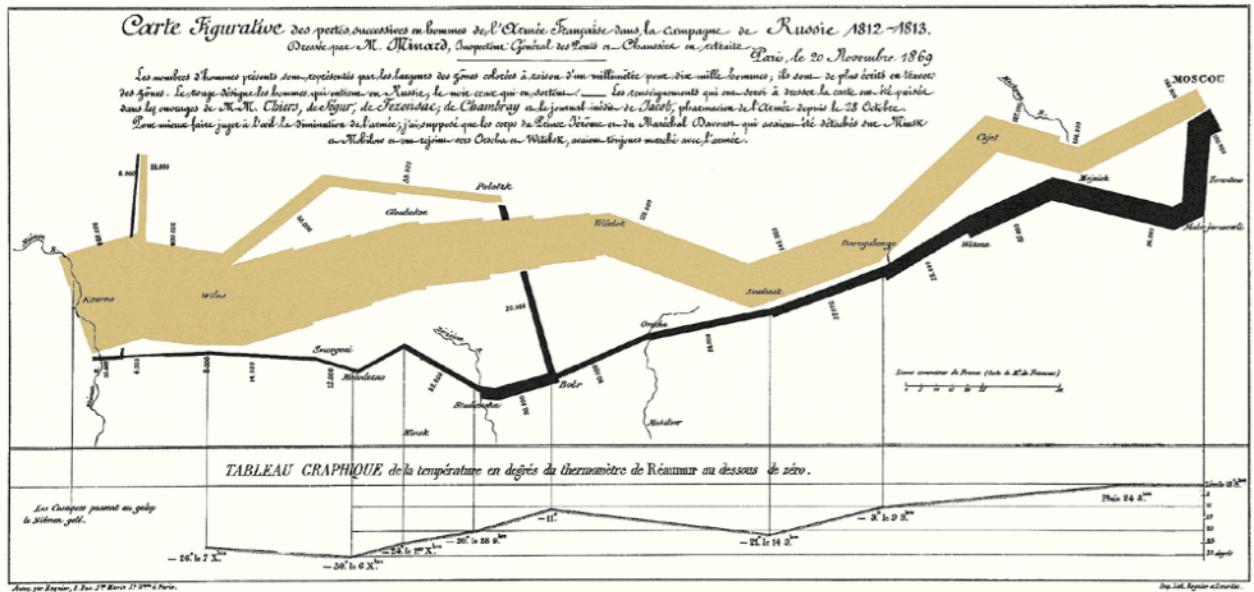


Data from NCES Digest 2017, Table 322.10.

Show Ponies



Unicorns



GETTING STARTED



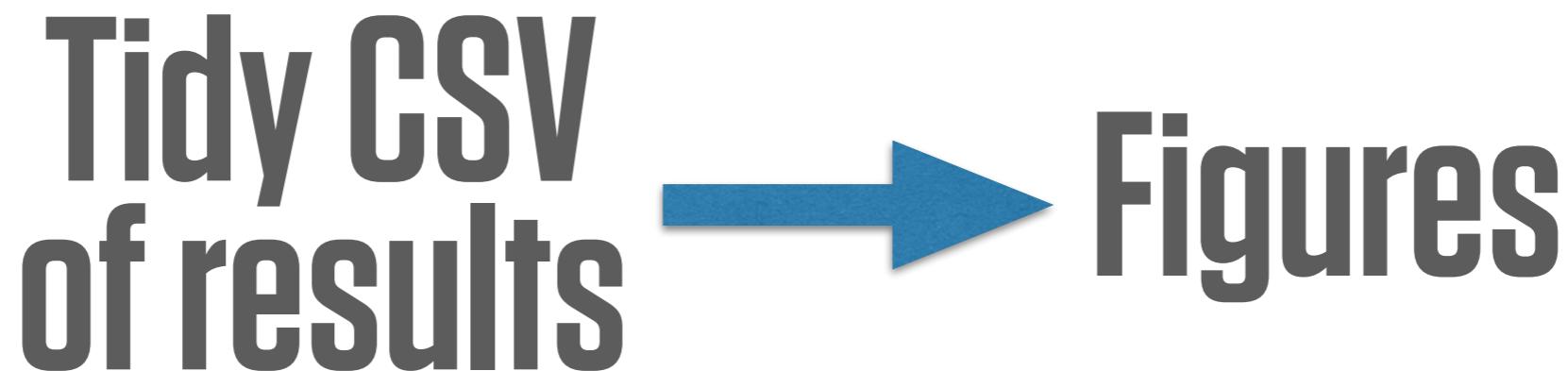
**WE WANT TO
DRAW GOOD
DATA GRAPHICS
REPRODUCIBLY**

Two ways to use R and ggplot

1. Do Everything in R



2. Use ggplot only



Abstraction in Software

Less → More

Easy things are awkward,
Hard things are doable,
Harder things are difficult

Easy things are trivial,
Hard things are difficult,
Harder things are impossible

D3
Grid

ggplot

Stata Excel

```
my_packages <- c("tidyverse", "broom", "coefplot", "dotwhisker",  
  "gapminder", "GGally", "ggrepel", "gridExtra",  
  "interplot", "margins", "maps", "mapproj",  
  "mapdata", "MASS", "quantreg", "scales",  
  "survey", "viridis", "viridisLite", "devtools")
```

```
install.packages(my_packages,  
  repos = "http://cran.rstudio.com")
```

```
devtools::install_github("kjhealy/socviz")
```

```
devtools::install_github("rstudio/radix")
```

REQUIRED LIBRARIES

LET'S GO

```
library(gapminder)
```

gapminder

```
# A tibble: 1,704 x 6
  country continent year lifeExp      pop gdpPercap
  <fctr>    <fctr> <int>   <dbl>    <int>     <dbl>
1 Afghanistan    Asia  1952 28.801 8425333 779.4453
2 Afghanistan    Asia  1957 30.332 9240934 820.8530
3 Afghanistan    Asia  1962 31.997 10267083 853.1007
4 Afghanistan    Asia  1967 34.020 11537966 836.1971
5 Afghanistan    Asia  1972 36.088 13079460 739.9811
6 Afghanistan    Asia  1977 38.438 14880372 786.1134
7 Afghanistan    Asia  1982 39.854 12881816 978.0114
8 Afghanistan    Asia  1987 40.822 13867957 852.3959
9 Afghanistan    Asia  1992 41.674 16317921 649.3414
10 Afghanistan   Asia  1997 41.763 22227415 635.3414
# ... with 1,694 more rows
```

named thing gets ...



... with these
arguments

```
p <- ggplot(data = gapminder,  
               mapping = aes(x = gdpPercap,  
                                 y = lifeExp))
```

... the output of
this function ...



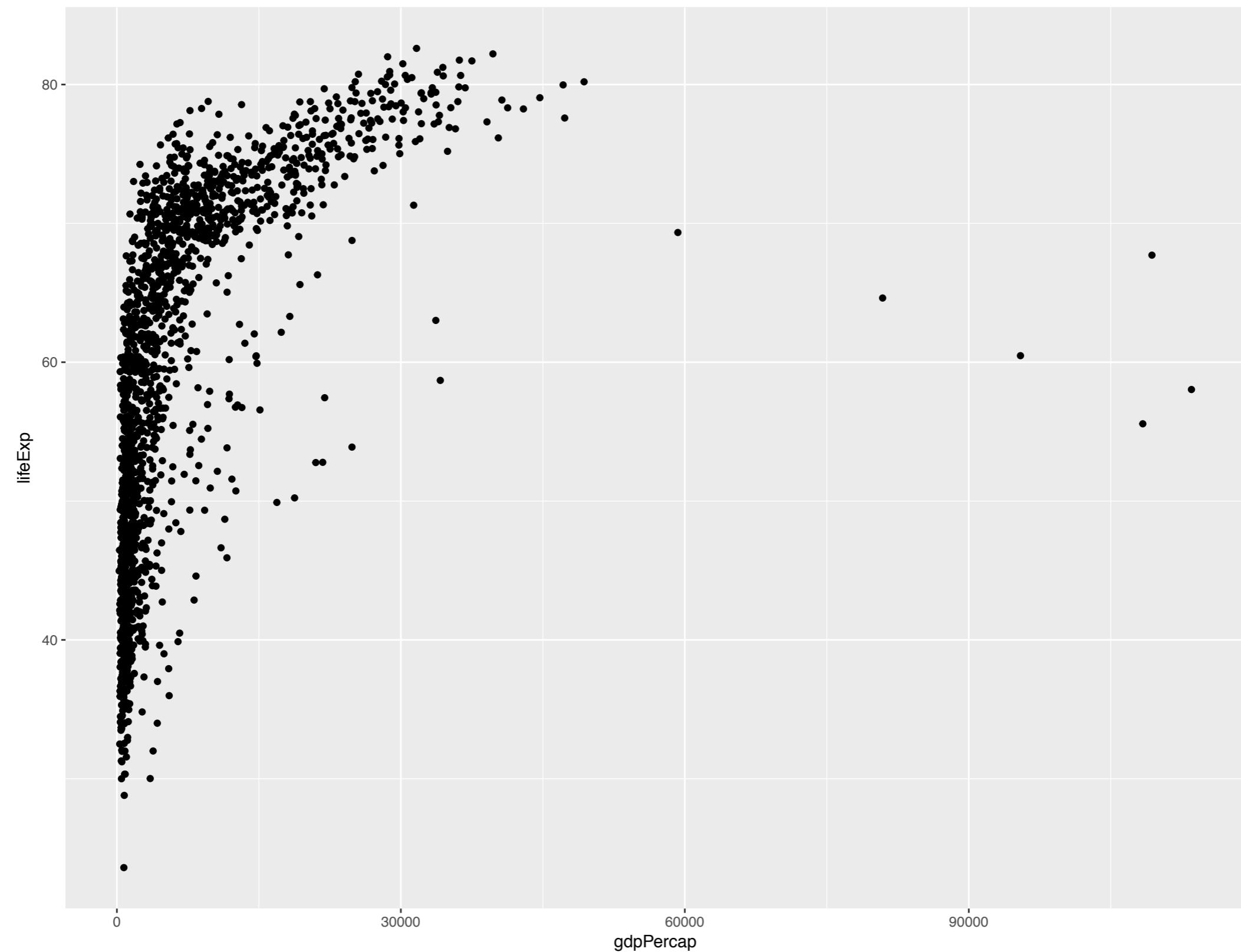
p

```
p + geom_point()
```



Objects created by
`ggplot()` are unusual in
that you can **add**
things to them, and
they will work as
though you wrote all
the code at once.

```
p <- ggplot(data = gapminder,  
               mapping = aes(x = gdpPercap,  
                                 y = lifeExp))  
  
p + geom_point()
```



R will be Frustrating

We're going to be adding a lot of objects together.

```
ggplot(data = mpg,  
       mapping = aes(x = displ, y = hwy)) +  
  geom_point()
```

"+"

goes
here

```
ggplot(data = mpg, mapping = aes(x = displ, y = hwy))  
+ geom_point()
```

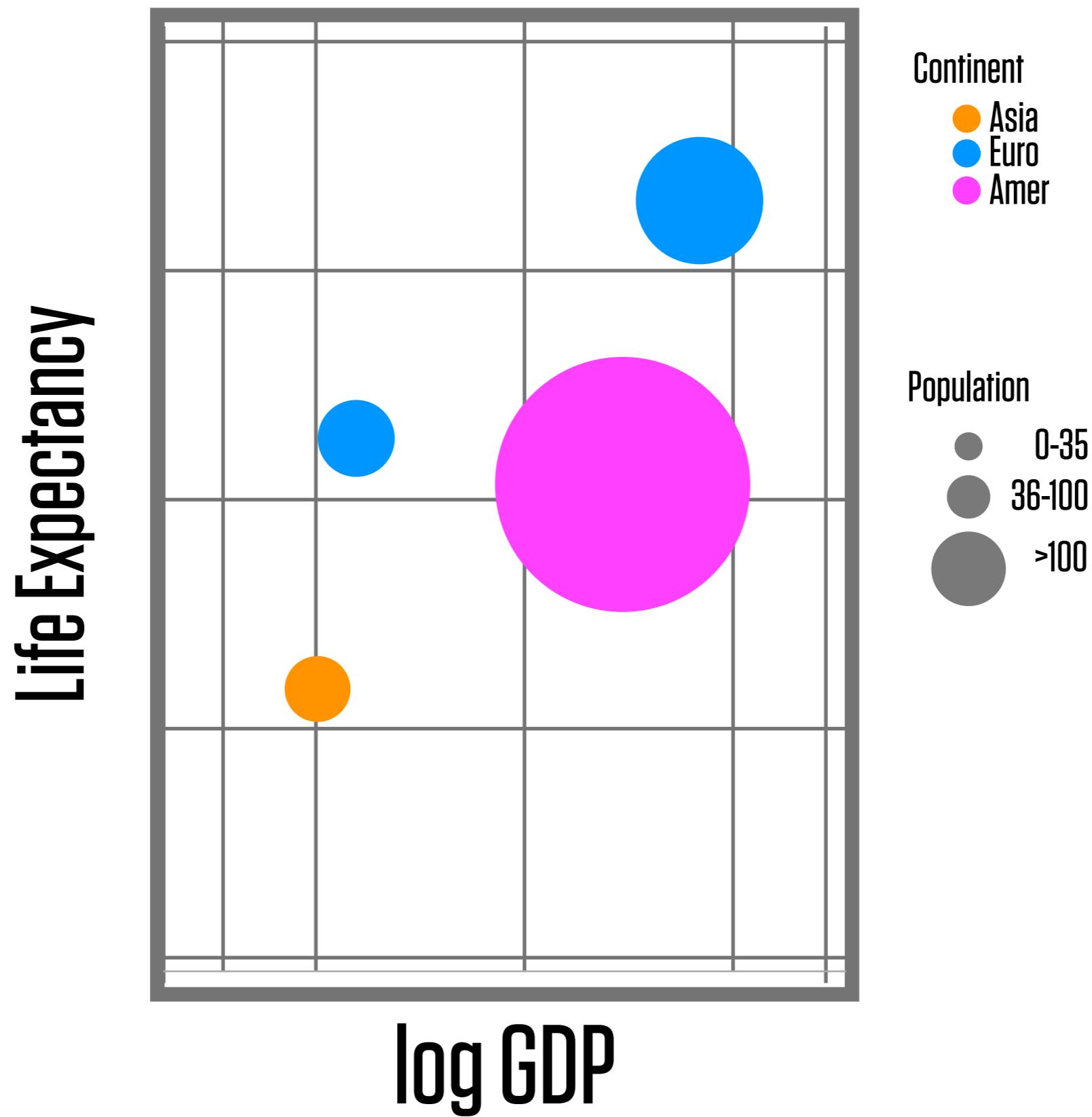


not here

HOW ggplot WORKS

gdp	lifexp	pop	continent
340	65	31	Euro
227	51	200	Amer
909	81	80	Euro
126	40	20	Asia

A Gapminder Plot



1. Tidy Data

gdp	lifexp	pop	continent
340	65	31	Euro
227	51	200	Amer
909	81	80	Euro
126	40	20	Asia

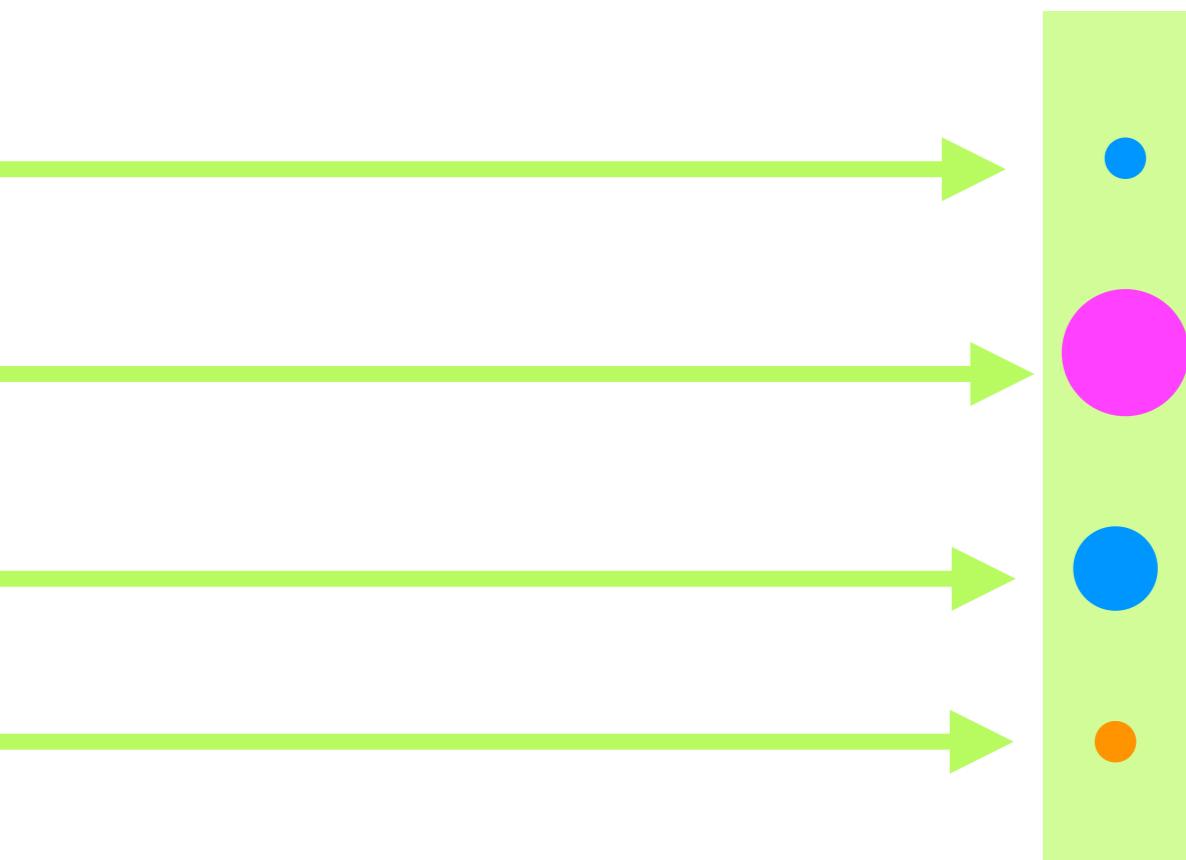
x=gdp y=lifexp size=pop color=continent

2. Mapping

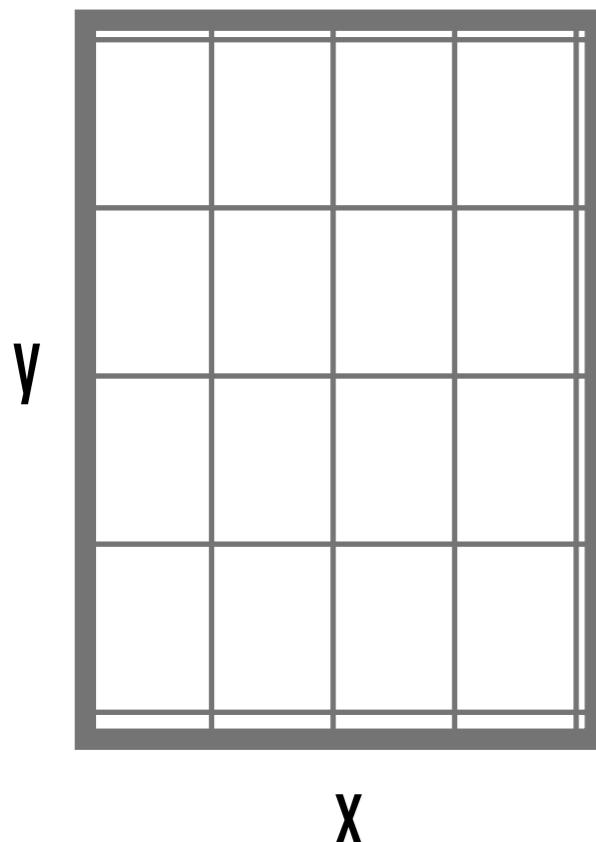
`ggplot(mapping = aes(x = ...))`

3. Geom

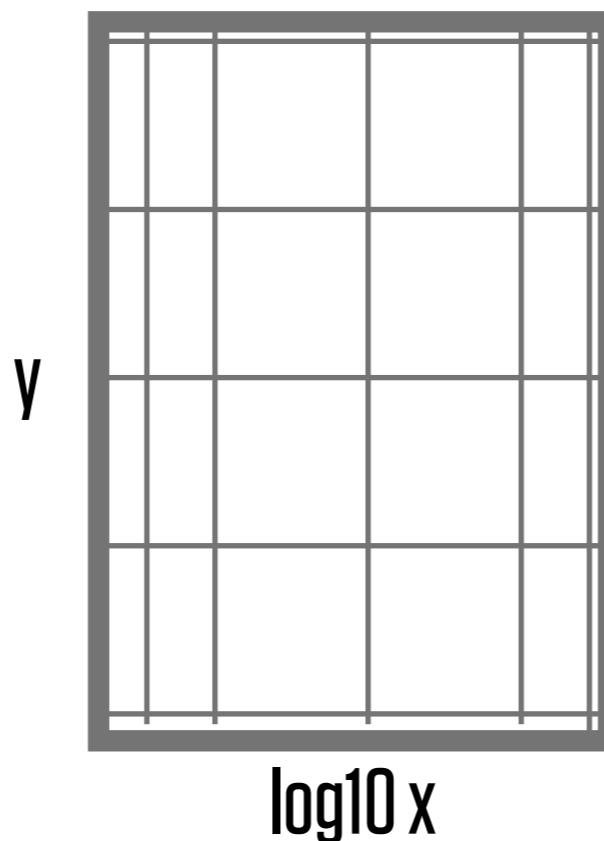
`geom_point()`



4. Coordinate System



5. Scales



6. Labels & Guides

A Gapminder Plot



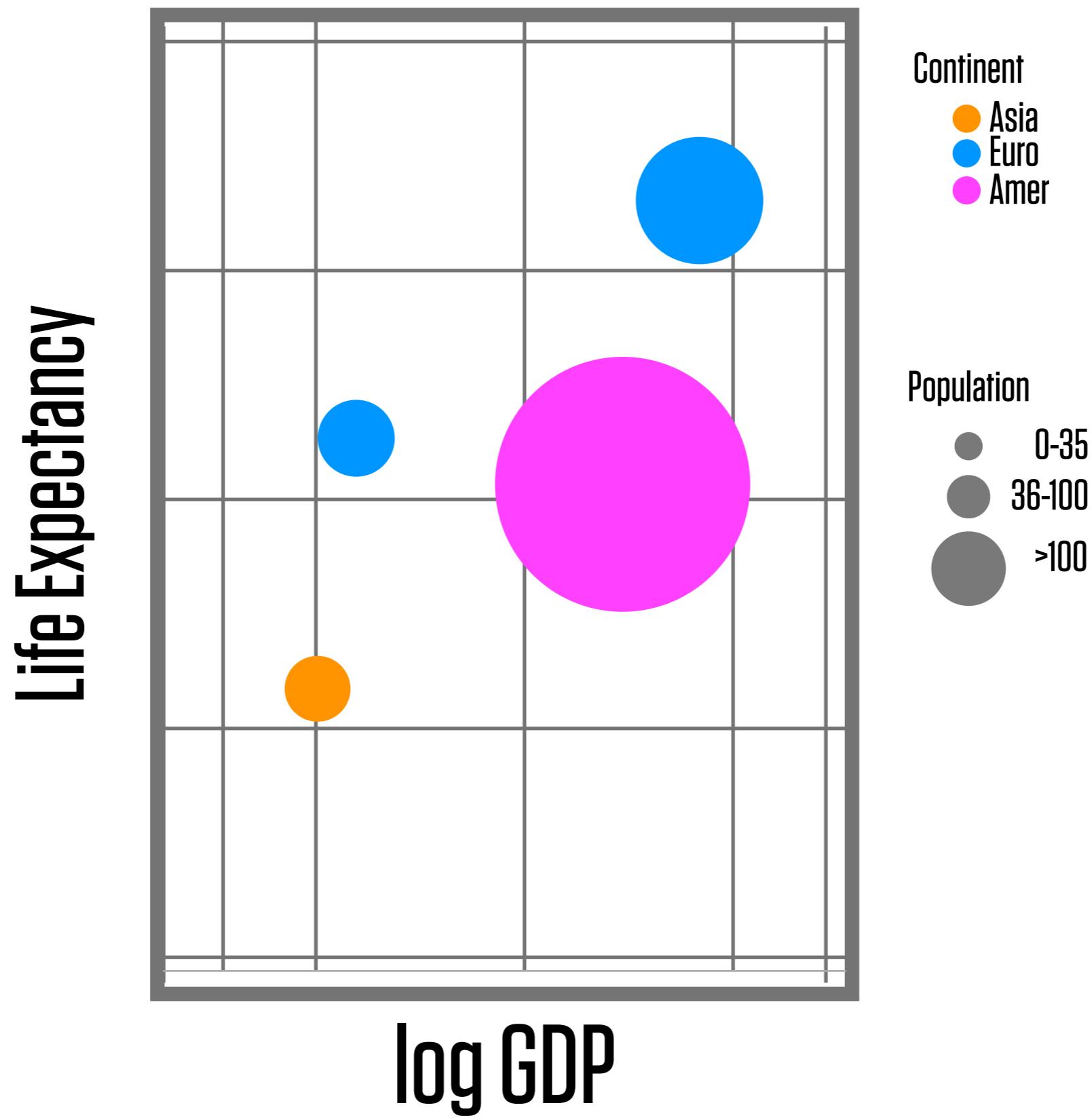
Life Expectancy

Continent

Population

log GDP

A Gapminder Plot



LAYER BY LAYER

```
head(gapminder)
```

```
## # A tibble: 6 × 6
##       country continent   year lifeExp      pop gdpPercap
##       <fctr>    <fctr> <int>   <dbl>     <int>     <dbl>
## 1 Afghanistan      Asia  1952 28.801 8425333 779.4453
## 2 Afghanistan      Asia  1957 30.332 9240934 820.8530
## 3 Afghanistan      Asia  1962 31.997 10267083 853.1007
## 4 Afghanistan      Asia  1967 34.020 11537966 836.1971
## 5 Afghanistan      Asia  1972 36.088 13079460 739.9811
## 6 Afghanistan      Asia  1977 38.438 14880372 786.1134
```

```
dim(gapminder)
```

```
## [1] 1704      6
```

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

Create a ggplot object
Data is gapminder table

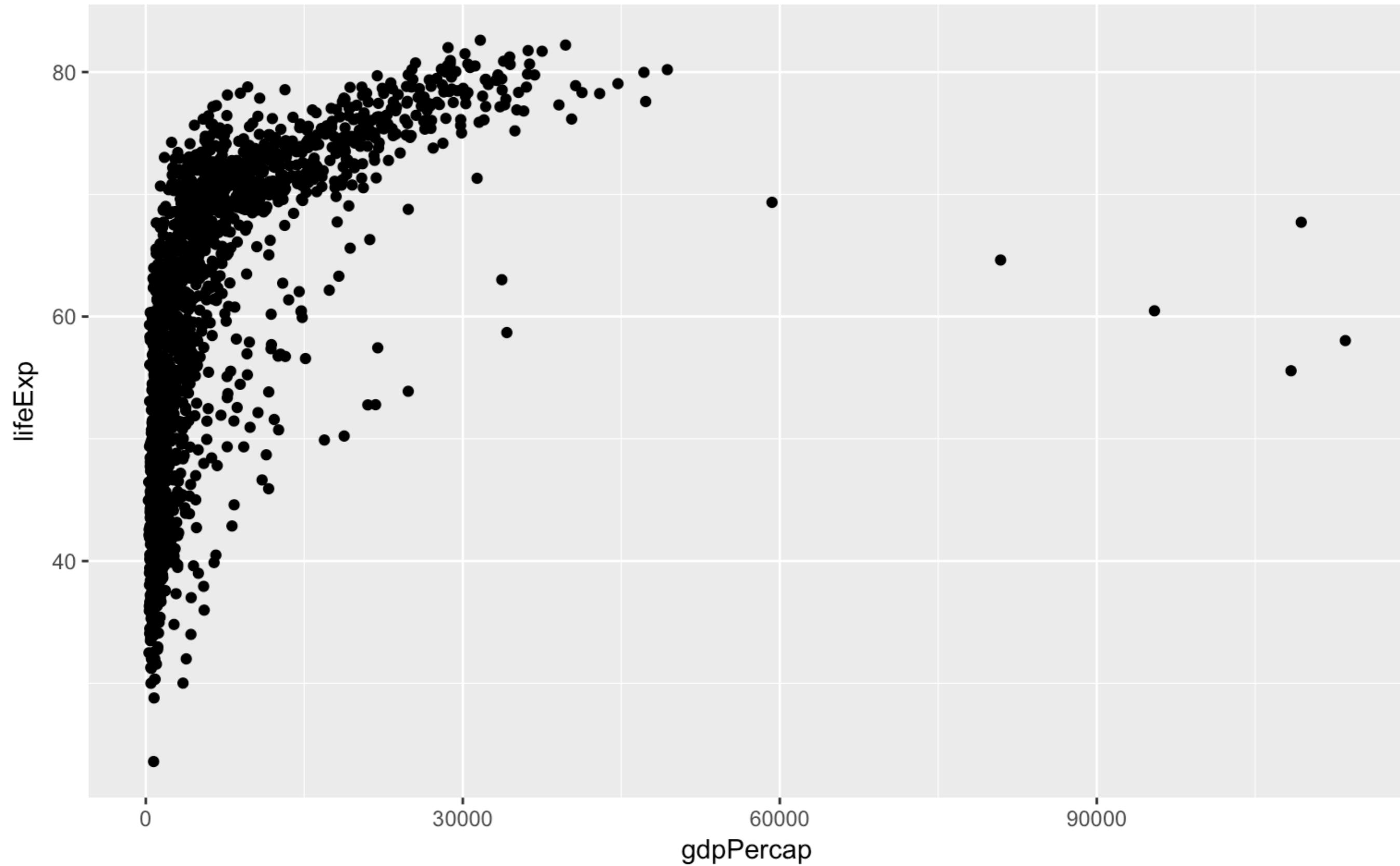
```
p <- ggplot(data = gapminder,  
             mapping = aes(x = gdpPercap,  
                            y = lifeExp))
```

Mapping: tell ggplot the
relationships you want to see

- The `mapping = aes(...)` instruction *links variables* to *things you will see* on the plot.
- The `x` and `y` values are the most obvious ones.
- Other aesthetic mappings can include, e.g., `color`, `shape`, and `size`.
- These mappings are not *directly* specifying what specific, e.g., colors or shapes will be on the plot. Rather they say which *variables* in the data will be *represented* by, e.g., colors and shapes on the plot.

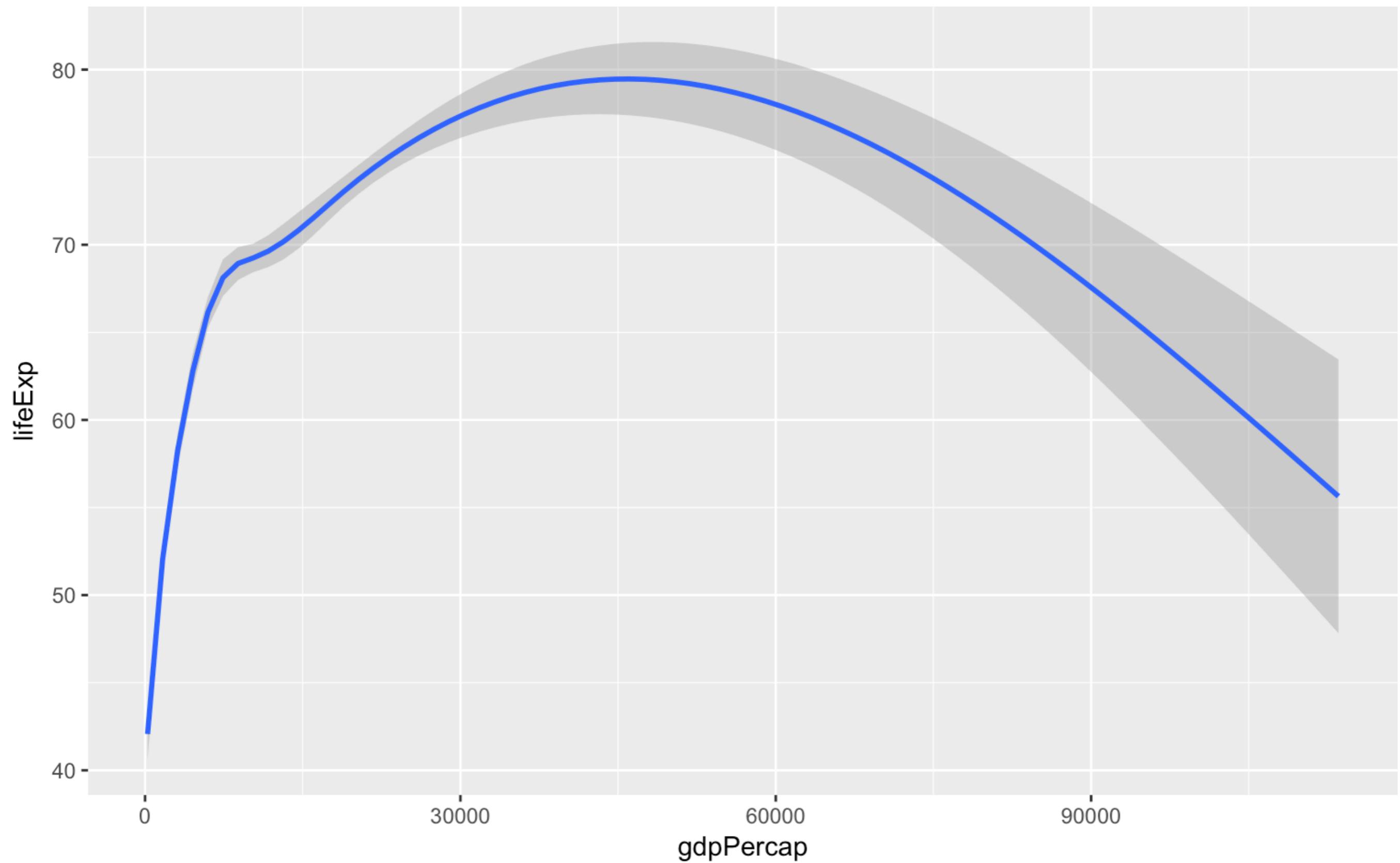
```
p + geom_point()
```

Add a geom layer
to the plot



```
p + geom_smooth()
```

Try a different geom

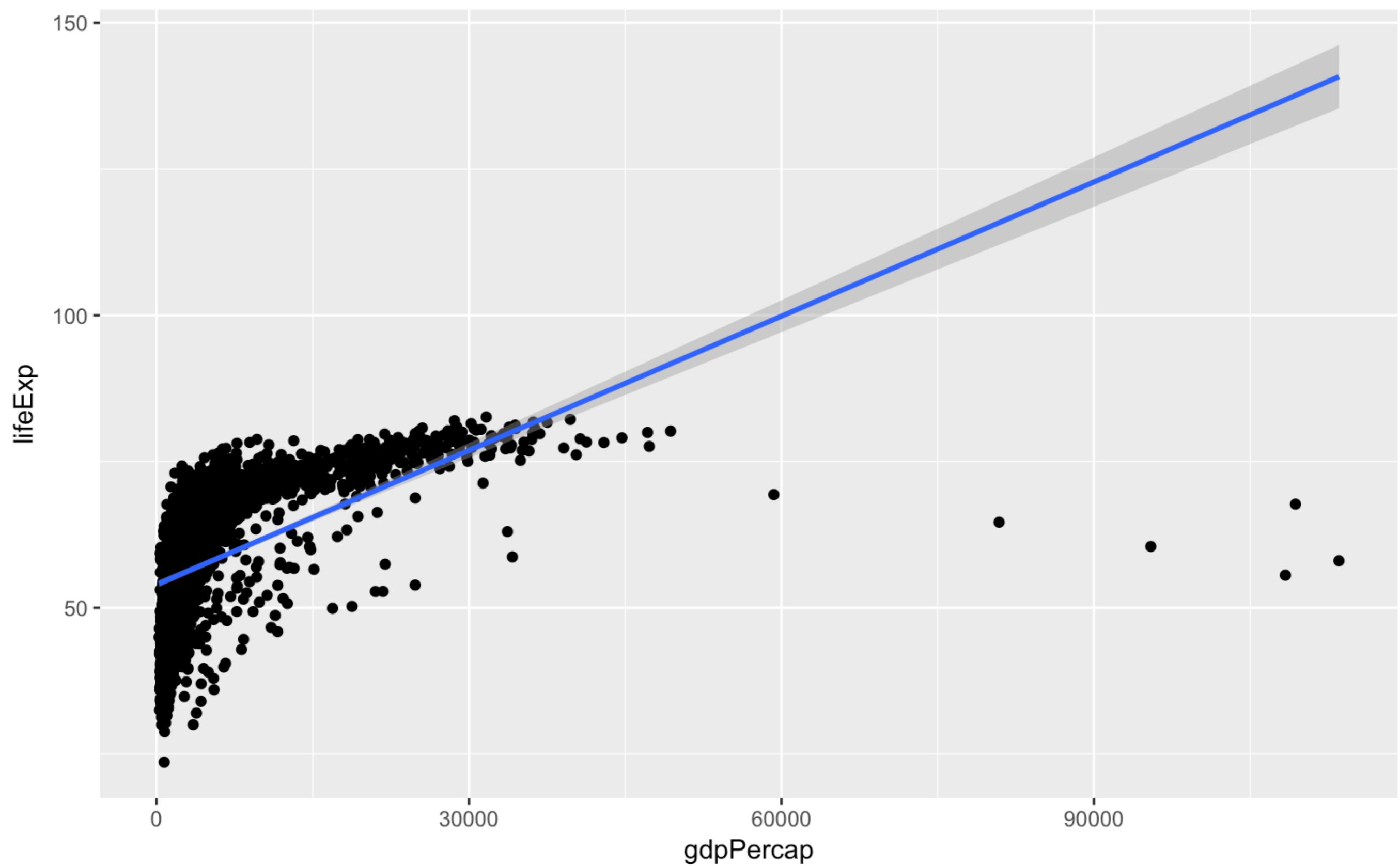


```
p + geom_point() +  
  geom_smooth(method = "gam") +  
  scale_x_log10(labels = scales::dollar)
```

This process is
literally additive

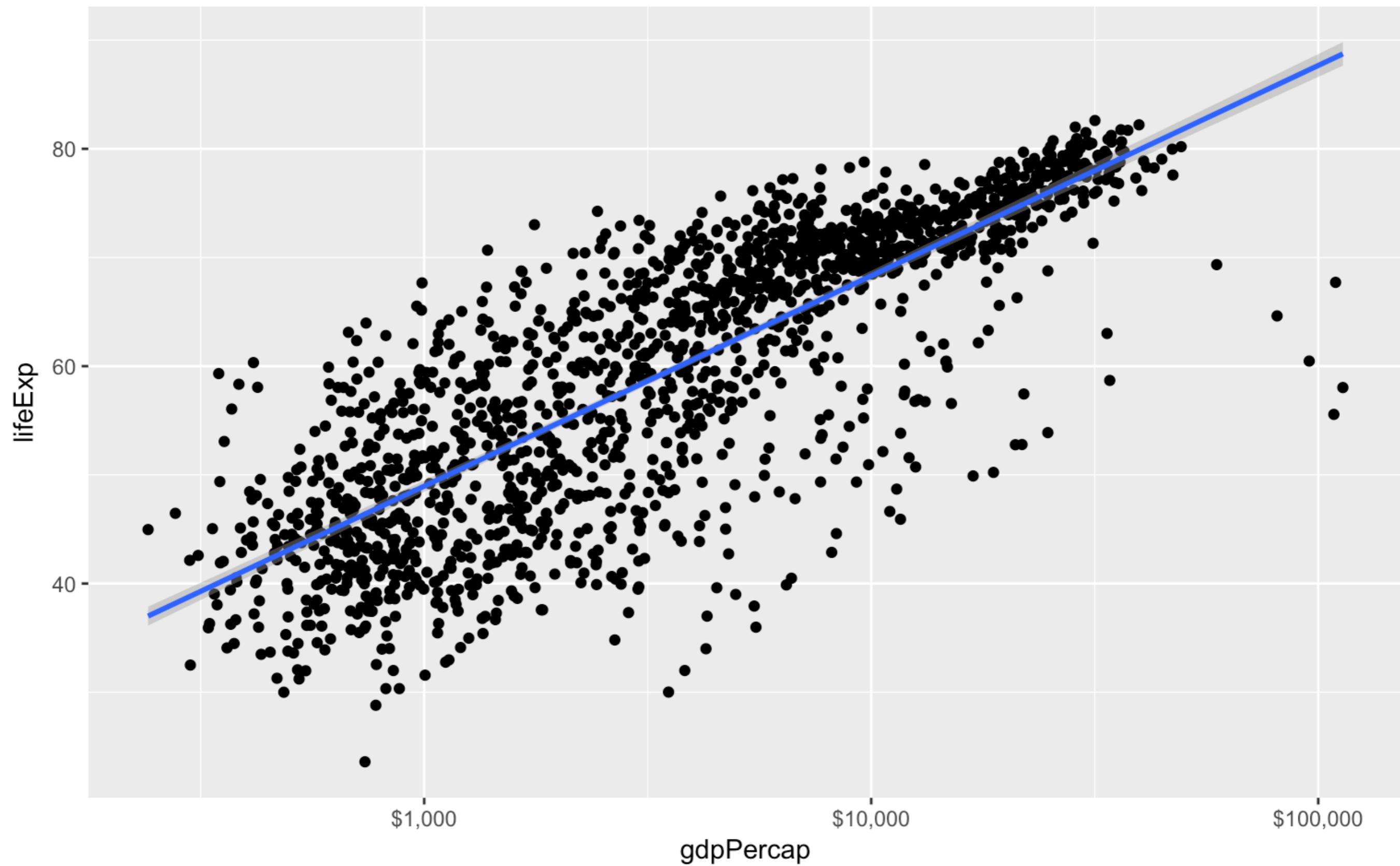
```
p + geom_point() +  
geom_smooth(method = "lm")
```

Every geom is a function.
It can take arguments.



```
p <- ggplot(data = gapminder,  
               mapping = aes(x = gdpPercap,  
                                 y = lifeExp))  
p + geom_point() +  
  geom_smooth(method = "lm") +  
  scale_x_log10(label = scales::dollar)
```

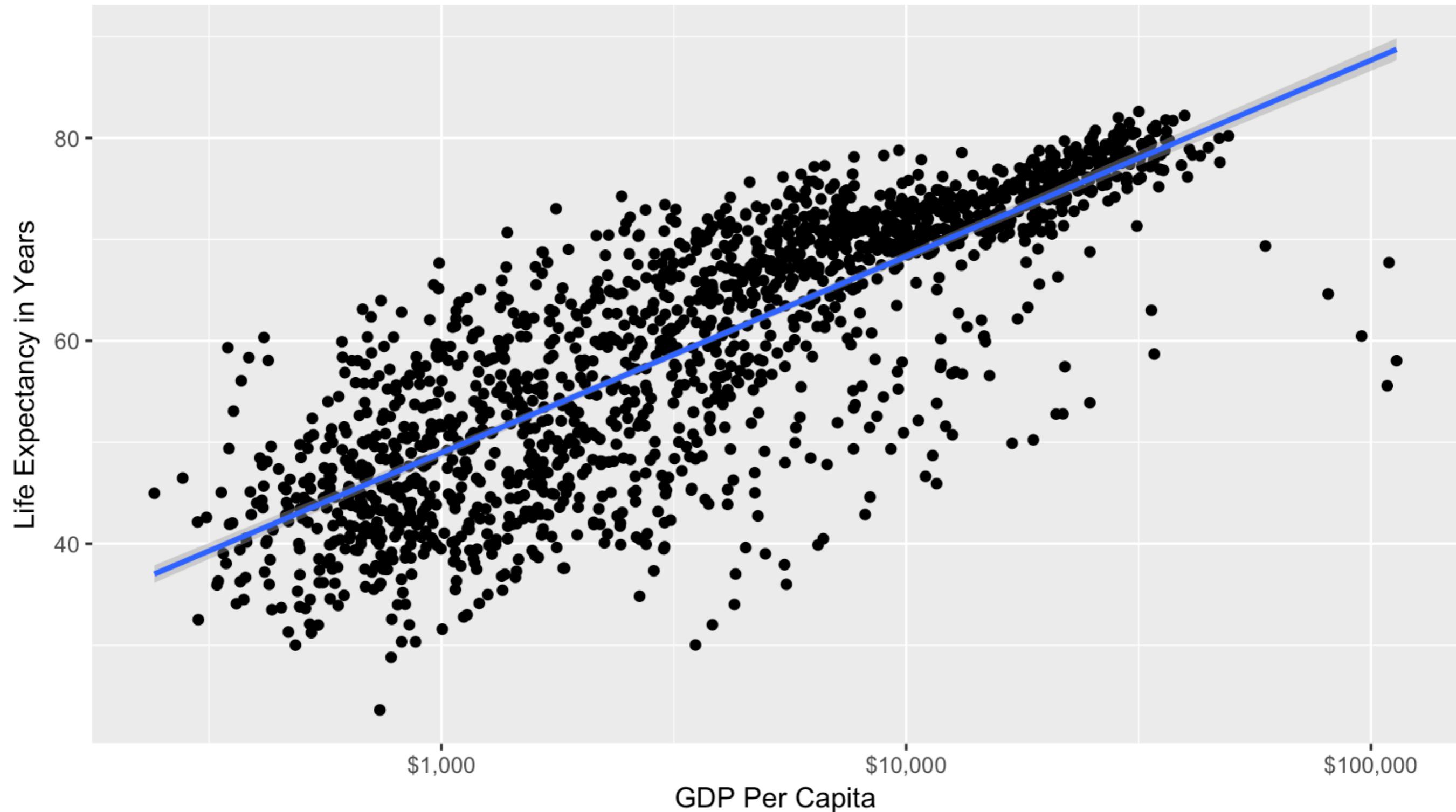
Keep Layering



```
p + geom_point() +
  geom_smooth(method = "gam") +
  scale_x_log10(labels = scales::dollar) +
  labs(x = "GDP Per Capita",
       y = "Life Expectancy in Years",
       title = "Economic Growth and Life Expectancy",
       subtitle = "Data points are country-years",
       caption = "Data source: Gapminder")
```

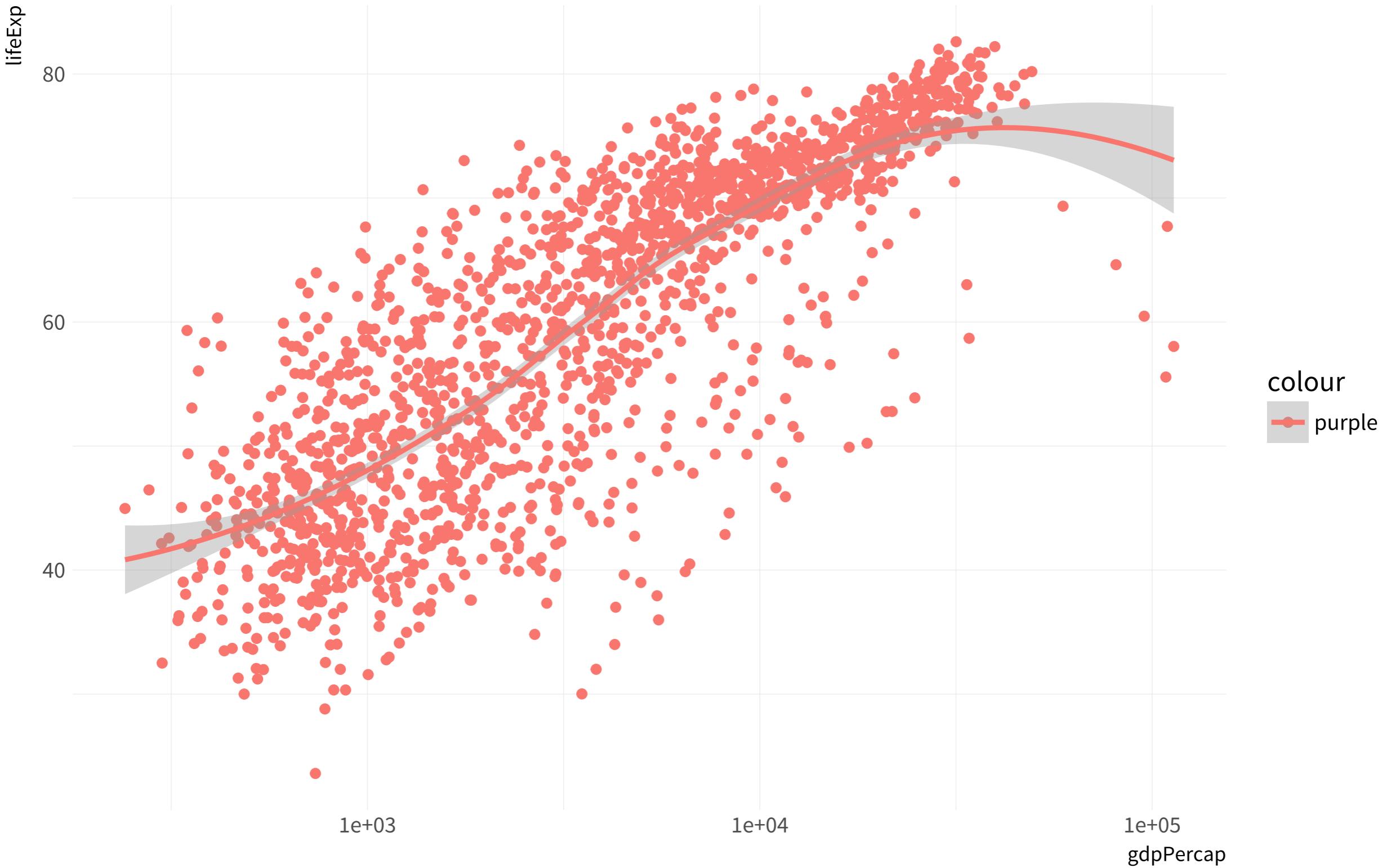
Economic Growth and Life Expectancy

Data points are country-years



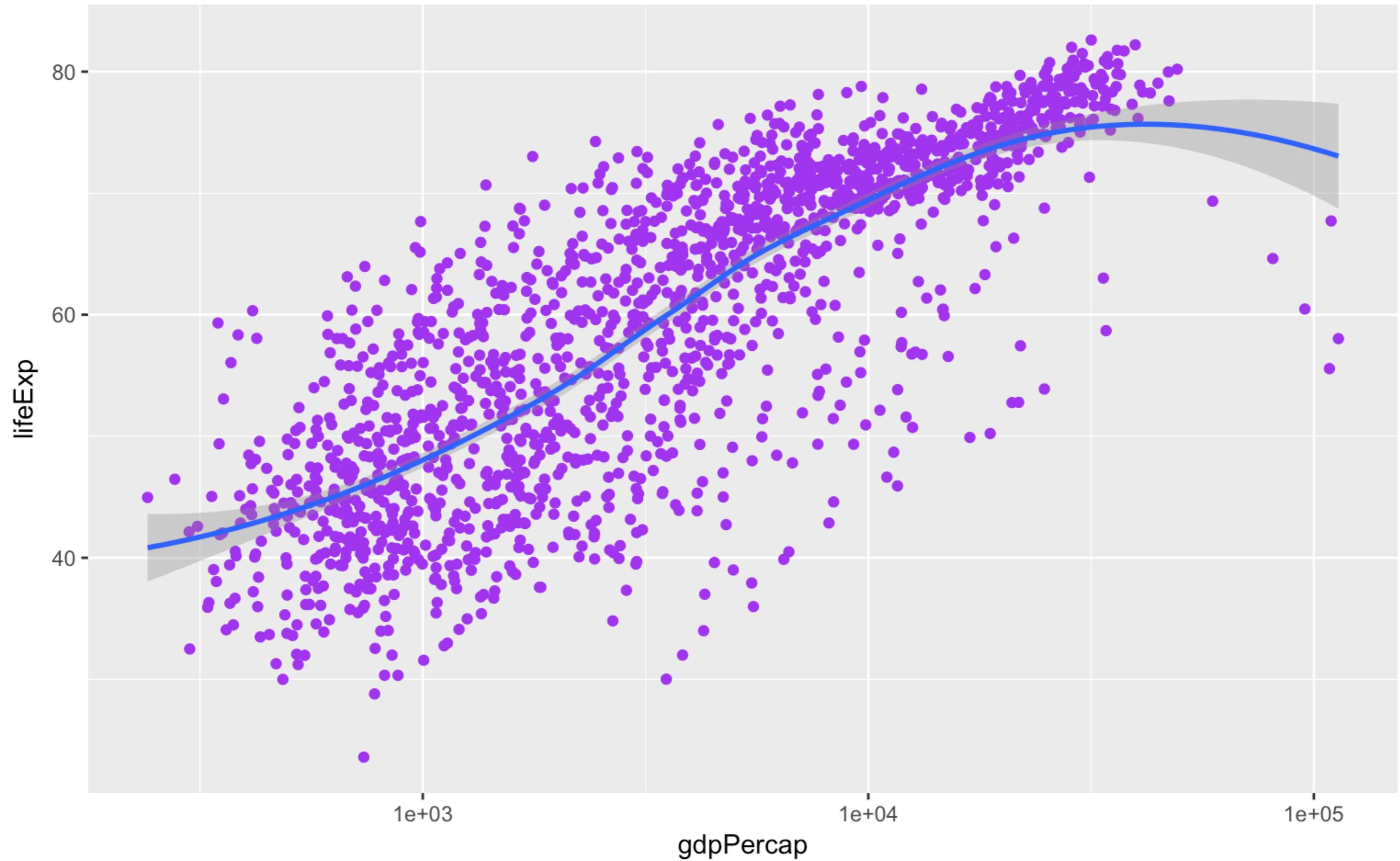
MAPPING
vs SETTING
AESTHETICS

```
p <- ggplot(data = gapminder,  
               mapping = aes(x = gdpPercap,  
                                 y = lifeExp,  
                                 color = "purple"))  
p + geom_point() +  
    geom_smooth(method = "loess") +  
    scale_x_log10()
```



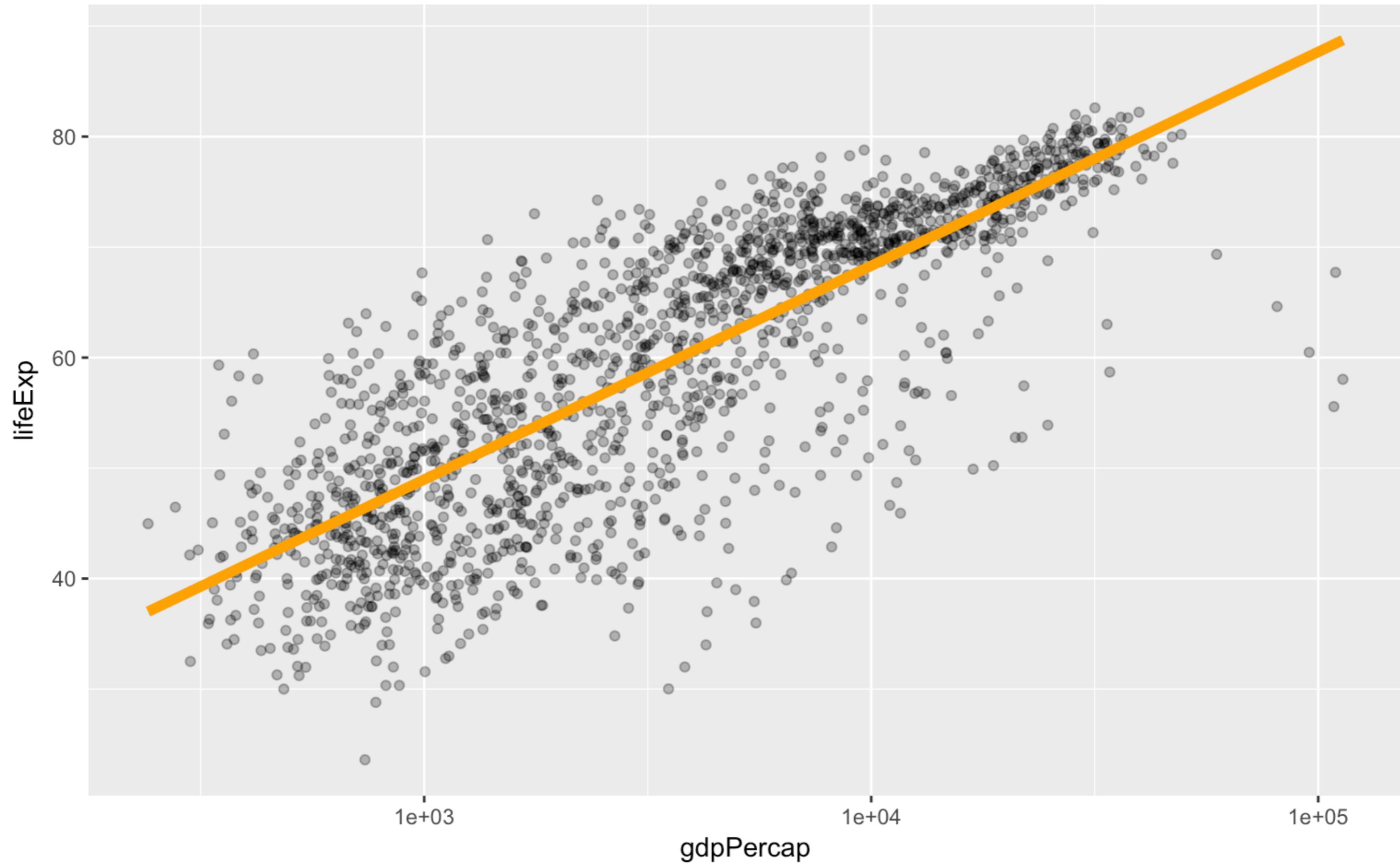
What has gone wrong here?

```
p <- ggplot(data = gapminder,  
             mapping = aes(x = gdpPercap,  
                            y = lifeExp))  
p + geom_point(color = "purple") +  
  geom_smooth(method = "loess") +  
  scale_x_log10()
```

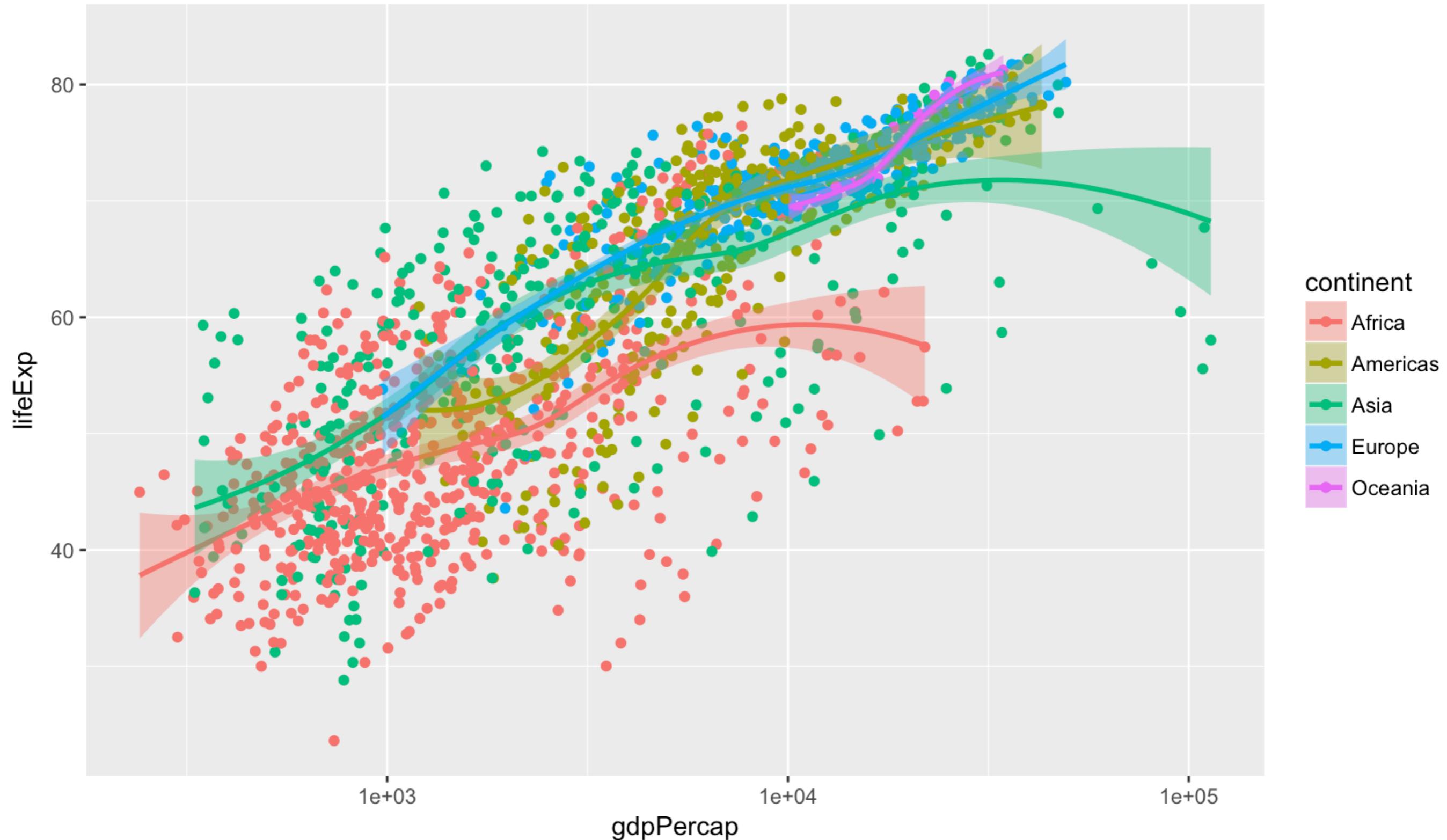


```
p <- ggplot(data = gapminder,  
             mapping = aes(x = gdpPercap,  
                            y = lifeExp))  
p + geom_point(alpha = 0.3) +  
  geom_smooth(color = "orange",  
              se=FALSE, size=2, method = "lm") +  
  scale_x_log10()
```

Here, some aesthetics are
mapped, and some are **set**



```
p <- ggplot(data = gapminder,  
             mapping = aes(x = gdpPercap,  
                            y = lifeExp,  
                            color = continent,  
                            fill = continent))  
p + geom_point() +  
  geom_smooth(method = "loess") +  
  scale_x_log10()
```



**MAP or SET
AESTHETICS
per geom**

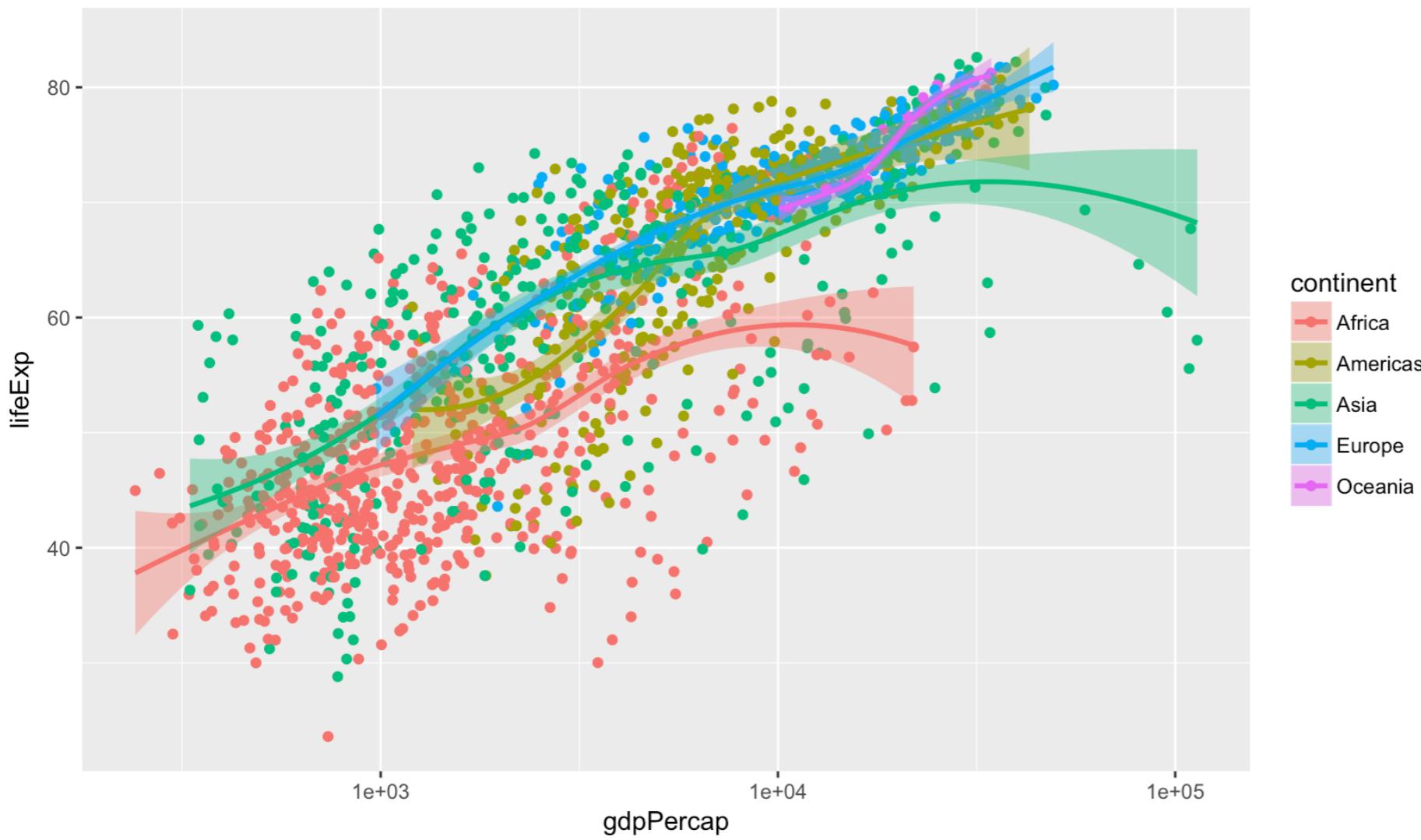
```
p <- ggplot(data = gapminder,  
               mapping = aes(x = gdpPercap,  
                                 y = lifeExp))  
p + geom_point(mapping =  
                  aes(color = continent)) +  
geom_smooth(method = "loess") +  
scale_x_log10()
```



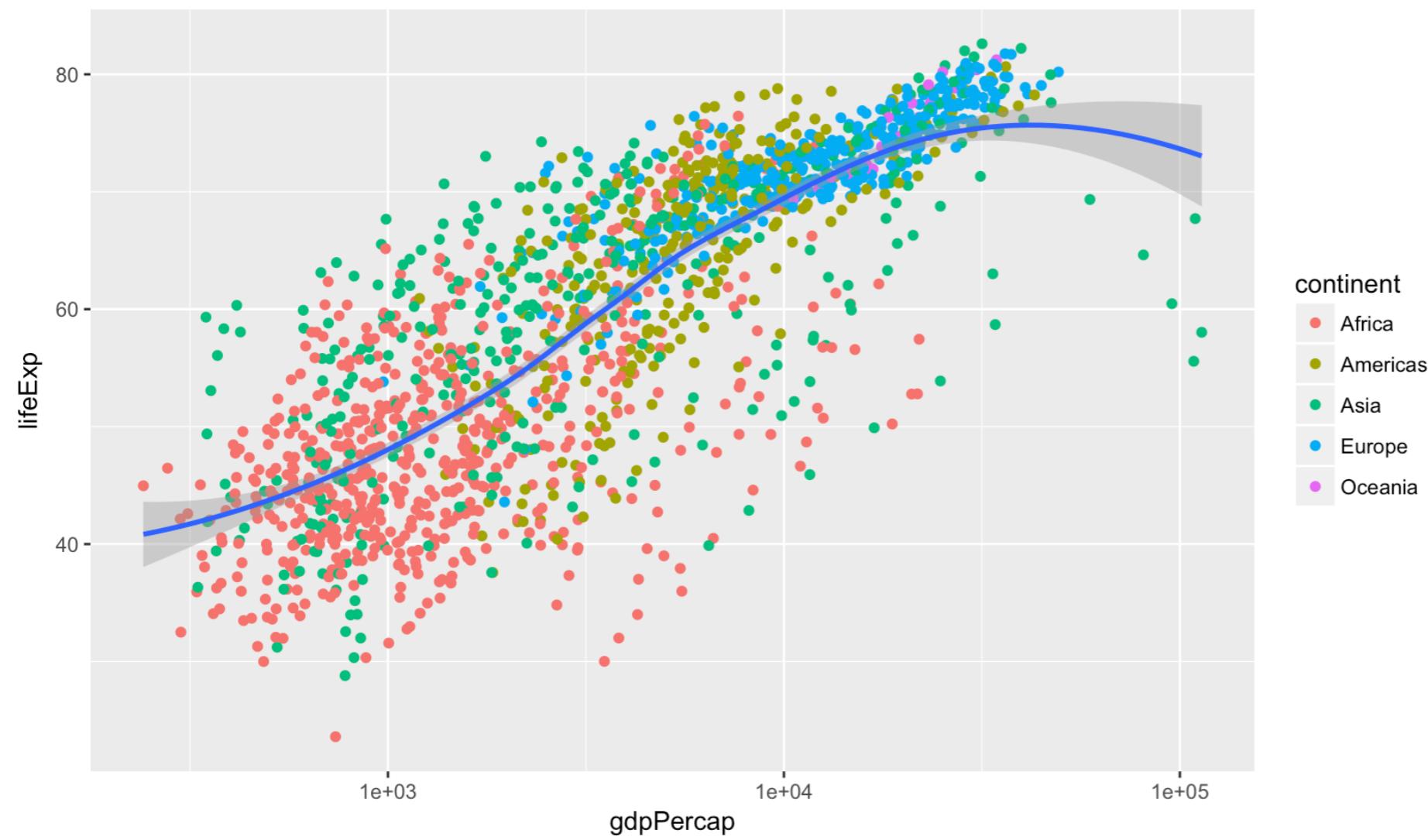
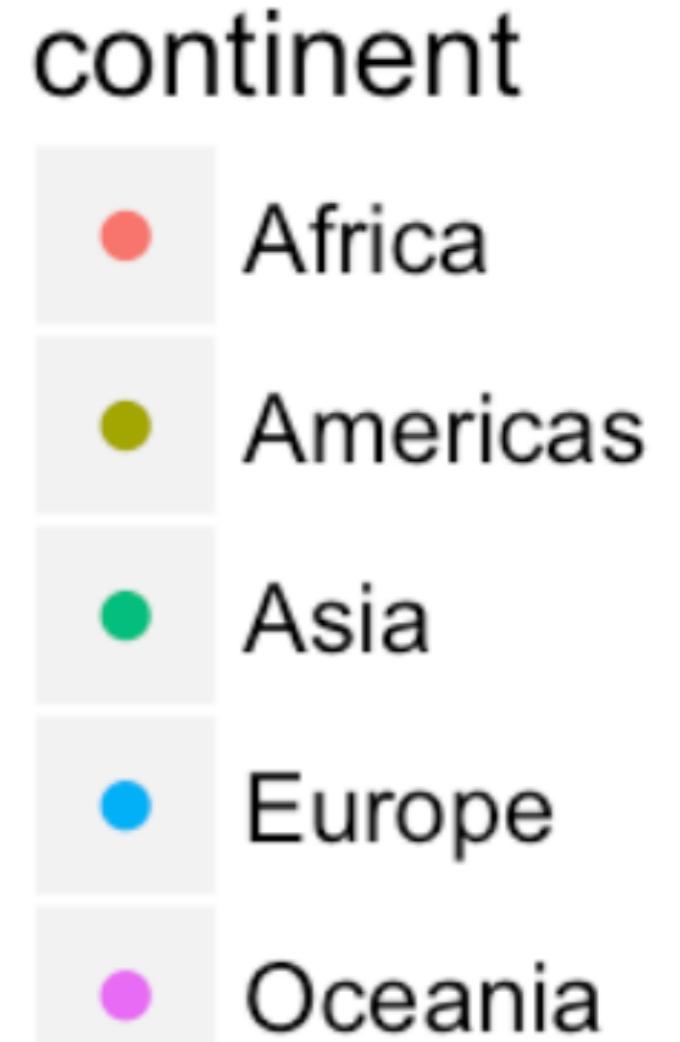
```

p <- ggplot(data = gapminder,
             mapping =
               aes(x = gdpPercap,
                   y = lifeExp,
                   color = continent,
                   fill = continent))
p + geom_point() +
  geom_smooth(method = "loess") +
  scale_x_log10()

```



```
p <- ggplot(data = gapminder,  
             mapping =  
               aes(x = gdpPercap,  
                   y = lifeExp))  
p + geom_point(mapping =  
               aes(color = continent)) +  
  geom_smooth(method = "loess") +  
  scale_x_log10()
```



Saving Your Work

ggsave()

With ggsave

ggsave("figures/my_figure.png")

ggsave("my_figure.pdf")

**ggsave("my_figure.pdf",
plot = p5,
scale = 1.2)**

**ggsave("figures/my-figure.pdf",
plot = p5,
width = 8,
height = 5)**

With `pdf()` or other graphics devices

```
pdf(file = "plot.pdf", height = 5in,  
     width = 5in)
```

▲
Open device ...

```
print(p5)
```

```
dev.off()
```



... and close when done

Within an Rmd chunk

```
```{r my_plot, fig.cap="My Plot", fig.width=9, fig.height=8}
```

```
p + geom_point(mapping =
 aes(color = continent)) +
 geom_smooth(method = "loess") +
 scale_x_log10()
```

```
...
```

## Set defaults in your first code chunk

```
knitr:::opts_chunk$set(fig.width=8, fig.height=5)
```

