

Visualization in R using ggplot2

Angela Zoss

April 2, 2018

<https://github.com/amzoss/ggplot2-DF18>

Set up environment

- R?
- RStudio?
- tidyverse?

If you haven't installed RStudio, try Docker:

<https://vm-manage.oit.duke.edu/containers>

Why visualize in R?

- Quickly explore data
- Save time switching to another tool
- Use charts to inspire new analyses and vice versa
- Reproducibility

Why care about reproducibility?

- Open science makes review easier
- Increasingly a requirement
- Saves you a lot of time trying to figure out what you did last time!

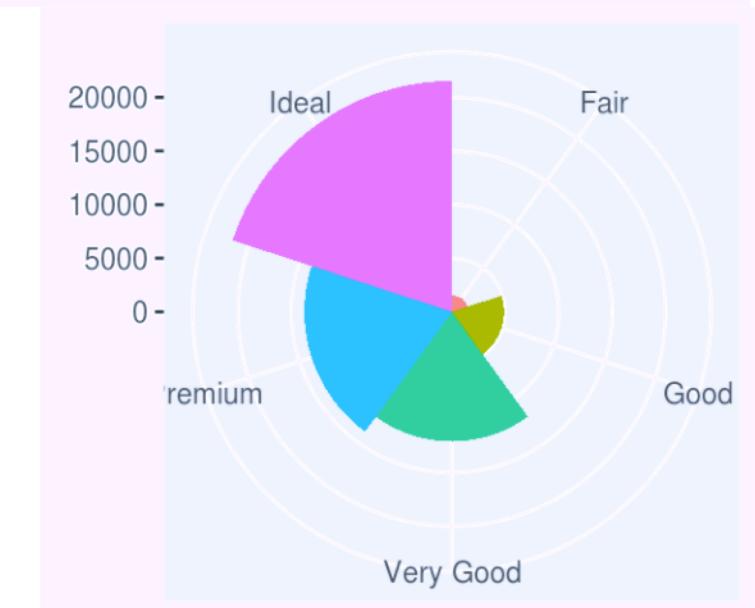
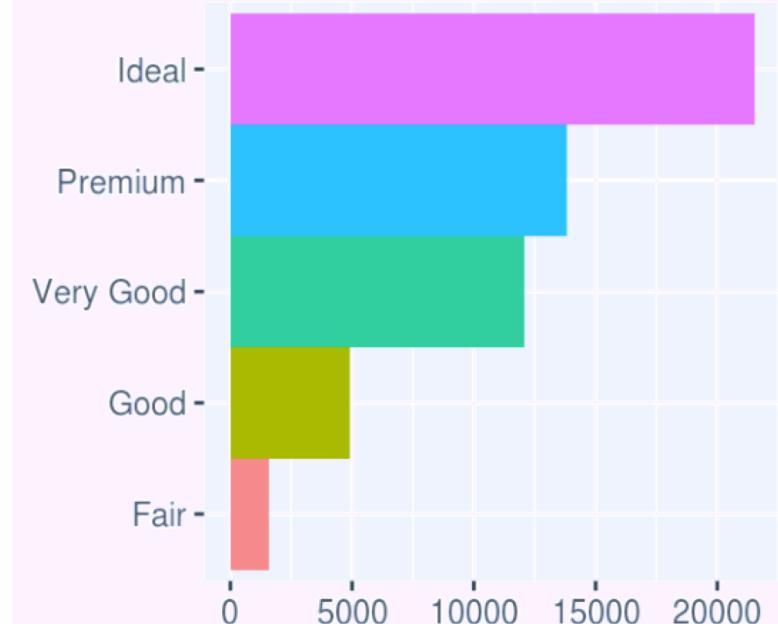
*“Your closest collaborator is **you** six months ago,
but you don’t reply to emails.”*

- *Mark Holder*

ggplot2

What is ggplot2?

an R package designed to create plots based on a theory of the grammar of graphics.



Why ggplot2 instead of base R?

- nice defaults
- easy faceting
- (arguably) more natural syntax
- can switch chart types more easily

“Why I use ggplot2”, David Robinson

<http://varianceexplained.org/r/why-i-use-ggplot2/>

Get workshop files

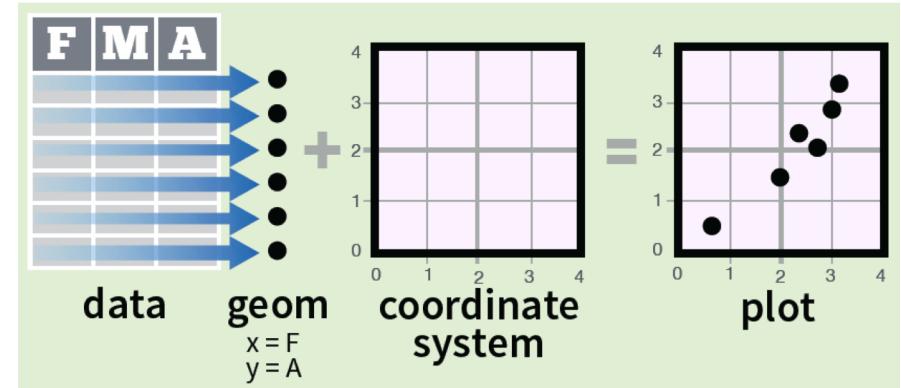
In RStudio:

- Project → New project
- Version Control
- Git
 - URL: <https://github.com/amzoss/ggplot2-DF18>
 - Project directory name: ggplot2-S18
 - Subdirectory: you choose
- Create Project

ggplot2: Elements

Basic elements in any ggplot2 visualization

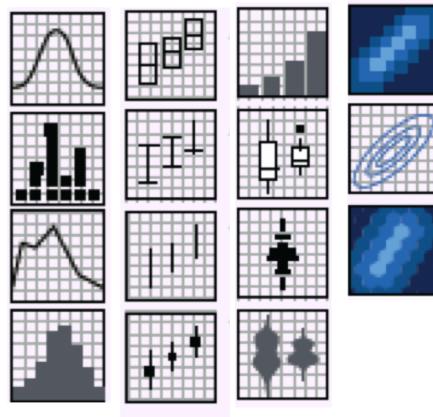
- **data**
- **aesthetics**
(variable mappings)
- **geom**
(chart type or shape)
- coordinate system
(the arrangement of the marks;
most geoms use default, cartesian)



<http://bit.ly/ggplot2-cheatsheet>

Types of geoms

- geom_bar()
- geom_point()
- geom_histogram()
- geom_map()
- etc.



<http://bit.ly/ggplot2-cheatsheet>

Note: some geoms also include data summary functions.
e.g., the “bar” geom will count data points in each category.

ggplot2: Basic syntax

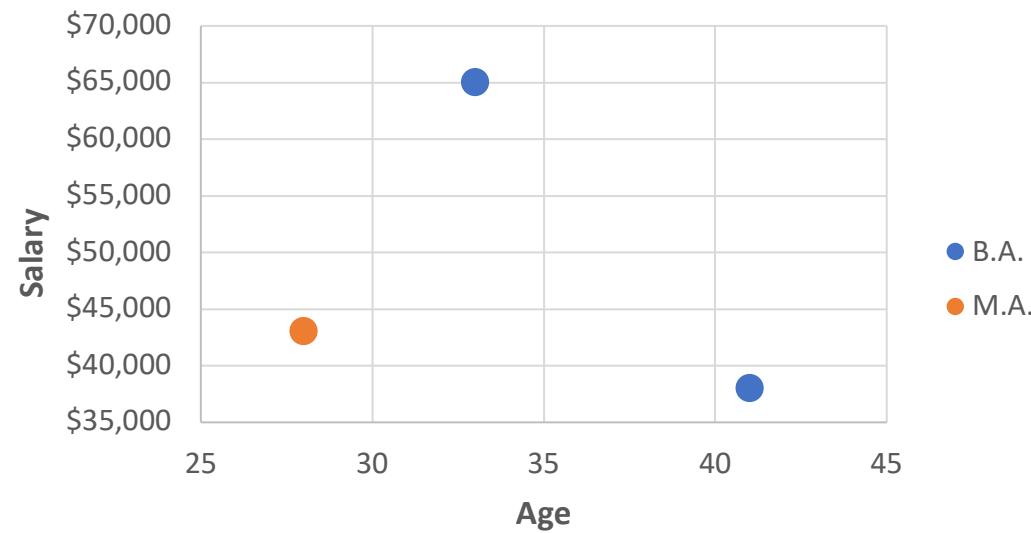
Template for a simple plot

```
ggplot() +
```

```
geom_... ( [data = data frame]  
          [aes(variable mappings)]  
          [non-variable adjustments] )
```

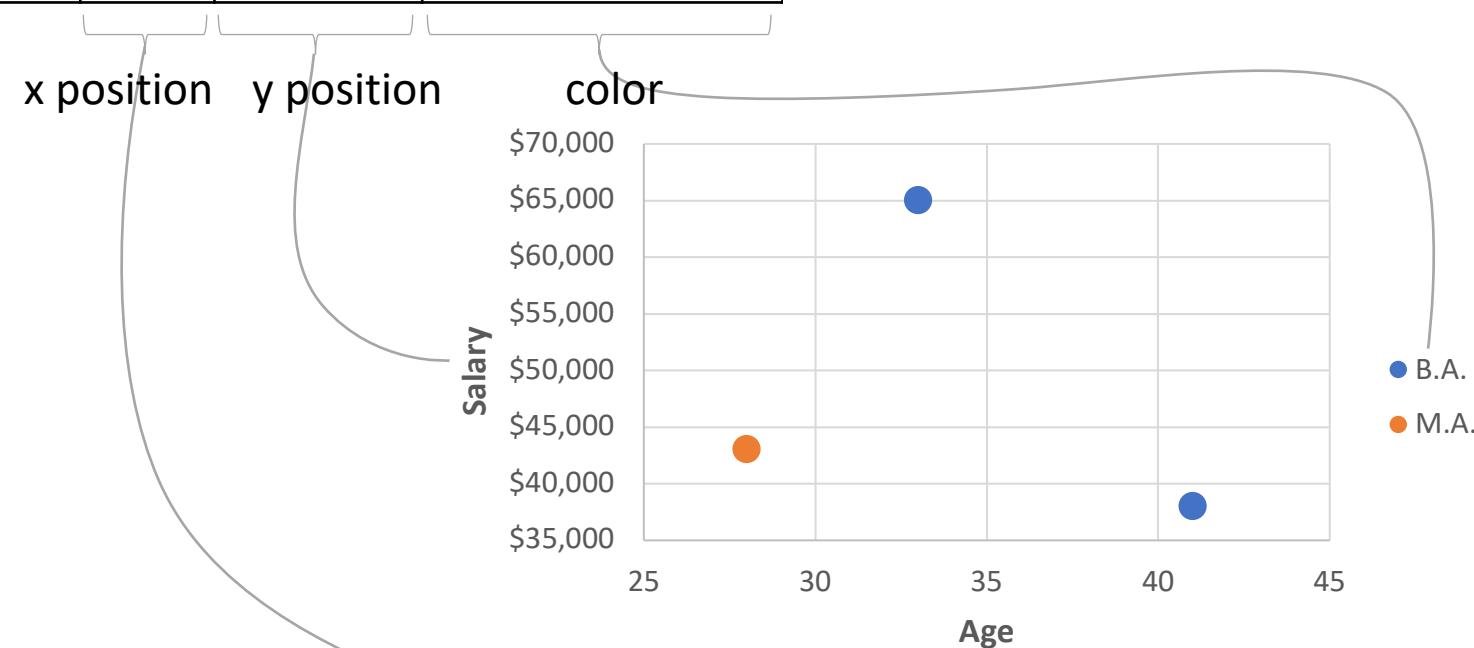
Aesthetic variable mappings

Name	Age	Salary	Highest Degree
Jane Smith	33	\$65,000	B.A.
Abby Jones	28	\$43,000	M.A.
Bridget Carden	41	\$38,000	B.A.



Aesthetic variable mappings

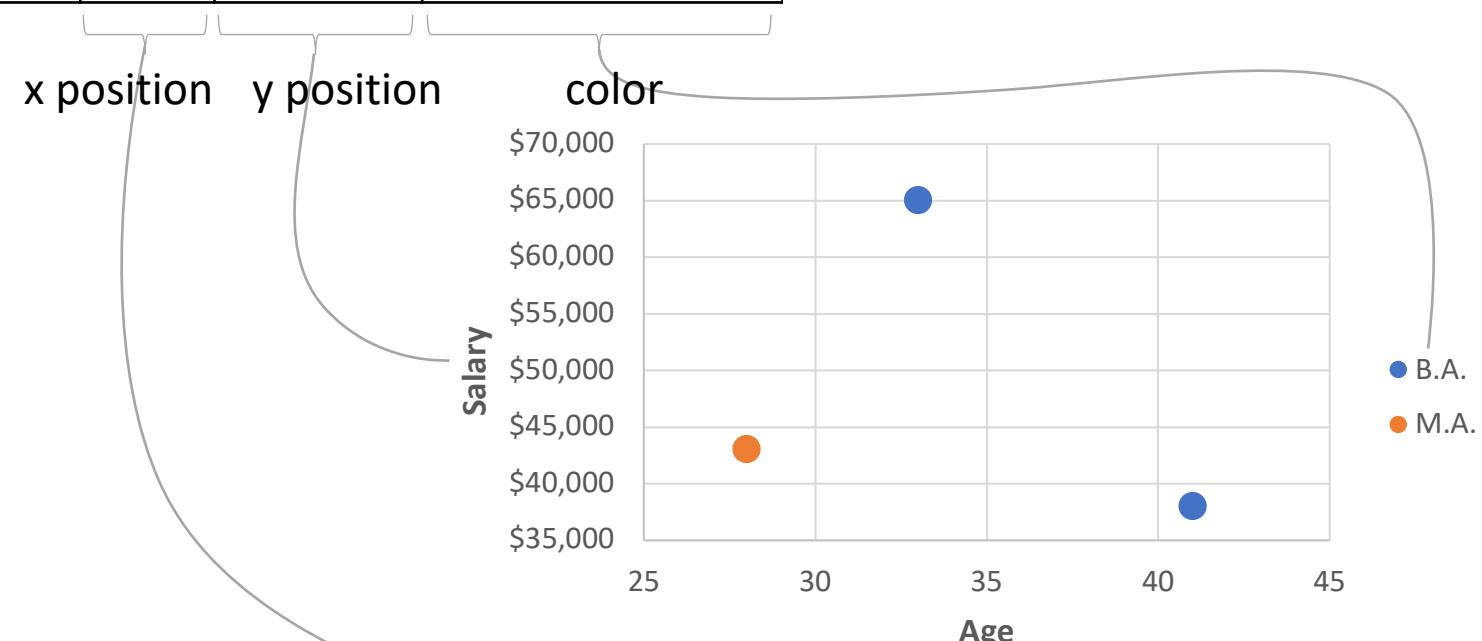
Name	Age	Salary	Highest Degree
Jane Smith	33	\$65,000	B.A.
Abby Jones	28	\$43,000	M.A.
Bridget Carden	41	\$38,000	B.A.



Aesthetic variable mappings

Name	Age	Salary	Highest Degree
Jane Smith	33	\$65,000	B.A.
Abby Jones	28	\$43,000	M.A.
Bridget Carden	41	\$38,000	B.A.

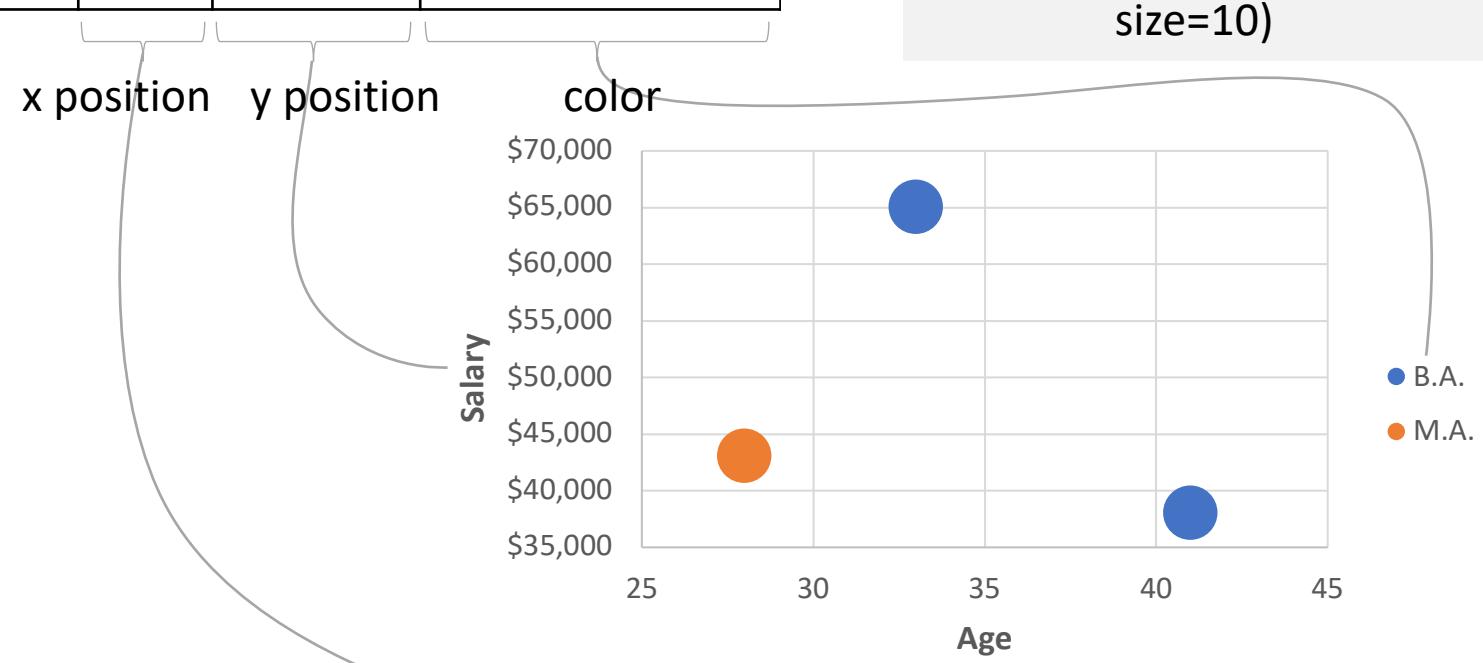
```
ggplot() +  
  geom_point(data,  
             aes(x=age,  
                  y=salary,  
                  color=degree))
```



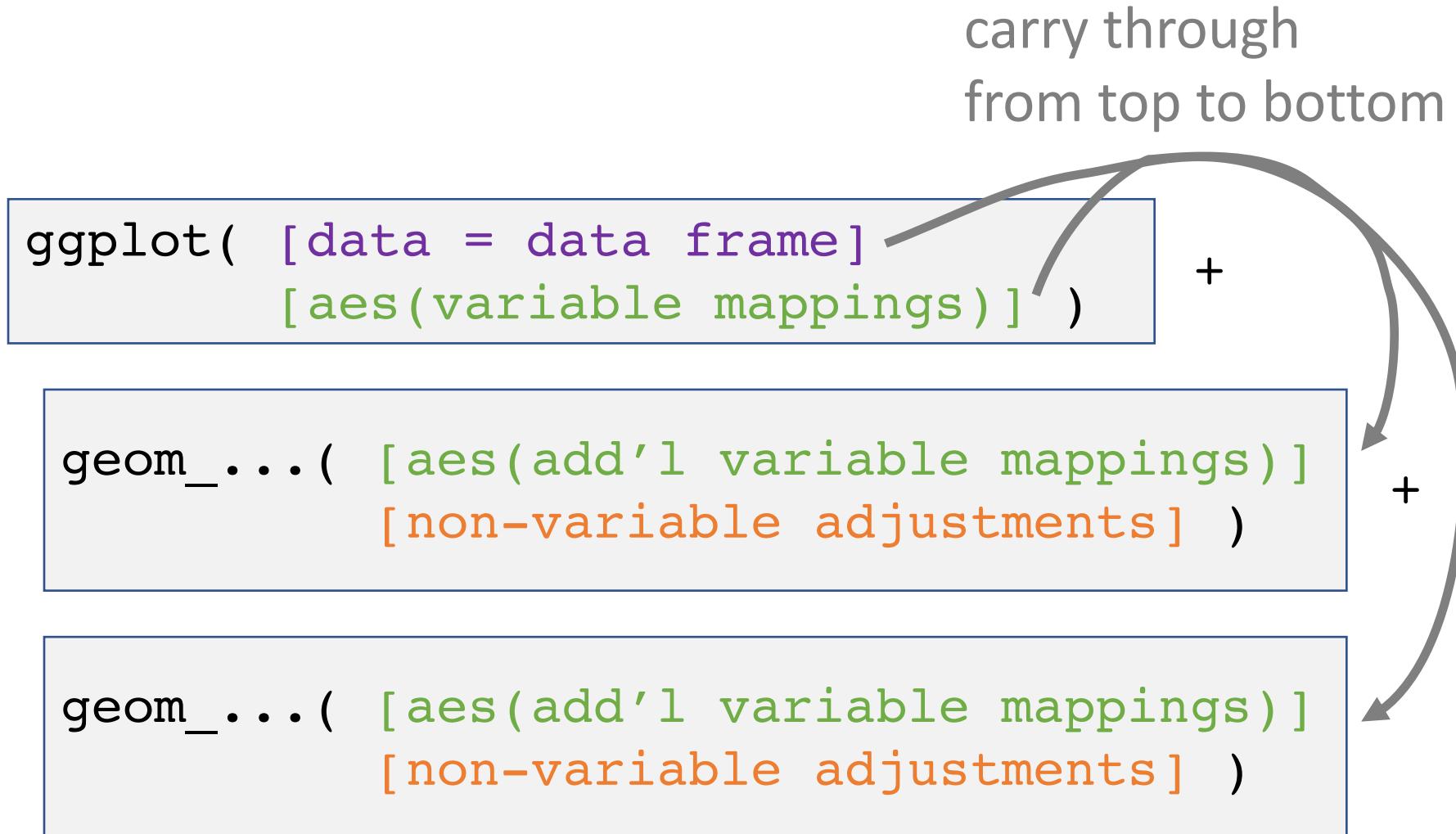
Non-variable adjustments

Name	Age	Salary	Highest Degree
Jane Smith	33	\$65,000	B.A.
Abby Jones	28	\$43,000	M.A.
Bridget Carden	41	\$38,000	B.A.

```
ggplot() +  
  geom_point(data,  
             aes(x=age,  
                  y=salary,  
                  color=degree),  
             size=10)
```



Template for a more complex plot

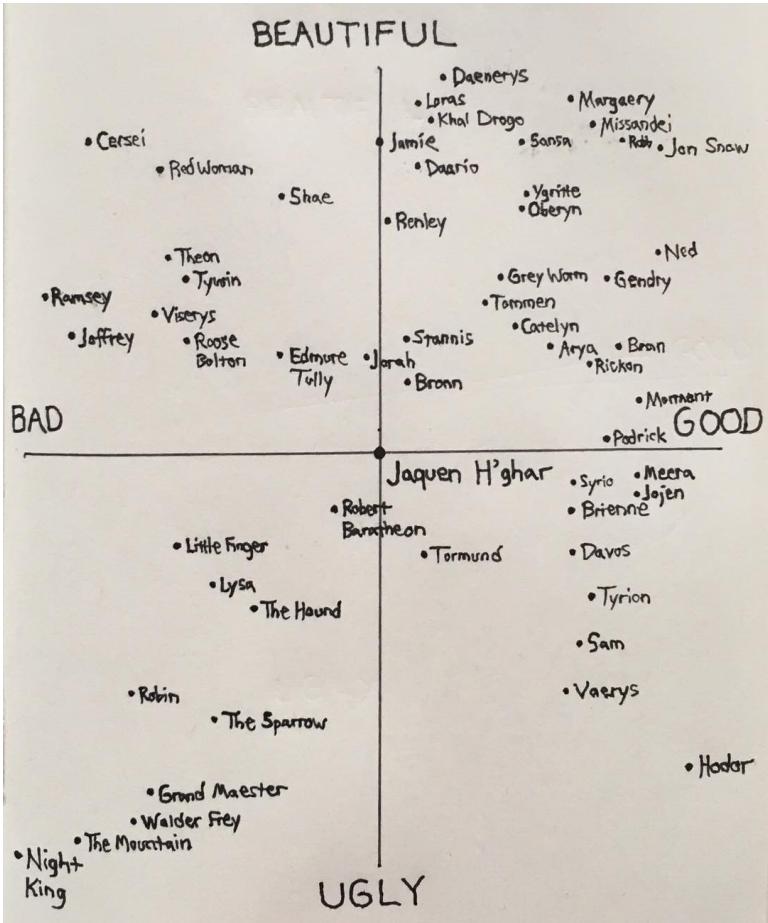


Exercise 1:

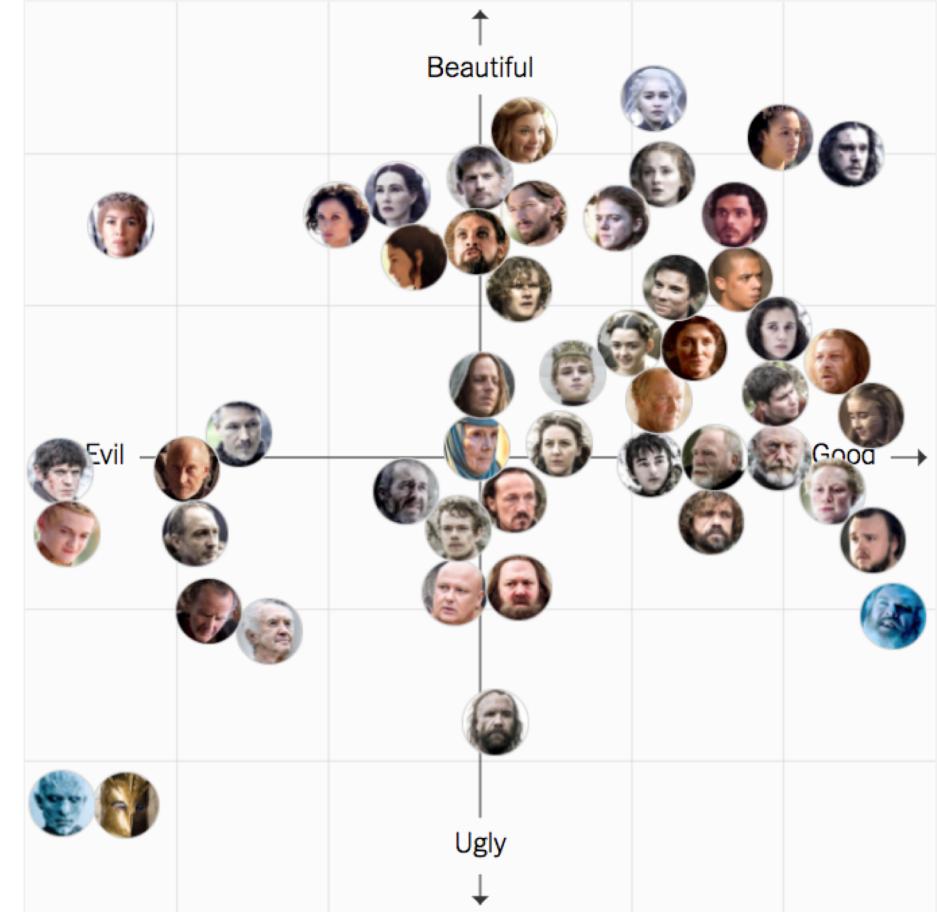
Game of Thrones character ratings

[https://www.nytimes.com/interactive/2017/08/09/upshot/game-of-thrones-
chart.html](https://www.nytimes.com/interactive/2017/08/09/upshot/game-of-thrones-chart.html)

Game of Thrones character ratings



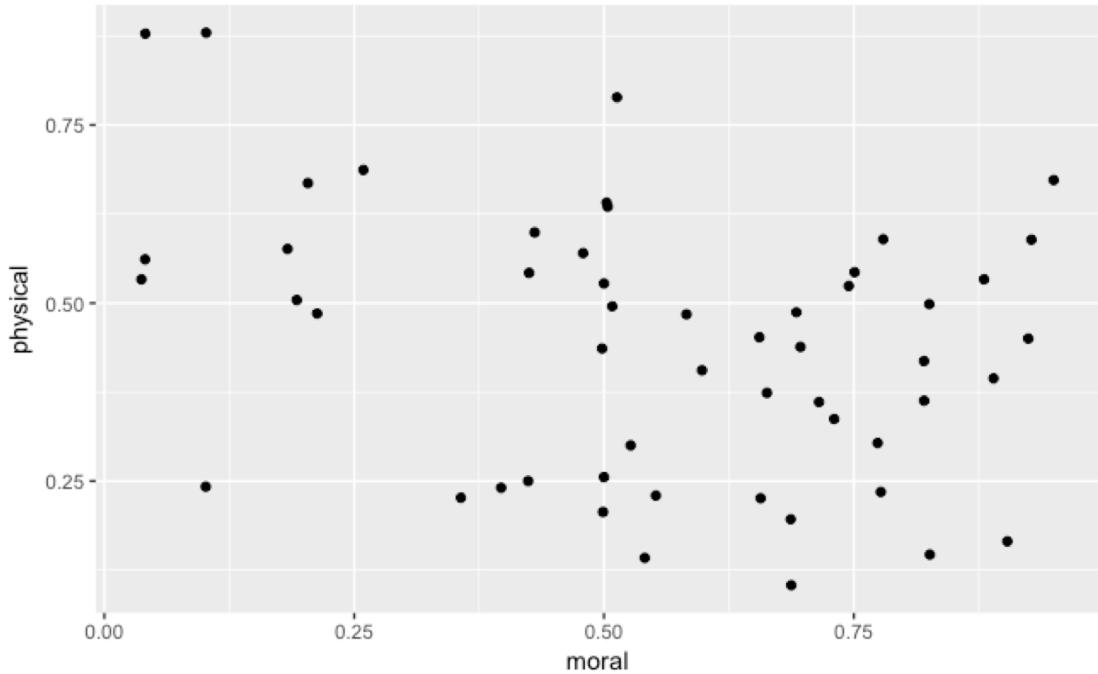
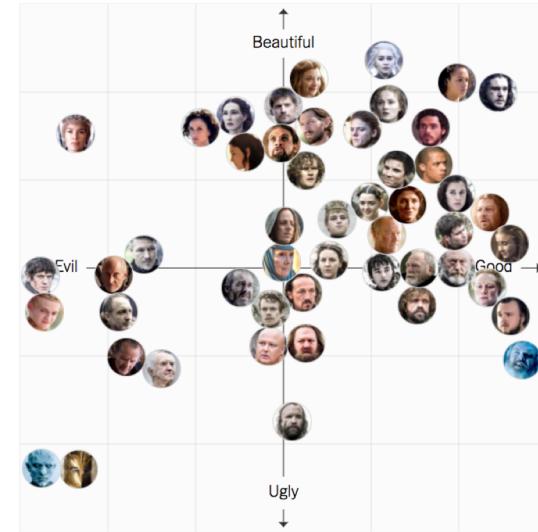
<https://www.instagram.com/p/BWnn-YogX1n/>



<https://www.nytimes.com/interactive/2017/08/09/upshot/game-of-thrones-chart.html>

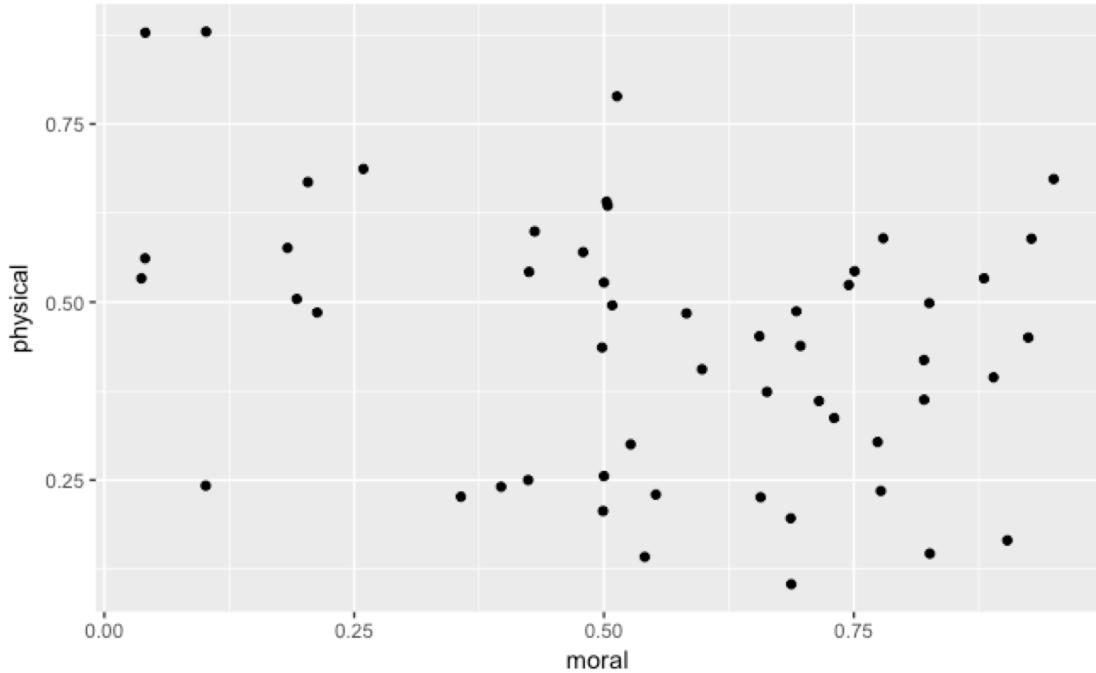
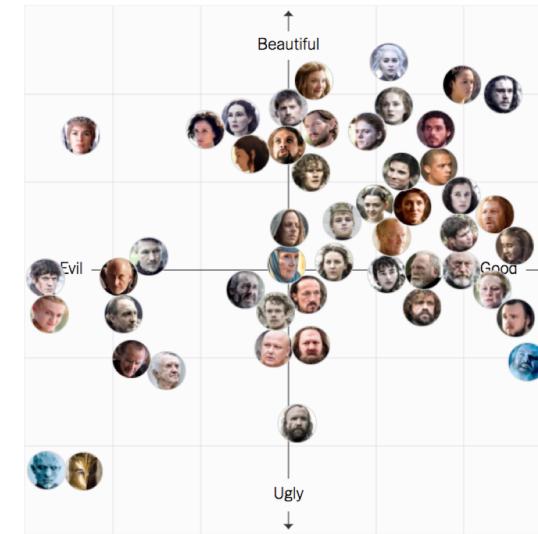
```
library(ggplot2)
```

...?



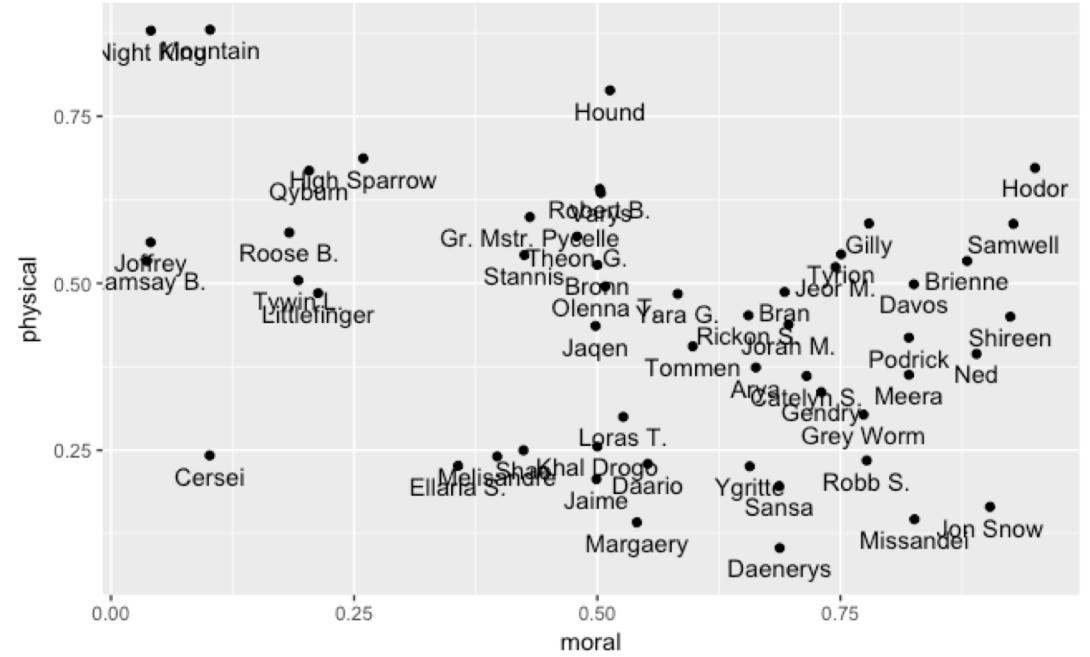
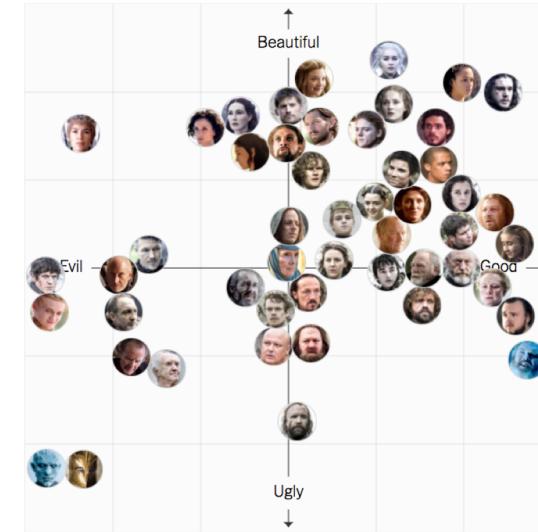
```
library(ggplot2)

ggplot(got,
       aes(x=moral,y=physical)) +
       geom_point()
```



```
library(ggplot2)
```

```
ggplot(got,  
       aes(x=moral,y=physical)) +  
       geom_point() +  
       ...?
```

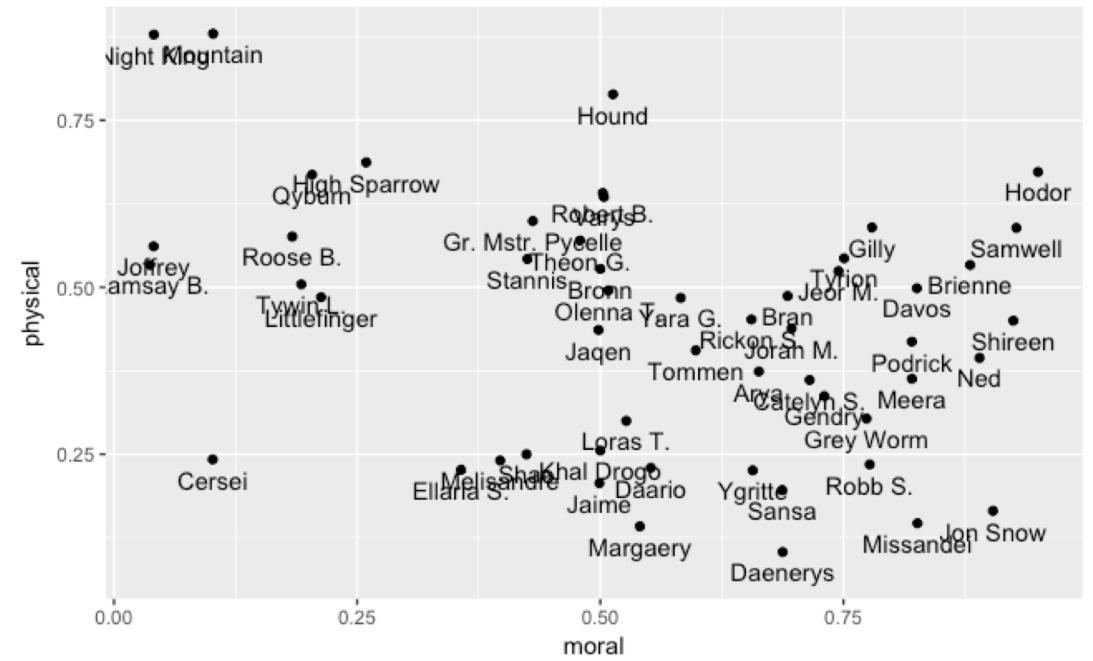
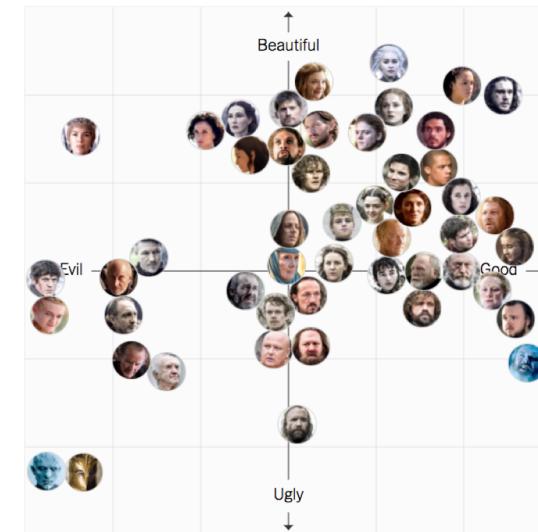


```

library(ggplot2)

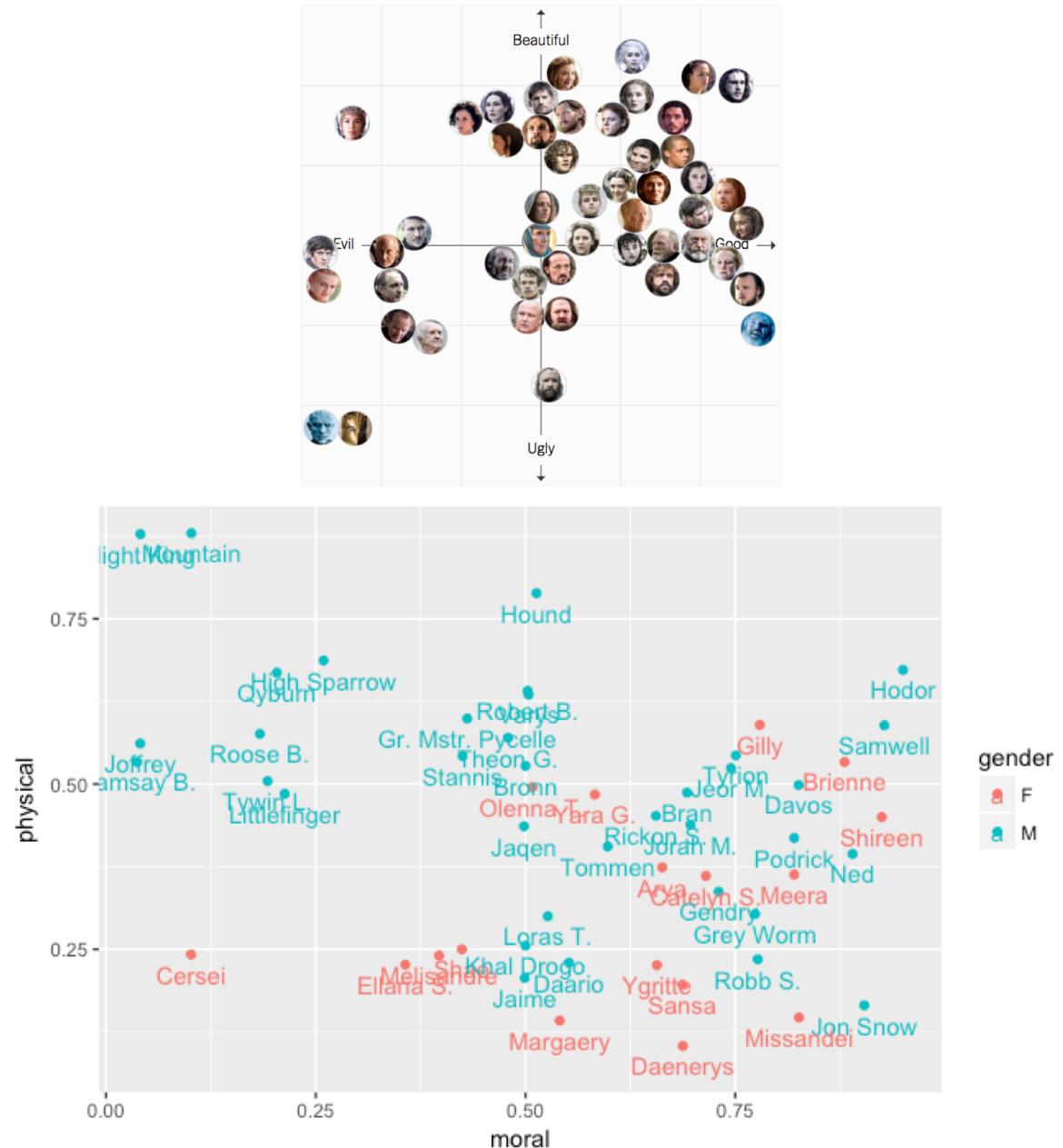
ggplot(got,
       aes(x=moral,y=physical)) +
  geom_point() +
  geom_text(aes(label=label),
            nudge_y = -.03)

```



```
library(ggplot2)

ggplot(got,
       aes(x=moral,y=physical)) +
  geom_point() +
  geom_text(aes(label=label),
            nudge_y = -.03)
```

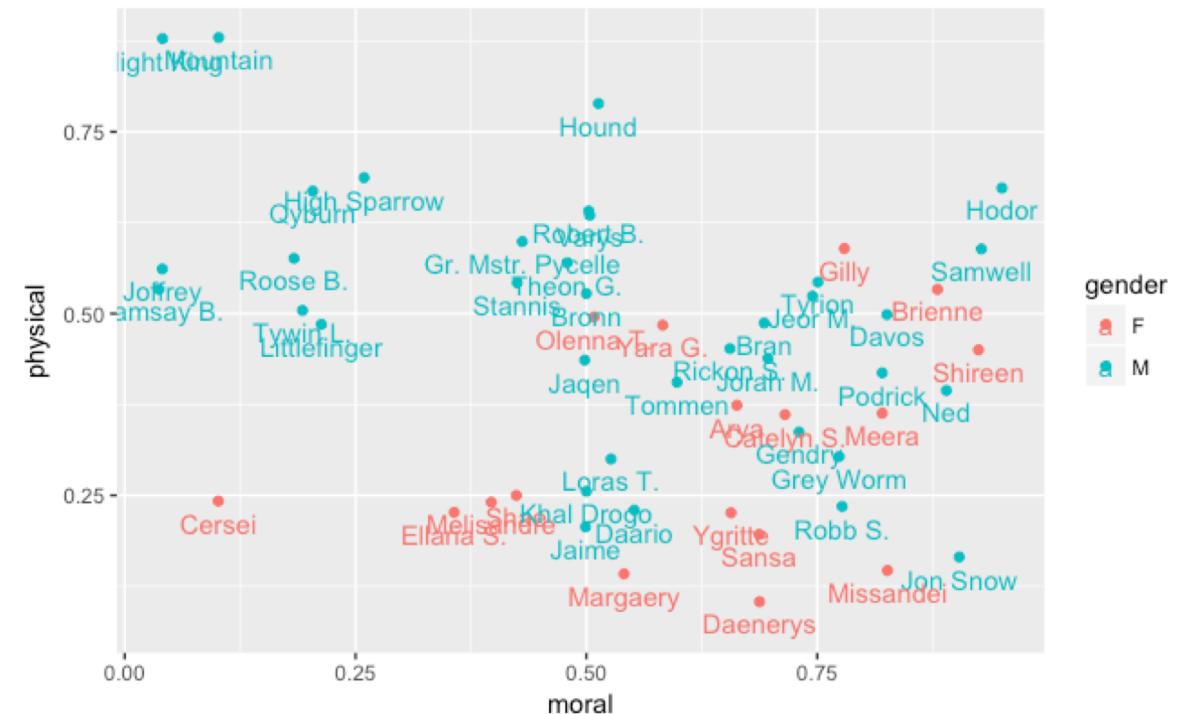
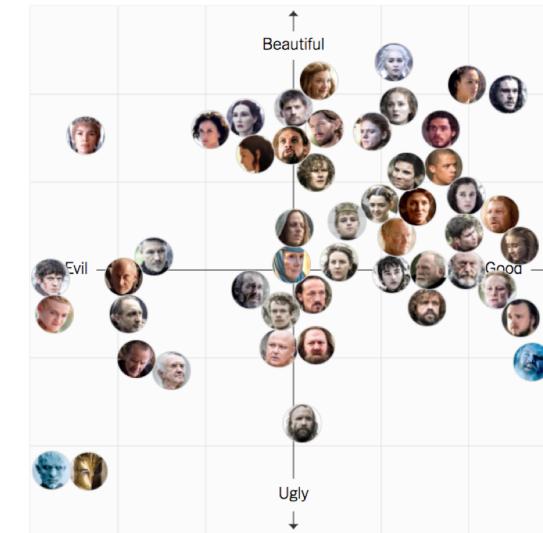


```

library(ggplot2)

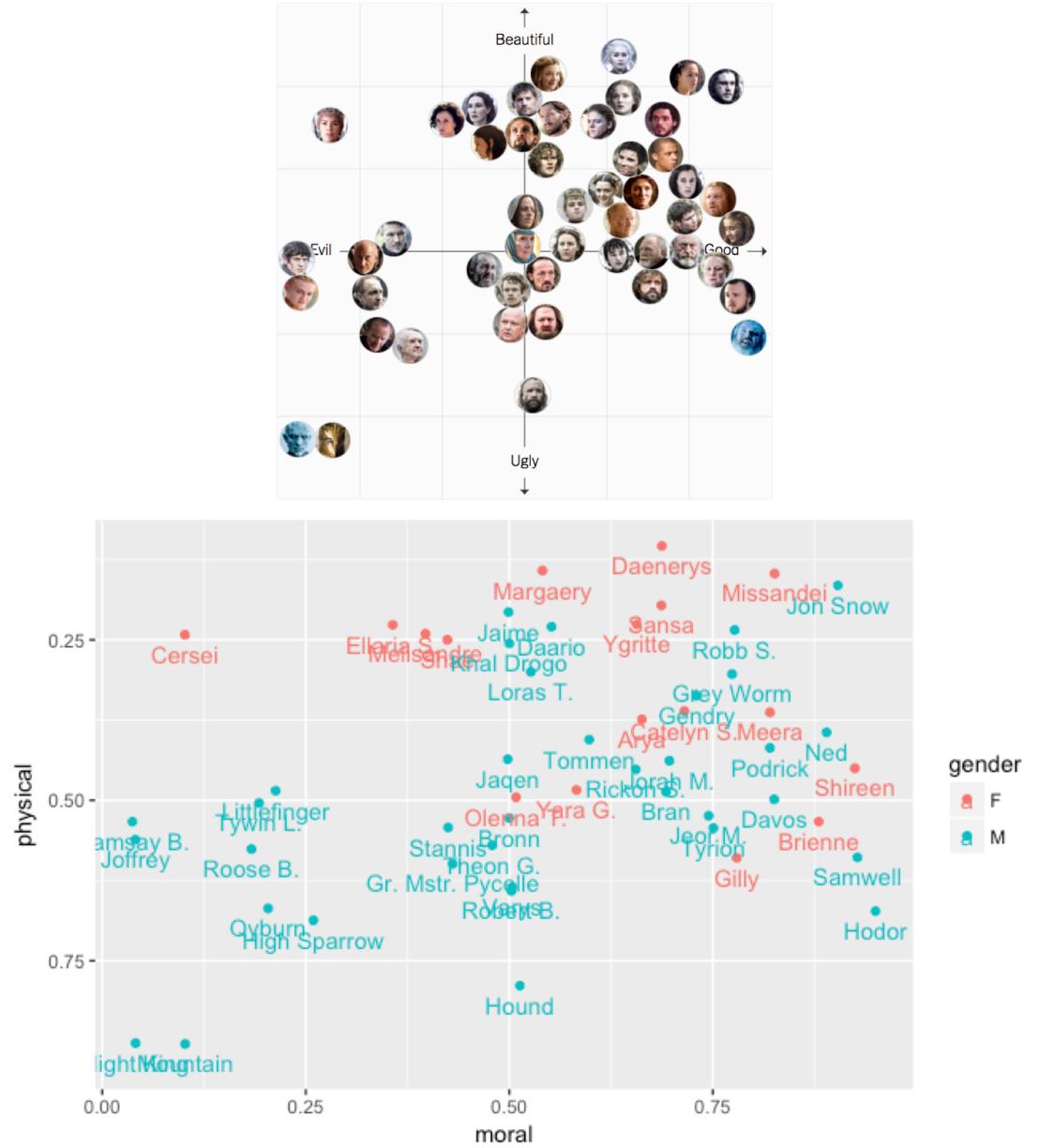
ggplot(got,
       aes(x=moral,y=physical,
            color=gender)) +
  geom_point() +
  geom_text(aes(label=label),
            nudge_y = -.03)

```



```
library(ggplot2)
```

```
ggplot(got,
       aes(x=moral,y=physical,
           color=gender)) +
  geom_point() +
  geom_text(aes(label=label),
            nudge_y = -.03) +
  ...?
```

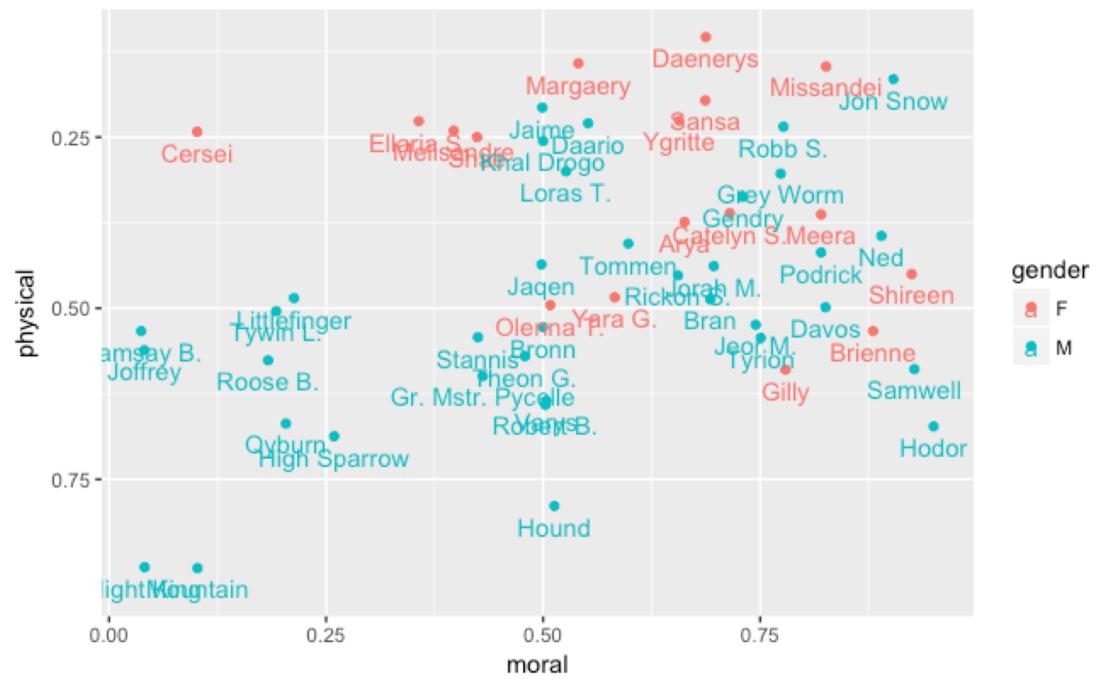
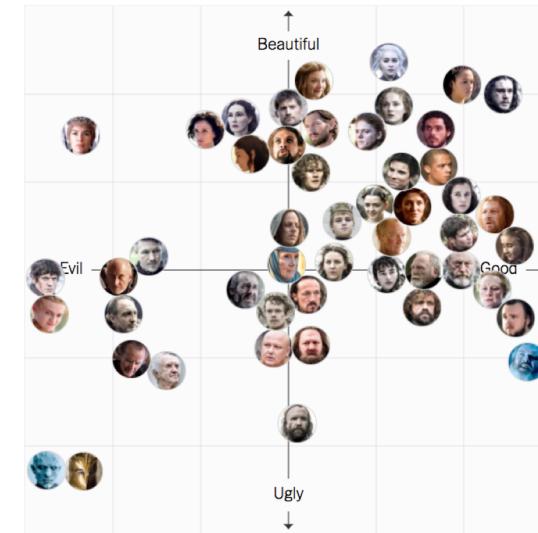


```

library(ggplot2)

ggplot(got,
       aes(x=moral,y=physical,
            color=gender)) +
  geom_point() +
  geom_text(aes(label=label),
            nudge_y = -.03) +
  scale_y_reverse()

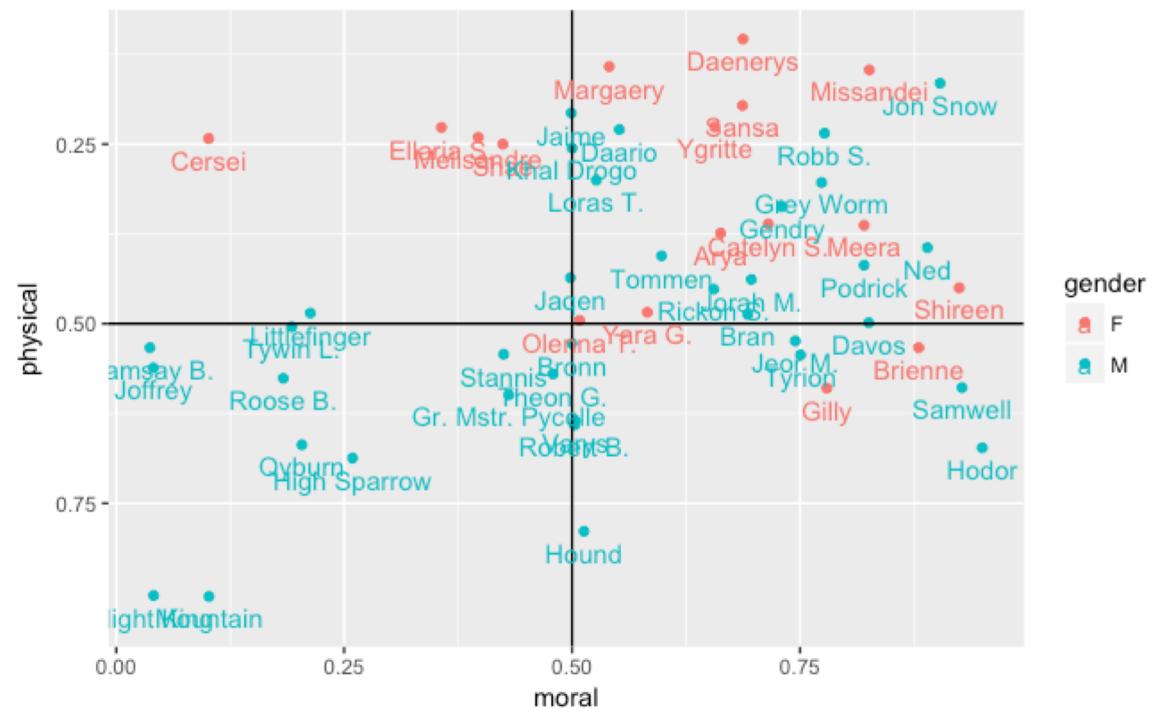
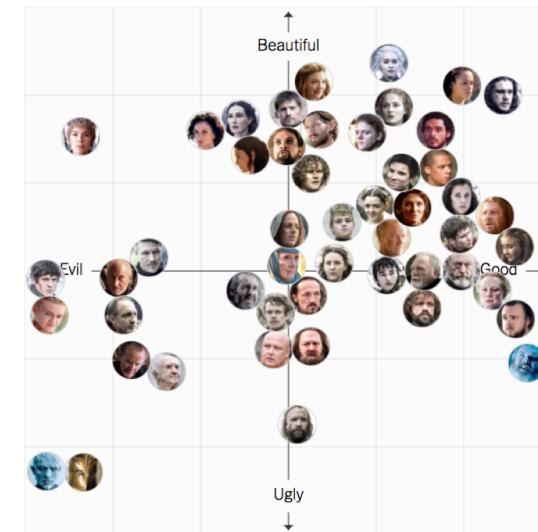
```



```

library(ggplot2)

ggplot(got,
       aes(x=moral,y=physical,
            color=gender)) +
  geom_point() +
  geom_text(aes(label=label),
            nudge_y = -.03) +
  scale_y_reverse() +
  ...
?
```

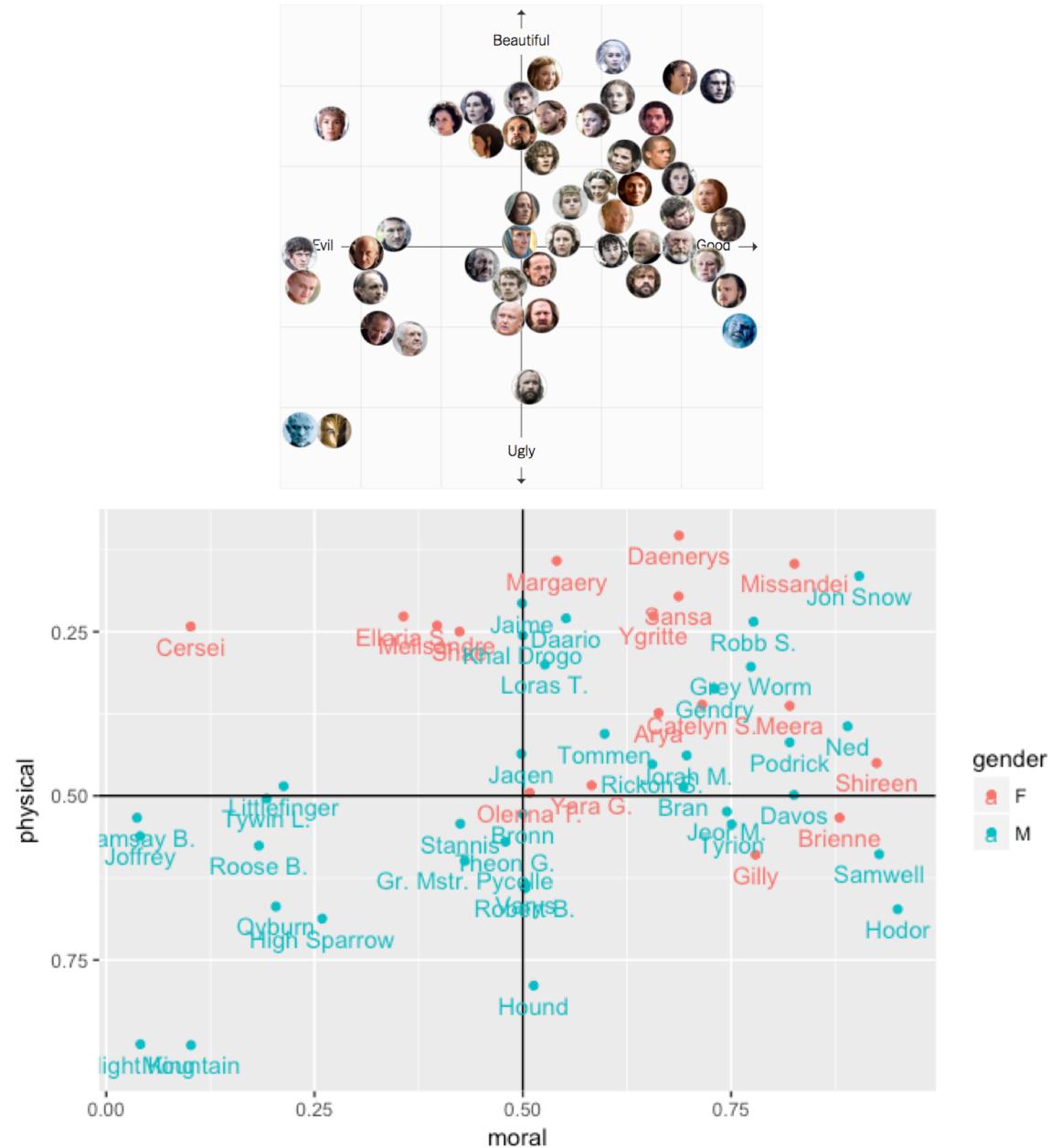


```

library(ggplot2)

ggplot(got,
       aes(x=moral,y=physical,
            color=gender)) +
  geom_point() +
  geom_hline(yintercept=.5) +
  geom_vline(xintercept=.5) +
  geom_text(aes(label=label),
            nudge_y = -.03) +
  scale_y_reverse()

```

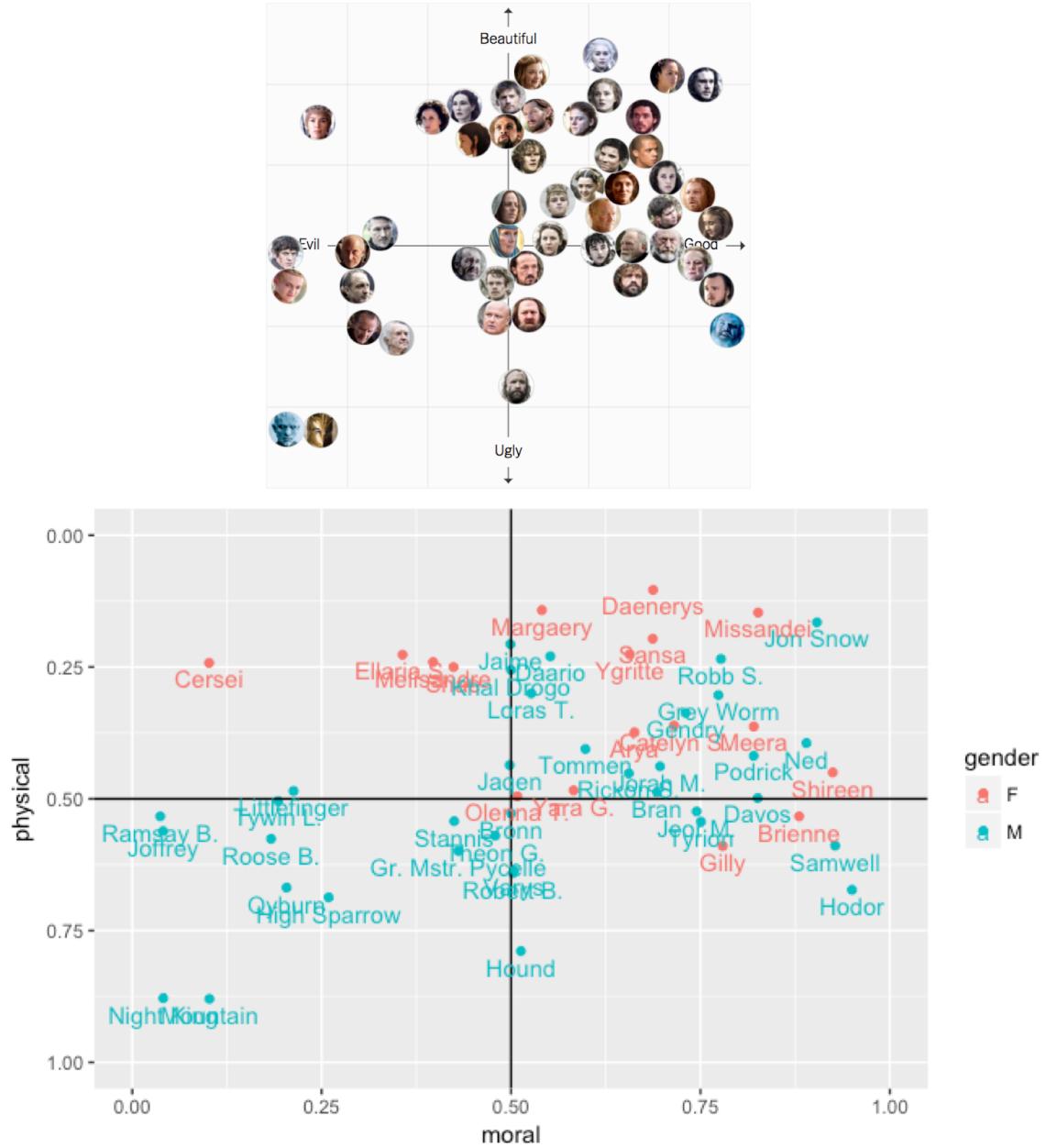


```

library(ggplot2)

ggplot(got,
       aes(x=moral,y=physical,
            color=gender)) +
  geom_point() +
  geom_hline(yintercept=.5) +
  geom_vline(xintercept=.5) +
  geom_text(aes(label=label),
            nudge_y = -.03) +
  scale_y_reverse() +
  ...
?

```



```

library(ggplot2)

ggplot(got,
       aes(x=moral,y=physical,
            color=gender)) +
  geom_point() +
  geom_hline(yintercept=.5) +
  geom_vline(xintercept=.5) +
  geom_text(aes(label=label),
            nudge_y = -.03) +
  scale_y_reverse(lim=c(1,0)) +
  scale_x_continuous(
    limits=c(0,1))

```

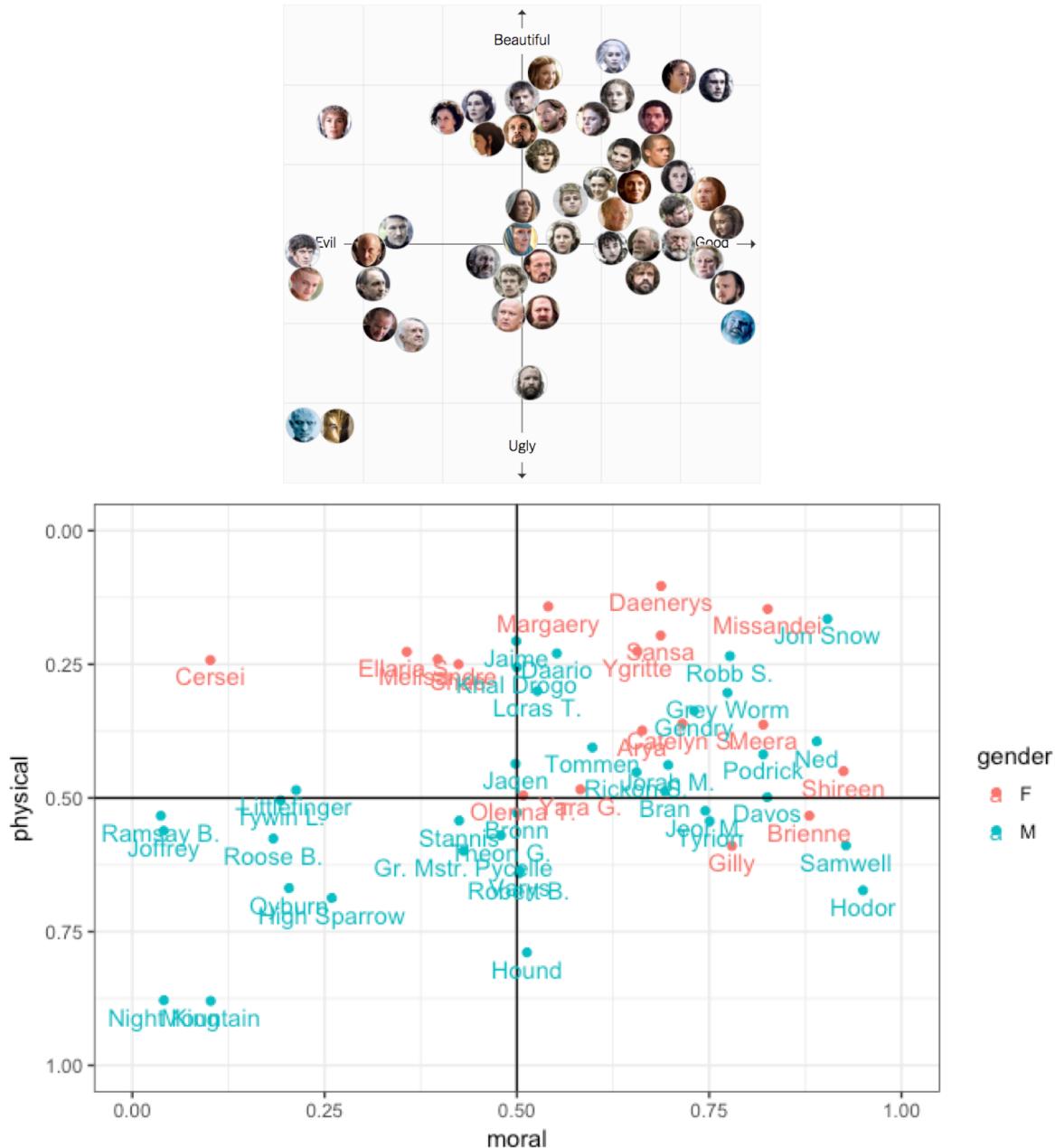


```

library(ggplot2)

ggplot(got,
       aes(x=moral,y=physical,
            color=gender)) +
  geom_point() +
  geom_hline(yintercept=.5) +
  geom_vline(xintercept=.5) +
  geom_text(aes(label=label),
            nudge_y = -.03) +
  scale_y_reverse(lim=c(1,0)) +
  scale_x_continuous(
    limits=c(0,1)) +
  ...
?

```



```
library(ggplot2)

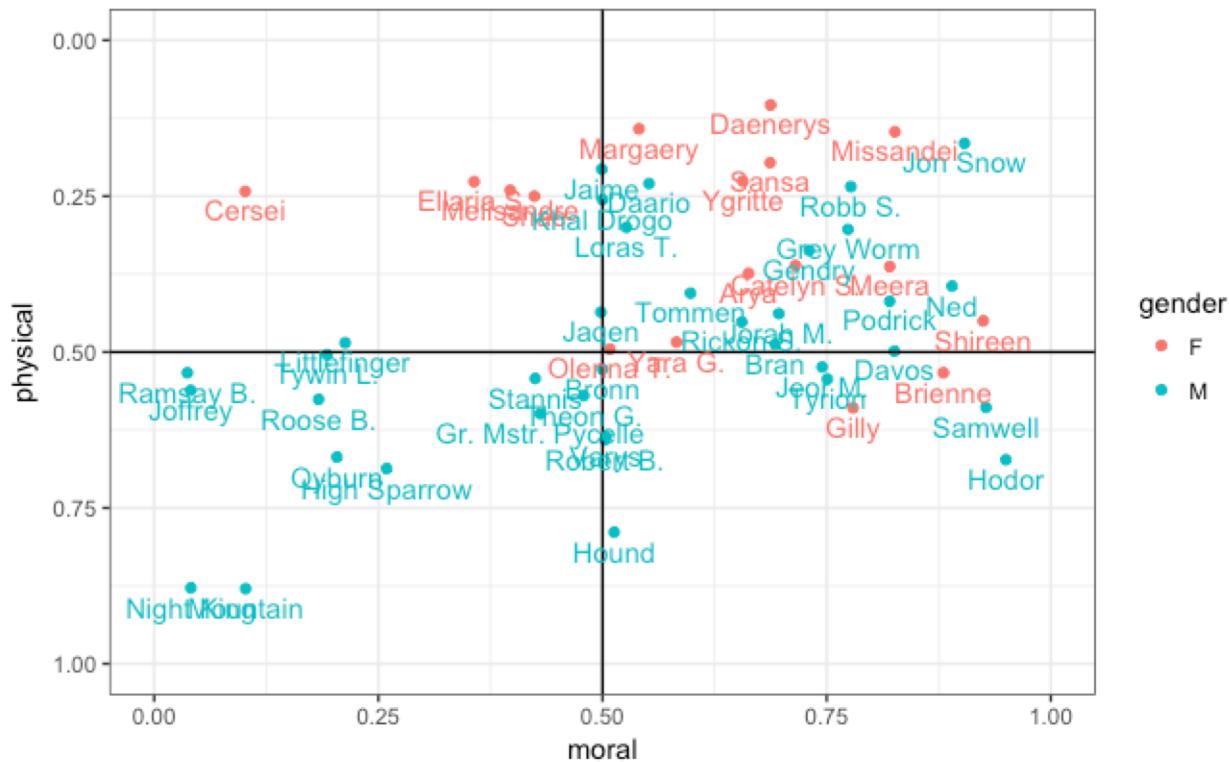
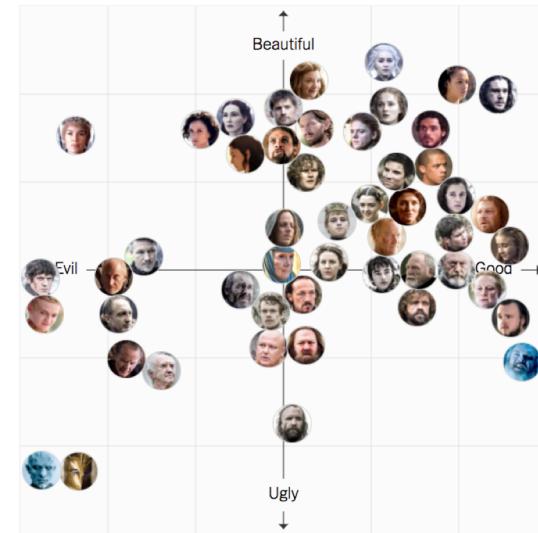
ggplot(got,
       aes(x=moral,y=physical,
           color=gender)) +
  geom_point() +
  geom_hline(yintercept=.5) +
  geom_vline(xintercept=.5) +
  geom_text(aes(label=label),
            nudge_y = -.03) +
  scale_y_reverse(lim=c(1,0)) +
  scale_x_continuous(
    limits=c(0,1)) +
  theme_bw()
```



```

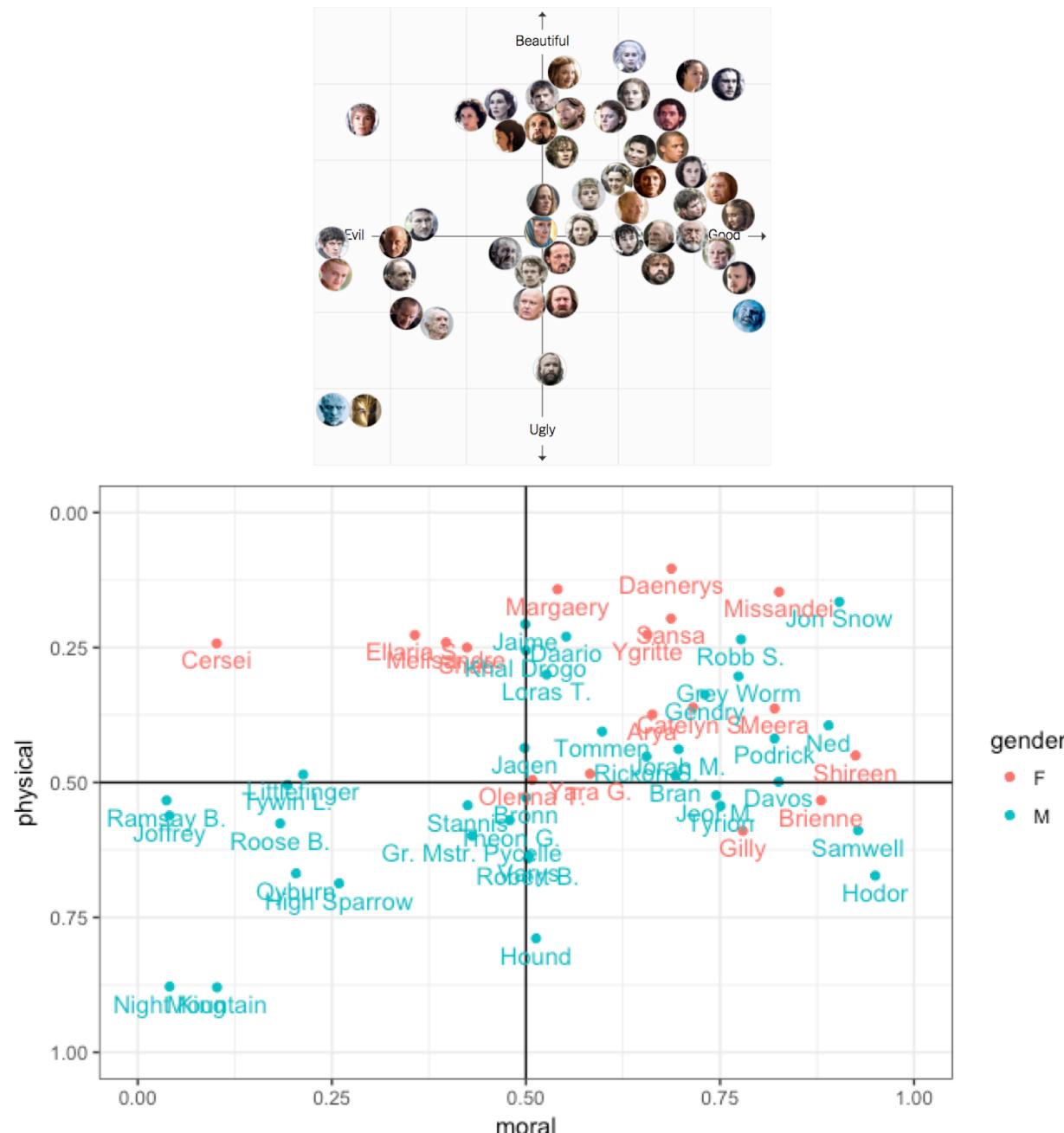
library(ggplot2)

ggplot(got,
       aes(x=moral,y=physical,
            color=gender)) +
  geom_point() +
  geom_hline(yintercept=.5) +
  geom_vline(xintercept=.5) +
  geom_text(aes(label=label),
            nudge_y = -.03) +
  scale_y_reverse(lim=c(1,0)) +
  scale_x_continuous(
    limits=c(0,1)) +
  theme_bw() +
  ...
?
```



```
library(ggplot2)

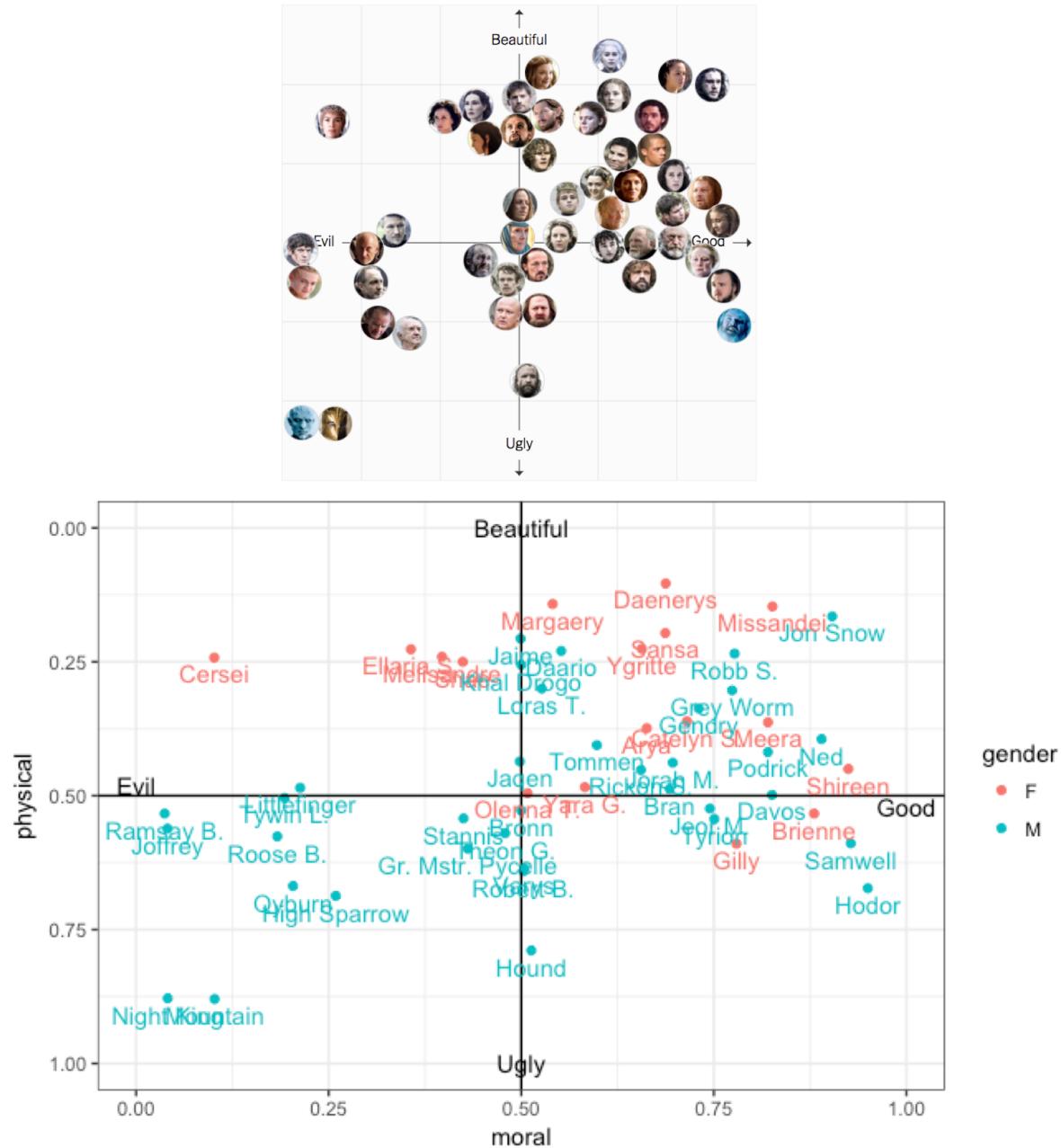
ggplot(got,
       aes(x=moral,y=physical,
            color=gender)) +
  geom_point() +
  geom_hline(yintercept=.5) +
  geom_vline(xintercept=.5) +
  geom_text(aes(label=label),
            nudge_y = -.03,
            show.legend = FALSE) +
  scale_y_reverse(lim=c(1,0)) +
  scale_x_continuous(
    limits=c(0,1)) +
  theme_bw()
```



```

library(ggplot2)

ggplot(got,
       aes(x=moral,y=physical,
            color=gender)) +
  geom_point() +
  geom_hline(yintercept=.5) +
  geom_vline(xintercept=.5) +
  geom_text(aes(label=label),
            nudge_y = -.03,
            show.legend = FALSE) +
  scale_y_reverse(lim=c(1,0)) +
  scale_x_continuous(
    limits=c(0,1)) +
  theme_bw() +
  ...
?
```

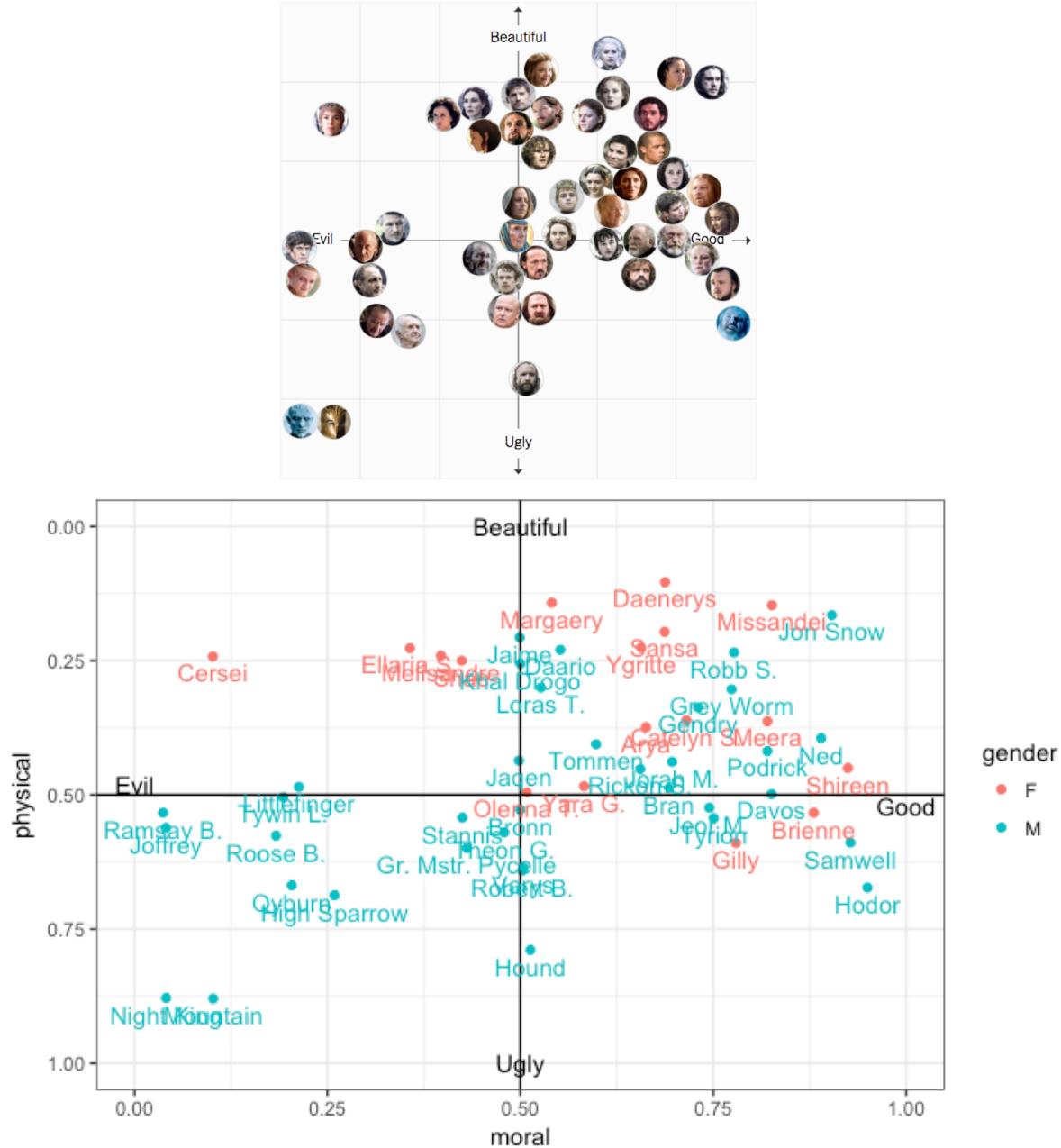


```

library(ggplot2)

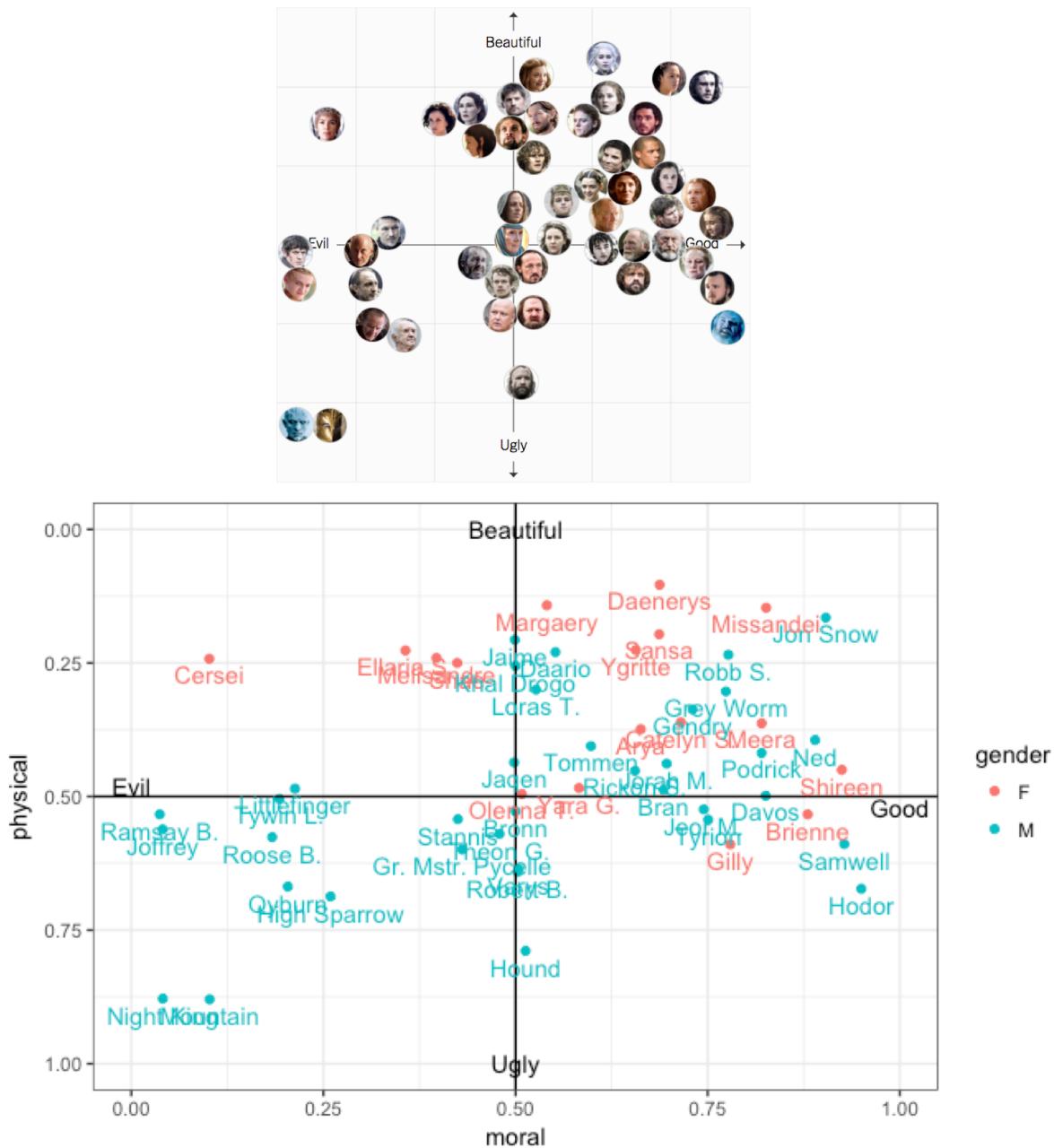
ggplot(got,
       aes(x=moral,y=physical,
            color=gender)) +
  geom_point() +
  geom_hline(yintercept=.5) +
  geom_vline(xintercept=.5) +
  geom_text(aes(label=label),
            nudge_y = -.03,
            show.legend = FALSE) +
  scale_y_reverse(lim=c(1,0)) +
  scale_x_continuous(
    limits=c(0,1)) +
  theme_bw() +
  annotate(geom = "text", x=.5,
          y=1, label="Ugly")

```



```
library(ggplot2)

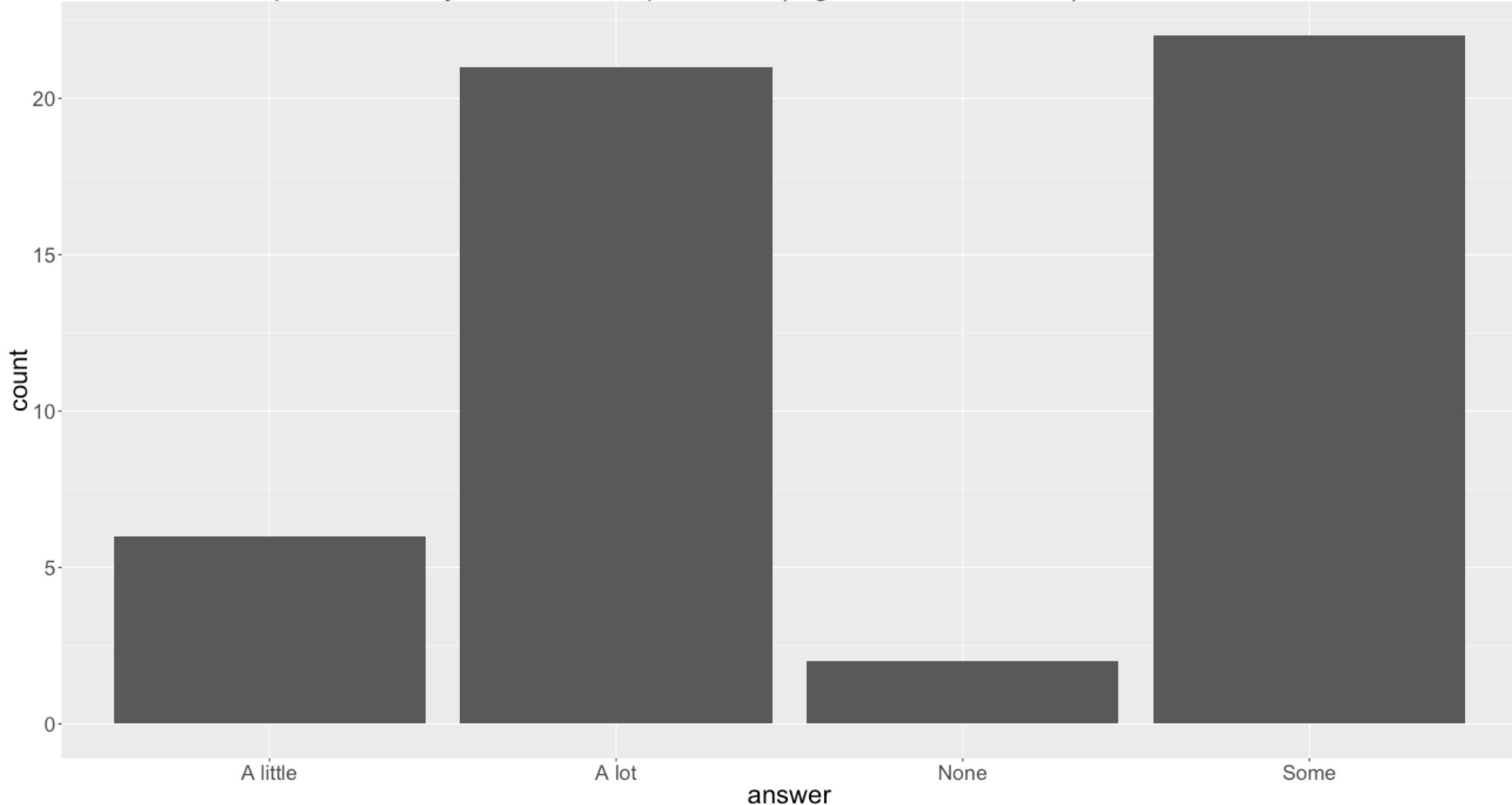
ggplot(got,
       aes(x=moral, y=physical,
           color=gender)) +
  geom_point() +
  geom_hline(yintercept=.5) +
  geom_vline(xintercept=.5) +
  geom_text(aes(label=label),
            nudge_y = - .03,
            show.legend = FALSE) +
  scale_y_reverse(lim=c(1,0)) +
  scale_x_continuous(
    limits=c(0,1)) +
  theme_bw() +
  annotate(geom = "text", x=.5,
           y=1, label="Ugly") +
  annotate(geom = "text", x=.5,
           y=0, label="Beautiful") +
  annotate(geom = "text", x=0,
           y=.48, label="Evil") +
  annotate(geom = "text", x=1,
           y=.52, label="Good")
```



Principles for Effective Visualizations

Principle 1: Order matters

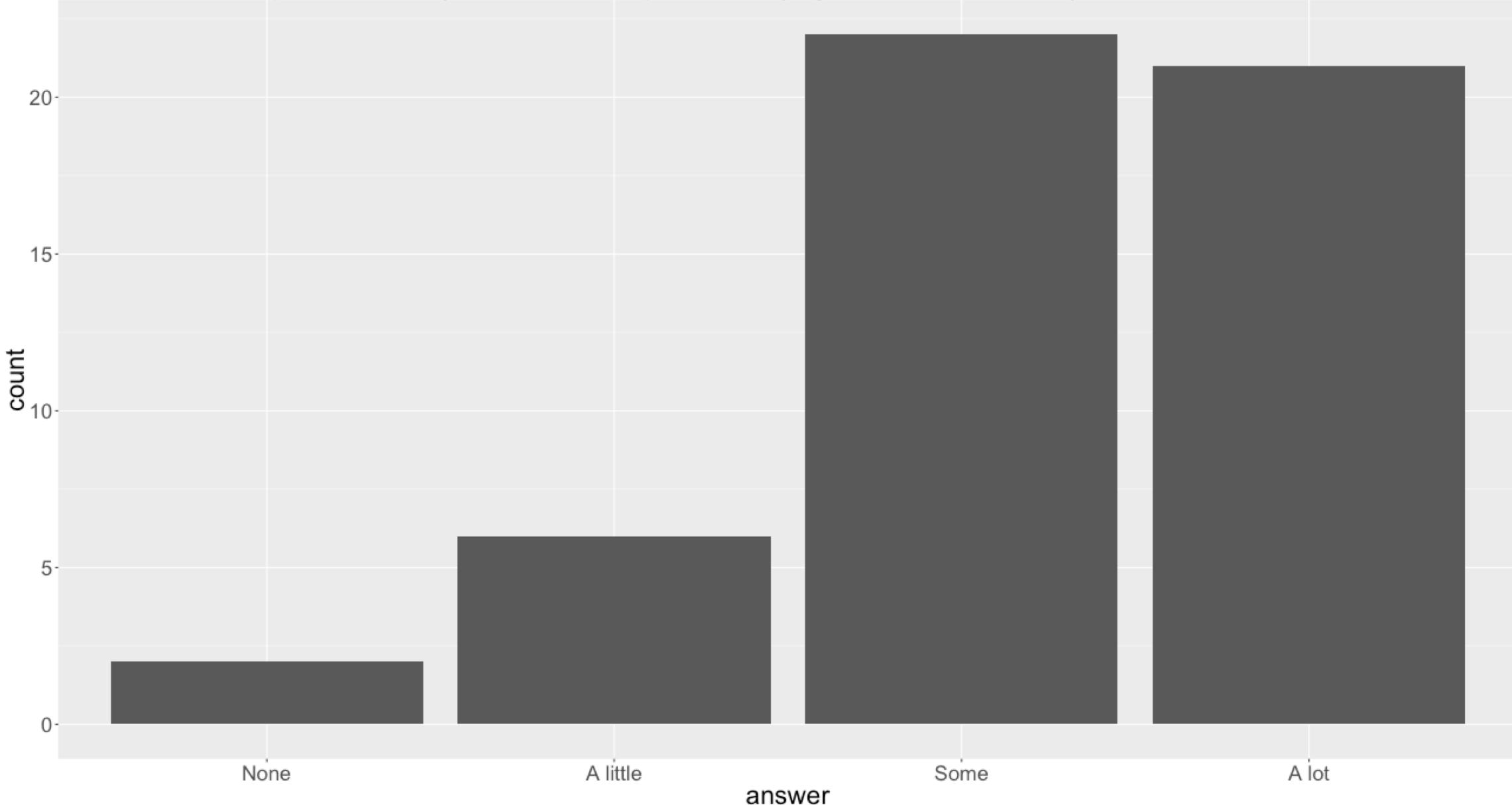
How much experience do you have as a producer (e.g., reader, follower) of network science research?



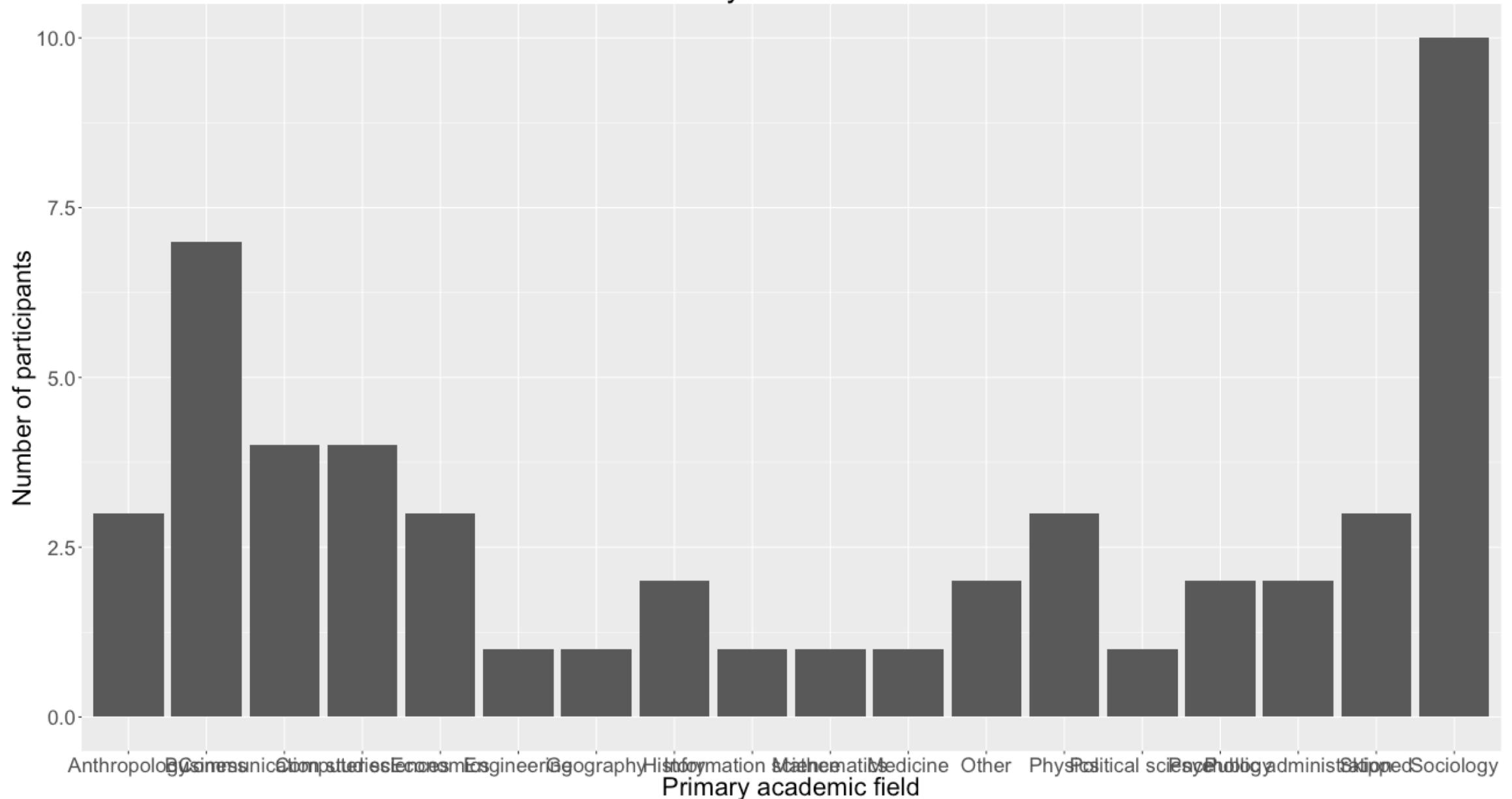
Order by meaning

```
data$answer <-  
  factor(data$answer,  
         levels=c("None", "A little", "Some", "A lot"),  
         ordered = TRUE)
```

How much experience do you have as a producer (e.g., reader, follower) of network science research?



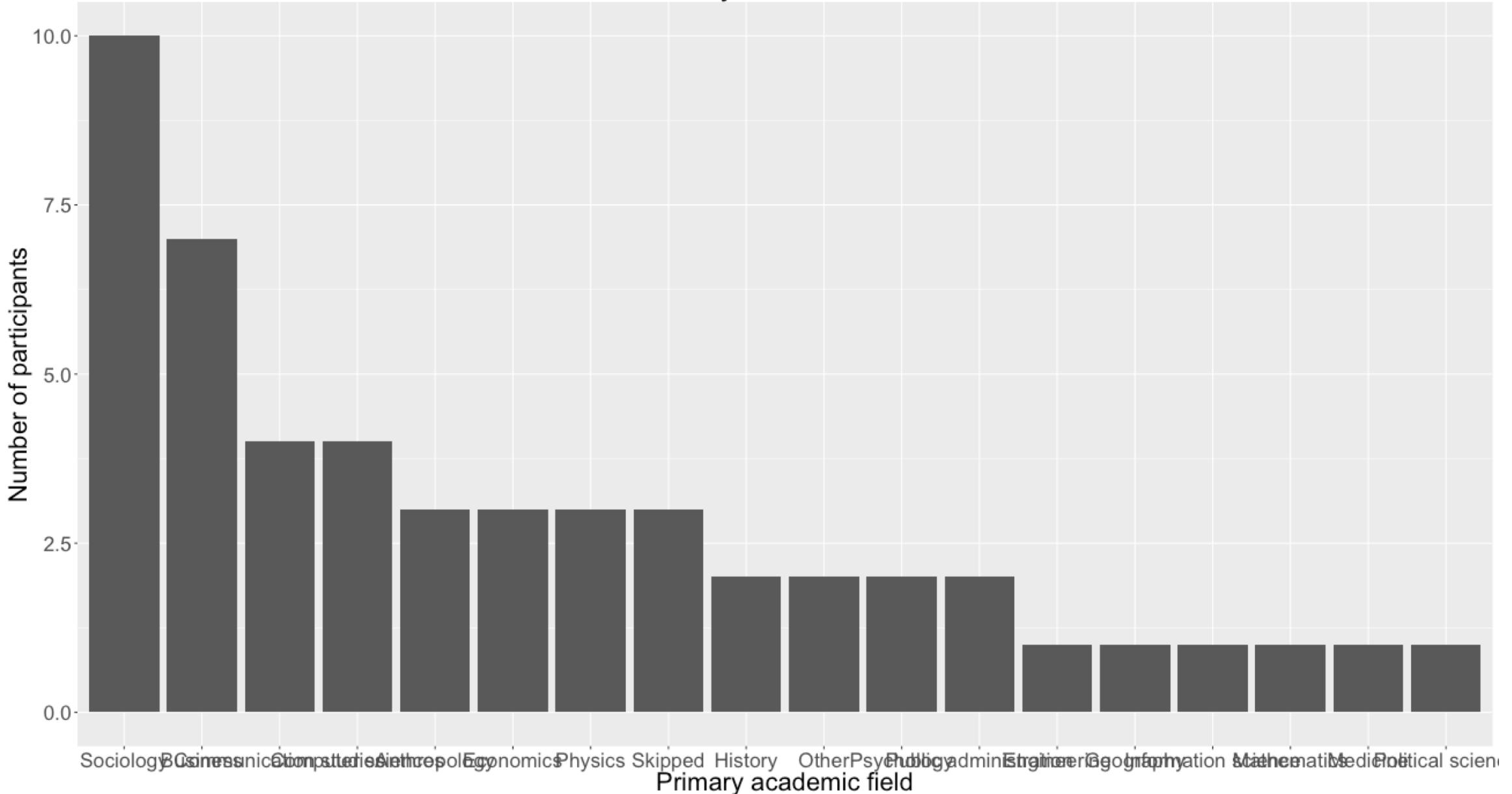
Primary academic field



Order by value

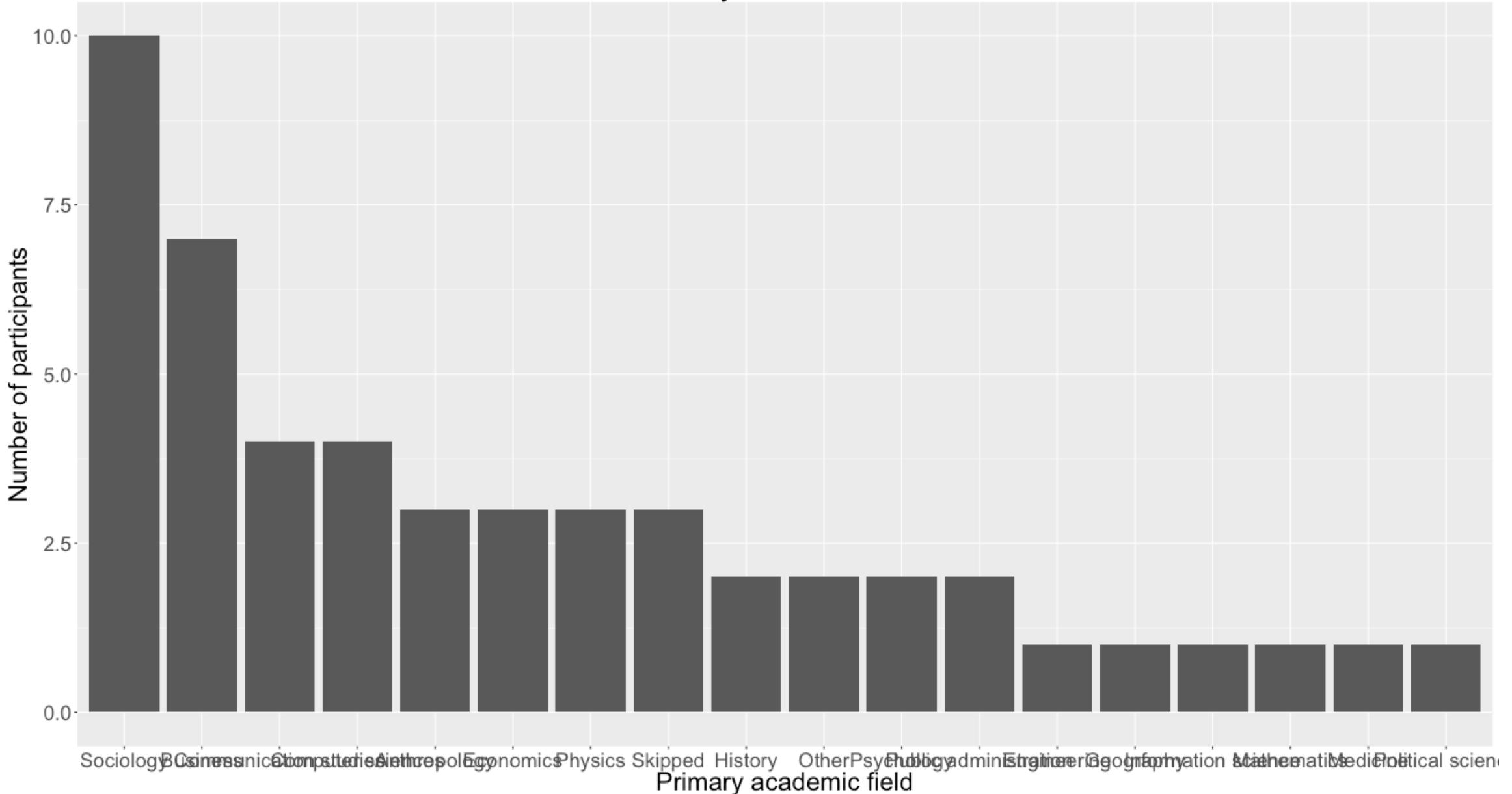
```
data$academic_field <-  
  factor(data$academic_field,  
         levels=names(  
           sort(  
             table(  
               data$academic_field),decreasing=TRUE))))
```

Primary academic field



Principle 2:
Put long categories on y-axis

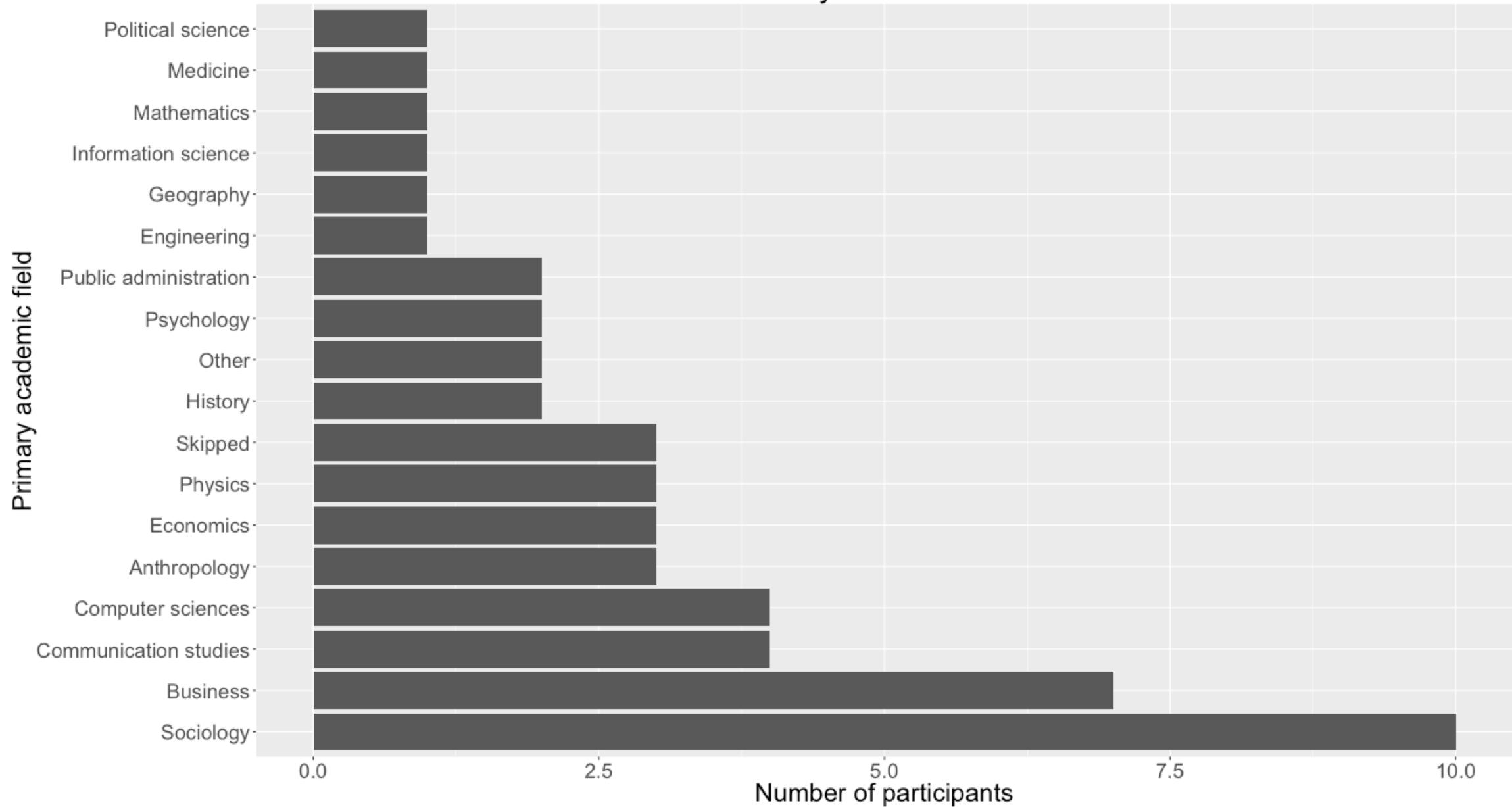
Primary academic field



Flip the axes

```
coord_flip()
```

Primary academic field

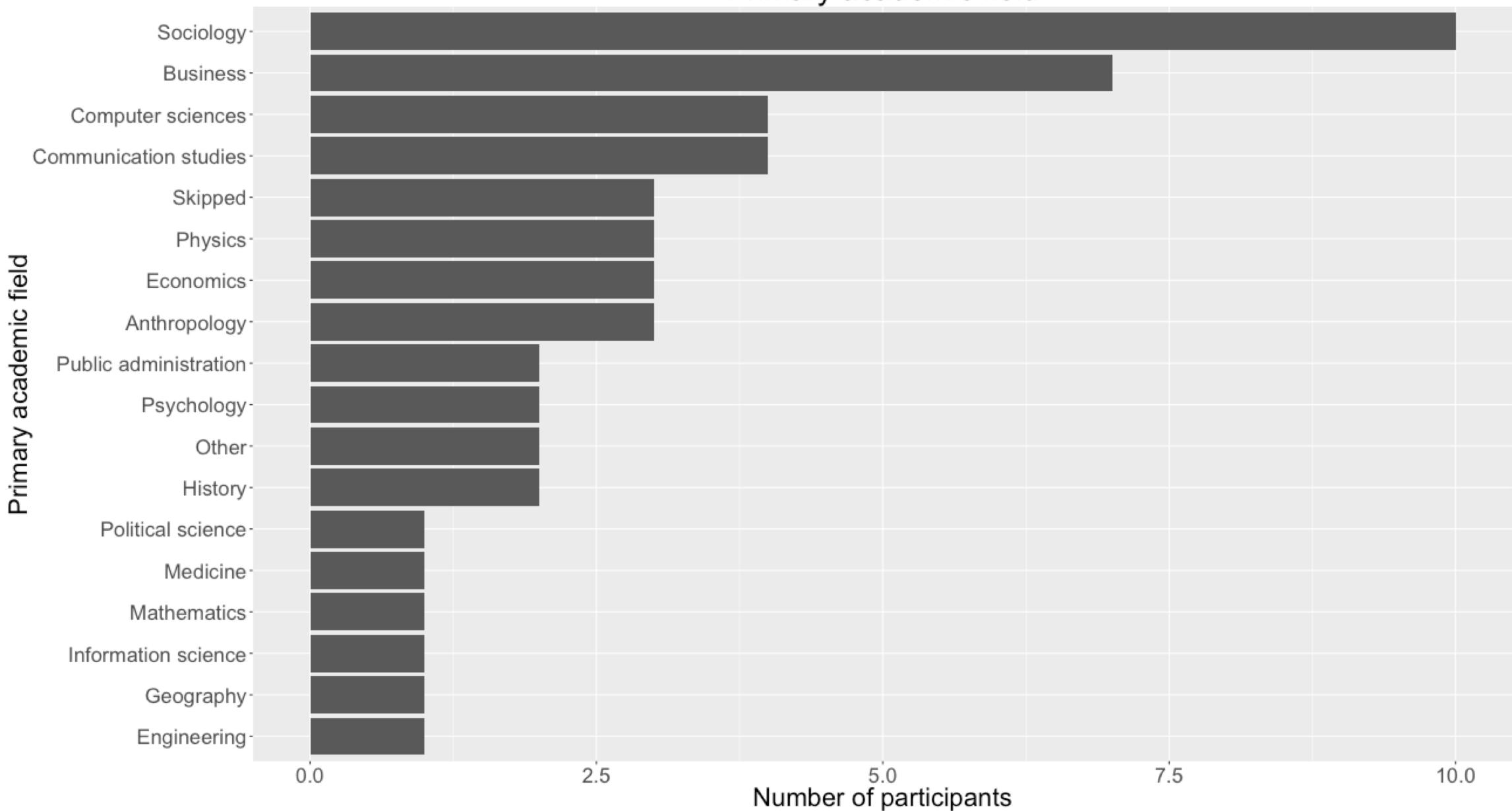


Oops!

```
data$academic_field <-  
  factor(data$academic_field,  
         levels=names(  
             sort(  
                 table(data$academic_field),  
                 decreasing=TRUE))))
```

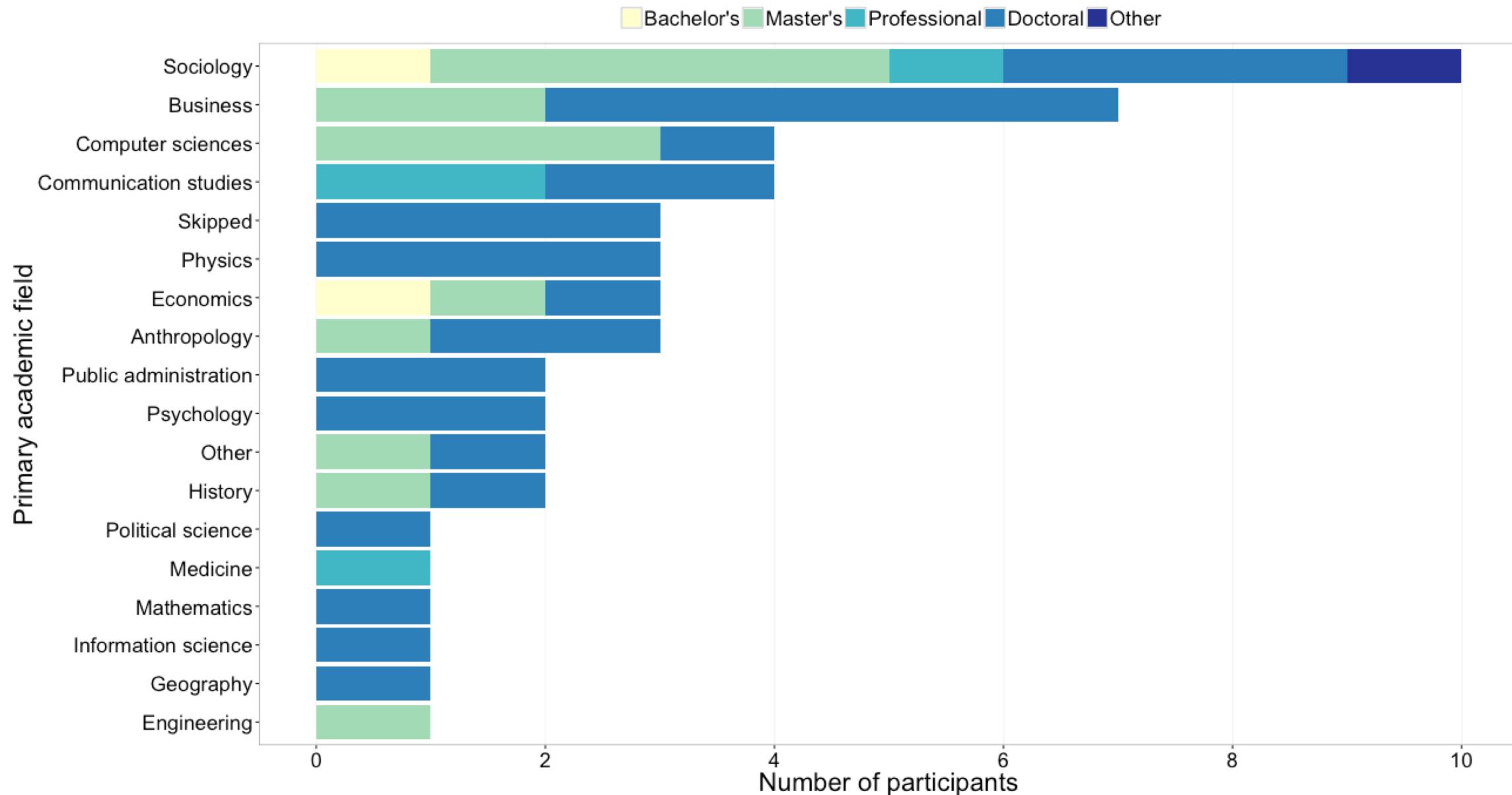
```
data$academic_field <-  
  factor(data$academic_field,  
         levels=names(  
             sort(  
                 table(data$academic_field)))))
```

Primary academic field

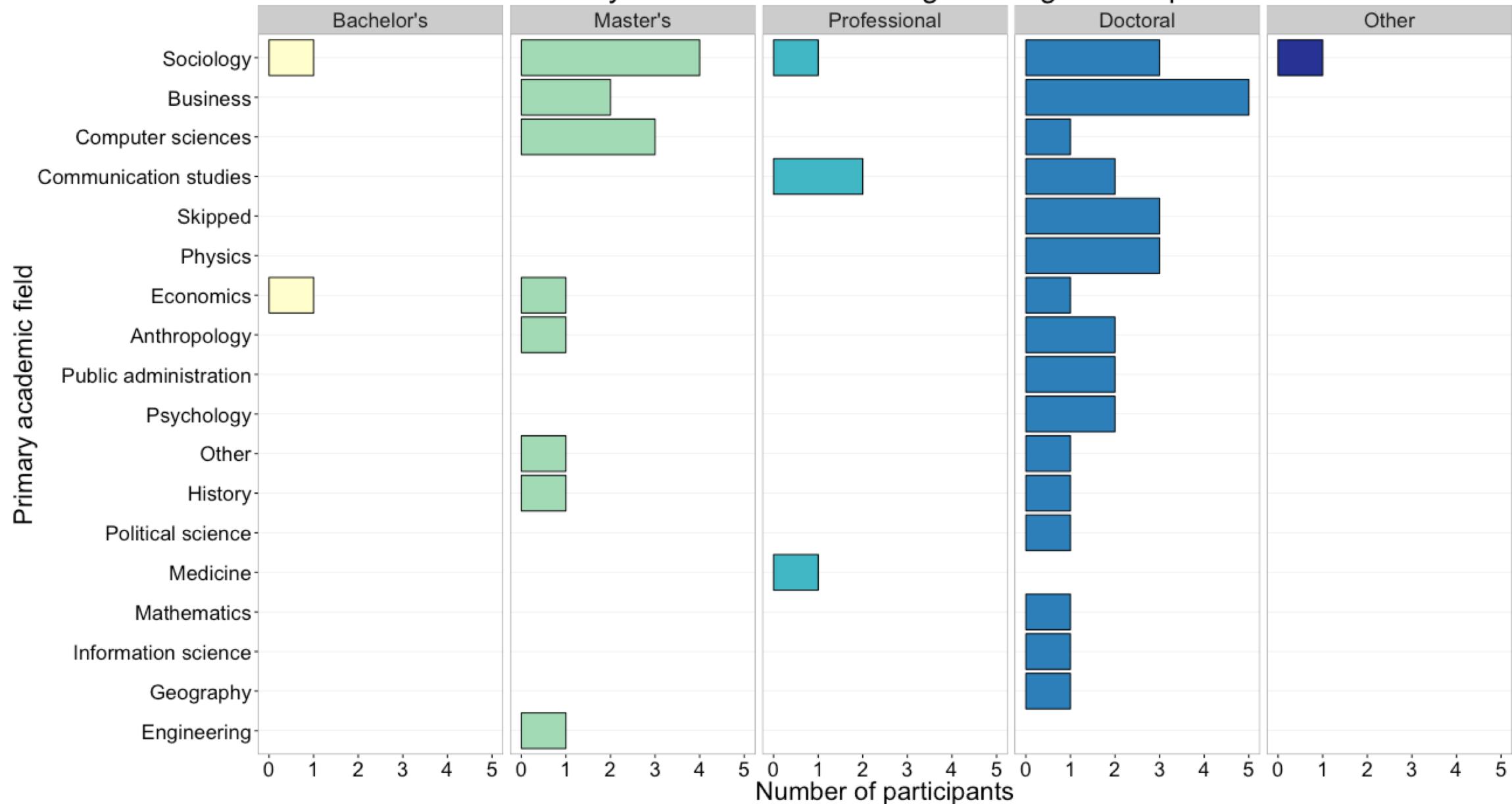


Principle 3: Pick a purpose

Primary academic field and highest degree completed



Primary academic field and highest degree completed



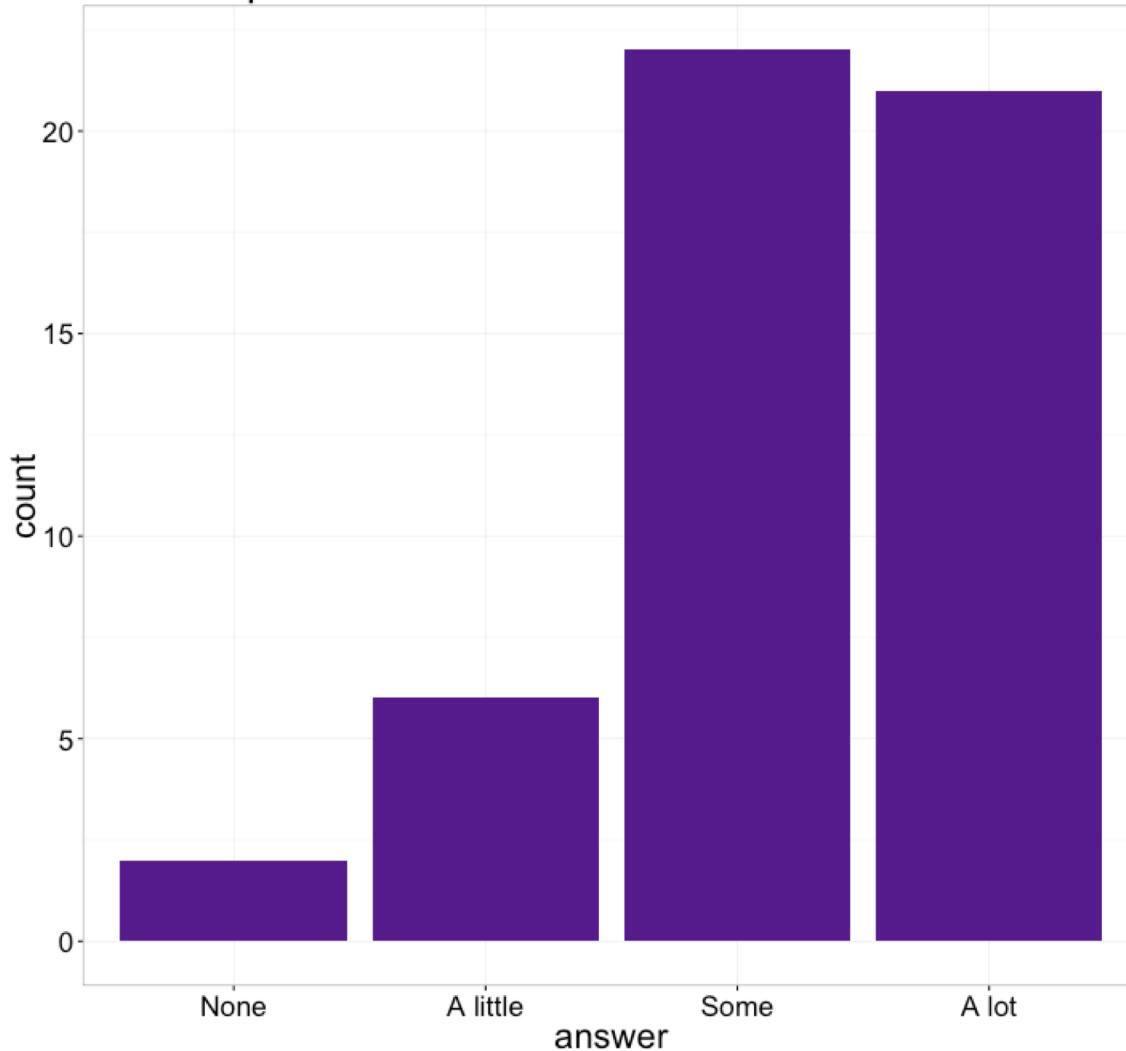
Different placement helps with different comparisons

```
fill=highest_degree
```

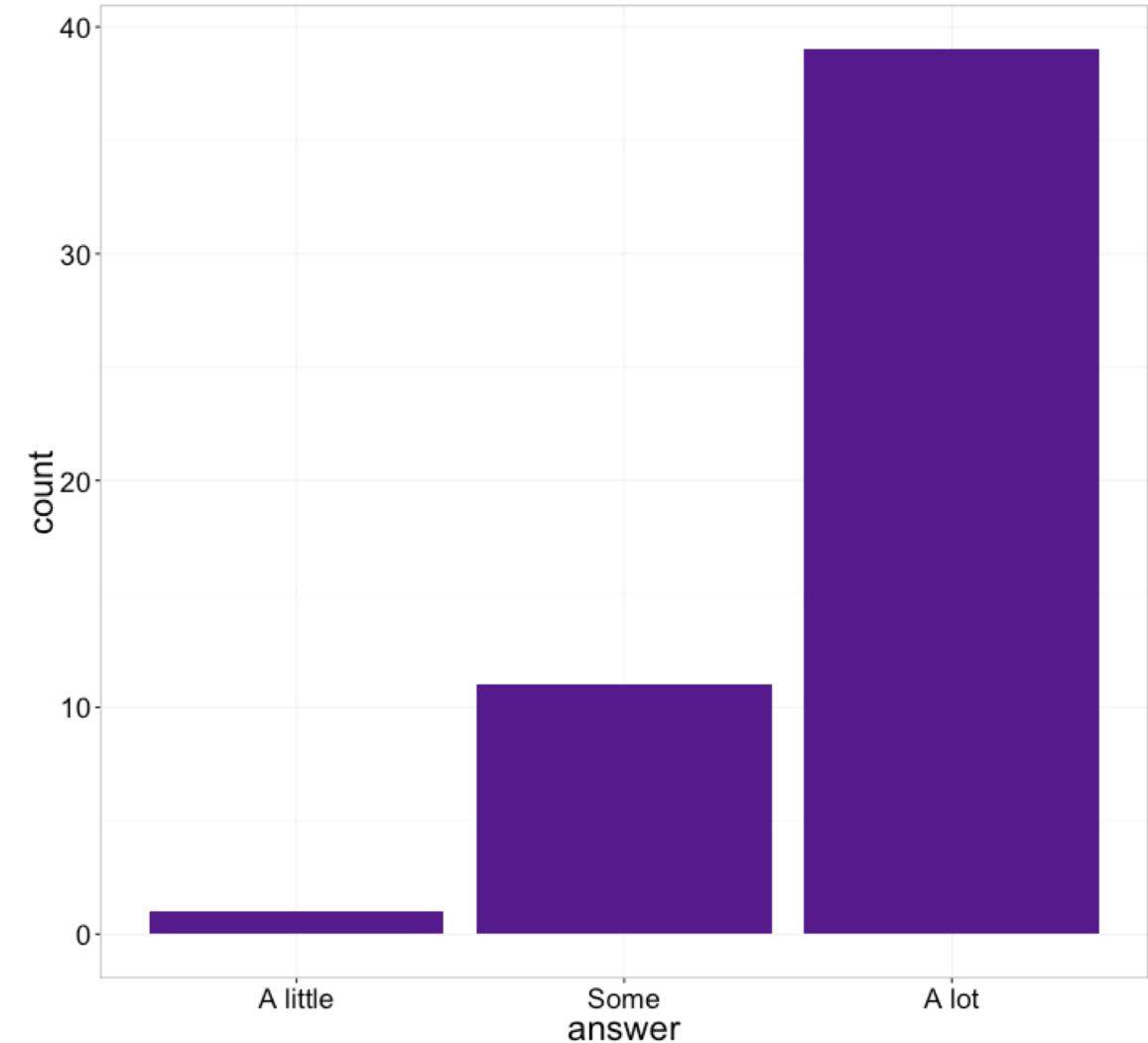
```
facet_grid(.~highest_degree)
```

Principle 4:
Keep scales consistent

How much experience do you have as a producer of network science research?



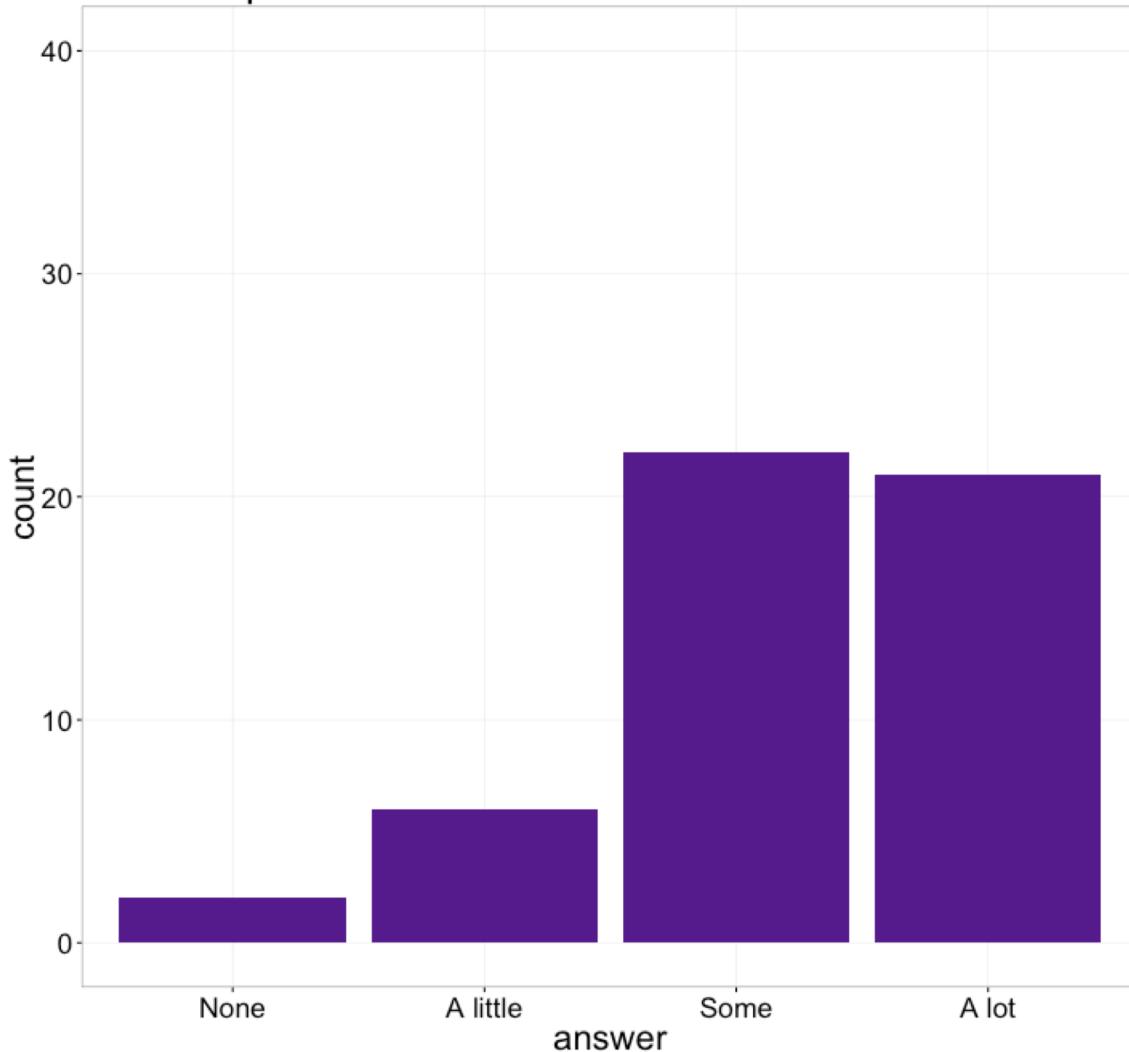
How much experience do you have as a consumer of network science research?



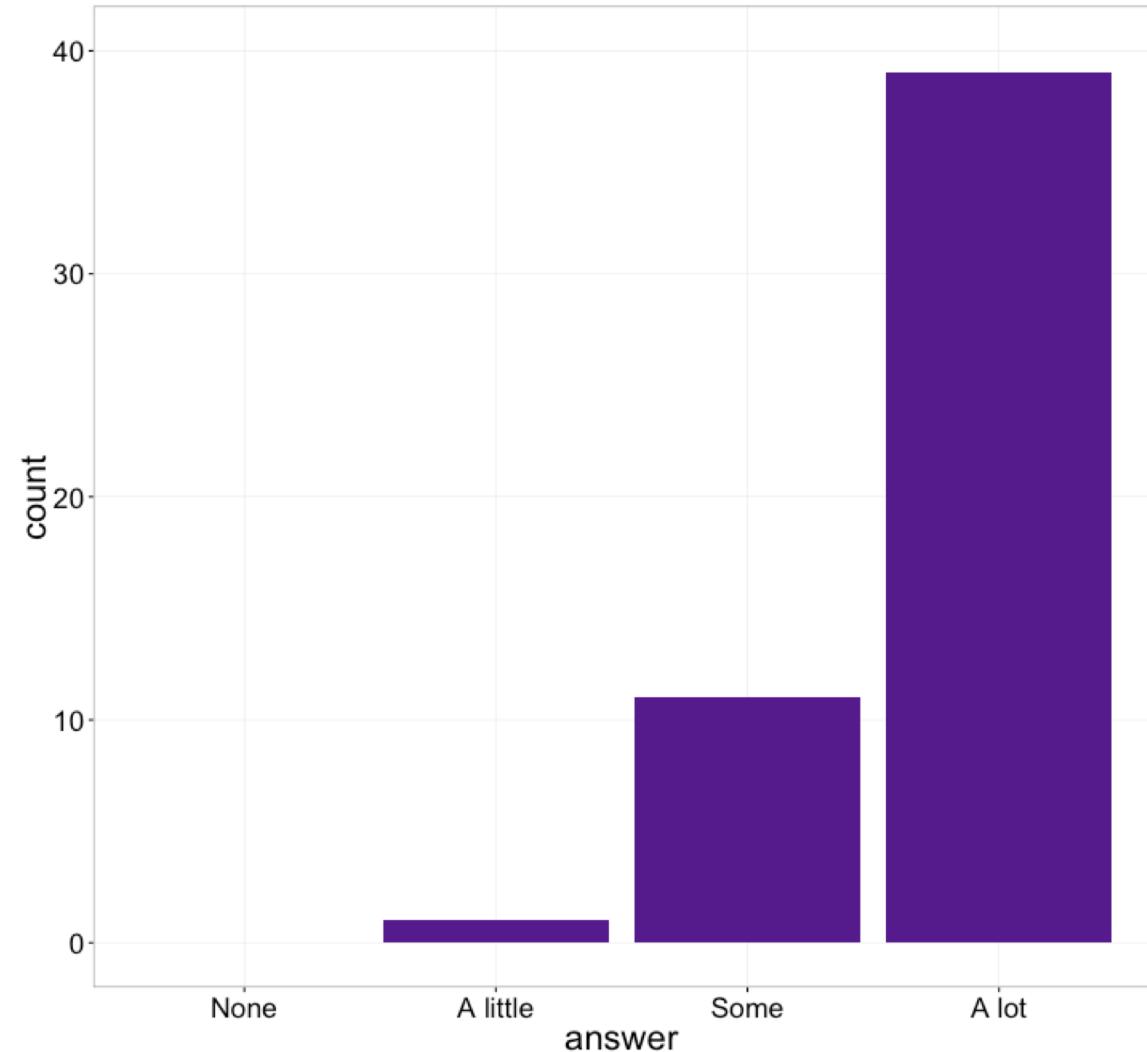
Keep all categories, manually set axes

```
scale_x_discrete(drop=FALSE)
scale_y_continuous(limits=c(0,40),
                   breaks=c(0,10,20,30,40),
                   minor_breaks=NULL)
```

How much experience do you have as a producer of network science research?

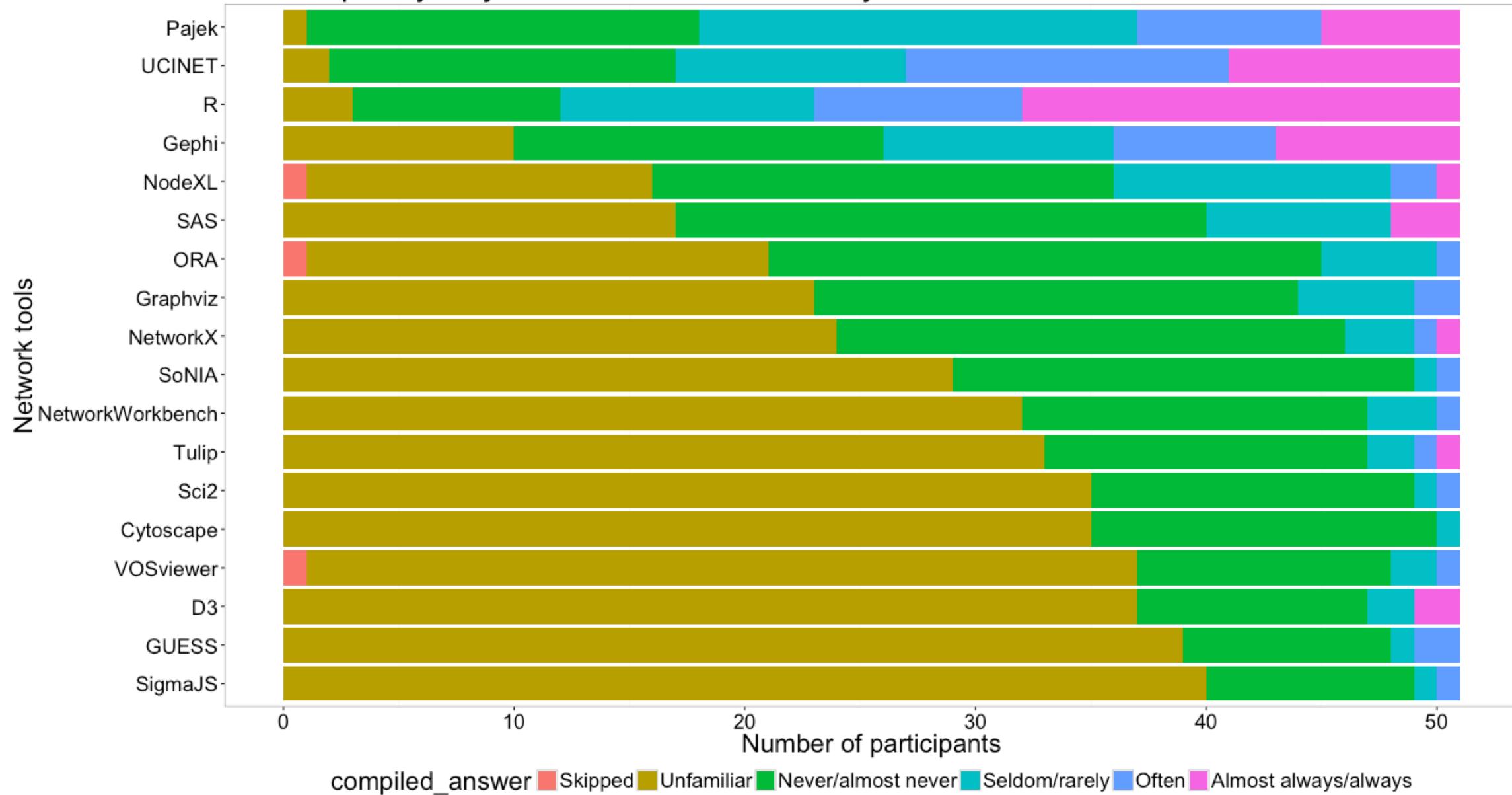


How much experience do you have as a consumer of network science research?



Principle 5:
Select meaningful colors

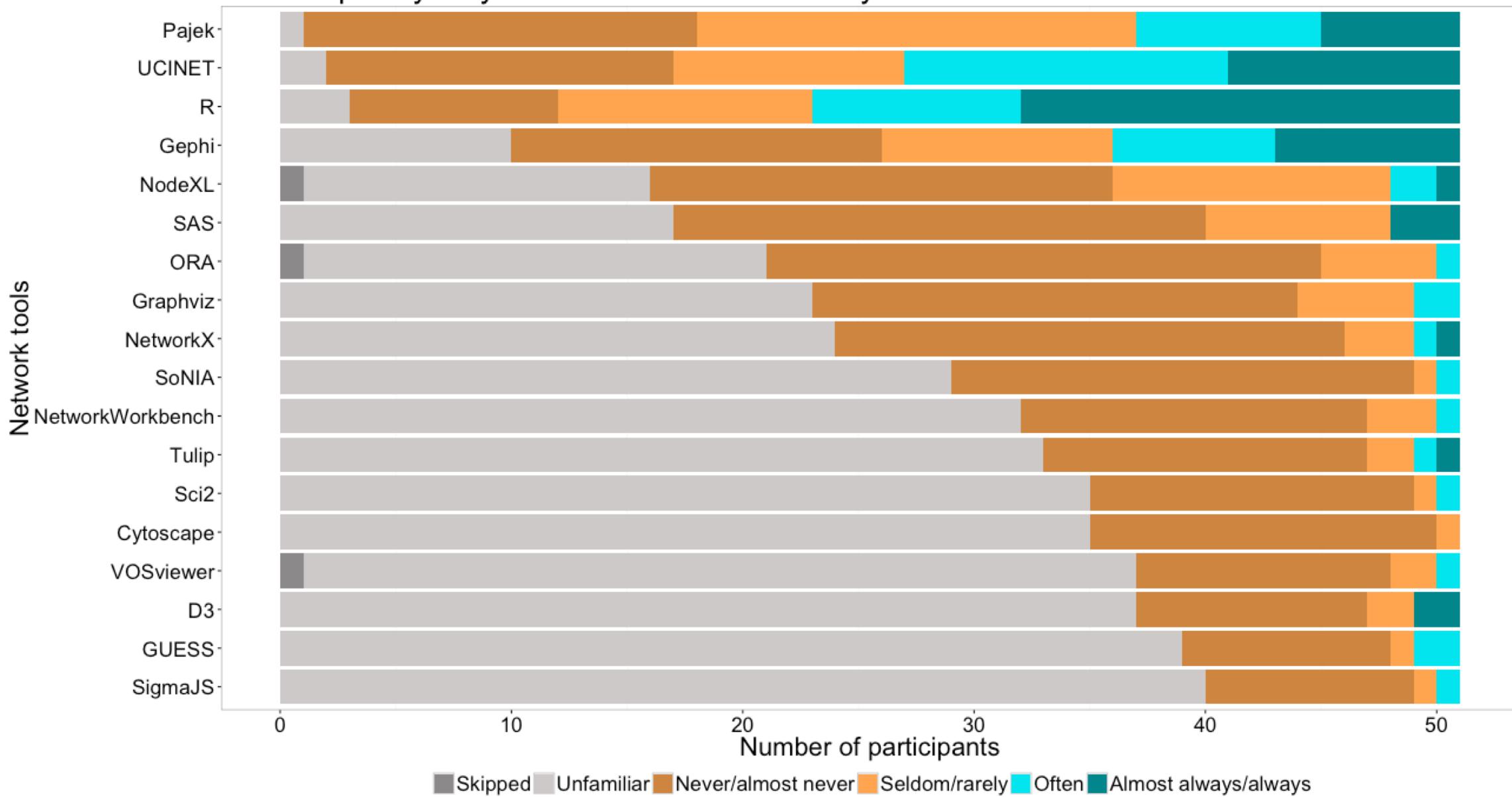
How frequently do you use these tools for analysis?



Select colors manually, or use alternate palette

```
scale_fill_manual(  
  values=c("snow4", "snow3",  
          "tan3", "tan1",  
          "turquoise2", "turquoise4"))  
  
scale_fill_manual(  
  values=c("#fee391", "#fe9929", "#cc4c02"))  
  
# Also see package RColorBrewer  
scale_fill_brewer(palette="BrBG")
```

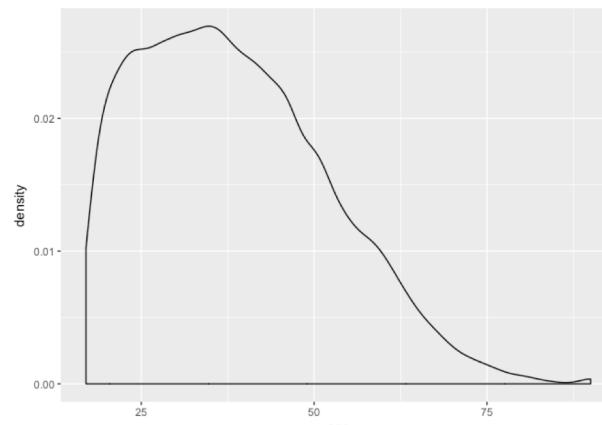
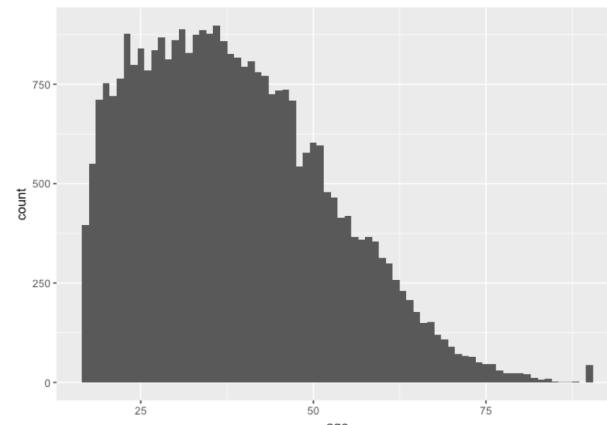
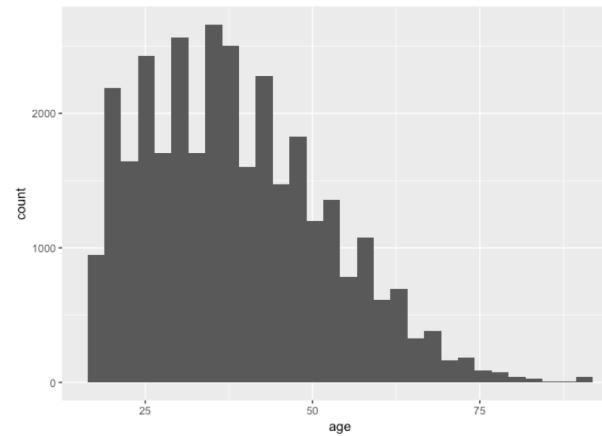
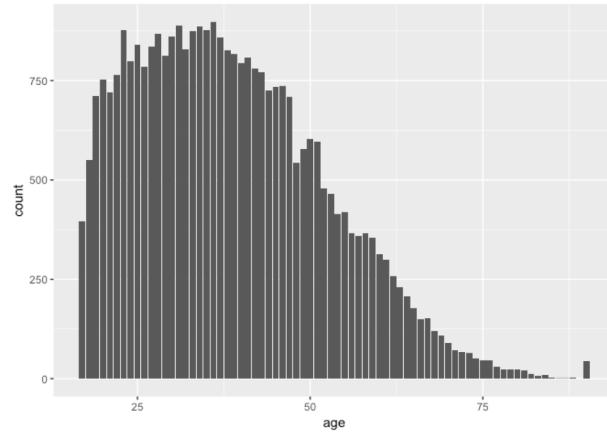
How frequently do you use these tools for analysis?



Exercise 2

Time to Statham Punch
<http://time.to.stathampun.ch/>

Templates



<https://github.com/amzoss/ggplot2-DF18/blob/master/5-templates.md>

Exercise 3

Gapminder Data

<http://www.gapminder.org/>

Averages across all years of the traditional Gapminder dataset

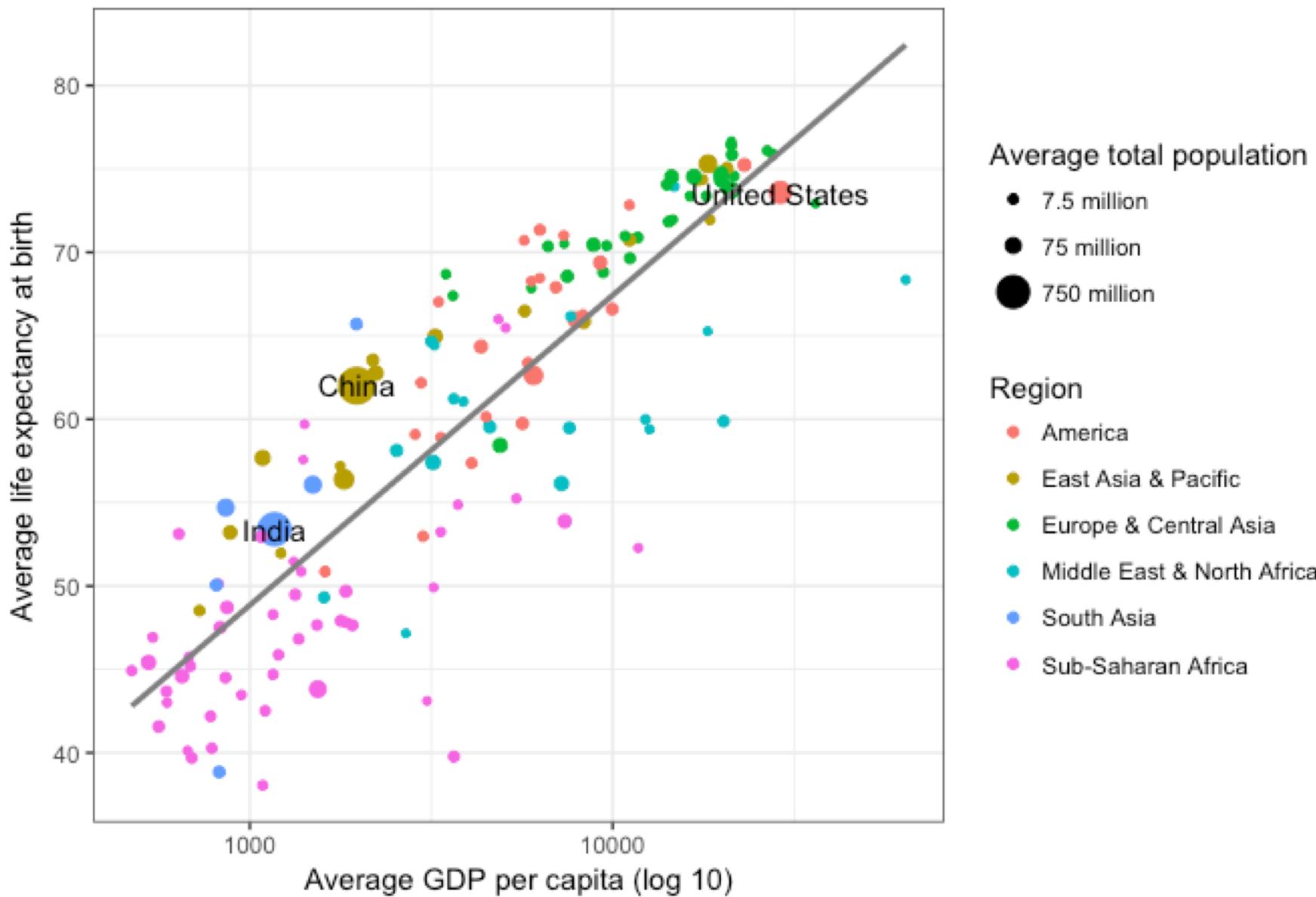


Chart choosing

data/2014Corollas.csv

Advanced

- Faceting
- Theme customization
- Data manipulations with dplyr, tidyr
- Custom stats functions

Getting help

ggplot2 Resources

- General ggplot2 information
<http://ggplot2.tidyverse.org/>
- R Graphics Cookbook (recipes for plots)
<http://www.cookbook-r.com/Graphs/index.html>
- R for Data Science (online book that includes ggplot2)
<http://r4ds.had.co.nz/>
- ggplot2: Elegant Graphs for Data Analysis (book by Hadley Wickham)
<http://ggplot2.org/book/>
- ggplot2 cheatsheet (also in RStudio)
<http://bit.ly/ggplot2-cheatsheet>

Data and Visualization Services



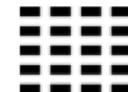
**Data and Visualization
Services Department**

<http://library.duke.edu/data>
askdata@duke.edu

Information about DVS

- Data collections, LibGuides, etc.
<http://library.duke.edu/data/>
- Blog (tutorials, announcements, etc.)
<http://blogs.library.duke.edu/data/>
- E-mail consultations
askdata@duke.edu
- Mailing list for announcements:
<https://lists.duke.edu/sympa/subscribe/dvs-announce>
- Twitter accounts
@duke_data, @duke_vis

Support Areas



Data Sources



Data Management



Data Cleaning



Data Analysis



Mapping and GIS



Data Visualization

Videos of past workshops

The image shows a screenshot of a Panopto video player interface. At the top left is the Panopto logo and the title "Figures and Posters". To the right are links for "Help" and "Sign in". Below the title, there is a thumbnail image of two people standing in front of a whiteboard in a classroom setting. The whiteboard has some handwritten text on it. To the right of the thumbnail, the main video area displays the title "Designing Academic Figures and Posters" in large bold letters, followed by the date "March 4, 2016" and a link "Slides: <http://duke.box.com/PostersSpring2016>". Below the title, two speakers are introduced: "Angela Zoss" and "Eric Monson", both described as Data Visualization Coordinators or Analysts from Data and Visualization Services. The video player includes a control bar at the bottom with a play button, a progress bar showing 0:03, a 10-second mark, and a 1:22:45 end time, along with buttons for "Speed" (set to 1X), "Quality", and "Hide". At the very bottom, there are four small thumbnail images related to poster design, each with a timestamp: 1:32, 4:32, 7:32, and 10:32.

<http://bit.ly/DVSvideos>

Questions?

askdata@duke.edu