# Visualization in R using ggplot2

Angela Zoss

9/18/19

https://github.com/amzoss/ggplot2-F18

# Set up environment

- R
- RStudio
- tidyverse package

# Get workshop files
URL: https://github.com/amzoss/ggplot2-F18

**With Git installed**

In RStudio:

- Project → New project

- Version Control

- Git
  - Paste in GitHub URL
  - Project directory name: ggplot2-F18
  - Subdirectory: you choose

- Create Project

**Without Git installed**

- Click green button to download ZIP

- Unzip files on your laptop

In RStudio:

- Project → New project…

- Existing directory

- Select unzipped folder

- Create Project

# 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

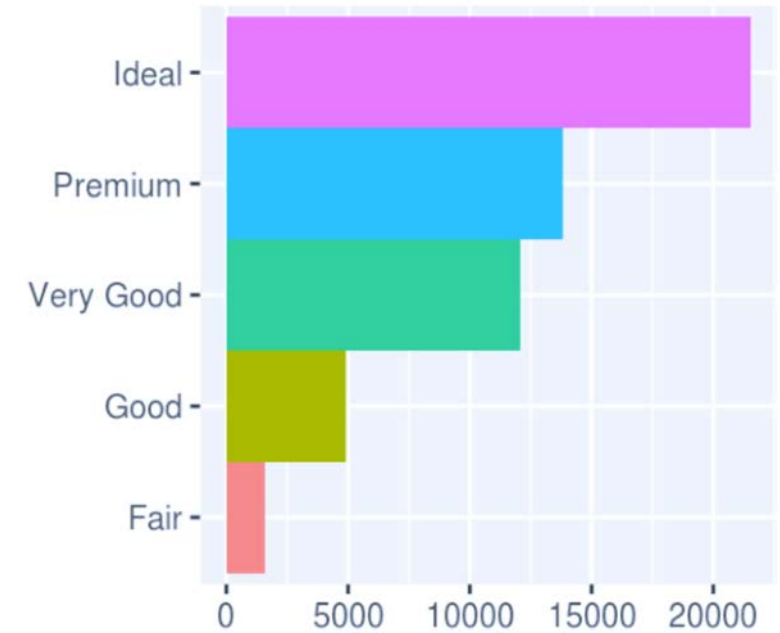- Saves you a lot of time trying to figure out what you did last time!

*"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.

# Grammar of graphics

1. DATA: a set of data operations that create variables from datasets
2. TRANS: variable transformations (e.g., rank)
3. SCALE: scale transformations (e.g., log)
4. COORD: a coordinate system (e.g., polar)
5. ELEMENT: graphs (e.g., points) and their aesthetic attributes (e.g., color)
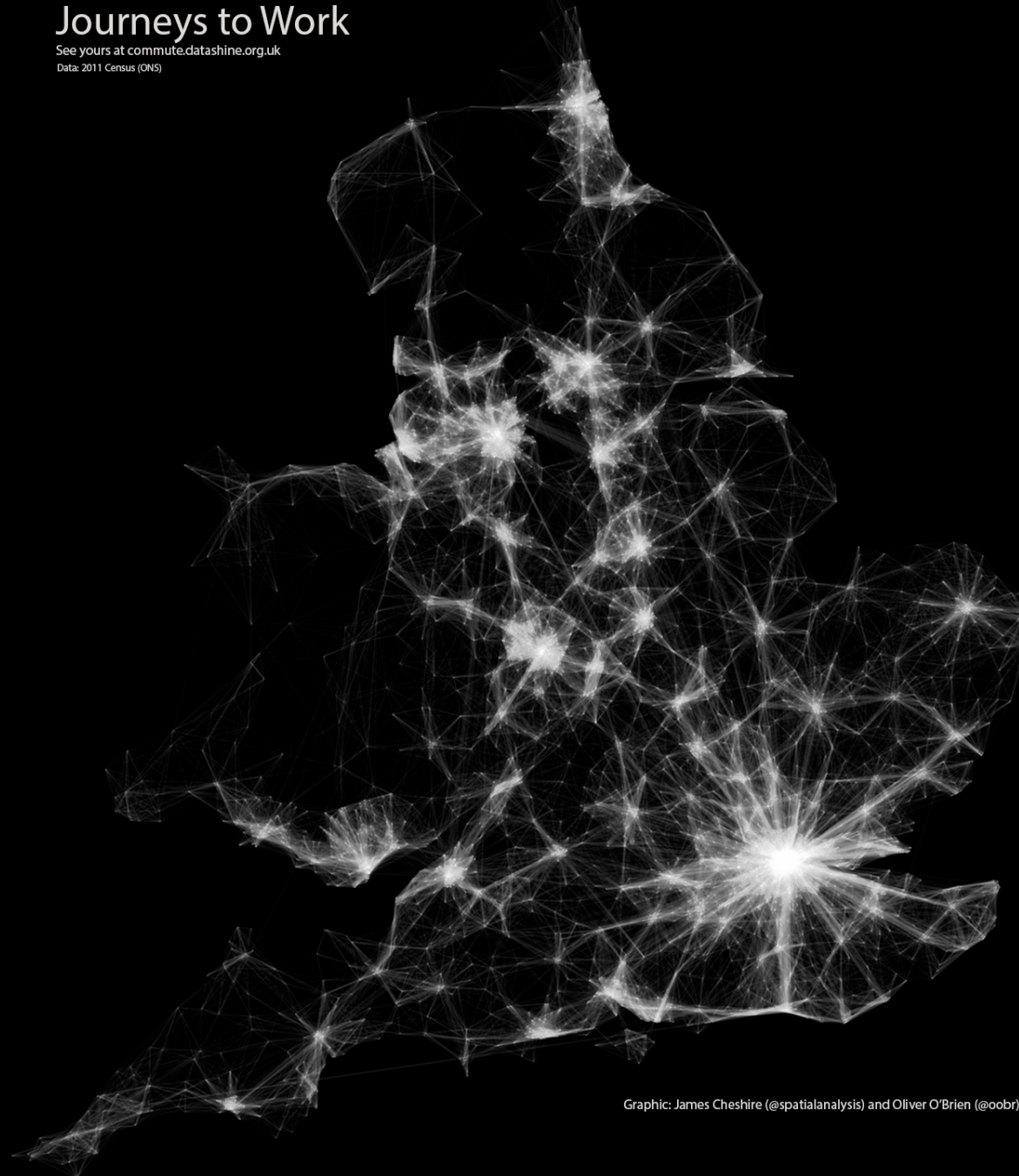6. GUIDE: one or more guides (axes, legends, etc.).

Wilkinson, Leland. (2005). *The grammar of graphics (2nd ed)*. New York: Springer.

# ggplot2 examples

# GDP per capita on Five Continents



http://socviz.co/groupfacettx.html

Journeys to Work

See yours at commute.datashine.org.uk

Data: 2011 Census (ONS)

Graphic: James Cheshire (@spatialanalysis) and Oliver O'Brien (@oobr)

http://spatial.ly/2015/03/mapping-flows/

http://sappingattention.blogspot.com/2017/05/a-brief-visual-history-of-marc.html

# 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/

# ggplot2: Elements

# Basic elements in any ggplot2 visualization

- **data**

- **aes**thetics
(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( data = data frame )
```
+

```
geom_...( 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. |

x position    y position    color

# 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. |

x position    y position    color

```
ggplot(data) +
    geom_point(
        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. |

x position    y position    color

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

# Template for a more complex plot

carry through
from top to bottom

```
ggplot( data = data frame,
        aes(variable mappings) )
```

+

```
geom_...( aes(add'l variable mappings),
          non-variable adjustments )
```

+

```
geom_...( aes(add'l variable mappings),
          non-variable adjustments )
```

# Using RStudio

- Projects
- Rmarkdown
- Cheat sheets

https://www.rstudio.com/resources/cheatsheets/#rmarkdown

# Why Rmarkdown?

- Plots show up inline
- Easier to incorporate explanatory text and materials
- Like to be able to easily run one chunk at a time

Caution: Running things out of order can mean your code won't work again later. Clear your environment often and run code chunks in order to be safe.

# Dataset 1: Game of Thrones character ratings

# Dataset 2:
# Star Wars character data

https://dplyr.tidyverse.org/reference/starwars.html

# 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

```r
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)))
```

```
data$academic_field <-
    fct_infreq(as_factor(data$academic_field))
```

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$academic_field <-
    fct_rev(fct_infreq(as_factor(data$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?

Legend: Skipped | Unfamiliar | Never/almost never | Seldom/rarely | Often | Almost always/always

# Dataset 3:
# Star Wars opinion survey

https://fivethirtyeight.com/features/americas-favorite-star-wars-movies-and-least-favorite-characters/

# Saving charts out

# Dataset 4: Gapminder Data

http://www.gapminder.org/

Averages across all years of the traditional Gapminder dataset

# ggplot2 Cheat Sheet

# 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

# Videos of past workshops



http://bit.ly/DVSvideos

# Questions?

[angela.zoss@duke.edu](mailto:angela.zoss@duke.edu)