





























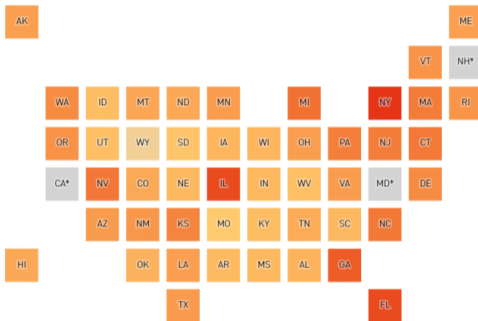




# Brady's Example, Comments by Evan

## New York, Illinois had highest rates of reported abortion in 2019

Number of reported abortions per 1,000 women ages 15-44 in hospitals and ambulatory care facilities, including clinics. The true number of abortions is likely larger, as not all abortions are reported.



\*Data not available for California, Maryland and New Hampshire  
Source: CDC  
Map: Sean McMinn / POLITICO

Goldberg, Dan. "Abortion Statistics by state: Maps, trigger laws, and possible bans," *Politico*, May 3, 2022.



















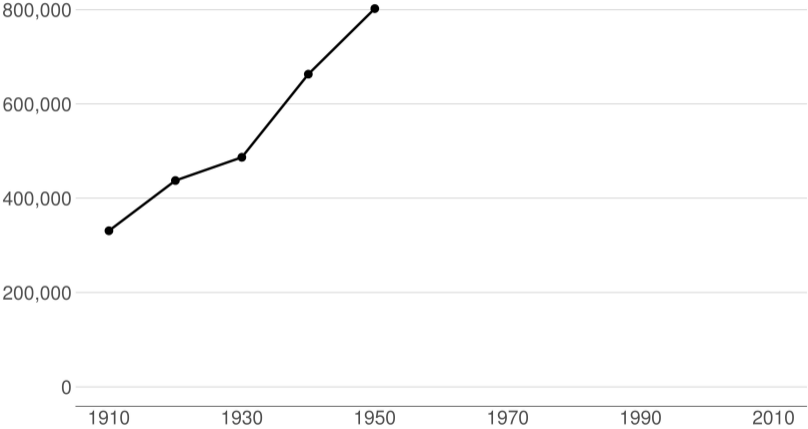




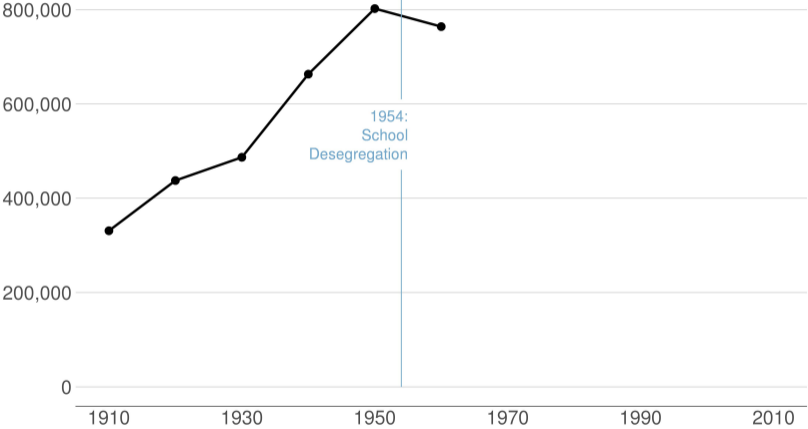




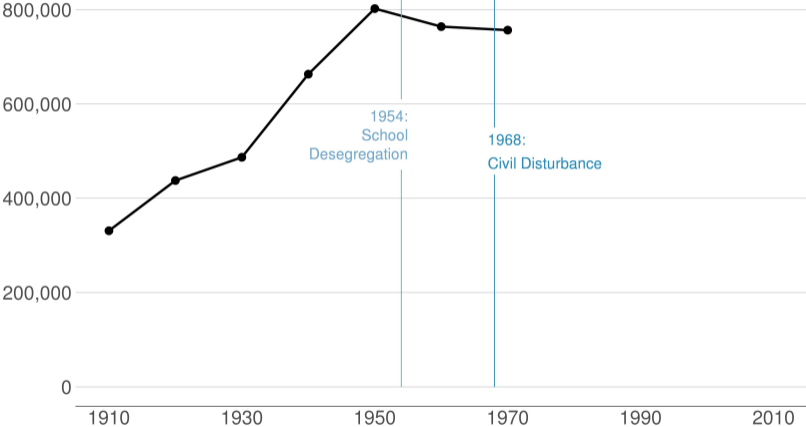
# DC Gains Population Through 1950



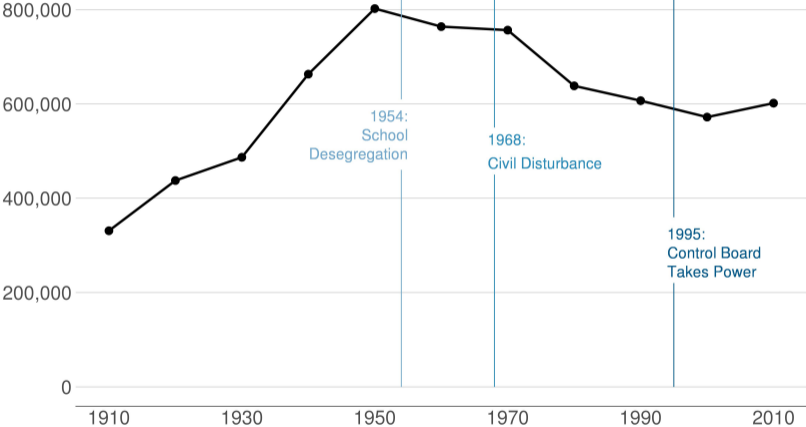
# Population Loses Start with Desegregation



# Continue After Civil Disturbance



# Population Turns Up After 2000



















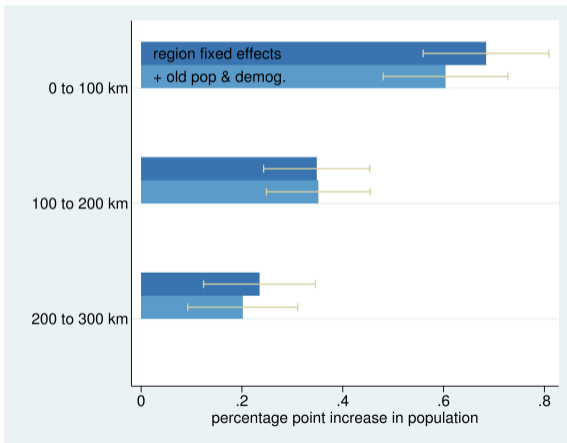




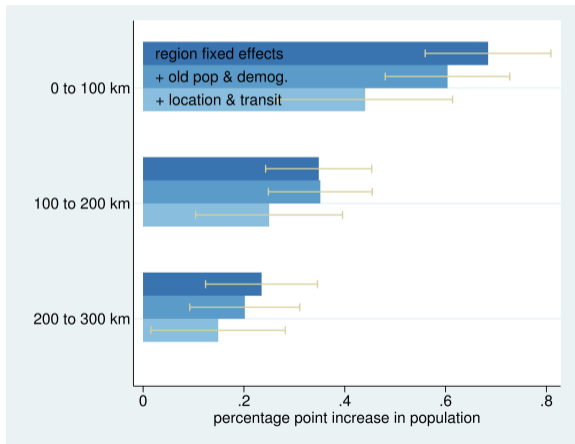




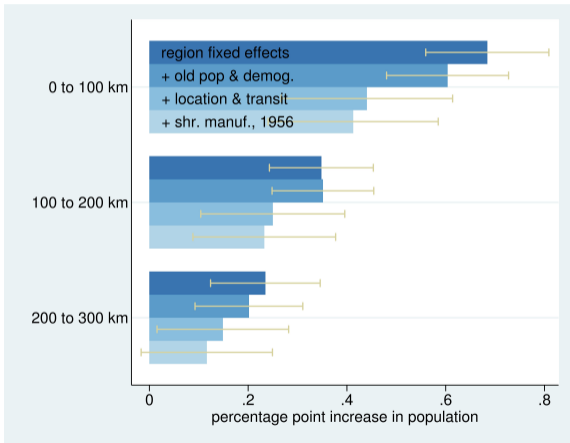
# Bars with Error Bars, Building



# Bars with Error Bars, Building



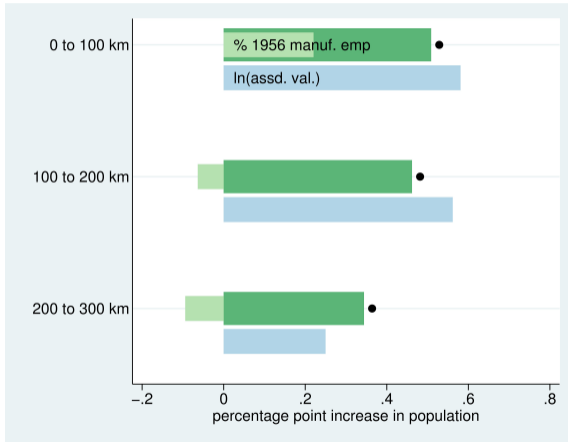
## Bars with Error Bars, Building







# Interaction Effects

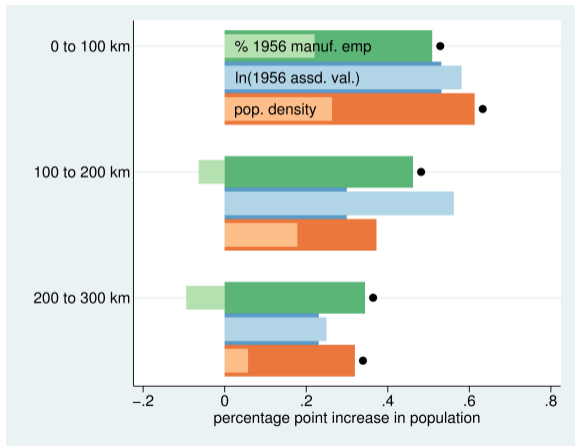








# Interaction Effects









# Today in R: Line Charts and De-Debugging

1. Line charts and `ggplot`
2. Summarizing data
3. Annotations
4. Making data long
5. De-debugging

## 1. Line charts

```
p1 <- ggplot() +  
  geom_line(data = polys,  
            mapping = aes(x = xvar, y = yvar))
```

## 1. Line charts

```
p1 <- ggplot() +  
  geom_line(data = polys,  
            mapping = aes(x = xvar, y = yvar))
```

- ▶ R does not require `xvar` to be time
- ▶ But your readers will assume it is



## Hurricanes by Year

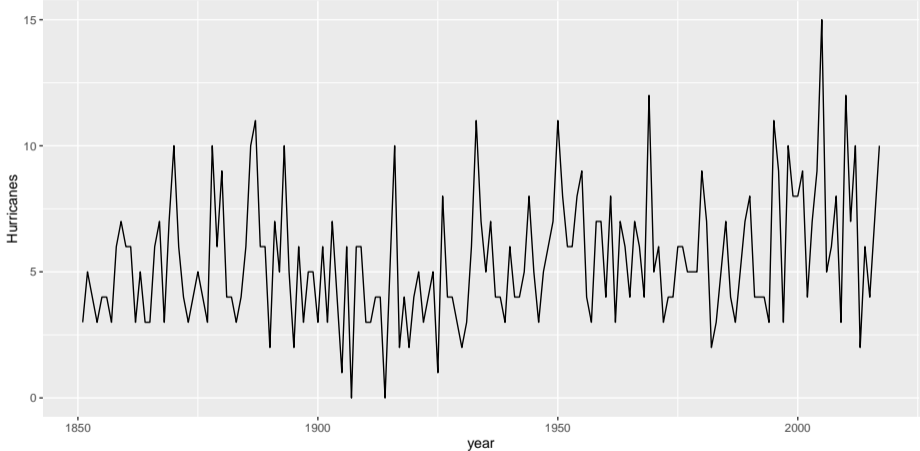
```
# remember the hurricanes  
hurr <- read.csv("H:/pppa_data_viz/2019/tutorial_data/lecture03/2019-02-02.  
# look at it  
head(hurr)
```

| ##   | year | Named.Storms | Hurricanes | Major |
|------|------|--------------|------------|-------|
| ## 1 | 1851 | 6            | 3          | 1     |
| ## 2 | 1852 | 5            | 5          | 1     |
| ## 3 | 1853 | 8            | 4          | 2     |
| ## 4 | 1854 | 5            | 3          | 1     |
| ## 5 | 1855 | 5            | 4          | 1     |
| ## 6 | 1856 | 6            | 4          | 2     |

## Graphing Hurricanes by Year

```
# plot it  
hurr.by.year <- ggplot() +  
  geom_line(data = hurr,  
            mapping = aes(x = year, y = Hurricanes))
```

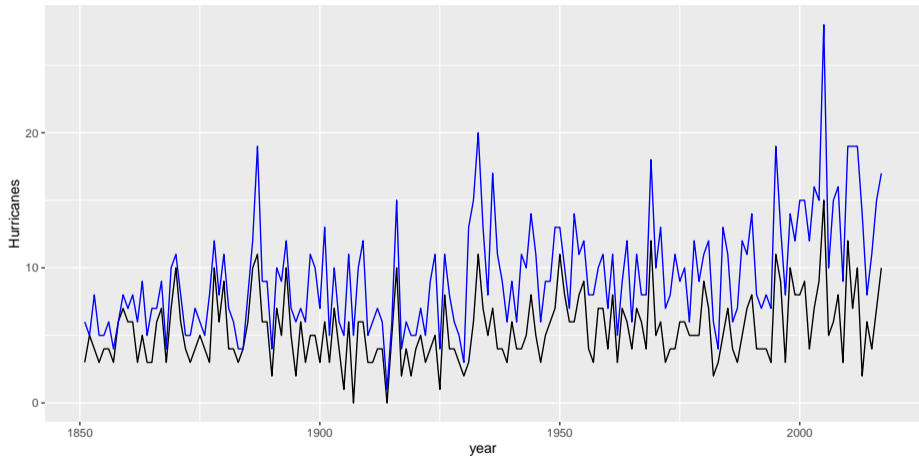
# Graphing Hurricanes by Year



## Graphing Hurricanes and Named Storms by Year

```
# plot it
hurr2.by.year <- ggplot() +
  geom_line(data = hurr,
            mapping = aes(x = year, y = Hurricanes)) +
  geom_line(data = hurr,
            mapping = aes(x = year, y = Named.Storms),
            color = "blue")
```

## Graphing Hurricanes and Named Storms by Year



## Using Group= to Make Multiple Lines

You need long data:

```
df <- data.frame(year = c(1, 2, 3, 1, 2, 3, 1, 2, 3),  
                 type = c(1, 1, 1, 2, 2, 2, 3, 3, 3),  
                 outcome = c(1, 2, 3, 2, 3, 3.1, 3, 4, 5))
```

df

```
##   year type outcome  
## 1     1     1     1.0  
## 2     2     1     2.0  
## 3     3     1     3.0  
## 4     1     2     2.0  
## 5     2     2     3.0  
## 6     3     2     3.1  
## 7     1     3     3.0  
## 8     2     3     4.0  
## 9     3     3     5.0
```

## Using Group= to Make Multiple Lines

You need long data:

```
# plot it  
df.by.year <- ggplot() +  
  geom_line(data = df,  
            mapping = aes(x = year, y = outcome,  
                          group = type))
```

## Using Group= to Make Multiple Lines

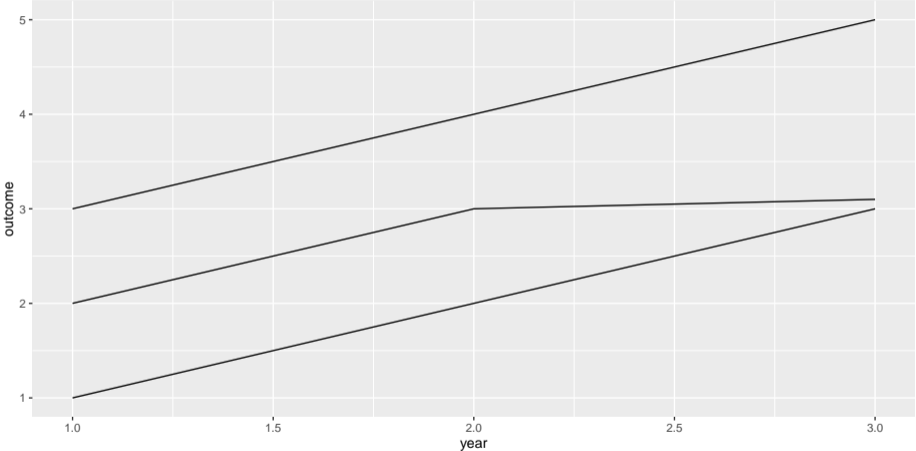
You need long data:

```
# plot it  
df.by.year <- ggplot() +  
  geom_line(data = df,  
            mapping = aes(x = year, y = outcome,  
                           group = type))
```

- ▶ Be wary of using too many lines



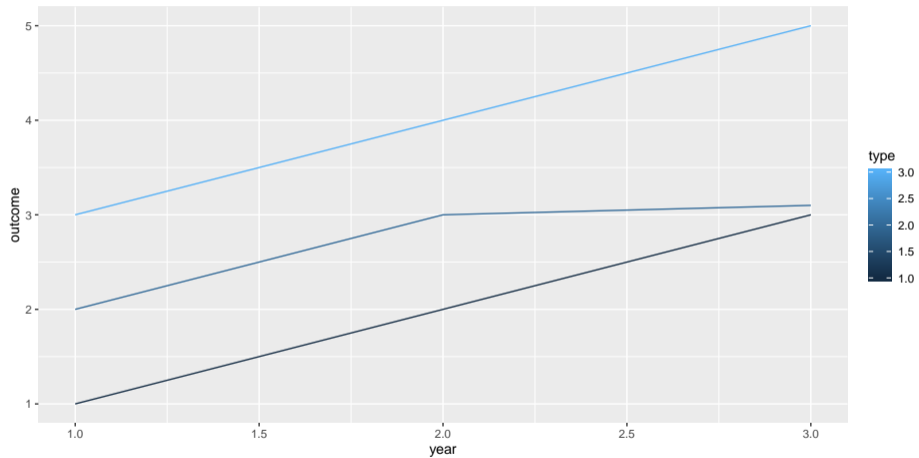
# Using Group= to Make Multiple Lines



## Using Group= to Make Multiple Lines and color = to color lines

```
df.by.year2 <- ggplot() +  
  geom_line(data = df,  
            mapping = aes(x = year, y = outcome,  
                          group = type, color = type))
```

## Using Group= to Make Multiple Lines and color = to color lines



## 2. Summarizing data

In today's tutorial, you'll use bikeshare data

- ▶ these come at the level of the individual ride
- ▶ we describe them by hour
  - ▶ → summarize by hour (`group_by first`)
- ▶ for the homework, you describe them by minute

## 2. Summarizing data

In today's tutorial, you'll use bikeshare data

- ▶ these come at the level of the individual ride
- ▶ we describe them by hour
  - ▶ → summarize by hour (`group_by` first)
- ▶ for the homework, you describe them by minute

Remember that `group_by` and then `summarize` take you from one unit of observation to another.

### 3. Annotations: Better than Legends

General logic is

```
np <- already.existing.plot +  
  annotate(geom = [name of geom],  
         x = [x location],  
         y = [y location],  
         [xmin = , xmax = ,  
          ymin = , ymax =]  
         other options -- size, color, etc)
```

### 3. Annotations: Better than Legends

General logic is

```
np <- already.existing.plot +  
  annotate(geom = [name of geom],  
         x = [x location],  
         y = [y location],  
         [xmin = , xmax = ,  
          ymin = , ymax =]  
         other options -- size, color, etc)
```

What you add is in `geom = [name of geom]` which can be

- ▶ text
- ▶ point
- ▶ segment (line)
- ▶ rectangle
- ▶ possibly others

### 3. Annotation example plan

- ▶ make a small dataframe to illustrate
- ▶ show the line graph of this small dataframe
- ▶ add an annotation



### 3. Small example dataframe

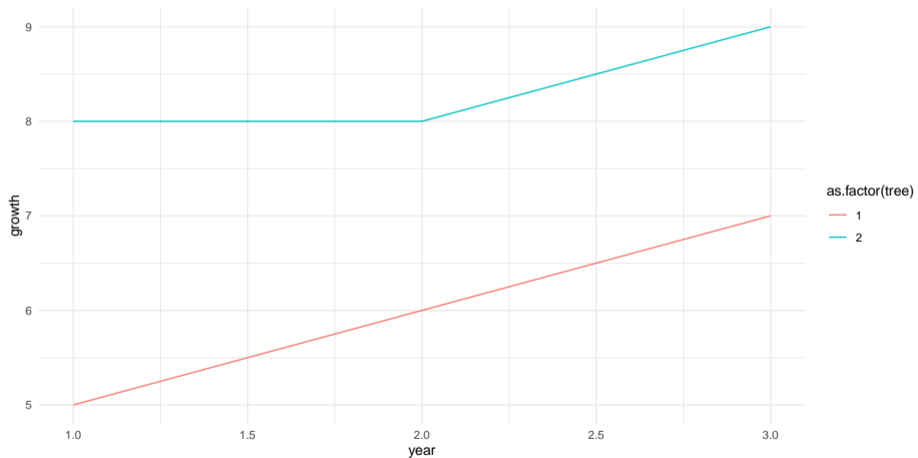
```
trees <- data.frame(year = c(1,2,3,1,2,3),  
                    tree = c(1,1,1,2,2,2),  
                    growth = c(5,6,7,8,8,9))
```

### 3. Line plot of trees

```
tp <- ggplot() +  
  geom_line(data = trees,  
            mapping = aes(x = year, y = growth,  
                          group = tree,  
                          color = as.factor(tree))) +  
  theme_minimal()
```

### 3. Line plot of trees

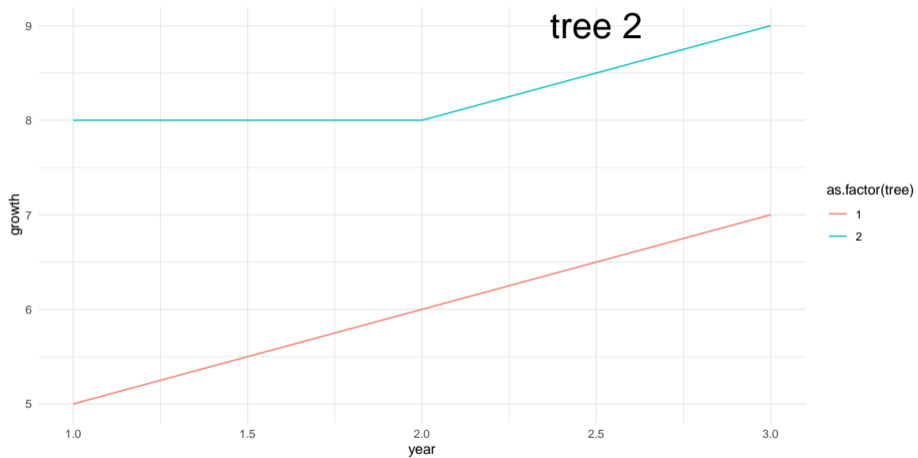
tp



### 3. Add text annotation

```
tp2 <- tp +  
  annotate(geom = "text",  
    x = 2.5,  
    y = 9,  
    label = "tree 2",  
    size = 10)
```

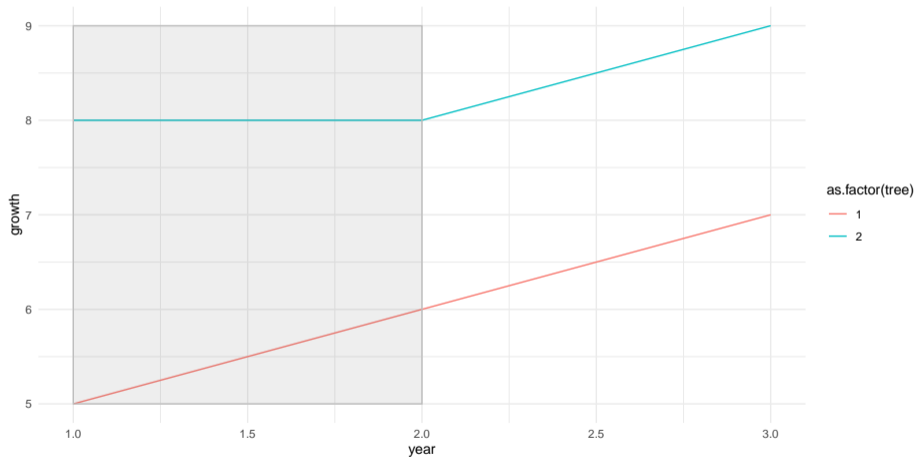
### 3. Line plot of trees with annotation



### 3. Add rectangle annotation

```
tp3 <- tp +  
  annotate(geom = "rect",  
          xmin = 1, xmax = 2,  
          ymin = 5, ymax = 9,  
          color = "grey",  
          alpha = 0.1)
```

### 3. Add rectangle annotation



## 4. Making Data Long

`ggplot()` prefers long data

To think about this we will

- ▶ show wide data
- ▶ show long data
- ▶ show how to go wide to long



## 4. Wide data

```
wide <- data.frame(state = c("6", "36", "48"),  
                  female_pop = c("10", "12", "14"),  
                  male_pop = c("11", "13", "12"))
```

wide

| ##   | state | female_pop | male_pop |
|------|-------|------------|----------|
| ## 1 | 6     | 10         | 11       |
| ## 2 | 36    | 12         | 13       |
| ## 3 | 48    | 14         | 12       |

## 4. Long data

```
long <- data.frame(state = c("6", "36", "48", "6", "36", "48"),  
                  pop = c("10", "12", "14", "11", "13", "12"),  
                  sex = c("female", "female", "female", "male", "male", "male"))
```

long

| ##   | state | pop | sex    |
|------|-------|-----|--------|
| ## 1 | 6     | 10  | female |
| ## 2 | 36    | 12  | female |
| ## 3 | 48    | 14  | female |
| ## 4 | 6     | 11  | male   |
| ## 5 | 36    | 13  | male   |
| ## 6 | 48    | 12  | male   |

## 4. Going from wide to long

```
long2 <- pivot_longer(data = wide,  
                      cols = c("female_pop", "male_pop"),  
                      names_to = "sex",  
                      values_to = "pop")
```

```
long2
```

```
## # A tibble: 6 x 3  
##   state sex      pop  
##   <fct> <chr>   <fct>  
## 1 6     female_pop 10  
## 2 6     male_pop   11  
## 3 36    female_pop 12  
## 4 36    male_pop   13  
## 5 48    female_pop 14  
## 6 48    male_pop   12
```

## 4. Additional notes

- ▶ you can clean up the sex variable with a `substr()` command
- ▶ or there is even a way to do set this up in `pivot_longer()` itself

## 4. Additional notes

- ▶ you can clean up the sex variable with a `substr()` command
- ▶ or there is even a way to do set this up in `pivot_longer()` itself
- ▶ and there is `pivot_wider()` for going the other way

## 4. Additional notes

- ▶ you can clean up the sex variable with a `substr()` command
- ▶ or there is even a way to do set this up in `pivot_longer()` itself
- ▶ and there is `pivot_wider()` for going the other way
- ▶ be careful with data in the dataframe that you are not pivoting – frequently wrongly organized
- ▶ → just keep what you need and pivot

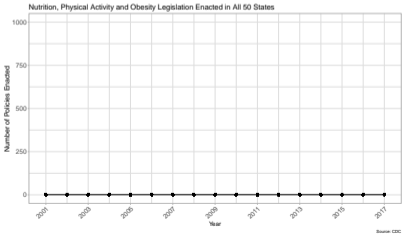
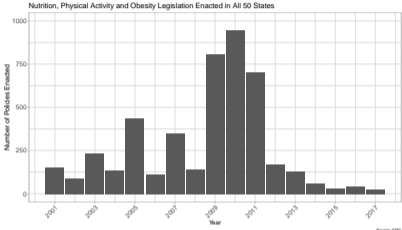
## 5. De-bugging

- ▶ Write a minimal reproducible example
- ▶ Doing this frequently solves your problem
- ▶ Two basic methods
  - ▶ A. start from scratch
  - ▶ B. Remove till problem disappears

Taken largely from Stack Overflow's [advice](#). For Hadley Wickham's official advice, see [here](#).

# 5.a. Start from scratch method

Problem: Jasmine is getting a wacky line graph





## 5.a Jasmine's code

```
enacted %>%  
  group_by(Year, ProvisionID, HealthTopic) %>%  
  slice(1) %>%  
  summarize(num.policies = n()) %>%  
  ggplot() + geom_col(aes(x = Year, y = num.policies))
```

## 5.a Jasmine's code

```
enacted %>%  
  group_by(Year, ProvisionID, HealthTopic) %>%  
  slice(1) %>%  
  summarize(num.policies = n()) %>%  
  ggplot() + geom_col(aes(x = Year, y = num.policies))
```

Ack! A lot of stuff together

## 5.a. My response

- (i) do the data prep before the graph, so the graph command has just the graphing, and not the summarizing or slicing or any of that. So prep the data and then show me the head of the head, with all the variables that go into the graph
- (ii) then the code for both `geom_col` and `geom_line` that provides different estimates.

It's possible that just doing (i) will figure out your problem, but we'll see.

LB

## 5.a. My response

- (i) do the data prep before the graph, so the graph command has just the graphing, and not the summarizing or slicing or any of that. So prep the data and then show me the head of the head, with all the variables that go into the graph
- (ii) then the code for both `geom_col` and `geom_line` that provides different estimates.

It's possible that just doing (i) will figure out your problem, but we'll see.

LB

→ (i) did figure out the problem – data not formatted correctly

## 5.a. How to implement start from scratch?

- ▶ Are data ok?
- ▶ Look at data by themselves
- ▶ Plot bar only
- ▶ Plot line only
- ▶ These should help you narrow down the problematic portion of the code

## 5.b. Remove till problem appears method

- ▶ This is for less obvious serious problems
- ▶ Method:
  - ▶ Get rid of bottom half of your code
  - ▶ Problem still exist?

## 5.b. Remove till problem appears method

- ▶ This is for less obvious serious problems
- ▶ Method:
  - ▶ Get rid of bottom half of your code
  - ▶ Problem still exist?
  - ▶ Get rid of bottom half of your code
  - ▶ Problem still exist?

## 5.b. Remove till problem appears method

- ▶ This is for less obvious serious problems
- ▶ Method:
  - ▶ Get rid of bottom half of your code
  - ▶ Problem still exist?
  - ▶ Get rid of bottom half of your code
  - ▶ Problem still exist?
  - ▶ etc..



## 5.b. Remove till problem appears method

- ▶ This is for less obvious serious problems
- ▶ Method:
  - ▶ Get rid of bottom half of your code
  - ▶ Problem still exist?
  - ▶ Get rid of bottom half of your code
  - ▶ Problem still exist?
  - ▶ etc..
- ▶ Surely a second-choice method
- ▶ But sometimes necessary
- ▶ I use this most frequently for R Markdown, which is buggy

## 5.c. Minimal Reproducible Example

- ▶ The smallest piece of code that generates your problem
- ▶ May need to include data
- ▶ Frequently, generating this solves your problem
- ▶ stackoverflow has a [great page](#) on this

