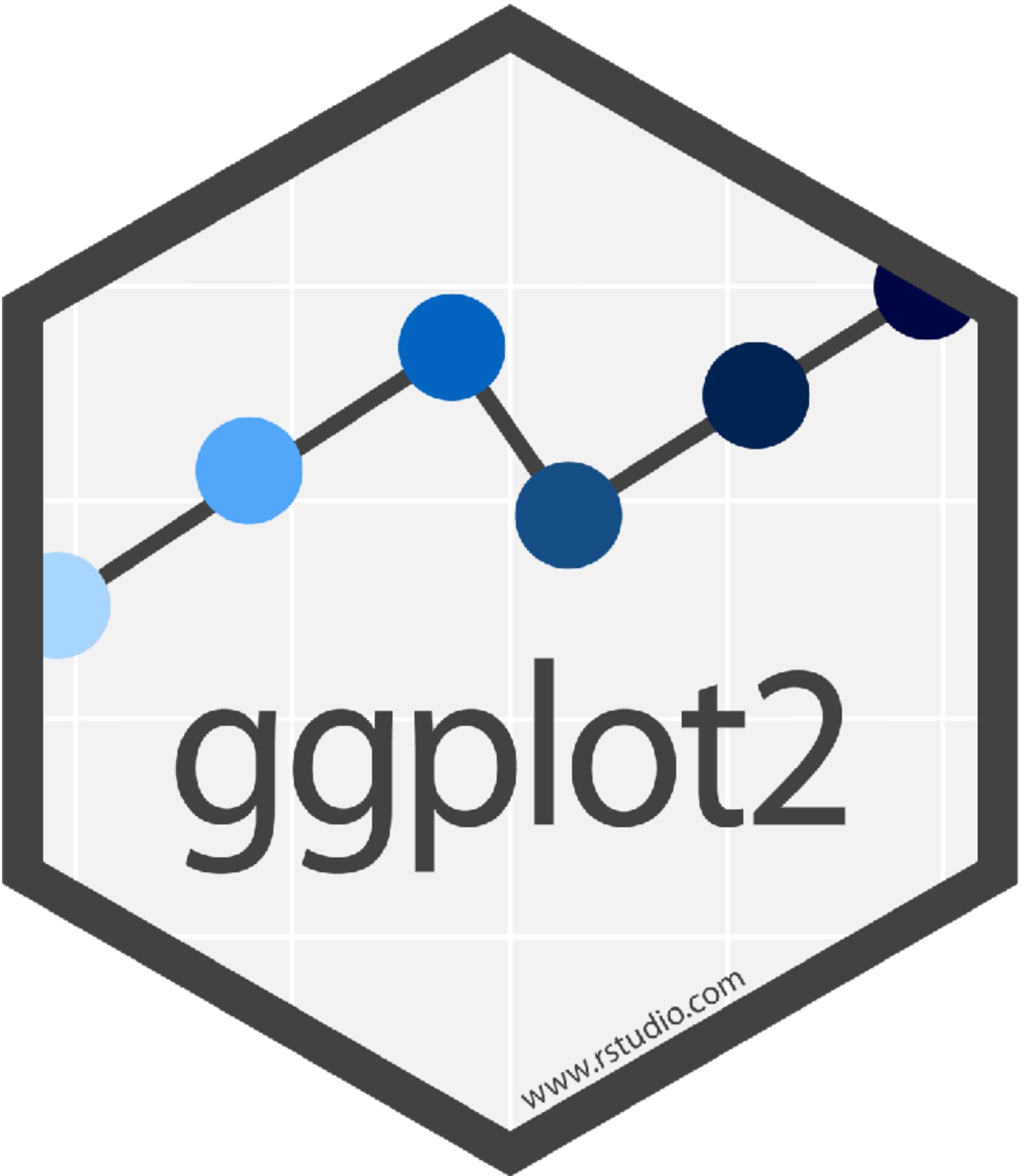


Data Visualization

Jake Thompson

 wjakethompson.com
  [@wjakethompson](https://twitter.com/wjakethompson)



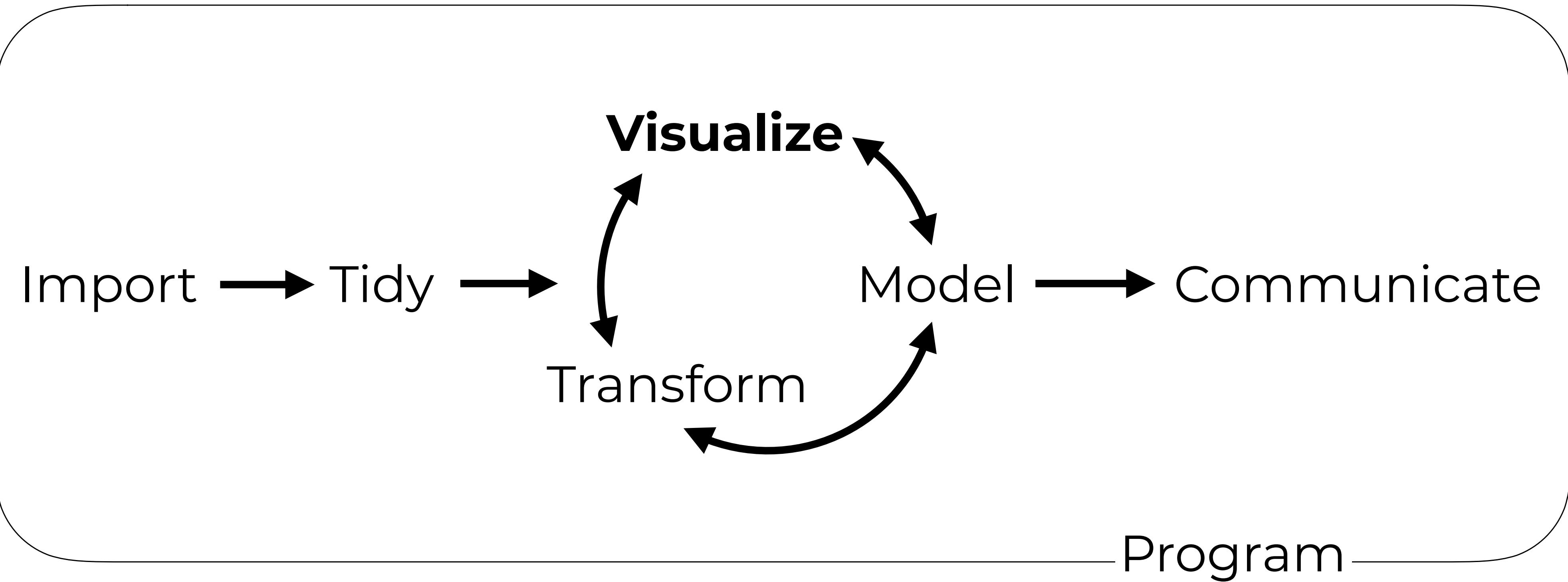


www.rstudio.com

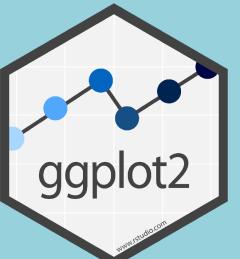


In R4DS
Data Visualisation





Adapted from Master the Tidyverse, CC BY RStudio





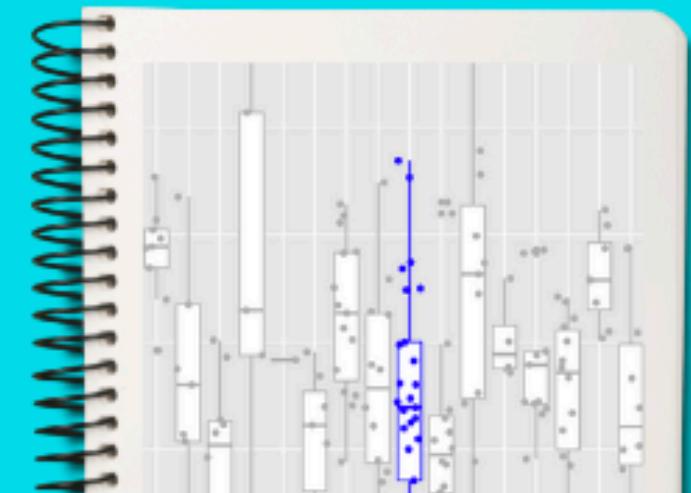
Lucy D'Agostino McGowan

Postdoctoral Researcher
Johns Hopkins University
 @LucyStats

BY LUCY D'AGOSTINO MCGOWAN

GGPLOT2 IN 2

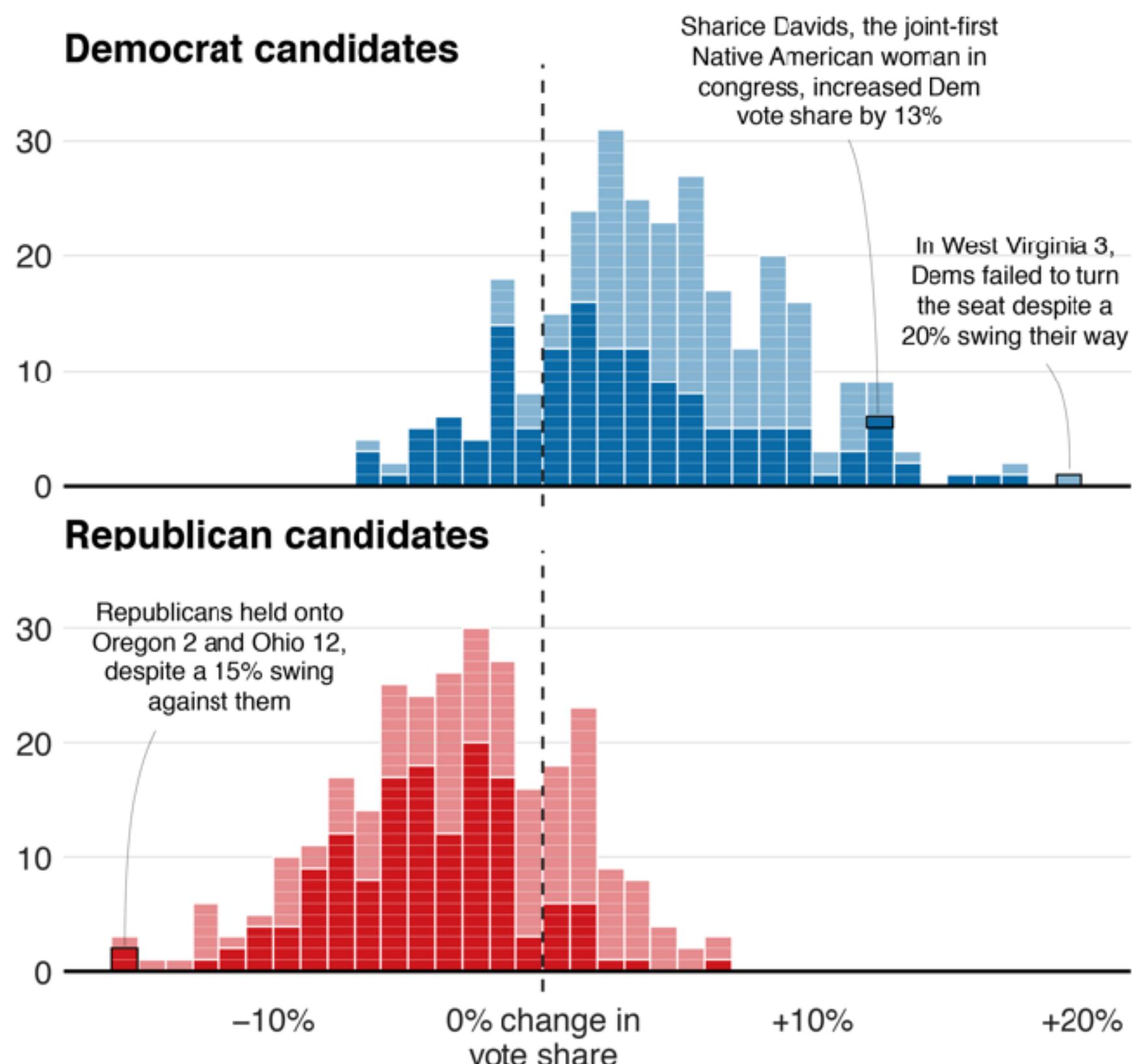
Learn the ggplot2 R
package in two hours!



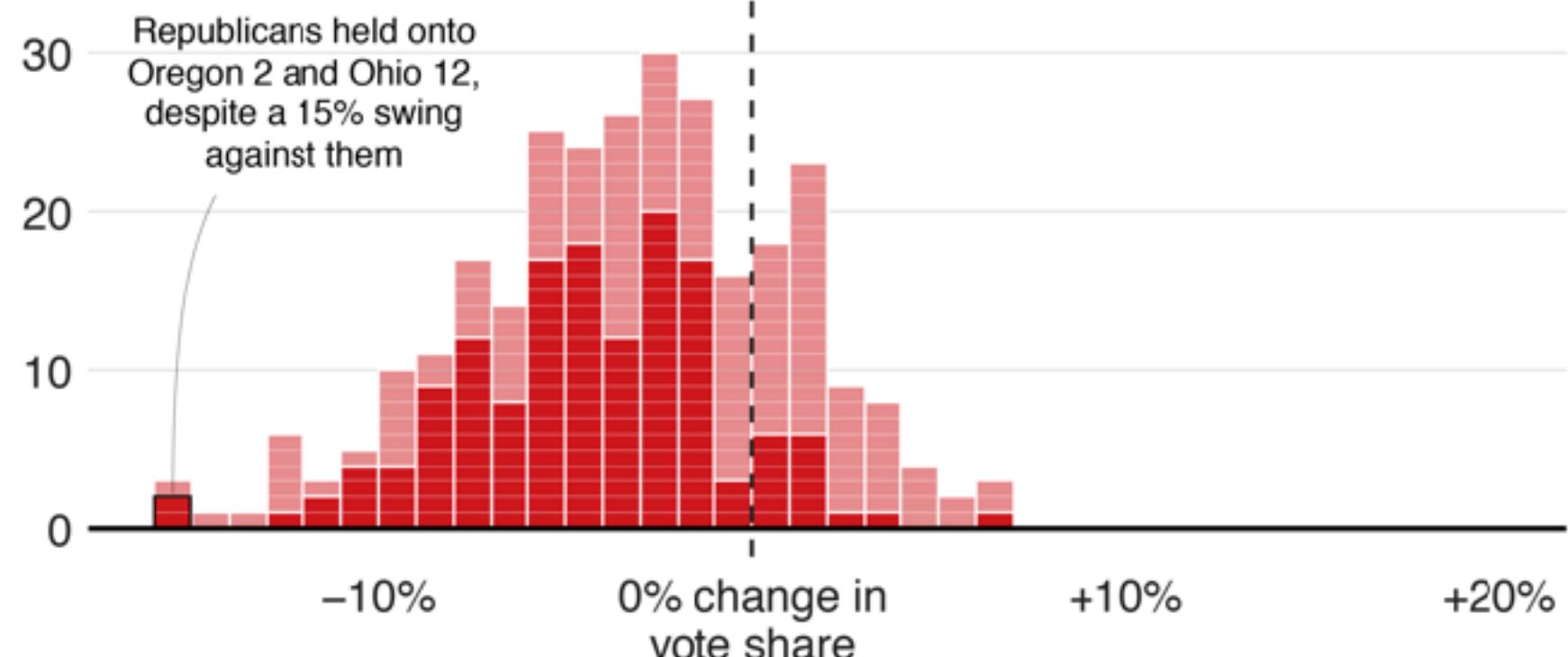
Blue wave

■ Won seat ■ Didn't win

Democrat candidates



Republican candidates



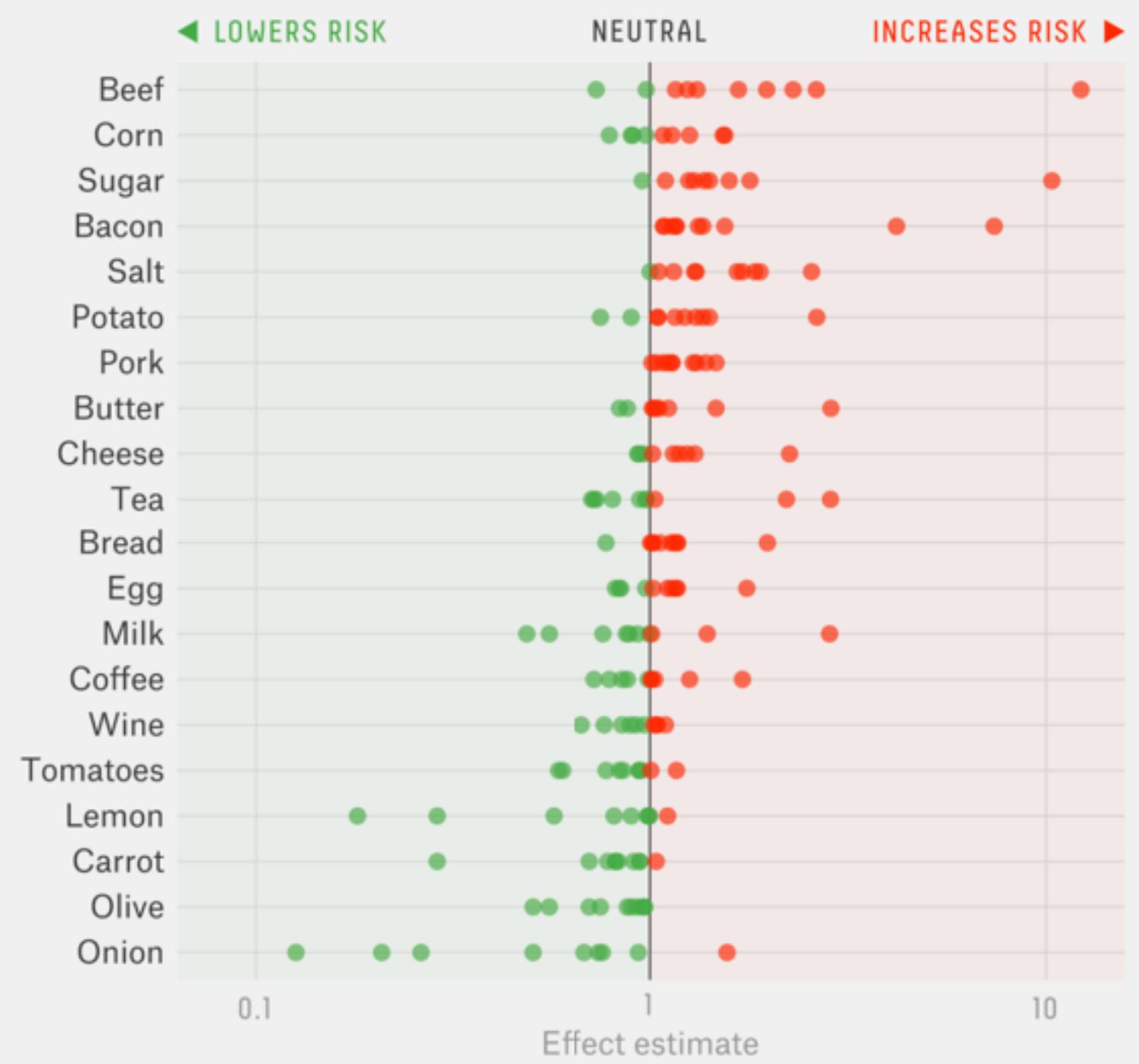
Source: AP, 19:01 ET

BBC

BBC News

Foods that may or may not give you cancer

Risk estimates for 20 foods (each studied at least 10 times) from a 2012 meta-analysis

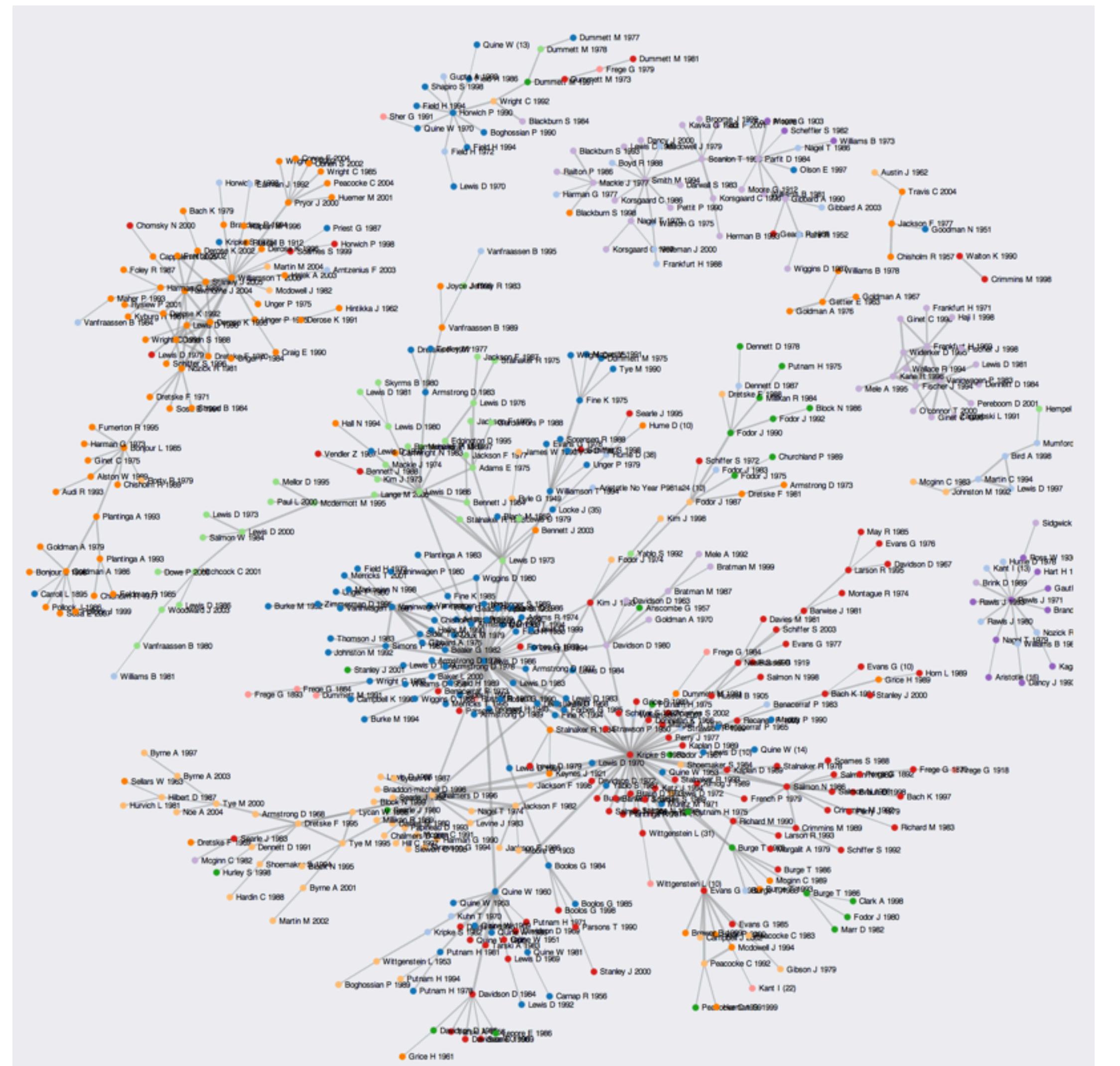


One outlier study not shown (corn, risk estimate of 19.43).

FIVETHIRTYEIGHT

SOURCE: AMERICAN JOURNAL OF CLINICAL NUTRITION

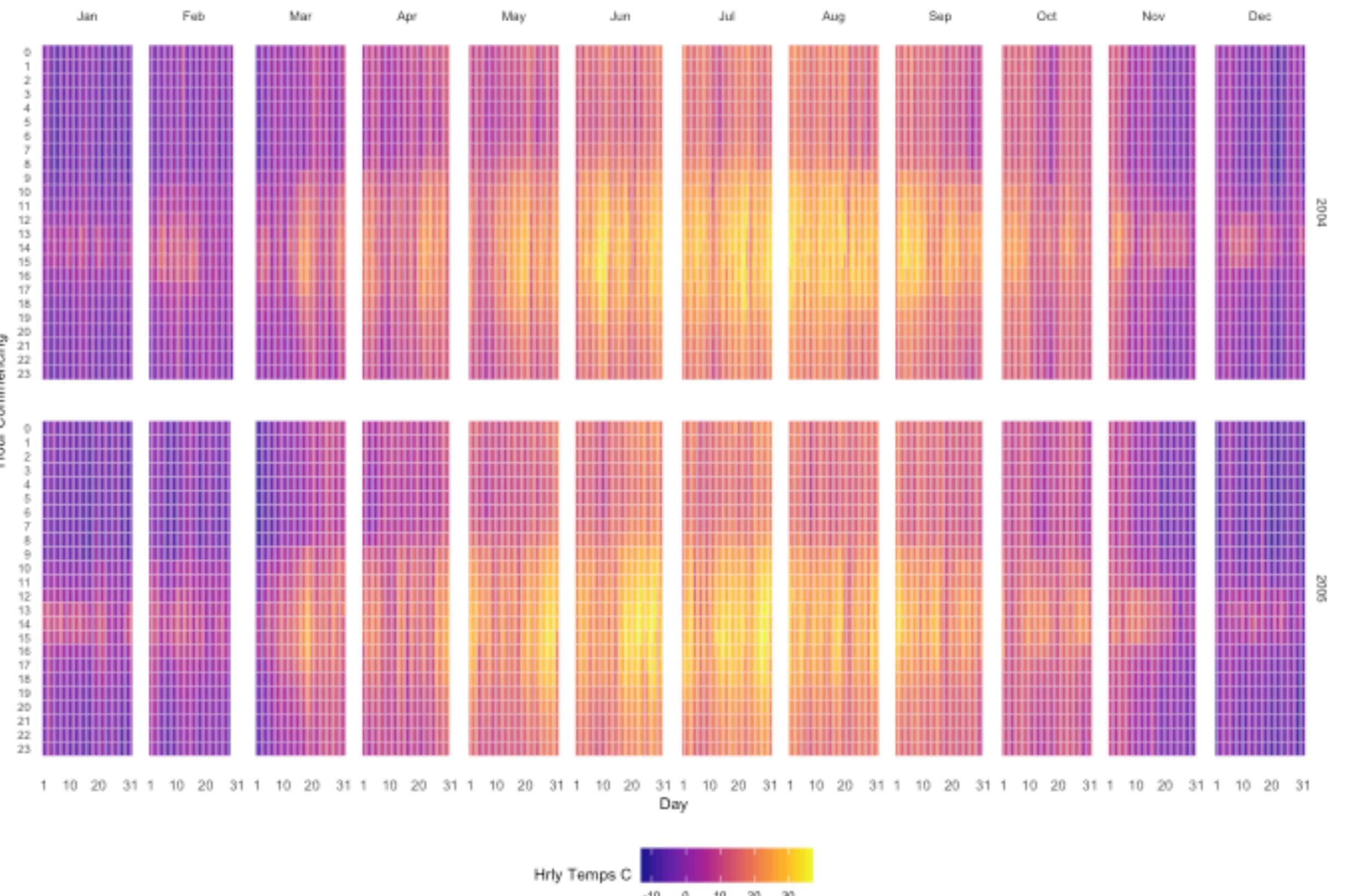
FiveThirtyEight



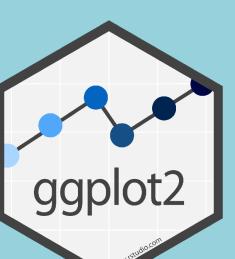
Philosophy Co-Citation Network, Kieran Healy



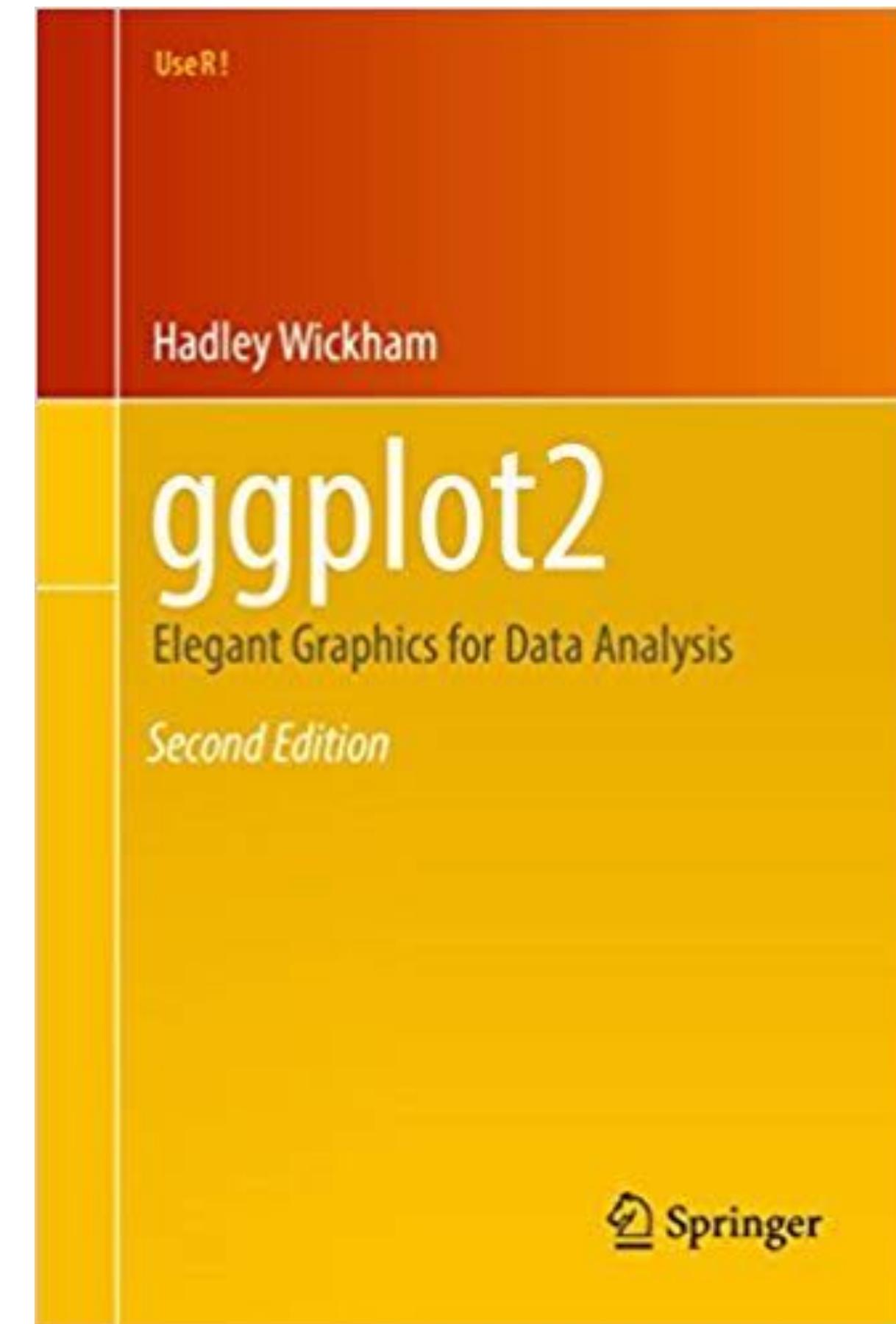
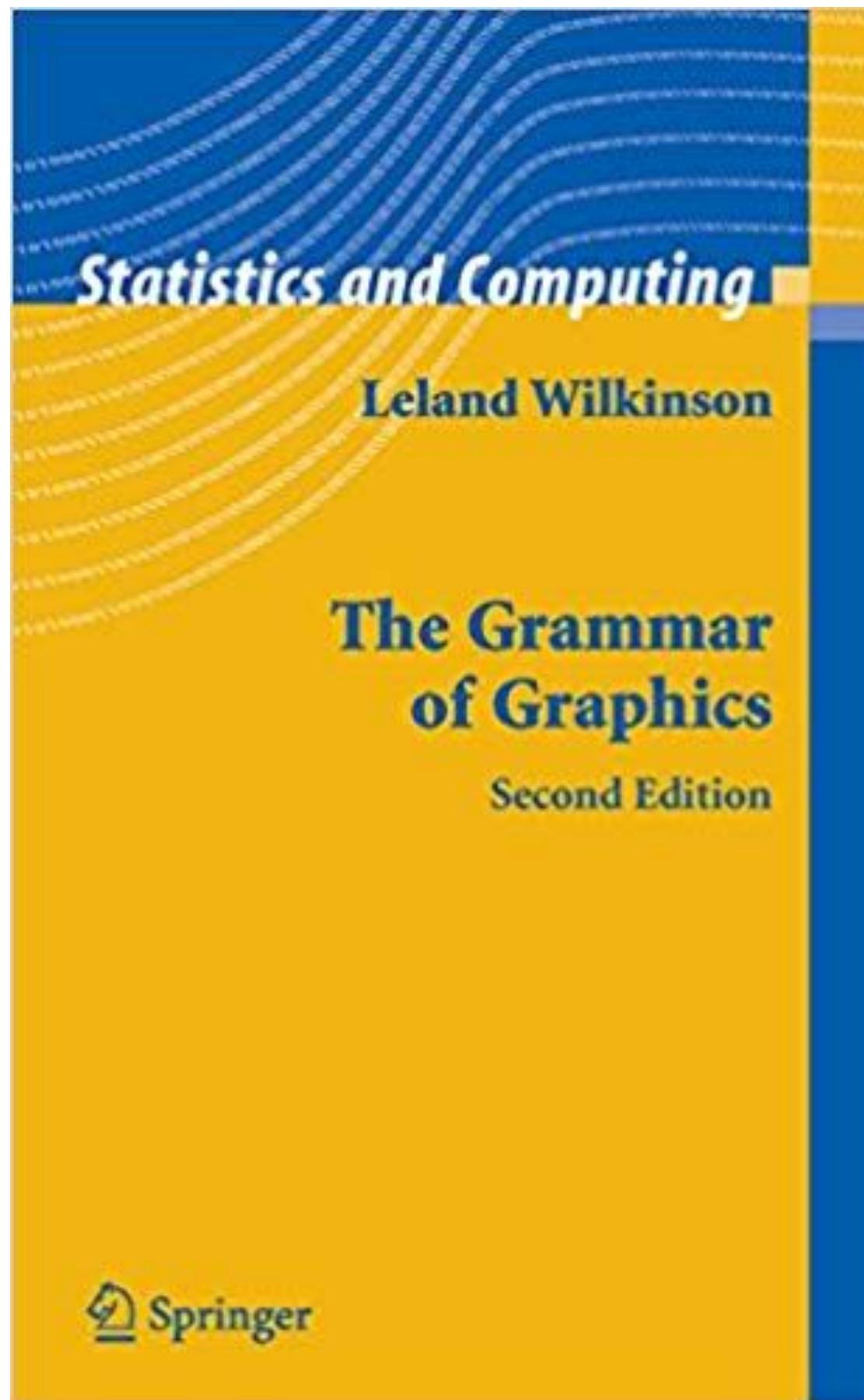
Hourly Temps - Station T0001



R Graph Gallery, John MacKintosh

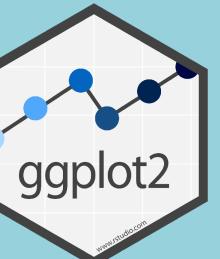


Grammar of Graphics

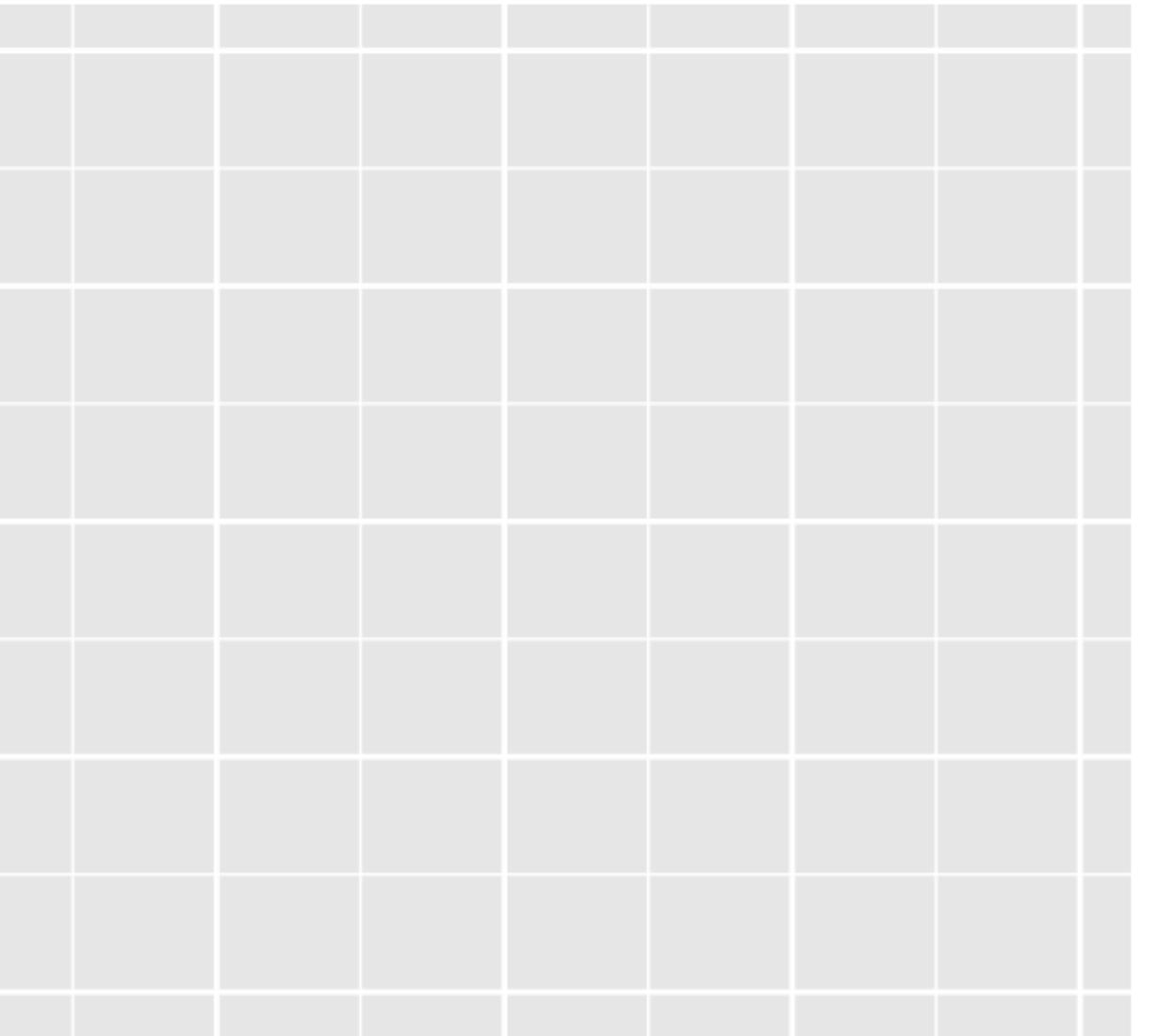


ggplot2

- **data** maps to
 - **aesthetics** in
 - layers

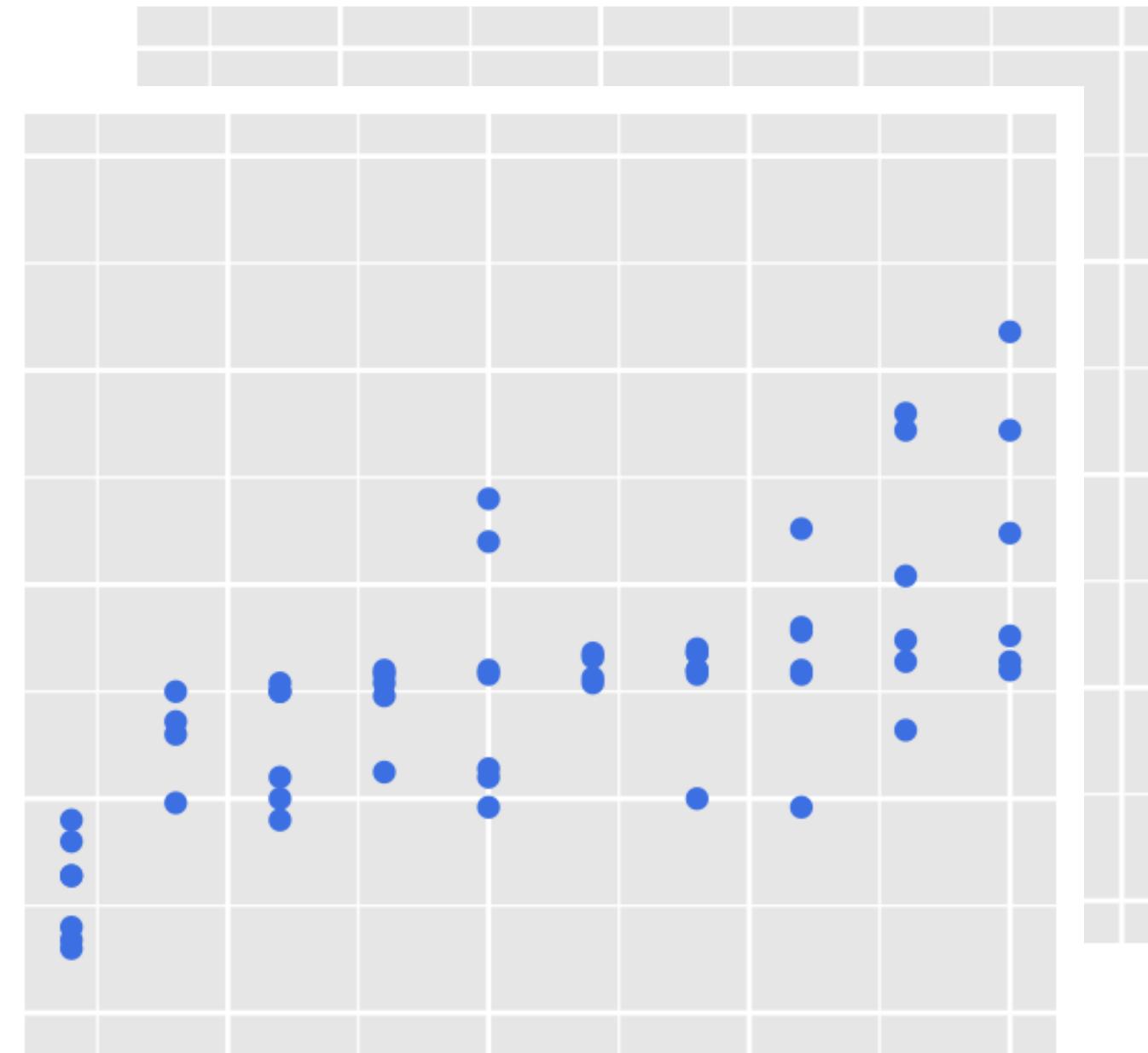


ggplot2 layers



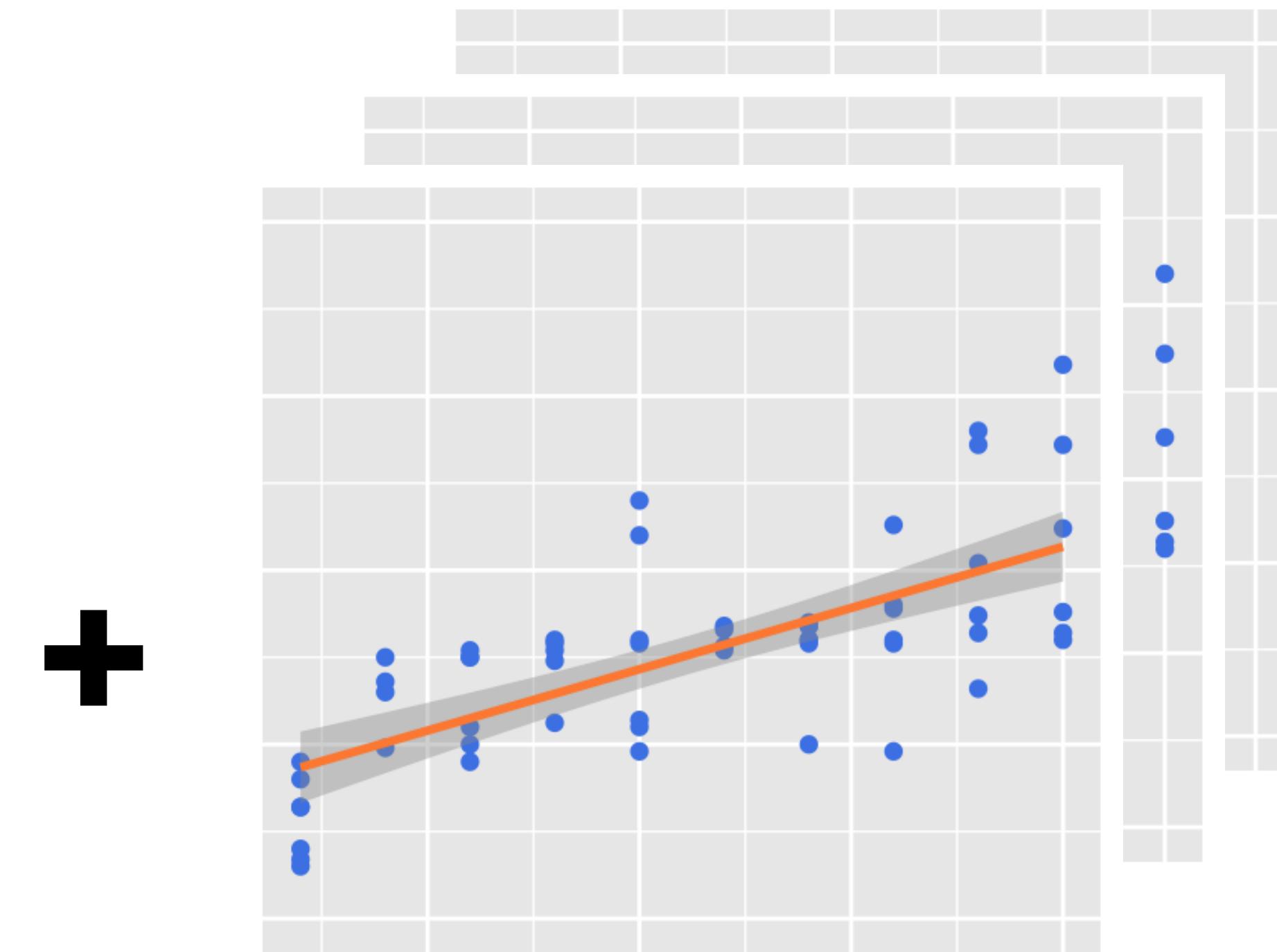
ggplot2 layers

+



geometric objects

ggplot2 layers

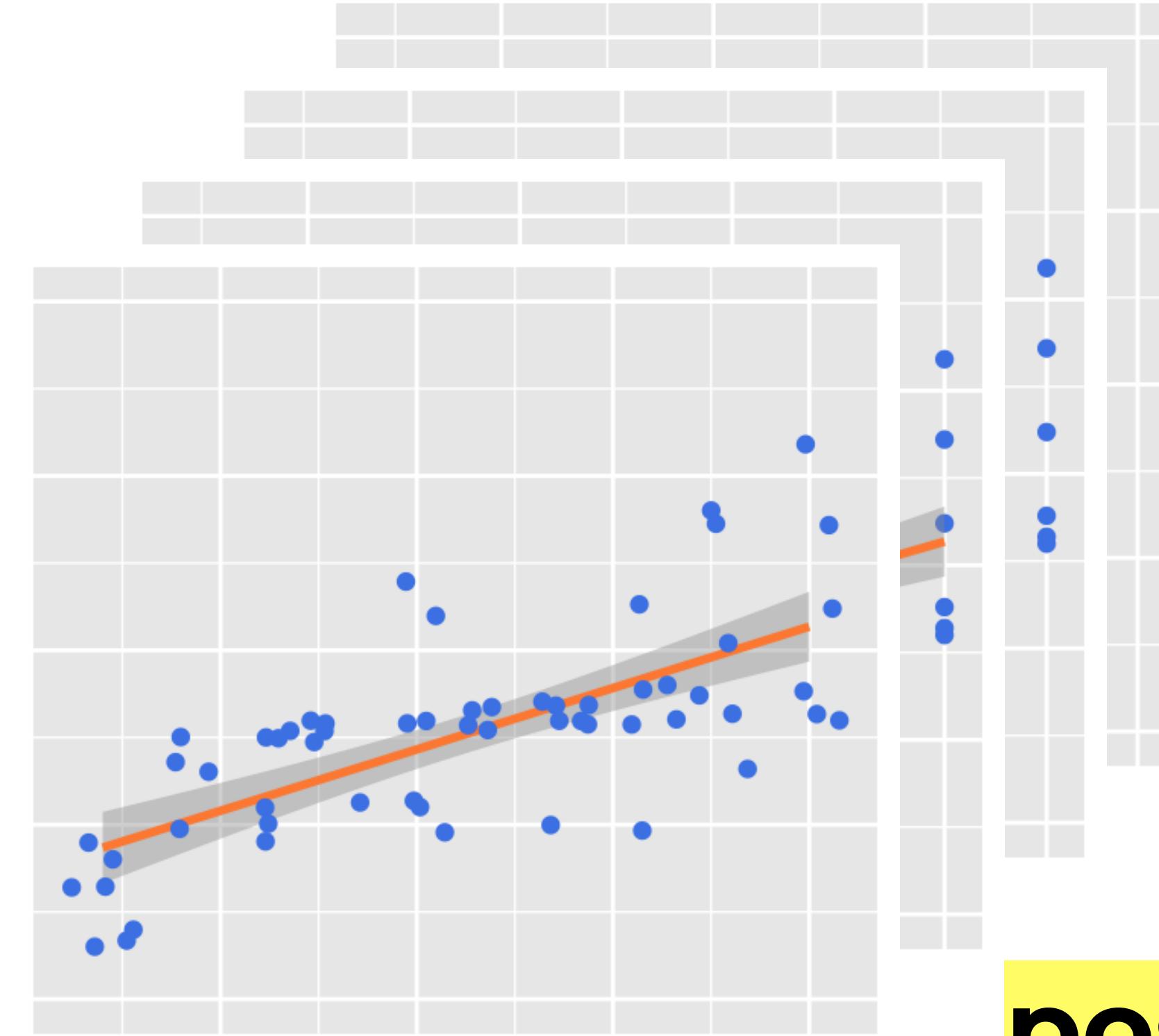


geometric objects

statistical transformations

ggplot2 layers

+

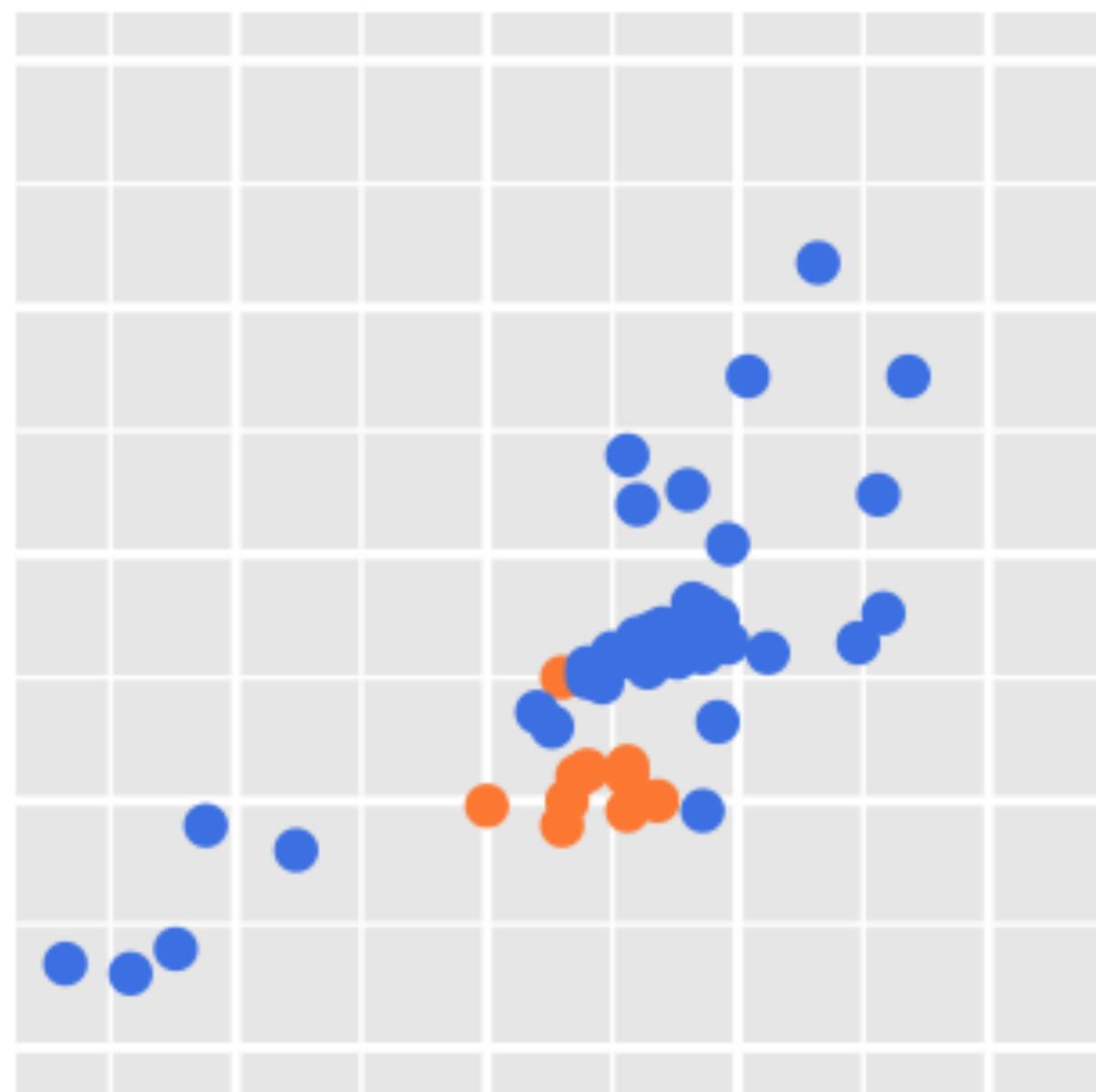


geometric objects

statistical transformations

position adjustments

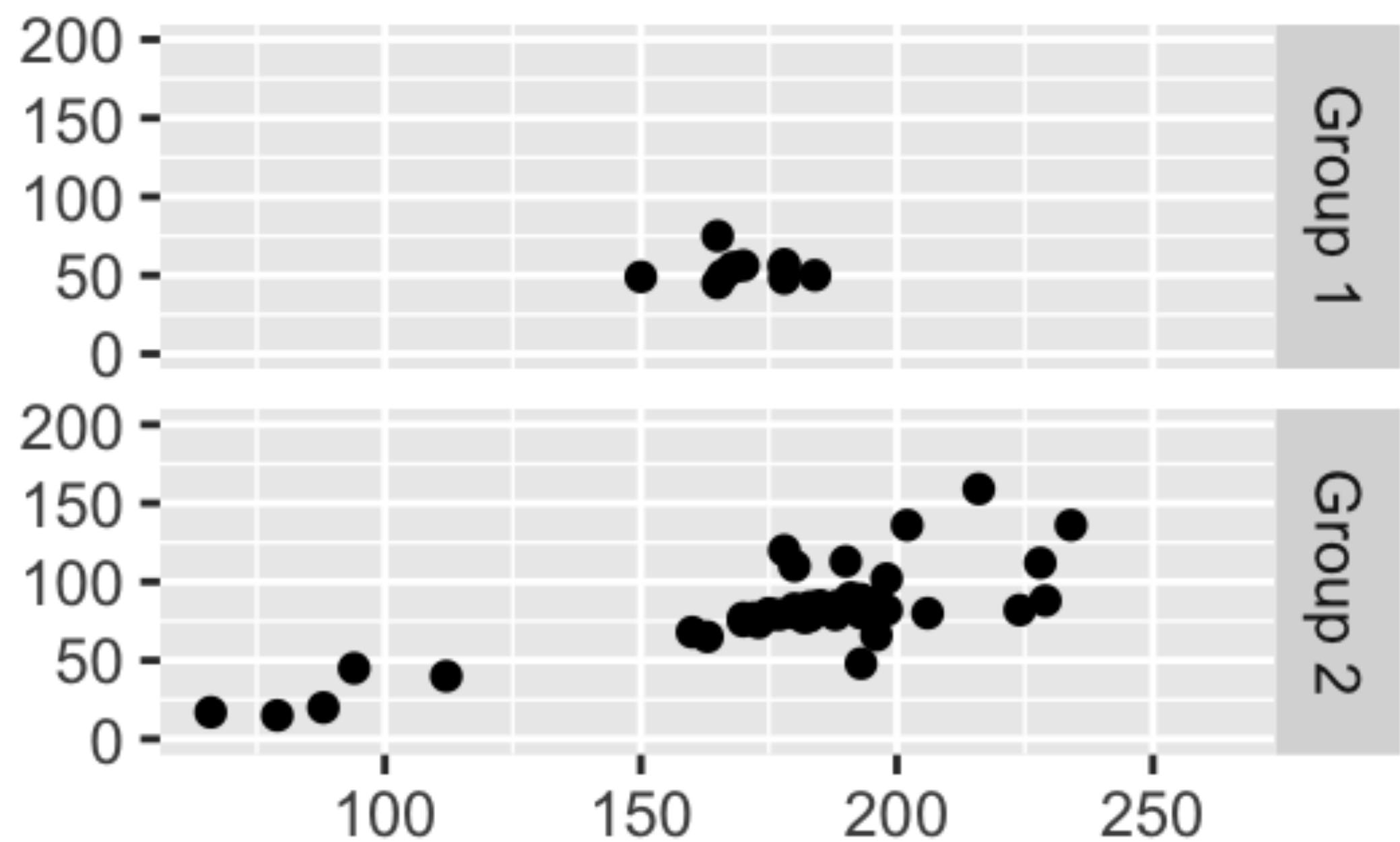
ggplot2 scales



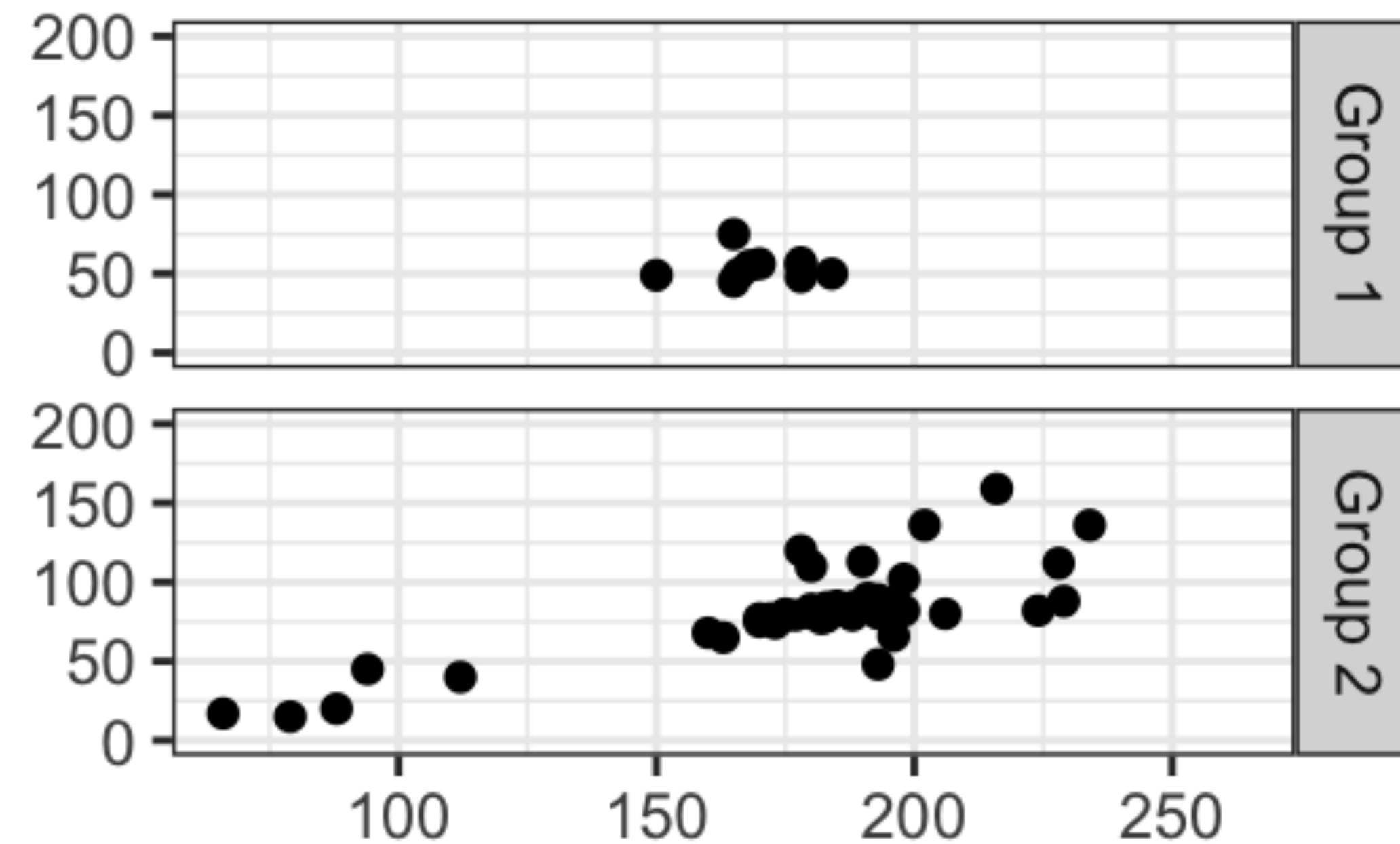
Groups

- Group 1
- Group 2

ggplot2 facets



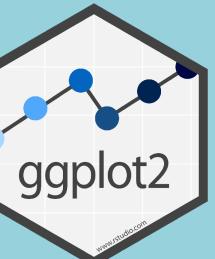
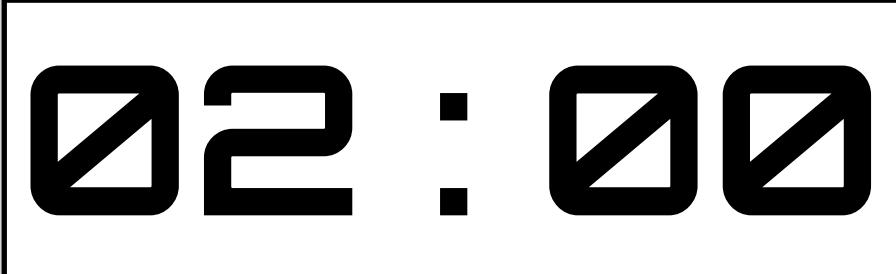
ggplot2 themes

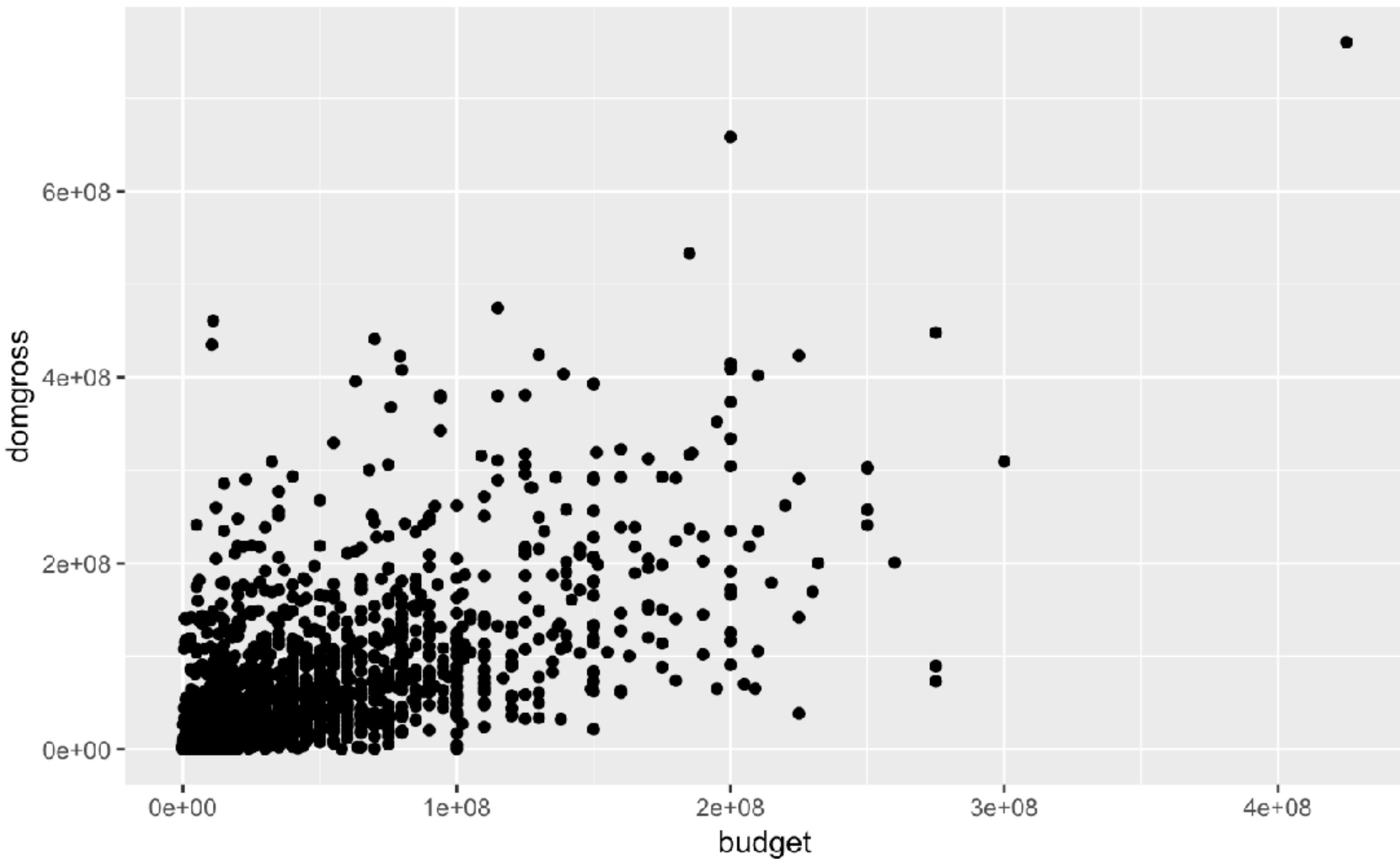


Your Turn 1

- Open the R Notebook **01-Visualize.Rmd**
- Let's look at the **bechdel** data set
- Run this code to make a graph:

```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross))
```





```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross))
```

```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross))
```

data

+ before new line

type of layer

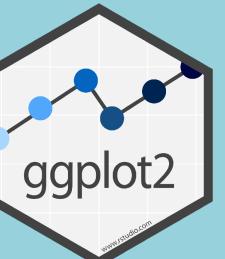
aes()

x variable

y variable

ggplot2 template

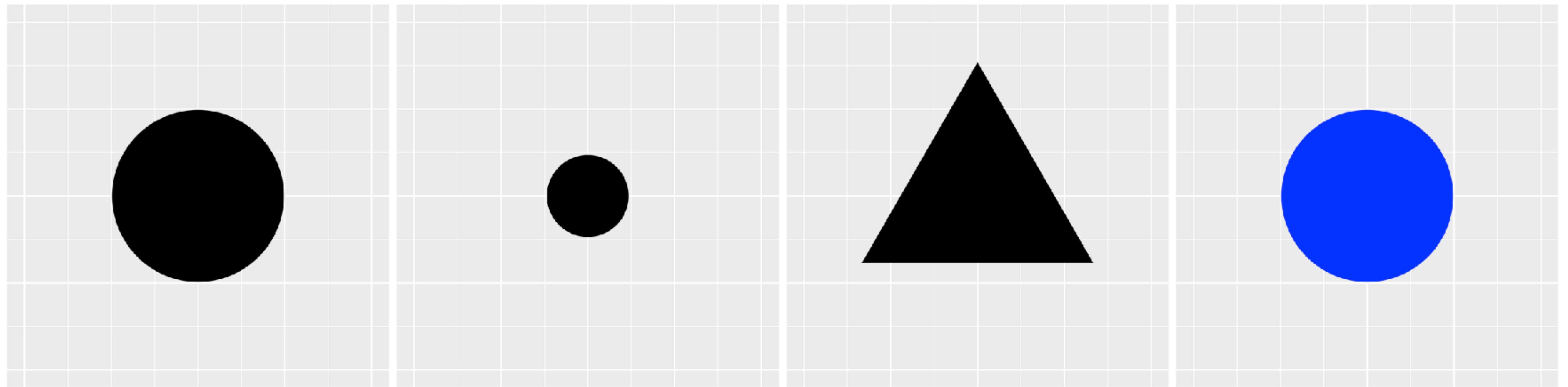
```
ggplot(data = <DATA>) +  
  <GEOM_FUNCTION>(mapping = aes(<MAPPINGS>))
```





Aesthetics

Aesthetics



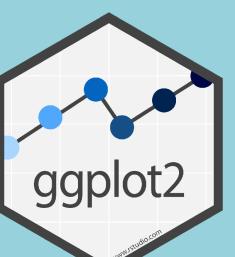
Aesthetics

color	clean_test
Purple	↔ nowomen
Blue	↔ notalk
Teal	↔ men
Lime	↔ dubious
Yellow	↔ ok

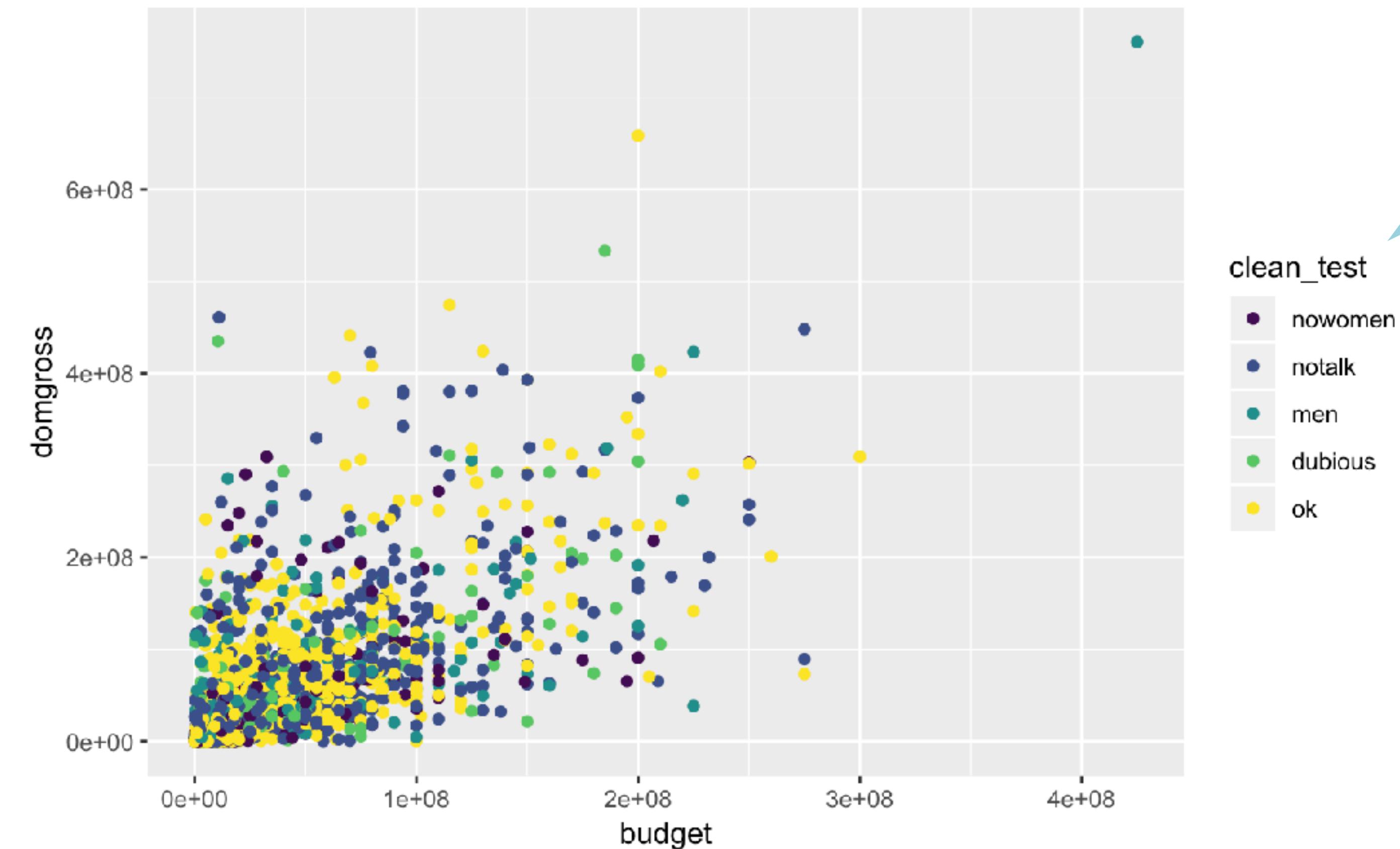
Aesthetic
property

Variable to
map it to

```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, color = clean_test))
```



```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, color = clean_test))
```



Legend added automatically

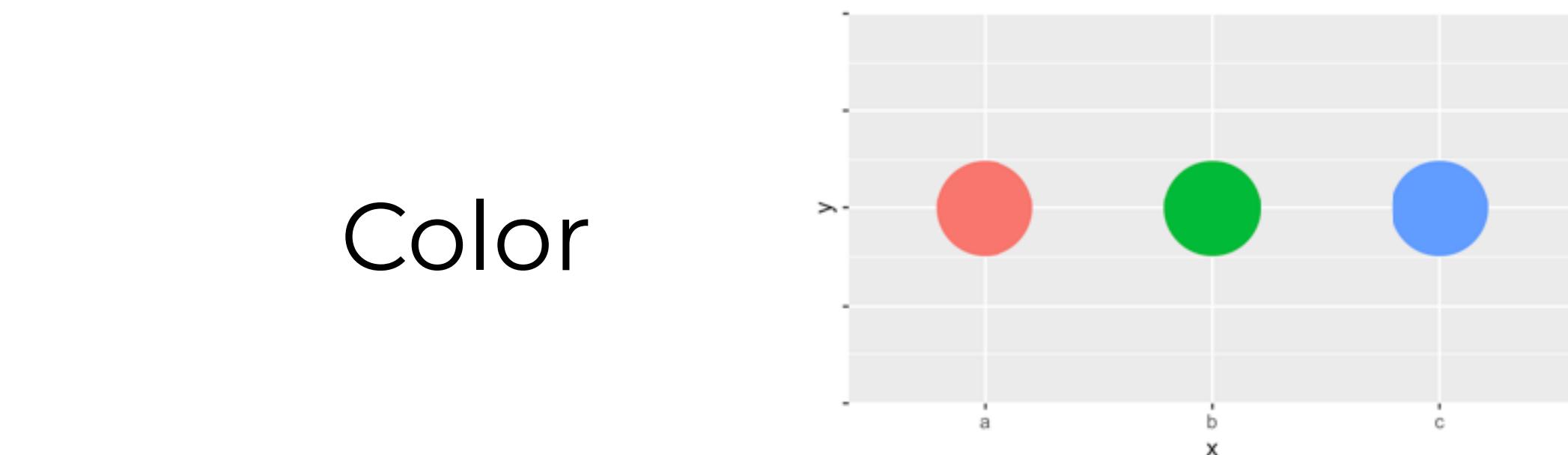
Your Turn 2

- Experiment with adding color, size, alpha, and shape aesthetics to your graph
- How do aesthetics behave different when mapped to discrete and continuous variables?
- What happens when you use more than one aesthetic?

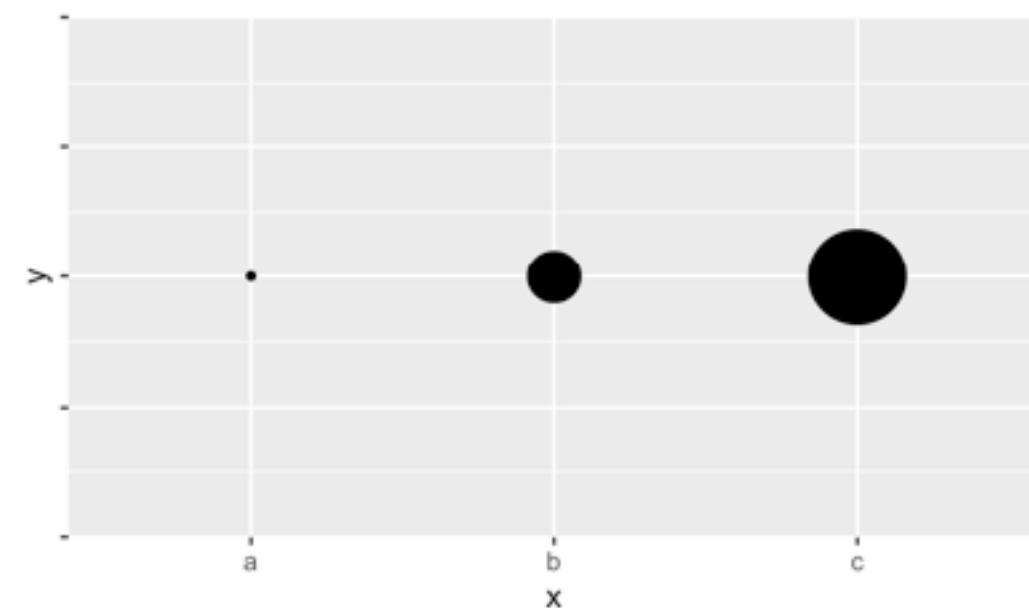
05 : 00



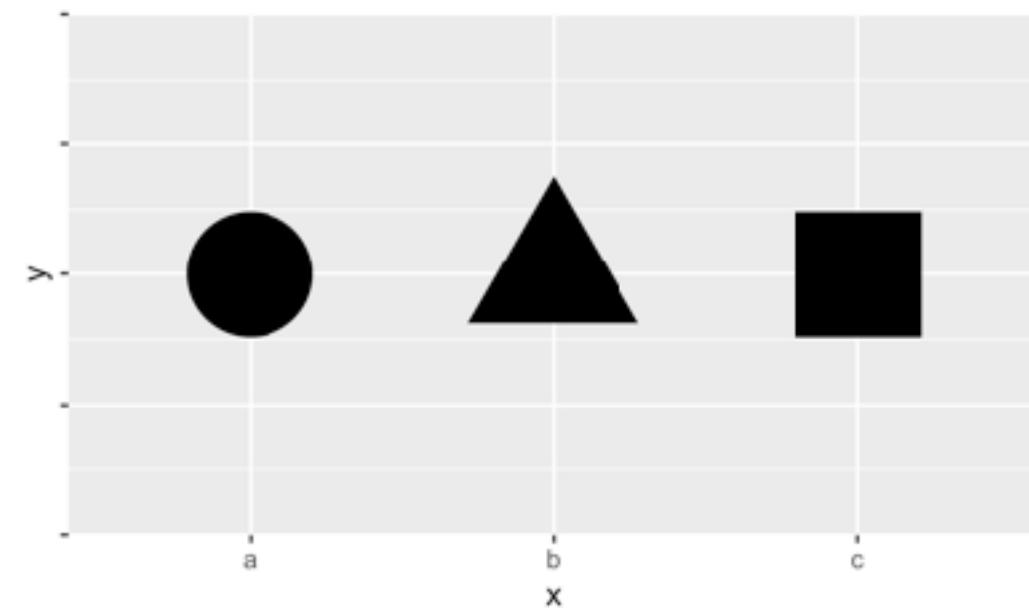
Color



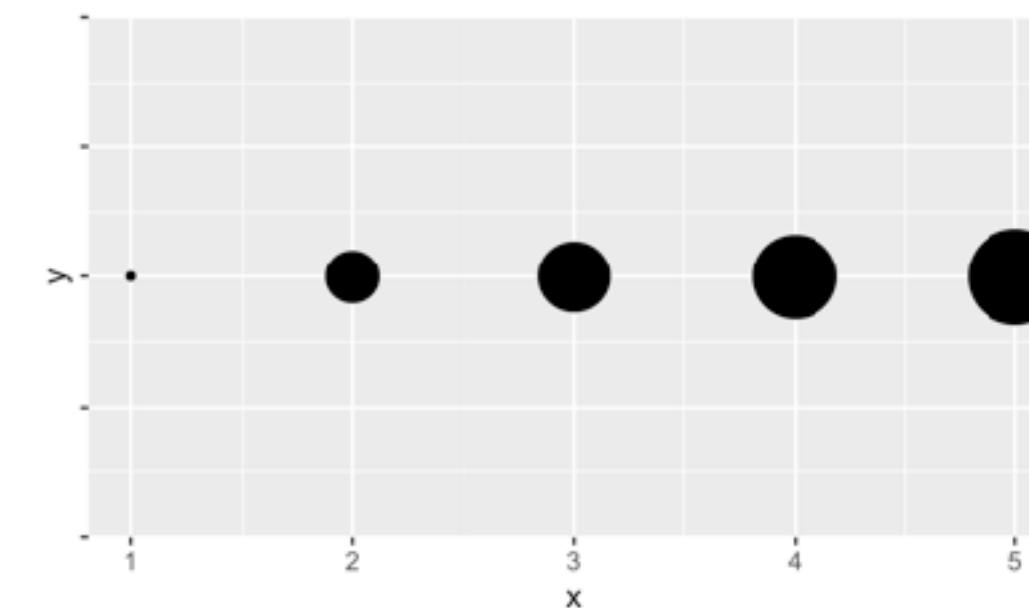
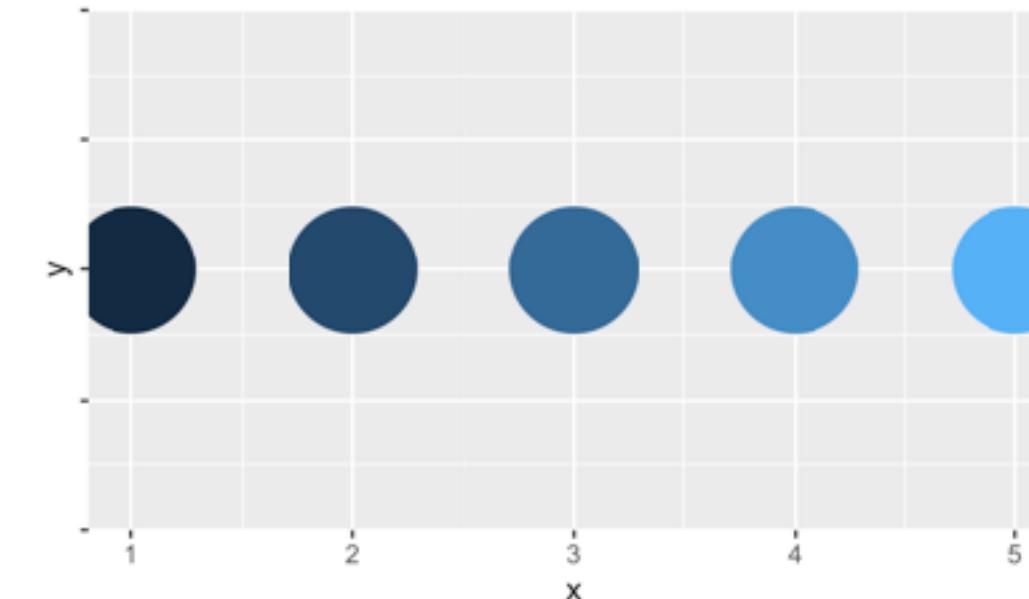
Size



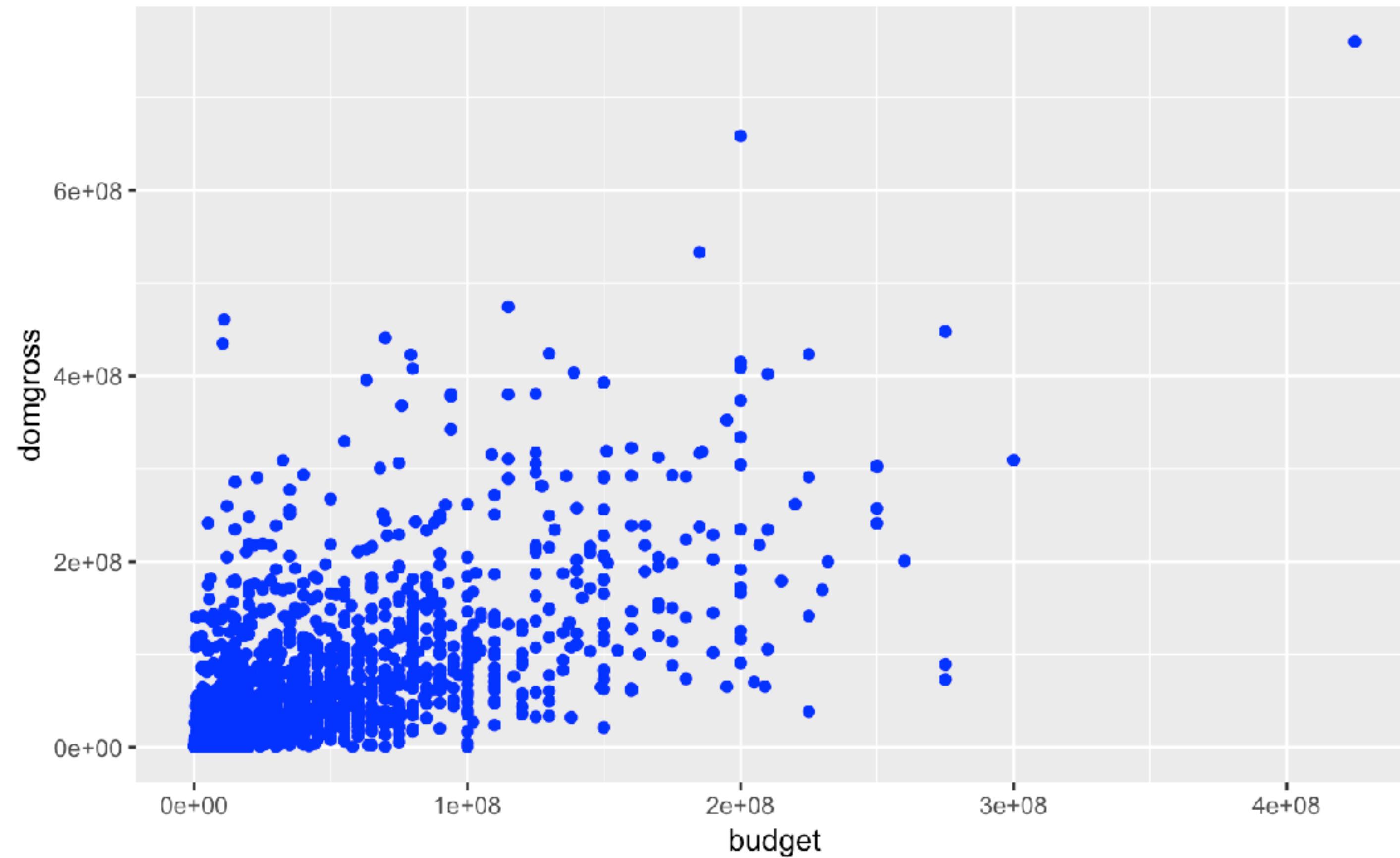
Shape

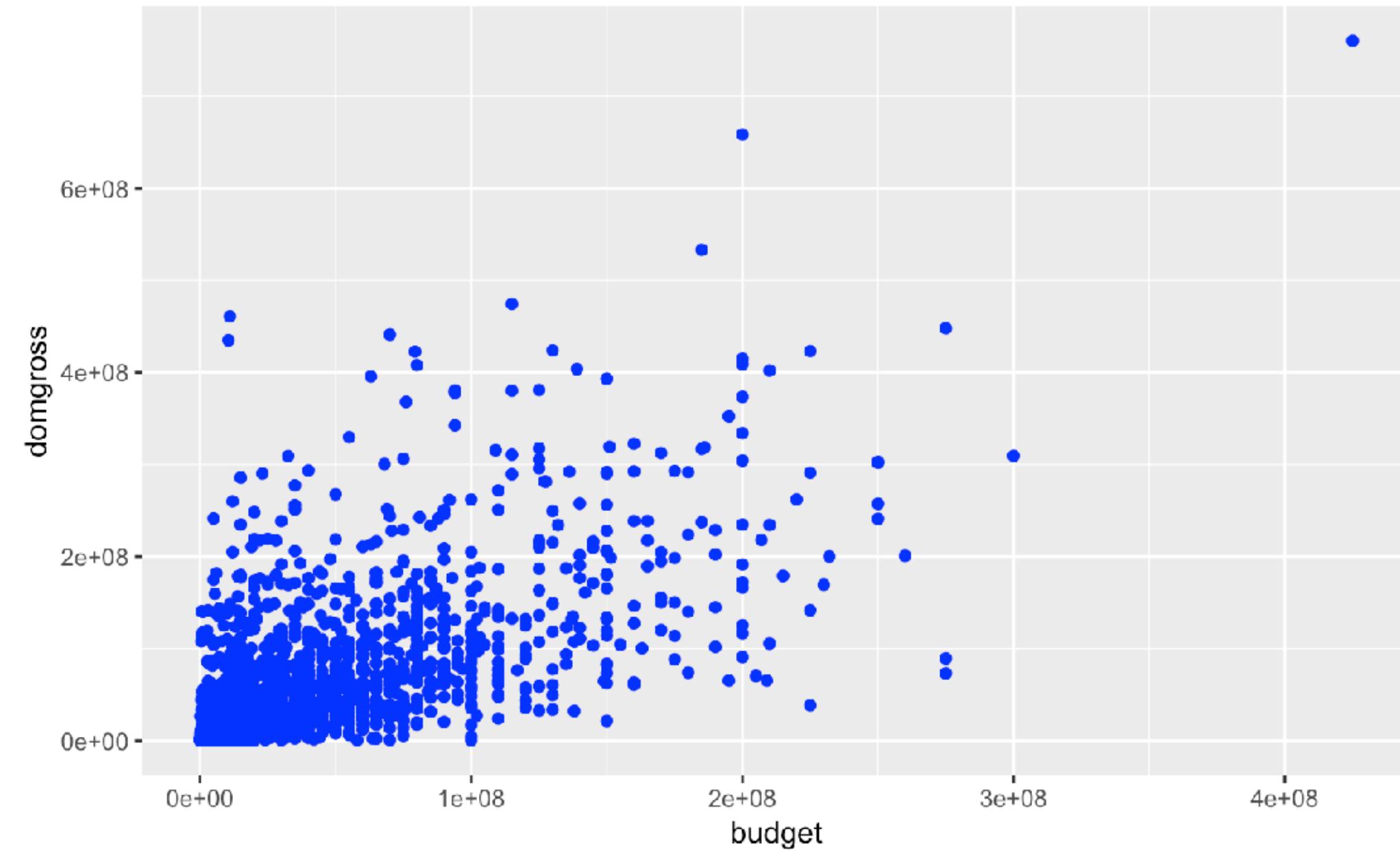


Continuous



Set vs. map

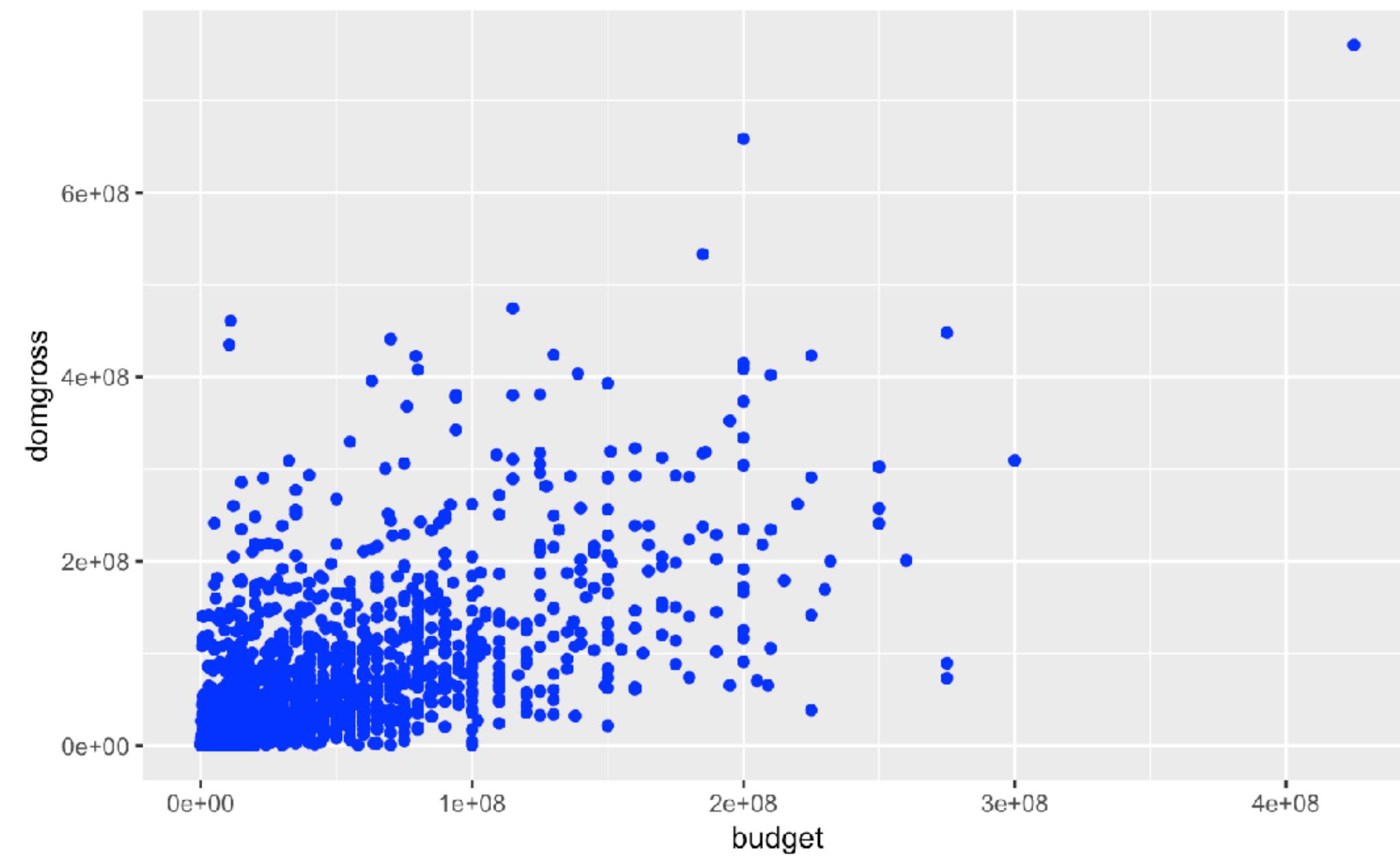
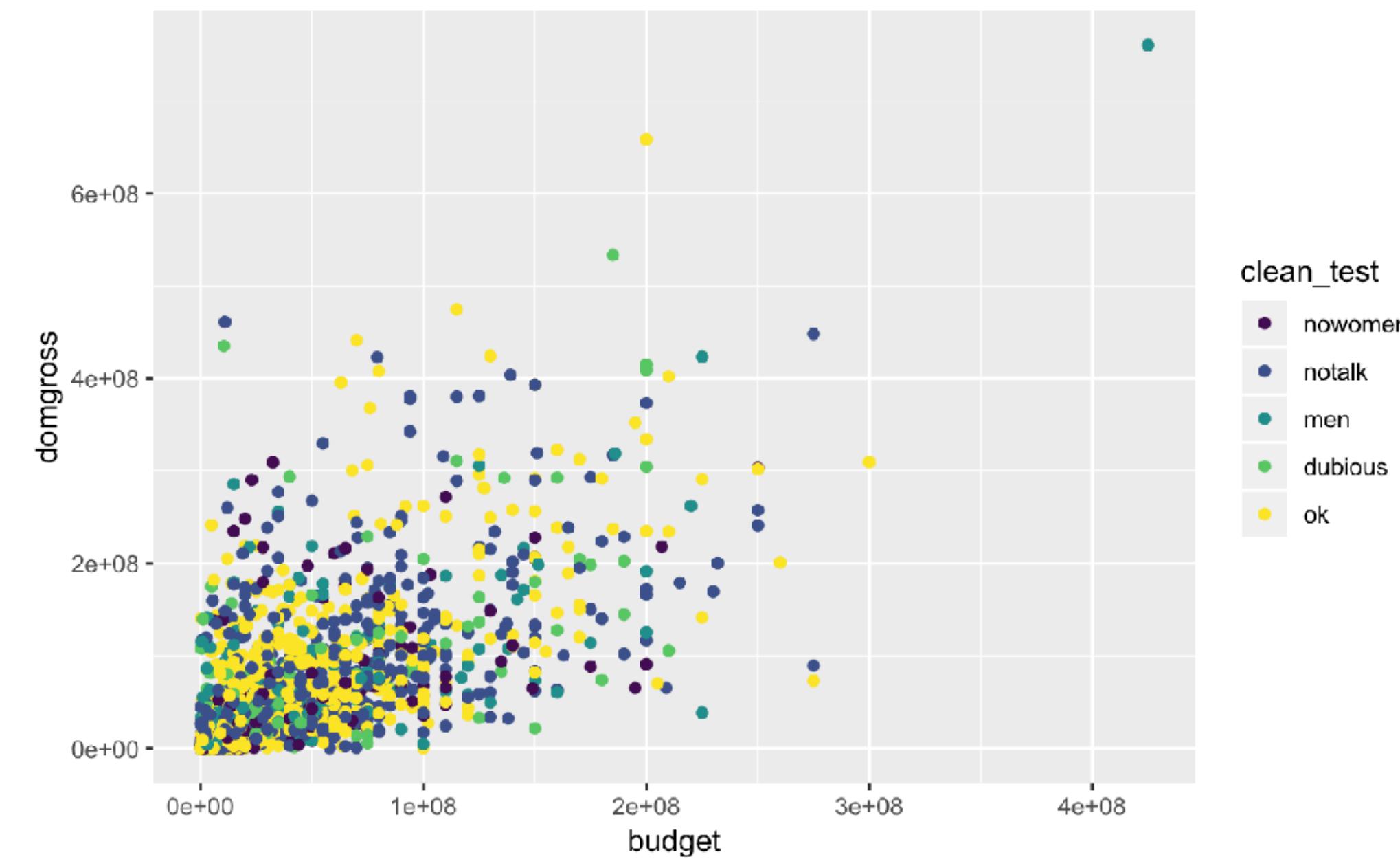




Outside of aes(): sets an aesthetic to a value

```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, color = clean_test))
```

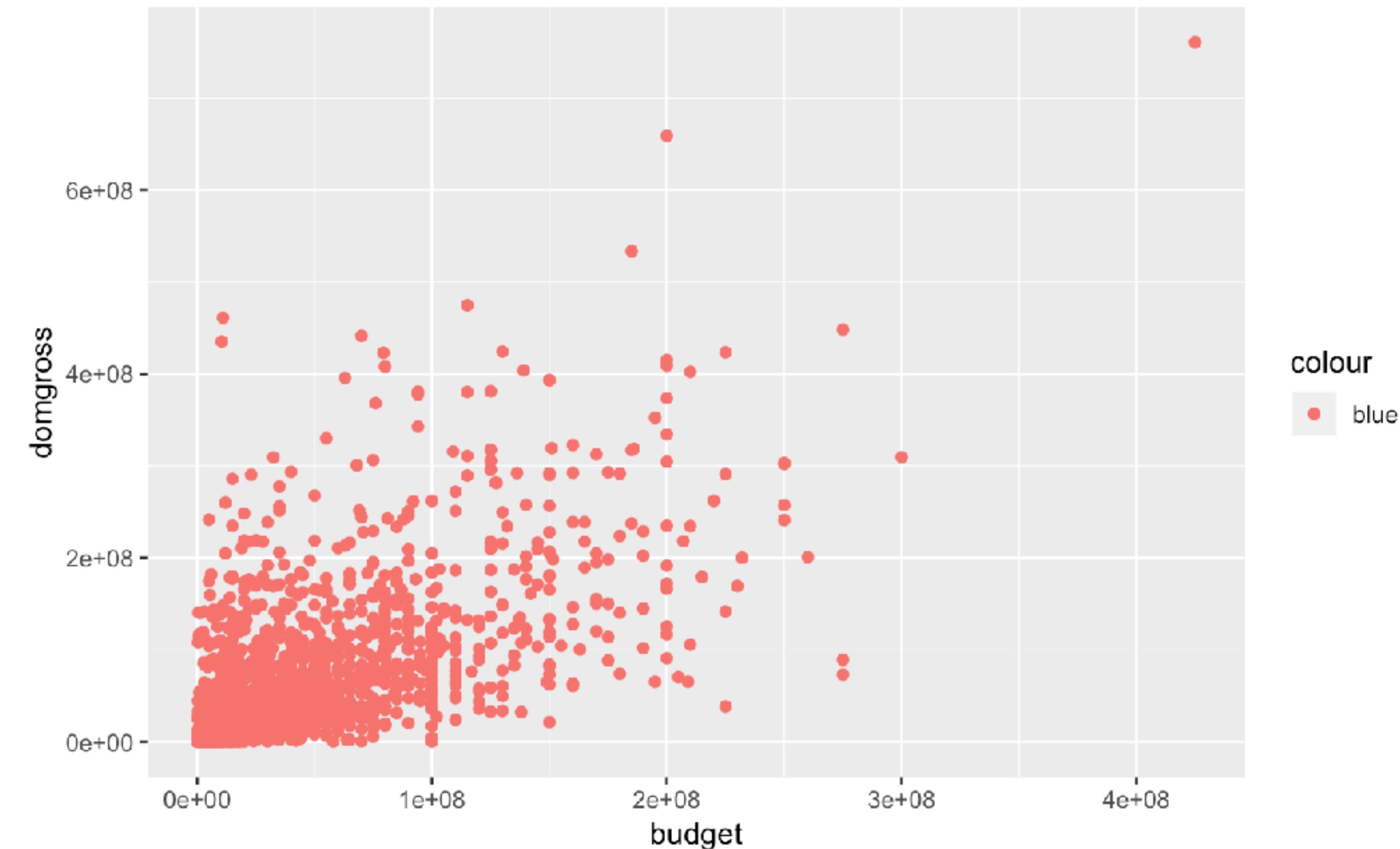
```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross), color = "blue")
```



```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, color = clean_test))
```

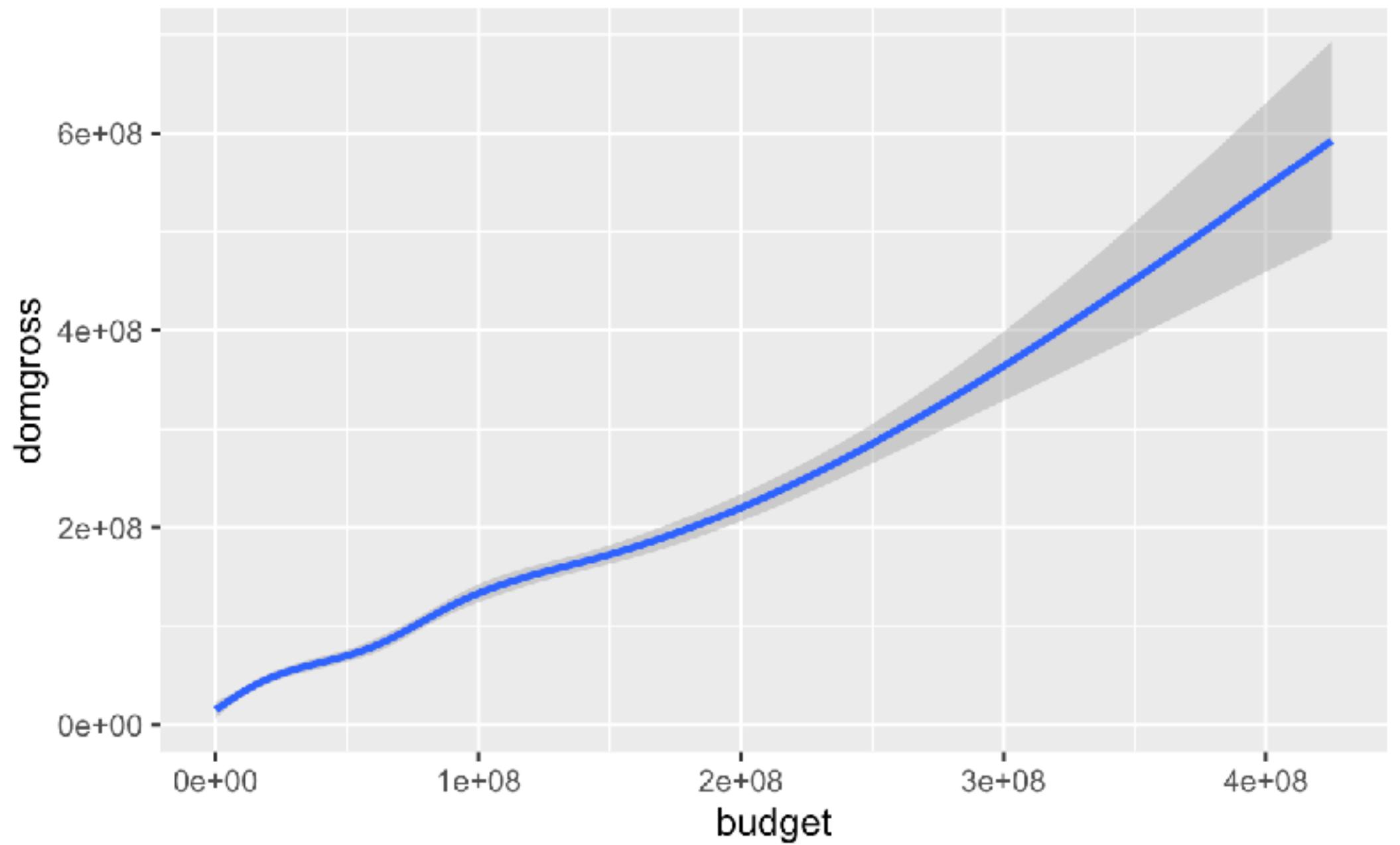
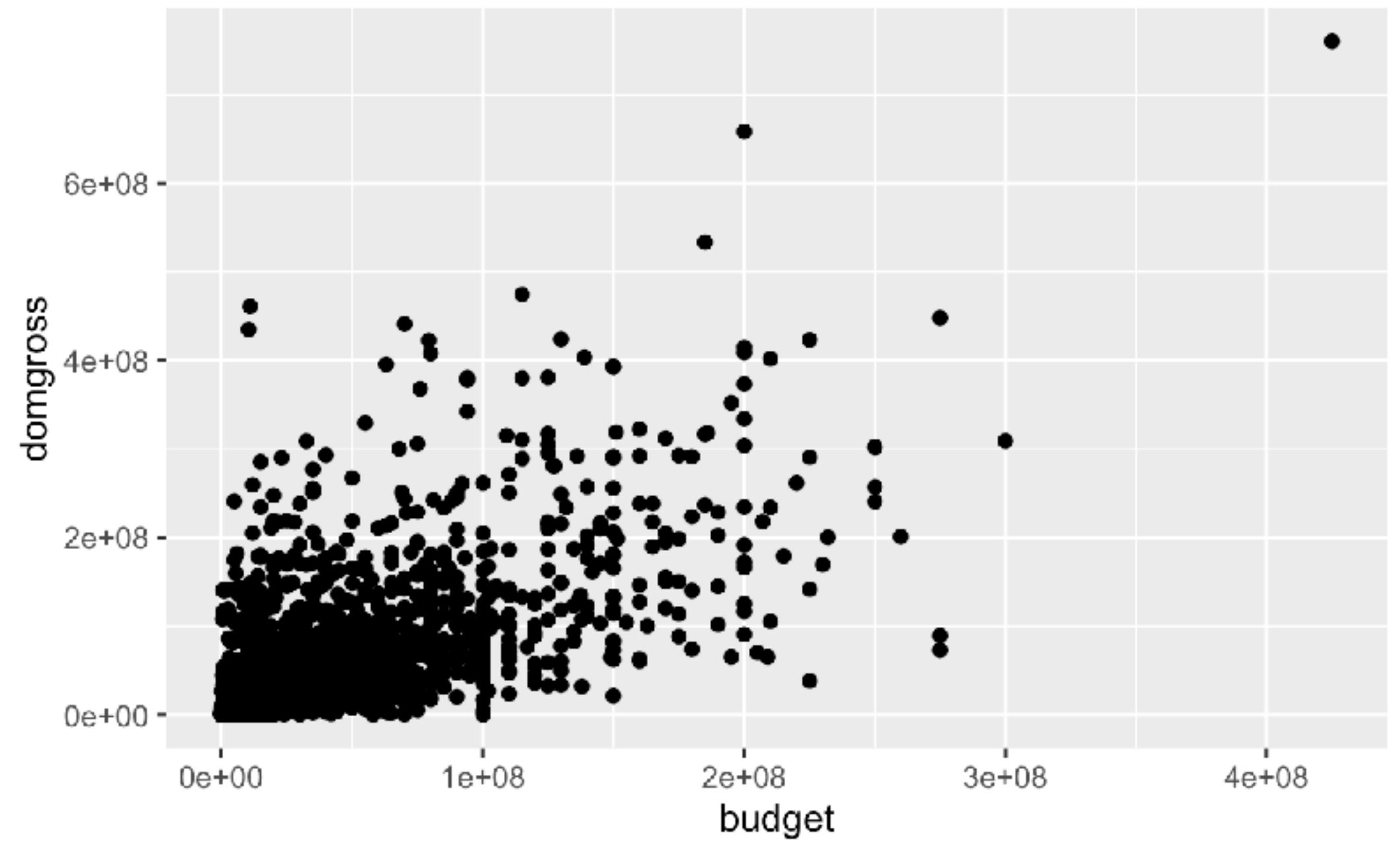
```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross), color = "blue")
```

```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, color = "blue"))
```

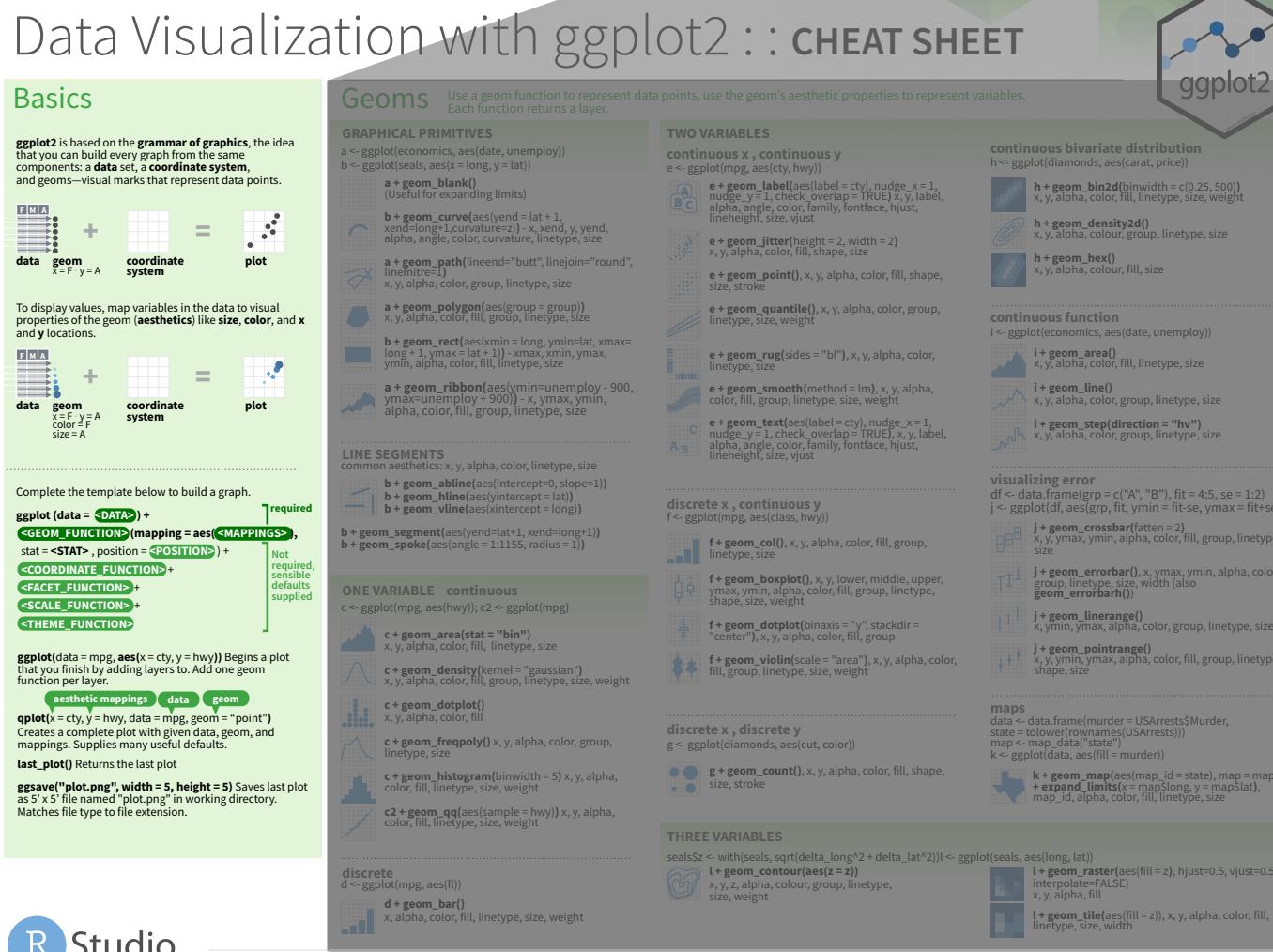




Geoms



ggplot(data = <DATA>) + <GEOM_FUNCTION>(mapping = aes(<MAPPINGS>))



Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables.
Each function returns a layer.

GRAPHICAL PRIMITIVES

```
a <- ggplot(economics, aes(date, unemploy))
b <- ggplot(seals, aes(x = long, y = lat))

a + geom_blank()
b + geom_curve(aes(yend = lat + 1,
  xend=long+1,curvature=z)) - x, yend, y, vend,
alpha, angle, color, curvature, linetype, size
e + geom_label(aes(label = cty), nudge_x = 1,
nudge_y = 1, check_overlap = TRUE) x, y, label,
alpha, angle, color, family, fontface, hjust,
lineheight, size, vjust
f + geom_rect(aes(xmin = long, ymin=lat, xmax=
long + 1, ymax = lat + 1)) - xmin, xmax, ymin,
ymax, alpha, color, fill, linetype, size
g + geom_hex()
h + geom_bin2d(binwidth = c(0.25, 500))
x, y, alpha, color, fill, linetype, size, weight
i + geom_density2d()
x, y, alpha, colour, group, linetype, size
j + geom_hex()
```

TWO VARIABLES

continuous x , continuous y

```
e + geom_label(aes(label = cty), nudge_x = 1,
nudge_y = 1, check_overlap = TRUE) x, y, label,
alpha, angle, color, family, fontface, hjust,
lineheight, size, vjust
f + geom_rect(aes(xmin = long, ymin=lat, xmax=
long + 1, ymax = lat + 1)) - xmin, xmax, ymin,
ymax, alpha, color, fill, linetype, size
g + geom_hex()
h + geom_bin2d(binwidth = c(0.25, 500))
x, y, alpha, color, fill, linetype, size, weight
i + geom_jitter(height = 2, width = 2)
x, y, alpha, color, fill, shape, size
j + geom_point(), x, y, alpha, color, fill, shape,
size, stroke
k + geom_quantile(), x, y, alpha, color, group,
linetype, size, weight
l + geom_rect(sides = "bl"), x, y, alpha, color,
linetype, size
m + geom_rug(sides = "bl"), x, y, alpha, color,
linetype, size
n + geom_ribbon(aes(ymin=unemploy - 900,
  ymax=unemploy + 900)) - x, ymax, ymin,
alpha, color, fill, group, linetype, size
o + geom_smooth(method = lm), x, y, alpha,
color, fill, group, linetype, size, weight
p + geom_text(aes(label = cty), nudge_x = 1,
nudge_y = 1, check_overlap = TRUE) x, y, label,
alpha, angle, color, family, fontface, hjust,
lineheight, size, vjust
q + geom_step(direction = "hv")
```

continuous bivariate distribution

```
h <- ggplot(diamonds, aes(carat, price))

h + geom_hex()
```

continuous function

```
i <- ggplot(economics, aes(date, unemploy))

i + geom_area()
x, y, alpha, color, fill, linetype, size
j + geom_line()
x, y, alpha, color, group, linetype, size
k + geom_step(direction = "hv")
```

visualizing error

```
df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)
j <- ggplot(df, aes(grp, fit, ymin = fit-se, ymax = fit+se))

j + geom_crossbar(fatten = 2)
```

discrete x , continuous y

```
f <- ggplot(mpg, aes(class, hwy))

f + geom_col(), x, y, alpha, color, fill, group,
linetype, size
g + geom_boxplot(), x, y, lower, middle, upper,
ymin, ymax, alpha, color, fill, group, linetype,
shape, size, weight
h + geom_errorbar(), x, y, min, alpha, color,
geom_errorbar()
i + geom_linerange()
x, ymin, ymax, alpha, color, group, linetype, size
j + geom_pointrange()
x, ymin, ymax, alpha, color, fill, group, linetype,
shape, size
```

ONE VARIABLE continuous

```
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)

c + geom_area(stat = "bin")
x, y, alpha, color, fill, linetype, size
c + geom_density(kernel = "gaussian")
x, y, alpha, color, fill, group, linetype, size, weight
c + geom_dotplot()
x, y, alpha, color, fill
c + geom_freqpoly()
x, y, alpha, color, group, linetype, size
c + geom_histogram(binwidth = 5) x, y, alpha,
color, fill, linetype, size, weight
c2 + geom_qq(aes(sample = hwy)) x, y, alpha,
color, fill, linetype, size, weight
```

maps

```
data <- data.frame(murder = USArrests$Murder,
state = tolower(rownames(USArrests)))
map <- map_data("state")
k <- ggplot(data, aes(fill = murder))

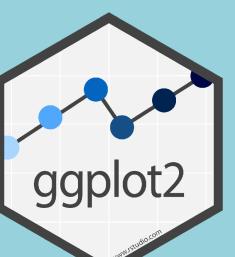
k + geom_map(aes(map_id = state), map = map)
+ expand_limits(x = map$long, y = map$lat),
map_id, alpha, color, fill, linetype, size
```

THREE VARIABLES

```
seals$z <- with(seals, sqrt(delta_long^2 + delta_lat^2)) l <- ggplot(seals, aes(long, lat))

l + geom_contour(aes(z = z))
x, y, z, alpha, colour, group, linetype,
size, weight
l + geom_raster(aes(fill = z), hjust=0.5, vjust=0.5,
interpolate=FALSE)
x, y, alpha, fill
l + geom_tile(aes(fill = z)), x, y, alpha, color, fill,
linetype, size, width
```

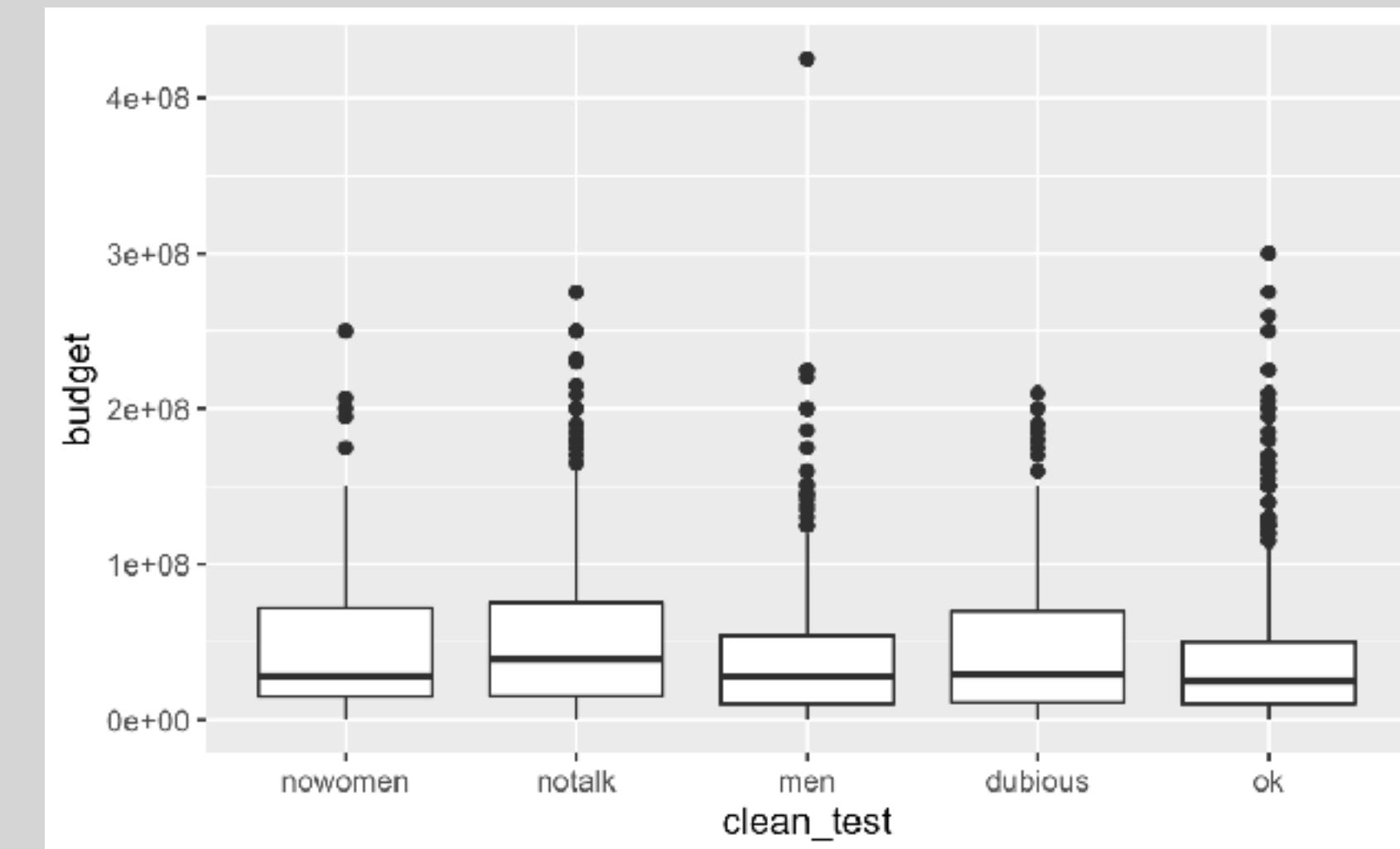
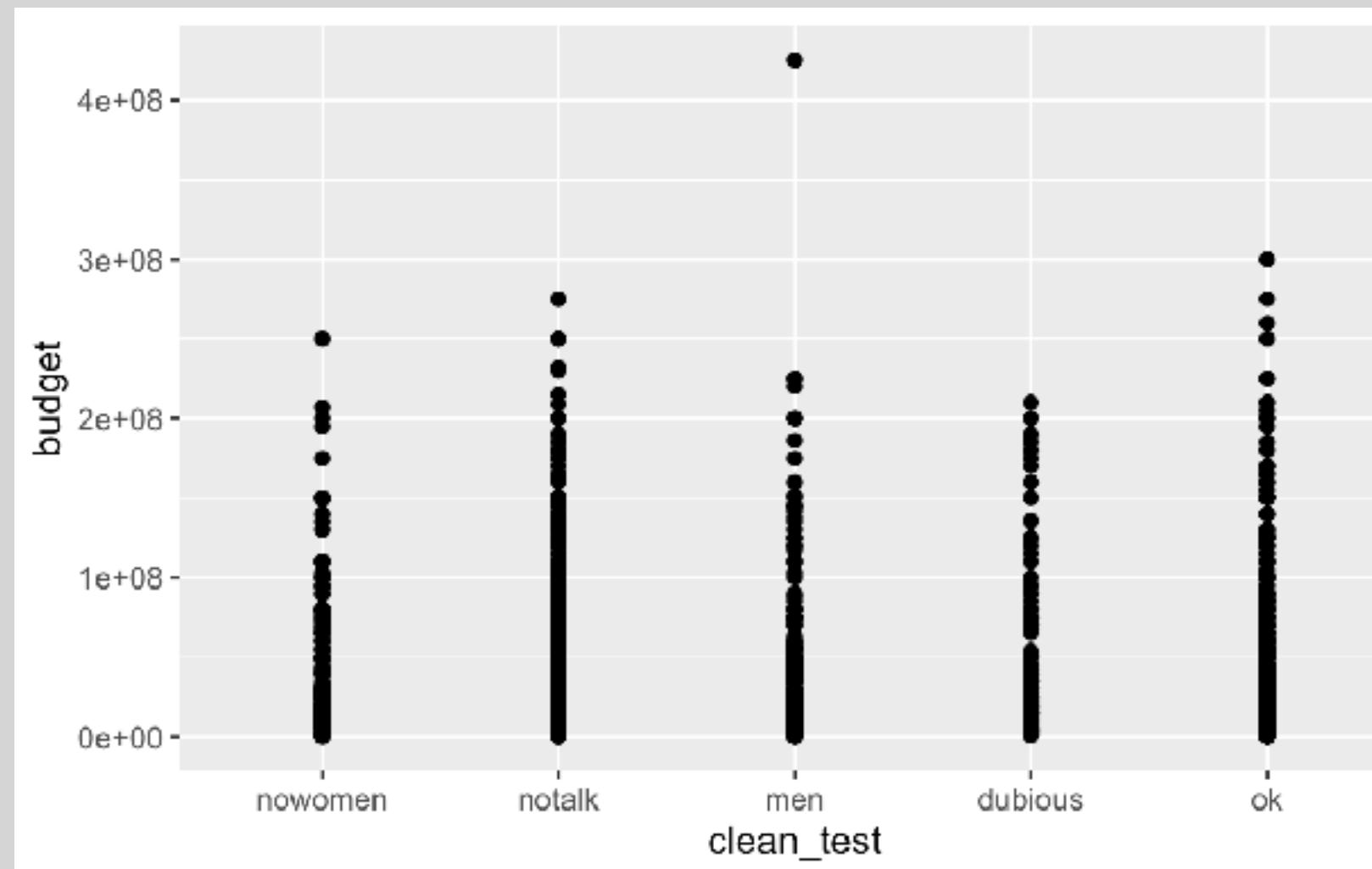
ggplot2



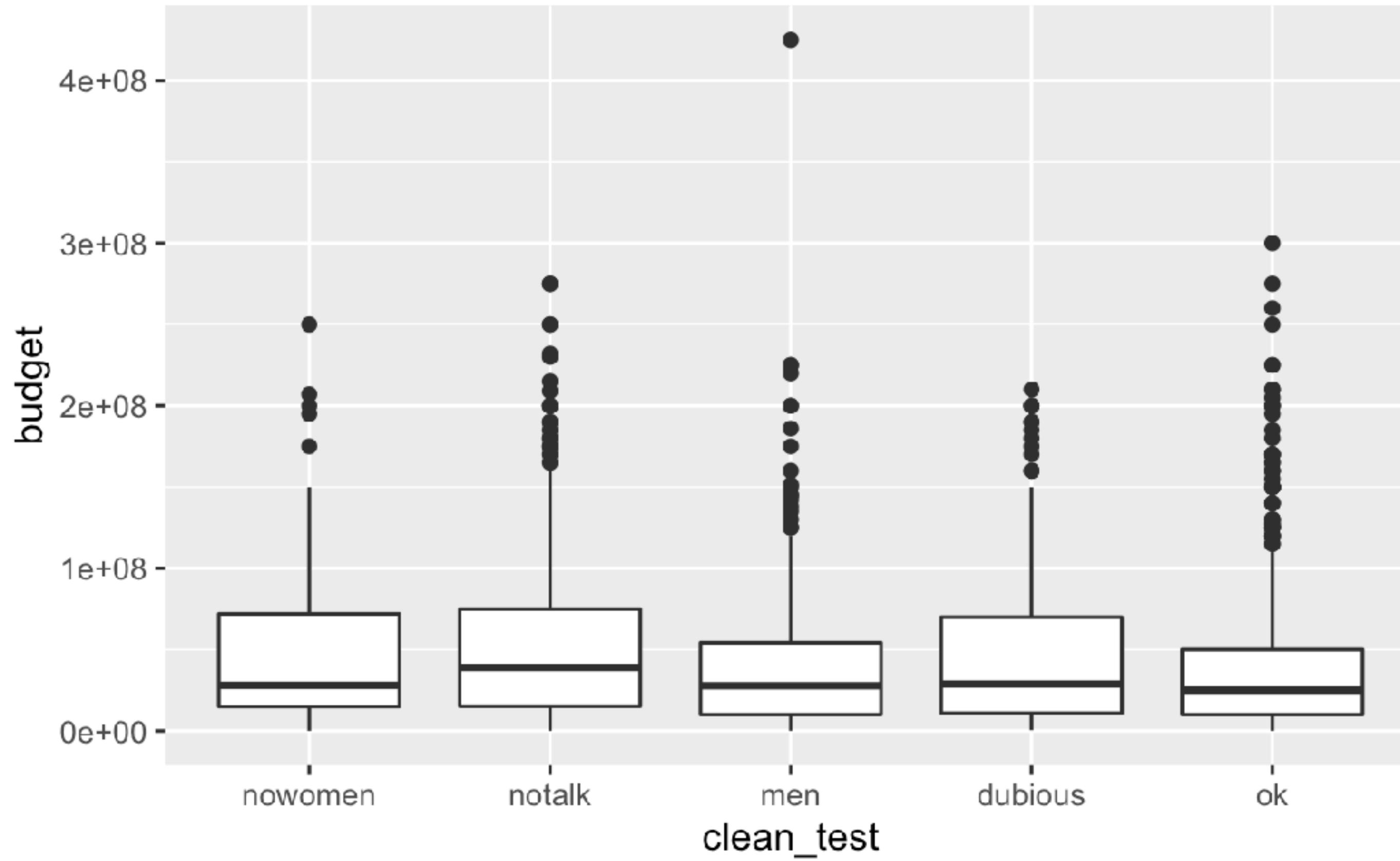
Your Turn 3

- Replace this scatterplot with one that draws boxplots

```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = clean_test, y = budget))
```



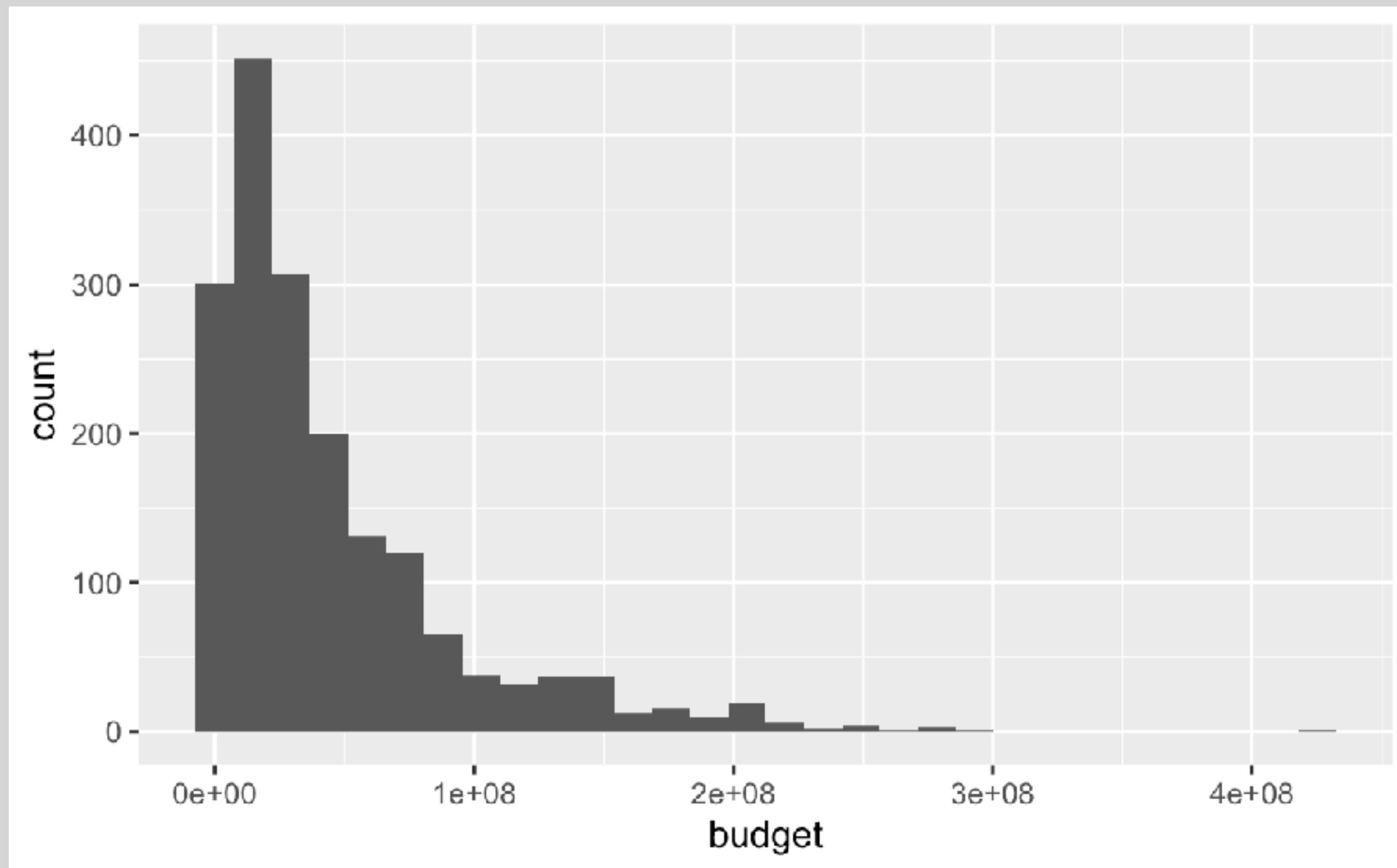
05 : 00

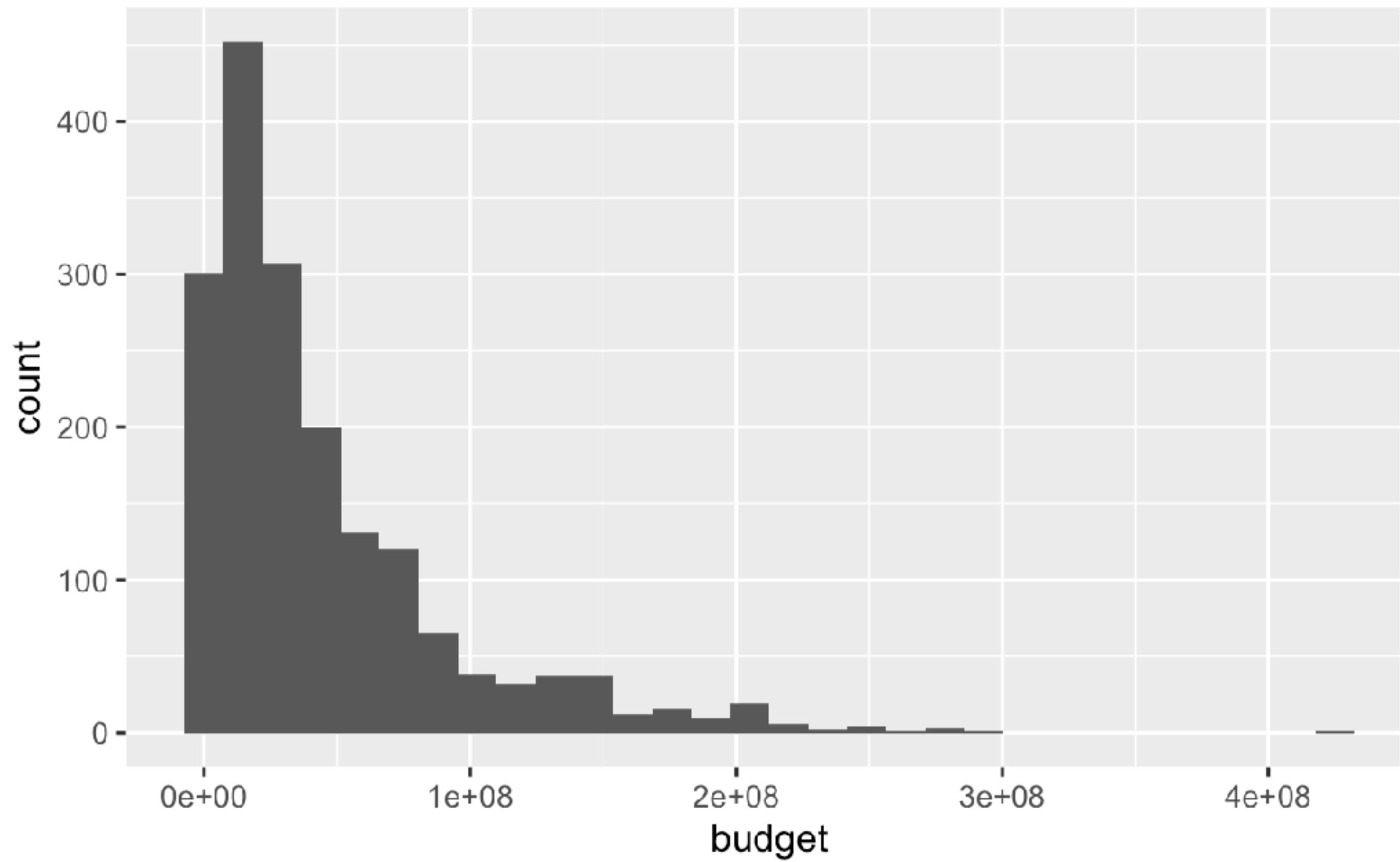


```
ggplot(data = bechdel) +  
  geom_boxplot(mapping = aes(x = clean_test, y = budget))
```

Your Turn 4

- Make the histogram of **budget** shown below

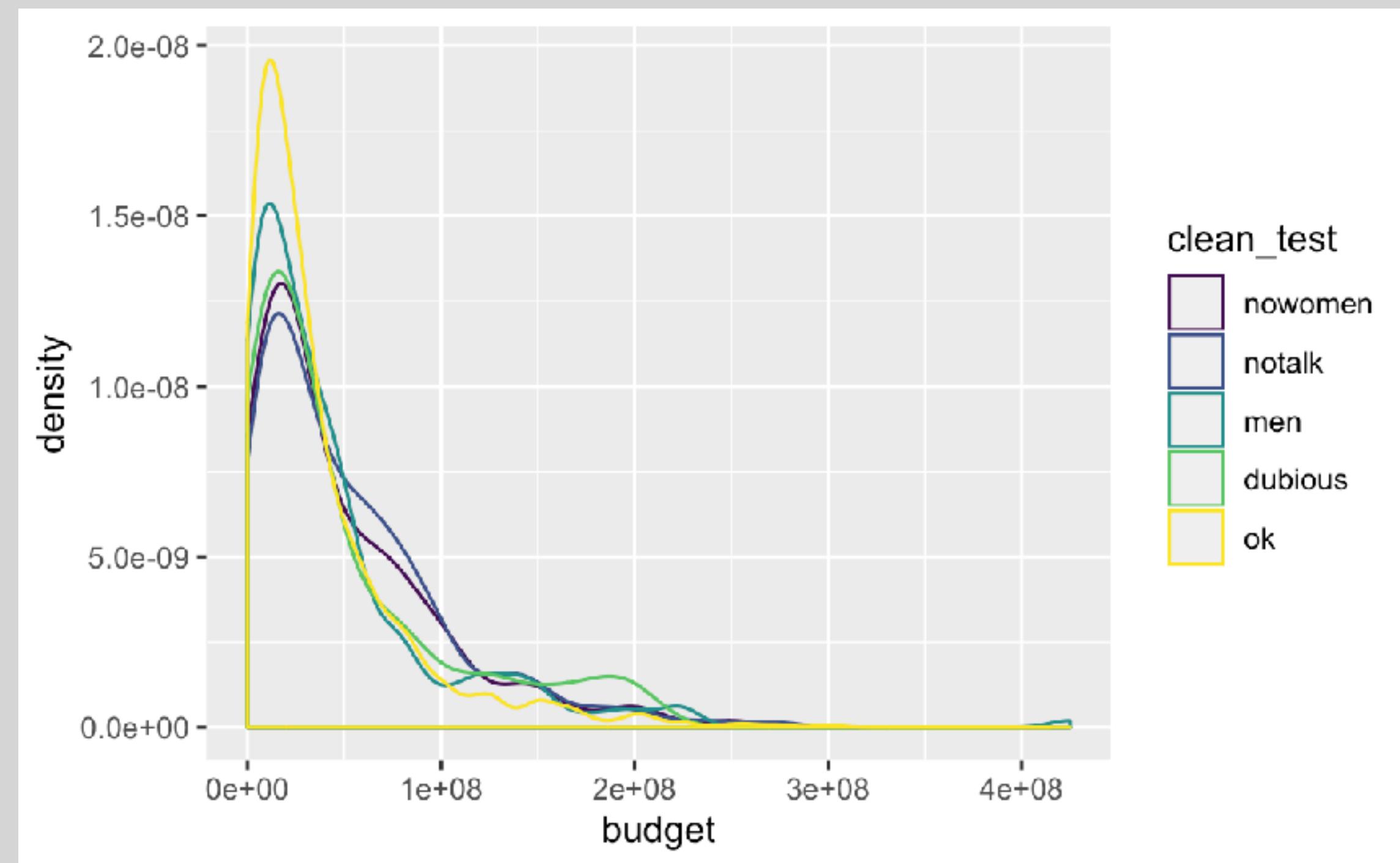




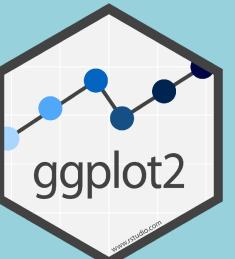
```
ggplot(data = bechdel) +  
  geom_histogram(mapping = aes(x = budget))
```

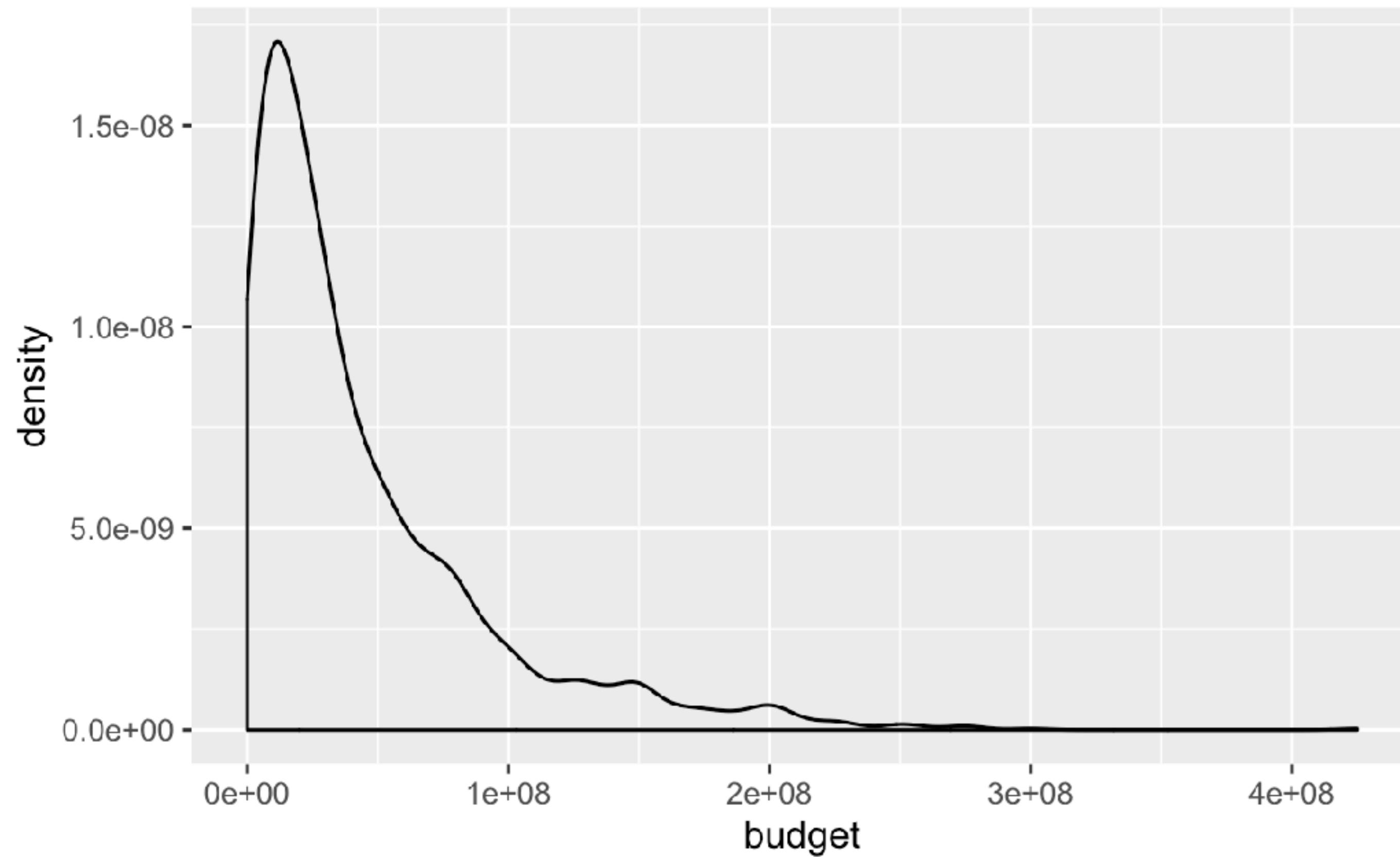
Your Turn 5

- Make the density plot of **budget** colored by **clean_test** shown below.



05 : 00

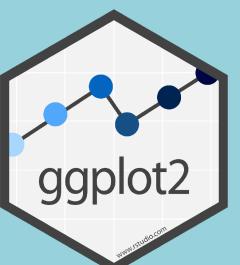


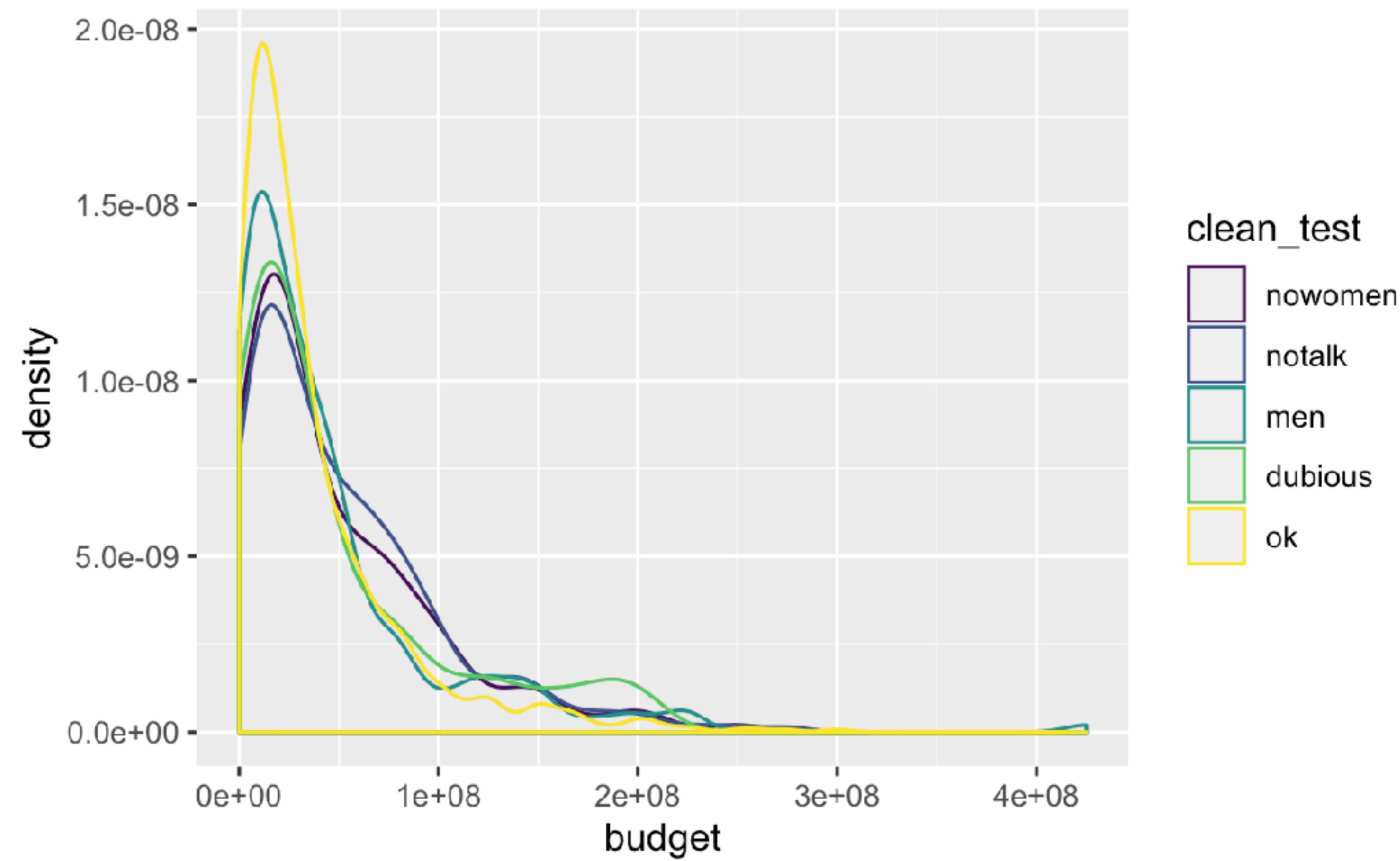


```
ggplot(data = bechdel) +  
  geom_density(mapping = aes(x = budget))
```



Adapted from Data Science in the tidyverse, CC BY Amelia McNamara

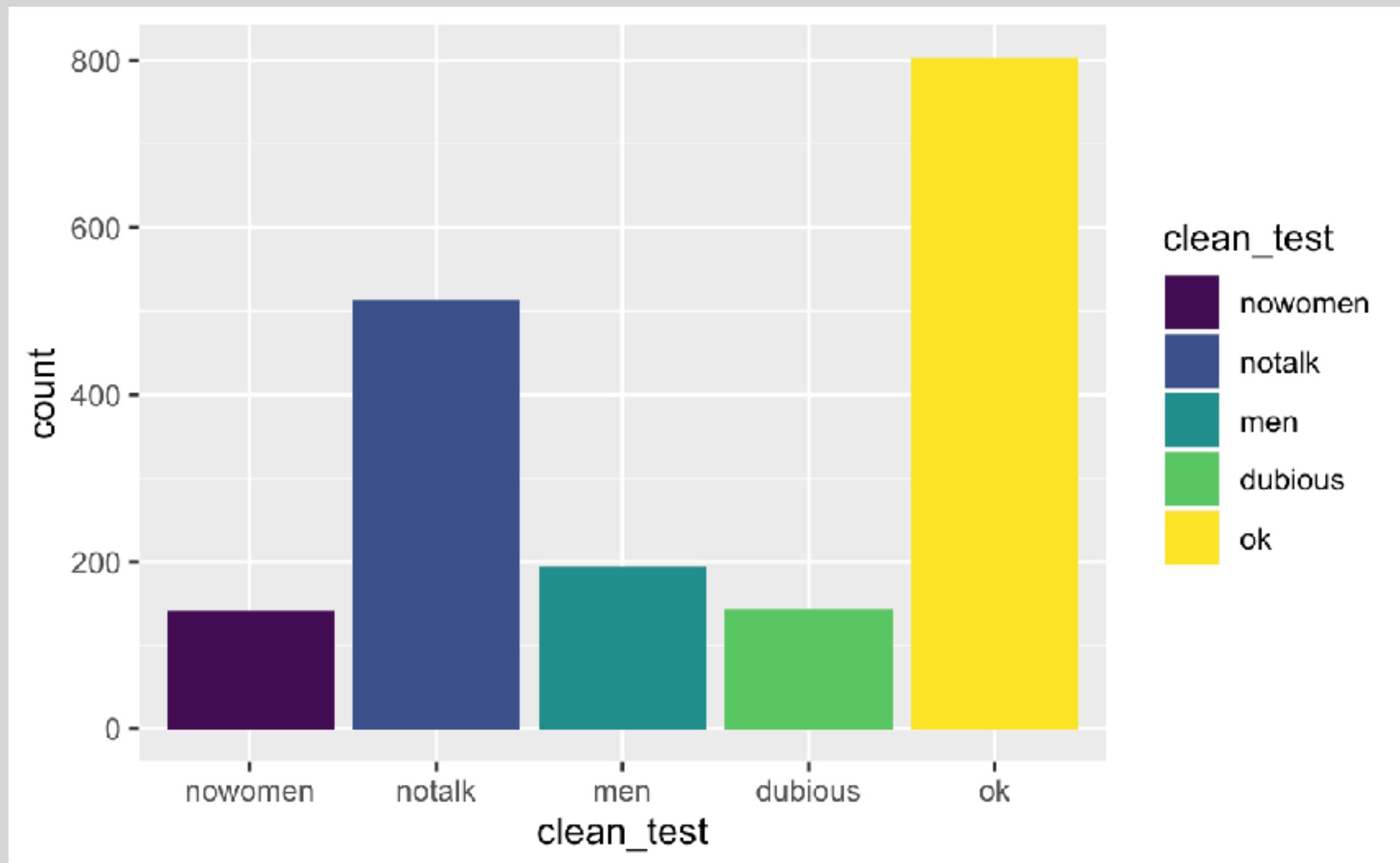




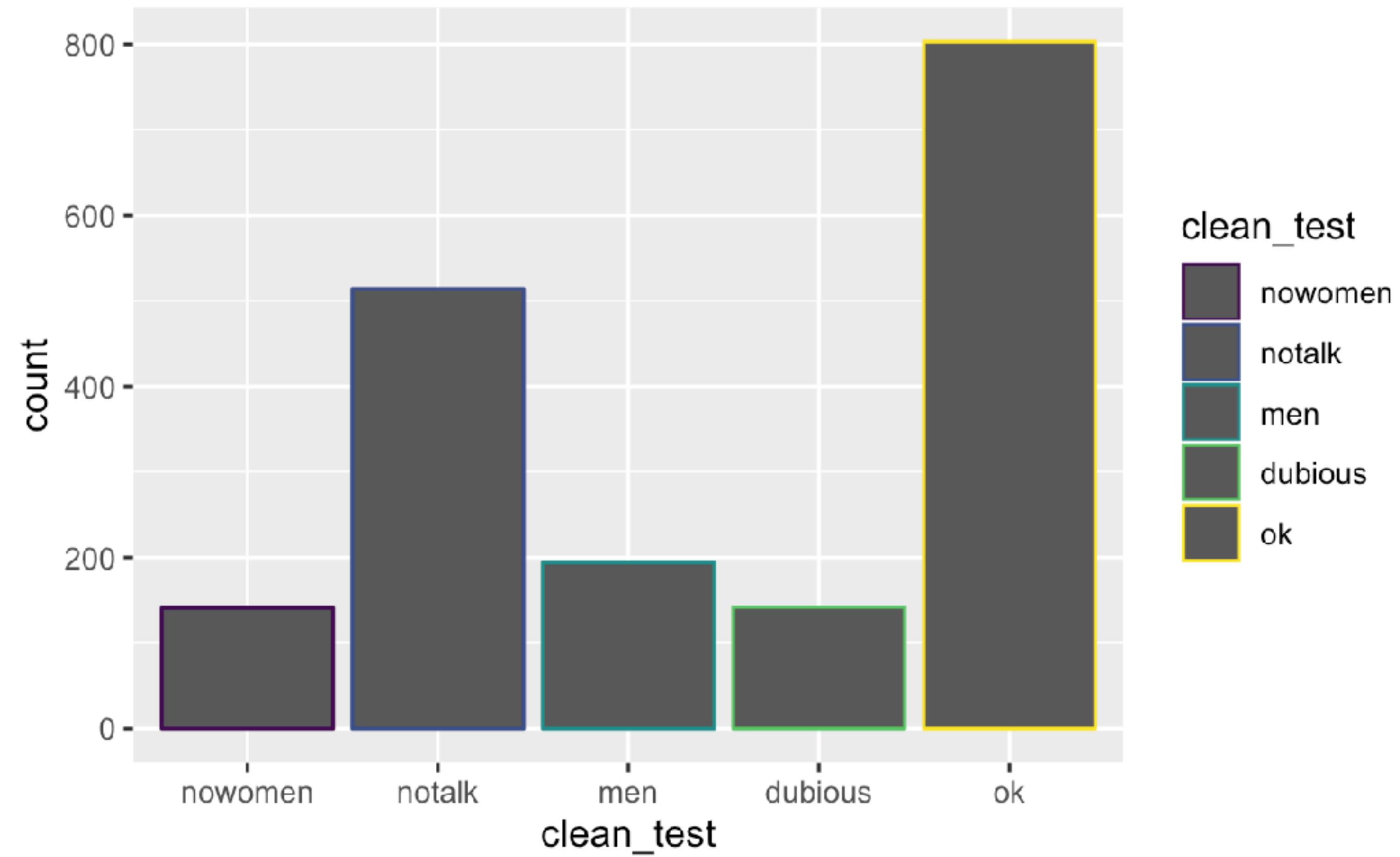
```
ggplot(data = bechdel) +  
  geom_density(mapping = aes(x = budget, color = clean_test))
```

Your Turn 6

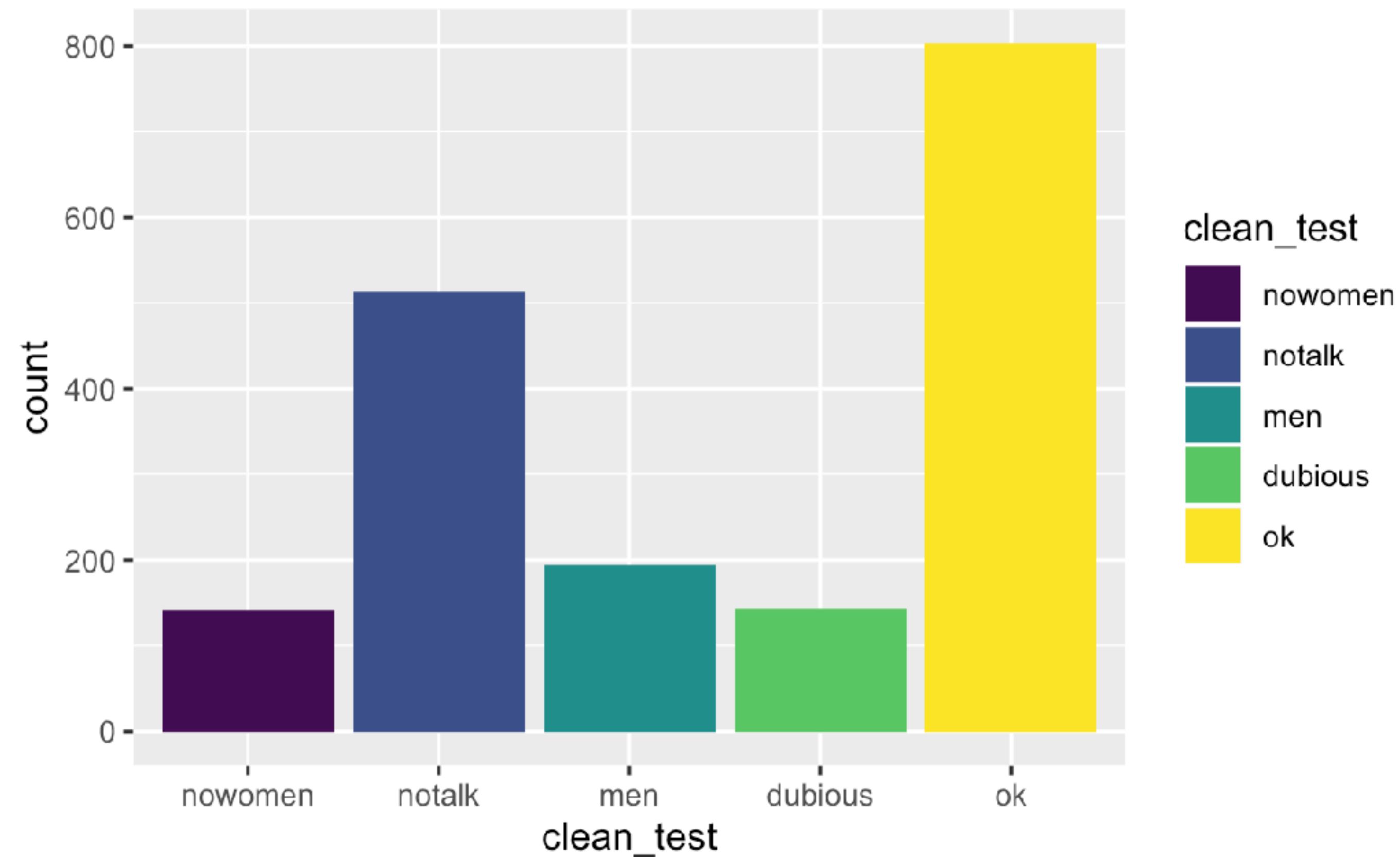
- Make the bar chart of **clean_test** colored by **clean_test** shown below.



05 : 00



```
ggplot(data = bechdel) +  
  geom_bar(mapping = aes(x = clean_test, color = clean_test))
```

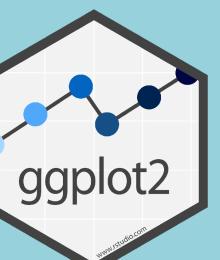


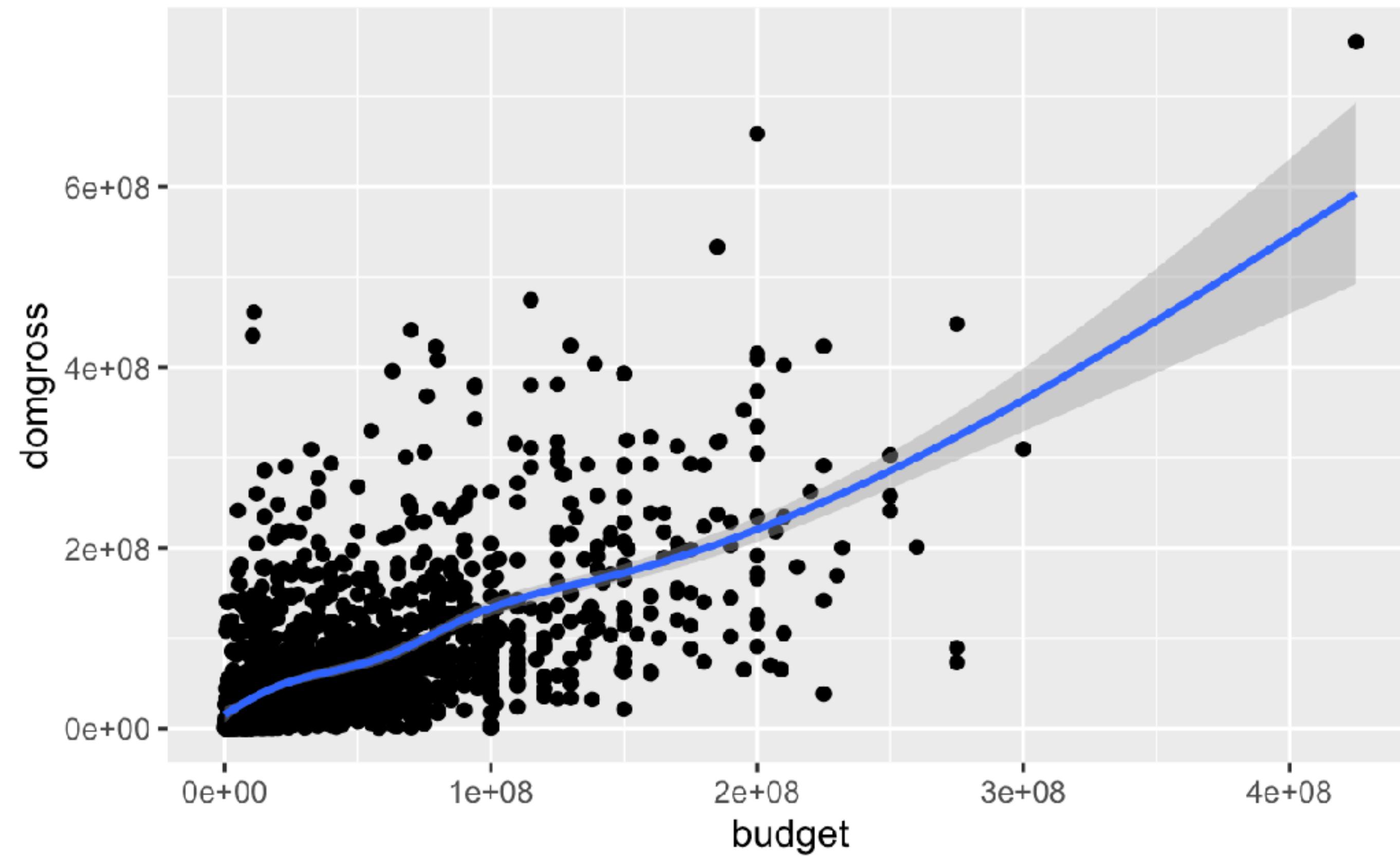
```
ggplot(data = bechdel) +  
  geom_bar(mapping = aes(x = clean_test, fill = clean_test))
```

Your Turn 7

- Predict what this code will do.
- Then run it.

```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross)) +  
  geom_smooth(mapping = aes(x = budget, y = domgross))
```



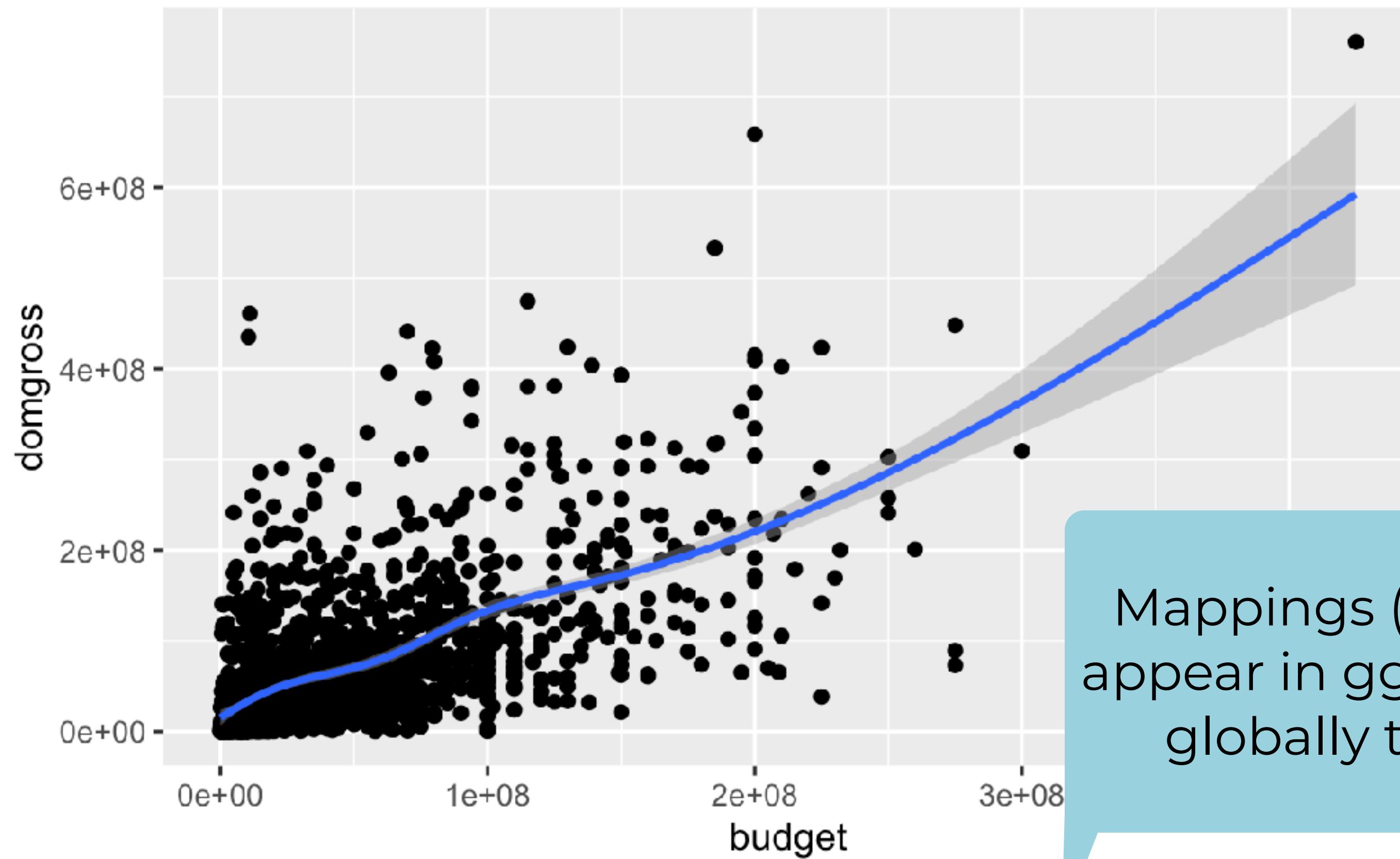


Each new
geom adds a
new layer

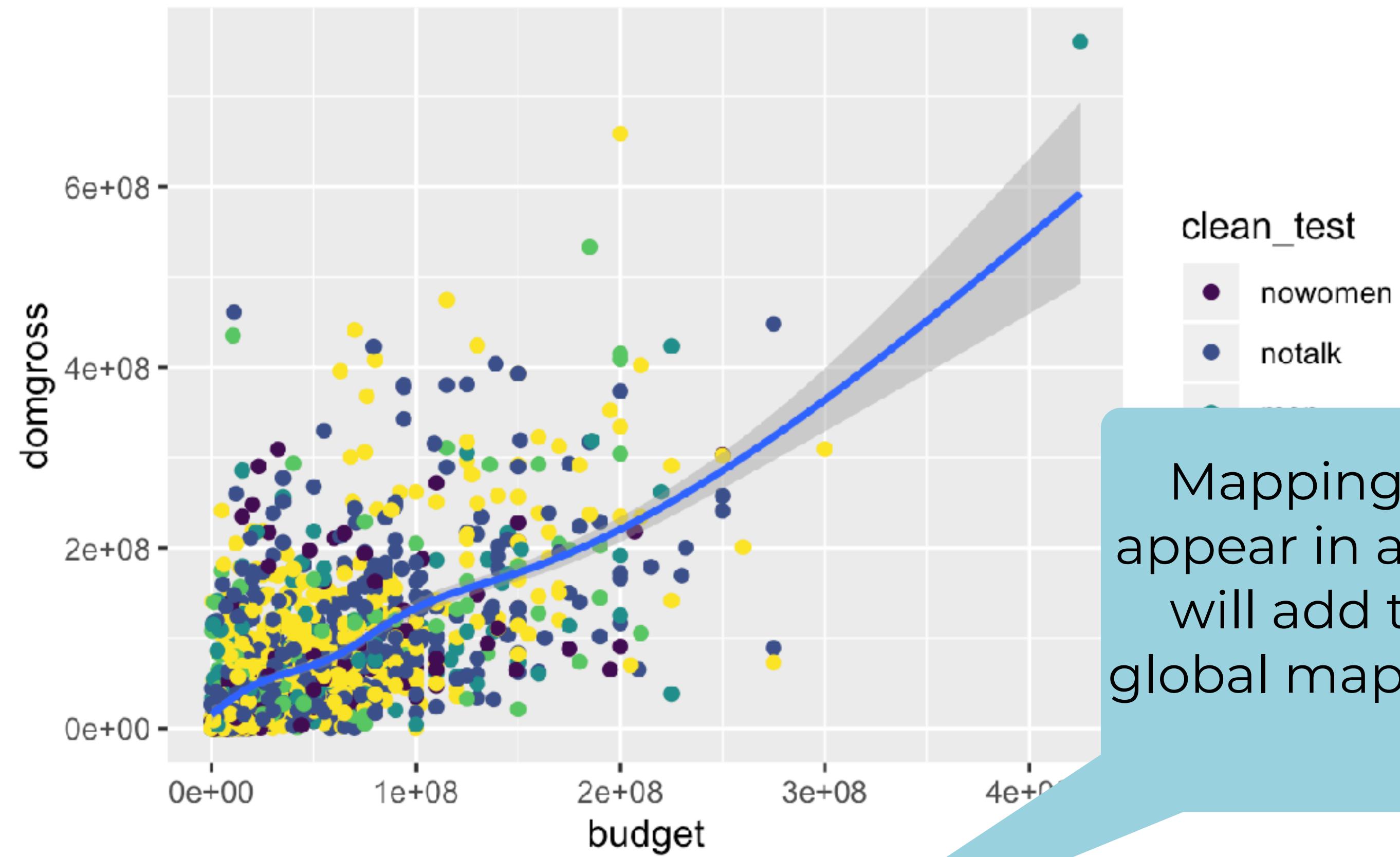
```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross)) +  
  geom_smooth(mapping = aes(x = budget, y = domgross))
```



Global vs. Local



```
ggplot(data = bechdel, mapping = aes(x = budget, y = domgross)) +  
  geom_point() +  
  geom_smooth()
```



Mappings (and data) that appear in a `geom_*`() function will add to or override the global mappings for that layer only

```
ggplot(data = bechdel, mapping = aes(x = budget, y = domgross)) +  
  geom_point(mapping = aes(color = clean_test)) +  
  geom_smooth()
```

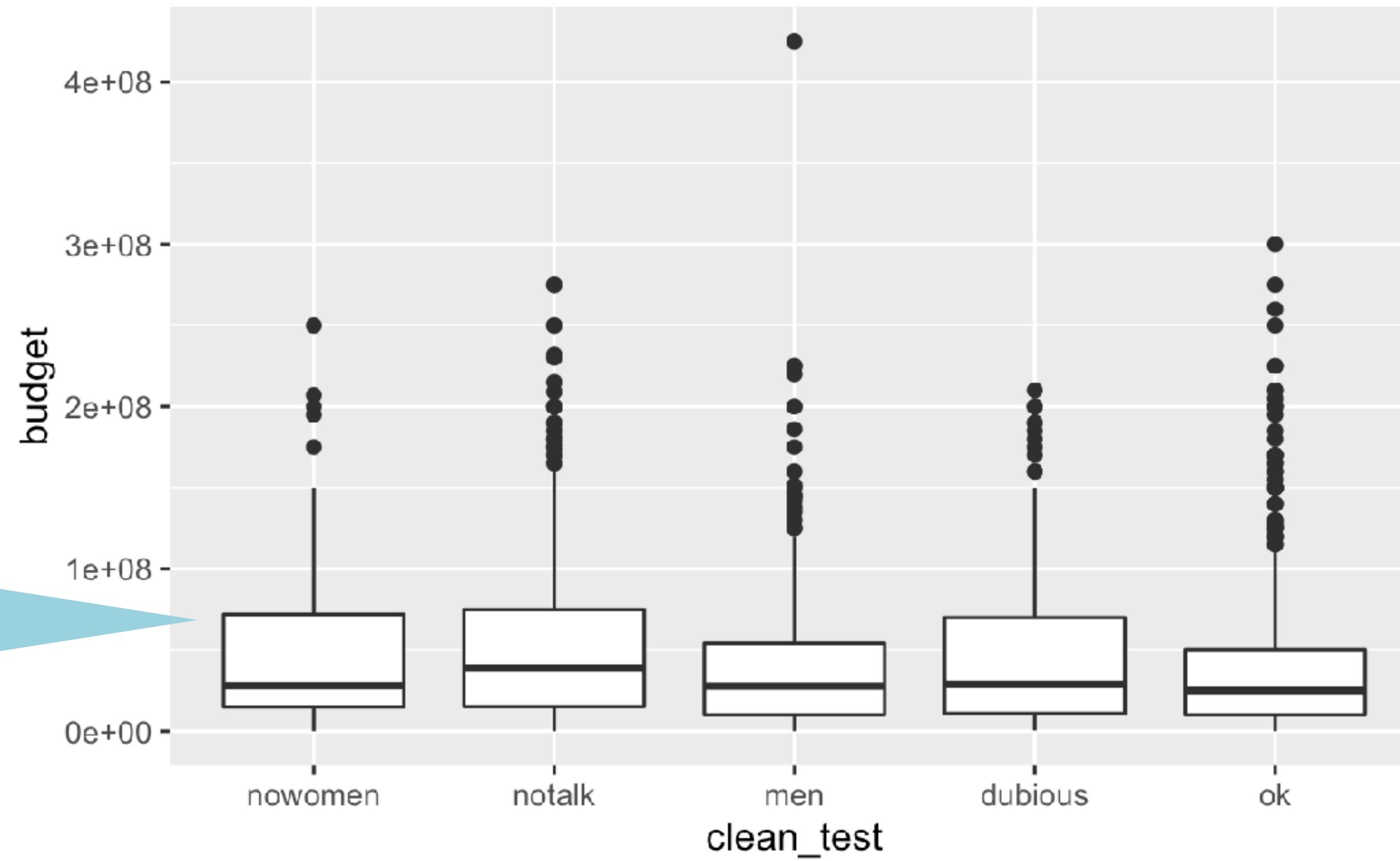


```
ggplot(data = bechdel, mapping = aes(x = budget, y = domgross)) +  
  geom_point(mapping = aes(color = clean_test)) +  
  geom_smooth(data = filter(bechdel, clean_test == "ok"))
```

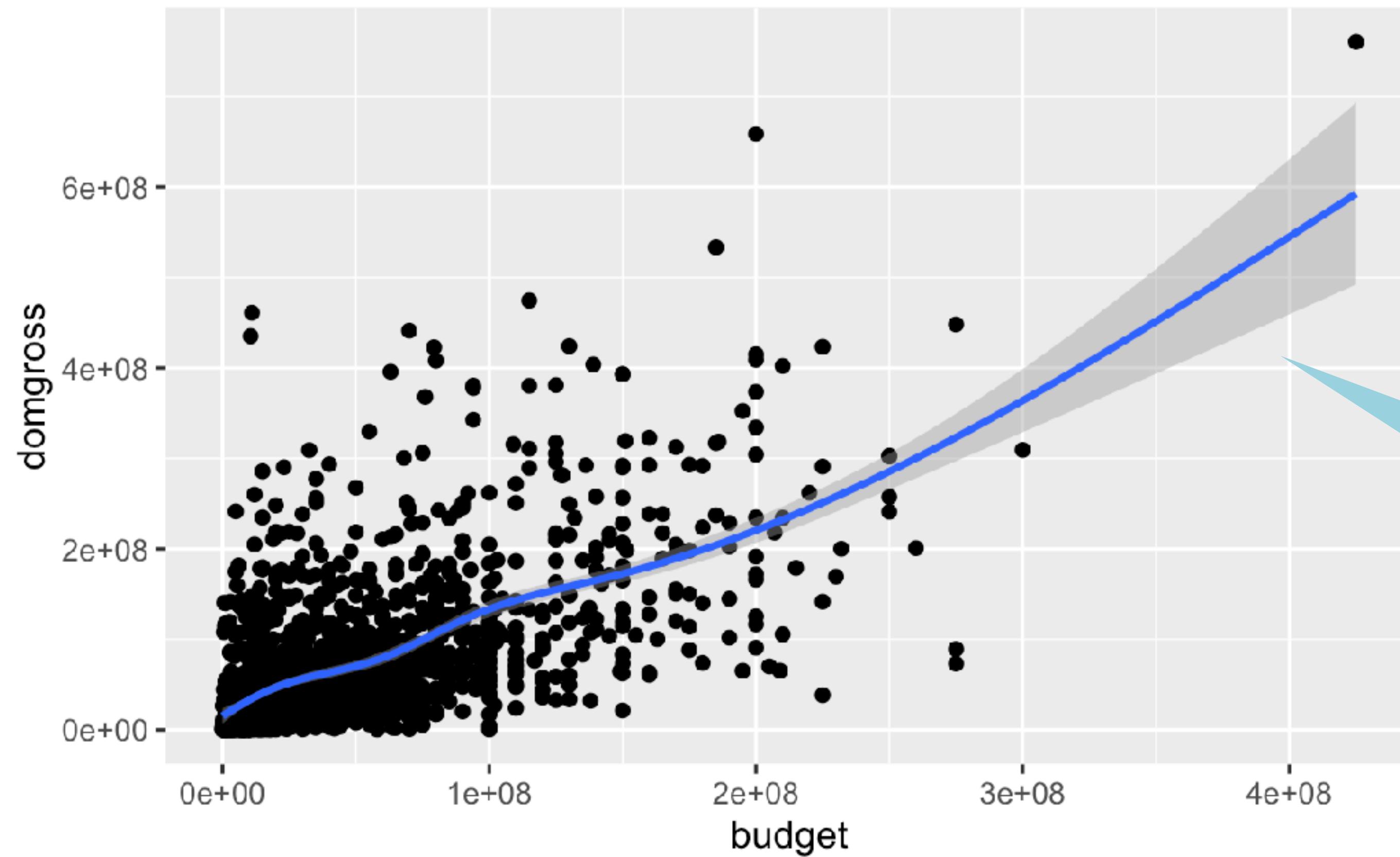


Stats

Where do
these values
come from?

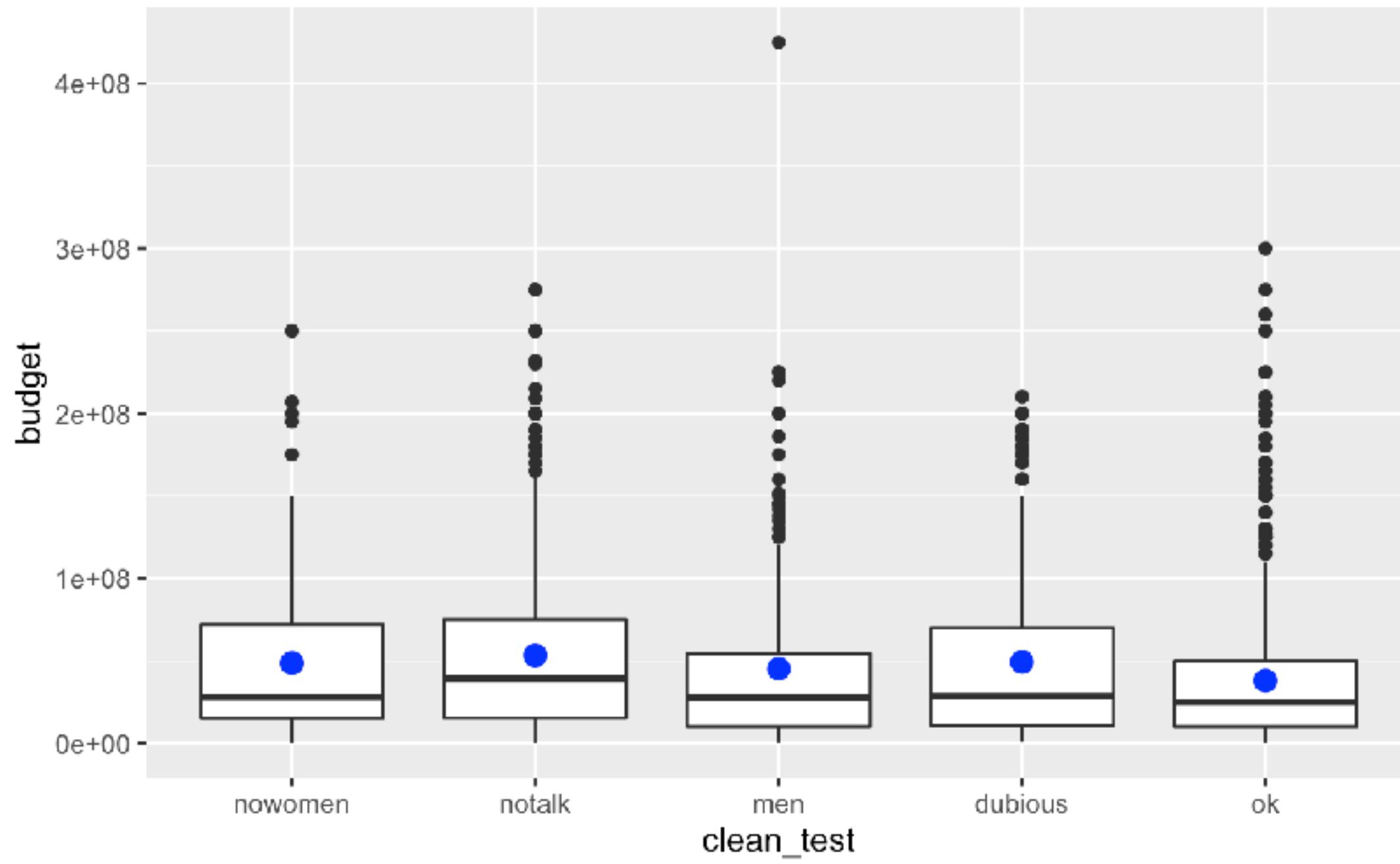


```
ggplot(data = bechdel) +  
  geom_boxplot(mapping = aes(x = clean_test, y = budget))
```



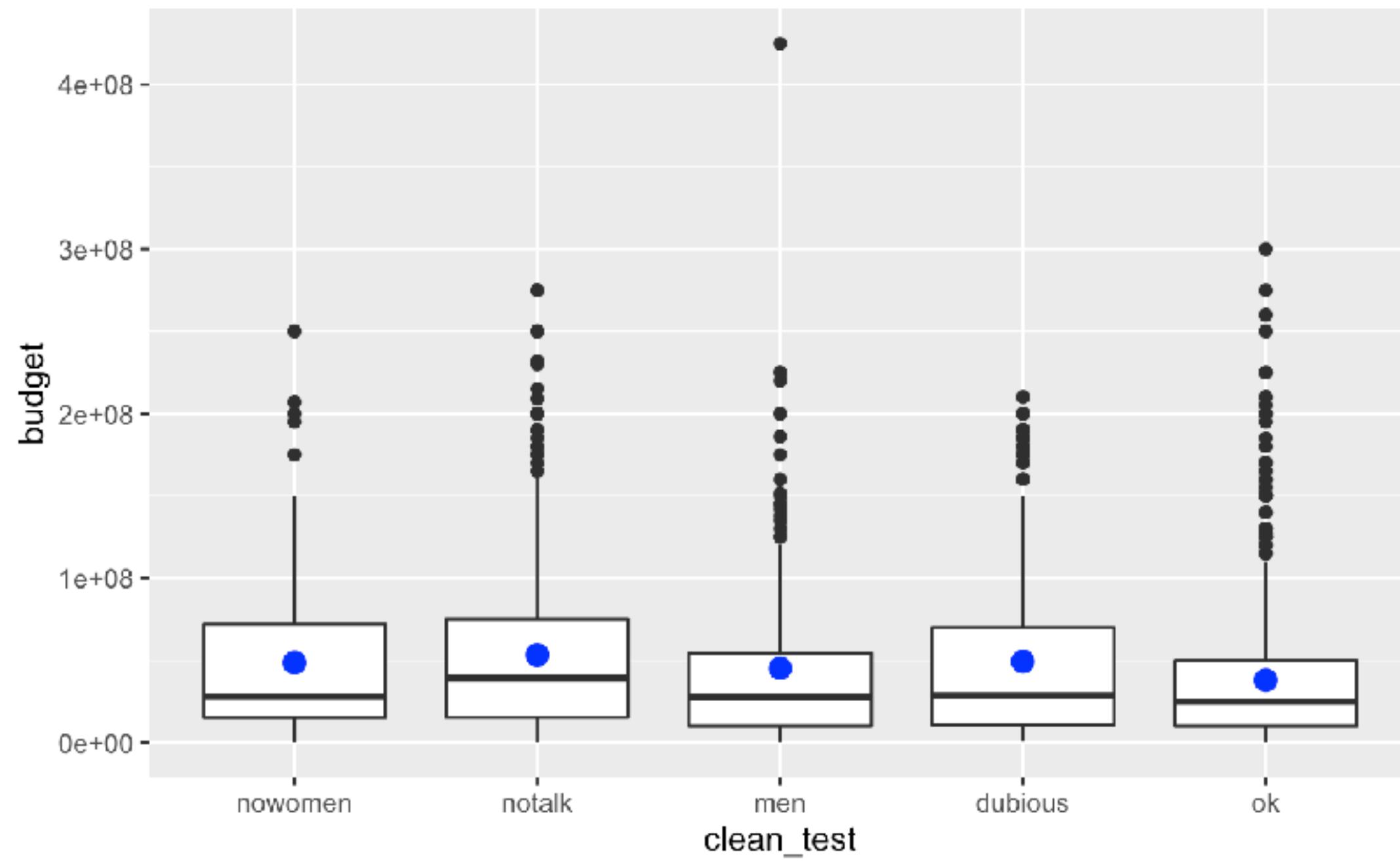
```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross)) +  
  geom_smooth(mapping = aes(x = budget, y = domgross))
```

Stats as layers



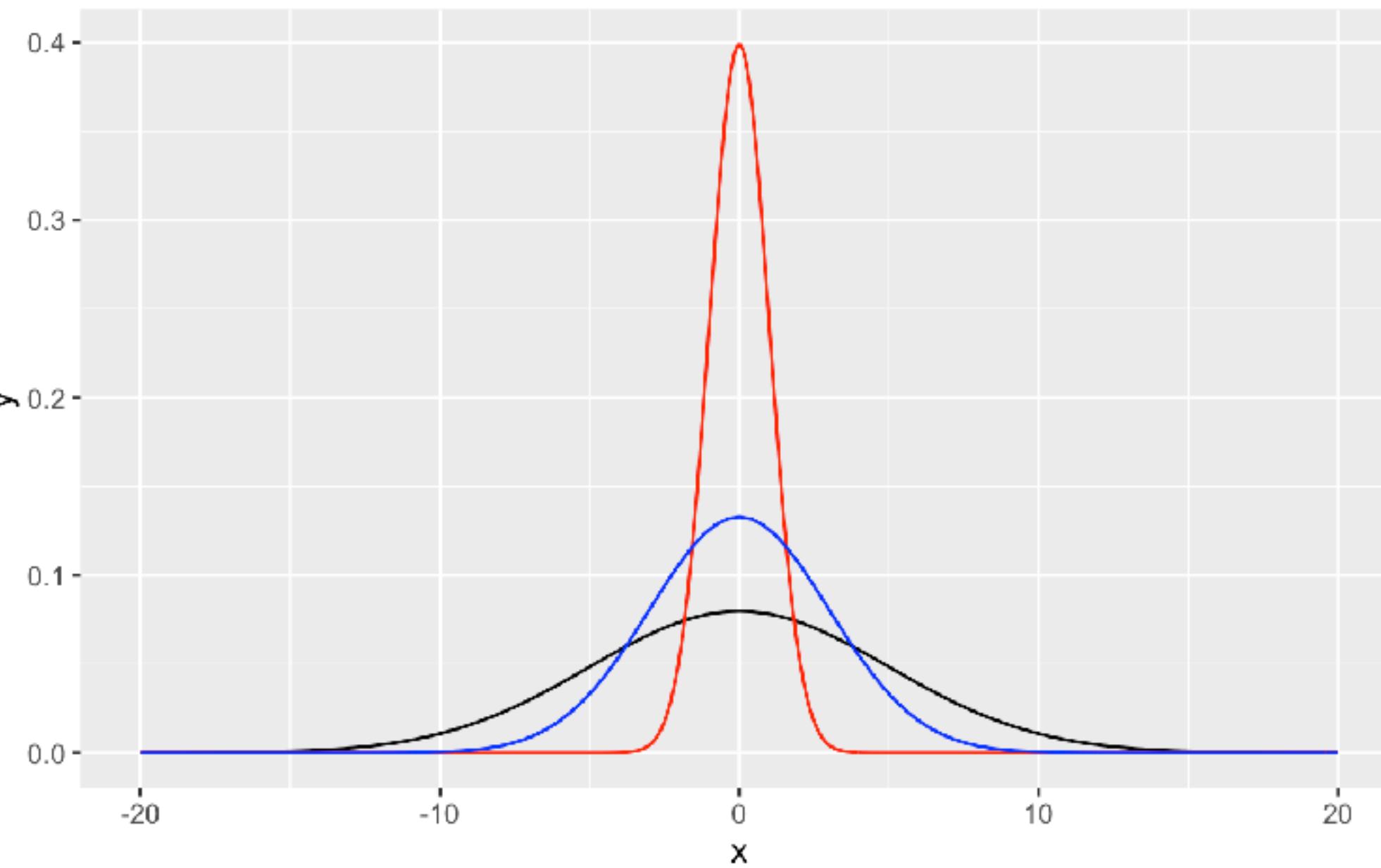
```
ggplot(data = bechdel, mapping = aes(x = clean_test, y = budget)) +  
  geom_boxplot() +  
  stat_summary(geom = "point", y.fun = "mean", color = "blue", size = 3)
```

Stats as layers



```
ggplot(data = bechdel, mapping = aes(x = clean_test, y = budget)) +  
  geom_boxplot() +  
  geom_point(stat = "summary", y.fun = "mean", color = "blue", size = 3)
```

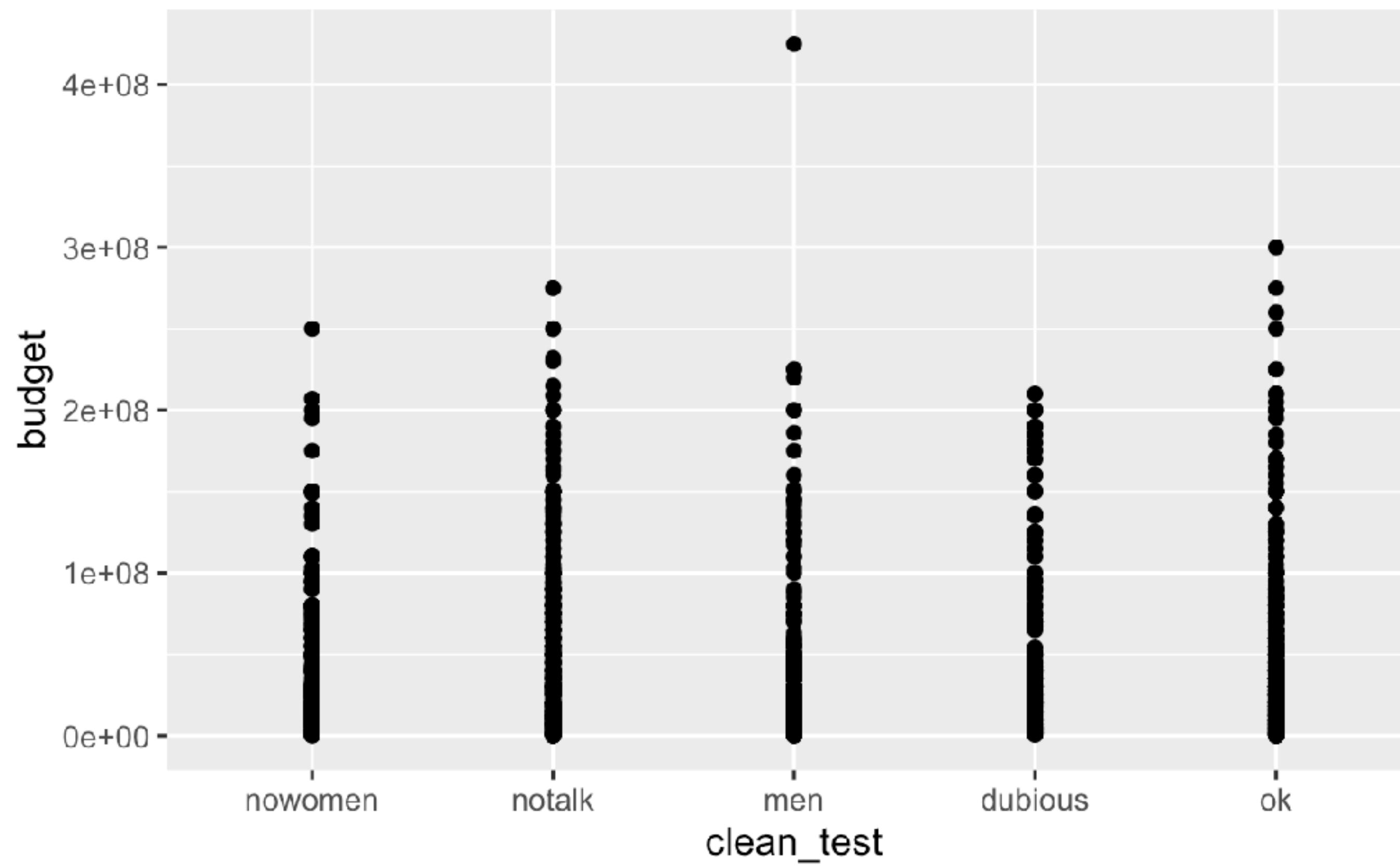
Distributions



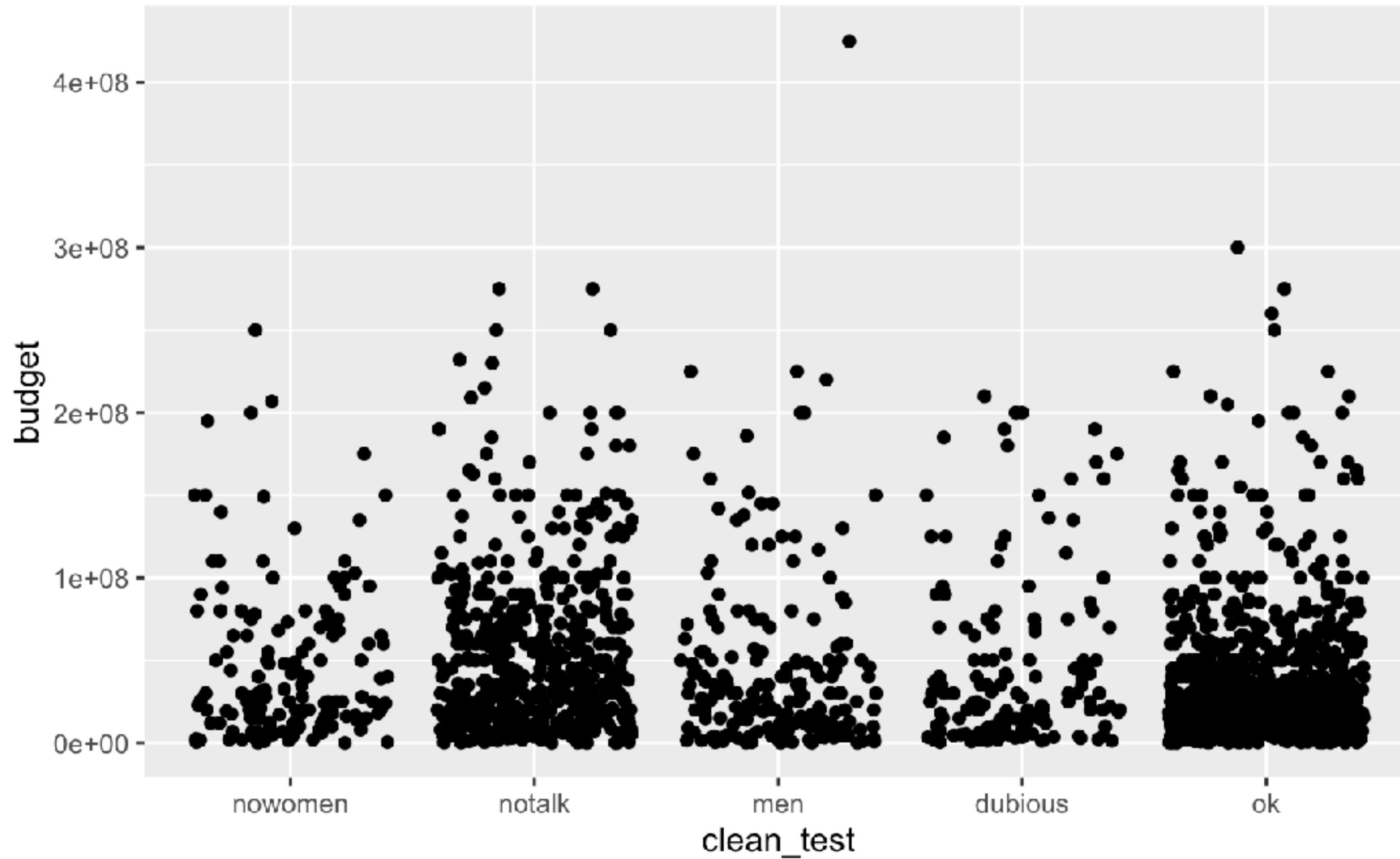
```
ggplot(data = tibble(x = c(-20, 20)), aes(x = x)) +  
  stat_function(fun = dnorm, args = list(mean = 0, sd = 5), color = "black") +  
  stat_function(fun = dnorm, args = list(mean = 0, sd = 1), color = "red") +  
  stat_function(fun = dnorm, args = list(mean = 0, sd = 3), color = "blue")
```



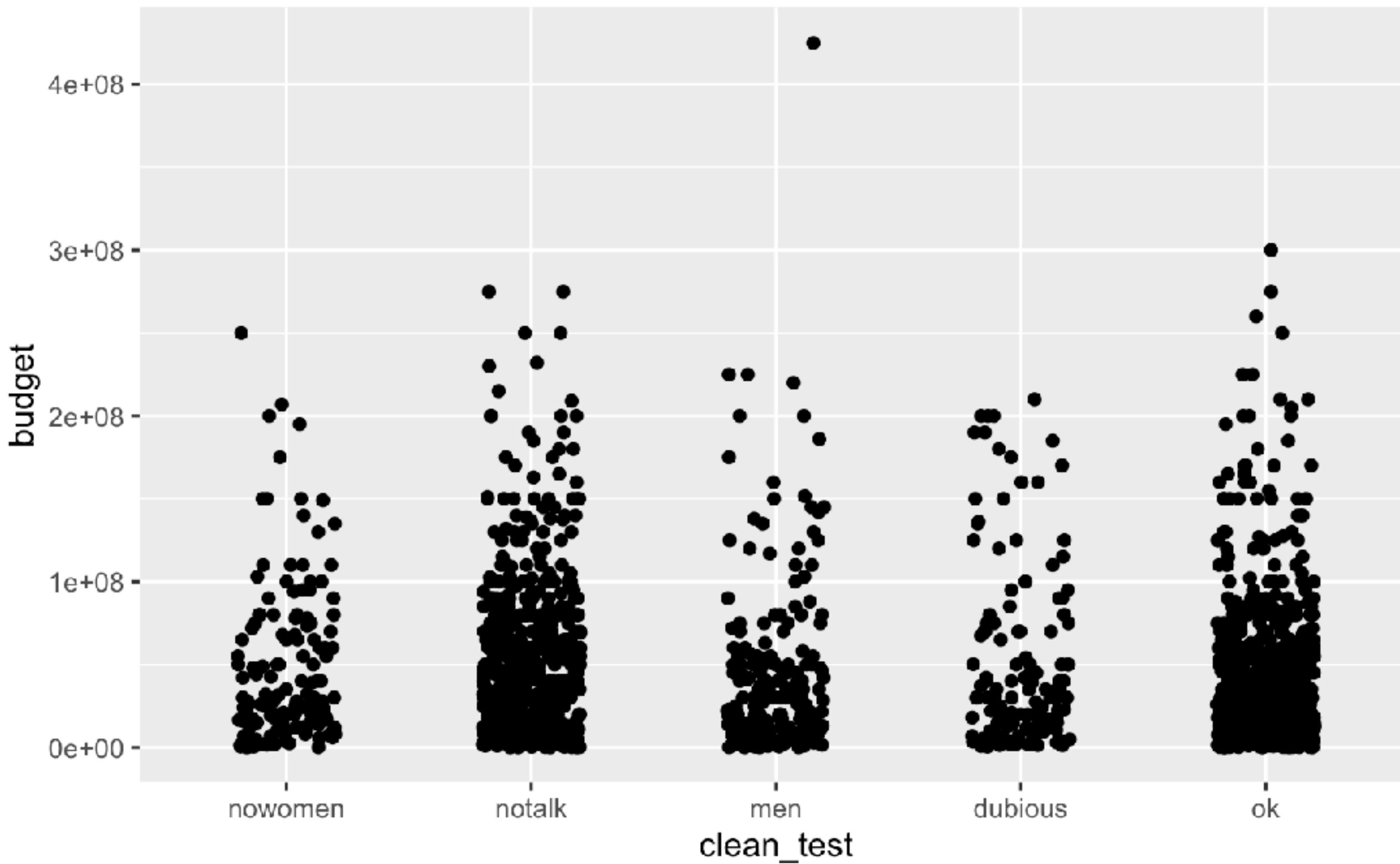
Position



```
ggplot(bechdel, aes(x = clean_test, y = budget)) +  
  geom_point()
```



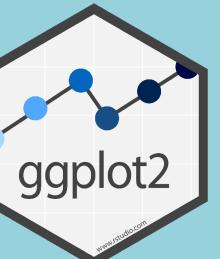
```
ggplot(bechdel, aes(x = clean_test, y = budget)) +  
  geom_point(position = "jitter")
```

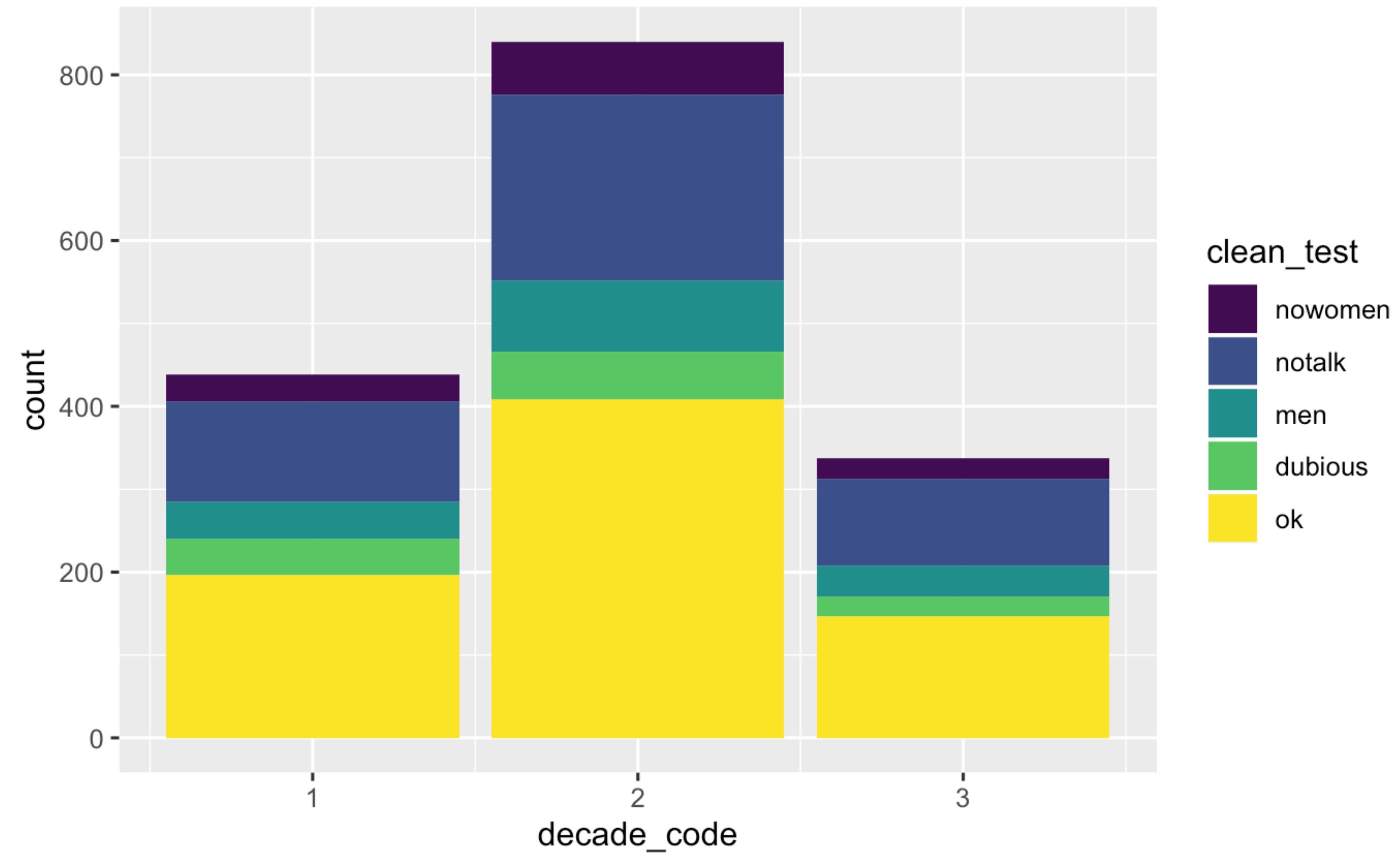


```
ggplot(bechdel, aes(x = clean_test, y = budget)) +  
  geom_point(position = position_jitter(width = 0.2, height = 0))
```

Positions

- `position_identity()`
- `position_jitter()`
- `position_dodge()`
- `position_fill()`
- And more...

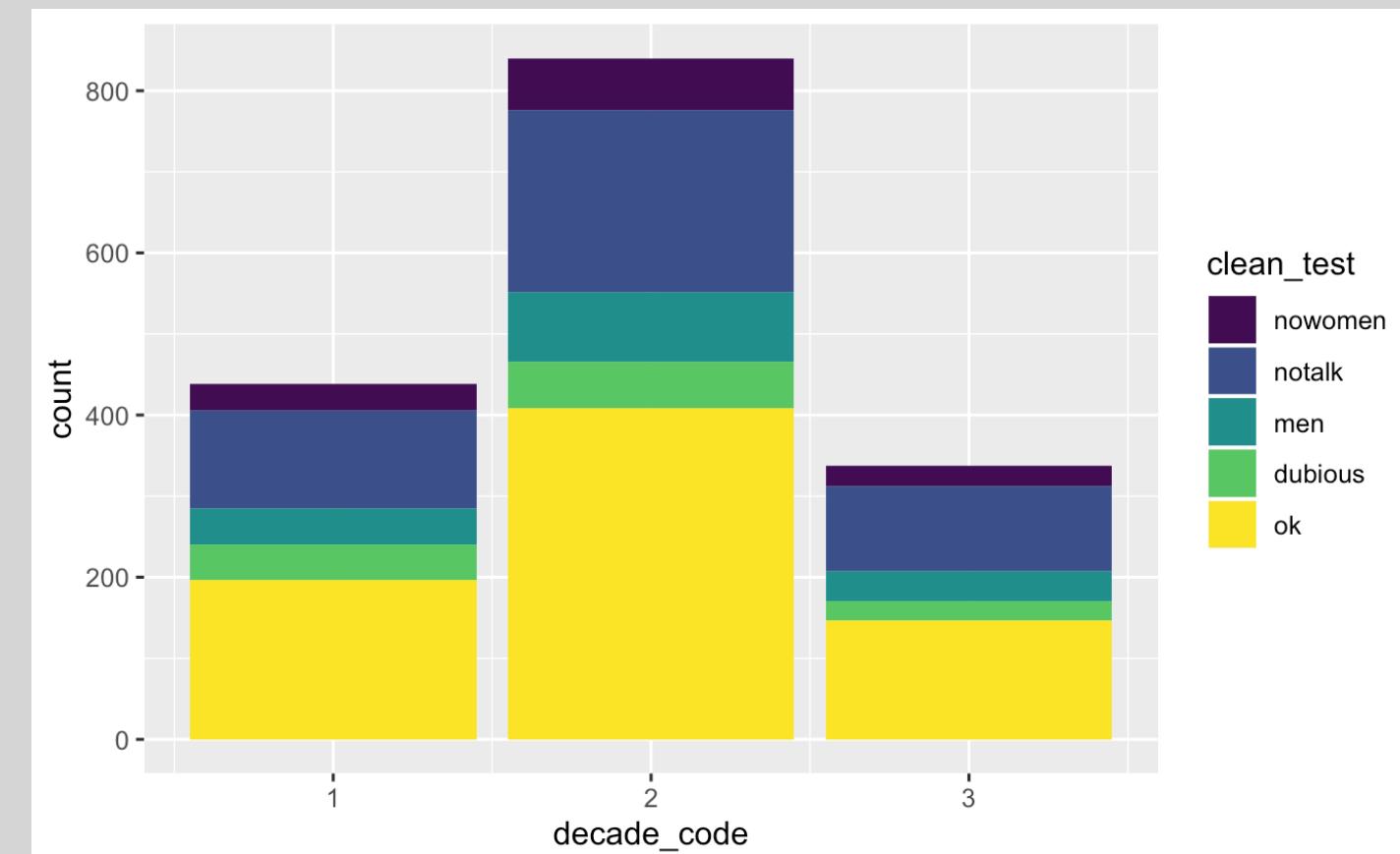


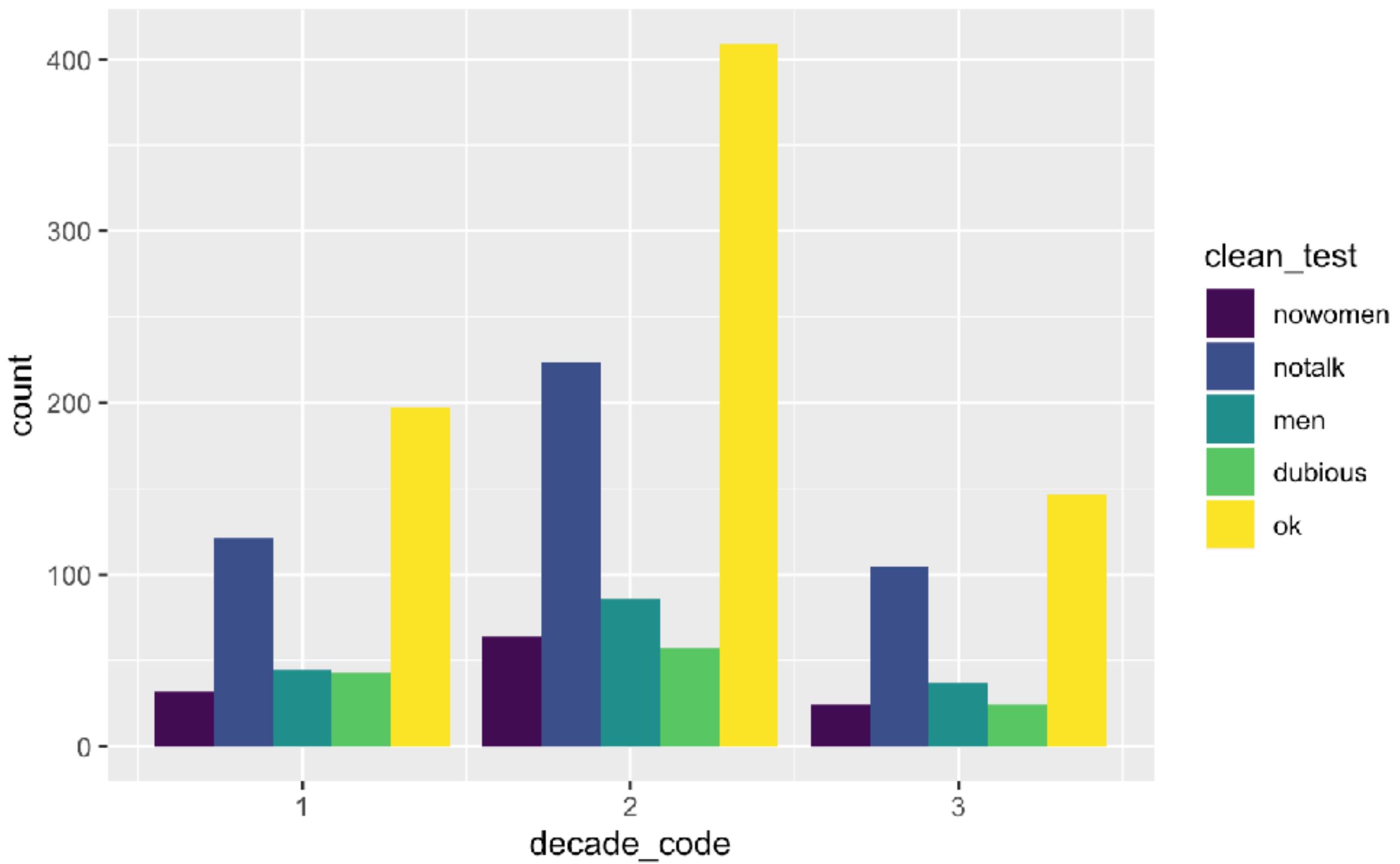


Your Turn 8

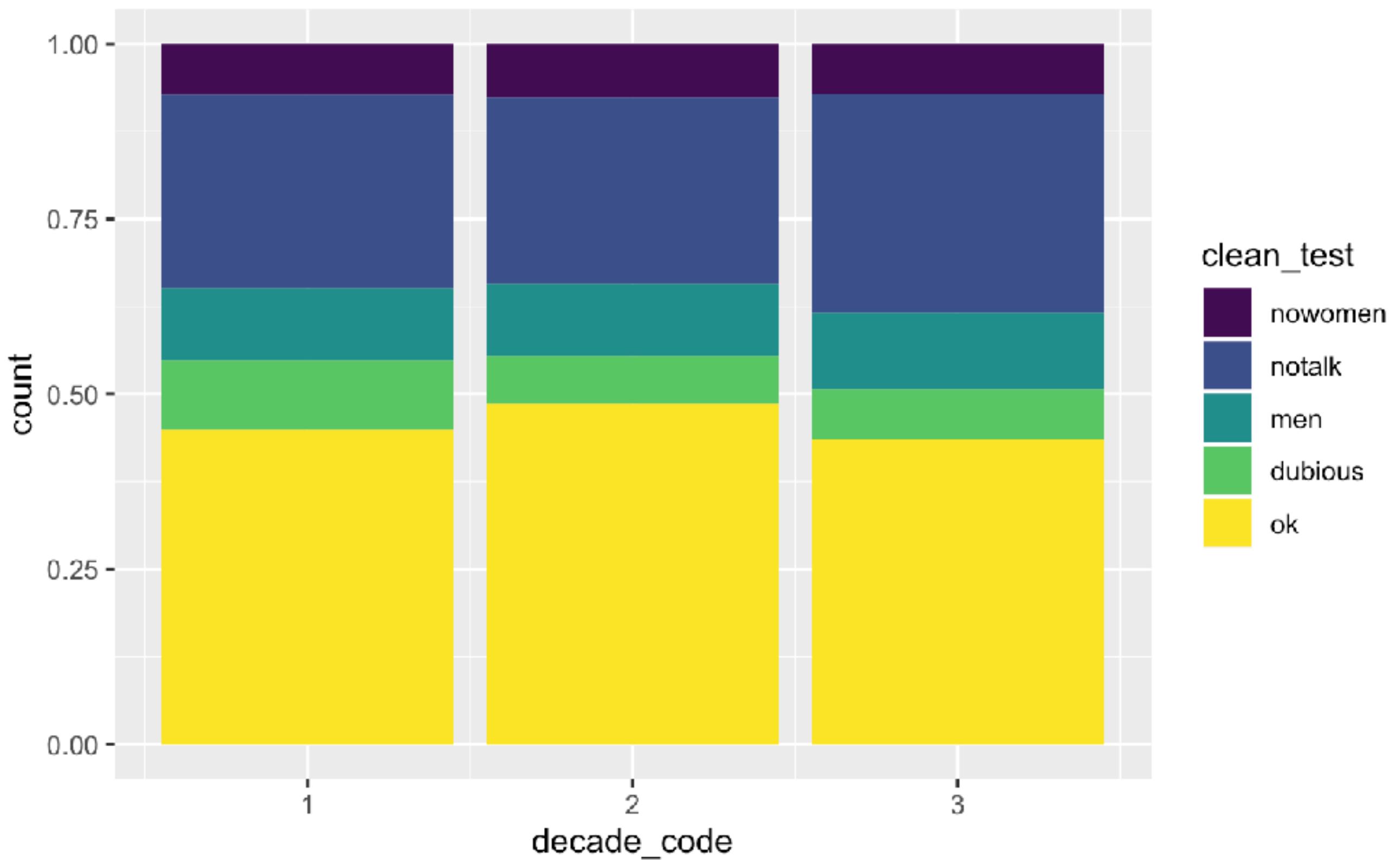
- Add a position adjustment to this plot to compare the frequency of test results across decades.

```
ggplot(data = bechdel, mapping = aes(x = decade_code)) +  
  geom_bar(mapping = aes(fill = clean_test))
```





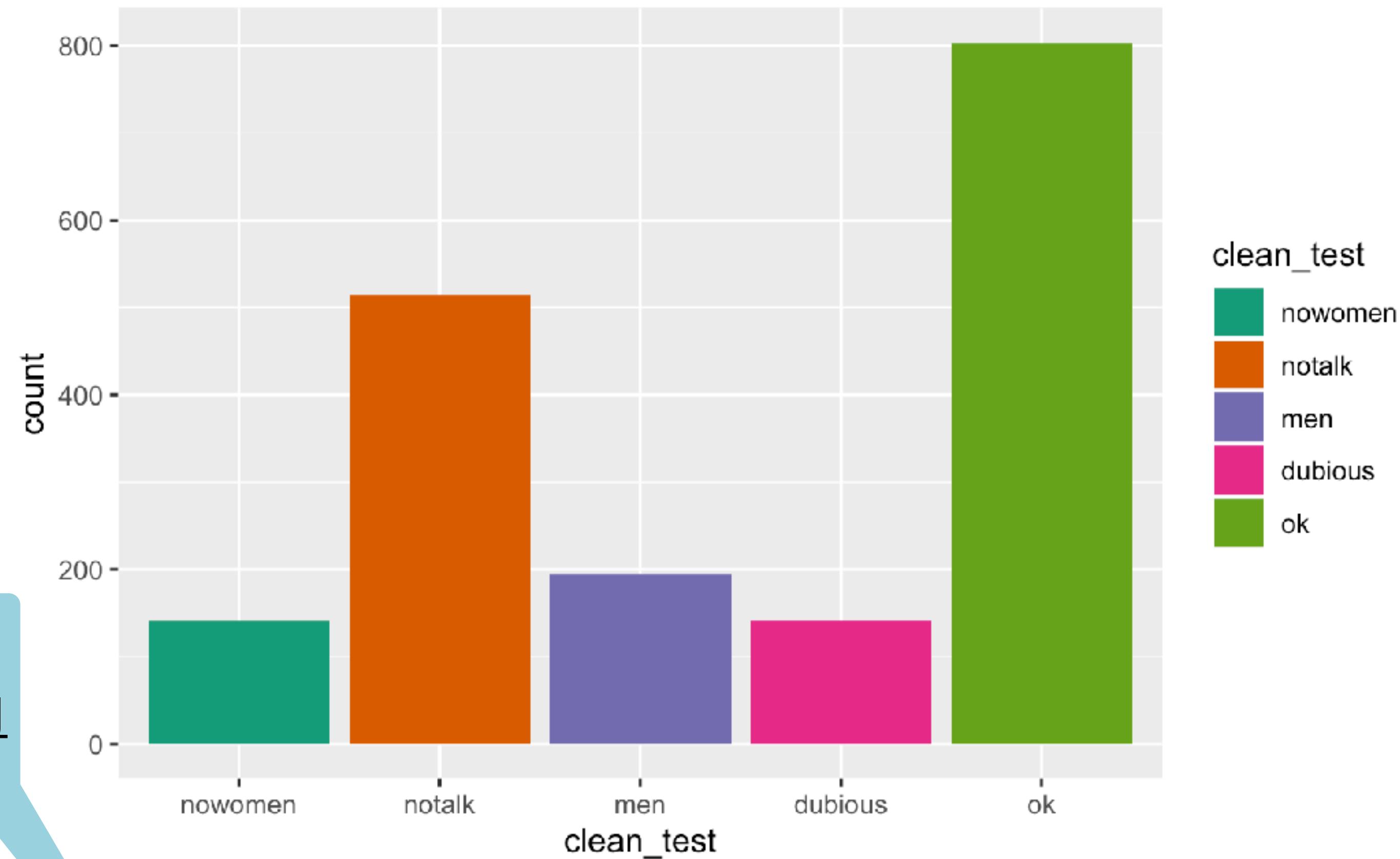
```
ggplot(bechdel, aes(x = decade_code)) +  
  geom_bar(aes(fill = clean_test), position = "dodge")
```



```
ggplot(bechdel, aes(x = decade_code)) +  
  geom_bar(aes(fill = clean_test), position = "fill")
```



Scales



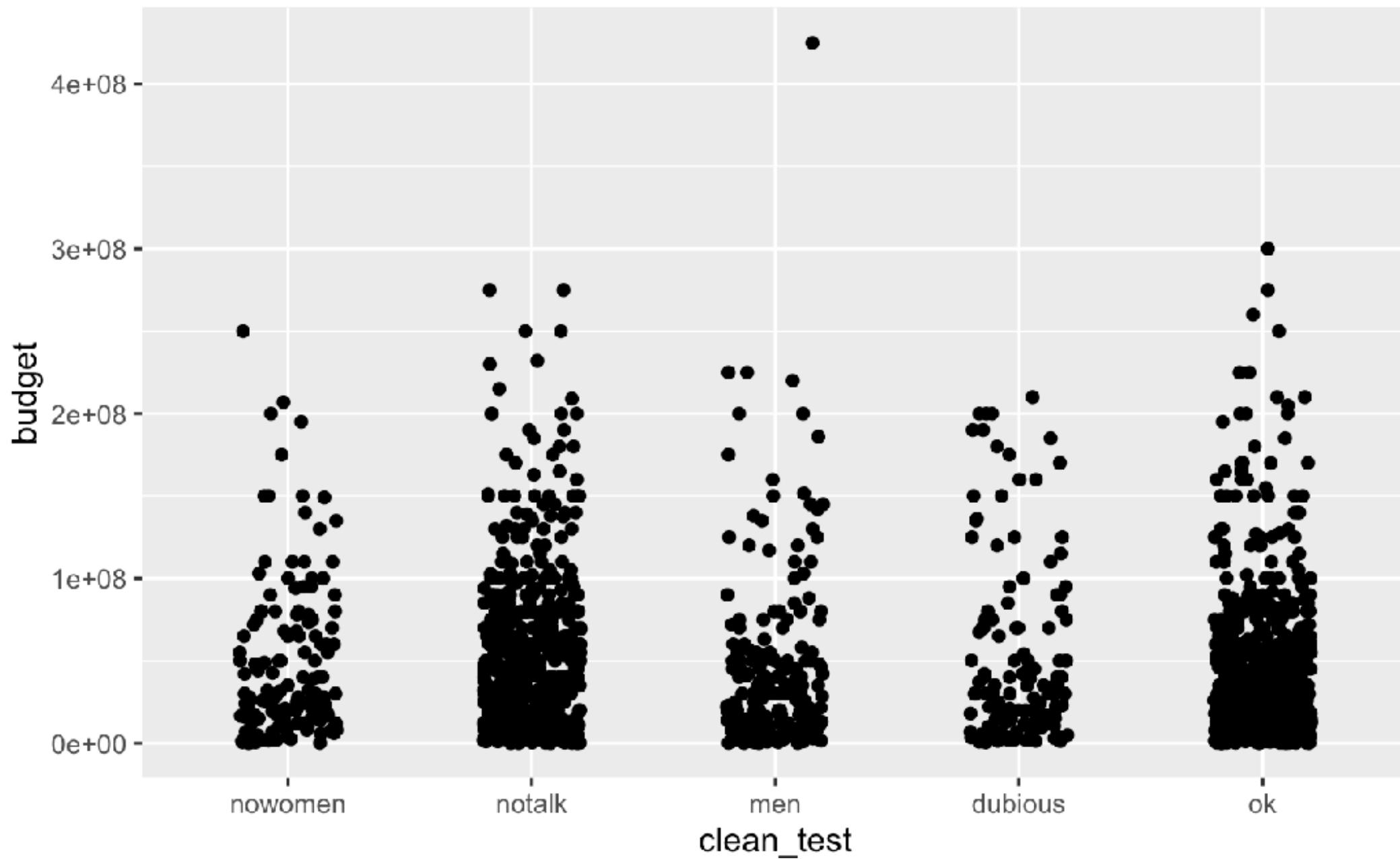
colorbrewer2.org

```
ggplot(data = bechdel) +  
  geom_bar(mapping = aes(x = clean_test, fill = clean_test)) +  
  scale_fill_brewer(palette = "Dark2")
```

Aesthetic scales

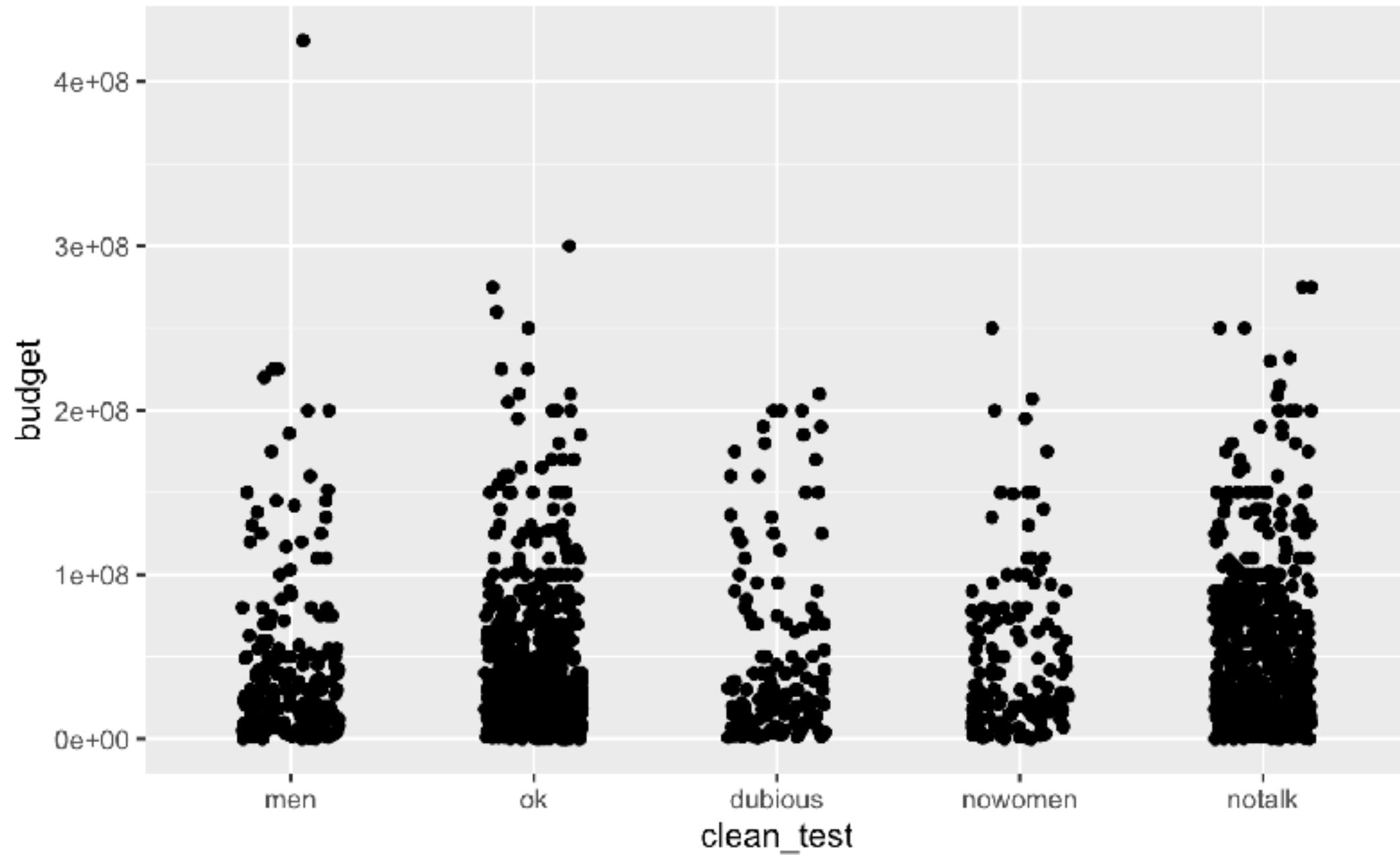
- `scale_*_continuous()`
- `scale_*_discrete()`
- `scale_*_ordinal()`
- `scale_{color/fill}_brewer()`
- `scale_{color/fill}_distiller()`
- `scale_{color/fill}_gradient()`

Coordinate scales



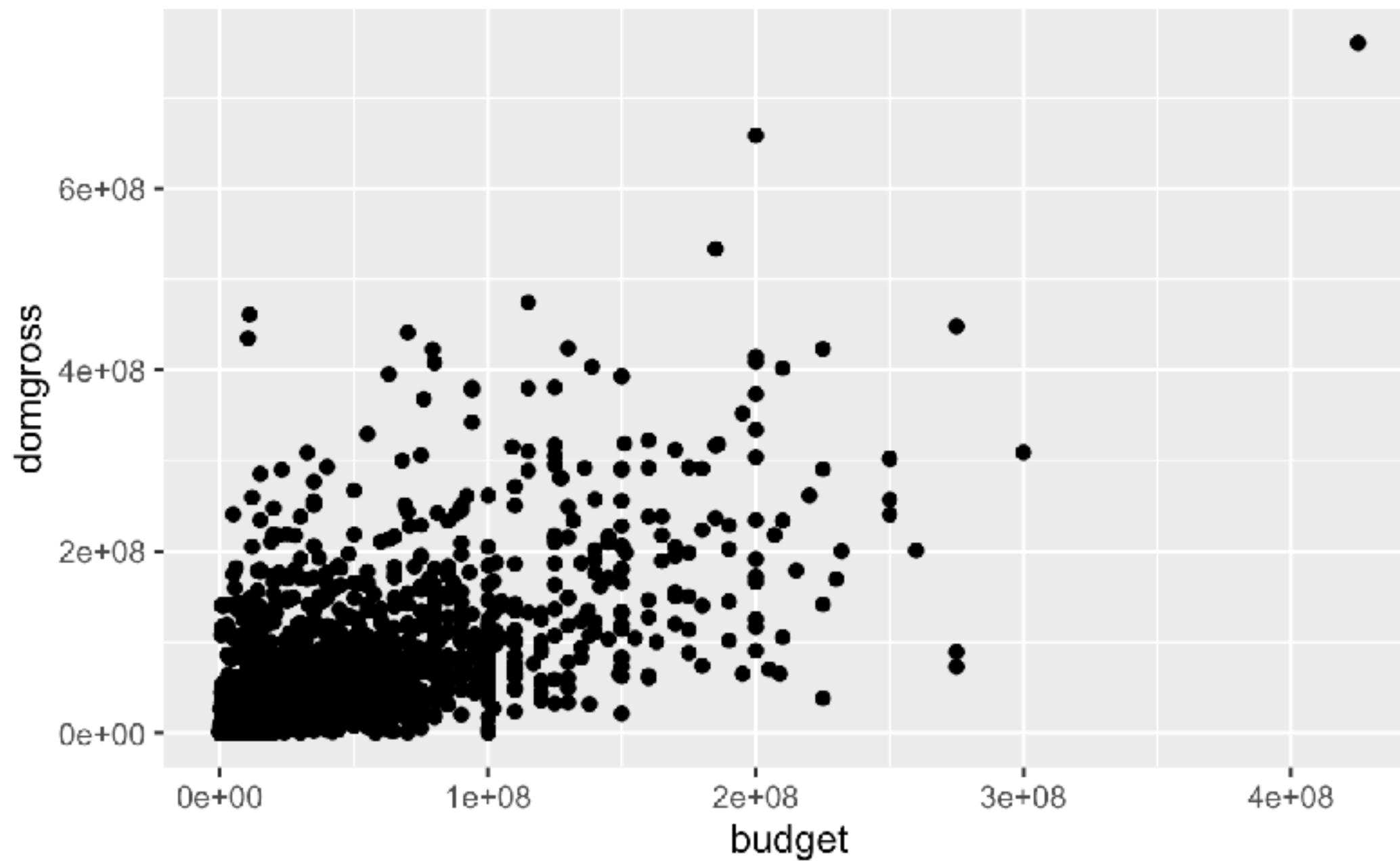
```
ggplot(bechdel, aes(x = clean_test, y = budget)) +  
  geom_point(position = position_jitter(width = 0.2, height = 0))
```

Coordinate scales



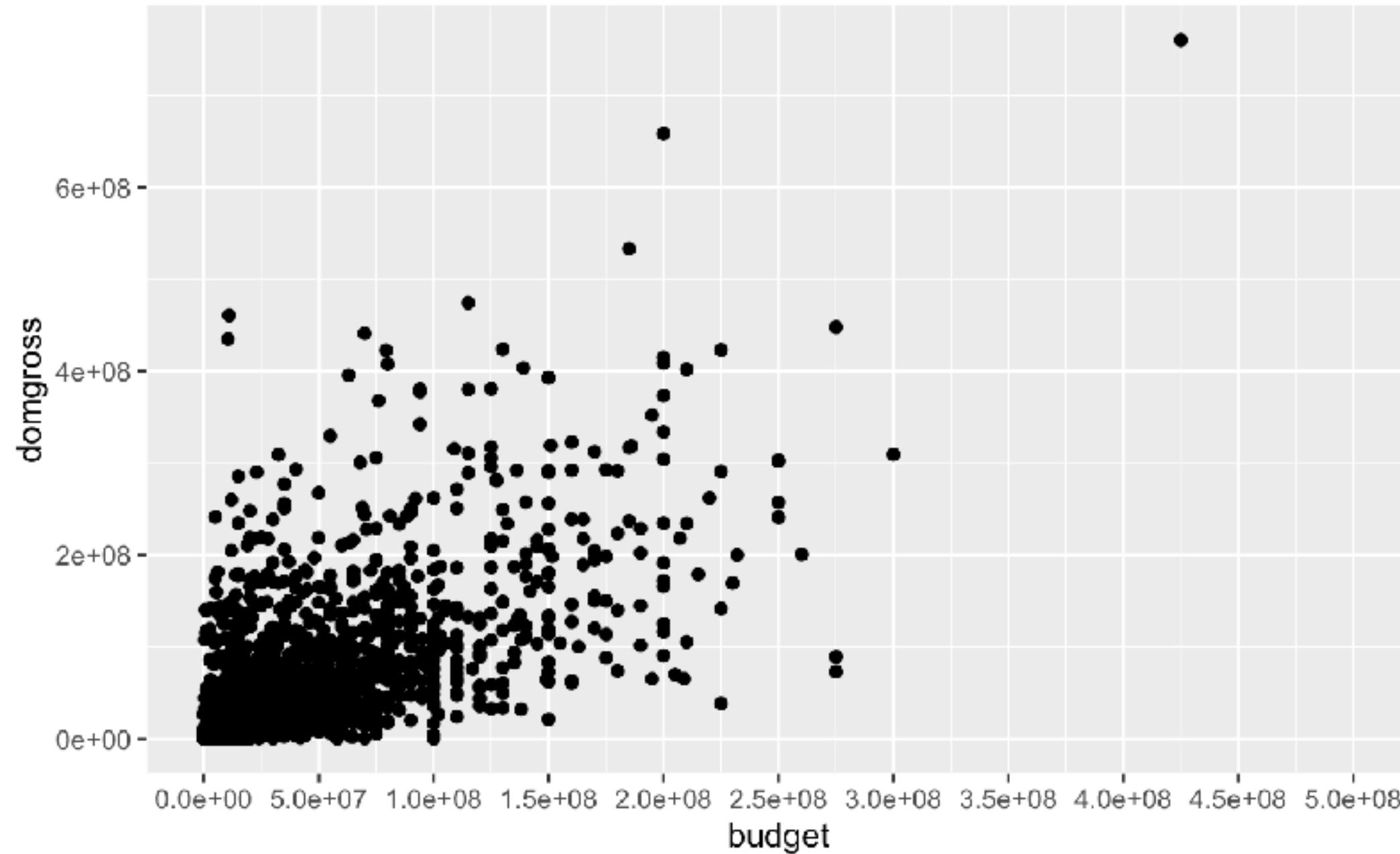
```
ggplot(bechdel, aes(x = clean_test, y = budget)) +  
  geom_point(position = position_jitter(width = 0.2, height = 0)) +  
  scale_x_discrete(limits = c("men", "ok", "dubious", "nowomen", "notalk"))
```

Coordinate scales



```
ggplot(bechdel, mapping = aes(x = budget, y = domgross)) +  
  geom_point()
```

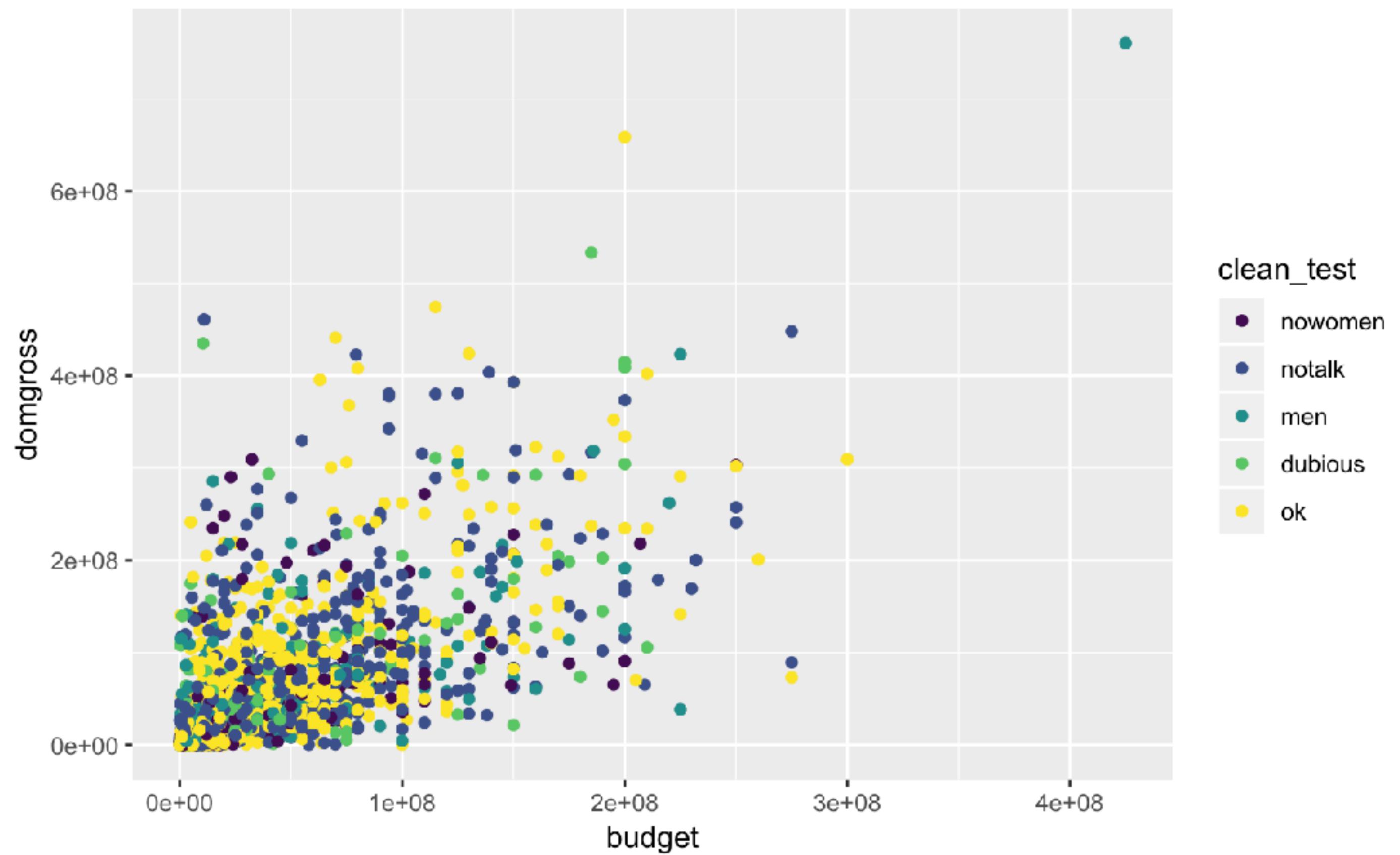
Coordinate scales



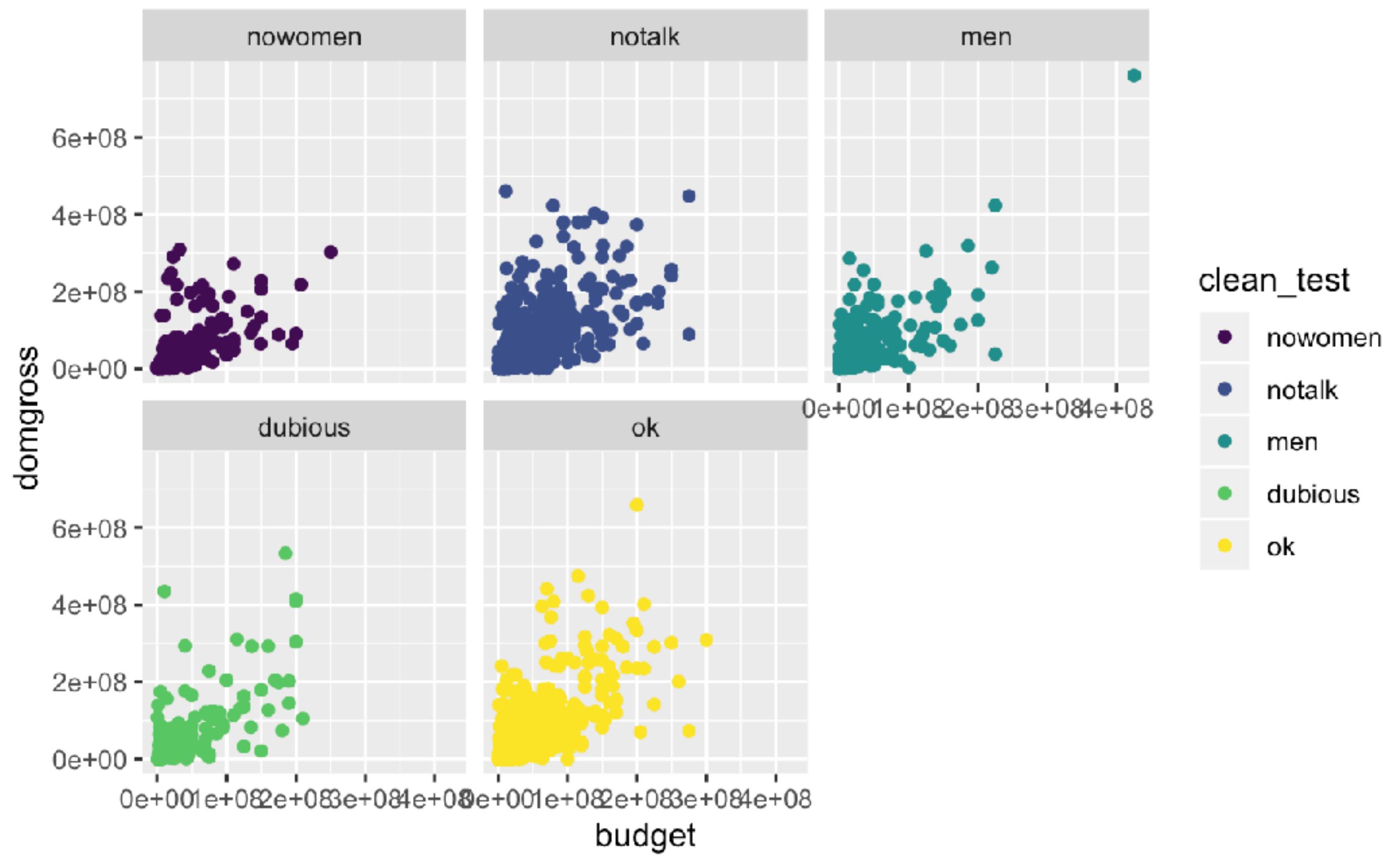
```
ggplot(bechdel, mapping = aes(x = budget, y = domgross)) +  
  geom_point() +  
  scale_x_continuous(limits = c(0, 5e+08), breaks = seq(0, 5e+08, 5e+07))
```



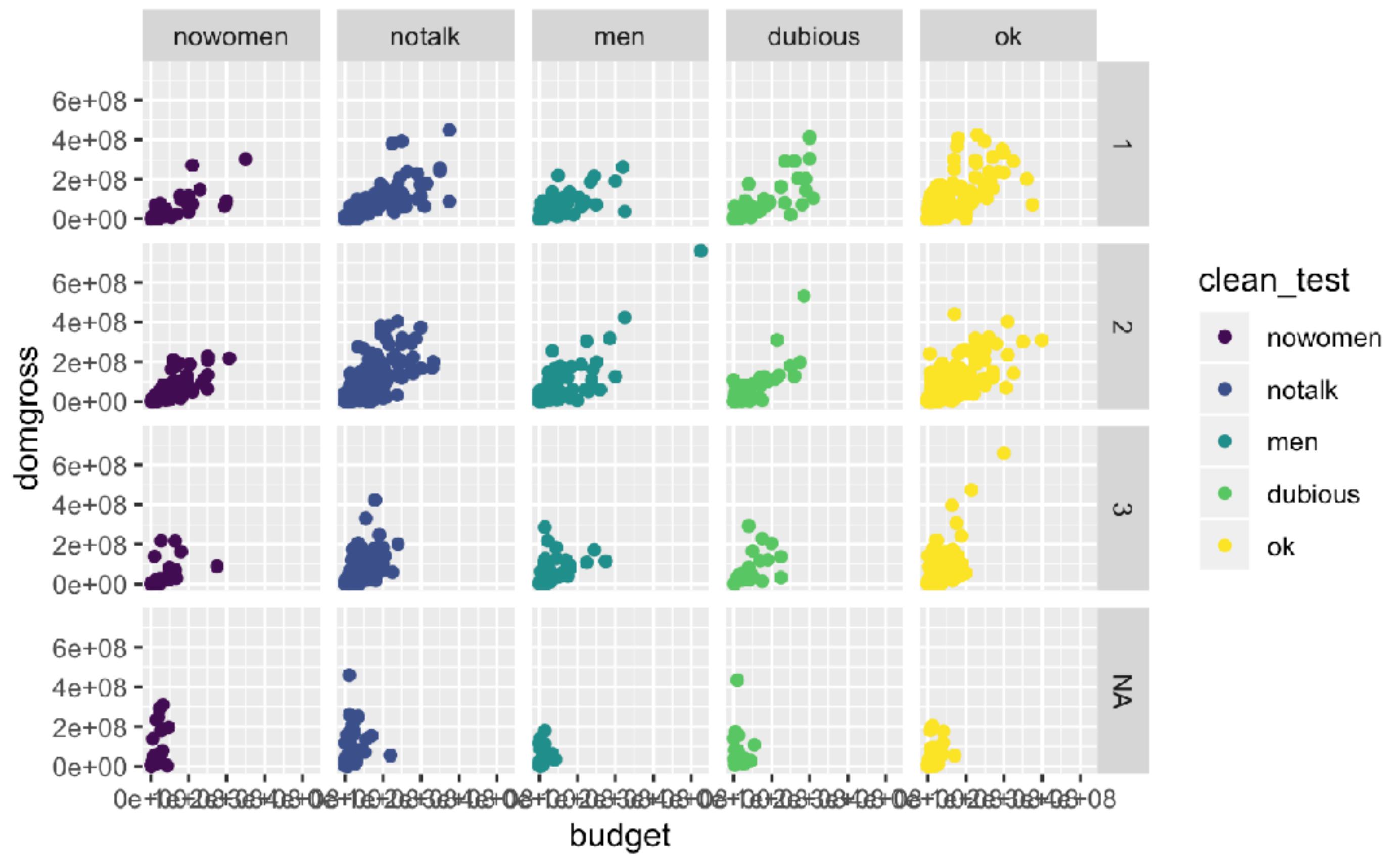
Facets



```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, color = clean_test))
```



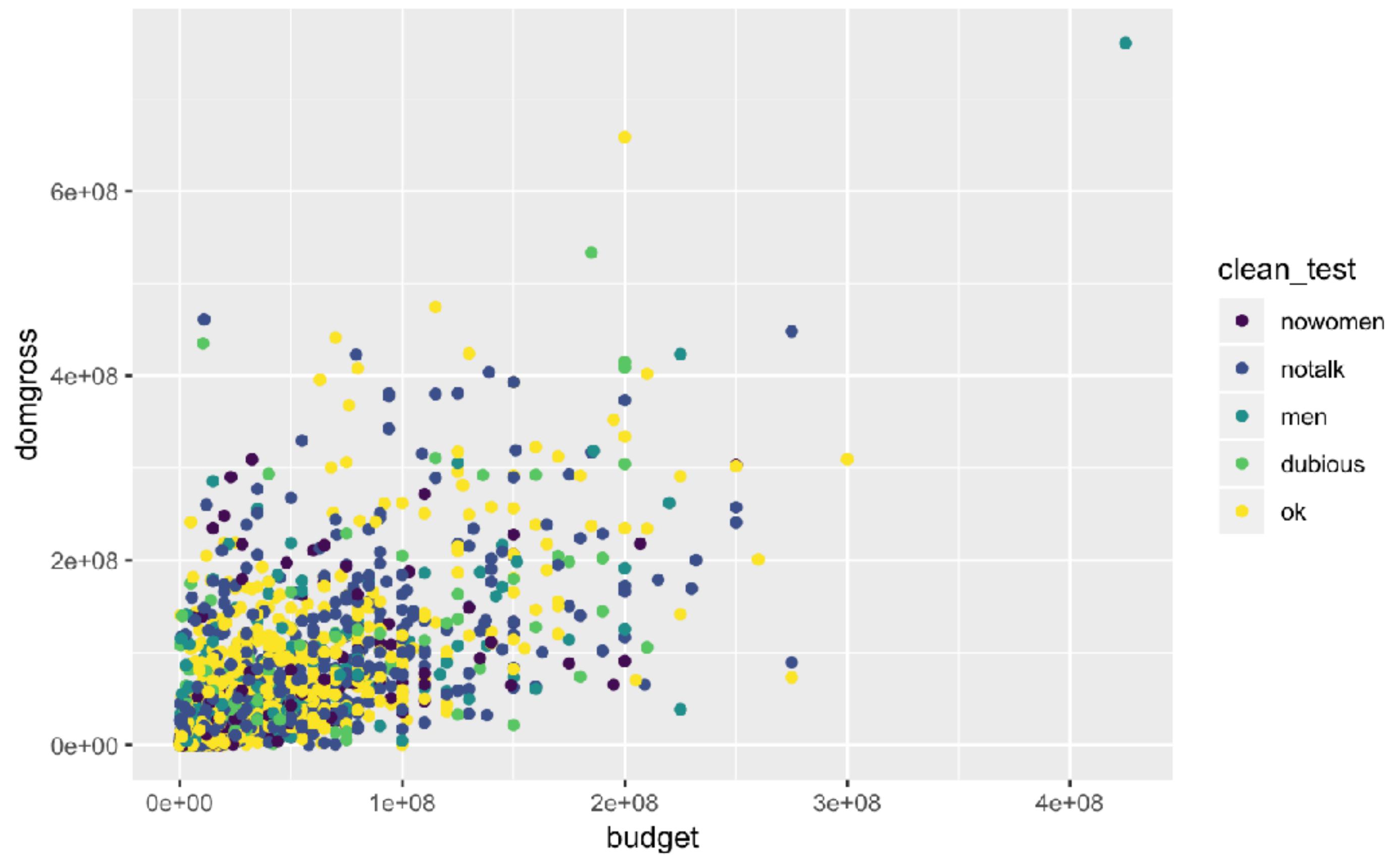
```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, color = clean_test)) +  
  facet_wrap(~ clean_test)
```



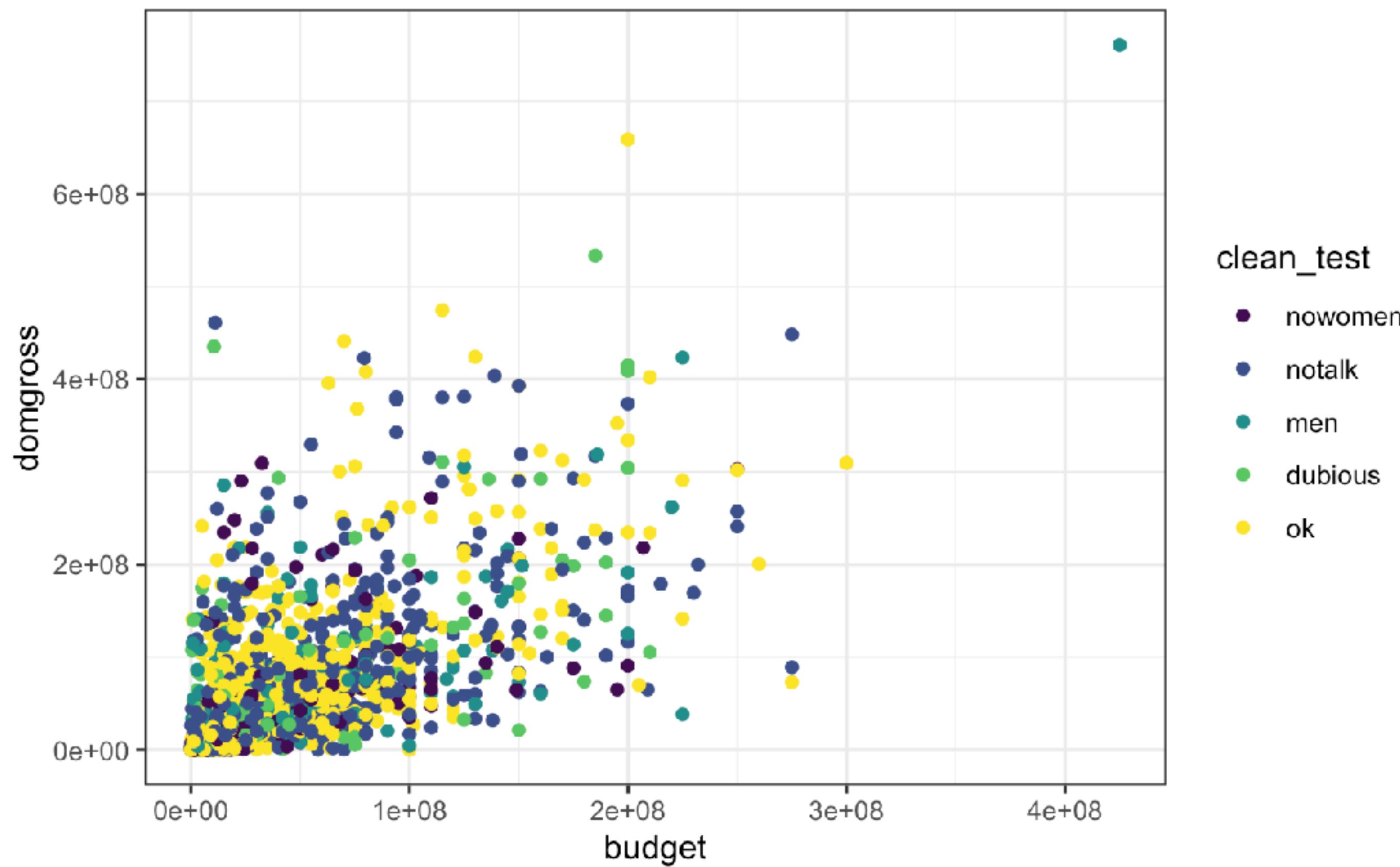
```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, color = clean_test)) +  
  facet_grid(decade_code ~ clean_test)
```



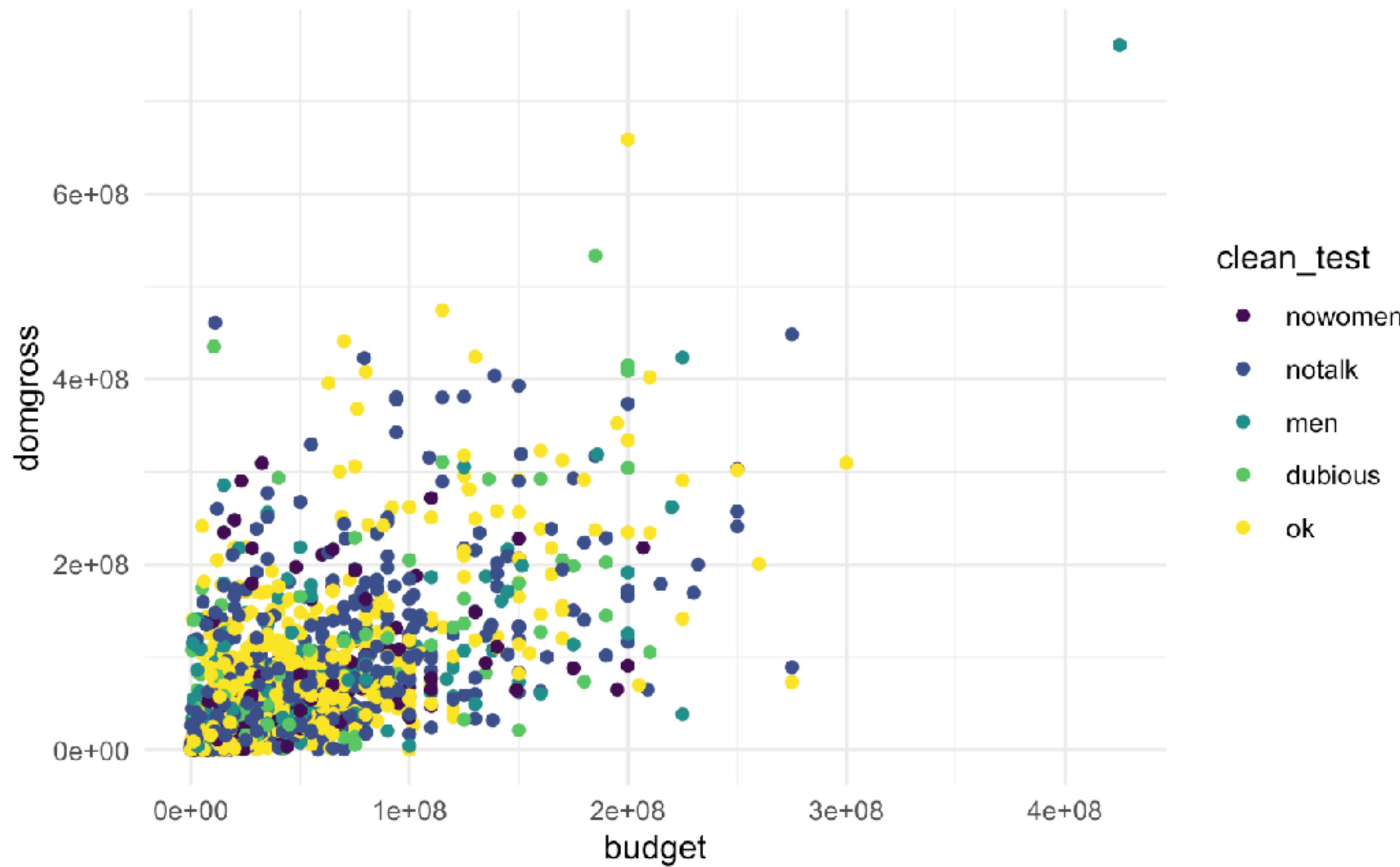
Themes



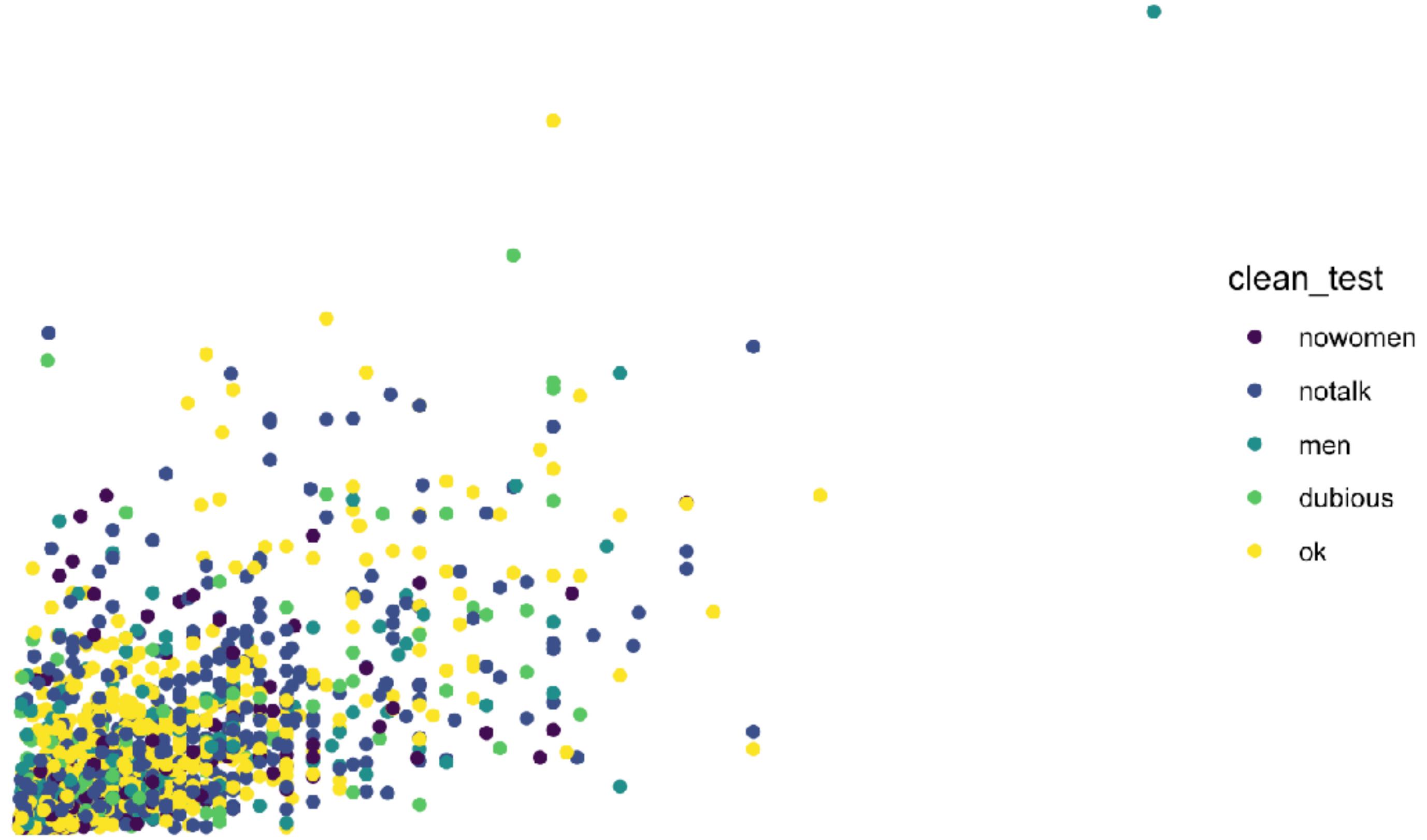
```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, color = clean_test))
```



```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, color = clean_test)) +  
  theme_bw()
```



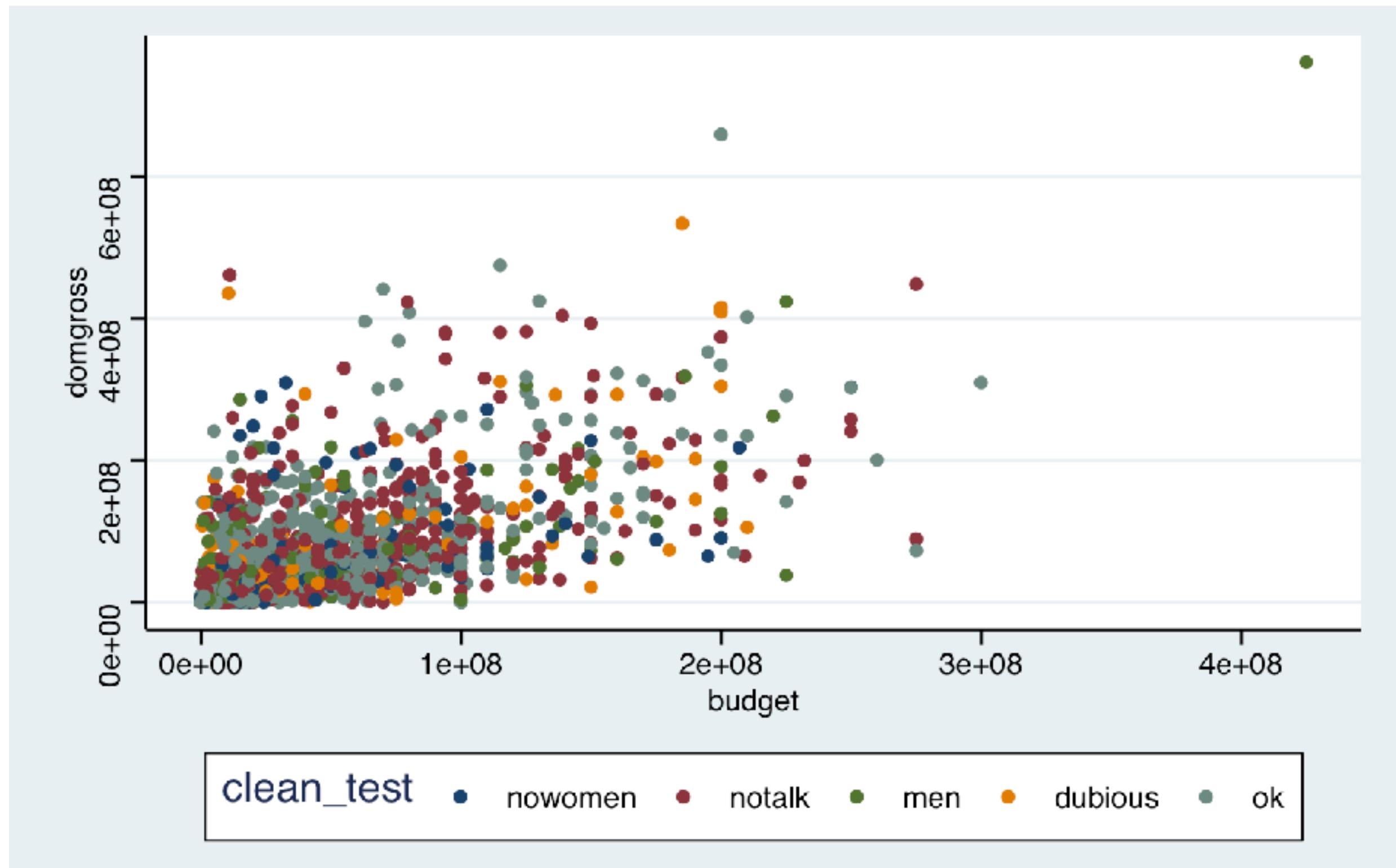
```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, color = clean_test)) +  
  theme_minimal()
```



```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, color = clean_test)) +  
  theme_void()
```

```
library(ggthemes)
```

```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, color = clean_test)) +  
  scale_color_stata() +  
  theme_stata()
```

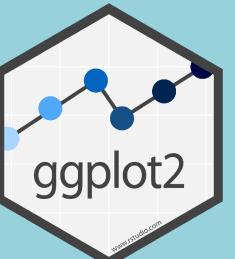


ggplot2 template

```
ggplot(data = <DATA>) +  
  <GEOM_FUNCTION>(mapping = aes(<MAPPINGS>),  
                    stat = <STAT>, position = <POSITION>) +  
  <FACET_FUNCTION> +  
  <SCALE_FUNCTION> +  
  <THEME_FUNCTION>
```

Required

Optional

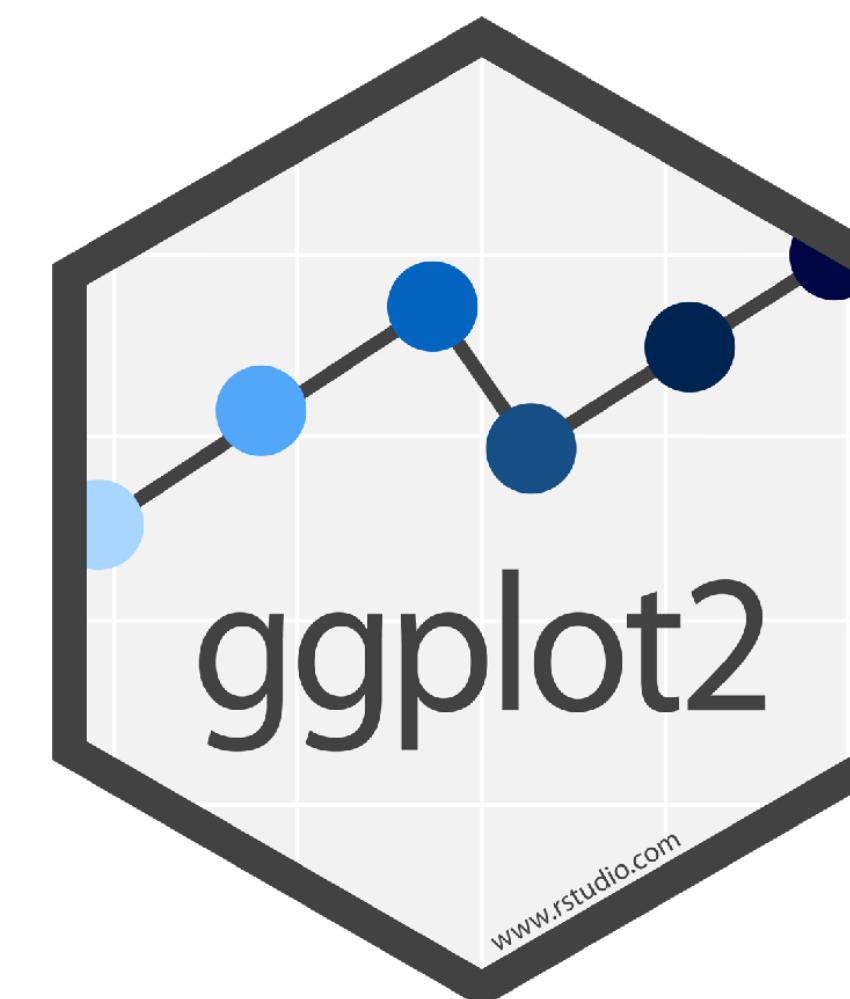


Saving plots

```
color_plot <- ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, color = clean_test))  
  
ggsave("filename.png", plot = color_plot, width = 8, height = 6,  
       units = "in", dpi = "retina")
```

Will save the last
plot created if not
specified

Data Visualization



wjakethompson.com

✉ wjakethompson@ku.edu

🐦 @wjakethompson

/github @wjakethompson