

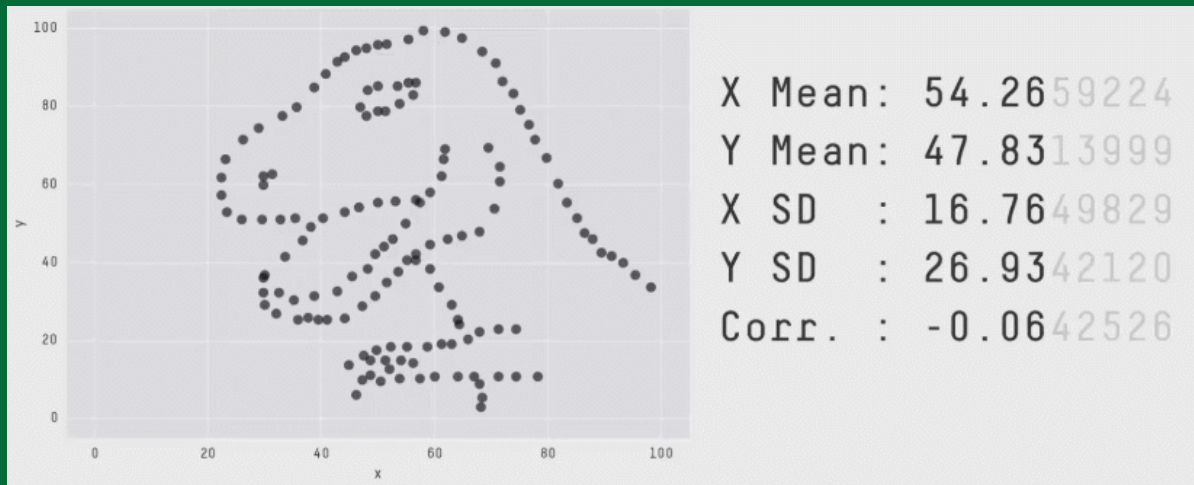
# DSBA 5122: Visual Analytics

## Class 5: Distributions and Uncertainty

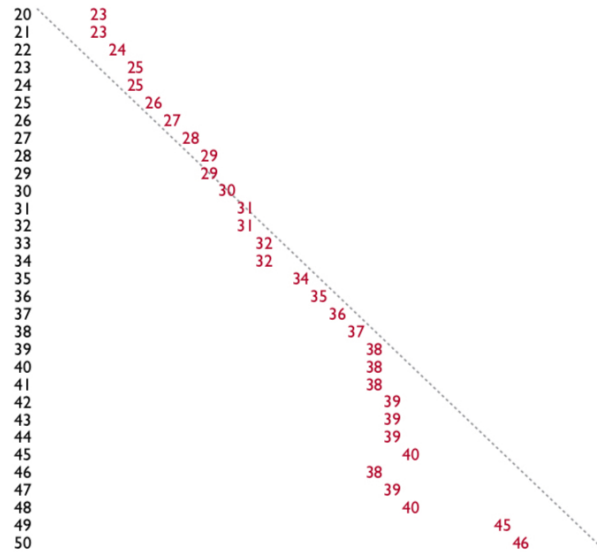
Ryan Wesslen

September 23, 2019

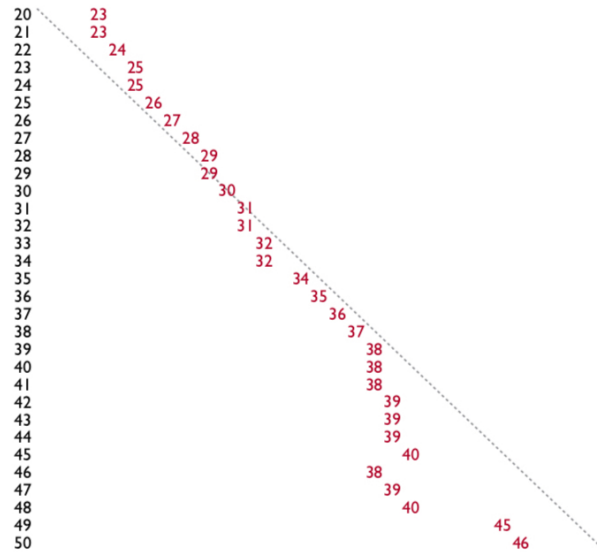
# Why view Distributions? Cairo Ch. 7 & Wilke Ch. 7 - 9



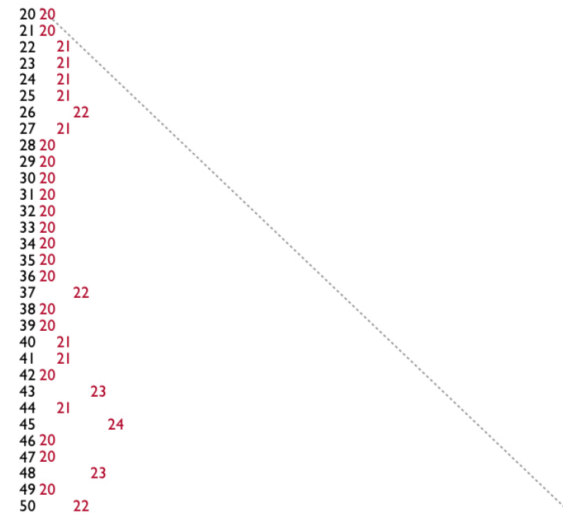
*a woman's age vs. the age of the men who look best to her*



*a woman's age vs. the age of the men who look best to her*



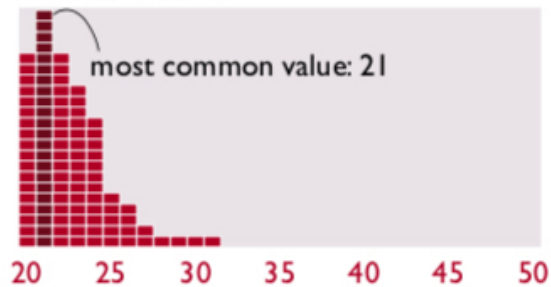
*a man's age vs. the age of the women who look best to him*



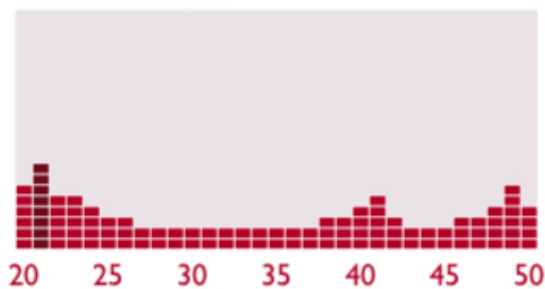
*men of 40 vs. the age  
of the women who look  
best to them*

■ = 1% of men

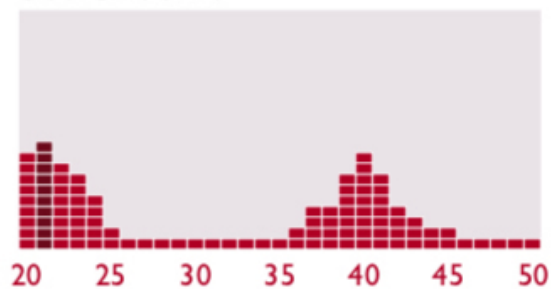
*distribution 1*



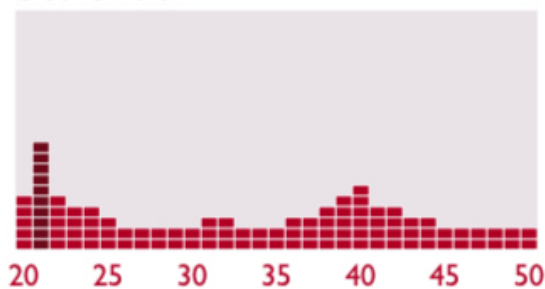
*distribution 2*



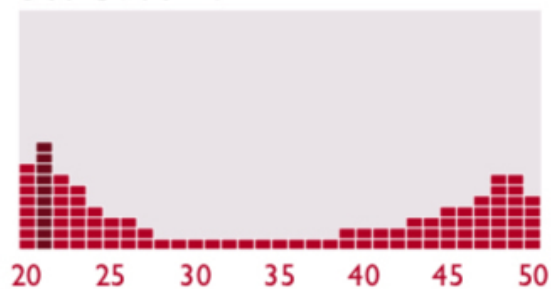
*distribution 3*



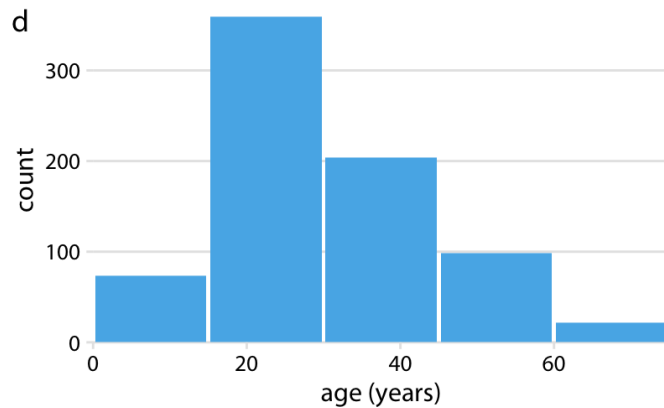
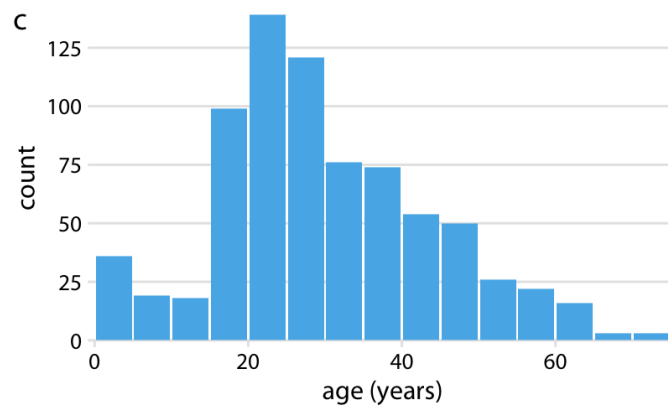
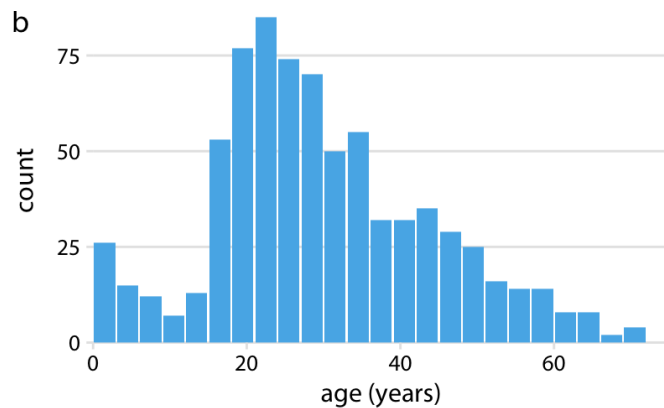
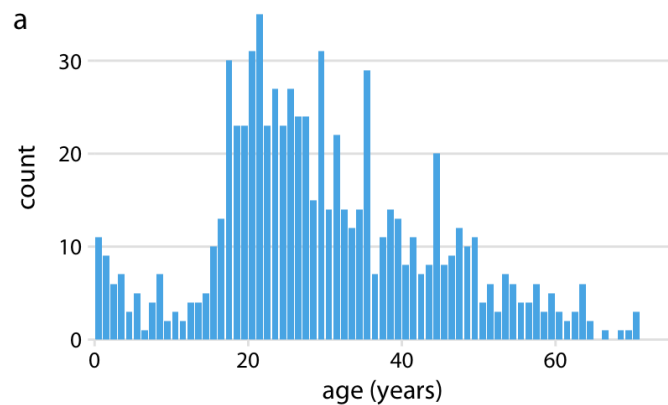
*distribution 4*

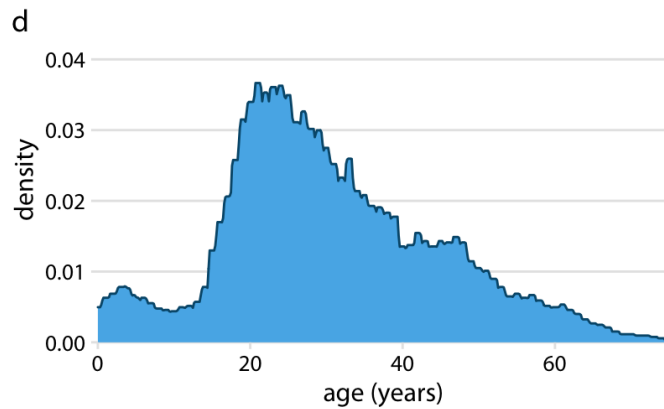
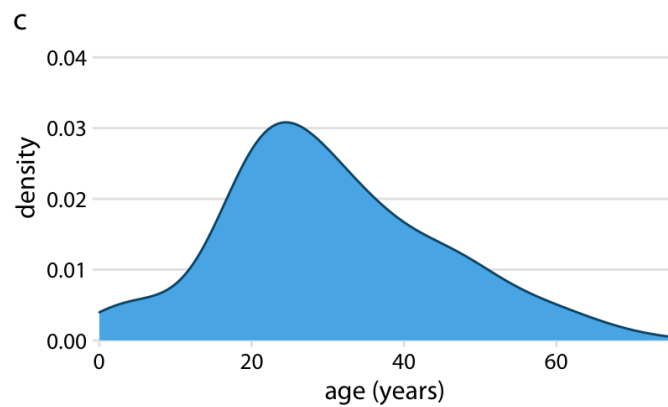
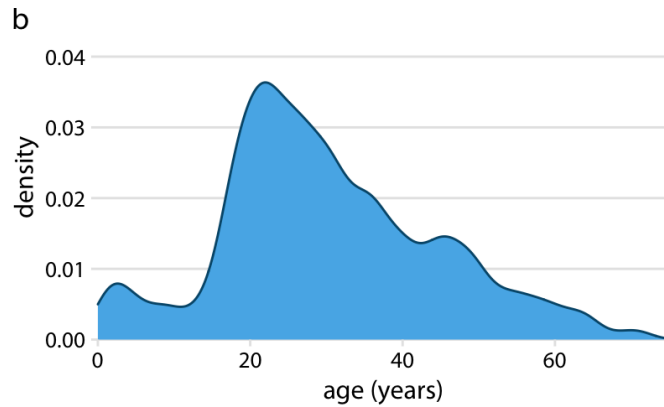
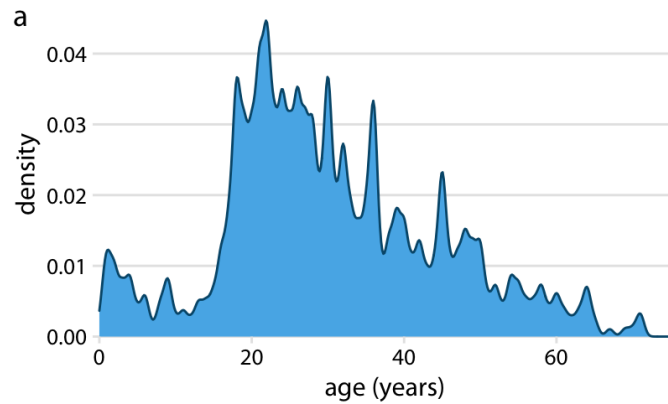


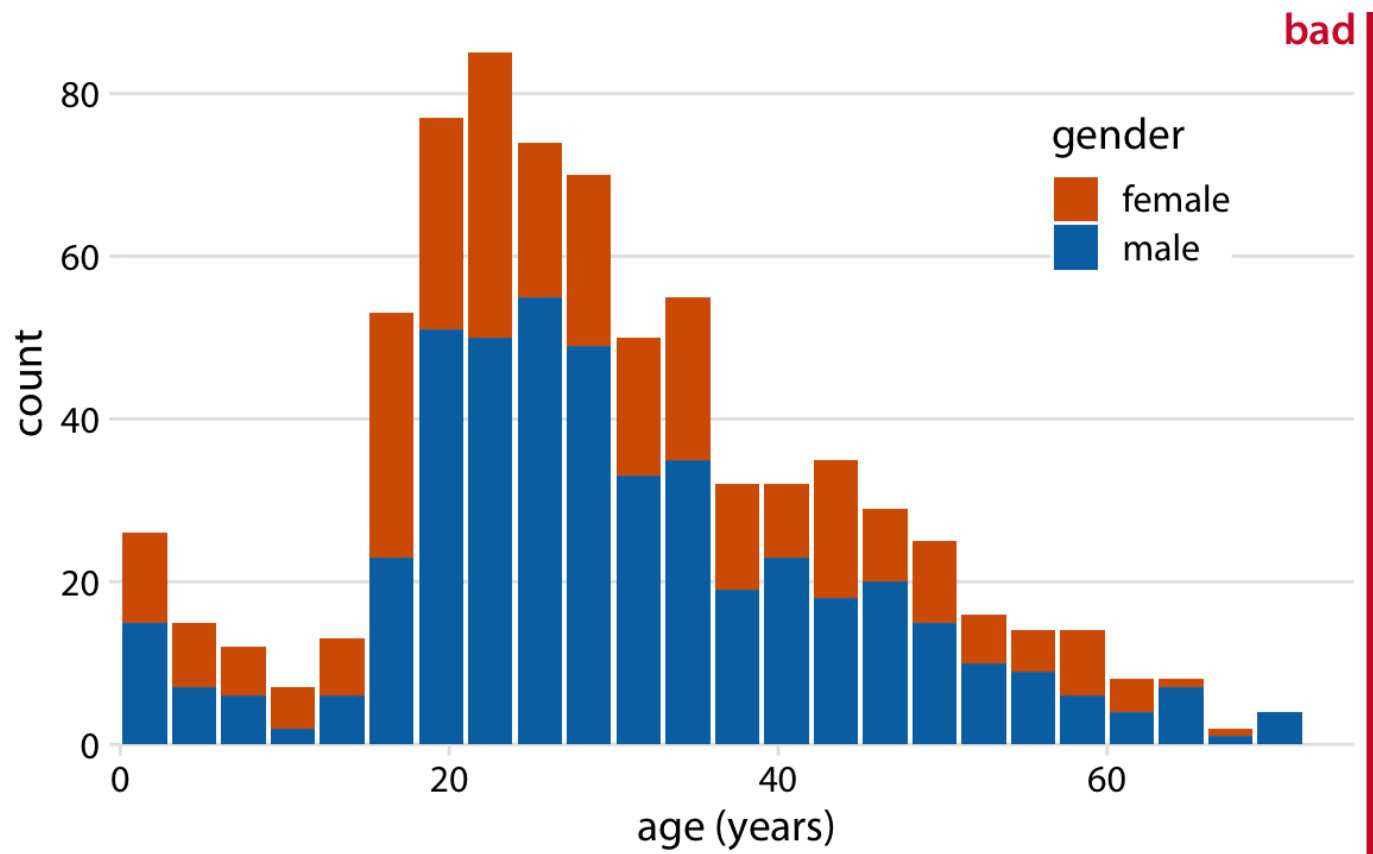
*distribution 5*



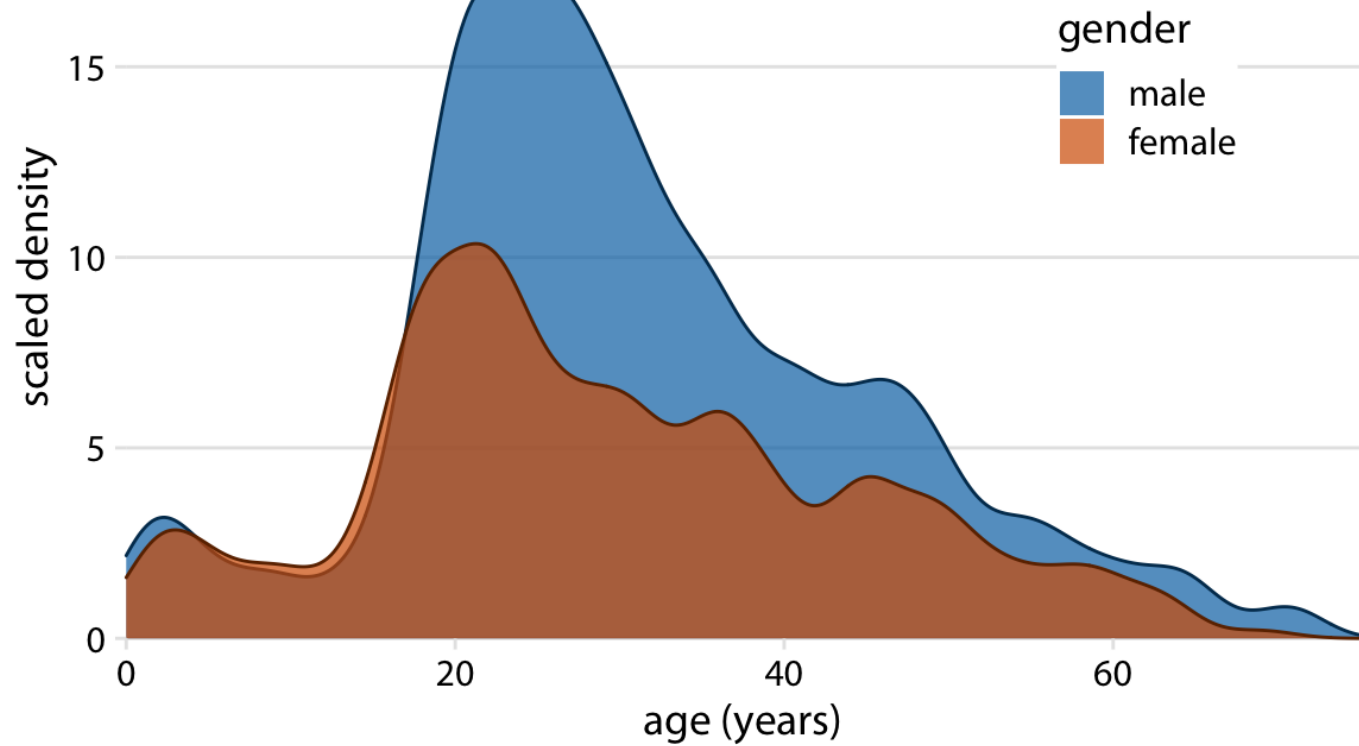
**Figure 7.2** Five possible (and fictitious) distributions for the data corresponding to the preferences of men of 40. All of them have the same mode: 21.

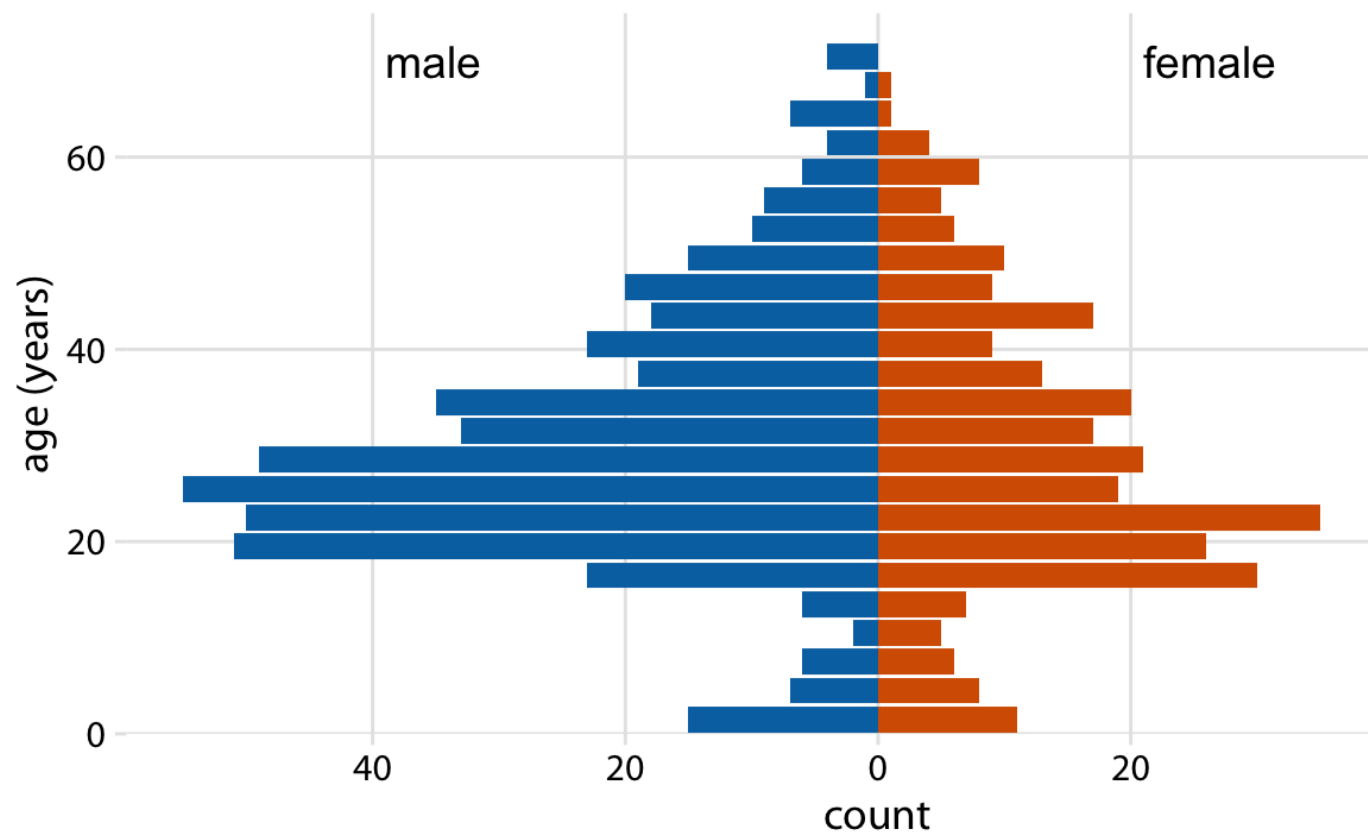


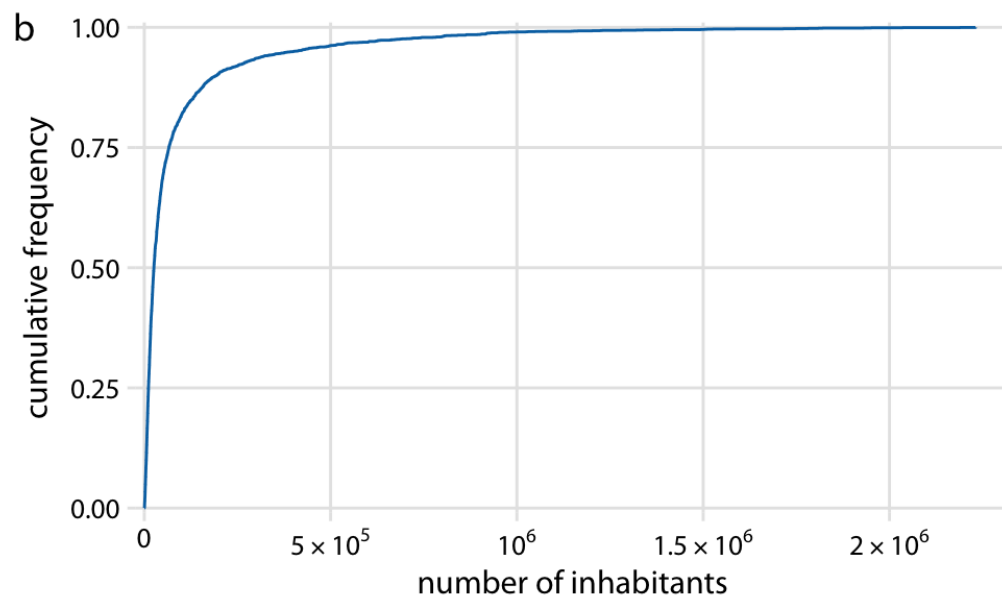
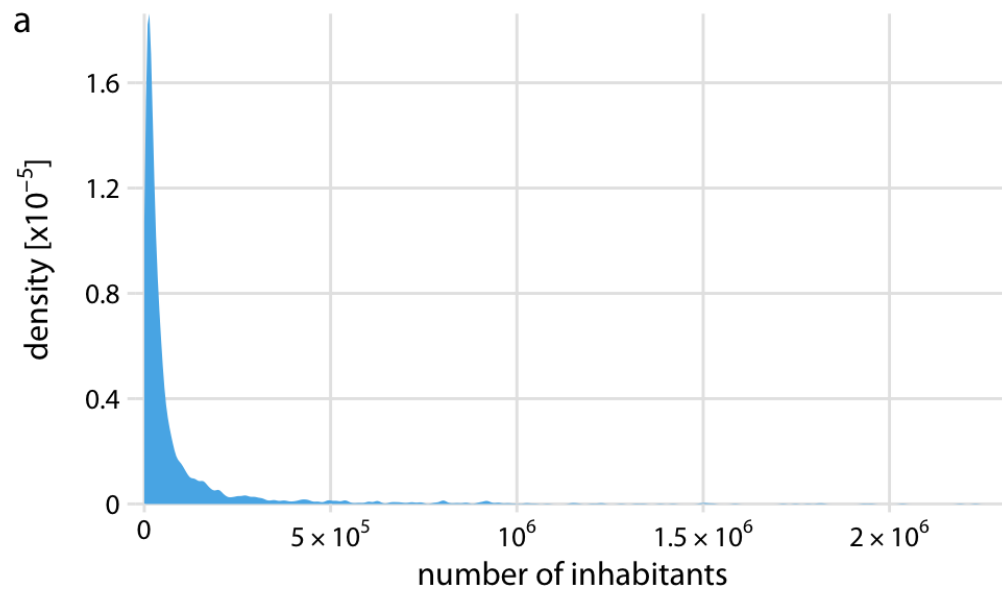


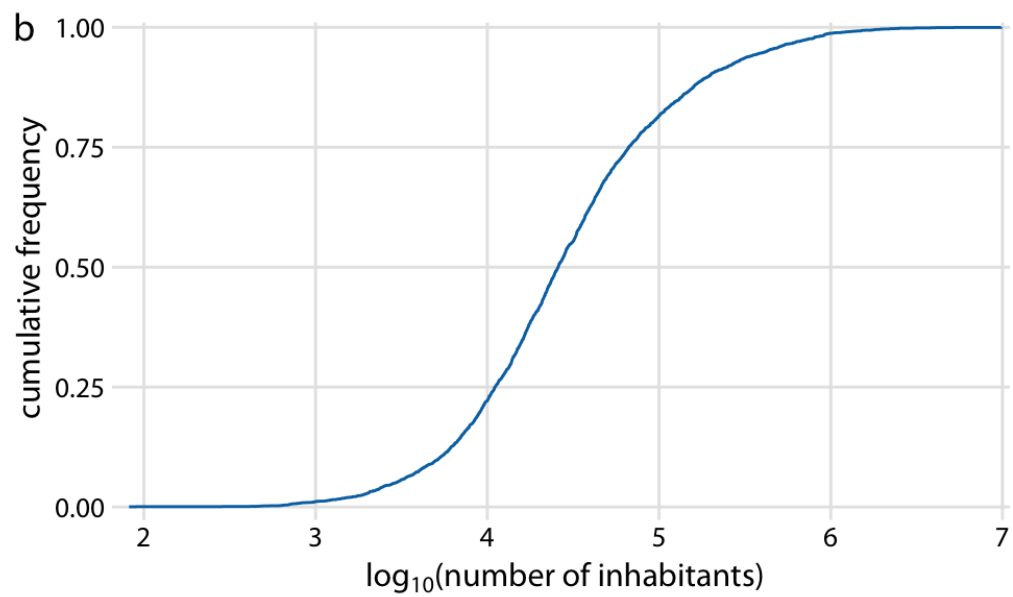
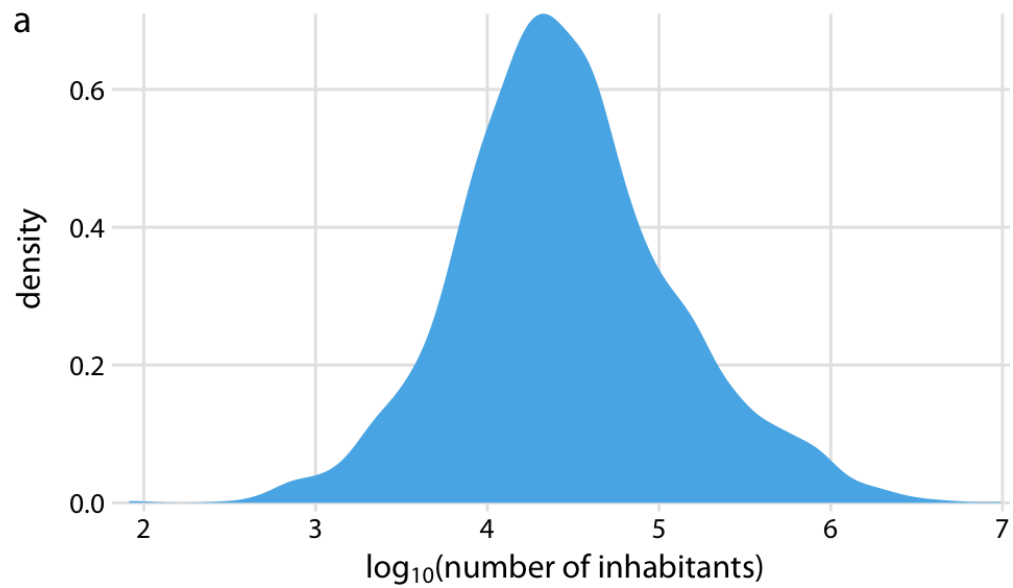




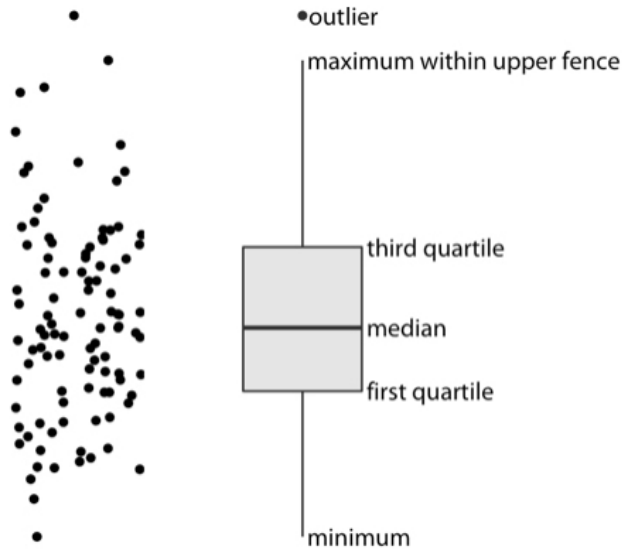




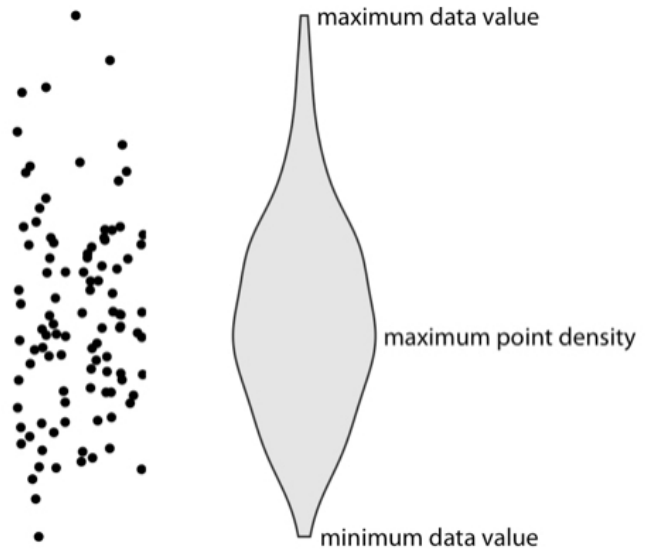


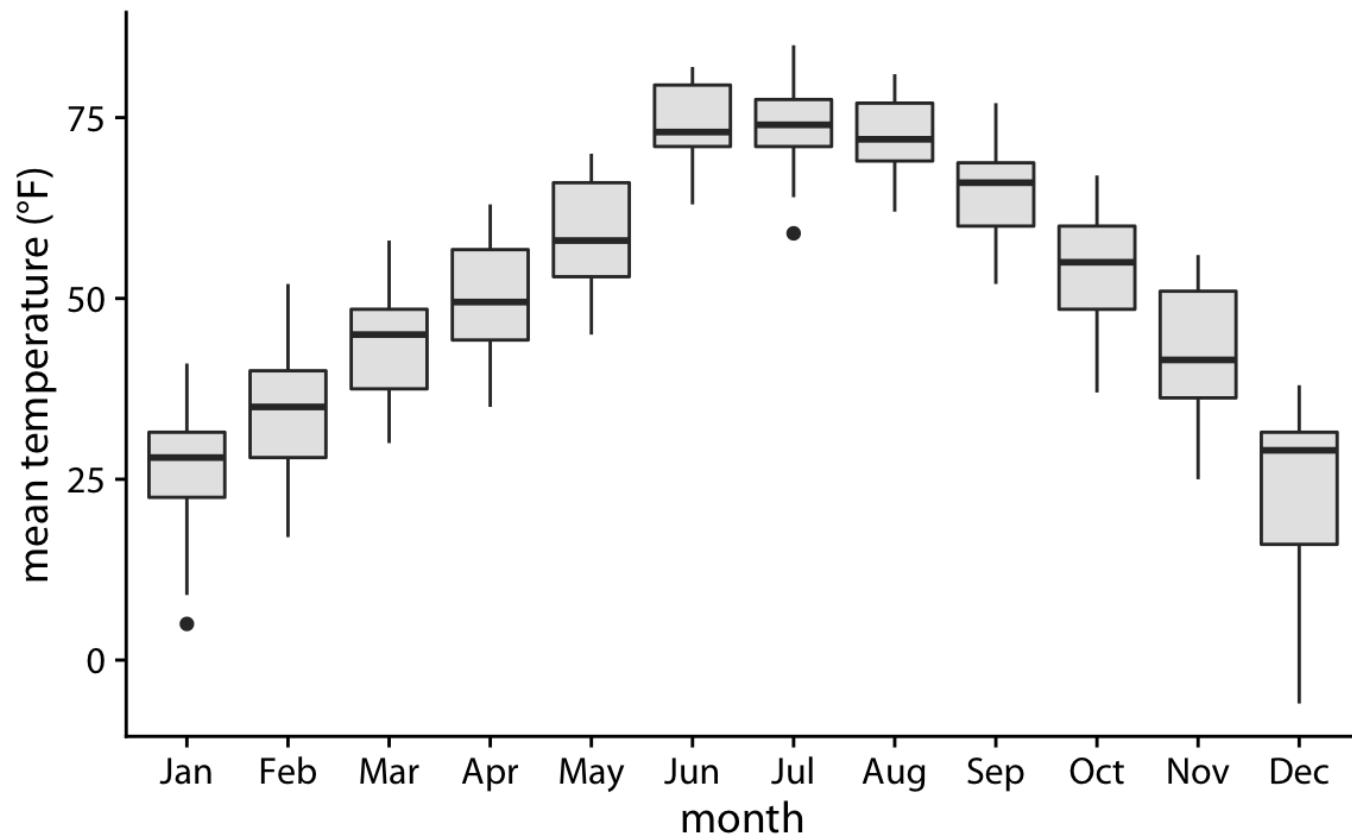


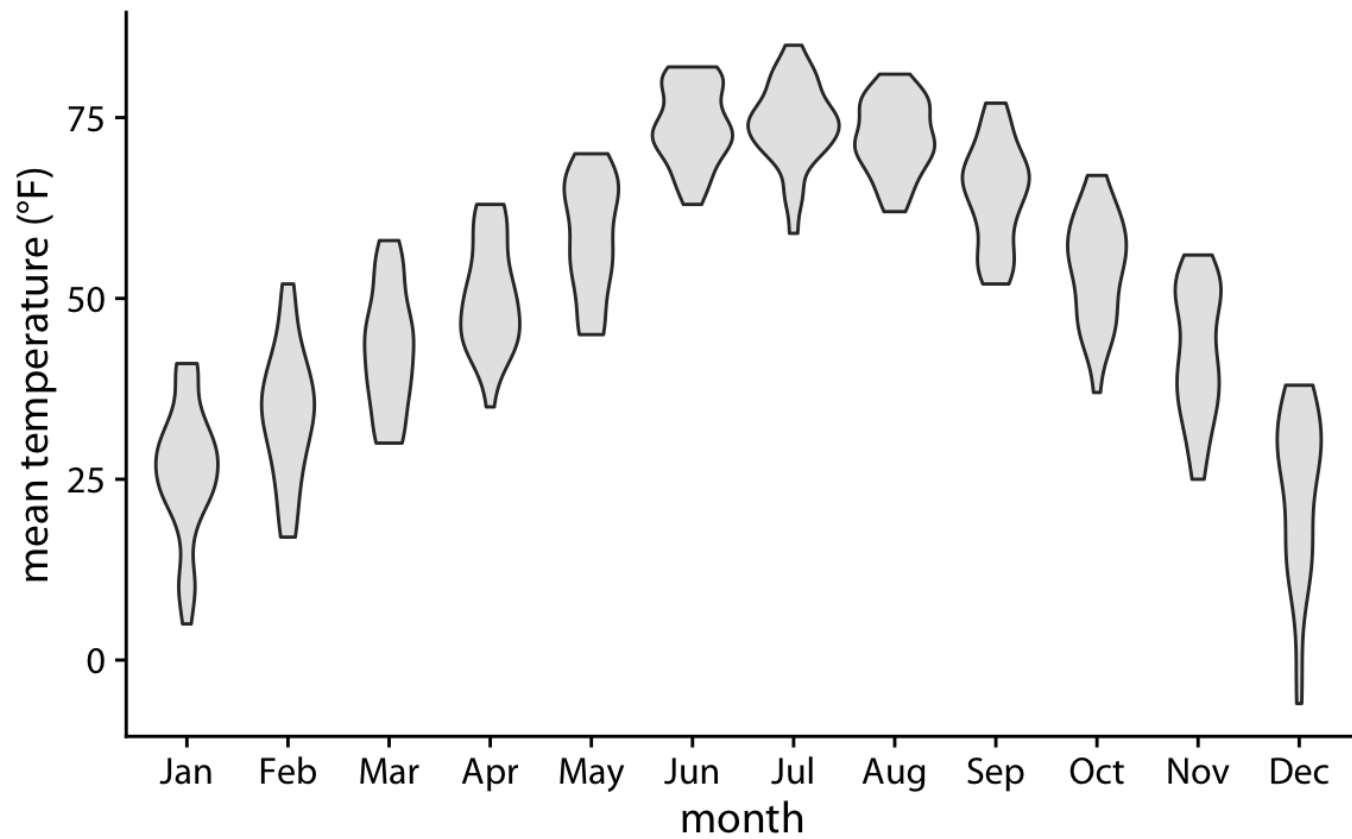
```
ggplot(df, aes(x, y)) + geom_boxplot()
```

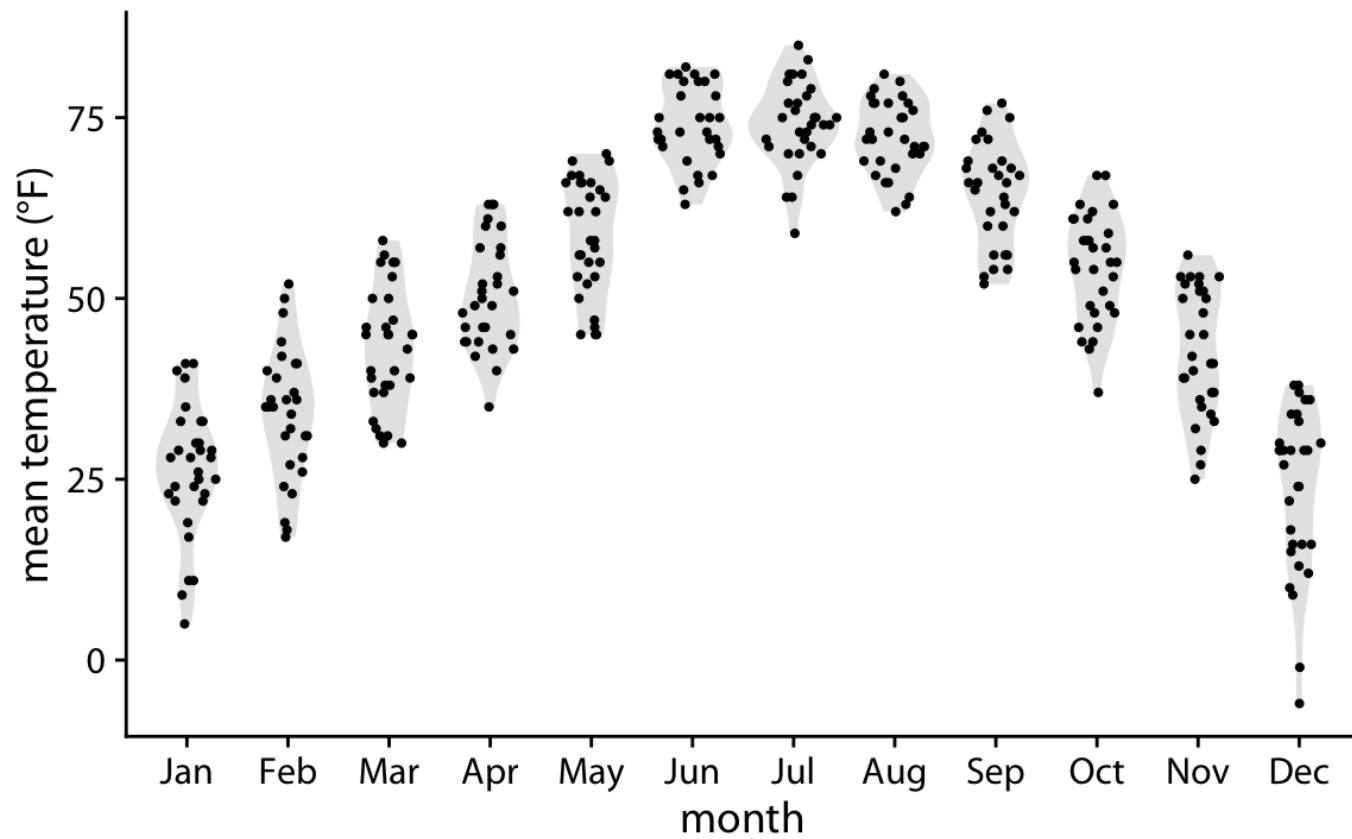


```
ggplot(df, aes(x, y)) + geom_violin()
```

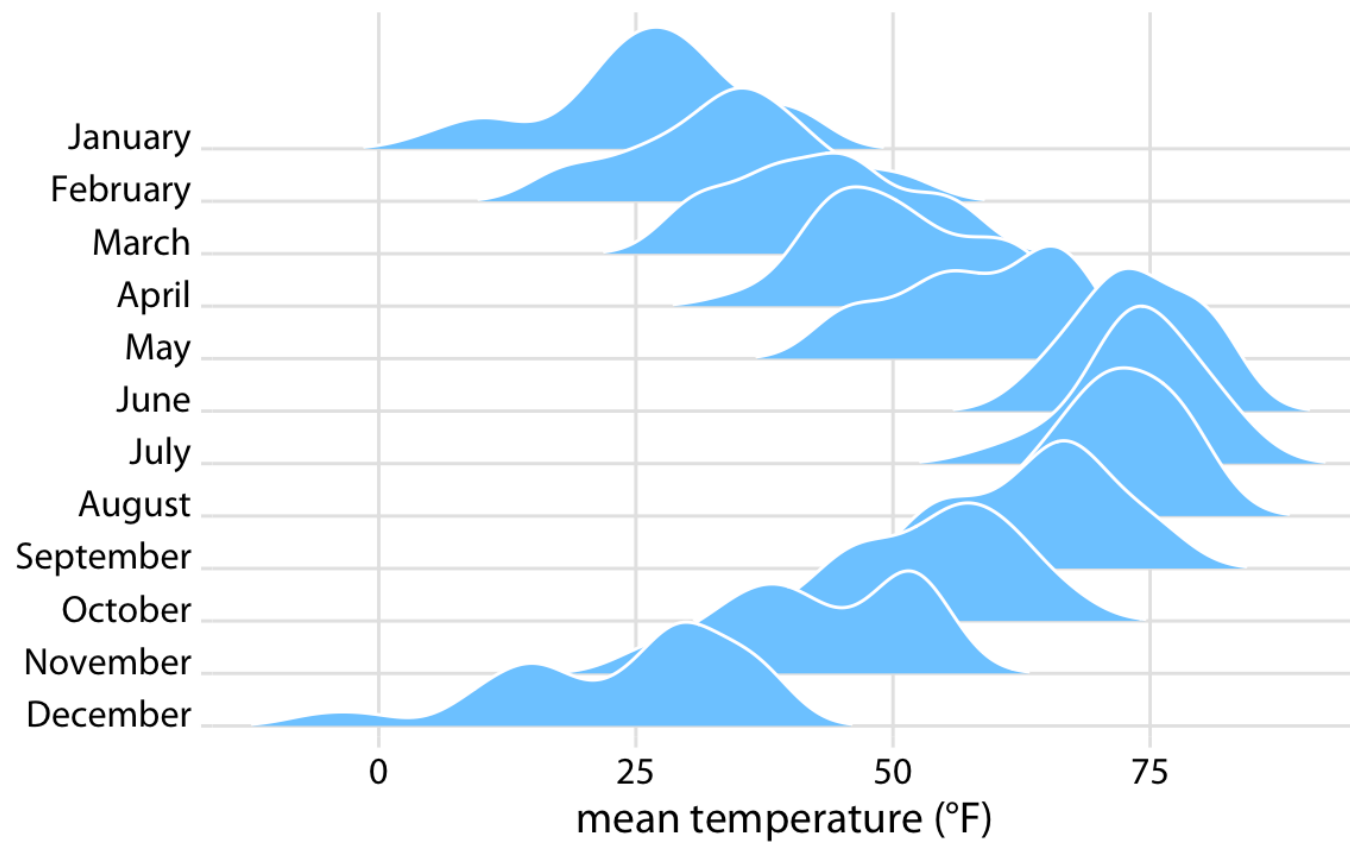


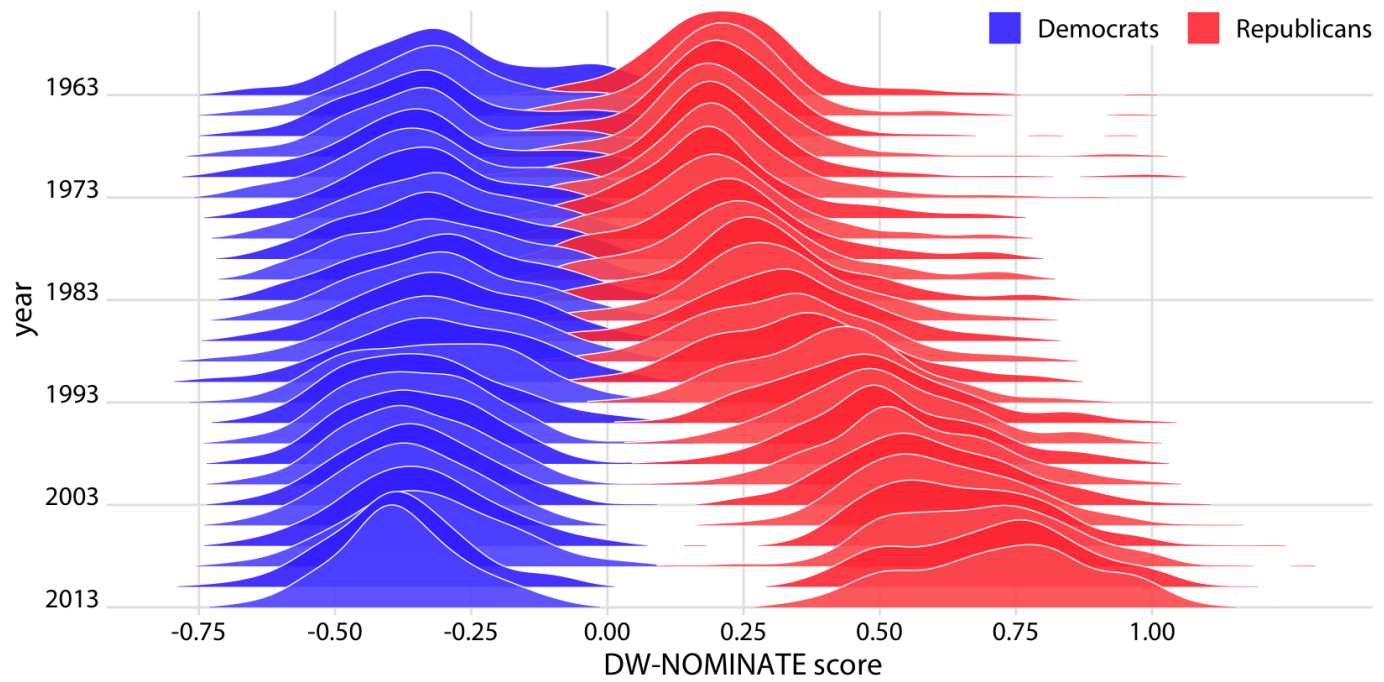




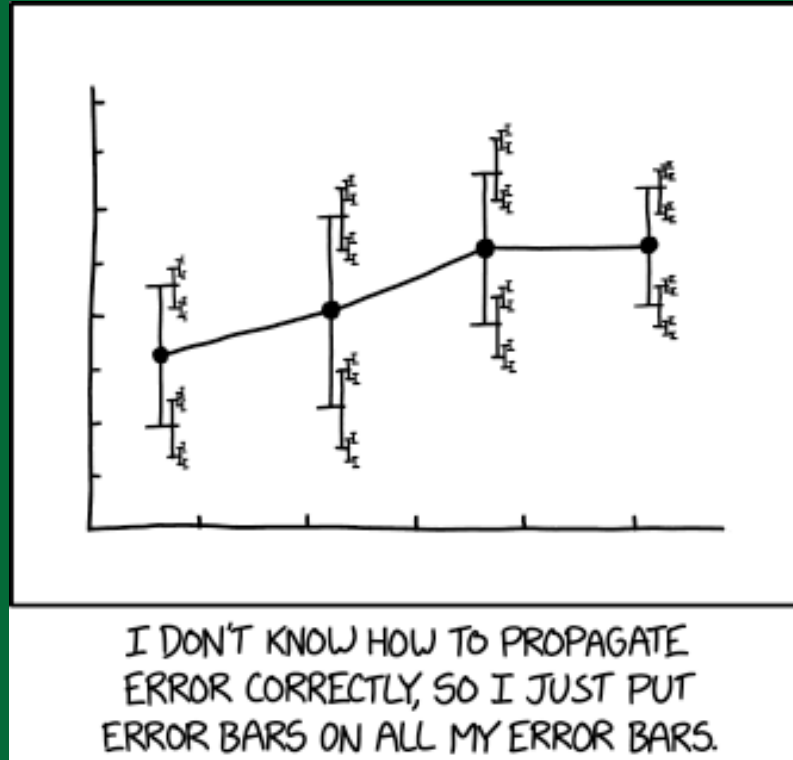




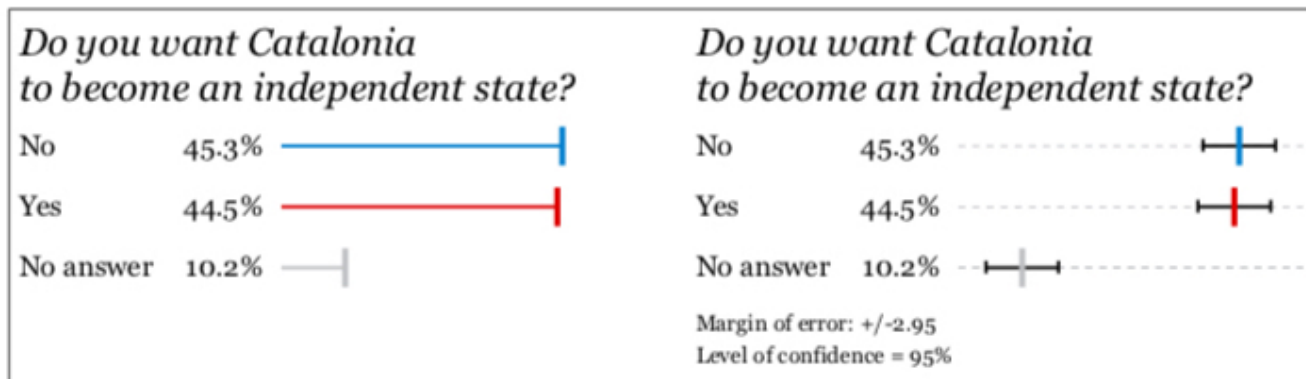




# Uncertainty: Cairo Ch. 10 & Wilke Ch. 16



xkcd



**Figure 11.1** Displaying the margin of error can change your view of the data.

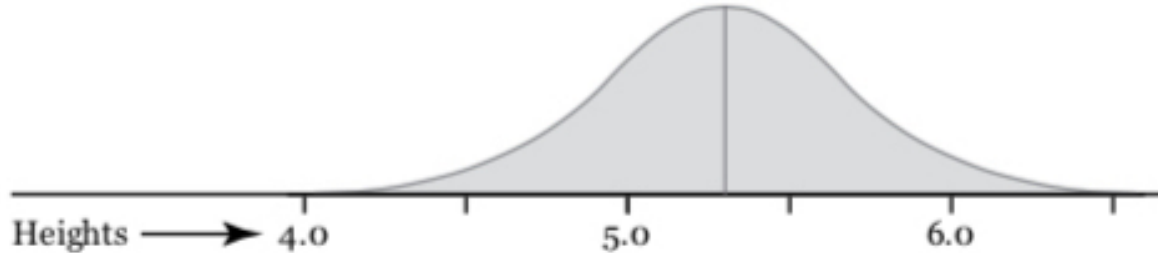
**POPULATION  
distribution  
(Unknown to you)**

I've drawn a normal distribution,  
but the actual population distribution  
could have a different shape



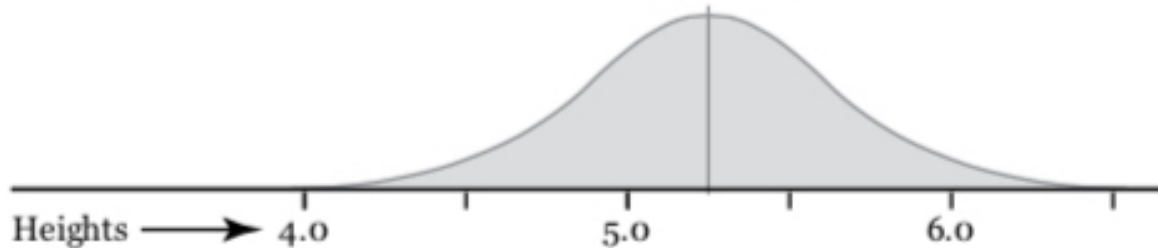
**SAMPLE A**

Mean: 5.3 feet



**SAMPLE B**

Mean: 5.25 feet



**Figure 11.2** Population and samples.

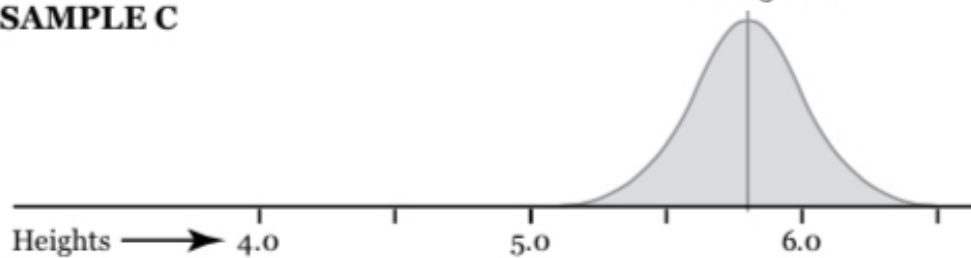
**POPULATION  
distribution  
(Unknown to you)**

I've drawn a normal distribution,  
but the actual population distribution  
could have a different shape

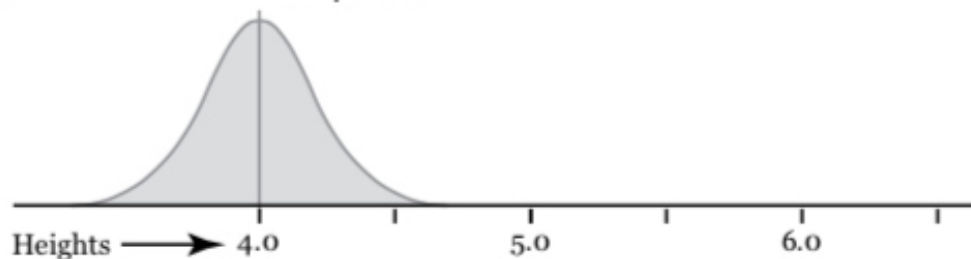


**SAMPLE C**

Mean: 5.8 feet



**SAMPLE D** Mean: 4.0 feet



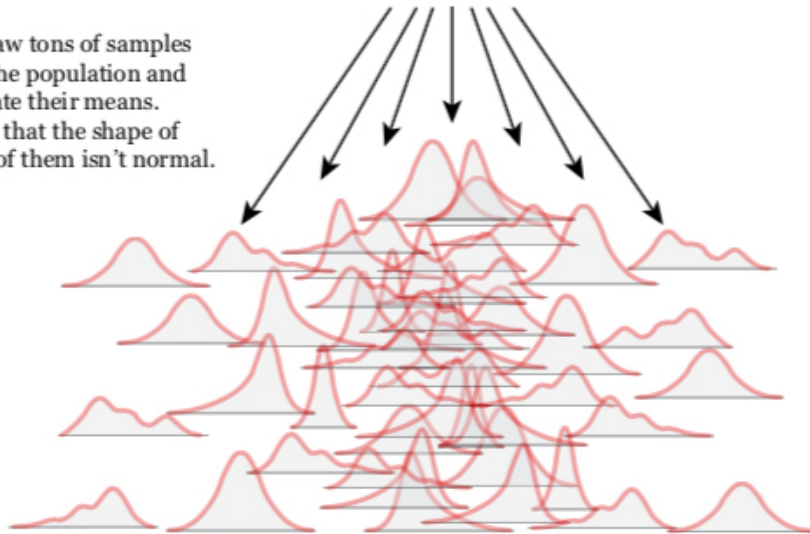
**Figure 11.3** When drawing many samples from a population, it is possible to obtain a few with means that greatly differ from the population mean.

**Population:**

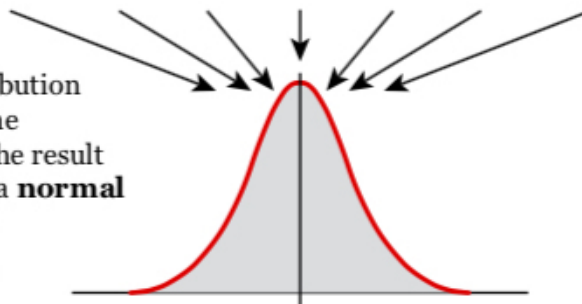
Unknown mean,  
standard deviation,  
and shape.



We draw tons of samples  
from the population and  
calculate their means.  
Notice that the shape of  
many of them isn't normal.



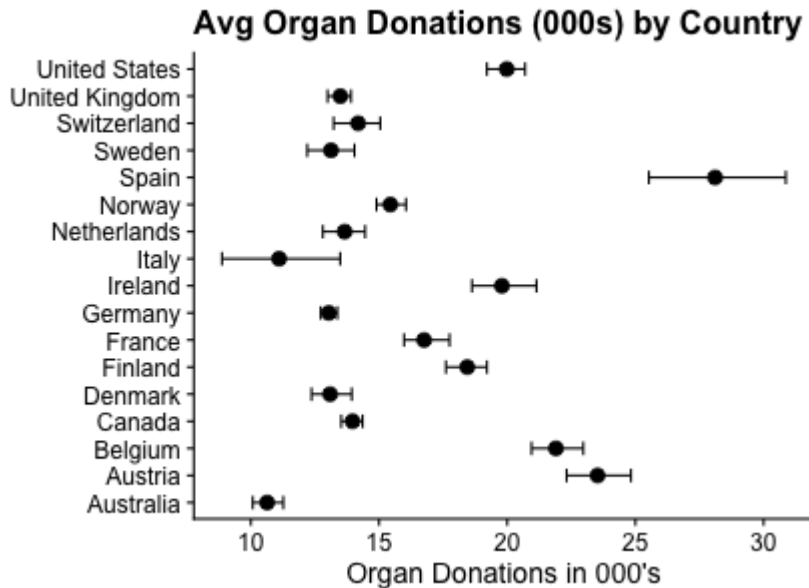
We plot the distribution  
of the means of the  
samples above. The result  
is approximately a **normal**  
**distribution of**  
**sample means.**



**Figure 11.5** The distribution of (imaginary) sample means.

# Bootstrapping: Within ggplot2

```
organdata %>%  
  ggplot(aes(x = country, y = donors)) +  
  stat_summary(fun.y = mean, geom = "point", size = 3) +  
  stat_summary(fun.data = mean_cl_boot, geom = "errorbar", width = 0.5) +  
  coord_flip() +  
  cowplot::theme_cowplot() +  
  labs(x = " ",  
       y = "Organ Donations in 000's",  
       title = "Avg Organ Donations (000s) by Country")
```

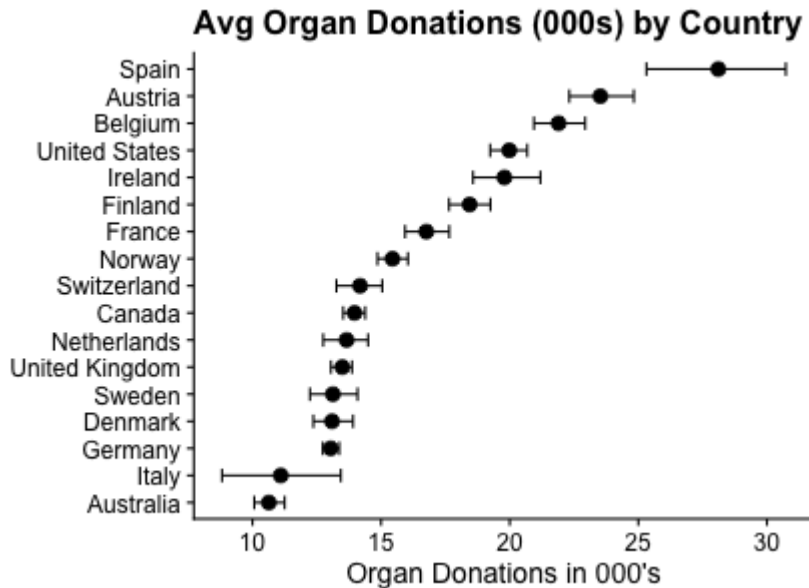


<https://rstudio.cloud/spaces/22733/project/527500>



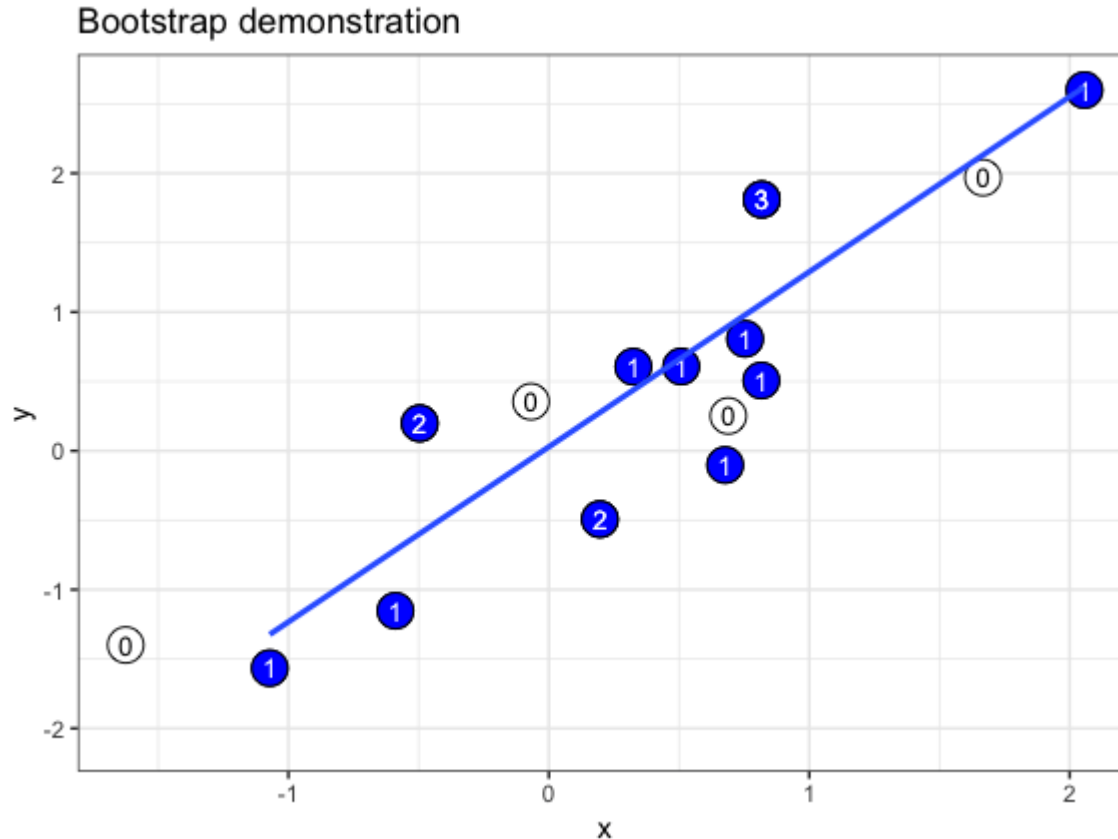
# Bootstrapping: Outside ggplot2

```
organdata %>%  
  group_by(country) %>%  
  do(as_tibble(bind_rows(Hmisc::smean.cl.boot(. $donors)))) %>%  
  ggplot(aes(x = reorder(country, Mean), y = Mean)) +  
  geom_point(size = 3) +  
  geom_errorbar(aes(ymin = Lower, ymax = Upper), width = 0.5) +  
  coord_flip() +  
  cowplot::theme_cowplot() +  
  labs(x = " ", y = "Avg Organ Donations (000's)", title = "Avg Organ Donations (000s) by Country")
```



<https://rstudio.cloud/spaces/22733/project/527500>

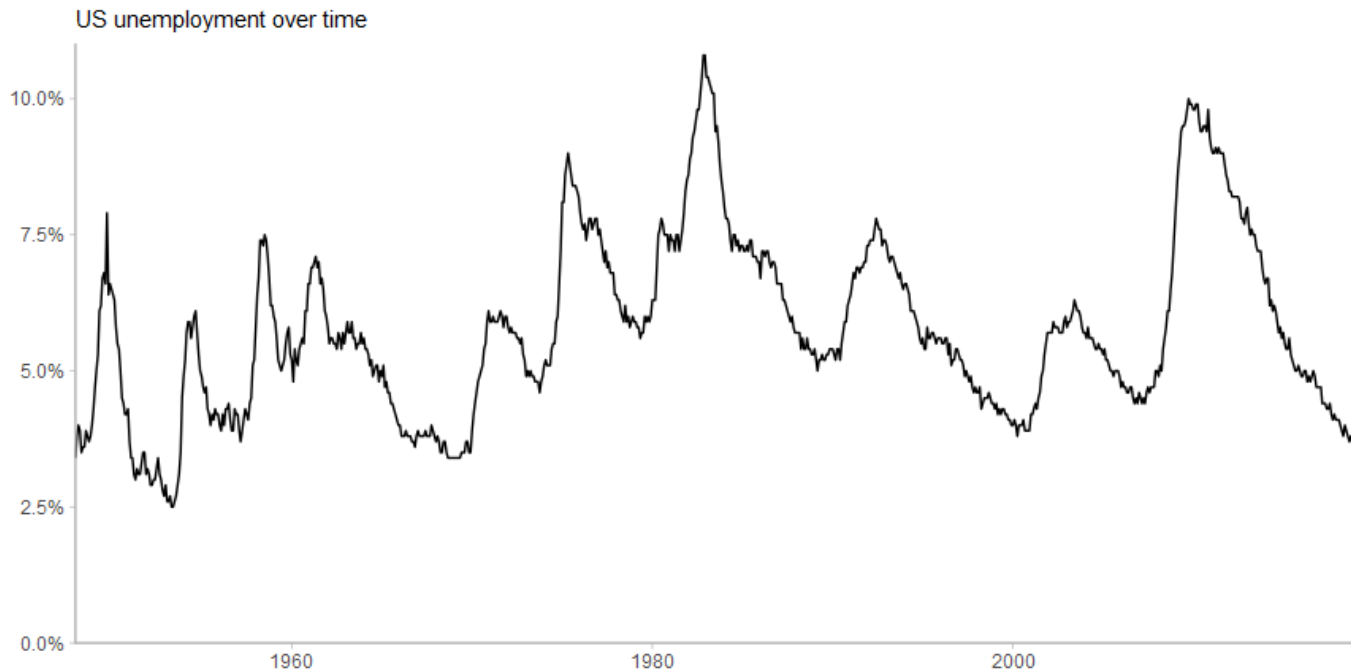
# Bootstrapping with HOPs + gganimate



ungeviz package by Claus Wilke

# Unemployment Rate

```
df %>%  
  ggplot(aes(x = date, y = unemployment)) +  
  geom_line() +  
  coord_cartesian(ylim = c(0, .11), expand = FALSE),  
  scale_y_continuous(labels = scales::percent) +  
  labs(x = NULL, y = NULL, subtitle = "US unemployment over time")
```



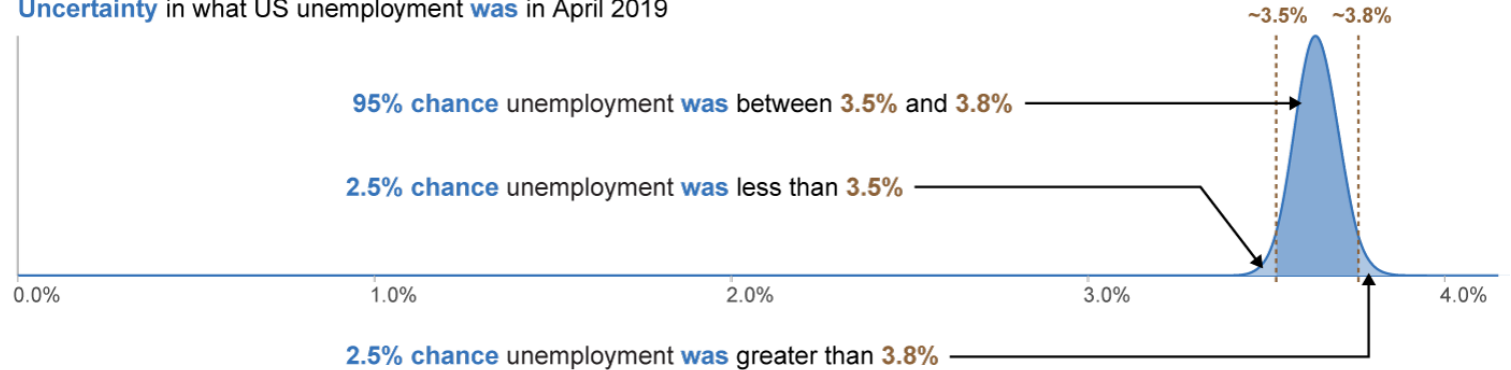
# Unemployment Rate



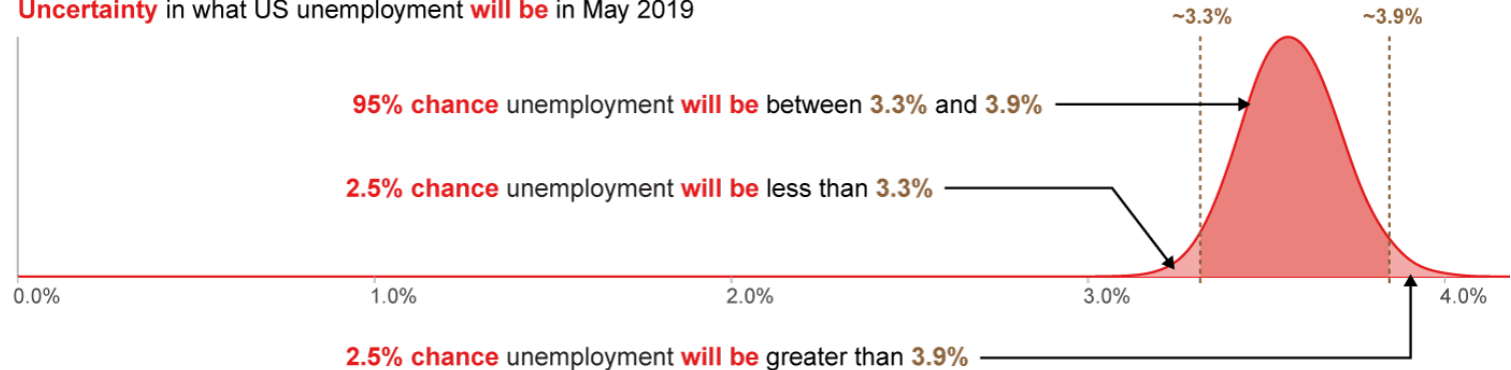
Kay and Hullman Multiple Views Blog 1

# Unemployment Rate

**Uncertainty** in what US unemployment **was** in April 2019

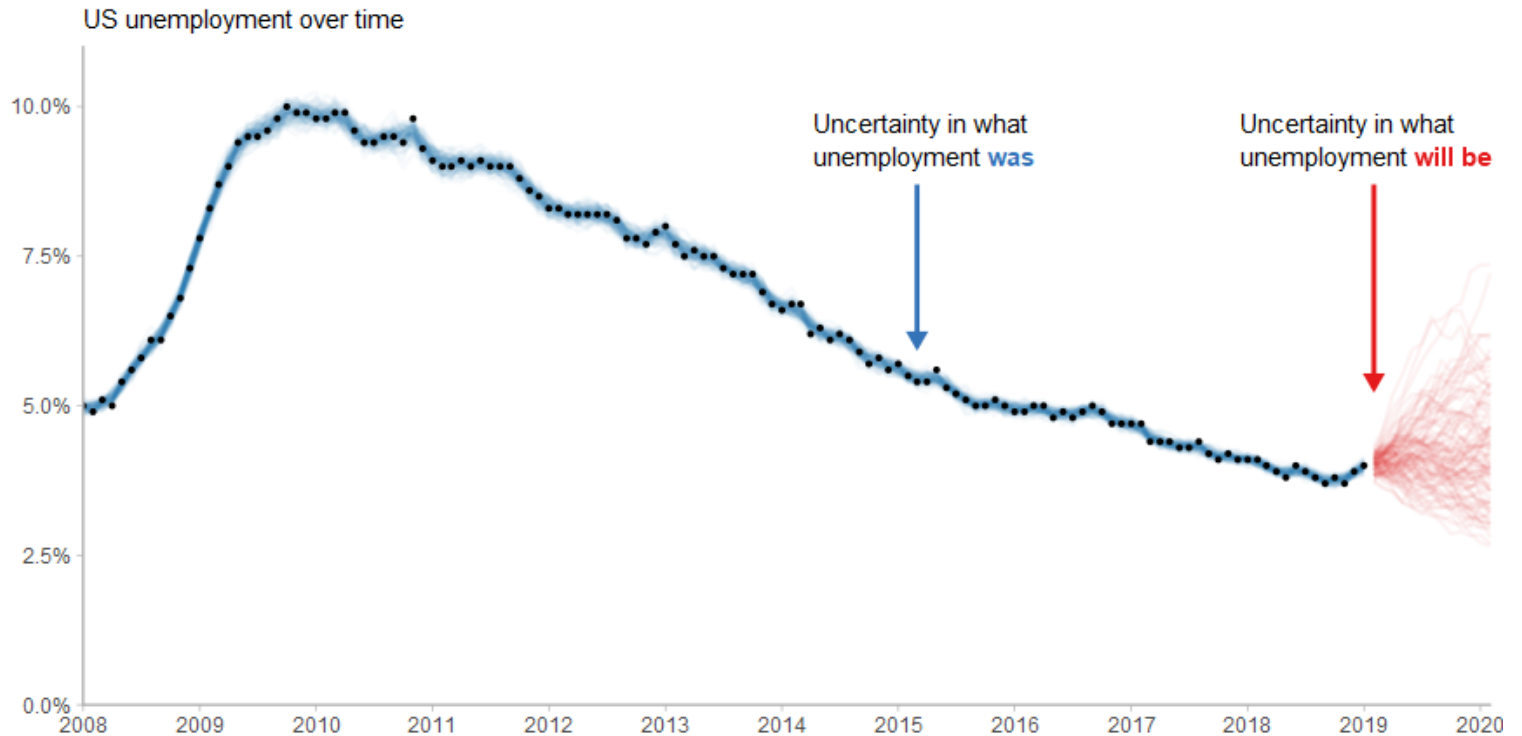


**Uncertainty** in what US unemployment **will be** in May 2019



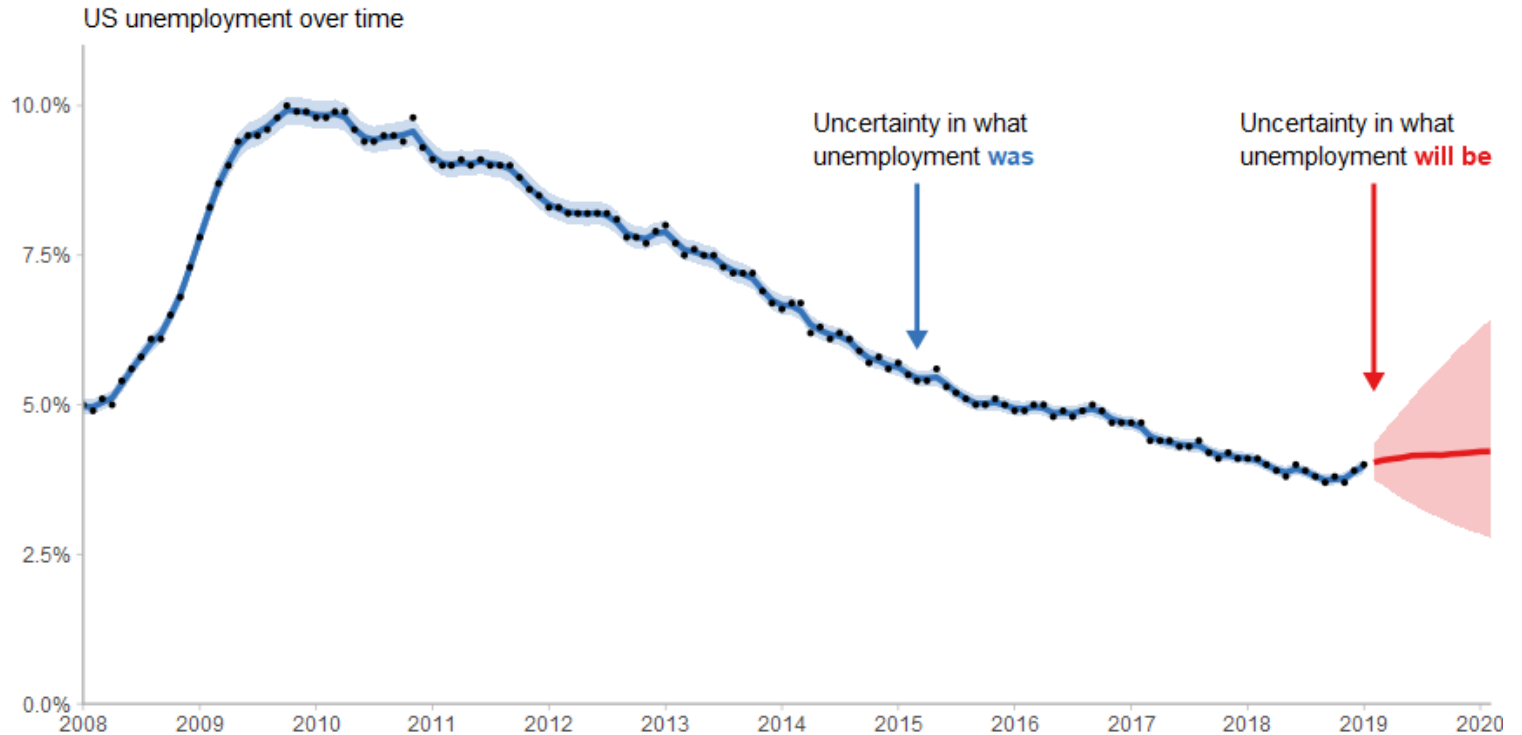
Kay and Hullman Multiple Views Blog 1

# Unemployment Rate



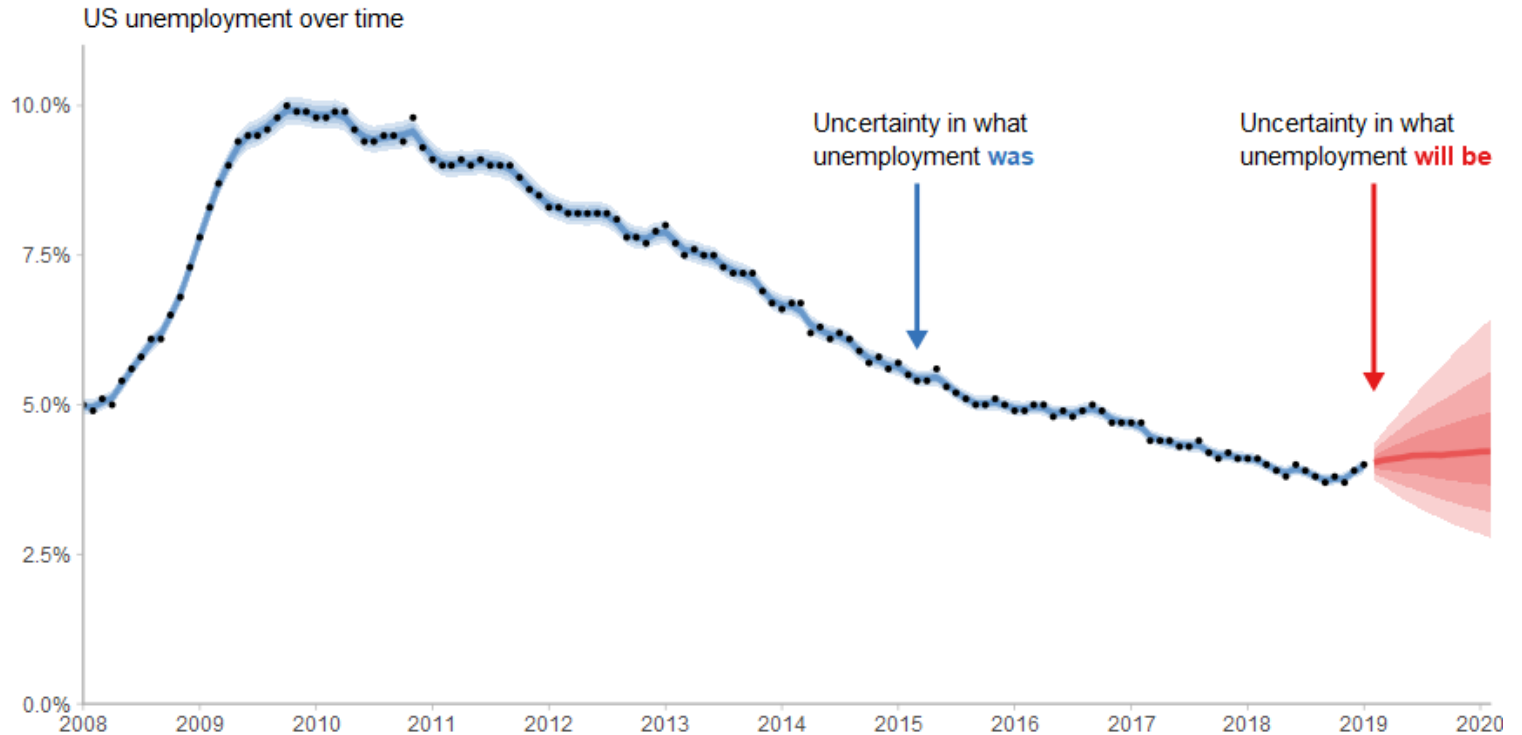
Source: Matthew Kay

# Unemployment Rate



Source: Matthew Kay

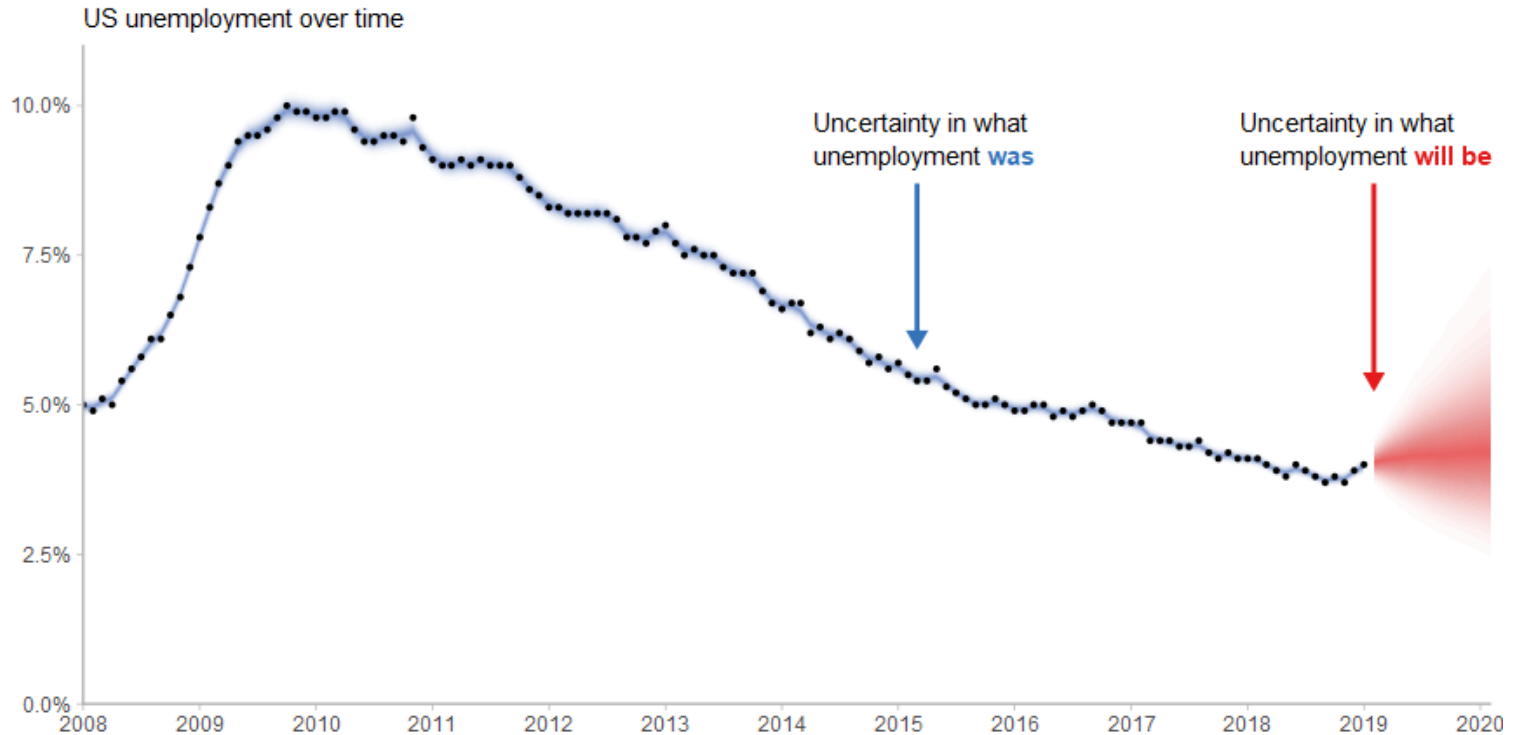
# Unemployment Rate



Source: Matthew Kay

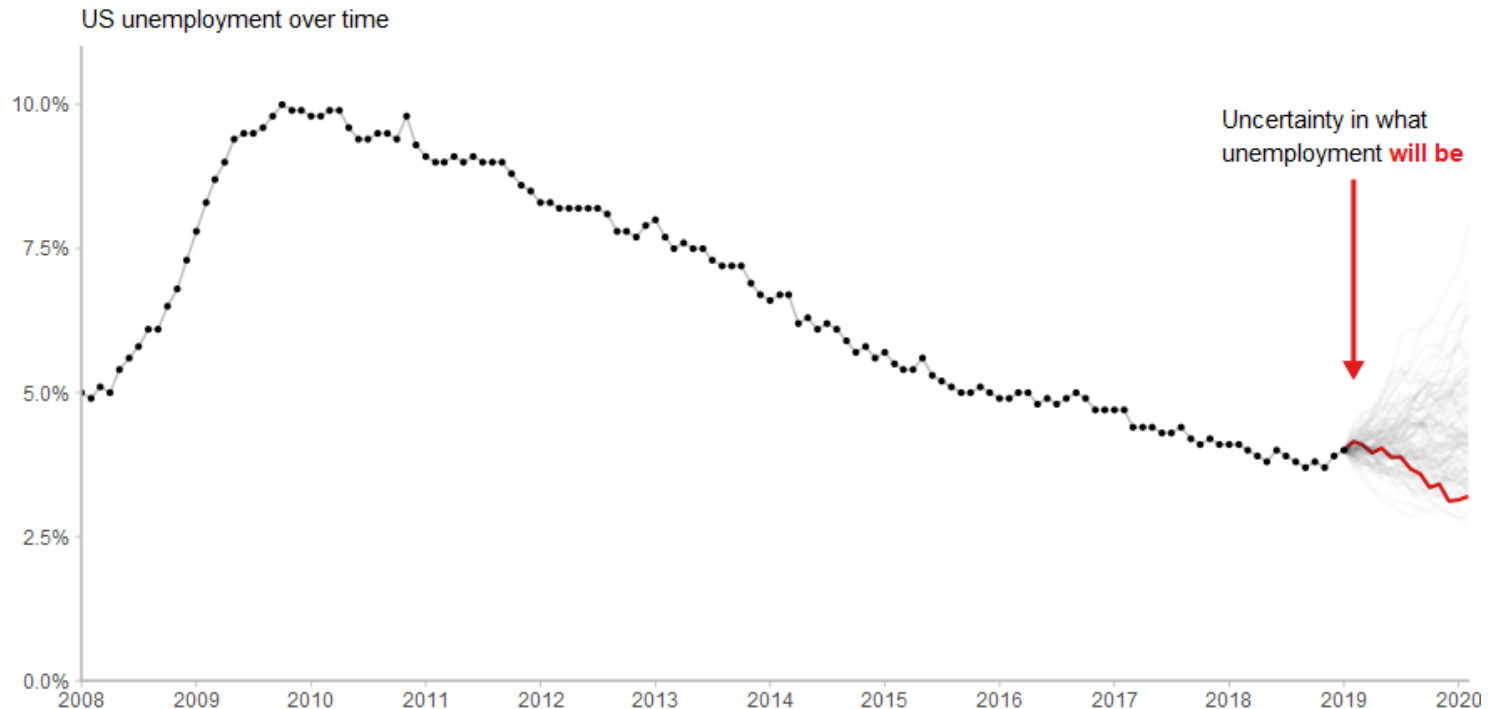


# Unemployment Rate



Source: Matthew Kay

# Unemployment Rate



Hypothetical Outcome Plots (HOP): `tidybayes` & `gganimate`

# Why is visualizing uncertainty hard?

- Efficient encodings for uncertainty can be hard to find.
- Make sure people understand encodings (what does the plot mean?).
- Perceptual models of probability (e.g., quantile dot plot, HOP).
- Decisions under uncertainty (e.g., Gigerenzer et al or Monty Hall problem).
- Findings may not apply in all contexts.
- Plus, you still have to actually build it!