

# Week 7: Factors, Amounts, & Proportions

m EMSE 4572 / 6572: Exploratory Data Analysis

2 John Paul Helveston

**October 08, 2025** 

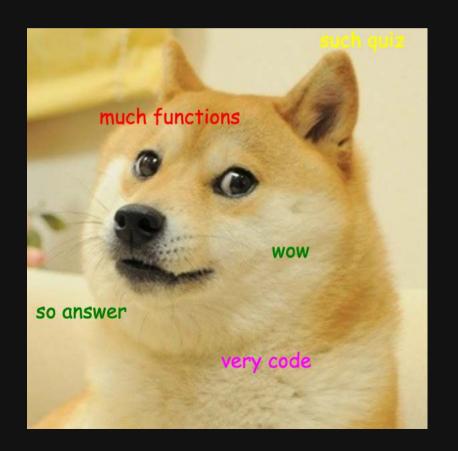
## Quiz 2

12:00

## Write your name on the quiz!

#### Rules:

- Work alone; no outside help of any kind is allowed.
- No calculators, no notes, no books, no computers, no phones.



# Next projects due:

- Mini project 2: Exploring Data (Due 10/14)
- Project Progress Report (Due 10/26)

### Today's data

## New packages

The {waffle} package

```
install.packages("waffle")
```

## Week 7: Factors, Amounts, & Proportions

- 1. Manipulating factors
- 2. Graphing amounts

**BREAK** 

3. Graphing proportions

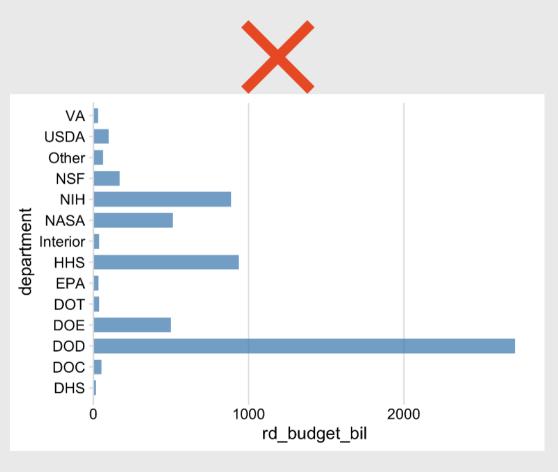
## Week 7: Factors, Amounts, & Proportions

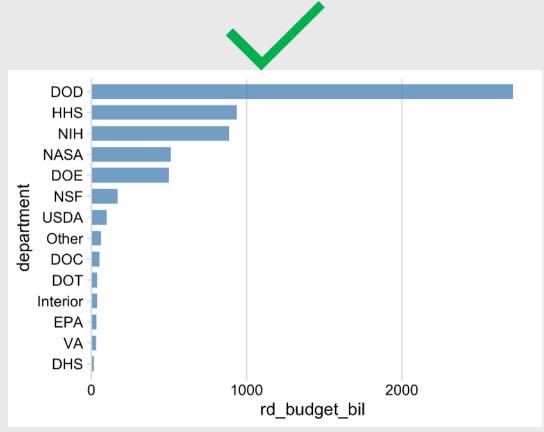
- 1. Manipulating factors
- 2. Graphing amounts

**BREAK** 

3. Graphing proportions

# Sorting in ggplot is done by reordering factors

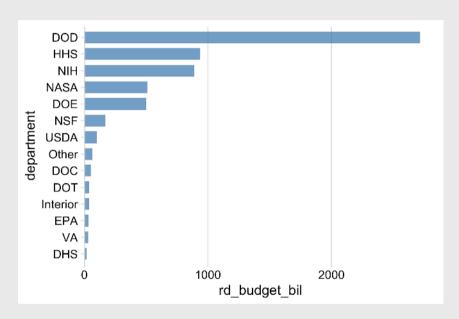




### Two ways to sort

#### Method 1: Use reorder() inside aesthetic mapping

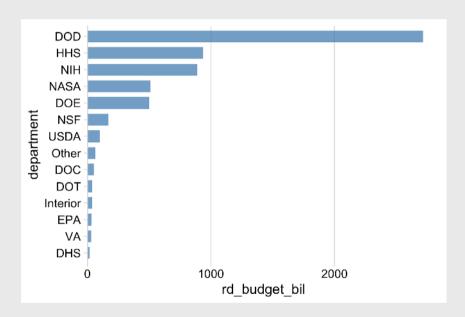
```
# Format the data frame
federal spending %>%
  group_by(department) %>%
  summarise(
    rd budget bil = sum(rd budget mil) / 10^3) %>%
# Make the chart
  ggplot() +
  geom_col(
    aes (
      x = rd budget bil,
      y = reorder(department, rd budget bil)
    width = 0.7, alpha = 0.8,
    fill = "steelblue"
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal vgrid()
```



### Two ways to sort

#### Method 2: Use fct\_reorder() when formatting the data frame

```
# Format the data frame
federal spending %>%
  group by(department) %>%
  summarise(
    rd budget bil = sum(rd budget mil) / 10^3) %>%
 mutate(
    department = fct_reorder(department, rd_budget_bil)
  ) %>%
# Make the chart
  ggplot() +
  geom col(
    aes(x = rd budget bil, y = department),
   width = 0.7, alpha = 0.8,
    fill = "steelblue"
  scale_x_continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme_minimal_vgrid()
```



# Reorder & modify factors with the **forcats** library

Loaded with library(tidyverse)



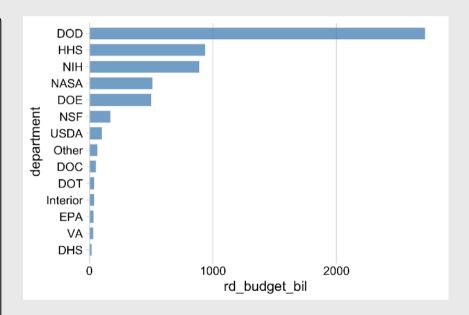
## Common situations for modifying / reording factors:

- 1. Reorder factors based on another numerical variable
- 2. Reorder factors manually
- 3. Modify factors manually
- 4. What if there are too many factor levels?

#### 1. Reorder factors based on another numerical variable

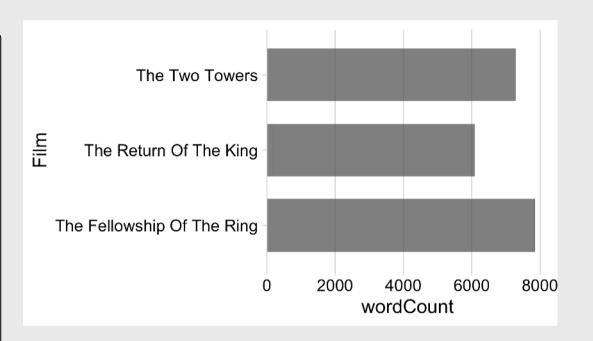
#### Use fct\_reorder()

```
# Format the data frame
federal_spending %>%
  group_by(department) %>%
  summarise(
    rd budget bil = sum(rd budget mil) / 10^3) %>%
  mutate(
    department = fct reorder(department, rd budget bil)
  ) %>%
# Make the chart
  ggplot() +
  geom col(
    aes(x = rd_budget_bil, y = department),
    width = 0.7, alpha = 0.8,
    fill = "steelblue"
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal vgrid()
```



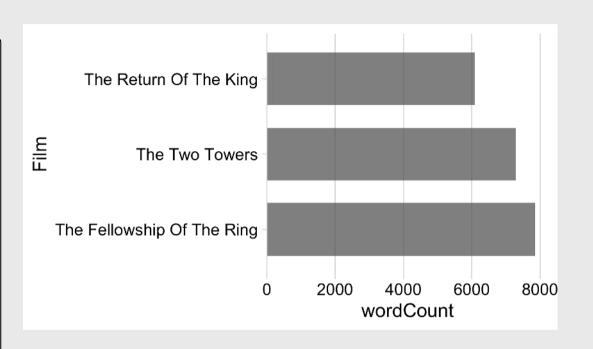
#### 2. Reorder factors **manually**

```
# Format the data frame
lotr words %>%
 pivot_longer(
      names_to = 'gender',
      values to = 'wordCount',
      cols = Female:Male) %>%
# Make the chart
  ggplot() +
  geom_col(
    aes(x = wordCount, y = Film),
    width = 0.7, alpha = 0.8
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal vgrid()
```



#### 2. Reorder factors **manually** with fct\_relevel()

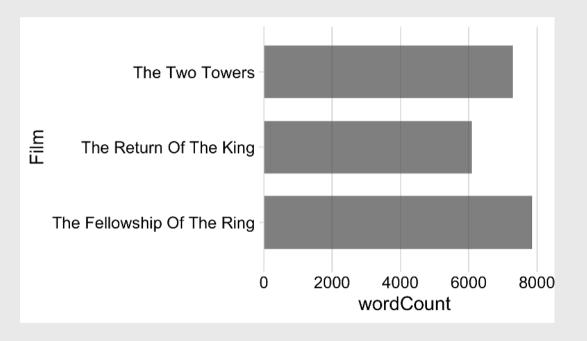
```
# Format the data frame
lotr words %>%
  pivot longer(
      names to = 'gender',
      values to = 'wordCount',
      cols = Female:Male) %>%
 mutate(
    Film = fct relevel(Film, levels = c(
      'The Fellowship Of The Ring',
      'The Two Towers',
      'The Return Of The King'))) %>%
# Make the chart
  ggplot() +
  geom col(
    aes(x = wordCount, y = Film),
    width = 0.7, alpha = 0.8
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal vgrid()
```



#### 3. Modify factors manually

```
# Format the data frame
lotr words %>%
  pivot_longer(
      names_to = 'gender',
      values to = 'wordCount',
      cols = Female:Male) %>%
# Make the chart
  ggplot() +
  geom col(
    aes(x = wordCount, y = Film),
    width = 0.7, alpha = 0.8
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme_minimal_vgrid()
```

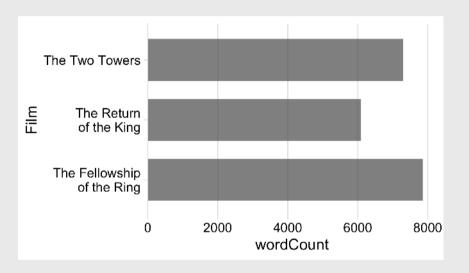
#### The film names here are too long



#### 3. Modify factors manually with fct\_recode()

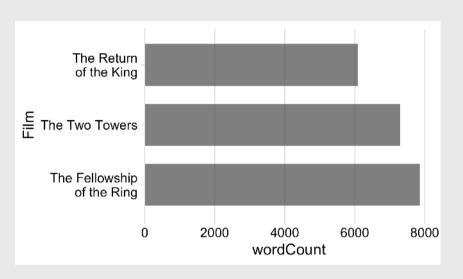
"new label" = "old label"

```
# Format the data frame
lotr words %>%
  pivot longer(
      names to = 'gender',
      values to = 'wordCount',
      cols = Female:Male) %>%
  mutate(
    Film = fct_recode(Film,
      'The Fellowship\nof the Ring' = 'The Fellowship Of
      'The Return\nof the King' = 'The Return Of The King'
# Make the chart
  qqplot() +
  geom col(
    aes(x = wordCount, y = Film),
    width = 0.7, alpha = 0.8
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal vgrid()
```



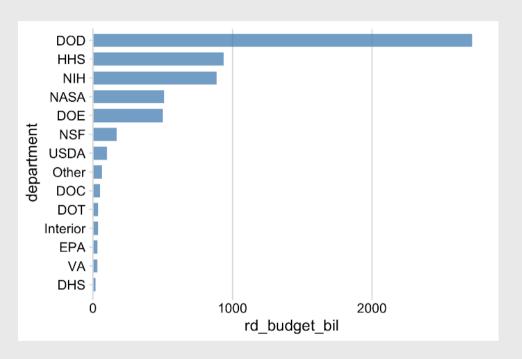
#### 2 & 3. Modify and reorder factors manually

```
# Format the data frame
lotr words %>%
  pivot longer(
      names to = 'gender',
      values to = 'wordCount',
      cols = Female:Male) %>%
 mutate(
    Film = fct relevel(Film, levels = c(
      'The Fellowship Of The Ring',
      'The Two Towers',
      'The Return Of The King')),
    Film = fct recode(Film,
      'The Fellowship\nof the Ring' = 'The Fellowship Of
      'The Return\nof the King' = 'The Return Of The King'
# Make the chart
  ggplot() +
  geom col(
      aes(x = wordCount, y = Film),
      width = 0.7, alpha = 0.8
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal vgrid()
```



#### 4. What if there are too many factor levels?

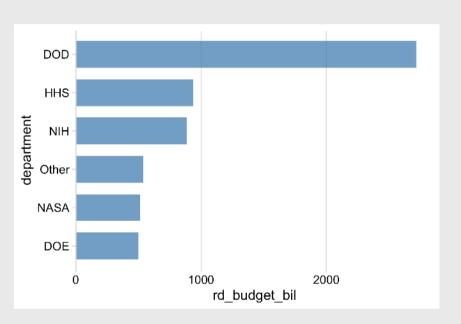
```
# Format the data frame
federal spending %>%
  group_by(department) %>%
  summarise(
    rd_budget_bil = sum(rd_budget_mil) / 10^3) %>%
  mutate(
    department = fct_reorder(department, rd_budget_
  ) %>%
# Make the chart
  ggplot() +
  geom col(
    aes(x = rd_budget_bil, y = department),
   width = 0.7, alpha = 0.8,
    fill = "steelblue"
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal vgrid()
```



#### 4. What if there are too many factor levels?

**Strategy**: Merge smaller factors into "Other" with fct\_other()

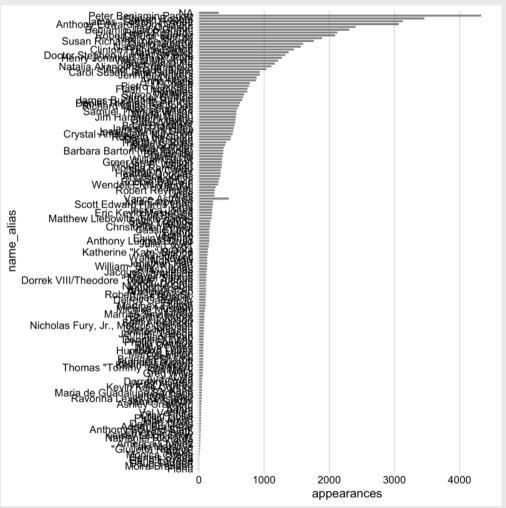
```
# Format the data frame
federal spending %>%
  mutate(
    department = fct other(department,
      keep = c('DOD', 'HHS', 'NIH', 'NASA', 'DOE'))) %>%
  group by(department) %>%
  summarise(
    rd budget bil = sum(rd budget mil) / 10^3) %>%
  mutate(
    department = fct reorder(department, rd budget bil))
# Make the chart
  aaplot() +
  geom col(
    aes(x = rd_budget_bil, y = department),
    width = 0.7, alpha = 0.8,
    fill = "steelblue"
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal vgrid()
```



#### 4. What if there are *really* too many factor levels?

```
# Format the data frame
avengers %>%
  mutate(
    name_alias = fct_reorder(name_alias, appear)

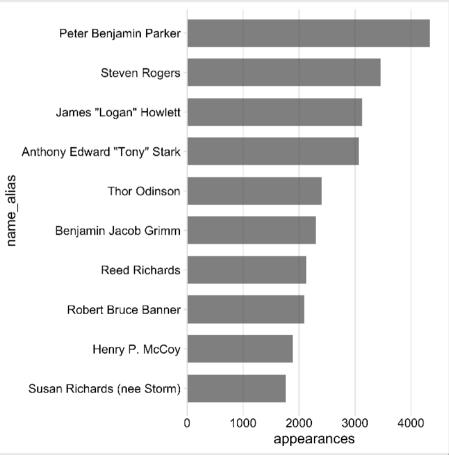
# Make the chart
  ggplot() +
  geom_col(
    aes(x = appearances, y = name_alias),
    width = 0.7, alpha = 0.8
) +
  scale_x_continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme_minimal_vgrid()
```



#### 4. What if there are *really* too many factor levels?

**Strategy**: Keep top N, drop the rest with slice()

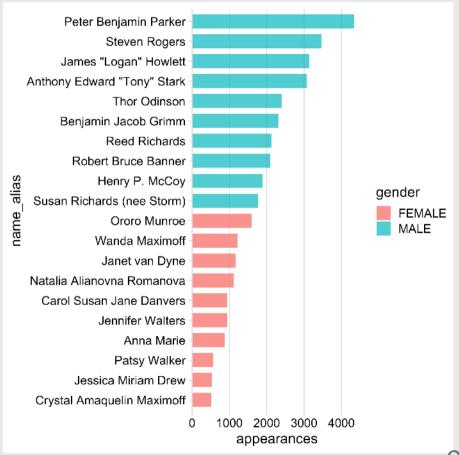
```
# Format the data frame
avengers %>%
  mutate(
    name alias = fct reorder(name alias, appear
  arrange(desc(appearances)) %>%
  slice(1:10) %>%
# Make the chart
  ggplot() +
  geom col(
    aes(x = appearances, y = name_alias),
    width = 0.7, alpha = 0.8
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal vgrid()
```



#### 4. What if there are *really* too many factor levels?

#### slice() works with grouping too!

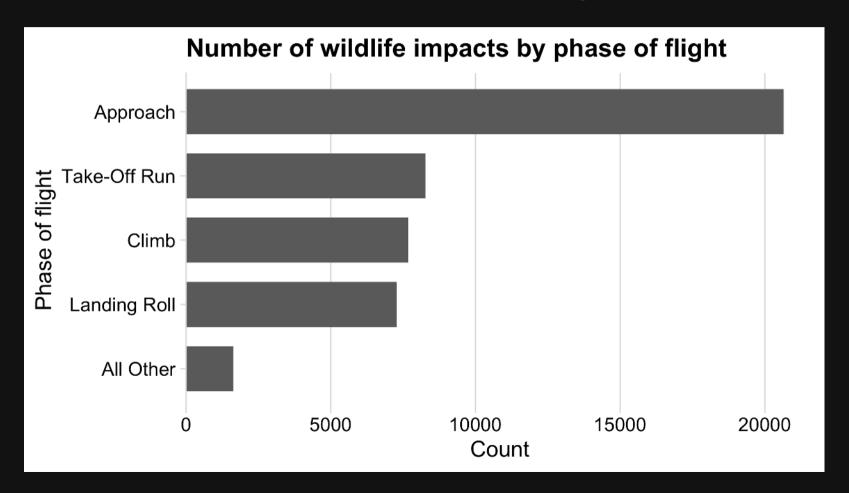
```
# Format the data frame
avengers %>%
  mutate(
    name alias = fct reorder(name alias, appear
  arrange(desc(appearances)) %>%
  group by(gender) %>%
  slice(1:10) %>%
# Make the chart
  ggplot() +
  geom col(
    aes (
      x = appearances,
      y = name_alias,
      fill = gender
    width = 0.7, alpha = 0.8
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal vgrid()
```



## Your turn - practice manipulating factors

15:00

Use the wildlife\_impacts data to create the following plot



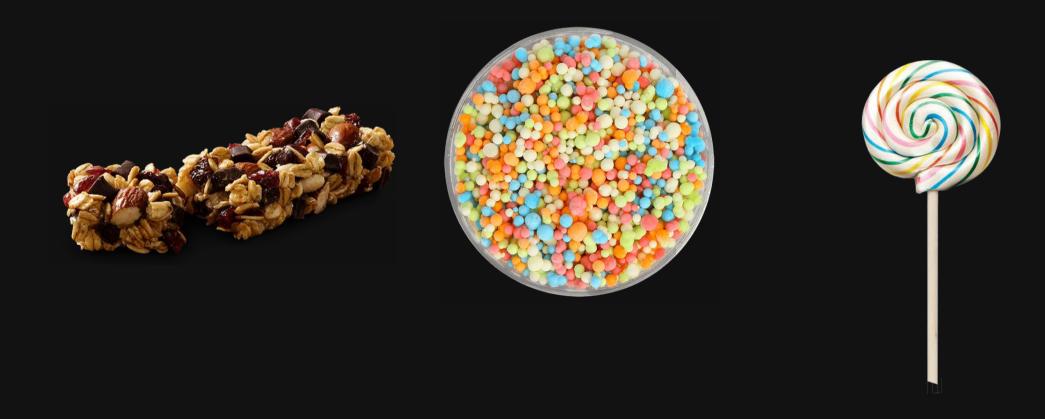
## Week 7: Factors, Amounts, & Proportions

- 1. Manipulating factors
- 2. Graphing amounts

**BREAK** 

3. Graphing proportions

# Show amounts with:





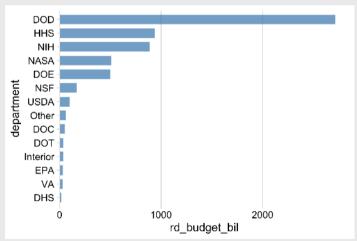


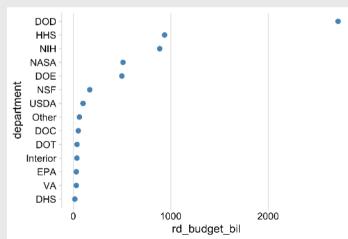


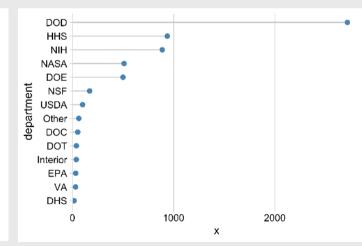
Bar chart

Dot chart

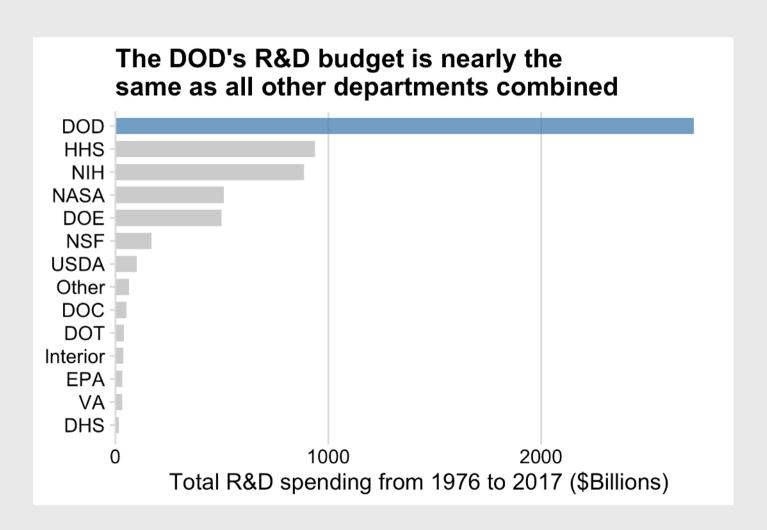
Lollipop chart





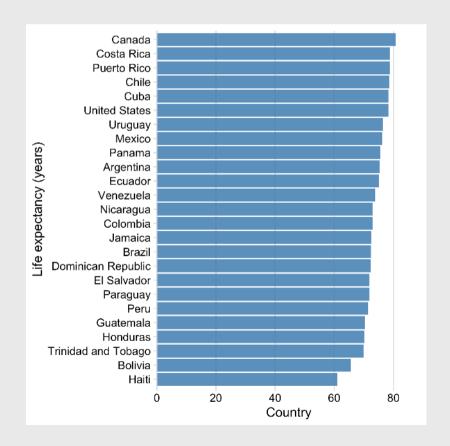


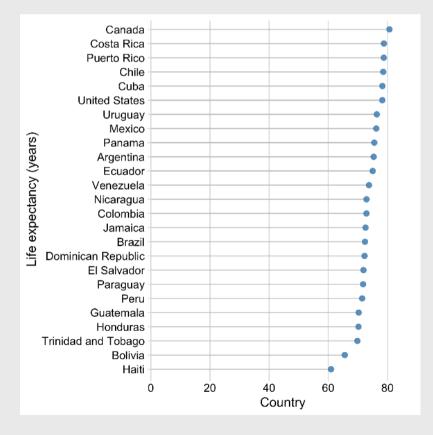
## Bars are good for highlighting specific categories



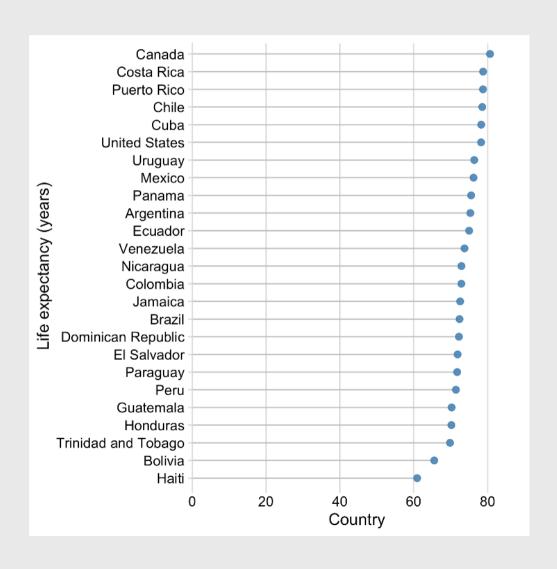
## Use lollipops when:

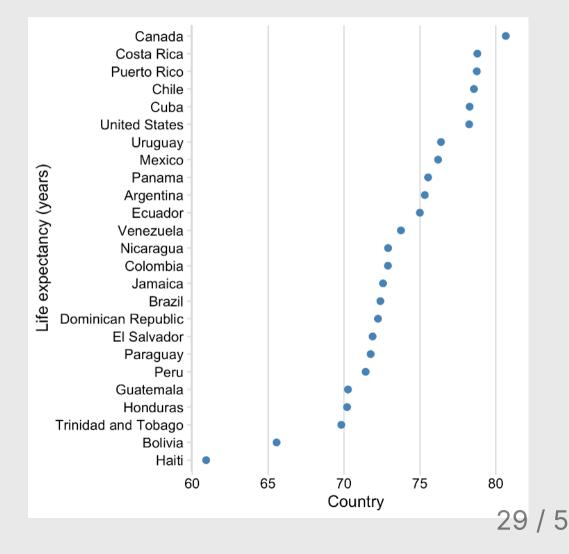
- The bars are overwhelming
- You're not highlighting categories





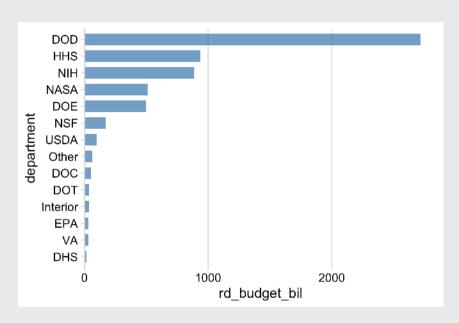
#### Or use dots and don't set axis to 0





#### How to make a **Bar chart**

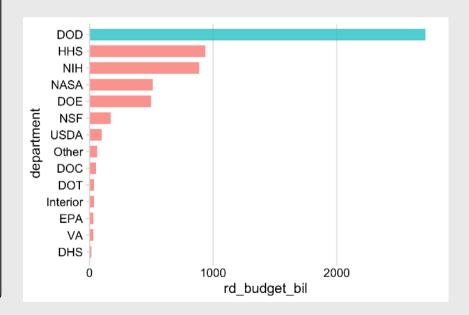
```
# Summarize the data
federal_spending %>%
    group_by(department) %>%
    summarise(rd_budget_bil = sum(rd_budget_mil) / 10^3) %;
    mutate(department = fct_reorder(department, rd_budget_
# Make chart
    ggplot() +
    geom_col(
        aes(x = rd_budget_bil, y = department),
        width = 0.7, alpha = 0.8,
        fill = 'steelblue') +
    scale_x_continuous(
        expand = expansion(mult = c(0, 0.05))) +
    theme_minimal_vgrid()
```



## Filling the bars with color

```
# Summarize the data
federal spending %>%
  group by(department) %>%
  summarise(rd budget bil = sum(rd budget mil) / 10^3) %
  mutate(
    department = fct_reorder(department, rd_budget_bil),
    is dod = if else(
      department == 'DOD', TRUE, FALSE)) %>%
# Make the chart
  ggplot() +
  geom col(
    aes(x = rd_budget_bil, y = department,
        fill = is dod),
   width = 0.7, alpha = 0.8) +
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal vgrid() +
  theme(legend.position = 'none')
```

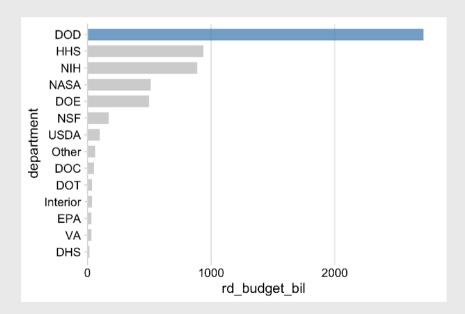
# The DOD's R&D budget is nearly the same as all other departments combined



## Filling the bars with color

```
# Summarize the data
federal spending %>%
  group by(department) %>%
  summarise(rd budget bil = sum(rd budget mil) / 10^3) %
 mutate(
    department = fct reorder(department, rd budget bil),
    is dod = if else(
      department == 'DOD', TRUE, FALSE)) %>%
# Make the chart
  ggplot() +
  geom col(
    aes(x = rd_budget_bil, y = department,
        fill = is dod),
    width = 0.7, alpha = 0.8) +
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  scale_fill_manual(values = c('grey', 'steelblue')) +
  theme_minimal_vgrid() +
  theme(legend.position = 'none')
```

# The DOD's R&D budget is nearly the same as all other departments combined



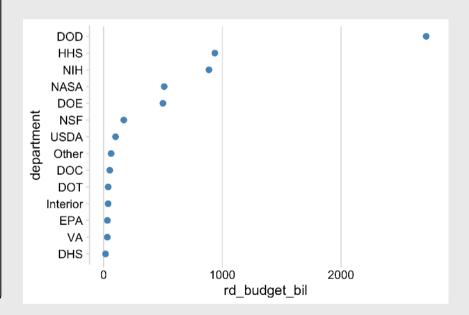
#### How to make a **Dot chart**

#### Summarize data frame:

```
# Summarize the data
federal_spending %>%
    group_by(department) %>%
    summarise(
    rd_budget_bil = sum(rd_budget_mil) / 10^3) %>%
    mutate(
    department = fct_reorder(department, rd_budget_bil))

# Make the chart
ggplot() +
geom_point(
    aes(x = rd_budget_bil, y = department),
    size = 2.5, color = 'steelblue') +
theme_minimal_vgrid()
```

# **Dot chart** of federal R&D spending by department

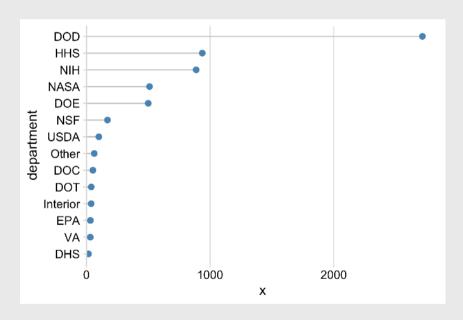


## How to make a **Lollipop chart**

#### Summarize data frame:

```
# Summarize the data
federal spending %>%
  group by(department) %>%
  summarise(
    rd budget bil = sum(rd budget mil) / 10^3) %>%
  mutate(
    department = fct_reorder(department, rd_budget_bil))
# Make the chart
  ggplot() +
  geom_segment(
    aes(x = 0, xend = rd_budget_bil,
        y = department, yend = department),
    color = 'grev') +
  geom point(
    aes(x = rd_budget_bil, y = department),
    size = 2.5, color = 'steelblue') +
  theme_minimal_vgrid()
```

# **Lollipop chart** of federal R&D spending by department

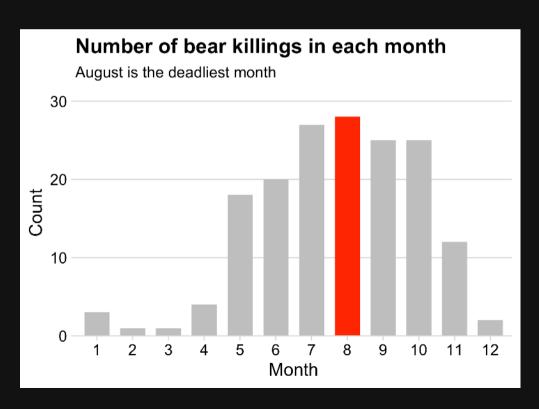


## Your turn - practice plotting amounts

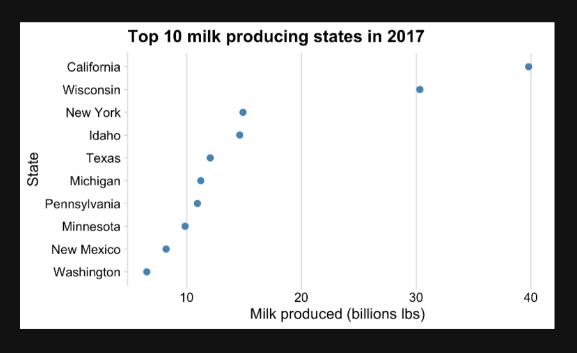


Create the following charts:

Data: bears



Data: milk\_production



## Break!

Stand up, Move around, Stretch!



# Week 7: Factors, Amounts, & Proportions

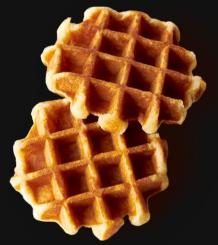
- 1. Manipulating factors
- 2. Graphing amounts

**BREAK** 

3. Graphing proportions

# Show proportions with:



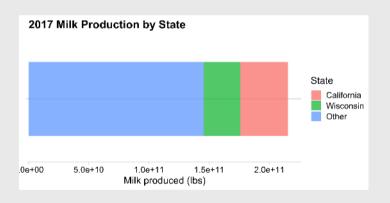




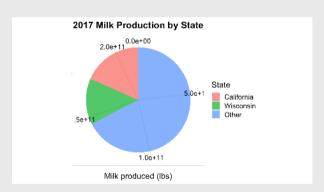




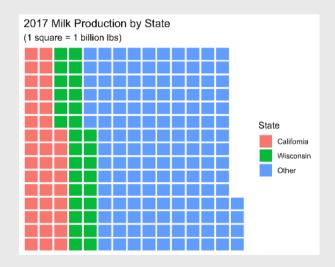
### Bar charts



## Pie charts

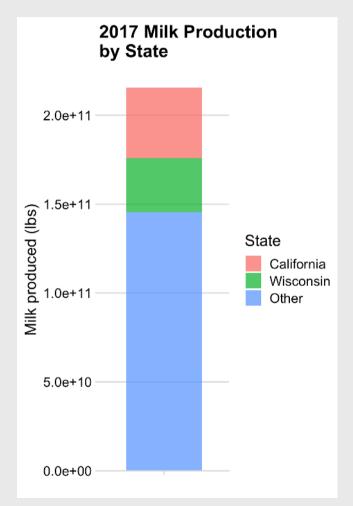


### Waffle charts



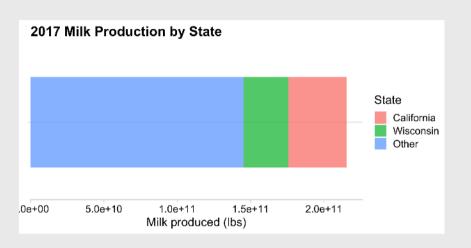
### Stacked bars

```
# Format the data
milk production %>%
  filter(year == 2017) %>%
 mutate(state = fct_other(state,
    keep = c('California', 'Wisconsin'))) %>%
  group by(state) %>%
  summarise(milk produced = sum(milk produced)) %>%
# Make the chart
  ggplot() +
  geom_col(
    aes(x = "", y = milk_produced, fill = state),
    width = 0.7, alpha = 0.8) +
  scale_y_continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal hgrid() +
  labs(x = NULL,
       y = 'Milk produced (lbs)',
       fill = 'State',
       title = '2017 Milk Production\nby State')
```



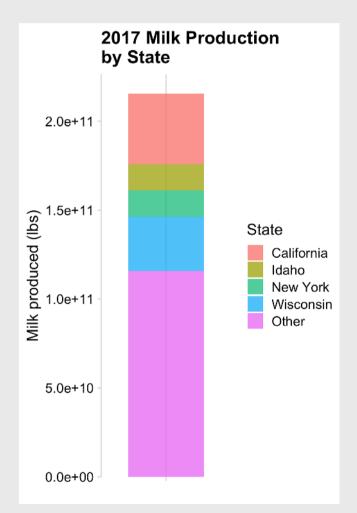
## Stacked bars - Rotated also looks good

```
# Format the data
milk production %>%
  filter(year == 2017) %>%
 mutate(state = fct other(state,
    keep = c('California', 'Wisconsin'))) %>%
  group by(state) %>%
  summarise(milk produced = sum(milk produced)) %>%
# Make the chart
  ggplot() +
  geom col(
    aes(x = milk_produced, y = "", fill = state),
    width = 0.7, alpha = 0.8) +
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal hgrid() +
  labs(y = NULL,
       x = 'Milk produced (lbs)',
       fill = 'State',
       title = '2017 Milk Production by State')
```



## Stacked bars - not great for more than a few categories

```
# Format the data
milk production %>%
  filter(year == 2017) %>%
 mutate(state = fct other(state,
    keep = c('California', 'Wisconsin',
             'New York', 'Idaho'))) %>%
  group by(state) %>%
  summarise(milk produced = sum(milk produced))
# Make the chart
  ggplot() +
  geom col(
    aes(x = "", y = milk_produced, fill = state),
    width = 0.7, alpha = 0.8) +
  scale y continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal vgrid() +
  labs(x = NULL,
       y = 'Milk produced (lbs)',
       fill = 'State',
       title = '2017 Milk Production\nby State')
```

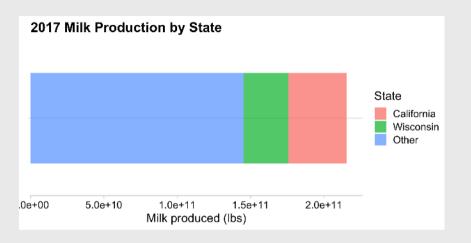


## Dodged bars

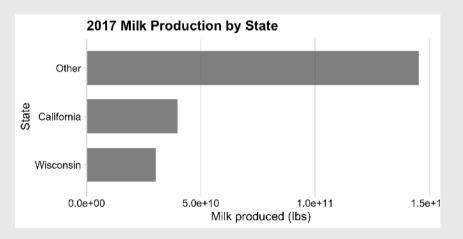
#### Better for part-to-whole comparison

```
# Format the data
milk production %>%
  filter(year == 2017) %>%
 mutate(state = fct other(state,
    keep = c('California', 'Wisconsin'))) %>%
  group by(state) %>%
  summarise(milk produced = sum(milk produced)) %>%
  mutate(state = fct reorder(state, milk produced)) %>%
# Make the chart
  ggplot() +
  geom col(
    aes(x = milk_produced, y = state),
    width = 0.7, alpha = 0.8) +
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal vgrid() +
  labs(x = 'Milk produced (lbs)',
       y = 'State',
       title = '2017 Milk Production by State')
```

#### Okay:



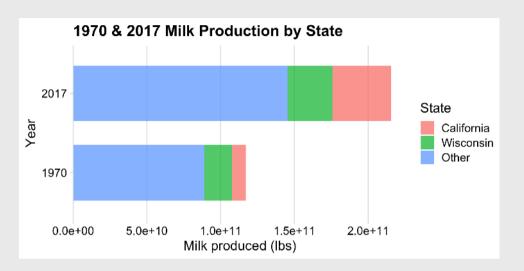
#### Better:



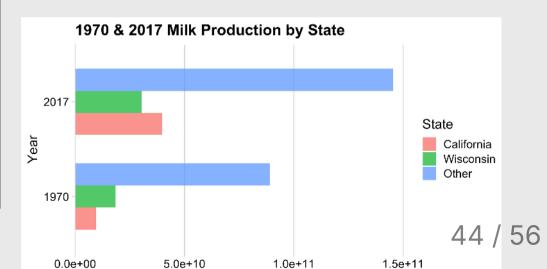
## Dodged bars

```
milk production %>%
  filter(year %in% c(1970, 2017)) %>%
  mutate(state = fct other(state,
    keep = c('California', 'Wisconsin'))) %>%
  group by(year, state) %>%
  summarise(milk produced = sum(milk produced)) %>%
# Make the chart
  ggplot() +
  geom col(
    aes(x = milk_produced,
        y = as.factor(year),
        fill = state),
    position = 'dodge',
   width = 0.7, alpha = 0.8) +
  scale x continuous(
    expand = expansion(mult = c(0, 0.05))) +
  theme minimal vgrid() +
  labs(x = 'Milk produced (lbs)',
       y = 'Year',
       fill = 'State',
       title = '1970 & 2017 Milk Production by State
```

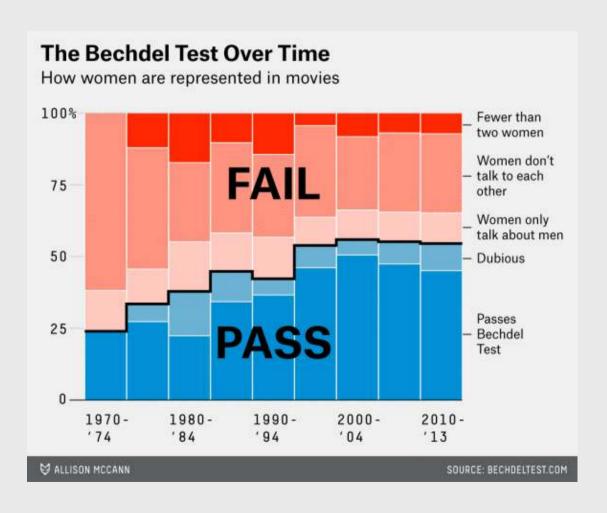
#### Better for comparing **total**:



#### Better for comparing **parts**:

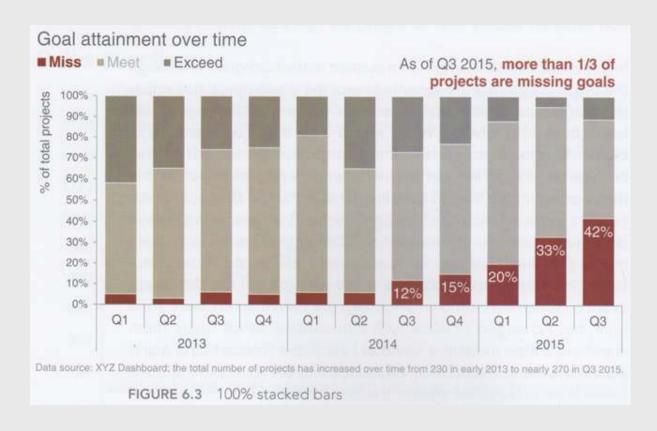


## Where stacking is useful



- 2 to 3 groups
- Proportions over time

## Where stacking is useful



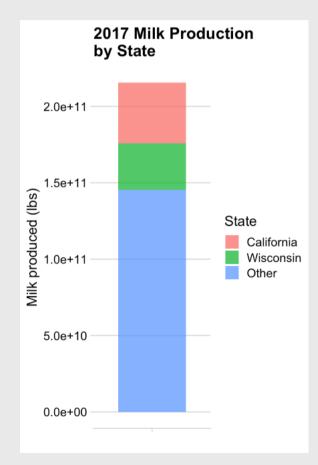
- 2 to 3 groups
- Proportions over time

https://www.perceptualedge.com/blog/?p=2239

### The Notorious P-I-E

#### Start with a bar chart

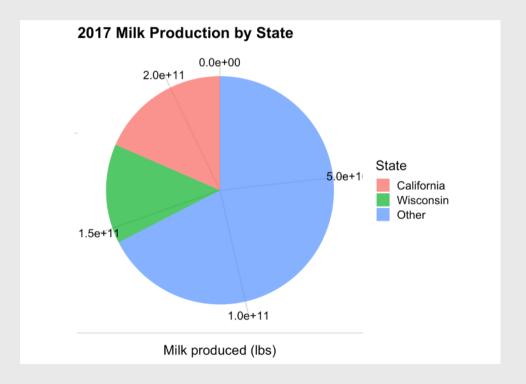
```
# Format the data
milk production %>%
  filter(year == 2017) %>%
  mutate(state = fct_other(state,
    keep = c('California', 'Wisconsin'))) %>%
  group by(state) %>%
  summarise(milk produced = sum(milk produced)) %>%
# Make the chart
  ggplot() +
  geom col(
    aes(x = "", y = milk_produced, fill = state),
    width = 0.7, alpha = 0.8) +
  theme minimal hgrid() +
  labs(x = NULL,
       y = 'Milk produced (lbs)',
       fill = 'State',
       title = '2017 Milk Production\nby State')
```



### The Notorious P-I-E

Convert bar to pie with coord\_polar()

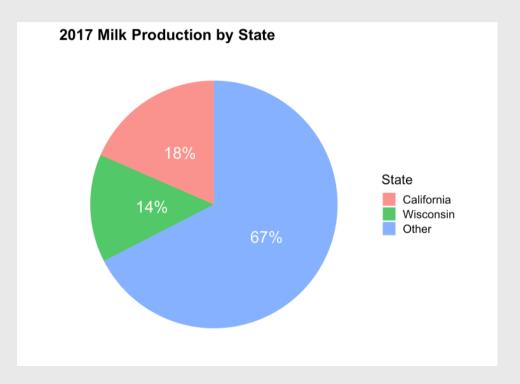
```
# Format the data
milk production %>%
  filter(year == 2017) %>%
  mutate(state = fct other(state,
    keep = c('California', 'Wisconsin'))) %>%
  group by(state) %>%
  summarise(milk_produced = sum(milk_produced)) %>%
# Make the chart
  ggplot() +
  geom col(
    aes(x = "", y = milk_produced, fill = state),
    width = 0.7, alpha = 0.8) +
  coord polar(theta = "y") +
  theme minimal hgrid() +
  labs(x = NULL,
       y = 'Milk produced (lbs)',
       fill = 'State',
       title = '2017 Milk Production by State')
```



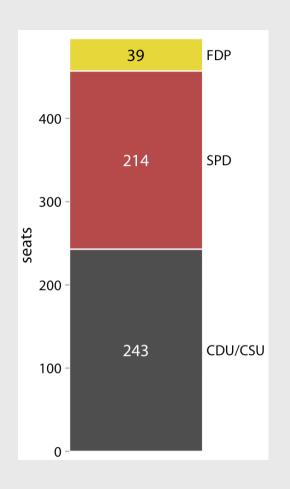
```
# Format the data
milk production %>%
  filter(year == 2017) %>%
  mutate(state = fct other(state,
    keep = c('California', 'Wisconsin'))) %>%
  group by(state) %>%
  summarise(milk produced = sum(milk produced)) %>%
  arrange(desc(state)) %>%
 mutate(p = 100*(milk produced / sum(milk produced)
         label = str c(round(p), '%')) %>%
# Make the chart
  ggplot() +
  geom col(
    aes(x = "", y = milk_produced, fill = state),
    width = 0.7, alpha = 0.8) +
  geom text(
    aes(x = "", y = milk_produced, label = label),
    color = "white", size = 6,
    position = position_stack(vjust = 0.5)) +
  coord_polar(theta = "y") +
  theme map() +
  labs(x = NULL,
       y = NULL
       fill = 'State',
       title = '2017 Milk Production by State')
```

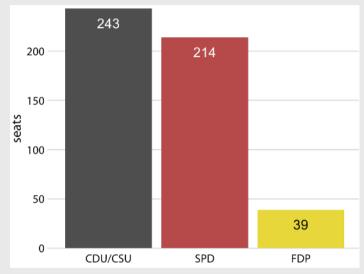
#### The Notorious P-I-E

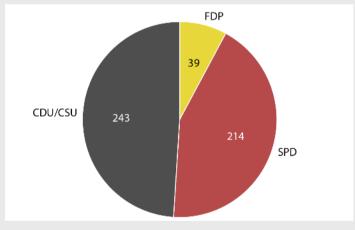
Final chart with labels & theme\_map()



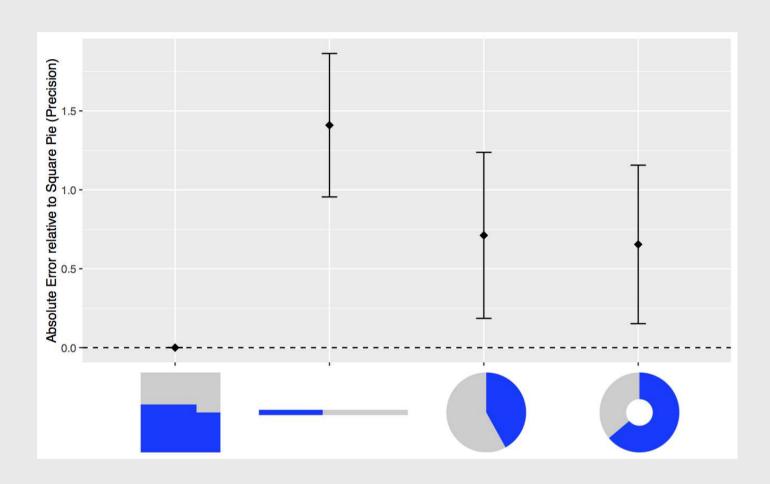
## Pies are still useful if the sum of components matters







## The best pies are square pies

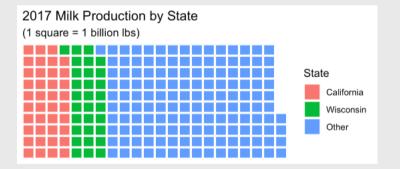


### Waffle plots

```
library(waffle)
# Format the data
milk production %>%
  filter(year == 2017) %>%
 mutate(state = fct other(state,
    keep = c('California', 'Wisconsin'))) %>%
  group by(state) %>%
  summarise(milk produced = sum(milk produced)) %>%
 mutate(milk produced = milk produced / 10^9) %>%
# Make the chart
  ggplot() +
  geom waffle(
    aes(fill = state, values = milk produced),
    color = "white", size = 1, n_rows = 10) +
  scale x discrete(expand = c(0, 0)) +
  scale y discrete(expand = c(0, 0)) +
  theme_minimal() +
  labs(fill = 'State',
       x = NULL, y = NULL,
       title = '2017 Milk Production by State',
       subtitle = '(1 square = 1 billion lbs)')
```

#### Use values between 100 - 1,000

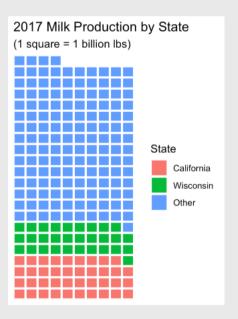
(You don't want 1,000,000,000 boxes!)



### Waffle plots

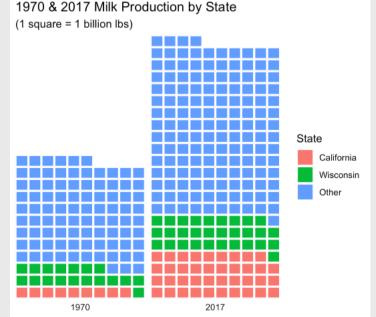
```
library(waffle)
# Format the data
milk production %>%
  filter(year == 2017) %>%
 mutate(state = fct other(state,
    keep = c('California', 'Wisconsin'))) %>%
  group by(state) %>%
  summarise(milk_produced = sum(milk_produced)) %>%
 mutate(milk produced = milk produced / 10^9) %>%
# Make the chart
  ggplot() +
  geom waffle(
    aes(fill = state, values = milk produced),
   color = "white", size = 1, n_cols = 10,
   flip = TRUE) +
  scale x discrete(expand = c(0, 0)) +
  scale y discrete(expand = c(0, 0)) +
  theme minimal() +
  labs(fill = 'State',
       x = NULL, y = NULL,
       title = '2017 Milk Production by State',
       subtitle = '(1 square = 1 billion lbs)')
```

#### Use flip = TRUE for vertical

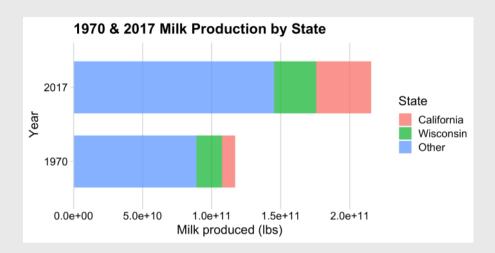


```
library(waffle)
# Format the data
milk production %>%
  filter(year %in% c(1970, 2017)) %>%
 mutate(state = fct other(state,
    keep = c('California', 'Wisconsin'))) %>%
 group by(year, state) %>%
  summarise(milk produced = sum(milk produced)) %>%
 mutate(milk produced = milk produced / 10^9) %>%
# Make the chart
  ggplot() +
  geom waffle(
    aes(fill = state, values = milk produced),
    color = "white", size = 1, n rows = 10,
    flip = TRUE) +
  facet_wrap(vars(year), strip.position = 'bottom')
  scale x discrete(expand = c(0, 0)) +
  scale y discrete(expand = c(0, 0)) +
  theme minimal() +
  labs(fill = 'State',
       x = NULL, y = NULL,
       title = '1970 & 2017 Milk Production by State
       subtitle = '(1 square = 1 billion lbs)')
```

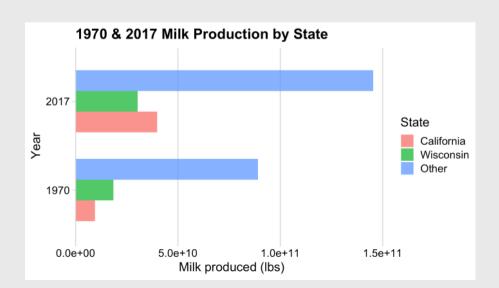
# Use facet\_wrap for side-by-side waffles



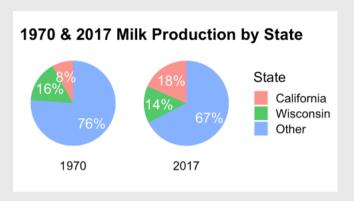
#### Stacked bars



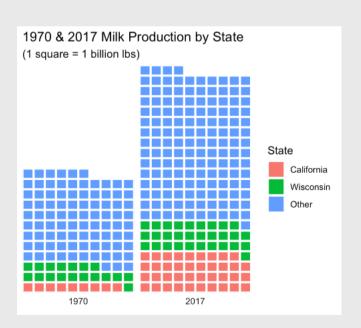
#### Dodged bars



#### Pie chart



#### Waffle chart



### Your turn



Using the wildlife\_impacts data, create plots that shows the proportion of incidents that occur at each different time of day.

For this exercise, you can remove NA values.

Try to create the following plots:

- Stacked bars
- Dodged bars
- Pie chart
- Waffle chart

To get started, you'll need to first summarize the data:

```
wildlife_summary <- wildlife_impacts %>%
  filter(!is.na(time_of_day)) %>%
  count(time_of_day)
wildlife_summary
```