

# Lecture 3: Bar Graphs

February 3, 2020

# Overview

Course Administration

Good, Bad and Ugly

General Graph Design, Few Ch. 9

What is a Bar Chart?

ggplot and Histograms



# Course Administration

1. Collect policy brief proposals
2. Make sure you're checking Piazza
3. Success rate for installing RMarkdown?

- Lauren G.
- Connor D.
- Basia D.

# This Week's Good Bad and Ugly

Finder	Commenter
Lindsay R.	Tereese S.
Kim W.	Danielle C.
Anna W.	David N.

**+** Bryant attempted **30,699** shots throughout his career.

**-**

**Legend:**

- Made
- Missed

April 13, 2016: Capping a 60-point performance, Kobe's final shot from the field was a 19-foot jumper.





Candidates are shown individually if they have at least 10 percent support in an average of national polls or if they've spent more than an estimated \$1 million on TV ads since Jan. 1, 2019.



1. Form
2. Color
3. Spatial Position

1. Maintain visual correspondence to quantity
2. Avoid 3D

# GRAPH CHOICE CHART

## Does your question ask you...

about the **variability** of a group of data points? (i.e. the range of the data, the shape of the distribution, or what the center of the data is)

1. Do all high tides rise to the same height?
2. How variable are wind speeds in Denmark?
3. What is the range and distribution of incomes in Sudan?

to **compare two or more groups** to decide if the groups are the same or different?

if **two numeric factors are correlated**?

1. Is the temperature inside the house correlated with the temperature outside?
2. How did electricity used by the kitchen circuit fluctuate during the past week?

how a **total is proportioned** into sub-groups? (Or what proportion a sub-group is of a total?)

1. What were Brazil's most significant exports in 2015?
2. What proportion of global electricity production comes from wind?
3. How do Parisians typically commute to work?

Do you want to compare the **variability of all data points** in each group to decide if any difference between the groups is meaningful?

1. Which of the two solar cars consistently goes the farthest?
2. Is there a meaningful difference in the heights of fertilized and unfertilized bean plants?

Are you comparing **single numbers** that summarize a group? (such as mean, median, or total...)

1. Was the total snowfall greater this winter than last winter?
2. Do cats and dogs have the same average body temperature?
3. How do the median incomes for the US and India compare?

Does it ask about how something changes through **linear TIME**?

N  
Y

1. Is the fuel efficiency of a car related to its weight?
2. Are smoking rates correlated with median income?
3. Given a fixed volume, how are temperature and pressure related?

1. Is sea level rising?
2. How did my weight change over the last 3 months?

Frequency Plot

MAKE EITHER

FOR EACH GROUP MAKE A

Histogram



Box Plots



Dot Plot



Bar Graph



Scatter Plot



Line Graph



Pie Chart



Stacked Bar Chart







## Bar Charts

Big idea: relative size

# Bars Outline

Big idea: relative size

- What do bars do?
- Few, Ch. 10, bits of bars
- Lollipops, esp from WSJ
- Giant numbers from WSJ



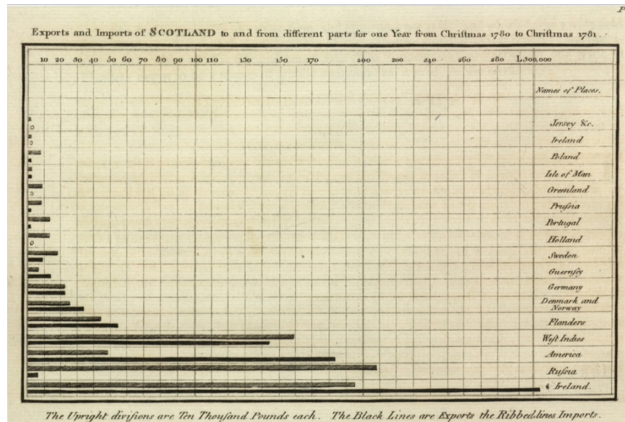
Bars compare quantities across categories

# What Does a Bar Chart Do?

Bars compare quantities across categories

- Levels can be shares
- Emphasize rank order of levels
- Highlight one level relative to others

# The First Bar Chart



Playfair, William, 1786. *The Commercial and Political Atlas*

## General Principles for Bar Charts

- Orientation
- Proximity
- Fills
- Borders
- Base value

Taken from Few Ch. 10, p. 210

Orientation: Bars horizontal or vertical?

- Horizontal better to fit in long lables
- Vertical better if axis is time

- You want mostly bars, not mostly white space
- But not touching bars
- Why not touching bars?

# Orientation & Proximity

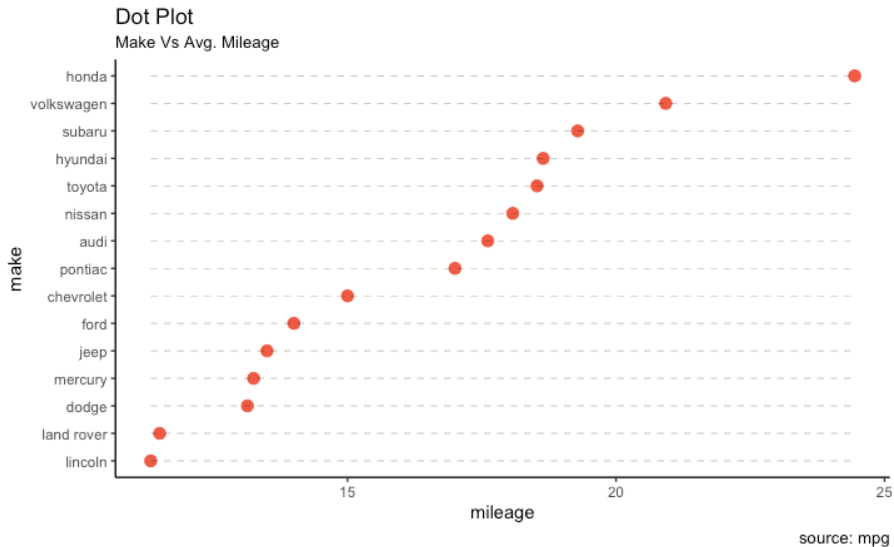
Orientation: Bars horizontal or vertical?

- Horizontal better to fit in long labels
- Vertical better if axis is time

Proximity – How close should the bars be?

- You want mostly bars, not mostly white space
- But not touching bars
- Why not touching bars?
- **Rank when you want to highlight ordering**

# Ranked Almost-Bars





# Fills

## Fills

### Do Not

- Use color as decoration
- Use hashed or lined fills

### Do

- As much as possible, put legend directly on the graph
- Highlight with color



## Borders

- Use sparingly to highlight
- Colors are better for highlighting

## Borders

- Use sparingly to highlight
- Colors are better for highlighting

Base Value

## Bars Must Start at Zero!

## Borders

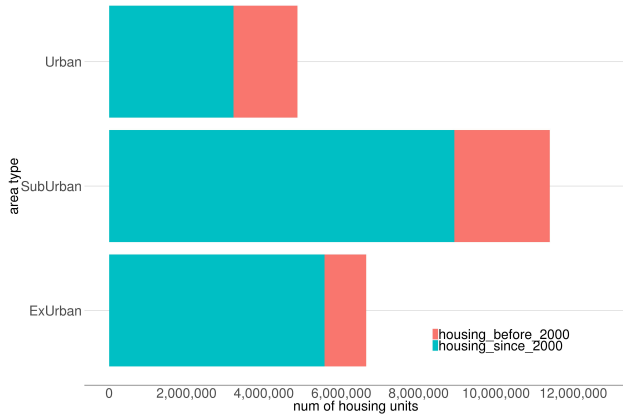
- Use sparingly to highlight
- Colors are better for highlighting

Base Value

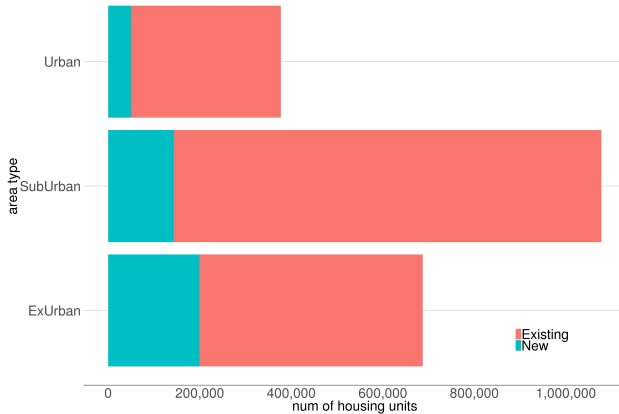
## Bars Must Start at Zero!

## Why?

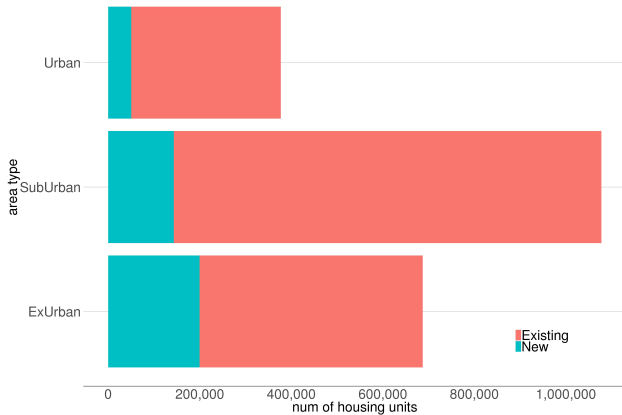
## Going From a Bad Bar Chart to a Decent One



## Going From a Bad Bar Chart to a Decent One

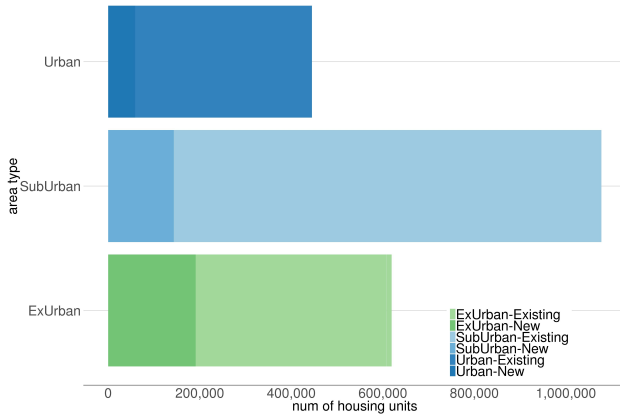


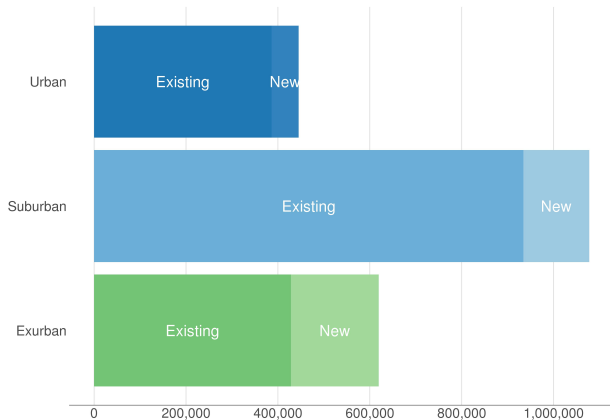
## Going From a Bad Bar Chart to a Decent One





## Going From a Bad Bar Chart to a Decent One



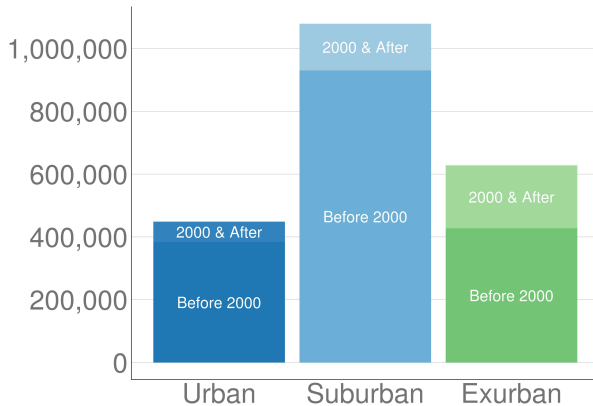


A stacked bar chart illustrating the number of existing and new jobs by location type. The Y-axis represents the number of jobs, ranging from 0 to 1,000,000 in increments of 200,000. The X-axis shows three location types: Urban, Suburban, and Exurban. Each bar is divided into two segments: 'Existing' (bottom) and 'New' (top). The 'New' segment is a lighter shade of the 'Existing' segment's color.

Location Type	Existing Jobs	New Jobs	Total Jobs
Urban	~380,000	~60,000	~440,000
Suburban	~940,000	~140,000	~1,080,000
Exurban	~430,000	~190,000	~620,000

Area	Existing Jobs	New Jobs	Total Jobs
Urban	380,000	60,000	440,000
Suburban	950,000	100,000	1,050,000
Exurban	430,000	190,000	620,000

## Going From a Bad Bar Chart to a Decent One



## When the Number is Too Big for a Bar

People really have trouble with big numbers

- is \$2 billion large part of a \$4 trillion budget?

## When the Number is Too Big for a Bar

People really have trouble with big numbers

- is \$2 billion large part of a \$4 trillion budget?

Microsoft plan for conveying big numbers, from *WSJ*

- attribute
- scaling factor
- reference

## When the Number is Too Big for a Bar

People really have trouble with big numbers

- is \$2 billion large part of a \$4 trillion budget?

Microsoft plan for conveying big numbers, from *WSJ*

- attribute
- scaling factor
- reference

“a conservation group that reclaimed about 100 million acres of land across the Earth. ... How big do you think that is?”



## When the Number is Too Big for a Bar

People really have trouble with big numbers

- is \$2 billion large part of a \$4 trillion budget?

Microsoft plan for conveying big numbers, from *WSJ*

- attribute
- scaling factor
- reference

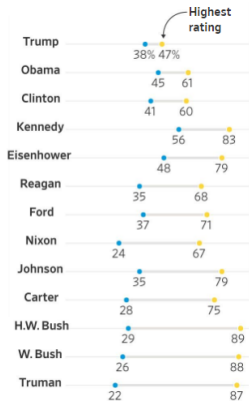
“a conservation group that reclaimed about 100 million acres of land across the Earth. ... How big do you think that is?”

About as big as (1.15x = scaling factor)  
California (reference)

# The Power of Lollipops

## Americans Locked In

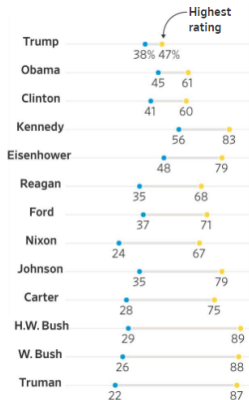
President Trump's job approval rating has held unusually steady throughout his term, moving up and down in the tightest range on record.



# The Power of Lollipops

## Americans Locked In

President Trump's job approval rating has held unusually steady throughout his term, moving up and down in the tightest range on record.

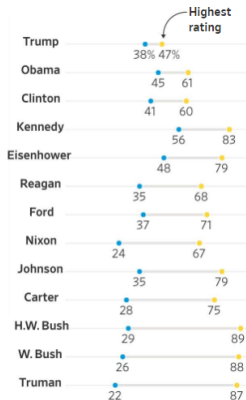


What is the point?

# The Power of Lollipops

## Americans Locked In

President Trump's job approval rating has held unusually steady throughout his term, moving up and down in the tightest range on record.

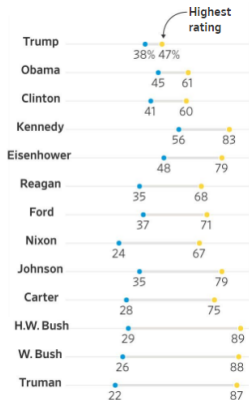


What is the point?  
What info does this convey?

# The Power of Lollipops

## Americans Locked In

President Trump's job approval rating has held unusually steady throughout his term, moving up and down in the tightest range on record.



What is the point?

What info does this convey?

- max and min
- approximately the variance
- by administration
- so a trend in variance!
- note the point in the title

- Admin
  -

G/B/U  
OOOOO

Graph Design  
○○○○

Bar Chart  
○○○○○○○○○○○○○○○○

R

●○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

R

# Today

- A. What is ggplot?
- B. The parts of ggplot
- C. Bars via ggplot
- D. Titles and axis scaling
- E. Factor re-ordering
- F. Summary statistics

## A. Graphing in R

- ▶ `ggplot` is the premier package for graphing in R
- ▶ There is a simple version of `ggplot` called `qplot`: we ignore it
- ▶ Developed in 2005 by Hadley Wickham
- ▶ In 2017, Wickham says “Ten years after `ggplot2`’s release, Wickham wonders how much longer his program will dominate chart making in R. ‘It really feels to me now like `ggplot2` is ripe for disruption,’ said Wickham. ‘I’m surprised some young gun hasn’t come along, and thought, ‘Wow this is crap,’ and done better. But so far, it hasn’t really happened.’”

Article link is [here](#).



## B. The Key Parts of a ggplot command

```
ggplot() +  
  geom_something(data = ,  
                 mapping = aes(x = xvar, [y = yvar]))
```

## B. The Key Parts of a ggplot command

```
ggplot() +  
  geom_something(data = ,  
                 mapping = aes(x = xvar, [y = yvar]))
```

This will pop up a graph in the plots window.

## Making your graph an object

```
leahs.graph <- ggplot() +  
  geom_something(data = ,  
                 mapping = aes(x = xvar, [y = yvar]))
```

Will show nothing, but creates `leahs.graph` to which you can refer.

## Making your graph an object

```
leahs.graph <- ggplot() +  
  geom_something(data = ,  
                 mapping = aes(x = xvar, [y = yvar]))
```

Will show nothing, but creates `leahs.graph` to which you can refer.

And

```
leahs.graph
```

Will pop up a graph.

You can add to an object

```
leahs.graph2 <- leahs.graph +  
  geom_another(data = ,  
               mapping = aes(x = xvar, [y = yvar]))
```

I usually name the object and call the named object.

## C.1. Bars

At its most basic

```
new.graph <- ggplot() +  
  geom_col(data = [your data],  
           mapping = aes(x = [categorical variable],  
                         y = [value]))
```

## C.2. Bar Chart Additions

Of course, there are many more things you can do

- ▶ make R create statistics for the value in the y variable: `geom_bar()`
- ▶ make stacked bars: `position = "stack"`
- ▶ make grouped bars: `position = "dodge"`
- ▶ change the bar width
- ▶ change bar colors
- ▶ put labels on bars
- ▶ and still oodles more

## D. Making Graphs Legible

```
new.graph <- ggplot() +  
  geom_col(data = [your data],  
           mapping = aes(x = [categorical variable],  
                         y = [value])) +  
  labs(title = "title here",  
        x = "x label",  
        y = "y label") +  
  [things about scales] +  
  theme([things you modify here])
```



## D. Making Graphs Legible

```
new.graph <- ggplot() +  
  geom_col(data = [your data],  
           mapping = aes(x = [categorical variable],  
                         y = [value])) +  
  labs(title = "title here",  
        x = "x label",  
        y = "y label") +  
  [things about scales] +  
  theme([things you modify here])
```

+ 1000s of more options

## E. Factor variables

- ▶ recall that R has a type of variable called a factor
- ▶ often created when a variable has a limited number of values
- ▶ useful to save memory space
- ▶ useful for making charts

## E.1. Factor levels

- ▶ we particularly care about factor levels this class
- ▶ R orders bar charts by the order of the factor
- ▶ to change the order, change the order of the factor

## E.2. Setting up a factor variable

```
states <- data.frame(state_abbrev = c("VA", "DC", "MD"),  
                      state_fips = c(51, 11, 24),  
                      av.feb.temp = c(50, 47, 45))  
  
str(states)
```

```
## 'data.frame':    3 obs. of  3 variables:  
## $ state_abbrev: Factor w/ 3 levels "DC","MD","VA": 3 1 2  
## $ state_fips  : num  51 11 24  
## $ av.feb.temp : num  50 47 45
```

- ▶ state is a factor variable
- ▶ has three levels: DC, MD, VA
- ▶ in that order – R auto-alphabetizes
- ▶ suppose we prefer it in another order: VA, DC, MD

### E.3. Re-ordering a factor

Change from [DC, MD, VA] to [VA, DC, MD]

```
levels(states$state_abbrev)
```

```
## [1] "DC" "MD" "VA"
```

```
states$state_abrev2 <- factor(states$state_abbrev,  
                             levels = c("VA", "DC", "MD"))
```

```
levels(states$state_abrev2)
```

```
## [1] "VA" "DC" "MD"
```

### E.3. Re-ordering a factor

Change from [DC, MD, VA] to [VA, DC, MD]

```
levels(states$state_abbrev)
```

```
## [1] "DC" "MD" "VA"
```

```
states$state_abrev2 <- factor(states$state_abbrev,  
                             levels = c("VA", "DC", "MD"))
```

```
levels(states$state_abrev2)
```

```
## [1] "VA" "DC" "MD"
```

Remember: you need this to re-order bars.

## F. Summary statistics are useful

- ▶ to check data
- ▶ to display data

## F.1 Call dplyr package

```
library(dplyr)
```

```
## Warning: package 'dplyr' was built under R version 3.6.2
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
##      filter, lag
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##      intersect, setdiff, setequal, union
```

Part of the tidyverse. If not installed, you'll need to do so.



## F.2. mutate()

- ▶ if you know stata's egen, it's like that
- ▶ create a new variable that has the average temperature for all three states

```
library(dplyr)
```

```
states <- mutate(.data = states,  
                  all.states.feb=mean(av.feb.temp,  
                                       na.rm = TRUE))
```

```
states
```

```
##   state_abbrev state_fips av.feb.temp state_abrev2 all.states.feb  
## 1          VA         51          50          VA      47.33333  
## 2          DC          11          47          DC      47.33333  
## 3          MD          24          45          MD      47.33333
```

## F.2. mutate()

- ▶ if you know stata's egen, it's like that
- ▶ create a new variable that has the average temperature for all three states

```
library(dplyr)
```

```
states <- mutate(.data = states,  
                 all.states.feb=mean(av.feb.temp,  
                                     na.rm = TRUE))
```

```
states
```

```
##   state_abbrev state_fips av.feb.temp state_abrev2 all.states.feb  
## 1          VA         51          50          VA      47.33333  
## 2          DC          11          47          DC      47.33333  
## 3          MD          24          45          MD      47.33333
```

Why not just `av.temp <- mean(states$av.feb.temp, na.rm = TRUE)?`

### F.3. More on `mutate()`

- ▶ it does many many other things as well
- ▶ you can use all kinds of functions in the second term
- ▶ and create more than one new variable
- ▶ can combine with `group_by()`

## Next Class

- Turn in PS 3
- Few, Chapter 6
- Chang, Chapter 6 (through 6.5)
- Linked Bloomberg article on quantities of land