

Lecture 6: Visualizing Numerical and Categorical Data

Chapter 1.6+1.7

1 / 25

Goals for Today

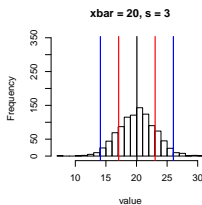
- ▶ Rule of thumb for standard deviations
- ▶ Population vs sample mean/variance/standard deviations
- ▶ Percentiles and Quartiles
- ▶ Boxplots
- ▶ Piecharts, barplots, mosaicplots

2 / 25

Rule of Thumb for Standard Deviations

3 / 25

Example

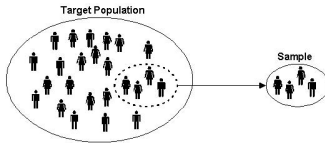


- ▶ black line is mean \bar{x}
- ▶ red lines mark about $\frac{2}{3}$:
 $[\bar{x} - s, \bar{x} + s] = [17, 23]$.
- ▶ blue lines mark about 95%:
 $[\bar{x} - 2s, \bar{x} + 2s] = [14, 26]$.

4 / 25

Population vs Sample Mean/Variance/Standard Deviation

Recall the notion of taking a **representative sample** from a **study/target population**. Say we are interested in the income of the individuals.



5 / 25

Population vs Sample Mean/Variance/Standard Deviation

6 / 25

Population vs Sample Mean/Variance/Standard Deviation

7 / 25

Percentiles

8 / 25

Quartiles and IQR

Quartiles and IQR

Example: `http:`

`//www.npr.org/sections/money/2015/03/19/394057221/
how-much-or-little-the-middle-class-makes-in-30-u-s-cities`

Robust Statistics (Chapter 1.6.6)

Boxplots

Boxplots are an alternative visualization of a sample x_1, \dots, x_n , but draw attention to outliers.

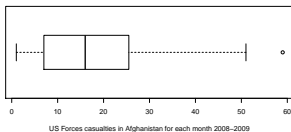
Example: # US Forces casualties in the war in Afghanistan for each month from 2008-2009:

7, 1, 7, 5, 16, 28, 20, 22, 27, 16, 1, 3, 14, 15, 13, 6, 12, 24, 44, 51, 37, 59, 17, 17

Boxplots

The summary values:

| Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. |
|------|---------|--------|-------|---------|-------|
| 1.00 | 7.00 | 16.00 | 19.25 | 24.75 | 59.00 |



Length of **whiskers**: they capture data that is no more than $1.5 \times IQR$ of both ends of the box.

13 / 25

Outliers Are Relatively Extreme

An **outlier** is an observation that appears extreme relative to the rest of the data.

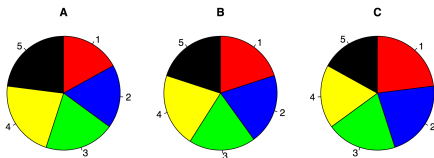
Why it is important to look for outliers? Examination of data for possible outliers serves many useful purposes, including

- ▶ Identifying strong skew in the distribution.
- ▶ Identifying data collection or entry errors.
- ▶ Providing insight into interesting properties of the data.

14 / 25

Piecharts

Say we have the following piecharts represent the polling from a local election with five candidates (1-5) at three different time points A, B, and C:

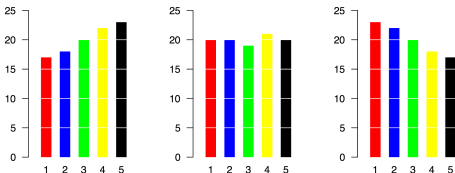


Answer the following questions:

- ▶ In the first race, is candidate 5 doing better than candidate 4?
- ▶ Who did better between time A and time B, candidate 2 or candidate 4?

15 / 25

Barplots Instead

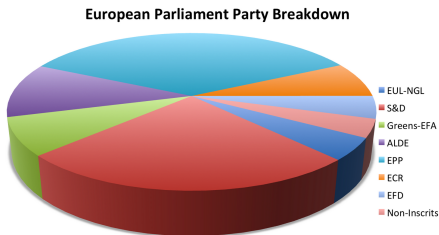


Answers:

- ▶ Candidate 5 is doing better than 4
- ▶ Between A and B, candidate 2 went from about 17% to 20% while candidate 4 went from about 22% to 21%. So candidate 2 did better

16 / 25

3D Piecharts Can Be Deceiving



EPP (teal) has 266 seats, whereas S&D (red) has 190 seats.

17 / 25

Titanic Survival Data

Typing `data(Titanic)` in R loads the survival and death counts, split by each of the following categories:

- ▶ Class: 1st, 2nd, 3rd, or crew (4 levels)
- ▶ Gender (2 levels)
- ▶ Age: Child or adult (2 levels)

i.e. $4 \times 2 \times 2 = 16$ possible groups to consider.

Questions

- ▶ What was the effect of class (1st, 2nd, 3rd, crew) on your chances of survival?
- ▶ Did the “women and children” first lifeboat policy hold?

18 / 25

Frequency Table

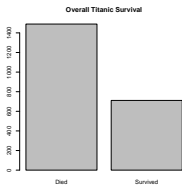
A table summarizing a single categorical variable is called a **frequency table**. Overall:

| | |
|----------|------|
| Died | 1490 |
| Survived | 711 |
| Total | 2201 |

19 / 25

Barplot

Barplots are ways to display categorical variables:



20 / 25

Contingency Table

A table that **cross-classifies** two categorical variables is a **contingency table**. Now let's split survival by class: 1st, 2nd, 3rd, and crew.

Before:

| | |
|----------|------|
| Died | 1490 |
| Survived | 711 |
| Total | 2201 |

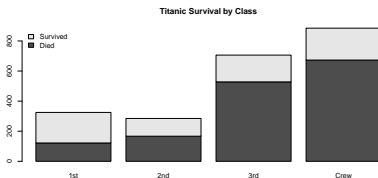
After:

| | 1st | 2nd | 3rd | Crew | Total |
|----------|-----|-----|-----|------|-------|
| Died | 122 | 167 | 528 | 673 | 1490 |
| Survived | 203 | 118 | 178 | 212 | 711 |
| Total | 325 | 285 | 706 | 885 | 2201 |

21 / 25

Stacked Barplot

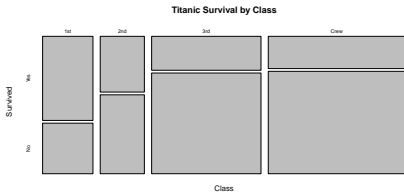
Stacked barplots are one way to display values from a contingency table:



22 / 25

Mosaic Plots

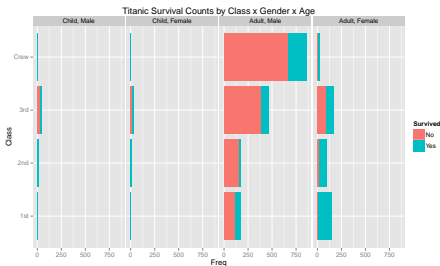
Mosaic plots are similar, but the widths of the bars now reflect proportions:



23 / 25

Stacked Barplots

Using the `ggplot2` package, we can plot survivals by class, age, and gender all at once.



24 / 25

Standardized/Normalized Stacked Barplots

Instead of raw counts, we can expand each bar to reflect proportions (i.e. standardize/normalize them).

