# Lecture 6: Visualizing Numerical and Categorical Data
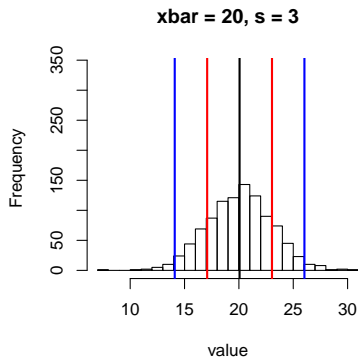
Chapter 1.6+1.7

# Goals for Today

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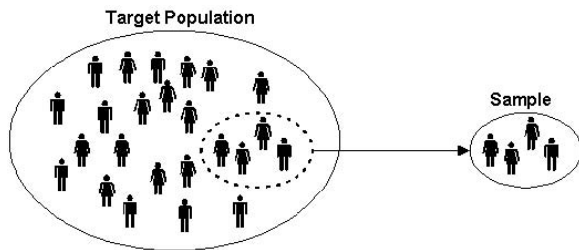- Piecharts, barplots, mosaicplots

# Rule of Thumb for Standard Deviations

# Example



xbar = 20, s = 3

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

# 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.

# Population vs Sample Mean/Variance/Standard Deviation

# Population vs Sample Mean/Variance/Standard Deviation

# Percentiles

# 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 $x_1, \ldots, x_n$, but draw attention to outliers.

# Boxplots

Boxplots are an alternative visualization of a sample $x_1, \ldots, x_n$, 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
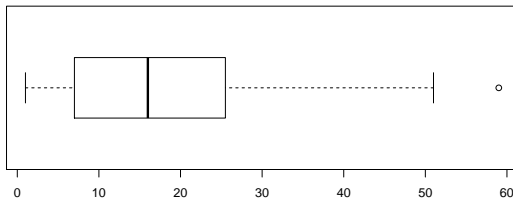
# Boxplots

The summary values:

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

# Boxplots

The summary values:

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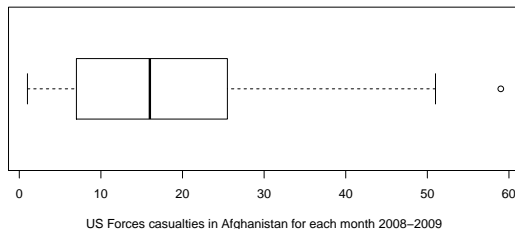


US Forces casualties in Afghanistan for each month 2008−2009

# Boxplots

The summary values:

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



US Forces casualties in Afghanistan for each month 2008–2009

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

# Outliers Are Relatively Extreme

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

# 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

# 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.

# 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
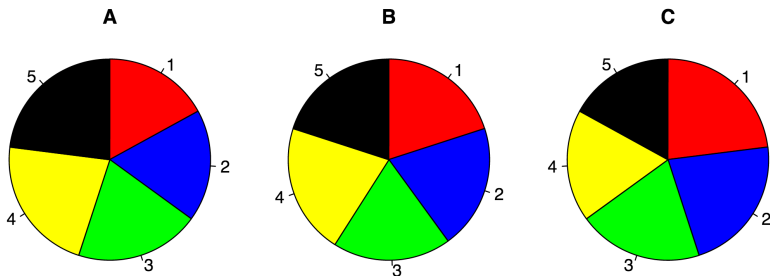
- Identifying strong skew in the distribution.
- Identifying data collection or entry errors.

# 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.
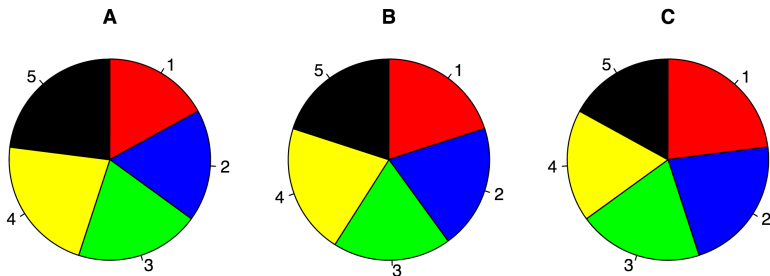- Providing insight into interesting properties of the data.

## 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, an C:

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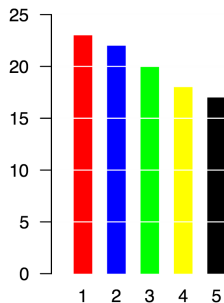
Say we have the following piecharts represent the polling from a local election with five candidates (1-5) at three different time points A, B, an C:
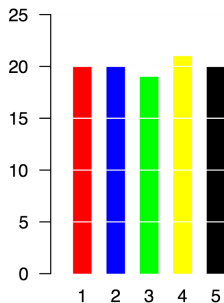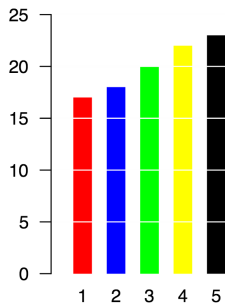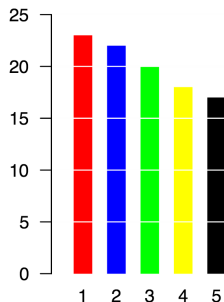
# Piecharts

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, an C:
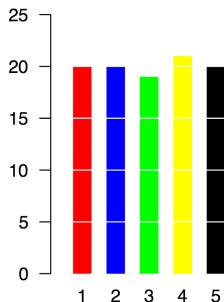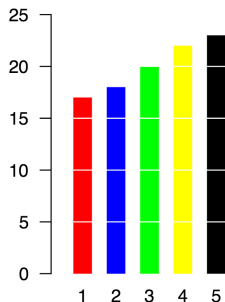


Answer the following questions:

- 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?

# Barplots Instead

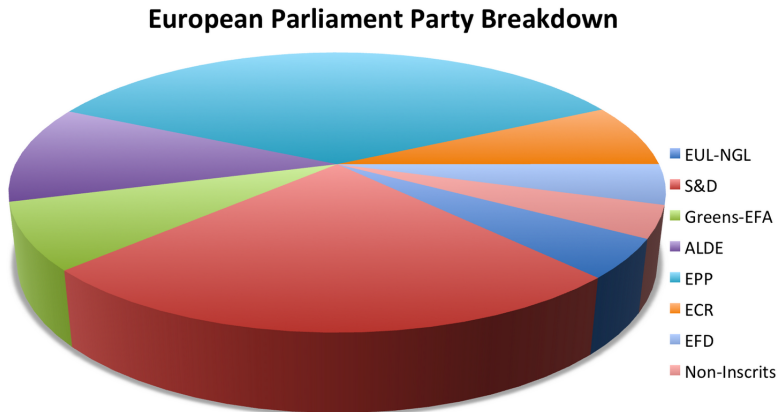# Barplots Instead



Answers:

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

# 3D Piecharts Can Be Deceiving



**European Parliament Party Breakdown**

- EUL-NGL
- S&D
- Greens-EFA
- ALDE
- EPP
- ECR
- EFD
- Non-Inscrits

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

# Titanic Survival Data

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

# 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)

# 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.

# 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?

# 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?
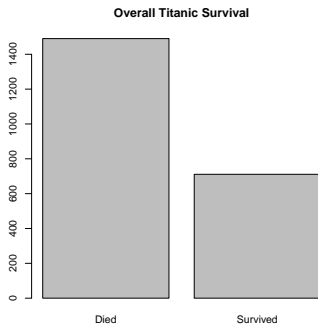- ▶ Did the "women and children" first lifeboat policy hold?

# Frequency Table

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

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

# Barplot

Barplots are ways to display categorial variables:



**Overall Titanic Survival**

# 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.

# 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.
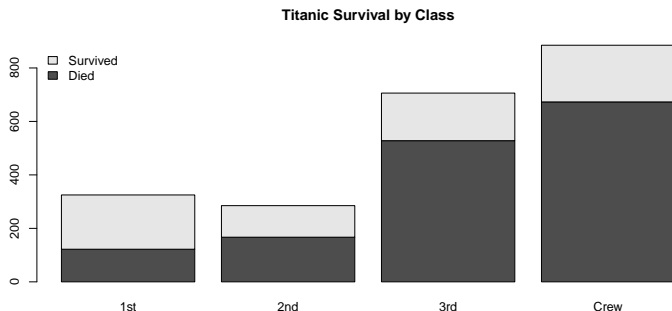
Before:

|          | Died |
|----------|------|
| 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  |

# Stacked Barplot

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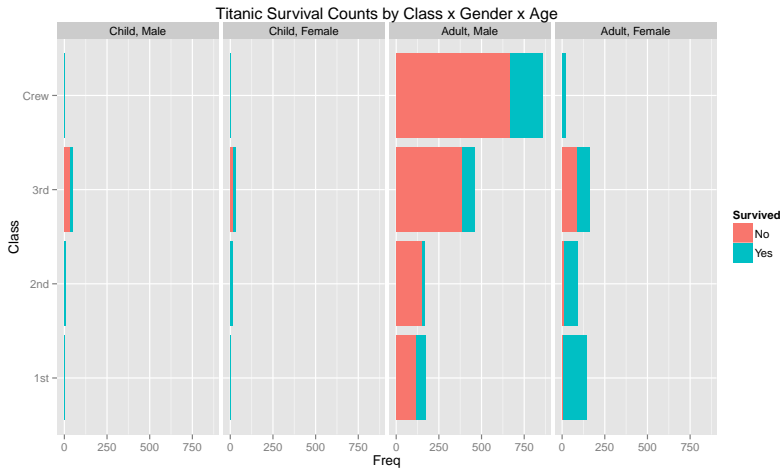


**Titanic Survival by Class**

# Mosaic Plots

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



**Titanic Survival by Class**

# Stacked Barplots

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

# Standardized/Normalized Stacked Barplots

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



Titanic Survival Proportions by Class x Gender x Age