

# Class 8: Statistical distributions

---

May 31, 2018



# General

# Announcements

- Homework 2 to be posted: <http://summer18.cds101.com/assignments/homework-2/>
- Reading 8 from **R for Data Science**, questions due on June 1st by 9:00am
  - From **chapter 12**: section 12.1 through to the end of section 12.5

# Practicing with real data: Chicago Towed Vehicles dataset

# Chicago Towing Data

- Data posted by the city of Chicago of Towed Vehicle information over the past 90 days: <https://data.cityofchicago.org/Transportation/Towed-Vehicles/ygr5-vcbg/>
- Use your RStudio project for the dplyr demos for this activity: <https://classroom.github.com/a/eF7HdfVO>
- Practice using dplyr by investigating this dataset!

# Towing Data Questions

Before we begin, let's clean up the dataset up a little so that it's easier to work with

Follow along in RStudio

# Towing Data Questions

What command would I run to group the towed vehicles by color and see a summary of how many cars of each color were impounded over the last 90 days?

What command would I run to see which day, over the last 90 days, saw the most cars impounded and the least cars impounded?

Come up with your own question to explore the dataset!

# Statistical distributions



# Variance

*Variance* is roughly the average squared deviation from the mean.

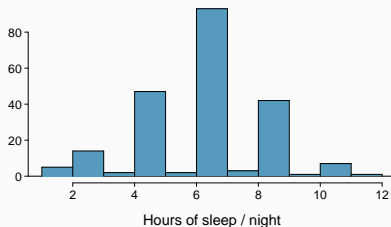
$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

# Variance

*Variance* is roughly the average squared deviation from the mean.

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

- The sample mean is  $\bar{x} = 6.71$ , and the sample size is  $n = 217$ .

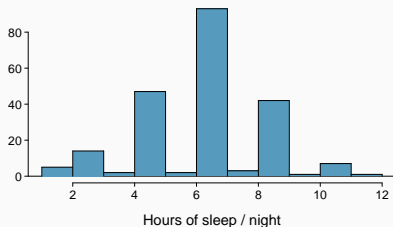


# Variance

*Variance* is roughly the average squared deviation from the mean.

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

- The sample mean is  $\bar{x} = 6.71$ , and the sample size is  $n = 217$ .
- The variance of amount of sleep students get per night can be calculated as:



$$s^2 = \frac{(5 - 6.71)^2 + (9 - 6.71)^2 + \dots + (7 - 6.71)^2}{217 - 1} = 4.11 \text{ hours}^2$$

## Variance (cont.)

Why do we use the squared deviation in the calculation of variance?

## Variance (cont.)

Why do we use the squared deviation in the calculation of variance?

- *To get rid of negatives so that observations equally distant from the mean are weighed equally.*
- *To weigh larger deviations more heavily.*

## Standard deviation

The *standard deviation* is the square root of the variance, and has the same units as the data.

$$s = \sqrt{s^2}$$

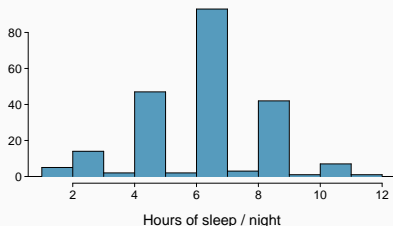
# Standard deviation

The *standard deviation* is the square root of the variance, and has the same units as the data.

$$s = \sqrt{s^2}$$

- The standard deviation of amount of sleep students get per night can be calculated as:

$$s = \sqrt{4.11} = 2.03 \text{ hours}$$



## Standard deviation

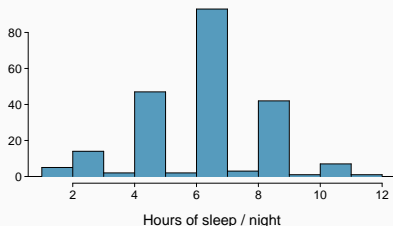
The *standard deviation* is the square root of the variance, and has the same units as the data.

$$s = \sqrt{s^2}$$

- The standard deviation of amount of sleep students get per night can be calculated as:

$$s = \sqrt{4.11} = 2.03 \text{ hours}$$

- We can see that all of the data are within 3 standard deviations of the mean.





# Median

- The *median* is the value that splits the data in half when ordered in ascending order.

0, 1, **2**, 3, 4

- If there are an even number of observations, then the median is the average of the two values in the middle.

$$0, 1, \underline{2}, 3, 4, 5 \rightarrow \frac{2 + 3}{2} = 2.5$$

- Since the median is the midpoint of the data, 50% of the values are below it. Hence, it is also the *50<sup>th</sup> percentile*.

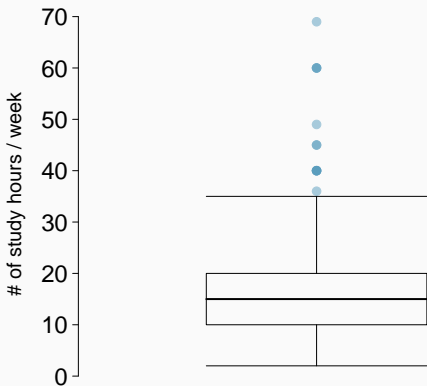
## Q1, Q3, and IQR

- The 25<sup>th</sup> percentile is also called the first quartile, **Q1**.
- The 50<sup>th</sup> percentile is also called the median.
- The 75<sup>th</sup> percentile is also called the third quartile, **Q3**.
- Between Q1 and Q3 is the middle 50% of the data. The range these data span is called the *interquartile range*, or the **IQR**.

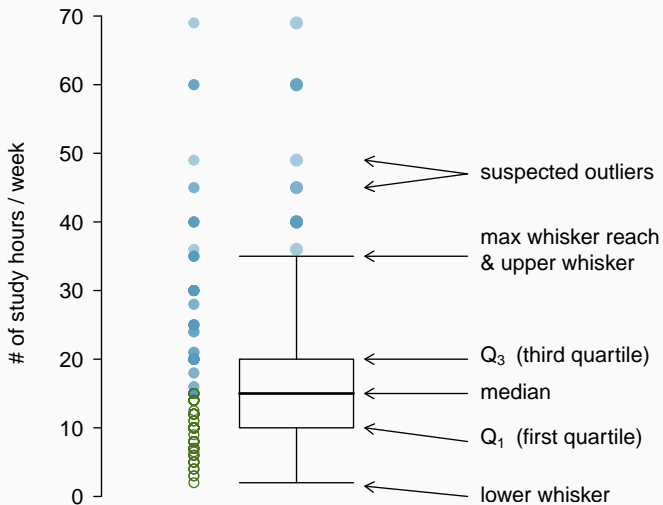
$$IQR = Q3 - Q1$$

## Box plot

The box in a *box plot* represents the middle 50% of the data, and the thick line in the box is the median.



# Anatomy of a box plot



## Whiskers and outliers

- *Whiskers*

of a box plot can extend up to  $1.5 \times IQR$  away from the quartiles.

$$\text{max upper whisker reach} = Q3 + 1.5 \times IQR$$

$$\text{max lower whisker reach} = Q1 - 1.5 \times IQR$$

## Whiskers and outliers

- *Whiskers*

of a box plot can extend up to  $1.5 \times IQR$  away from the quartiles.

$$\text{max upper whisker reach} = Q3 + 1.5 \times IQR$$

$$\text{max lower whisker reach} = Q1 - 1.5 \times IQR$$

$$IQR : 20 - 10 = 10$$

$$\text{max upper whisker reach} = 20 + 1.5 \times 10 = 35$$

$$\text{max lower whisker reach} = 10 - 1.5 \times 10 = -5$$

## Whiskers and outliers

- *Whiskers*

of a box plot can extend up to  $1.5 \times IQR$  away from the quartiles.

$$\text{max upper whisker reach} = Q3 + 1.5 \times IQR$$

$$\text{max lower whisker reach} = Q1 - 1.5 \times IQR$$

$$IQR : 20 - 10 = 10$$

$$\text{max upper whisker reach} = 20 + 1.5 \times 10 = 35$$

$$\text{max lower whisker reach} = 10 - 1.5 \times 10 = -5$$

- A potential *outlier* is defined as an observation beyond the maximum reach of the whiskers. It is an observation that appears extreme relative to the rest of the data.

## Outliers (cont.)

Why is it important to look for outliers?



## Outliers (cont.)

Why is it important to look for outliers?

- *Identify extreme skew in the distribution.*
- *Identify data collection and entry errors.*
- *Provide insight into interesting features of the data.*

## Robust statistics

Median and IQR are more robust to skewness and outliers than mean and SD. Therefore,

- for skewed distributions it is often more helpful to use median and IQR to describe the center and spread
- for symmetric distributions it is often more helpful to use the mean and SD to describe the center and spread

## Robust statistics

Median and IQR are more robust to skewness and outliers than mean and SD. Therefore,

- for skewed distributions it is often more helpful to use median and IQR to describe the center and spread
- for symmetric distributions it is often more helpful to use the mean and SD to describe the center and spread

If you would like to estimate the typical household income for a student, would you be more interested in the mean or median income?

## Robust statistics

Median and IQR are more robust to skewness and outliers than mean and SD. Therefore,

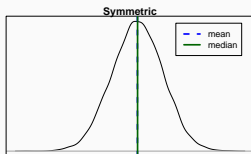
- for skewed distributions it is often more helpful to use median and IQR to describe the center and spread
- for symmetric distributions it is often more helpful to use the mean and SD to describe the center and spread

If you would like to estimate the typical household income for a student, would you be more interested in the mean or median income?

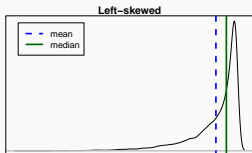
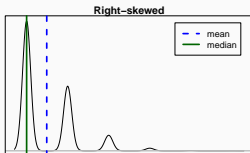
*Median*

# Mean vs. median

- If the distribution is symmetric, center is often defined as the mean:  $\text{mean} \approx \text{median}$



- If the distribution is skewed or has extreme outliers, center is often defined as the median
  - Right-skewed:  $\text{mean} > \text{median}$
  - Left-skewed:  $\text{mean} < \text{median}$



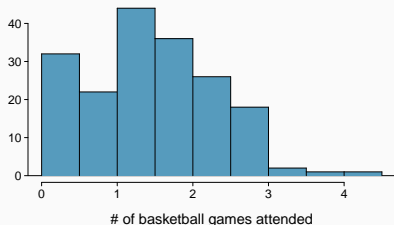
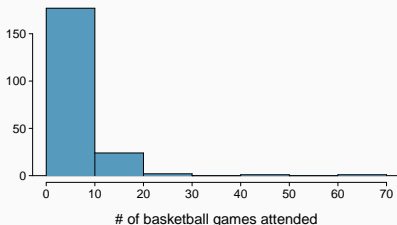
## Extremely skewed data

When data are extremely skewed, transforming them might make modeling easier. A common transformation is the *log transformation*.

## Extremely skewed data

When data are extremely skewed, transforming them might make modeling easier. A common transformation is the *log transformation*.

The histograms on the left shows the distribution of number of basketball games attended by students. The histogram on the right shows the distribution of log of number of games attended.



# Credits

The slides with blue headers originate from the following source:

- The Chapter 1 [OpenIntro Statistics slides](#) developed by Mine Çetinkaya-Rundel and made available under the [CC BY-SA 3.0 license](#).