

Unit 1: Introduction to data
2. Exploratory data analysis

Sta 101 – Spring 2019

Duke University, Department of Statistical Science

How old were you when you had your first kiss?

Dr. Ellison

Slides posted at <https://www2.stat.duke.edu/courses/Spring19/sta101.001/index.html>

Register your Clicker!

To do now:

- Turn on your clicker (**orange button**)
- Wait about 6 seconds.
- IF your clicker says "READY,":**
 - Look for when your name appears on the slides. (*If you don't see your name and you are officially enrolled in the course let me know!*)
 - When you see your name, type the 4 letters you see under it (you have 15 seconds).
 - If your name box turned another color, your clicker should now be registered to the class!**
- IF your clicker does not say "READY":**
 - Hold down on the **orange button** until the clicker screen changes.
 - Quickly press AA.
 - Wait ~6 seconds.
 - Your screen should now say "READY". (*If not, ask a TA for help!*)

We will be registering clickers 1/14, 1/16, 1/23... clicker grading begins on 1/28!

Readiness assessment

► **Individual:** 15 minutes, using clickers

► **Team:** 10 minutes, using scratch off sheets (1 per team)









1









Summary of main ideas

To Do:





- Getting to Know you Survey + Pretest **due tomorrow 1/15**
- Start working on Problem Set 1 (**Due Friday 1/26**)

16

Outline	
I.	Readiness assessment
II.	Housekeeping
III.	Main ideas
1.	<u>Analysis work flow:</u>
1.	  Always start your exploration with a visualization
2.	<u>Single Numerical Variable</u>
1.	   When describing numerical distributions discuss shape, center, spread, and unusual observations
2.	 Robust statistics are not easily affected by outliers and extreme skew
3.	  Use box plots to display quartiles, median, and outliers
IV.	Application exercises
V.	Summary

Outline	
I.	Readiness assessment
II.	Housekeeping
III.	Main ideas
1.	<u>Analysis work flow:</u>
1.	  Always start your exploration with a visualization
2.	<u>Single Numerical Variable</u>
1.	   When describing numerical distributions discuss shape, center, spread, and unusual observations
2.	 Robust statistics are not easily affected by outliers and extreme skew
3.	  Use box plots to display quartiles, median, and outliers
IV.	Application exercises
V.	Summary

Outline	
I.	Readiness assessment
II.	Housekeeping
III.	Main ideas
1.	<u>Analysis work flow:</u>
1.	  Always start your exploration with a visualization
2.	<u>Single Numerical Variable</u>
1.	   When describing numerical distributions discuss shape, center, spread, and unusual observations
2.	 Robust statistics are not easily affected by outliers and extreme skew
3.	  Use box plots to display quartiles, median, and outliers
IV.	Application exercises
V.	Summary

Outline	
I.	Readiness assessment
II.	Housekeeping
III.	Main ideas
1.	<u>Analysis work flow:</u>
1.	  Always start your exploration with a visualization
2.	<u>Single Numerical Variable</u>
1.	   When describing numerical distributions discuss shape, center, spread, and unusual observations
2.	 Robust statistics are not easily affected by outliers and extreme skew
3.	  Use box plots to display quartiles, median, and outliers
IV.	Application exercises
V.	Summary

Outline

Should we calculate a summary statistic or make a data visualization first?

```

    graph LR
      A[Form Research Question  
(about Population)] --> B[Collect Data  
(Sample)]
      B --> C[Exploratory Data Analysis  
(of Sample)  
*Summary Statistics  
*Visualizations]
      C --> D[Make Inferences and Conclusions  
(about Population)]
      C --> A
      C --> B
  
```

From a past Sta 101 survey...

Do you see anything out of the ordinary?

Age Bin	Frequency
0-1	1
1-2	0
2-3	0
3-4	2
4-5	2
5-6	2
6-7	2
7-8	2
8-9	2
9-10	6
10-11	15
11-12	17
12-13	19
13-14	20
14-15	18
15-16	9
16-17	0
17-18	0
18-19	0
19-20	0

age at first kiss

3

From a past Sta 101 survey...

Do you see anything out of the ordinary?

Age Bin	Frequency
0-1	1
1-2	0
2-3	0
3-4	2
4-5	2
5-6	2
6-7	2
7-8	2
8-9	2
9-10	6
10-11	15
11-12	17
12-13	19
13-14	20
14-15	18
15-16	9
16-17	0
17-18	0
18-19	0
19-20	0

age at first kiss

Some people reported very low ages, which might suggest the survey question wasn't clear: romantic kiss or any kiss?

3

Outline

We should start our exploratory data analysis with a visualization first!

```

    graph LR
      A[Form Research Question  
(about Population)] --> B[Collect Data  
(Sample)]
      B --> C[Exploratory Data Analysis  
(of Sample)  
*Summary Statistics  
*Visualizations]
      C --> D[Make Inferences and Conclusions  
(about Population)]
      C --> A
      C --> B
  
```

Outline

Lab Hint: When asked to describe a visualization of a single numerical variable, there are four things we should always discuss. What are they?



Outline

- I. Readiness assessment
- II. Housekeeping
- III. Main ideas
 1. Analysis work flow:
 1. NEW Always start your exploration with a visualization
 2. Single Numerical Variable
 1. Q Q Q When describing numerical distributions discuss shape, center, spread, and unusual observations
 2. Q Robust statistics are not easily affected by outliers and extreme skew
 3. Q NEW Use box plots to display quartiles, median, and outliers
- IV. Application exercises
- V. Summary



Describing distributions of numerical variables

- ▶ **Shape:** skewness, modality
- ▶ **Center:** an estimate of a *typical* observation in the distribution (mean, median, mode, etc.)
 - Notation: μ : population mean, \bar{x} sample mean
- ▶ **Spread:** measure of variability in the distribution (standard deviation, IQR, range, etc.)
- ▶ **Unusual observations:** observations that stand out from the rest of the data that may be suspected outliers

7

Outline

What are some things to think about when guessing the distribution of any variable?



Clicker question

Which of these is most likely to have a roughly symmetric distribution?

- (a) salaries of a random sample of people from North Carolina
- (b) weights of adult females
- (c) scores on an well-designed exam
- (d) last digits of phone numbers

8

Clicker question

Which of these is most likely to have a roughly symmetric distribution?

- (a) salaries of a random sample of people from North Carolina
- (b) *weights of adult females*
- (c) scores on an well-designed exam
- (d) *last digits of phone numbers*



8

Outline

Think about natural boundaries!



Application exercise: 1.1 Distributions of numerical variables

See the course website for instructions.

15

Outline

How do mean, median, and skewness relate?

median mean

Mean vs. median

Clicker question

How do the mean and median of the following two datasets compare?

Dataset 1: 30, 50, 70, 90
Dataset 2: 30, 50, 70, 1000

- (a) $\bar{x}_1 = \bar{x}_2, median_1 = median_2$
- (b) $\bar{x}_1 < \bar{x}_2, median_1 = median_2$
- (c) $\bar{x}_1 < \bar{x}_2, median_1 < median_2$
- (d) $\bar{x}_1 > \bar{x}_2, median_1 < median_2$
- (e) $\bar{x}_1 > \bar{x}_2, median_1 = median_2$

9

Mean vs. median

Clicker question

How do the mean and median of the following two datasets compare?

Dataset 1: 30, 50, 70, 90
Dataset 2: 30, 50, 70, 1000

- (a) $\bar{x}_1 = \bar{x}_2, median_1 = median_2$
- (b) $\bar{x}_1 < \bar{x}_2, median_1 = median_2$**
- (c) $\bar{x}_1 < \bar{x}_2, median_1 < median_2$
- (d) $\bar{x}_1 > \bar{x}_2, median_1 < median_2$
- (e) $\bar{x}_1 > \bar{x}_2, median_1 = median_2$

9

Mean vs. median

Clicker question

How do the mean and median of the following two datasets compare?

Dataset 1: 30, 50, 70, 90
Dataset 2: 30, 50, 70, 1000

- (a) $\bar{x}_1 = \bar{x}_2, median_1 = median_2$
- (b) $\bar{x}_1 < \bar{x}_2, median_1 = median_2$**
- (c) $\bar{x}_1 < \bar{x}_2, median_1 < median_2$
- (d) $\bar{x}_1 > \bar{x}_2, median_1 < median_2$
- (e) $\bar{x}_1 > \bar{x}_2, median_1 = median_2$

9

Mean vs. median

Clicker question

How do the mean and median of the following two datasets compare?

Dataset 1: 30, 50, 70, 90
 Dataset 2: 30, 50, 70, 1000

(a) $\bar{x}_1 = \bar{x}_2, median_1 = median_2$
 (b) $\bar{x}_1 < \bar{x}_2, median_1 = median_2$
 (c) $\bar{x}_1 < \bar{x}_2, median_1 < median_2$
 (d) $\bar{x}_1 > \bar{x}_2, median_1 < median_2$
 (e) $\bar{x}_1 > \bar{x}_2, median_1 = median_2$

The figure contains two dotplots. The first, titled 'Dotplot of Dataset 1', shows a symmetric distribution with data points (blue dots) at 30, 50, 70, and 90. The average (orange 'x') and median (green triangle) are both at 60. The second, titled 'Dotplot of Dataset 2', shows a right-skewed distribution with data points at 30, 50, 70, and 1000. The average (orange 'x') is at 60, while the median (green triangle) is at 30.

Outline

Why are the ways we calculate population standard deviation and sample standard deviation different?

The diagram shows a large blue circle labeled 'all graduates' containing many small blue dots. A smaller red circle labeled 'sample' is shown with arrows pointing to it from the population. A large blue Greek letter sigma (σ) is on the left, and a large green letter S is on the right.

Standard deviation and variance

- ▶ Most commonly used measure of variability is the *standard deviation*, which roughly measures the average deviation from the mean
 - Notation: σ : population standard deviation, s : sample standard deviation
- ▶ Calculating the standard deviation, for a population (rarely, if ever) and for a sample:

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{n}} \quad s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

- ▶ Square of the standard deviation is called the *variance*.

10

Standard deviation and variance

- ▶ Most commonly used measure of variability is the *standard deviation*, which roughly measures the average deviation from the mean
 - Notation: σ : population standard deviation, s : sample standard deviation
- ▶ Calculating the standard deviation, for a population (rarely, if ever) and for a sample:

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{n}} \quad s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

- ▶ Square of the standard deviation is called the *variance*.

10

More on SD

Why divide by $n - 1$ instead of n when calculating the sample standard deviation?

11

More on SD

Why divide by $n - 1$ instead of n when calculating the sample standard deviation?

Lose a “degree of freedom” for using an estimate (the sample mean, \bar{x}) (in estimating the sample variance/standard deviation.

11

More on SD

Why divide by $n - 1$ instead of n when calculating the sample standard deviation?

Lose a “degree of freedom” for using an estimate (the sample mean, \bar{x}) (in estimating the sample variance/standard deviation.

- More uncertainty introduced by using \bar{x} instead of μ .
- All else held equal (ie. $\bar{x} = \mu$) :
 - $n \rightarrow n-1$... gets smaller!

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{n}} \rightarrow s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

...gets larger.... incorporates more variation/uncertainty.


11

More on SD

Why divide by $n - 1$ instead of n when calculating the sample standard deviation?

Lose a “degree of freedom” for using an estimate (the sample mean, \bar{x}) (in estimating the sample variance/standard deviation.)

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



11


More on SD

Why divide by $n - 1$ instead of n when calculating the sample standard deviation?

Lose a “degree of freedom” for using an estimate (the sample mean, \bar{x} (in estimating the sample variance/standard deviation.)

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

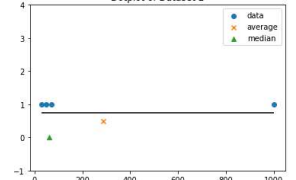
- ▶ To get rid of negatives so that observations equally distant from the mean are weighed equally.
- ▶ To weigh larger deviations more heavily.



11

Outline

Center:
When should we use mean vs. median?



Spread:
When should we use standard deviation vs. IQR vs. range?

Outline

- I. Readiness assessment
- II. Housekeeping
- III. Main ideas
 1. Analysis work flow:
 1. Always start your exploration with a visualization
 2. Single Numerical Variable
 1. When describing numerical distributions discuss shape, center, spread, and unusual observations
 2. Robust statistics are not easily affected by outliers and extreme skew
 3. Use box plots to display quartiles, median, and outliers
- IV. Application exercises
- V. Summary

Range and IQR

Clicker question

True / False: The range is always at least as large as the IQR for a given dataset.

(a) Yes
 (b) No

12

Range and IQR

Clicker question

True / False: The range is always at least as large as the IQR for a given dataset.

(a) Yes
(b) No

$Range = max - min, IQR = Q3 - Q1$

12

Range and IQR

Clicker question

True / False: The range is always at least as large as the IQR for a given dataset.

(a) Yes
(b) No

$Range = max - min, IQR = Q3 - Q1$

Is the range or the IQR more robust to outliers?

12

Range and IQR

Is the range or the IQR more robust to outliers?

IQR

12

Robust statistics


- ▶ Mean and standard deviation are **easily affected by extreme observations** since the value of each data point contributes to their calculation.
- ▶ Median and IQR are **more robust to outliers**.
- ▶ Therefore we choose **median & IQR** (over mean&SD) when describing skewed distributions.
- ▶ We choose **mean & SD** when describing symmetric distributions, as they are more useful in using mathematical theory to make inferences.

13

Outline

How do we determine if a data point is an outlier in a numerical distribution?

How do we construct a boxplot?



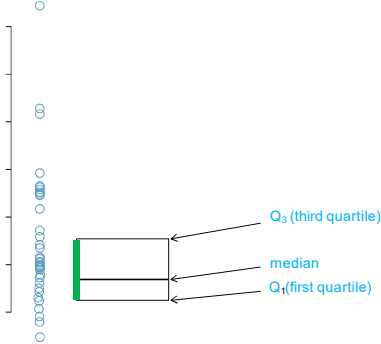
<https://www.kdnuggets.com/2017/01/3-methods-deal-outliers.htm>

Outline

- I. Readiness assessment
- II. Housekeeping
- III. Main ideas
 1. Analysis work flow:
 1. Always start your exploration with a visualization
 2. Single Numerical Variable
 1. When describing numerical distributions discuss shape, center, spread, and unusual observations
 2. Robust statistics are not easily affected by outliers and extreme skew
 3. Use box plots to display quartiles, median, and outliers
- IV. Application exercises
- V. Summary

Box plot

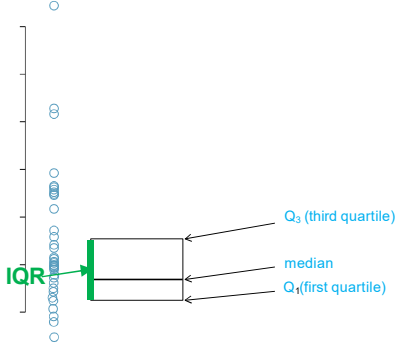
A *box plot* visualizes the median, the quartiles, and suspected outliers. An **outlier** is defined as an observation more than $1.5 \times \text{IQR}$ away from the quartiles (Q1 and Q3).



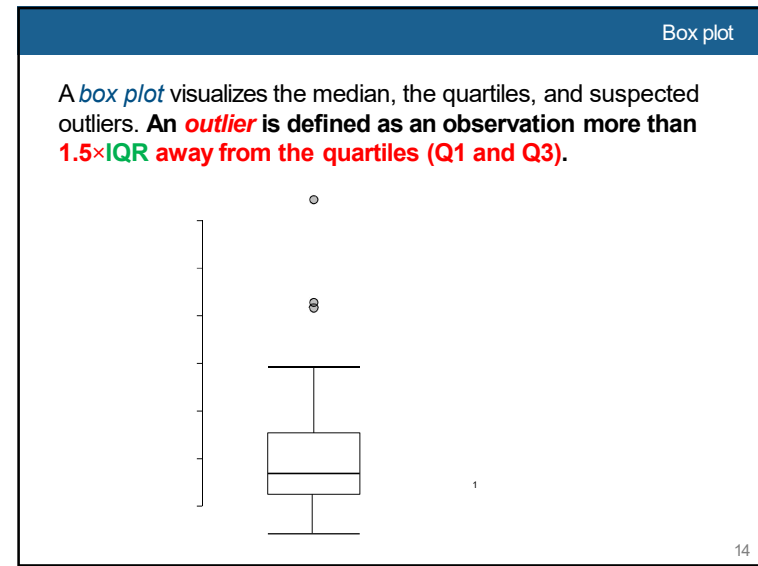
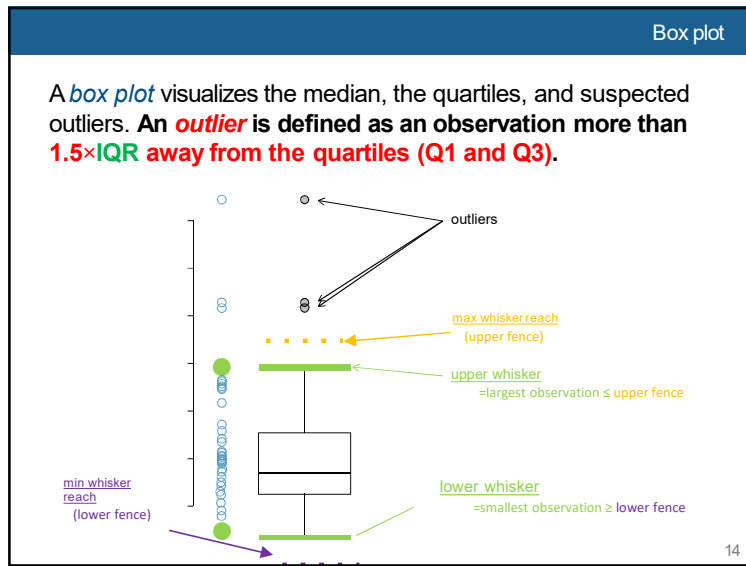
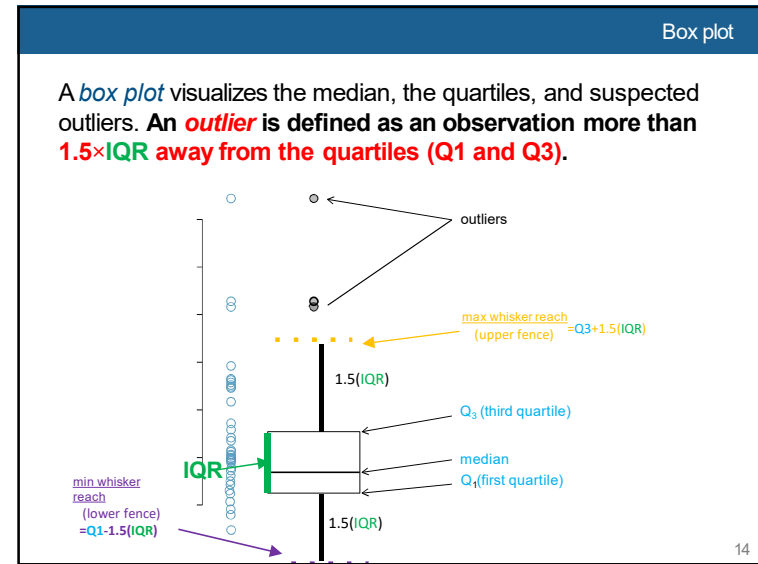
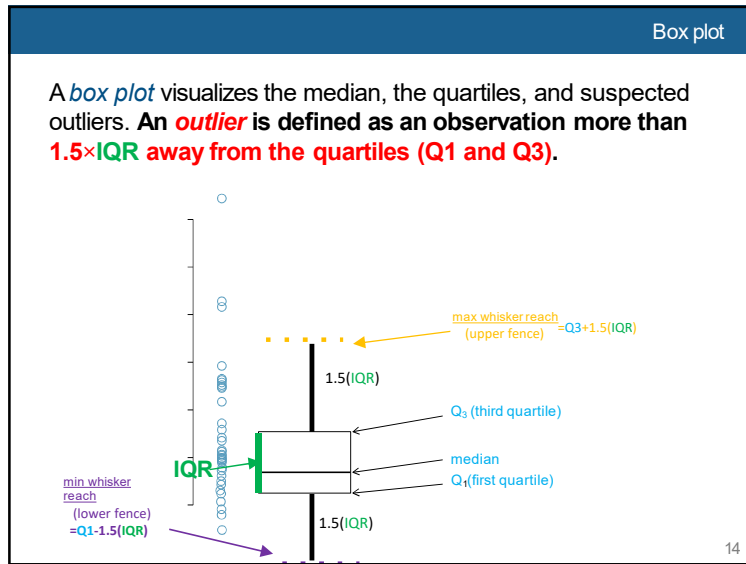
14

Box plot








A *box plot* visualizes the median, the quartiles, and suspected outliers. An **outlier** is defined as an observation more than $1.5 \times \text{IQR}$ away from the quartiles (Q1 and Q3).



14



Outline

- I. Readiness assessment
- II. Housekeeping
- III. Main ideas
 - 1. Analysis work flow:
 - 1.   Always start your exploration with a visualization
 - 2. Single Numerical Variable
 - 1.   When describing numerical distributions discuss shape, center, spread, and unusual observations
 - 2.  Robust statistics are not easily affected by outliers and extreme skew
 - 3.   Use box plots to display quartiles, median, and outliers
- IV. Application exercises
- V. Summary