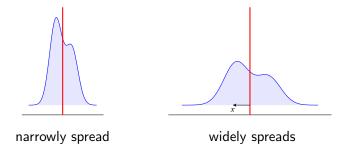
# American University of Beirut STAT 210: Elementary Statistics for Sciences 2022–2023 Fall

Siamak Taati

Chapter 3
Measures of center and variation
Part 3

# Measures of center and variation

## Measures of spread: variance and standard deviation



The variance and standard deviation measure the "typical" deviation of the data values from the mean.

deviation from the population mean  $= x - \mu$ deviation from the sample mean  $= x - \overline{x}$ 

## Measures of spread: variance

Sample vs. population variance For a numerical variable x,

population variance of 
$$x$$
: 
$$\sigma^2 \coloneqq \frac{\sum (x-\mu)^2}{N}$$
 sample variance of  $x$ : 
$$s^2 \coloneqq \frac{\sum (x-\overline{x})^2}{n-1}$$

where N is the population size and n is the sample size.

## Measures of spread: variance

Sample vs. population variance For a numerical variable x,

population variance of 
$$x$$
: 
$$\sigma^2 \coloneqq \frac{\sum (x-\mu)^2}{N}$$
 sample variance of  $x$ : 
$$s^2 \coloneqq \frac{\sum (x-\bar{x})^2}{n-1}$$

where N is the population size and n is the sample size.

An alternative way to calculate the variances:

population variance of 
$$x$$
: 
$$\sigma^2 = \frac{\sum x^2 - (\sum x)^2/N}{N}$$
 sample variance of  $x$ : 
$$s^2 := \frac{\sum x^2 - (\sum x)^2/n}{n-1}$$

## Measures of spread: variance

Sample vs. population variance For a numerical variable x,

population variance of 
$$x$$
: 
$$\sigma^2 \coloneqq \frac{\sum (x - \mu)^2}{N}$$
 sample variance of  $x$ : 
$$s^2 \coloneqq \frac{\sum (x - \overline{x})^2}{n - 1}$$

where N is the population size and n is the sample size.

An alternative way to calculate the variances:

population variance of 
$$x$$
: 
$$\sigma^2 = \frac{\sum x^2 - (\sum x)^2 / N}{N}$$
 sample variance of  $x$ : 
$$s^2 := \frac{\sum x^2 - (\sum x)^2 / n}{n-1}$$

## Example

A list of values and its corresponding frequency table:

1	4	5	5	5	5	1	5	3	3

value	frequency
1	2
3	2
4	1
5	5

## Example

A list of values and its corresponding frequency table:

1	4	5	5	5	5	1	5	3	3

value	frequency
1	2
3	2
4	1
5	5

Q What is the sample variance?

## Example

A list of values and its corresponding frequency table:

1	4	5	5	5	5	1	5	3	3

We calculated before: | mean = 3.7 |.

value	frequency
1	2
3	2
4	1
5	5

Q) What is the sample variance?

## Example

A list of values and its corresponding frequency table:

1	4	5	5	5	5	1	5	3	3

value	frequency
1	2
3	2
4	1
5	5

We calculated before: mean = 3.7.

Q What is the sample variance?

## Method 1 (definition):

$$s^2 = \frac{(1-3.7)^2 + (4-3.7)^2 + (5-3.7)^2 + (5-3.7)^2 + (5-3.7)^2}{+(5-3.7)^2 + (1-3.7)^2 + (5-3.7)^2 + (3-3.7)^2 + (3-3.7)^2}{10-1} \approx \boxed{2.6778}$$

#### Example

A list of values and its corresponding frequency table:

1	4	5	5	5	5	1	5	3	3

value	frequency
1	2
3	2
4	1
5	5

We calculated before: mean = 3.7.

Q What is the sample variance?

## Method 1 (definition):

$$s^2 = \frac{(1-3.7)^2 + (4-3.7)^2 + (5-3.7)^2 + (5-3.7)^2 + (5-3.7)^2}{10-1} \approx \boxed{2.6778}$$

## Method 2 (definition + frequency table):

$$s^2 = \frac{2 \times (1 - 3.7)^2 + 2 \times (3 - 3.7)^2 + 1 \times (4 - 3.7)^2 + 5 \times (5 - 3.7)^2}{10 - 1} \approx \boxed{2.6778} \; .$$

## Example

A list of values and its corresponding frequency table:

1	4	5	5	5	5	1	5	3	3

value	frequency
1	2
3	2
4	1
5	5

We calculated before: mean = 3.7.

Q What is the sample variance?

Method 1 (definition):

$$s^{2} = \frac{\frac{(1-3.7)^{2}}{(1-3.7)^{2}} + (4-3.7)^{2} + (5-3.7)^{2} + (5-3.7)^{2} + (5-3.7)^{2} + (5-3.7)^{2} + (5-3.7)^{2} + (5-3.7)^{2} + (3-3.7)^{2} + (3-3.7)^{2}}{10-1} \approx \boxed{2.6778}.$$

 $\underline{\mathsf{Method}\ 2}\ (\mathsf{definition} + \mathsf{frequency\ table})$ :

$$s^2 = \frac{2 \times (1 - 3.7)^2 + 2 \times (3 - 3.7)^2 + 1 \times (4 - 3.7)^2 + 5 \times (5 - 3.7)^2}{10 - 1} \approx \boxed{2.6778} \, .$$

#### Example

A list of values and its corresponding frequency table:

1	4	5	5	5	5	1	5	3	3

value	frequency
1	2
3	2
4	1
5	5

We calculated before: mean = 3.7.

Q What is the sample variance?

Method 1 (definition):

$$s^2 = \frac{(1-3.7)^2 + (4-3.7)^2 + (5-3.7)^2 + (5-3.7)^2 + (5-3.7)^2}{+(5-3.7)^2 + (1-3.7)^2 + (5-3.7)^2 + \frac{(3-3.7)^2}{(3-3.7)^2} + \frac{(3-3.7)^2}{(3-3.7)^2}}{10-1} \approx \boxed{2.6778}.$$

Method 2 (definition + frequency table):

$$s^2 = \frac{2 \times (1 - 3.7)^2 + \frac{2 \times (3 - 3.7)^2 + 1 \times (4 - 3.7)^2 + 5 \times (5 - 3.7)^2}{10 - 1} \approx \boxed{2.6778}.$$

#### Example

A list of values and its corresponding frequency table:

1 4 5 5 5 5 1 5 3 3

value	frequency
1	2
3	2
4	1
5	5

We calculated before: mean = 3.7.

Q What is the sample variance?

## Method 1 (definition):

$$s^2 = \frac{(1-3.7)^2 + \frac{(4-3.7)^2}{(4-3.7)^2 + (5-3.7)^2 + (5-3.7)^2 + (5-3.7)^2}{+(5-3.7)^2 + (1-3.7)^2 + (5-3.7)^2 + (3-3.7)^2 + (3-3.7)^2}{10-1} \approx \boxed{2.6778}$$

## Method 2 (definition + frequency table):

$$s^2 = \frac{2 \times (1 - 3.7)^2 + 2 \times (3 - 3.7)^2 + \frac{1 \times (4 - 3.7)^2}{10 - 1} + 5 \times (5 - 3.7)^2}{10 - 1} \approx \boxed{2.6778}.$$

#### Example

A list of values and its corresponding frequency table:

1	4	5	5	5	5	1	5	3	3

value	frequency
1	2
3	2
4	1
5	5

We calculated before: mean = 3.7.

Q What is the sample variance?

Method 1 (definition):

$$s^{2} = \frac{(1 - 3.7)^{2} + (4 - 3.7)^{2} + (5 - 3.7)^{2} + (5 - 3.7)^{2} + (5 - 3.7)^{2}}{10 - 1} \approx \boxed{2.6778}$$

 $\underline{\mathsf{Method}\ 2}\ (\mathsf{definition}\ +\ \mathsf{frequency}\ \mathsf{table})$ :

$$s^2 = \frac{2 \times (1 - 3.7)^2 + 2 \times (3 - 3.7)^2 + 1 \times (4 - 3.7)^2 + \frac{5 \times (5 - 3.7)^2}{10 - 1} \approx \boxed{2.6778} \,.$$

## Example

A list of values and its corresponding frequency table:

1	4	5	5	5	5	1	5	3	3

We calculated before: mean = 3.7.

value	frequency
1	2
3	2
4	1
5	5

Q What is the sample variance?

#### Example

A list of values and its corresponding frequency table:

1 4 5 5 5 5 1 5 3 3

value	frequency
1	2
3	2
4	1
5	5

We calculated before: mean = 3.7.

Q What is the sample variance?

$$\sum x = 1 + 4 + 5 + 5 + 5 + 5 + 1 + 5 + 3 + 3 = 37$$

$$(\sum x)^2 = \sum x^2 = \frac{\sum x^2 - (\sum x)^2 / n}{n - 1} = \frac{\sum x}{n} = \frac{\sum x}{n} = \frac{n}{n}$$

#### Example

A list of values and its corresponding frequency table:

1 4 5 5 5 5 1 5 3 3

value	frequency
1	2
3	2
4	1
5	5

We calculated before: mean = 3.7.

Q What is the sample variance?

$$\sum x = 1 + 4 + 5 + 5 + 5 + 5 + 1 + 5 + 3 + 3 = 37$$
$$(\sum x)^2 = 37^2 = 1369$$
$$\sum x^2 = \frac{\sum x^2 - (\sum x)^2 / n}{n - 1} = \frac{\sum x}{n -$$

#### Example

A list of values and its corresponding frequency table:

1	4	5	5	5	5	1	5	3	3

value	frequency
1	2
3	2
4	1
5	5

We calculated before: mean = 3.7.

Q What is the sample variance?

$$\sum x = 1 + 4 + 5 + 5 + 5 + 5 + 1 + 5 + 3 + 3 = 37$$

$$(\sum x)^2 = 37^2 = 1369$$

$$\sum x^2 = 1^2 + 4^2 + 5^2 + 5^2 + 5^2 + 5^2 + 1^2 + 5^2 + 3^2 + 3^2 = 161$$

$$s^2 = \frac{\sum x^2 - (\sum x)^2 / n}{n - 1} =$$

#### Example

A list of values and its corresponding frequency table:

1	4	5	5	5	5	1	5	3	3	

value	frequency
1	2
3	2
4	1
5	5

We calculated before: mean = 3.7.

Q What is the sample variance?

$$\sum x = 1 + 4 + 5 + 5 + 5 + 5 + 1 + 5 + 3 + 3 = 37$$

$$(\sum x)^2 = 37^2 = 1369$$

$$\sum x^2 = 1^2 + 4^2 + 5^2 + 5^2 + 5^2 + 5^2 + 1^2 + 5^2 + 3^2 + 3^2 = 161$$

$$s^2 = \frac{\sum x^2 - (\sum x)^2 / n}{n - 1} = \frac{161 - 1369 / 10}{9} \approx \boxed{2.6778}.$$

#### Example

A list of values and its corresponding frequency table:

1	4	5	5	5	5	1	5	3	3

We calculated before: mean = 3.7.

value	frequency
1	2
3	2
4	1
5	5

Q What is the sample variance?

$$\sum x = (\sum x)^2 = \sum x^2 =$$

$$\sum x^2 = \frac{\sum x^2 - (\sum x)^2 / n}{n - 1} =$$

## Example

A list of values and its corresponding frequency table:

1	4	5	5	5	5	1	5	3	3

value	frequency
1	2
3	2
4	1
5	5

We calculated before: mean = 3.7.

Q What is the sample variance?

$$\sum x = 2 \times 1 + 2 \times 3 + 1 \times 4 + 5 \times 5 = 37$$
$$(\sum x)^{2} = \sum x^{2} = s^{2} = \frac{\sum x^{2} - (\sum x)^{2} / n}{n - 1} = s^{2}$$

#### Example

A list of values and its corresponding frequency table:

1	4	5	5	5	5	1	5	3	3

We calculated before:  $\boxed{\text{mean} = 3.7}$ .

value	frequency
1	2
3	2
4	1
5	5

(Q) What is the sample variance?

$$\sum x = 2 \times 1 + 2 \times 3 + 1 \times 4 + 5 \times 5 = 37$$
$$(\sum x)^2 = 37^2 = 1369$$
$$\sum x^2 = s^2 = \frac{\sum x^2 - (\sum x)^2 / n}{n - 1} = s^2$$

## Example

A list of values and its corresponding frequency table:

1	4	5	5	5	5	1	5	3	3

value	frequency
1	2
3	2
4	1
5	5

We calculated before: mean = 3.7.

Q What is the sample variance?

$$\sum x = 2 \times 1 + 2 \times 3 + 1 \times 4 + 5 \times 5 = 37$$
$$(\sum x)^2 = 37^2 = 1369$$
$$\sum x^2 = 2 \times 1^2 + 2 \times 3^2 + 1 \times 4^2 + 5 \times 5^2 = 161$$
$$s^2 = \frac{\sum x^2 - (\sum x)^2 / n}{n - 1} =$$

#### Example

A list of values and its corresponding frequency table:

val										
	3	3	5	1	5	5	5	5	4	1
		3.7	n = 3	near	e: r	efore	ed b	ulate	calc	We

 value
 frequency

 1
 2

 3
 2

 4
 1

 5
 5

Q What is the sample variance?

$$\sum x = 2 \times 1 + 2 \times 3 + 1 \times 4 + 5 \times 5 = 37$$

$$(\sum x)^2 = 37^2 = 1369$$

$$\sum x^2 = 2 \times 1^2 + 2 \times 3^2 + 1 \times 4^2 + 5 \times 5^2 = 161$$

$$s^2 = \frac{\sum x^2 - (\sum x)^2 / n}{n - 1} = \frac{161 - 1369 / 10}{9} \approx \boxed{2.6778}.$$

#### Scenario

A group of marine biologists are studying two species of fish. The population mean and standard deviation of the length of the fish in the two species are:

Species 1	Species 2
$\mu_1=5\mathrm{cm}$	$\mu_2=20\mathrm{cm}$
$\sigma_1=0.7\mathrm{cm}$	$\sigma_2=2\mathrm{cm}$

#### Scenario

A group of marine biologists are studying two species of fish. The population mean and standard deviation of the length of the fish in the two species are:

Species 1	Species 2
$\mu_1=5\mathrm{cm}$	$\mu_2=20\mathrm{cm}$
$\sigma_1=0.7\mathrm{cm}$	$\sigma_2=2{ m cm}$

Q The length of which species of fish has more variations?



#### Scenario

A group of marine biologists are studying two species of fish. The population mean and standard deviation of the length of the fish in the two species are:

Species 1	Species 2
$\mu_1=5\mathrm{cm}$	$\mu_2=20\mathrm{cm}$
$\sigma_1=0.7\mathrm{cm}$	$\sigma_2=2{ m cm}$

- Q The length of which species of fish has more variations?
- $\longrightarrow$  The absolute standard deviation of the length is larger for the 2nd species,

#### Scenario

A group of marine biologists are studying two species of fish. The population mean and standard deviation of the length of the fish in the two species are:

Species 1	Species 2
$\mu_1=5\mathrm{cm}$	$\mu_2=20\mathrm{cm}$
$\sigma_1=0.7\mathrm{cm}$	$\sigma_2=2{ m cm}$

- Q The length of which species of fish has more variations?
- → The absolute standard deviation of the length is larger for the 2nd species, but relative to the mean length, there is more variation within the 1st species:

$$\frac{\sigma_1}{\mu_1} = \frac{0.7}{5} = \boxed{14\%}$$
  $\frac{\sigma_2}{\mu_2} = \frac{2}{20} = \boxed{10\%}$ 



The coefficient of variation of a variable x is the ratio of its standard deviation over its mean:

population 
$$\mathsf{CV} \coloneqq \frac{\sigma}{\mu}$$
 sample  $\mathsf{CV} \coloneqq \frac{s}{\bar{x}}$ 

and is often written in percentage format.

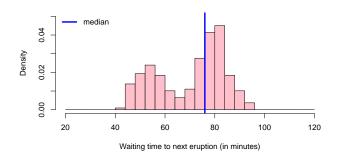
The coefficient of variation of a variable x is the ratio of its standard deviation over its mean:

population 
$$CV := \frac{\sigma}{\mu}$$
 sample  $CV := \frac{s}{\bar{x}}$ 

and is often written in percentage format.

population CV for species 1=14% population CV for species 2=10%

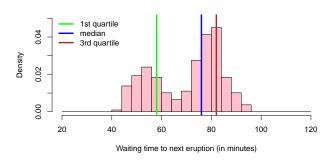
Recall that median of a list of data values divides the values into two halves: those <u>larger than</u> the median and those <u>smaller than</u> the median.



Old Faithful geyser

Recall that median of a list of data values divides the values into two halves: those <u>larger than</u> the median and those <u>smaller than</u> the median.

The quartiles further divide the data values into quarters.



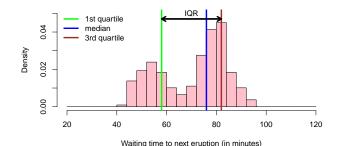
Old Faithful geyser



Recall that median of a list of data values divides the values into two halves: those <u>larger than</u> the median and those <u>smaller than</u> the median.

The quartiles further divide the data values into quarters.

The difference between the 1st and the 3rd quartiles is a measure of <u>variation</u>. It is called the <u>inter-quartile range</u> (IQR).



Old Faithful geyser



The quartiles partition the data values into three equal parts.

- ▶ The 2nd quartile (denoted by  $Q_2$ ) is the same as the median.
- ▶ The 1st quartile (denoted by  $Q_1$ ) is the median of the data values that are smaller than the median.
- The 3rd quartile (denoted by  $Q_3$ ) is the median of the data values that are larger than the median.

The inter-quartile range is the difference between the 1st and the 3rd quartiles:

$$IQR := Q_3 - Q_1$$
.

It is a measure of spread.

Back to the birds data set . . .

## Weight of females

```
98.8
      92.5
            87.3
                   103.5
                         98.3
                               99.1
                                    96.9
                                          103.0
                                                96.5
                                                     93.4
102.9
      101.1
            100.9
                   95.6
                         96.1
                               99.3
                                    93.3
                                          95.4
                                                88.5
                                                     97.6
92.2
      98.0
```

Back to the birds data set . . .

```
87.3
       88.5
             92.2
                  92.5 93.3
                              93.4
                                    95.4
                                          95.6
                                                 96.1
                                                       96.5
96.9 97.6
             98.0
                  98.3 98.8
                              99.1
                                    99.3
                                         100.9
                                                101.1
                                                       102.9
103.0
      103.5
```

Back to the birds data set . . .

```
87.3
       88.5
             92.2
                  92.5 93.3
                              93.4
                                    95.4
                                          95.6
                                                 96.1
                                                       96.5
96.9 97.6
             98.0
                  98.3 98.8
                              99.1
                                    99.3
                                         100.9
                                                101.1
                                                       102.9
103.0
      103.5
```

Back to the birds data set . . .

```
87.3 88.5 92.2 92.5 93.3 93.4 95.4 95.6 96.1 96.5 96.9 97.6 98.0 98.3 98.8 99.1 99.3 100.9 101.1 102.9 103.0 103.5 Q_2 = \text{median} = \frac{96.9 + 97.6}{2} = \boxed{97.25} \text{ grams}
```

Back to the birds data set . . .

87.3 88.5 92.2 92.5 93.3 93.4 95.4 95.6 96.1 96.5 96.9 97.6 98.0 98.3 98.8 99.1 99.3 100.9 101.1 102.9 103.0 103.5 
$$Q_2 = \text{median} = \frac{96.9 + 97.6}{2} = \boxed{97.25} \text{ grams}$$
 
$$Q_1 = \boxed{93.4} \text{ grams}$$

Back to the birds data set . . .

87.3 88.5 92.2 92.5 93.3 93.4 95.4 95.6 96.1 96.5 96.9 97.6 98.0 98.3 98.8 99.1 99.3 100.9 101.1 102.9 103.0 103.5 
$$Q_2 = \text{median} = \frac{96.9 + 97.6}{2} = \boxed{97.25} \text{ grams}$$
 
$$Q_1 = \boxed{93.4} \text{ grams}$$
 
$$Q_3 = \boxed{99.3} \text{ grams}$$

Back to the birds data set . . .

87.3 88.5 92.2 92.5 93.3 93.4 95.4 95.6 96.1 96.5 96.9 97.6 98.0 98.3 98.8 99.1 99.3 100.9 101.1 102.9 103.0 103.5 
$$Q_2 = \text{median} = \frac{96.9 + 97.6}{2} = \boxed{97.25} \text{ grams}$$
 
$$Q_1 = \boxed{93.4} \text{ grams}$$
 
$$Q_3 = \boxed{99.3} \text{ grams}$$
 
$$IQR = Q_3 - Q_1 = 99.3 - 93.4 = \boxed{5.9} \text{ grams}$$

Back to the birds data set . . .

#### Weight of females (sorted)

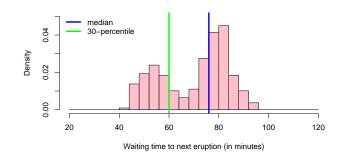
87.3 88.5 92.2 92.5 93.3 93.4 95.4 95.6 96.1 96.5 96.9 97.6 98.0 98.3 98.8 99.1 99.3 100.9 101.1 102.9 103.0 103.5 
$$Q_2 = \text{median} = \frac{96.9 + 97.6}{2} = \boxed{97.25} \text{ grams}$$
 
$$Q_1 = \boxed{93.4} \text{ grams}$$
 
$$Q_3 = \boxed{99.3} \text{ grams}$$
 
$$IQR = Q_3 - Q_1 = 99.3 - 93.4 = \boxed{5.9} \text{ grams}$$

Exercise: Compare the IQR for the sampled male and female birds.



The percentiles are similar to quartiles, but with other ratios.

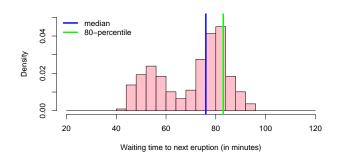
The 30-percentile divides the data values into the lower 30% and the upper 70%.



Old Faithful geyser

The percentiles are similar to quartiles, but with other ratios.

The 30-percentile divides the data values into the lower 30% and the upper 70%.



The 30-percentile of a list of data values is a value q with the following property:

▶ 30-percent of the data values are smaller than (or equal to) q 70-percent of them are larger (or equal to) than q.

The 30-percentile of a list of data values is a value q with the following property:

▶ 30-percent of the data values are smaller than (or equal to) q 70-percent of them are larger (or equal to) than q.

The 30-percentile of a list of n data values can be calculated as follows:

The 30-percentile of a list of data values is a value q with the following property:

▶ 30-percent of the data values are smaller than (or equal to) q 70-percent of them are larger (or equal to) than q.

The 30-percentile of a list of n data values can be calculated as follows:

Sort all the data values.

The 30-percentile of a list of data values is a value q with the following property:

▶ 30-percent of the data values are smaller than (or equal to) q 70-percent of them are larger (or equal to) than q.

The 30-percentile of a list of n data values can be calculated as follows:

- Sort all the data values.
- ▶ Let k be  $\frac{30}{100} \times n$  rounded up to the next whole number.

The 30-percentile of a list of data values is a value q with the following property:

▶ 30-percent of the data values are smaller than (or equal to) q 70-percent of them are larger (or equal to) than q.

The 30-percentile of a list of n data values can be calculated as follows:

- Sort all the data values.
- ▶ Let k be  $\frac{30}{100} \times n$  rounded up to the next whole number.
- ▶ The 30-percentile is the *k*-th element in the sorted list.

The 30-percentile of a list of data values is a value q with the following property:

▶ 30-percent of the data values are smaller than (or equal to) q 70-percent of them are larger (or equal to) than q.

The 30-percentile of a list of n data values can be calculated as follows:

- Sort all the data values.
- ▶ Let k be  $\frac{30}{100} \times n$  rounded up to the next whole number.
- ▶ The 30-percentile is the *k*-th element in the sorted list.

#### Remarks

- ▶ 25-percentile  $\approx Q_1$ .
- ▶ 50-percentile  $\approx Q_2 = \text{median}$ .
- ▶ 75-percentile  $\approx Q_3$ .

The 30-percentile of a list of data values is a value q with the following property:

▶ 30-percent of the data values are smaller than (or equal to) q 70-percent of them are larger (or equal to) than q.

The 30-percentile of a list of n data values can be calculated as follows:

- Sort all the data values.
- ▶ Let k be  $\frac{30}{100} \times n$  rounded up to the next whole number.
- ▶ The 30-percentile is the *k*-th element in the sorted list.

#### Remarks

- ▶ 25-percentile  $\approx Q_1$ .
- ▶ 50-percentile  $\approx Q_2 = \text{median}$ .
- ▶ 75-percentile  $\approx Q_3$ .
- ▶ In our convention, the values might be somewhat different!



Back to the birds data set . . .

#### Weight of females

```
98.8
      92.5
             87.3
                   103.5 98.3
                               99.1
                                    96.9
                                          103.0
                                                96.5
                                                     93.4
102.9
      101.1
            100.9
                   95.6
                         96.1
                               99.3
                                    93.3
                                          95.4
                                                88.5
                                                     97.6
92.2
      98.0
```

Q What is the 30-percentile of the data values?

Back to the birds data set . . .

#### Weight of females (sorted)

```
87.3
      88.5
           92.2
                  92.5 93.3
                             93.4
                                  95.4
                                        95.6
                                               96.1
                                                     96.5
96.9 97.6
           98.0
                  98.3 98.8
                             99.1
                                  99.3
                                        100.9
                                              101.1
                                                     102.9
103.0
      103.5
```

What is the 30-percentile of the data values?

 $\longrightarrow$  Sort the values!

Back to the birds data set . . .

```
87.3 88.5 92.2 92.5 93.3 93.4 95.4 95.6 96.1 96.5 96.9 97.6 98.0 98.3 98.8 99.1 99.3 100.9 101.1 102.9 103.0 103.5
```

- What is the 30-percentile of the data values?
  - $\longrightarrow$  Sort the values!
  - $\longrightarrow \frac{30}{100} \times 22 = 6.6$ , which is rounded up to 7.

Back to the birds data set . . .

```
87.3 88.5 92.2 92.5 93.3 93.4 <mark>95.4</mark> 95.6 96.1 96.5 96.9 97.6 98.0 98.3 98.8 99.1 99.3 100.9 101.1 102.9 103.0 103.5
```

- What is the 30-percentile of the data values?
  - $\longrightarrow$  Sort the values!
  - $\longrightarrow \frac{30}{100} \times 22 = 6.6$ , which is rounded up to 7.
  - $\rightarrow$  30-percentile = 95.4 grams.

Back to the birds data set . . .

```
88.5
           92.2
                                  95.4 95.6
                                               96.1
87.3
                  92.5 93.3 93.4
                                                     96.5
96.9 97.6
           98.0
                  98.3 98.8
                             99.1
                                  99.3
                                        100.9
                                              101.1
                                                     102.9
103.0
      103.5
```

- What is the 30-percentile of the data values?
  - $\longrightarrow$  Sort the values!
  - $\longrightarrow \frac{30}{100} \times 22 = 6.6$ , which is rounded up to 7.
  - $\rightarrow$  30-percentile = 95.4 grams.
- Q What is the percentile rank of 98.3?

Back to the birds data set . . .

```
88.5
           92.2
                                  95.4 95.6
                                              96.1
87.3
                  92.5 93.3 93.4
                                                     96.5
96.9 97.6
           98.0
                  98.3 98.8
                            99.1
                                  99.3
                                       100.9
                                              101.1
                                                     102.9
103.0
      103.5
```

- What is the 30-percentile of the data values?
  - $\longrightarrow$  Sort the values!
  - $\longrightarrow \frac{30}{100} \times 22 = 6.6$ , which is rounded up to 7.
  - $\rightarrow$  30-percentile = 95.4 grams.
- Q What is the percentile rank of 98.3?
  - $\longrightarrow$  There are 13 values smaller than 98.3.

Back to the birds data set . . .

```
87.3 88.5 92.2 92.5 93.3 93.4 95.4 95.6 96.1 96.5 96.9 97.6 98.0 98.3 98.8 99.1 99.3 100.9 101.1 102.9 103.0 103.5
```

- $\mathbb{Q}$  What is the 30-percentile of the data values?
  - $\longrightarrow$  Sort the values!
  - $\longrightarrow \frac{30}{100} \times 22 = 6.6$ , which is rounded up to 7.
  - $\rightarrow$  30-percentile = 95.4 grams.
- Q What is the percentile rank of 98.3?
  - $\longrightarrow$  There are 13 values smaller than 98.3.
  - $\longrightarrow$  percentile rank of  $98.3 = \frac{13}{22} \approx \boxed{59.1\%}$ .

The percentile rank of a value q in a list of n data values is the percentage of the values in the list that are smaller than q:

```
(percentile rank of q) := (percentage of data values < q)
```

The percentile rank of a value q in a list of n data values is the percentage of the values in the list that are smaller than q:

```
(percentile rank of q) := (percentage of data values < q)
```

In other words, the percentage rank of q is simply the cumulative percentage at q.

Remark (slight inconsistency)

Consider the following (sorted) list of 10 data values:

$$-5$$
  $-4$   $-3$   $-1$   $-1$  0 0 1 2 4

#### Remark (slight inconsistency)

Consider the following (sorted) list of 10 data values:

$$-5$$
  $-4$   $-3$   $-1$   $-1$  0 0 1 2 4

 $\bigcirc$  What is the percentile rank of -3?

#### Remark (slight inconsistency)

Consider the following (sorted) list of 10 data values:

$$-5$$
  $-4$   $-3$   $-1$   $-1$  0 0 1 2 4

 $\bigcirc$  What is the percentile rank of -3?

$$\longrightarrow \frac{2}{10} = 20\%$$
.

#### Remark (slight inconsistency)

Consider the following (sorted) list of 10 data values:

$$-5$$
  $-4$   $-3$   $-1$   $-1$  0 0 1 2 4

 $\bigcirc$  What is the percentile rank of -3?

$$\longrightarrow \frac{2}{10} = 20\%$$
.

 $\bigcirc$  What is the percentile rank of -4?

#### Remark (slight inconsistency)

Consider the following (sorted) list of 10 data values:

$$-5$$
  $-4$   $-3$   $-1$   $-1$  0 0 1 2 4

 $\bigcirc$  What is the percentile rank of -3?

$$\longrightarrow \ \tfrac{2}{10} = \boxed{20\%}.$$

 $\bigcirc$  What is the percentile rank of -4?

$$\longrightarrow \frac{1}{10} = \boxed{10\%}.$$

#### Remark (slight inconsistency)

Consider the following (sorted) list of 10 data values:

$$-5$$
  $-4$   $-3$   $-1$   $-1$  0 0 1 2 4

 $\bigcirc$  What is the percentile rank of -3?

$$\longrightarrow \frac{2}{10} = 20\%$$
.

 $\bigcirc$  What is the percentile rank of -4?

$$\longrightarrow \frac{1}{10} = \boxed{10\%}$$
.

Q What is the 20-percentile?

#### Remark (slight inconsistency)

Consider the following (sorted) list of 10 data values:

$$-5$$
  $-4$   $-3$   $-1$   $-1$  0 0 1 2 4

 $\bigcirc$  What is the percentile rank of -3?

$$\longrightarrow \frac{2}{10} = 20\%$$
.

 $\bigcirc$  What is the percentile rank of -4?

$$\longrightarrow \frac{1}{10} = \boxed{10\%}$$
.

Q What is the 20-percentile?

 $\longrightarrow \frac{20}{100} \times 10 = 2$ , which does not require rounding.

#### Remark (slight inconsistency)

Consider the following (sorted) list of 10 data values:

$$-5$$
  $-4$   $-3$   $-1$   $-1$  0 0 1 2 4

 $\bigcirc$  What is the percentile rank of -3?

$$\longrightarrow \frac{2}{10} = 20\%$$
.

 $\bigcirc$  What is the percentile rank of -4?

$$\longrightarrow \frac{1}{10} = \boxed{10\%}$$
.

Q What is the 20-percentile?

 $\longrightarrow \frac{20}{100} \times 10 = 2$ , which does not require rounding.

 $\longrightarrow$  Hence, the 20-percentile is |-4|.

#### Remark (slight inconsistency)

Consider the following (sorted) list of 10 data values:

$$-5$$
  $-4$   $-3$   $-1$   $-1$  0 0 1 2 4

 $\bigcirc$  What is the percentile rank of -3?

$$\longrightarrow \frac{2}{10} = 20\%$$
.

 $\bigcirc$  What is the percentile rank of -4?

$$\longrightarrow \frac{1}{10} = \boxed{10\%}$$
.

Q What is the 20-percentile?

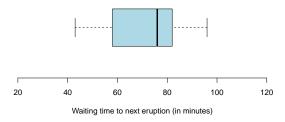
 $\longrightarrow \frac{20}{100} \times 10 = 2$ , which does not require rounding.

 $\longrightarrow$  Hence, the 20-percentile is  $\left| -4 \right|$ .

This inconsistency is partly due to our (simplifying) convention, and partly inevitable.

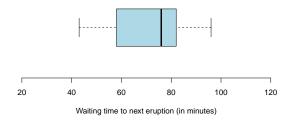
# Summarizing and visualizing data: box plots

#### Old Faithful geyser

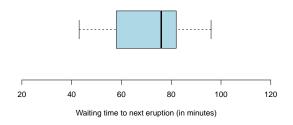


# Summarizing and visualizing data: box plots

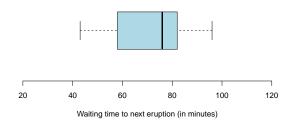
#### Old Faithful geyser



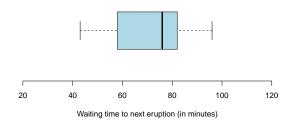
Q What does the thick vertical line stand for?



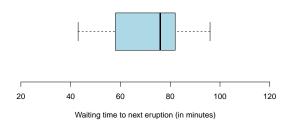
- What does the thick vertical line stand for?
  - $\longrightarrow$  Median



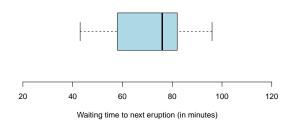
- Q What does the thick vertical line stand for?
  - --> Median
- Q What does the box stand for?



- Q What does the thick vertical line stand for?
  - --> Median
- What does the box stand for?
  - → The 1st and the 3rd quartiles

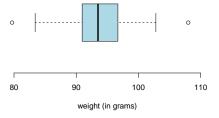


- Q What does the thick vertical line stand for?
  - --> Median
- Q What does the box stand for?
  - $\longrightarrow$  The 1st and the 3rd quartiles
- Q What do the whiskers stand for?

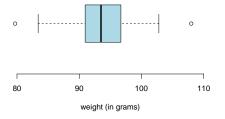


- Q What does the thick vertical line stand for?
  - → Median
- Q What does the box stand for?
  - $\longrightarrow$  The 1st and the 3rd quartiles
- Q What do the whiskers stand for?
  - → The minimum and the maximum

### Sampled male birds

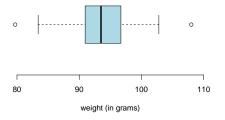


### Sampled male birds



What do the small circles stand for?

#### Sampled male birds

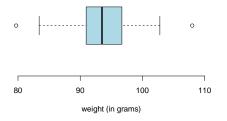


Q What do the small circles stand for?

→ The outliers

[Mild outliers; Tukey's convention]

#### Sampled male birds

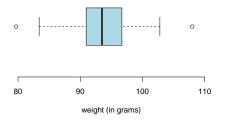


- Q What do the small circles stand for?
  - $\longrightarrow$  The outliers

[Mild outliers; Tukey's convention]

Q What do the whiskers stand for?

### Sampled male birds

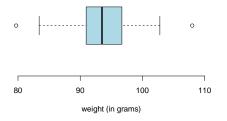


- Q What do the small circles stand for?
  - $\longrightarrow$  The outliers

[Mild outliers; Tukey's convention]

- What do the whiskers stand for?
  - $\longrightarrow$  The minimum and the maximum excluding the outliers

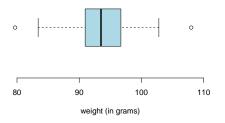
### Sampled male birds



### Identifying outliers (Tukey's convension)

Outliers are those observations that are farther than 1.5 IQR from the box.

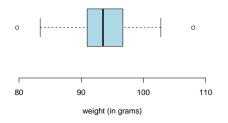
#### Sampled male birds



### Identifying outliers (Tukey's convension)

- Outliers are those observations that are farther than 1.5 IQR from the box.
- ► The observations that are farther than 3 IQR from the box are considered extreme outliers.

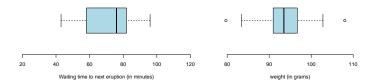
#### Sampled male birds

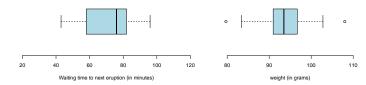


### Identifying outliers (Tukey's convension)

- Outliers are those observations that are farther than 1.5 IQR from the box.
- ► The observations that are farther than 3 IQR from the box are considered extreme outliers.
- The non-extreme outliers are considered mild outliers.

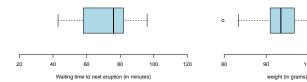






The box plot, visually summarizes:

► The center (median)

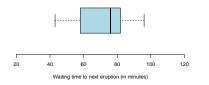


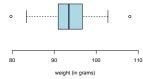
The box plot, visually summarizes:

- ▶ The center (median)
- The spread (IQR)

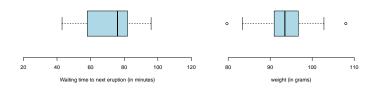
100

110

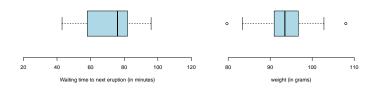




- The center (median)
- ▶ The spread (IQR)
- The skew



- The center (median)
- The spread (IQR)
- The skew
  - $\circ~Q_3-Q_2$  vs.  $Q_2-Q_1$
  - o Distances of the left and right whiskers from the box



- The center (median)
- ▶ The spread (IQR)
- ▶ The skew
  - $\circ \ Q_3 Q_2 \ \text{vs.} \ Q_2 Q_1$
  - o Distances of the left and right whiskers from the box
- ▶ The outliers