


4. Measure of Central Tendency and Dispersion

In the case of quantitative data, the central tendency and dispersion of the data are measured and analyzed.

- Central tendency : mean / median
- Dispersion : variance / standard deviation

Data for two quantitative variables are measured using a covariance and a correlation coefficient.

4.1 Measure of Central Tendency – Mean and Median

 Think	<p>The data obtained by taking a sample of 5 middle school students and surveying their weight are as follows:</p> <p style="text-align: center;">(Data 4.1) Weights of five middle school students (kg)</p> <div style="border: 1px solid black; padding: 5px; margin: 10px auto; width: fit-content;"> 63 60 65 55 77 </div>
Explore	<p>1) What kind of graph is used to find a representative value of these data?</p> <p>2) What would be a representative value for the weight of 5 students?</p>

- The **average** (or **mean**) is a measure of central tendency of quantitative data and is widely used as a representative value for the data. The mean is the sum of all data and divided by the number of data, which implies the center of gravity of the data. The mean is expressed as μ (read mu), and the mean of (Data 4.1) is obtained as follows:

$$\text{Mean} = \mu = \frac{63 + 60 + 65 + 55 + 77}{5} = \frac{320}{5} = 64$$

- When n number of data is expressed as x_1, x_2, \dots, x_n , mean is expressed by the following formula.

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

- In general, the mean is very appropriate as a representative value of the data, but when there is a very large or a small value in the data, it is greatly affected by this extreme value. In this case, a **median** can be used. The median is the value in the middle when the data are sorted in order. In (Data 4.1), the number of data is 5 which is an odd number, and the 3rd ($\frac{5+1}{2}$) number when data is sorted in ascending order is the median as follows:

(Data 4.1) is sorted in ascending order.

55 60 63 65 77

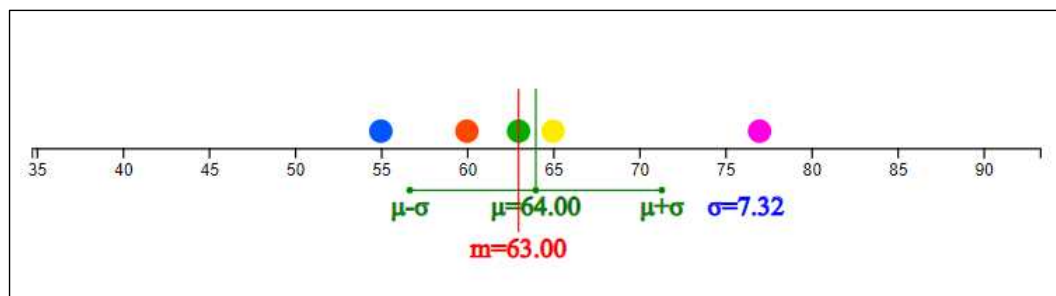
Median is the 3rd number in these sorted data which is 63.

- If the number of data is 6 which is an even number, how do we find the median? In this case, the median of the data is calculated as the average of the 3rd ($= \frac{6}{2}$) and 4th ($= \frac{6+2}{2}$) of the sorted data.
- Generally, a median is expressed as m , and, if the number of data is n , it is calculated as follows:

- 1) Data are sorted in ascending order.
- 2) Check whether the number of data n is an odd number or an even number.
- 3) If n is odd, $m = (\frac{n+1}{2})^{th}$ data in sorted data.

If n is even, $m = \text{Average of } (\frac{n}{2})^{th} \text{ data } (\frac{n+2}{2})^{th} \text{ data in sorted data.}$

- In order to see the overall distribution of the weight data, a stem and leaf plot or a histogram discussed in Chapter 3 can be considered, but a dot graph is more useful. In a dot graph, after obtaining the minimum and maximum values of data, the position of each data is calculated on the horizontal axis, and displayed as a dot.
- <Figure 4.1> is a dot graph for (Data 4.1). In proportion to the minimum value of 55 and the maximum value of 76, each data is displayed by a dot. The green line is the mean μ and the red line is the median m . In this data, the mean is located slightly to the right of the median because 77 of the data is located to the right of the other four data. That is, the mean is more sensitive to an extreme than the median.



<Figure 4.1> Dot graph of weight data

- If there are lots of data, it is time-consuming and difficult to obtain the mean and median manually as above. Let's find the representative value of the data using 『eStatH』 software.

Practice 4.1

Using 『eStatH』, draw a dot graph for the weights of 5 students (Data 4.1) and find the mean and median.

Solution

- Using the QR on the left, select 'Dot Graph – Mean / Standard Deviation' from the 『eStatH』 menu, then a window like <Figure 4.2> appears.
- Enter students' weight data in 'Data input'. (You can also copy and paste the data from the e-book)

Dot Graph - Mean/StdDev Menu

[Enter Data]

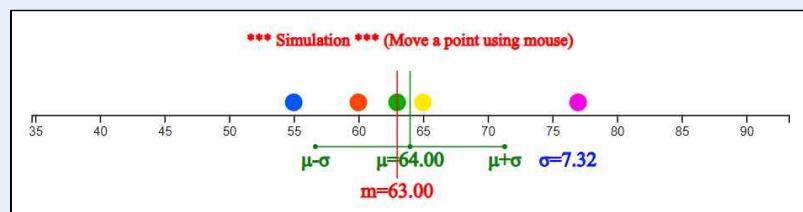
Sorted Data

Number of Data	<i>n</i>	<input type="text" value="5"/>	Minimum	<i>min</i>	<input type="text" value="55.00"/>
Total		<input type="text" value="320.00"/>	Maximum	<i>max</i>	<input type="text" value="77.00"/>
Mean	μ	<input type="text" value="64.00"/>	Range	<i>range</i>	<input type="text" value="22.00"/>
Median	<i>m</i>	<input type="text" value="63.00"/>	Variance	σ^2	<input type="text" value="53.60"/>
Mode	<i>mode</i>	<input type="text"/>	Std Deviation	σ	<input type="text" value="7.32"/>

<Figure 4.2> Data input of weight data for a dot graph



- When data are entered, the number of data, minimum, maximum, mean, and median are calculated immediately. If you click the [Execute] button, a dot graph appears as shown in <Figure 4.1> and the mean and median values are displayed.
- Below <Figure 4.1>, a simulation window as shown in <Figure 4.3> appears. In this simulation window, you move a point with the mouse to see the changes in the mean and median. For example, if you drag the rightmost point and move it to the right, the mean changes but the median does not. That is, the median is not affected by the extreme points.



<Figure 4.3> Simulation window to see a change in mean and median if you move a point

Practice 4.2

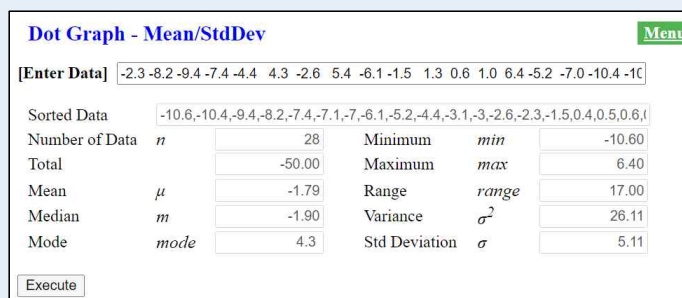
Using 『eStatH』, let's find the mean and median of the daily minimum temperature ([Practice 3.2]) in Seoul in February (Data 3.2).

(Data 3.2) Daily minimum temperature in February 2021 in Seoul (unit: degree in Celcius)

-2.3	-8.2	-9.4	-7.4	-4.4	4.3	-2.6	5.4	-6.1	-1.5
1.3	0.6	1.0	6.4	-5.2	-7.0	-10.4	-10.6	-7.1	5.5
4.7	0.4	-3.1	-3.0	0.7	0.5	4.3	3.2		

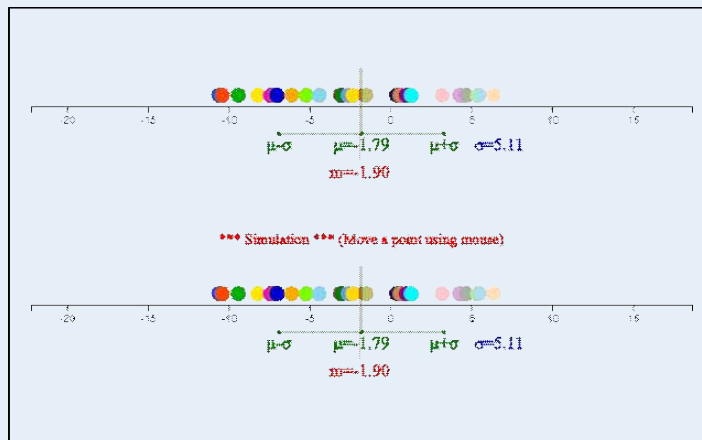
Solution

- If you select 'Dot Graph - Mean/Standard Deviation' from the 『eStatH』 menu using the QR on the left, the data input window as shown in <Figure 4.4> appears.




<Figure 4.4> Temperature data input for a dot graph


- When the daily minimum temperature data are entered in 'Data input' (you can copy and paste the data from the e-book), as shown in <Figure 4.4>, it shows immediately that the number of data is 28, mean -1.79 , median -1.90 , minimum -10.6 , maximum 6.4 degrees.
- If you click the [Execute] button, a dot graph as shown in <Figure 4.5> appears and the mean (μ) and median (m) are displayed. Below this dot graph, a simulation window appears where you can change a point with the mouse and observe the changes in the mean and median values.



<Figure 4.5> Dot graph of daily minimum temperature and its simulation window

Practice 4.2 Solution (Continued)	<ul style="list-style-type: none"> Looking at this dot graph, it can be seen that there is almost no difference between the mean and the median due to the absence of extreme value.
--	---

Exercise 4.1	<p>The following is data on the length of bicycle-only roads by 25 administrative districts in Seoul as of 2019 ([Exercise 3.1]). Use 『eStatH』 to draw a dot graph, and to find and analyze the representative values of data.</p>
	<p>(Data 3.3) Length of bicycle-only roads by 25 administrative districts in Seoul as of 2019 (unit: km)</p> <div style="border: 1px solid black; padding: 5px; margin: 5px 0;"> 24 15 23 20 30 24 7 8 7 12 28 27 19 35 41 42 11 8 37 13 20 29 53 93 42 </div>

Exercise 4.2	<p>The following is data on the maximum wind speed of typhoons that passed through Korea in 2020 ([Exercise 3.2]). Use 『eStatH』 to draw a dot graph, and to find and analyze the representative values of data.</p>
	<p>(Data 3.4) Maximum wind speed of typhoons that passed through Korea in 2020 (unit: m/sec)</p> <div style="border: 1px solid black; padding: 5px; margin: 5px 0;"> 40 22 21 29 19 22 24 45 49 55 24 27 29 35 19 24 35 40 56 24 21 43 18 </div>

A. Calculation of mean using frequency table

<p> Think</p>	<p>Assume that a frequency table of the academic achievement test scores of a middle school class is given as follows:</p> <p style="text-align: center; margin: 10px 0;">[Table 4.1] Frequency table of the academic achievement test scores of a middle school</p> <table border="1" style="margin: auto; border-collapse: collapse; text-align: center;"> <thead> <tr style="background-color: #d3d3d3;"> <th style="padding: 5px;">Class</th> <th style="padding: 5px;">Frequency</th> </tr> </thead> <tbody> <tr> <td style="padding: 5px;">$60 \leq \sim < 70$</td> <td style="padding: 5px;">2</td> </tr> <tr> <td style="padding: 5px;">$70 \sim 80$</td> <td style="padding: 5px;">10</td> </tr> <tr> <td style="padding: 5px;">$80 \sim 90$</td> <td style="padding: 5px;">15</td> </tr> <tr> <td style="padding: 5px;">$90 \sim 100$</td> <td style="padding: 5px;">3</td> </tr> <tr style="background-color: #d3d3d3;"> <td style="padding: 5px;">Total</td> <td style="padding: 5px;">30</td> </tr> </tbody> </table>	Class	Frequency	$60 \leq \sim < 70$	2	$70 \sim 80$	10	$80 \sim 90$	15	$90 \sim 100$	3	Total	30
Class	Frequency												
$60 \leq \sim < 70$	2												
$70 \sim 80$	10												
$80 \sim 90$	15												
$90 \sim 100$	3												
Total	30												
<p>Explore</p>	<p>How do we find the mean using this frequency table?</p>												

- When a frequency table is given rather than the raw data, the mean can be obtained approximately as follows using the middle values of each class interval.
- First, find the middle value of each class interval. Then, it is assumed that each class has the middle value as many as the frequency, and the mean is obtained using this approximated data.

[Table 4.2] Approximated data using the middle value of each class interval in the academic achievement test scores

Class	Middle value	Frequency	Approximated data
60 ≤ ~ < 70	65	2	65 65
70 ~ 80	75	5	75 75 75 75 75
80 ~ 90	85	10	85 85 85 85 85 85 85 85 85 85
90 ~ 100	95	3	95 95 95
Total		20	

- Mean is calculated as follows:

$$\begin{aligned}
 \text{Mean} &= \frac{65 + 65 + 75 + 75 + 75 + 75 + 75 + 85 + 85 + 85 + 85 + 85 + 85 + 85 + 85 + 85 + 85 + 85 + 95 + 95 + 95}{20} \\
 &= \frac{65 \times 2 + 75 \times 5 + 85 \times 10 + 95 \times 3}{20} \\
 &= \frac{1640}{20} = 82
 \end{aligned}$$

- Using ‘Frequency Distribution Polygon - Relative Frequency Comparison’ of 『eStatH』, the approximate mean of the frequency table can be obtained as shown in <Figure 4.6>. After entering the left value of the class interval and ‘Frequency 1’, click the [Execute] button.

<Figure 4.6> Mean calculation using a frequency table

4.2 Measure of Dispersion – Standard Deviation

👉 Think	<p>The quiz scores (out of 10) of five middle school students are as follows:</p> <p>(Data 4.2) The quiz scores (out of 10) of five middle school students</p> <div style="border: 1px solid #ccc; padding: 5px; display: inline-block;">6 8 7 4 10</div>
Explore	Is there a way to measure how scattered these data are?

- The degree to which data are scattered is called a dispersion. A simple measure of the dispersion is a **range** which is the maximum minus the minimum.

$$\text{Range} = \text{max} - \text{min}$$

In (Data 4.2), the maximum value is 10 and the minimum value is 4, so the range is 22.

$$\text{Range} = 10 - 4 = 6$$

- Since the range is too sensitive to extreme values, a variance or a standard deviation is generally used to measure the dispersion. The **variance** is obtained by squaring the distance between each data value and the mean, and dividing it by the number of data. Therefore, when the data are scattered far from the mean, the variance is large, and when the data are clustered around the mean, the variance is small. The variance is expressed as σ^2 (read as sigma squared).
- The mean of the data in (Data 4.2) is as follows:

$$\text{Mean} = \mu = \frac{6 + 8 + 7 + 4 + 10}{5} = \frac{35}{5} = 7$$

- The variance is calculated by squaring the distances from the mean to each data value to find the sum, and then finding the mean. That is, it is the average of squared distances.

$$\text{Variance} = \sigma^2 = \frac{(6-7)^2 + (8-7)^2 + (7-7)^2 + (4-7)^2 + (10-7)^2}{5} = \frac{20}{5} = 4$$

- When n number of data is expressed as x_1, x_2, \dots, x_n and the mean is expressed as μ , the variance can be expressed by the following formula.


Variance
$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \quad (n: \text{number of data})$$

- The **standard deviation** is defined as the square root of the variance and denoted by σ . The variance is not easy to interpret practically because it is the average of the squared distances, but the standard deviation is the square root of the variance, so it can be interpreted as a measure of the average distance between each data value and the mean.

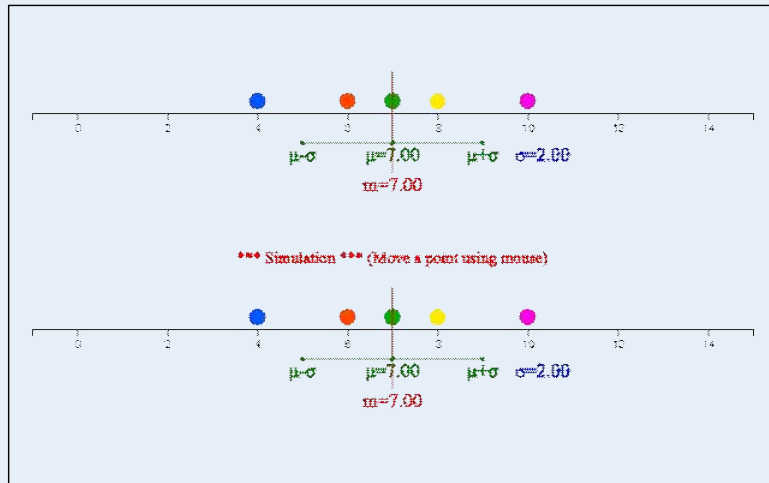
Standard deviation
$$\sigma = \sqrt{\sigma^2}$$

The standard deviation of (Data 4.2) is $\sigma = \sqrt{\sigma^2} = \sqrt{4} = 2$.

<p style="text-align: center;">Practice 4.3</p>	<p>Using 『eStatH』, draw a dot graph for the quiz scores of 5 sample students (Data 4.2) and find the mean and standard deviation.</p>																														
<p style="text-align: center;">Solution</p>	<ul style="list-style-type: none"> Using the QR on the left, select 'Dot Graph - Mean / Standard Deviation' from the 『eStatH』 menu. Then a window like <Figure 4.7> appears. Enter students' quiz scores in 'Data input'. (You can also copy and paste the material from the e-book) <div style="border: 1px solid black; padding: 10px; margin: 10px 0;"> <p style="text-align: center;">Dot Graph - Mean/StdDev Menu</p> <p>[Enter Data] <input style="width: 100%;" type="text" value="6 8 7 4 10"/></p> <p>Sorted Data <input style="width: 100%;" type="text" value="4,6,7,8,10"/></p> <table style="width: 100%; border-collapse: collapse;"> <tr> <td style="width: 30%;">Number of Data</td> <td style="width: 10%;"><i>n</i></td> <td style="width: 15%;"><input style="width: 50%;" type="text" value="5"/></td> <td style="width: 15%;">Minimum</td> <td style="width: 10%;"><i>min</i></td> <td style="width: 15%;"><input style="width: 50%;" type="text" value="4.00"/></td> </tr> <tr> <td>Total</td> <td></td> <td><input style="width: 50%;" type="text" value="35.00"/></td> <td>Maximum</td> <td><i>max</i></td> <td><input style="width: 50%;" type="text" value="10.00"/></td> </tr> <tr> <td>Mean</td> <td>μ</td> <td><input style="width: 50%;" type="text" value="7.00"/></td> <td>Range</td> <td><i>range</i></td> <td><input style="width: 50%;" type="text" value="6.00"/></td> </tr> <tr> <td>Median</td> <td><i>m</i></td> <td><input style="width: 50%;" type="text" value="7.00"/></td> <td>Variance</td> <td>σ^2</td> <td><input style="width: 50%;" type="text" value="4.00"/></td> </tr> <tr> <td>Mode</td> <td><i>mode</i></td> <td><input style="width: 50%;" type="text"/></td> <td>Std Deviation</td> <td>σ</td> <td><input style="width: 50%;" type="text" value="2.00"/></td> </tr> </table> <p style="text-align: center;"><input type="button" value="Execute"/></p> </div> <p style="text-align: center;"><Figure 4.7> Quiz data input for a dot graph</p> <ul style="list-style-type: none"> When the data are entered, the number of data, minimum, maximum, mean, median, variance and standard deviation are calculated. If you click the [Execute] button, a dot graph as shown in <Figure 4.8> appears and the mean, median, standard deviation, and a line of mean \pm standard deviation are displayed. Using the simulation window below the figure, you can check the change in the standard deviation by moving a data point with the mouse. The standard deviation is also affected by an extreme point. 	Number of Data	<i>n</i>	<input style="width: 50%;" type="text" value="5"/>	Minimum	<i>min</i>	<input style="width: 50%;" type="text" value="4.00"/>	Total		<input style="width: 50%;" type="text" value="35.00"/>	Maximum	<i>max</i>	<input style="width: 50%;" type="text" value="10.00"/>	Mean	μ	<input style="width: 50%;" type="text" value="7.00"/>	Range	<i>range</i>	<input style="width: 50%;" type="text" value="6.00"/>	Median	<i>m</i>	<input style="width: 50%;" type="text" value="7.00"/>	Variance	σ^2	<input style="width: 50%;" type="text" value="4.00"/>	Mode	<i>mode</i>	<input style="width: 50%;" type="text"/>	Std Deviation	σ	<input style="width: 50%;" type="text" value="2.00"/>
Number of Data	<i>n</i>	<input style="width: 50%;" type="text" value="5"/>	Minimum	<i>min</i>	<input style="width: 50%;" type="text" value="4.00"/>																										
Total		<input style="width: 50%;" type="text" value="35.00"/>	Maximum	<i>max</i>	<input style="width: 50%;" type="text" value="10.00"/>																										
Mean	μ	<input style="width: 50%;" type="text" value="7.00"/>	Range	<i>range</i>	<input style="width: 50%;" type="text" value="6.00"/>																										
Median	<i>m</i>	<input style="width: 50%;" type="text" value="7.00"/>	Variance	σ^2	<input style="width: 50%;" type="text" value="4.00"/>																										
Mode	<i>mode</i>	<input style="width: 50%;" type="text"/>	Std Deviation	σ	<input style="width: 50%;" type="text" value="2.00"/>																										



Practice 4.3
Solution
(continued)



<Figure 4.8> Dot graph with a line of mean \pm standard deviation

Practice 4.4

Using 『eStatH』, let's draw a dot graph for the daily minimum temperature ([Practice 3.2]) in Seoul in February (Data 3.2) and find the mean and standard deviation.

(Data 3.2) Daily minimum temperature ([Practice 3.2]) in Seoul in February in 2021 (degree in Celsius)

-2.3	-8.2	-9.4	-7.4	-4.4	4.3	-2.6	5.4	-6.1	-1.5
1.3	0.6	1.0	6.4	-5.2	-7.0	-10.4	-10.6	-7.1	5.5
4.7	0.4	-3.1	-3.0	0.7	0.5	4.3	3.2		

Solution

- If you select 'Dot Graph - Mean / Standard Deviation' from the 『eStatH』 menu that appears using the QR on the left, the data input window as shown in <Figure 4.9> appears.



Menu

Dot Graph - Mean/StdDev

[Enter Data]

Sorted Data	<input type="text" value="-10.6,-10.4,-9.4,-8.2,-7.4,-7.1,-7,-6.1,-5.2,-4.4,-3.1,-3,-2.6,-2.3,-1.5,0.4,0.5,0.6,0.7,0.5,4.3,3.2"/>				
Number of Data	<i>n</i>	<input type="text" value="28"/>	Minimum	<i>min</i>	<input type="text" value="-10.60"/>
Total		<input type="text" value="-50.00"/>	Maximum	<i>max</i>	<input type="text" value="6.40"/>
Mean	μ	<input type="text" value="-1.79"/>	Range	<i>range</i>	<input type="text" value="17.00"/>
Median	<i>m</i>	<input type="text" value="-1.90"/>	Variance	σ^2	<input type="text" value="26.11"/>
Mode	<i>mode</i>	<input type="text" value="4.3"/>	Std Deviation	σ	<input type="text" value="5.11"/>

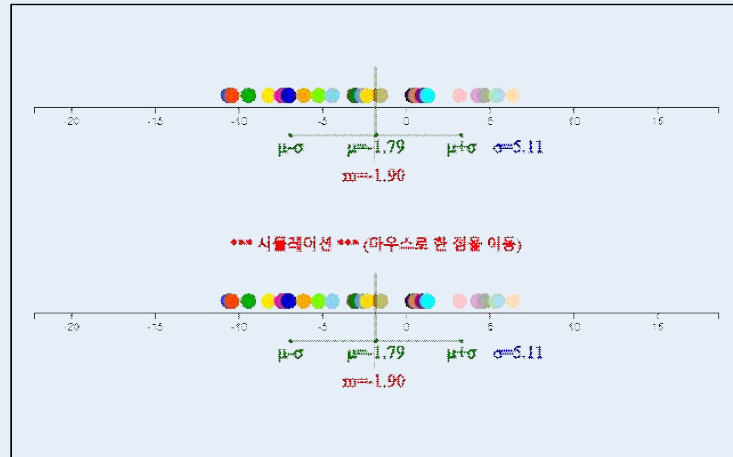
<Figure 4.9> Temperature data input for a dot graph

- When data are entered, the number of data, minimum, maximum, mean, median, variance and standard deviation are calculated. If you click the [Execute] button, a dot graph as shown in <Figure 4.10> appears and the mean, median, standard deviation, and a line of mean \pm standard deviation are displayed.

**Practice 4.4
Solution
(continued)**



- Using the simulation window below the figure, you can check the change in the standard deviation by moving a point with the mouse. The standard deviation is also affected by an extreme point.



<Figure 4.10> Dot graph of daily minimum temperature and a simulation window

Exercise 4.3



The following is data on the length of bicycle-only roads by 25 administrative districts in Seoul as of 2019 ([Exercise 3.1]). Use 『eStatH』 to draw a dot graph and to find and analyze the mean and standard deviation of the data.

(Data 3.3) Length of bicycle-only roads by 25 administrative districts in Seoul as of 2019. (unit km)

24 15 23 20 30 24 7 8 7 12 28 27 19 35 41 42 11 8 37 13
20 29 53 93 42

Exercise 4.4




The following is data on the maximum wind speed of typhoons that passed through Korea in 2020 ([Exercise 3.2]). Use 『eStatH』 to draw a dot graph and find and analyze the mean and standard deviation of data.

(Data 3.4) Maximum wind speed of typhoons that passed through Korea in 2020 (unit m/sec)

40 22 21 29 19 22 24 45 49 55 24 27 29 35 19 24 35 40 56 24
21 43 18

A. Calculation of standard deviation using frequency table

 Think	<p>Assume that the frequency table of the academic achievement test scores of a middle school class is given as follows:</p> <p style="text-align: center;">[Table 4.3] Frequency table of the academic achievement test scores of a middle school</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th style="text-align: center;">Class</th> <th style="text-align: center;">Frequency</th> </tr> </thead> <tbody> <tr> <td style="text-align: center;">$60 \leq \sim < 70$</td> <td style="text-align: center;">2</td> </tr> <tr> <td style="text-align: center;">$70 \sim 80$</td> <td style="text-align: center;">10</td> </tr> <tr> <td style="text-align: center;">$80 \sim 90$</td> <td style="text-align: center;">15</td> </tr> <tr> <td style="text-align: center;">$90 \sim 100$</td> <td style="text-align: center;">3</td> </tr> <tr> <td style="text-align: center;">Total</td> <td style="text-align: center;">30</td> </tr> </tbody> </table>	Class	Frequency	$60 \leq \sim < 70$	2	$70 \sim 80$	10	$80 \sim 90$	15	$90 \sim 100$	3	Total	30
Class	Frequency												
$60 \leq \sim < 70$	2												
$70 \sim 80$	10												
$80 \sim 90$	15												
$90 \sim 100$	3												
Total	30												
Explore	How to find the standard deviation of the data in this frequency table?												

- In the previous section, when a frequency table was given rather than the raw data, the mean was approximated using the middle value of each class interval. The standard deviation is calculated in a similar way.
- First, find the middle value of each class. Then, it is assumed that each class has the middle value as many as the frequency, and the average is obtained using this approximate data.

[Table 4.4] Approximated data using the middle value of each class interval in the academic achievement test scores

Class	Middle value	Frequency	Approximated data
$60 \leq \sim < 70$	65	2	65 65
$70 \sim 80$	75	5	75 75 75 75 75
$80 \sim 90$	85	10	85 85 85 85 85 85 85 85 85 85
$90 \sim 100$	95	3	95 95 95
Total		20	

- Mean is calculated as follows:

$$\begin{aligned}
 \text{Mean} &= \frac{65 + 65 + 75 + 75 + 75 + 75 + 75 + 85 + 85 + 85 + 85 + 85 + 85 + 85 + 85 + 85 + 85 + 95 + 95 + 95}{20} \\
 &= \frac{65 \times 2 + 75 \times 5 + 85 \times 10 + 95 \times 3}{20} \\
 &= \frac{1640}{20} = 82
 \end{aligned}$$

- The variance and standard deviation are calculated in a similar way.

$$\text{Variance} = \frac{(65 - 82)^2 \times 2 + (75 - 82)^2 \times 5 + (85 - 82)^2 \times 10 + (95 - 82)^2 \times 3}{20}$$

$$= \frac{1420}{20} = 71$$

$$\text{Standard deviation} = \sqrt{71} = 8.43$$

- Using 'Frequency Distribution Polygon - Relative Frequency Comparison' of 『eStatH』, the approximate mean and standard deviation of the frequency table can be obtained as shown in <Figure 4.6>. After entering the left value of the class interval and 'Frequency 1', click the [Execute] button.

Category	Frequency 1	Frequency 2	Relative Freq 1	Relative Freq 2
1 60 ≤ ~ < 70.00	2		0.100	
2 70 ≤ ~ < 80.00	5		0.250	
3 80 ≤ ~ < 90.00	10		0.500	
4 90 ≤ ~ < 100.00	3		0.150	
5				
6				
7				
8				
9				
Total	20		1.000	
Mean	82.00			
Std Deviation	8.43			

<Figure 4.11> Standard deviation calculation using a frequency table

4.3 Covariance and Correlation Coefficient

Think	The height and weight of 7 male middle school students were investigated as follows: (Data 4.3) Height and weight of 7 male middle school students <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th></th> <th>1</th> <th>2</th> <th>3</th> <th>4</th> <th>5</th> <th>6</th> <th>7</th> </tr> </thead> <tbody> <tr> <td>Height (cm)</td> <td>162</td> <td>164</td> <td>170</td> <td>158</td> <td>175</td> <td>168</td> <td>172</td> </tr> <tr> <td>Weight (kg)</td> <td>54</td> <td>60</td> <td>64</td> <td>52</td> <td>65</td> <td>60</td> <td>67</td> </tr> </tbody> </table>		1	2	3	4	5	6	7	Height (cm)	162	164	170	158	175	168	172	Weight (kg)	54	60	64	52	65	60	67
	1	2	3	4	5	6	7																		
Height (cm)	162	164	170	158	175	168	172																		
Weight (kg)	54	60	64	52	65	60	67																		
Explore	Is there a measure to determine the correlation between two quantitative variables?																								


- Just as the variance is used as a measure of dispersion in one quantitative variable, the following covariance is used in two quantitative variables. When n number of x and y data are expressed as $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, and the mean is expressed as (μ_x, μ_y) , the **covariance** σ_{xy} can be expressed by the following formula.

$$\text{Covariance} \quad \sigma_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y) \quad (n: \text{number of data})$$

- Covariance implies the total average of the values obtained by multiplying the x-axis distance and the y-axis distance between each data point and the mean point (μ_x, μ_y) of data. Therefore, if there are many points on the upper right and lower left of the mean point, the covariance has a positive value, indicating a positive correlation. If there are many points on the upper left and lower right of the mean point, the covariance has a negative value, indicating a negative correlation. However, since covariance can increase in value depending on the unit of data, the following correlation coefficient is used as a measure of correlation.

$$\text{Correlation coefficient} \quad \rho = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

- The correlation coefficient is a variation of the covariance and can only have values between -1 and +1. When the correlation coefficient is close to +1, the two variables are said to have a strong positive correlation, and when the correlation coefficient is close to -1, it is said to have a strong negative correlation. When the correlation coefficient is close to 0, there is no correlation between the two variables.

Practice 4.5	Using 『eStatH』, calculate the covariance and correlation coefficient of the height and weight of 7 students.
Solution 	<ul style="list-style-type: none"> • If you select 'Scatter Plot – Correlation Coefficient' from the 『eStatH』 menu using the QR on the left, a window like <Figure 4.12> appears. • Enter students' height in 'Enter X data' and their weight in 'Enter Y data'. (You can also copy and paste the material from the e-book)

Practice 4.5
Solution
(continued)



Scatterplot Menu

Main Title: Height and weight of 7 male middle school students

y title: Weight x title: Height

X Enter Data: 162 164 170 158 175 168 172

Y Enter Data: 54 60 64 52 65 60 67

Number of Data	n_x	7	n_y	7		
Mean	μ_x	167.00	μ_y	60.29		
Variance	σ_x^2	30.57	σ_y^2	27.06	Covariance	σ_{xy} 27.00
Std Deviation	σ_x	5.53	σ_y	5.20	Correlation Coefficient	ρ 0.94

Execute

<Figure 4.12> Calculation of covariance and correlation coefficient using height and weight data

- After the data input, click the [Execute] button. Then the number of data, mean, variance, standard deviation, covariance and correlation coefficient are calculated as in <Figure 4.12> and a scatter plot is displayed.
- If you check the 'regression line' under the scatterplot, a regression line that explains the relationship between height and weight is drawn.
- As shown in <Figure 4.12>, the covariance of height and weight in (Data 4.3) is 27 and the correlation coefficient is 0.94, indicating a strong positive correlation.

Exercise 4.5


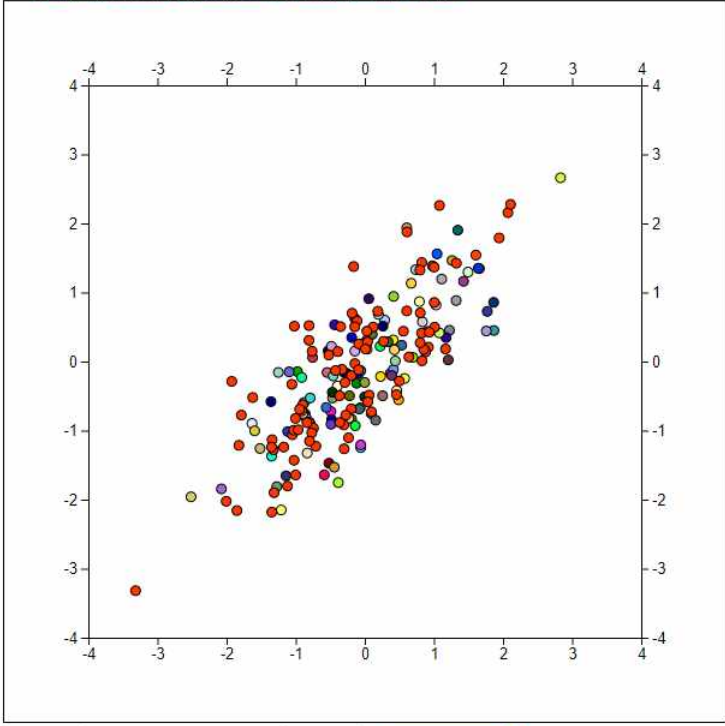


The following data are weekly study hours and test scores of 10 middle school students. Find the covariance and correlation coefficient using 『eStatH』 and analyze what kind of correlation is.

(Data 4.4) Weekly study hours and test score of 10 students

	1	2	3	4	5	6	7	8	9	10
Study hours	10	25	15	16	20	5	18	21	12	20
Test score	75	95	82	85	97	65	87	88	76	90

- By using 『eStatH』, you can examine examples of data for various correlation coefficients.

Practice 4.6	Let's simulate various correlation coefficients using 『eStat H』 .
<p data-bbox="379 398 488 432">Solution</p> 	<ul data-bbox="555 387 1385 495" style="list-style-type: none"> • If you select 'Correlation Coefficient' from the QR on the left or 『eStatH』 menu, an initial scatter plot as shown in <Figure 4.13> appears. <div data-bbox="592 539 1347 1384" style="border: 1px solid black; padding: 5px;"> <p data-bbox="616 555 858 580">Correlation Coefficient Menu</p> <p data-bbox="616 600 1062 622">**Enter Correlation Coefficient and click Execute button</p>  <p data-bbox="608 1352 1326 1375">Execute Correlation Coefficient 0.8 <input type="range" value="0.8"/> 1 <input type="checkbox"/> Regression Line</p> </div> <p data-bbox="643 1391 1299 1417"><Figure 4.13> Simulation window for a correlation coefficient</p> <ul data-bbox="555 1451 1385 1666" style="list-style-type: none"> • Change the 'Correlation Coefficient' under the initial scatter plot to the desired value and click the [Execute] button to display a scatter plot for the corresponding correlation coefficient. • If 'Regression Line' is checked, a regression line representing the points appears.

- If the correlation is strong, a straight line that can explain the relationship between the variables is obtained, which is called a regression line. A detailed explanation of the regression line is covered in university level of statistics.

<p>Practice 4.7</p>	<p>Using 『eStatH』, make a point on the plane and observe the correlation coefficient and regression line while moving.</p>
<p>Solution</p>	<ul style="list-style-type: none"> If you select 'Correlation Coefficient - Regression Line Experiment' from the QR on the left or the 『eStatH』 menu, a screen for testing the correlation coefficient and regression line as shown in <Figure 4.14> appears. <div data-bbox="676 573 1182 1160" data-label="Figure"> </div> <p><Figure 4.14> Simulation window of a correlation coefficient and a regression line</p> <ul style="list-style-type: none"> If you place a dot on this blank screen with the mouse, the regression line and correlation coefficient appear as shown in <Figure 4.15>. You can observe the change of the regression line and the correlation coefficient by moving the point after clicking it with the mouse. <div data-bbox="687 1462 1171 1939" data-label="Figure"> </div> <p><Figure 4.15> Example of dots and its regression line and correlation coefficient</p>

