

CS 1173: Hand computation of statistical indicators

In the following, $\mathbf{A} = [2, 4, 2, 3, 5, 1, 4, 3]$ and $\mathbf{x} = [x_1, x_2, x_3, \dots, x_n]$ is a general set of n measurements used in the formulas.

1. The average or mean measures the central tendency of a set of numbers. To calculate the average, add the values and divide by the number of values. The average of \mathbf{x} is:

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k$$

Example 1: The average or mean of \mathbf{A} is 3:

- Sum the elements of \mathbf{A} : $2 + 4 + 2 + 3 + 5 + 1 + 4 + 3 = 24$.
- Divide by the number of elements in \mathbf{A} : $24/8 = 3$

2. The median is another measure of central tendency. In contrast to the average, the median is not heavily influenced by outliers. When the average and the median are far apart in value, you can surmise that there must be some outliers in the data. To calculate the median, sort the values and take the middle one. If the number of items is even, average the middle two values.

Example 2: The median of \mathbf{A} is 3:

- Sort \mathbf{A} in increasing order: $[1, 2, 2, 3, 3, 4, 4, 5]$.
- Average of the middle two values since \mathbf{A} has an even number of elements: $(3 + 3)/2 = 3$.

3. The mode is the value most frequently appearing in the data. This value may not be unique.

Example 3: \mathbf{A} has three most frequent values: 2, 3 and 4. The mode may not be a single value.

- Compute the frequencies of the elements of \mathbf{A} : 1:1, 2:2, 3:2, 4:2, 5:1
- Find the element(s) with the largest count: 2, 3, and 4 all appear twice.

4. The variance measures the mean-squared error (MSE) of the values from the mean. The variance is the average of the squares of the distance of each value from the mean. The variance \mathbf{x} is:

$$\sigma^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})^2$$

where \bar{x} is the mean of \mathbf{x} .

Example 4: The variance of \mathbf{A} is 1.5:

- First compute the sum of the squares of the differences with the mean:
 $(2 - 3)^2 + (4 - 3)^2 + (2 - 3)^2 + (3 - 3)^2 + (5 - 3)^2 + (1 - 3)^2 + (4 - 3)^2 + (3 - 3)^2 =$
 $1 + 1 + 1 + 0 + 4 + 4 + 1 + 0 = 12$
- Divide by n to obtain the variance: $12/8 = 1.5$

5. The standard deviation measures the spread or dispersion of values from the mean. The standard deviation is calculated as the square root of the average of the squares of the distance of each value from the mean. In other words, the standard deviation is the square root of the variance. The standard deviation of \mathbf{x} is:

$$\sigma = \sqrt{\frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})^2}$$

where \bar{x} is the mean of \mathbf{x} .

Example 5: The standard deviation \mathbf{A} :

- First compute the sum of the squares of the differences with the mean:
 $(2 - 3)^2 + (4 - 3)^2 + (2 - 3)^2 + (3 - 3)^2 + (5 - 3)^2 + (1 - 3)^2 + (4 - 3)^2 + (3 - 3)^2 =$
 $1 + 1 + 1 + 0 + 4 + 4 + 1 + 0 = 12$
- Divide by n and take the square root: $\sqrt{12/8} = 1.225$

6. The p^{th} percentile is the data value such that $p\%$ of the data fall below that value. The 50% percentile is the median. The 25th percentile is sometimes called the first quartile.

Example 6: The 25th percentile of \mathbf{A} is 2, and the 75th percentile of \mathbf{A} is 4.

- Sort \mathbf{A} in increasing order: [1, 2, 2, 3, 3, 4, 4, 5]
- Compute the percentiles for the data: $0.5*(100/8), 1.5*(100/8), \dots, 7.5*(100/8) =$
 $1:6.25, 2:18.75, 3:31.25, 4:43.75, 5:56.25, 6:68.75, 7:81.25, 8:93.75$
- Linearly interpolate the two closest values. For the 25th percentile, we interpolate between 2 at the 18.75th percentile and 3 at the 31.25th percentile, so the answer is 2.

7. The inter-quartile range (IQR) is the 75th percentile minus the 25th percentile.

Example 7: What is the inter-quartile range of \mathbf{A} ?

Since the 25th percentile of \mathbf{A} is 2 and the 75th percentile of \mathbf{A} is 4, the inter-quartile range is 2.

Example 8: What is the position of the fence for a box plot of \mathbf{A} ?

The top of the fence is at the 75th percentile plus $1.5 \times \text{IQR} = 4 + 1.5 \times 2 = 4 + 3 = 7$

The bottom of the fence is at the 25th percentile minus $1.5 \times \text{IQR} = 2 - 3 = -1$

Example 8: What are the positions of the whiskers for a box plot of \mathbf{A} ?

All of the values of \mathbf{A} are inside the fence, so the top whisker is at 5 and the bottom whisker is at 1.