

In this activity, we investigate relationships and differences among the notions of “(overall, or population) mean,” “sample means,” and “mean of the sample means.”

Consider the population

$$X = \{-93, -42, -33, 51, 81\}.$$

Think of the points in X as being measurements on some overall population (of size 5) that we’re interested in studying.

1. Compute the mean, let’s call it μ , and the standard deviation, call it σ , of this population. Use the formulas for mean and standard deviation that we studied in class. Write your answers to two decimal places. (Note: we’re calling these things μ and σ here, to emphasize that we’re thinking of X as a population. But to compute then, you should use the formulas for \bar{x} and s .)

$$\mu = \underline{-7.2000} \qquad \sigma = \underline{71.4227}$$

2. Write down all three-element subsets of X . Hint: there are $\binom{5}{3} = 10$ of these subsets; call them S_1, S_2, \dots, S_{10} . The first two have been done for you, to get you started. (You can write down the remaining ones in any order.)

Remark. Think of each of the subsets S_1, S_2, \dots, S_{10} as being a size-3 *sample* from the overall population X .

$$S_1 = \underline{\{-93, -42, -33\}} \qquad S_2 = \underline{\{-93, -42, 51\}}$$

$$S_3 = \underline{\{-93, -42, 81\}} \qquad S_4 = \underline{\{-93, -33, 51\}}$$

$$S_5 = \underline{\{-93, -33, 81\}} \qquad S_6 = \underline{\{-93, 51, 81\}}$$

$$S_7 = \underline{\{-42, -33, 51\}} \qquad S_8 = \underline{\{-42, -33, 81\}}$$

$$S_9 = \underline{\{-42, 51, 81\}} \qquad S_{10} = \underline{\{-33, 51, 81\}}$$

3. Each set S_k has a mean \bar{x}_k : that is, S_1 has mean \bar{x}_1 , S_2 has mean \bar{x}_2 , and so on. Each \bar{x}_k is called a (size-3) *sample mean* from the overall population X .

Compute the \bar{x}_k 's and write your answers in the spaces below (the first two have been done for you). You can just write the final answer (to two decimal places) in each case.

$$\bar{x}_1 = \frac{-93 - 42 - 33}{3} = -56$$

$$\bar{x}_2 = \frac{-93 - 42 + 51}{3} = -28$$

$$\bar{x}_3 = \frac{-93 - 42 + 81}{3} = -18$$

$$\bar{x}_4 = \frac{-93 - 33 + 51}{3} = -25$$

$$\bar{x}_5 = \frac{-93 - 33 + 81}{3} = -15$$

$$\bar{x}_6 = \frac{-93 + 51 + 81}{3} = 13$$

$$\bar{x}_7 = \frac{-42 - 33 + 51}{3} = -8$$

$$\bar{x}_8 = \frac{-42 - 33 + 81}{3} = 2$$

$$\bar{x}_9 = \frac{-42 + 51 + 81}{3} = 30$$

$$\bar{x}_{10} = \frac{-33 + 51 + 81}{3} = 33$$

4. Let \bar{X} denote the population of all of the above sample means; that is,

$$\bar{X} = \{\bar{x}_1, \bar{x}_2, \bar{x}_3, \bar{x}_4, \bar{x}_5, \bar{x}_6, \bar{x}_7, \bar{x}_8, \bar{x}_9, \bar{x}_{10}\}.$$

Compute the mean, let's call it $\bar{\mu}$, and the standard deviation, call it $\bar{\sigma}$, of \bar{X} . Again, use the formulas for mean and standard deviation that we studied in class. Write your answers to two decimal places.

$$\bar{\mu} = -7.2000$$

$$\bar{\sigma} = 27.4906$$

5. How do μ and $\bar{\mu}$ compare? That is: which is smaller, or are they equal?

$$\mu = \bar{\mu}.$$

6. How do σ and $\bar{\sigma}$ compare? That is: which is smaller, or are they equal?

$$\bar{\sigma} < \sigma.$$

7. Fill in the blanks. Use your answers above to guide you. Each blank should be filled with one of the words/phrases “smaller than,” “equal to,” “average,” “extreme,” or “spread.”

Consider a set X of data points corresponding to some overall population; suppose X has mean μ and standard deviation s .

Fix a sample size n , and compute the mean of each size- n sample from X . Let \bar{X} denote the set of all of these sample means. Then:

- The mean $\bar{\mu}$ of \bar{X} is equal to the mean μ of X . This reflects the intuitively plausible fact that “the average of the averages equals the average.”
- The standard deviation $\bar{\sigma}$ of \bar{X} is smaller than the standard deviation σ of X . Why should this be true? It’s because the process of taking averages tends to mitigate the effect of outliers, or extreme values, in your data. That is: an extreme data value, meaning one that’s far away from the mean, can result in a relatively large standard deviation, or spread, in your data. But if this extreme value is averaged against other data values, then its impact on the spread in the data will not be as extreme, since the other data values involved in the computation will tend to pull things back towards the average. Consequently, averaging data will tend to yield numbers with a smaller spread, and consequently a smaller standard deviation, than the original data itself. Or, to summarize (and repeat): the standard deviation $\bar{\sigma}$ of the set \bar{X} of size- n sample means from a population X is smaller than the standard deviation σ of the data set X itself.

8. Fill in the blank: According to the SDM theorem from the class of 11/13 (this past Wednesday), under certain conditions, the standard deviation $\bar{\sigma}$ of a population \bar{X} of size- n sample means should be related to the the standard deviation σ of the original population X by the formula $\bar{\sigma} = \sigma / \underline{\sqrt{n}}$.

Question: does this formula hold for the standard deviations σ and $\bar{\sigma}$ that you actually computed above? If not, why not? Hint: are any conditions of SDM not being met in the present case?

The condition not being met is $n \geq 30$. The approximations described in SDM are better and better the larger n is. For small n , they may not be very good approximations at all.