

Sample Variance (3 pages; 8/7/21)

(1) Strictly speaking, the variance of a sample is defined as the average squared deviation from the sample mean;

$$\text{ie } s^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$$

(and then the standard deviation is $s = \sqrt{\frac{1}{n} \sum (x_i - \bar{x})^2}$).

However, if the sample variance is intended to be an estimate for the population variance, then it can be shown that an unbiased estimate of the population variance is

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

This means that, if we define S^2 to be the random variable

$$\frac{1}{n-1} \sum (X_i - \bar{X})^2, \text{ then } E(S^2) = \sigma^2, \text{ the population variance.}$$

[See the Appendix for a proof of this. The fact that we are using \bar{x} (the sample mean) in the formula, instead of the population mean μ , means that the n deviations $x_i - \bar{x}$ are not independent (for example, $x_n - \bar{x}$ can be determined, if the other deviations are known). This suggests that an average of n deviations isn't appropriate, but doesn't constitute a proof that a divisor of $n - 1$ gives the right value.]

For exam purposes, $s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$ is usually preferred, even when there is no mention of it being an estimate for the population variance.

(2) Alternative, and generally more convenient formulae are:

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

and $s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2 = \frac{1}{n-1} \{(\sum x_i^2) - n\bar{x}^2\}$ when the divisor of $n - 1$ is being used

Proof (for $s^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$)

$$\begin{aligned} s^2 &= \frac{1}{n} \sum (x_i - \bar{x})^2 \\ &= \frac{1}{n} \{ \sum (x_i^2 - 2x_i\bar{x} + \bar{x}^2) \} \\ &= \frac{1}{n} \{ (\sum x_i^2) - 2\bar{x}(\sum x_i) + n\bar{x}^2 \} \\ &= \frac{1}{n} \{ (\sum x_i^2) - 2\bar{x}(n\bar{x}) + n\bar{x}^2 \} \\ &= \frac{1}{n} \{ (\sum x_i^2) - n\bar{x}^2 \} \end{aligned}$$

A useful check is as follows:

If all the n data items are the same, then each $x_i = \bar{x}$,

and $\sum x_i^2 = n\bar{x}^2$, so that $s^2 = 0$; as expected, since there is no variance amongst the x_i .

Notes

(i) $(\sum x_i^2) - n\bar{x}^2$ is often denoted S_{xx}

(ii) $(\sum x_i^2) - n\bar{x}^2$ can also be written as $(\sum x_i^2) - \frac{(\sum x_i)^2}{n}$

(iii) It is tempting to write $\frac{1}{n}\{(\sum x_i^2) - n\bar{x}^2\}$ as $\frac{1}{n}(\sum x_i^2) - \bar{x}^2$,

but $\frac{1}{n}\{(\sum x_i^2) - n\bar{x}^2\}$ has the advantage that it can easily be converted to $\frac{1}{n-1}\{(\sum x_i^2) - n\bar{x}^2\}$ if necessary.

Appendix

$S^2 = \frac{1}{n-1}([\sum X^2] - n\bar{X}^2)$ is an unbiased estimator for the population variance

Proof

$$E(S^2) = \frac{1}{n-1}\{[\sum E(X^2)] - nE(\bar{X}^2)\}$$

$$= \frac{1}{n-1}\{nE(X^2) - nE(\bar{X}^2)\}$$

$$= \frac{n}{n-1}\{E(X^2) - E(\bar{X}^2)\}$$

Now $\sigma^2 = E(X^2) - \mu^2$, so that $E(X^2) = \sigma^2 + \mu^2$

$$\text{Also, } \text{Var}(\bar{X}) = E(\bar{X}^2) - \mu^2,$$

$$\text{and } \text{Var}(\bar{X}) = \text{Var}\left(\frac{X_1 + \dots + X_n}{n}\right) = \frac{1}{n^2}\text{Var}(X_1 + \dots + X_n)$$

$$= \frac{1}{n^2}(n\text{Var}(X_i)) = \frac{\sigma^2}{n},$$

$$\text{so that } E(\bar{X}^2) = \frac{\sigma^2}{n} + \mu^2$$

$$\text{Then } E(S^2) = \frac{n}{n-1}([\sigma^2 + \mu^2] - \left[\frac{\sigma^2}{n} + \mu^2\right])$$

$$= \frac{n\sigma^2}{n-1}\left(1 - \frac{1}{n}\right)$$

$$= \frac{n\sigma^2}{n-1}\left(\frac{n-1}{n}\right) = \sigma^2$$