
Why Sample Variance has $n - 1$ on denominator?

Recall the population variance formula:

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (X_i - \mu)^2$$

Sample variance:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

Rewrite $\mathbb{E}[\sum_{i=1}^n (X_i - \bar{X})^2]$ as:

$$\mathbb{E}\left[\sum_{i=1}^n (X_i - \mu + \mu - \bar{X})^2\right]$$

Expand the squared term:

$$\mathbb{E}\left[\sum_{i=1}^n ((X_i - \mu)^2 + 2(X_i - \mu)(\mu - \bar{X}) + (\mu - \bar{X})^2)\right]$$

Simplify:

$$\mathbb{E}\left[\sum_{i=1}^n (X_i - \mu)^2\right] + \mathbb{E}\left[\sum_{i=1}^n 2(X_i - \mu)(\mu - \bar{X})\right] + \mathbb{E}\left[\sum_{i=1}^n (\mu - \bar{X})^2\right]$$

The middle term is

$$\begin{aligned} \mathbb{E}\left[\sum_{i=1}^n 2(X_i - \mu)(\mu - \bar{X})\right] &= 2\mathbb{E}\left[\sum_{i=1}^n (X_i - \mu)(\mu - \bar{X})\right] \\ &= 2\mathbb{E}[(n\bar{X} - n\mu)(\mu - \bar{X})] \\ &= -2n\mathbb{E}[(\mu - \bar{X})^2] \end{aligned}$$

For the third item

$$\mathbb{E}\left[\sum_{i=1}^n (\mu - \bar{X})^2\right] = \mathbb{E}[n(\mu - \bar{X})^2] = n\mathbb{E}[(\mu - \bar{X})^2]$$

So we have

$$\mathbb{E}\left[\sum_{i=1}^n (X_i - \bar{X})^2\right] = \mathbb{E}\left[\sum_{i=1}^n (X_i - \mu)^2\right] - n\mathbb{E}[(\mu - \bar{X})^2]$$

We know that $\mathbb{E}[\sum_{i=1}^n (X_i - \mu)^2] = n\sigma^2$

$$\mathbb{E}[(\mu - \bar{X})^2] = \mathbb{E}[\mu^2 - 2\mu\bar{X} + \bar{X}^2] = \mathbb{E}[\bar{X}^2] - \mu^2$$

And

$$\begin{aligned}\text{Var}(X) &= \mathbb{E}[(X - \mathbb{E}[X])^2] \\ &= \mathbb{E}[X^2 - 2X\mathbb{E}[X] + \mathbb{E}[X]^2] \\ &= \mathbb{E}[X^2] - 2\mathbb{E}[X]\mathbb{E}[X] + \mathbb{E}[X]^2 \\ &= \mathbb{E}[X^2] - 2\mathbb{E}[X]^2 + \mathbb{E}[X]^2 \\ &= \mathbb{E}[X^2] - \mathbb{E}[X]^2\end{aligned}$$

So

$$\mathbb{E}[\bar{X}^2] = \text{Var}[\bar{X}] + \mathbb{E}[\bar{X}]^2$$

So

$$\begin{aligned}\mathbb{E}[(\mu - \bar{X})^2] &= \mathbb{E}[\bar{X}^2] - \mu^2 \\ &= \text{Var}[\bar{X}] \\ &= \text{Var}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] \\ &= \frac{1}{n^2} \text{Var}\left[\sum_{i=1}^n X_i\right] \\ &= \frac{1}{n^2} \times n \times \text{Var}[X_i] \\ &= \frac{\sigma^2}{n}\end{aligned}$$

So we have

$$\mathbb{E}\left[\sum_{i=1}^n (X_i - \bar{X})^2\right] = n\sigma^2 - n\left(\frac{\sigma^2}{n}\right) = n\sigma^2 - \sigma^2 = (n-1)\sigma^2$$

If we use n in the denominator:

$$\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right] = \frac{n-1}{n} \sigma^2$$

This is biased because it's less than the true population variance σ^2 .

But if we use $(n-1)$ in the denominator:

$$\mathbb{E}\left[\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2\right] = \sigma^2$$

This gives us an unbiased estimator of the population variance.