# Proof $\frac{\sum(X-\bar{X})^2}{n}$ Is Biased Estimator Of $\sigma^2$

Preliminaries:

$\mu$ is the population mean.

$\sigma^2$ is the population variance.

$\bar{X}$ is the sample mean and is a random variable.

It is a given that $\mathbb{E}(X) = \mu$ and $\text{Var}(X) = \sigma^2$.

We will make extensive use of the results $\mathbb{E}(aX + bY) = a\mathbb{E}(X) + b\mathbb{E}(Y)$ and $\text{Var}(aX + bY) = a^2\text{Var}(X) + b^2\text{Var}(Y)$. Note that the *second* result requires *independence* between $X$ and $Y$. We need to prove a few results before we start.

$$
\begin{aligned}
\mathbb{E}(\bar{X}) &= \mathbb{E}\left(\frac{X_1 + X_2 + X_3 \cdots + X_n}{n}\right) \\
&= \mathbb{E}\left(\frac{X_1}{n}\right) + \mathbb{E}\left(\frac{X_2}{n}\right) + \mathbb{E}\left(\frac{X_3}{n}\right) + \cdots + \mathbb{E}\left(\frac{X_n}{n}\right) \\
&= \frac{\mathbb{E}(X_1)}{n} + \frac{\mathbb{E}(X_2)}{n} + \frac{\mathbb{E}(X_3)}{n} + \cdots + \frac{\mathbb{E}(X_n)}{n} \\
&= \frac{\mu}{n} + \frac{\mu}{n} + \frac{\mu}{n} + \cdots + \frac{\mu}{n} \\
&= n \times \frac{\mu}{n} \\
&= \mu.
\end{aligned}
$$

In other words $\bar{X}$ is an unbiased estimator of $\mu$ which is what our intuition would say; the mean of the sample is our best guess of what the population mean is.

$$
\begin{aligned}
\text{Var}(\bar{X}) &= \text{Var}\left(\frac{X_1 + X_2 + \cdots + X_n}{n}\right) \\
&= \text{Var}\left(\frac{X_1}{n}\right) + \text{Var}\left(\frac{X_2}{n}\right) + \text{Var}\left(\frac{X_3}{n}\right) + \cdots + \text{Var}\left(\frac{X_n}{n}\right) \\
&= \frac{\text{Var}(X_1)}{n^2} + \frac{\text{Var}(X_2)}{n^2} + \frac{\text{Var}(X_3)}{n^2} + \cdots + \frac{\text{Var}(X_n)}{n^2} \\
&= \frac{\sigma^2}{n^2} + \frac{\sigma^2}{n^2} + \frac{\sigma^2}{n^2} + \cdots + \frac{\sigma^2}{n^2} \\
&= n \times \frac{\sigma^2}{n^2} \\
&= \frac{\sigma^2}{n}.
\end{aligned}
$$

This is much less obvious, but shows that the *larger* the sample size the *smaller* the variance in the sample mean so we can be rather more confident that $\bar{X}$ is close to $\mu$.

Also since $\text{Var}(X) \equiv \mathbb{E}(X^2) - (\mathbb{E}(X))^2$ we have $\mathbb{E}(X^2) = \text{Var}(X) + (\mathbb{E}(X))^2 = \sigma^2 + \mu^2$.

Similarly $\text{Var}(\bar{X}) \equiv \mathbb{E}(\bar{X}^2) - (\mathbb{E}(\bar{X}))^2$, so $\mathbb{E}(\bar{X}^2) = \text{Var}(\bar{X}) + (\mathbb{E}(\bar{X}))^2 = \frac{\sigma^2}{n} + \mu^2$.

We want to show that if we take a sample from a population then calculating the variance of the sample using $\frac{\sum(X-\bar{X})^2}{n}$ (which is equivalent to $\frac{\sum X^2}{n} - \bar{X}^2$) does not (on average) give $\sigma^2$; i.e. it is a *biased* estimator.

$$\mathbb{E}\left(\frac{\sum(X - \bar{X})^2}{n}\right) = \mathbb{E}\left(\frac{\sum X^2}{n} - \bar{X}^2\right)$$

$$= \mathbb{E}\left(\frac{X_1^2}{n} + \frac{X_2^2}{n} + \frac{X_3^2}{n} + \cdots + \frac{X_n^2}{n} - \bar{X}^2\right)$$

$$= \mathbb{E}\left(\frac{X_1^2}{n}\right) + \mathbb{E}\left(\frac{X_2^2}{n}\right) + \mathbb{E}\left(\frac{X_3^2}{n}\right) + \cdots + \mathbb{E}\left(\frac{X_n^2}{n}\right) - \mathbb{E}\left(\bar{X}^2\right)$$

$$= \frac{\mathbb{E}(X_1^2)}{n} + \frac{\mathbb{E}(X_2^2)}{n} + \frac{\mathbb{E}(X_3^2)}{n} + \cdots + \frac{\mathbb{E}(X_n^2)}{n} - \mathbb{E}(\bar{X}^2)$$

$$= \frac{\sigma^2 + \mu^2}{n} + \frac{\sigma^2 + \mu^2}{n} + \frac{\sigma^2 + \mu^2}{n} + \cdots + \frac{\sigma^2 + \mu^2}{n} - \left(\frac{\sigma^2}{n} + \mu^2\right)$$

$$= \sigma^2 + \mu^2 - \frac{\sigma^2}{n} - \mu^2$$

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

So we don't (on average) get the desired $\sigma^2$. Which is why when estimating $\sigma^2$ we multiply by

$$\frac{n}{n-1}.$$