

The Maximum Likelihood Estimation For The Variance Is Biased

BIO210 Biostatistics

Extra Reading Material for Lecture 18

Xi Chen

School of Life Sciences

Southern University of Science and Technology

Fall 2024

The random variable \mathbf{X} denotes certain metric (*e.g.* height, weight) we are interested in from a population, and $\mathbf{X} \sim \mathcal{N}(\mu, \sigma^2)$. We draw a random sample of size n from the population. Like we discussed during the lecture, a random sample of size n can be thought as n **i.i.d.** random variables. That is:

$$\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3, \dots, \mathbf{X}_n \sim \mathcal{N}(\mu, \sigma^2)$$

We have seen that the maximum likelihood estimator for σ^2 is:

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})^2$$

Then, what is $\mathbb{E}[\hat{\sigma}^2]$? If $\mathbb{E}[\hat{\sigma}^2] = \sigma^2$, it is an unbiased estimator. Otherwise, it is a biased one.

Now let's have a look.

$$\begin{aligned} \mathbb{E}[\hat{\sigma}^2] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})^2\right] = \frac{1}{n} \mathbb{E}\left[\sum_{i=1}^n (\mathbf{X}_i^2 - 2\bar{\mathbf{X}}\mathbf{X}_i + \bar{\mathbf{X}}^2)\right] \\ &= \frac{1}{n} \mathbb{E}\left[\sum_{i=1}^n \mathbf{X}_i^2 - 2\bar{\mathbf{X}} \sum_{i=1}^n \mathbf{X}_i + \sum_{i=1}^n \bar{\mathbf{X}}^2\right] \end{aligned}$$

Note that: $\sum_{i=1}^n \mathbf{X}_i = n\bar{\mathbf{X}}$. Since $\bar{\mathbf{X}}$ remains the same for each i , we have $\sum_{i=1}^n \bar{\mathbf{X}}^2 = n\bar{\mathbf{X}}^2$. Replacing the blue terms above, we have:

$$\begin{aligned} \mathbb{E}[\hat{\sigma}^2] &= \frac{1}{n} \mathbb{E}\left[\sum_{i=1}^n \mathbf{X}_i^2 - 2\bar{\mathbf{X}} \cdot n\bar{\mathbf{X}} + n\bar{\mathbf{X}}^2\right] = \frac{1}{n} \mathbb{E}\left[\sum_{i=1}^n \mathbf{X}_i^2 - n\bar{\mathbf{X}}^2\right] \\ &= \frac{1}{n} \left(\mathbb{E}\left[\sum_{i=1}^n \mathbf{X}_i^2\right] - \mathbb{E}[n\bar{\mathbf{X}}^2] \right) \end{aligned} \tag{1}$$

Since $\text{Var}(\mathbf{X}) = \mathbb{E}[\mathbf{X}^2] - (\mathbb{E}[\mathbf{X}])^2$, so we have $\mathbb{E}[\mathbf{X}^2] = \text{Var}(\mathbf{X}) + (\mathbb{E}[\mathbf{X}])^2$,

then,

$$\begin{aligned}
 \mathbb{E} \left[\sum_{i=1}^n \mathbf{X}_i^2 \right] &= \mathbb{E} [\mathbf{X}_1^2] + \mathbb{E} [\mathbf{X}_2^2] + \mathbb{E} [\mathbf{X}_3^2] + \cdots + \mathbb{E} [\mathbf{X}_n^2] \\
 &= \text{Var}(\mathbf{X}_1) + (\mathbb{E}[\mathbf{X}_1])^2 + \text{Var}(\mathbf{X}_2) + (\mathbb{E}[\mathbf{X}_2])^2 + \cdots \\
 &\quad + \text{Var}(\mathbf{X}_n) + (\mathbb{E}[\mathbf{X}_n])^2 \\
 &= \sigma^2 + \mu^2 + \sigma^2 + \mu^2 + \cdots + \sigma^2 + \mu^2 \\
 &= n\sigma^2 + n\mu^2
 \end{aligned} \tag{2}$$

Putting equation (2) into equation (1), we have:

$$\begin{aligned}
 \mathbb{E} [\hat{\sigma}^2] &= \sigma^2 + \mu^2 - \frac{1}{n} \cdot \mathbb{E} [n\bar{\mathbf{X}}^2] = \sigma^2 + \mu^2 - \mathbb{E} [\bar{\mathbf{X}}^2] \\
 &= \sigma^2 + \mu^2 - (\sigma_{\bar{\mathbf{X}}}^2 + \mu_{\bar{\mathbf{X}}}^2)
 \end{aligned} \tag{3}$$

According to the central limit theorem, we have $\mu_{\bar{\mathbf{X}}} = \mu$ and $\sigma_{\bar{\mathbf{X}}}^2 = \frac{\sigma^2}{n}$. Therefore, equation (3) becomes:

$$\mathbb{E} [\hat{\sigma}^2] = \sigma^2 + \mu^2 - \frac{\sigma^2}{n} - \mu^2 = \frac{n-1}{n} \sigma^2 \neq \sigma^2$$

Hence, it is not an unbiased estimator.