# Artificial Intelligence and Machine Learning

Derivation of the Empirical Variance Formula

Let $N$ independent and identically distributed random variables $X_1, X_2, \ldots, X_N$ be given. We assume they have mean $\mathbb{E}\{X_n\} := \mu$ and variance $\mathbb{V}\{X_n\} := \sigma^2$ ($1 \le n \le N$). Our goal is to find an **unbiased estimator** for the variance parameter. *(The estimator $\mu^{ML} := \frac{1}{N}\sum_{n=1}^{N} X_n$ for the mean – which we have derived in the lecture notes – is an unbiased estimator.)*

First, we show that the maximum likelihood estimator for the variance

$$(\sigma^2)^{\mathrm{ML}} := \frac{1}{N}\sum_{n=1}^{N}(X_n - \mu^{\mathrm{ML}})^2$$

is biased. For this we determine the expected value of $(\sigma^2)^{\mathrm{ML}}$. We start by computing:

$$\mathbb{E}\left\{\sum_{n=1}^{N}(X_n - \mu^{\mathrm{ML}})^2\right\} = \mathbb{E}\left\{\sum_{n=1}^{N}\left(X_n^2 - 2X_n\mu^{\mathrm{ML}} + (\mu^{\mathrm{ML}})^2\right)\right\}$$

[**Pull sum inside**]

$$= \mathbb{E}\left\{\sum_{n=1}^{N} X_n^2 - 2\mu^{\mathrm{ML}}\sum_{n=1}^{N} X_n + N(\mu^{\mathrm{ML}})^2\right\}$$

[**Plug in the definition of $\mu^{\mathbf{ML}}$**]

$$= \mathbb{E}\left\{\sum_{n=1}^{N} X_n^2 - \frac{2}{N}\sum_{n=1}^{N} X_n \sum_{n=1}^{N} X_n + N\left(\frac{1}{N}\sum_{n=1}^{N} X_n\right)^2\right\}$$

$$= \mathbb{E}\left\{\sum_{n=1}^{N} X_n^2 - \frac{2}{N}\left(\sum_{n=1}^{N} X_n\right)^2 + \frac{1}{N}\left(\sum_{n=1}^{N} X_n\right)^2\right\}$$

$$= \mathbb{E}\left\{\sum_{n=1}^{N} X_n^2 - \frac{1}{N}\left(\sum_{n=1}^{N} X_n\right)^2\right\}$$

[**Make use of the linearity of $\mathbb{E}$**]

$$= \sum_{n=1}^{N}\mathbb{E}\left\{X_n^2\right\} - \frac{1}{N}\mathbb{E}\left\{\left(\sum_{n=1}^{N} X_n\right)^2\right\}$$

[**Plug in definitions: For any random variable $Y$ we have $\mathbb{V}\{Y\} := \mathbb{E}\left\{Y^2\right\} - \mathbb{E}\{Y\}^2$. Moreover, by definition of the random variables $X_n$ ($1 \le n \le N$) we have that $\mathbb{E}\{X_n\} := \mu$ and $\mathbb{V}\{X_n\} := \sigma^2$**]

$$= \sum_{n=1}^{N}(\mathbb{V}\{X_n\} + \mu^2) - \frac{1}{N}\left(\mathbb{V}\left\{\sum_{n=1}^{N} X_n\right\} + (N\mu)^2\right)$$

$$= N(\sigma^2 + \mu^2) - \frac{1}{N}(N\sigma^2 + N^2\mu^2)$$

$$= N\sigma^2 + N\mu^2 - \sigma^2 - N\mu^2$$

$$= (N - 1)\sigma^2 \tag{1}$$

Using the result we obtained in (1) we are now able to show that the maximum likelihood estimator for the variance is biased:

$$\mathbb{E}\left\{(\sigma^2)^{\mathrm{ML}}\right\} = \mathbb{E}\left\{\frac{1}{N}\sum_{n=1}^{N}(X_n - \mu^{\mathrm{ML}})^2\right\}$$

**[Linearity of $\mathbb{E}$]**

$$= \frac{1}{N}\mathbb{E}\left\{\sum_{n=1}^{N}(X_n - \mu^{\mathrm{ML}})^2\right\}$$

$$\stackrel{(1)}{=} \frac{N-1}{N}\sigma^2$$

Since $\frac{N-1}{N} < 1$, we see that $(\sigma^2)^{\mathrm{ML}}$ **systematically underestimates** the true variance of the data. We can correct for this bias by defining the **empirical variance** according to:

$$(\sigma^2)^{\mathrm{Emp}} := \frac{N}{N-1}(\sigma^2)^{\mathrm{ML}}$$

$$= \frac{N}{N-1}\left(\frac{1}{N}\sum_{n=1}^{N}(X_n - \mu^{\mathrm{ML}})^2\right)$$

$$= \boxed{\frac{1}{N-1}\sum_{n=1}^{N}(X_n - \mu^{\mathrm{ML}})^2} \tag{2}$$

Finally, let us verify that the empirical variance is indeed unbiased:

$$\mathbb{E}\left\{(\sigma^2)^{\mathrm{Emp}}\right\} = \mathbb{E}\left\{\frac{1}{N-1}\sum_{n=1}^{N}(X_n - \mu^{\mathrm{ML}})^2\right\}$$

**[Linearity of $\mathbb{E}$]**

$$= \frac{1}{N-1}\mathbb{E}\left\{\sum_{n=1}^{N}(X_n - \mu^{\mathrm{ML}})^2\right\}$$

$$\stackrel{(1)}{=} \frac{1}{N-1}(N-1)\sigma^2$$

$$= \sigma^2$$

$\square$

---