

# Notes On Estimation

Alon Jacobson

August 5, 2021

## 1 Notation

Say we have a fixed, unknown quantity of interest  $\theta$  which is some statistic or functional of an underlying statistical model. We gather scalar or vector-valued data which is generated from this statistical model, and this data that we obtain is described by a random variable  $X$ . From a realization of  $X$ , we construct a point estimate or “best guess” of  $\theta$  by  $\hat{\theta} = g(X)$  which is some fixed function  $g$  of  $X$ . Thus  $\hat{\theta}$  is a random variable.

The bias of an estimator is defined by

$$\text{bias}(\hat{\theta}) = \mathbb{E}(\hat{\theta}) - \theta = \mathbb{E}(\hat{\theta} - \theta).$$

We say that  $\hat{\theta}$  is *unbiased* when  $\text{bias}(\hat{\theta}) = 0$ , or equivalently when  $\mathbb{E}(\hat{\theta}) = \theta$ . We use  $\mathbb{E}(\hat{\theta})$  to mean the expected value over the distribution of data  $X$  (i.e., averaging over all possible observations of  $X$ ):  $\mathbb{E}(r(X)) = \int r(x)f_X(x) dx$ . Similarly, we write  $\mathbb{V}$  for the variance over the distribution of  $X$ .

A way to assess the quality of a point estimate is by the mean squared error, or MSE, defined by:

$$\text{MSE}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \theta)^2].$$

The MSE can be written as

$$\text{MSE}(\hat{\theta}) = \text{bias}^2(\hat{\theta}) + \mathbb{V}(\hat{\theta}).$$

*Proof.*

$$\begin{aligned}\text{MSE}(\hat{\theta}) &= \mathbb{E}[(\hat{\theta} - \theta)^2] \\ &= (\mathbb{E}[\hat{\theta} - \theta])^2 + \mathbb{V}(\hat{\theta} - \theta) \\ &= \text{bias}^2(\hat{\theta}) + \mathbb{V}(\hat{\theta}).\end{aligned}$$

where we have used the fact that, since for a random variable  $X$ ,  $\mathbb{V}(X) = \mathbb{E}(X^2) - (\mathbb{E}X)^2$ , we have  $\mathbb{E}(X^2) = (\mathbb{E}X)^2 + \mathbb{V}(X)$ .  $\square$

## 2 Bias Correction

Suppose you have a point estimate  $\hat{\theta}$  of a quantity  $\theta$  and you want to correct for its bias to make it unbiased. If you know the bias, you could make a new estimator  $\hat{\theta}_1$ :

$$\hat{\theta}_1 := \hat{\theta} - \text{bias}(\hat{\theta}).$$

This estimator is unbiased:

$$\begin{aligned}
\text{bias}(\hat{\theta}_1) &= \mathbb{E}(\hat{\theta}_1 - \theta) \\
&= \mathbb{E}(\hat{\theta} - \text{bias}(\hat{\theta}) - \theta) \\
&= \mathbb{E}(\hat{\theta} - \theta) - \mathbb{E}(\text{bias}(\hat{\theta})) \\
&= \text{bias}(\hat{\theta}) - \mathbb{E}(\text{bias}(\hat{\theta})) \\
&= \text{bias}(\hat{\theta}) - \text{bias}(\hat{\theta}) = 0.
\end{aligned}$$

Now, suppose instead that you know that  $\text{bias}(\hat{\theta}) = K\theta$  for some constant  $K$ . (Or, alternatively, you only estimate  $K$  with  $\hat{K}$ , so you have  $\text{bias}(\hat{\theta}) = K\theta \approx \hat{K}\theta$ . But this will not affect the calculations.) But of course you don't know  $\theta$  itself, so you don't know  $\text{bias}(\hat{\theta}) = K\theta$ . So you plug in  $\hat{\theta}$  for  $\theta$  to estimate  $\text{bias}(\hat{\theta})$ :

$$\widehat{\text{bias}}(\hat{\theta}) := K\hat{\theta}.$$

So you construct the estimator

$$\begin{aligned}
\hat{\theta}_2 &:= \hat{\theta} - \widehat{\text{bias}}(\hat{\theta}) \\
&= \hat{\theta} - K\hat{\theta} = (1 - K)\hat{\theta}.
\end{aligned}$$

Is this estimator unbiased? Let's check.

$\hat{\theta}_2$  is an estimator of the form  $\alpha\hat{\theta}$  for some constant  $\alpha$ ; in this case  $\alpha = 1 - K$ . Let's first calculate the bias of an estimator of this more general form (take note - this formula will be used a lot):

$$\begin{aligned}
\text{bias}(\alpha\hat{\theta}) &= \mathbb{E}(\alpha\hat{\theta}) - \theta \\
&= \alpha\mathbb{E}(\hat{\theta}) - \alpha\theta - (1 - \alpha)\theta \\
&= \alpha(\mathbb{E}(\hat{\theta}) - \theta) - (1 - \alpha)\theta \\
&= \alpha\text{bias}(\hat{\theta}) - (1 - \alpha)\theta.
\end{aligned}$$

Now we can compute  $\text{bias}(\hat{\theta}_2)$ :

$$\begin{aligned}
\text{bias}(\hat{\theta}_2) &= \text{bias}((1 - K)\hat{\theta}) \\
&= (1 - K)\text{bias}(\hat{\theta}) - (1 - (1 - K))\theta \\
&= \text{bias}(\hat{\theta}) - K\text{bias}(\hat{\theta}) - K\theta \\
&= \text{bias}(\hat{\theta}) - K\text{bias}(\hat{\theta}) - \text{bias}(\hat{\theta}) \\
&= -K\text{bias}(\hat{\theta}) = -K^2\theta.
\end{aligned}$$

So, assuming  $K \neq 0$  and  $\theta \neq 0$ ,  $\hat{\theta}_2$  is biased!

OK, let's take a different approach. We seek an estimator  $\hat{\theta}_3$  that is unbiased. Let's directly solve for  $\hat{\theta}_3$ . No matter the value of  $\hat{\theta}_3$ , as long as  $\hat{\theta} \neq 0$ , we have

$$\hat{\theta}_3 = \alpha\hat{\theta}$$

for some  $\alpha \in \mathbb{R}$ . Let's solve for  $\alpha$  by setting  $\text{bias}(\hat{\theta}_3) = 0$ .

$$\begin{aligned}
\text{bias}(\hat{\theta}_3) &= \text{bias}(\alpha\hat{\theta}) = \alpha\text{bias}(\hat{\theta}) - (1 - \alpha)\theta \\
&= \alpha K\theta - (1 - \alpha)\theta \\
&= (\alpha K - (1 - \alpha))\theta \\
&= ((1 + K)\alpha - 1)\theta
\end{aligned}$$

Set this to equal 0 and solve for  $\alpha$ :

$$\begin{aligned} ((1+K)\alpha - 1)\theta &= 0 \\ (1+K)\alpha - 1 &= 0 \\ \alpha &= \frac{1}{1+K} \end{aligned}$$

where we assumed that  $\theta \neq 0$  and  $K \neq -1$ . So the unbiased estimator is

$$\hat{\theta}_3 = \frac{1}{1+K}\hat{\theta}.$$

So instead of multiplying by  $1-K$  to get an unbiased estimate, we should divide by  $1+K$ . However, for  $K$  near 0,  $\hat{\theta}_2 = (1-K)\hat{\theta}$  is almost unbiased, since  $\text{bias}(\hat{\theta}_2) = -K^2\theta$  which gets small in magnitude relative to  $\theta$  for  $K$  near 0. In fact,  $\hat{\theta}_2$  is a local linear approximation to  $\hat{\theta}_3$  about  $K = 0$ . One way to see this analytically is to write the Maclaurin series for  $\frac{1}{1+K}$  using the well-known geometric series formula:

$$\frac{1}{1+K} = \frac{1}{1-(-K)} = \sum_{n=0}^{\infty} (-K)^n = \sum_{n=0}^{\infty} (-1)^n K^n = 1 - K + K^2 - K^3 + K^4 - \dots$$

Therefore, for  $K$  near 0,  $1-K \approx \frac{1}{1+K}$  and hence  $\hat{\theta}_2 \approx \hat{\theta}_3$ . This means that if someone derives a bias-correction estimator of the form of  $\hat{\theta}_2$ , and when calculating the value of the  $\hat{\theta}_2$ ,  $K$  is generally small in magnitude, they might not notice that the estimator is actually slightly unbiased.

### 3 Minimum MSE

But what if, rather than an unbiased estimator, you want an estimator with minimum mean squared error? Loosely speaking, you want an estimator that is as close as possible on average to the true value. Call this minimum MSE estimator  $\hat{\theta}_4$ , and again write this estimator in the form

$$\hat{\theta}_4 = \alpha\hat{\theta}.$$

We will choose  $\alpha$  such that  $\hat{\theta}_4$  achieves minimum mean squared error. Let's calculate:

$$\begin{aligned} \text{MSE}(\hat{\theta}_4) &= \text{bias}^2(\hat{\theta}_4) + \mathbb{V}(\hat{\theta}_4) \\ &= (\text{bias}(\alpha\hat{\theta}))^2 + \mathbb{V}(\alpha\hat{\theta}) \\ &= [((1+K)\alpha - 1)\theta]^2 + \alpha^2\mathbb{V}(\hat{\theta}) \\ &= \theta^2((1+K)\alpha - 1)^2 + \mathbb{V}(\hat{\theta})\alpha^2. \end{aligned}$$

We're treating  $\text{MSE}(\hat{\theta}_4)$  as a function of  $\alpha$  and seeking to find  $\alpha$  that minimizes it. So let's find the derivative and set it to 0.

$$\begin{aligned} \frac{d\text{MSE}(\hat{\theta}_4)}{d\alpha} &= 2\theta^2((1+K)\alpha - 1)(1+K) + 2\mathbb{V}(\hat{\theta})\alpha \\ &= 2\theta^2(1+K)^2\alpha - 2\theta^2(1+K) + 2\mathbb{V}(\hat{\theta})\alpha \\ &= [2\theta^2(1+K)^2 + 2\mathbb{V}(\hat{\theta})]\alpha - 2\theta^2(1+K) = 0 \end{aligned}$$

Solving for  $\alpha$  yields

$$\alpha^* = \frac{\theta^2(1+K)}{\theta^2(1+K)^2 + \mathbb{V}(\hat{\theta})} = \frac{1}{1+K + \frac{\mathbb{V}(\hat{\theta})}{\theta^2(1+K)}}.$$

But we need to check that  $\text{MSE}(\hat{\theta}_4)$  is a global minimum at the critical point  $\alpha^*$ . Applying the second derivative test,

$$\frac{d^2 \text{MSE}(\hat{\theta}_4)}{d\alpha^2} = 2\theta^2(1+K)^2 + 2\mathbb{V}(\hat{\theta})$$

which must be positive, as long as  $K \neq -1$  and  $\theta \neq 0$ . In particular the second derivative is positive at  $\alpha^*$ , so by the second-derivative test, the critical point is at a local minimum. Also, since the second derivative is always positive,  $\text{MSE}(\hat{\theta}_4)$  is a convex (parabola) function of  $\alpha$ , hence the critical point is the global minimum.

Notice the form of  $\alpha^*$ ; it has  $\theta$  in it, but  $\theta$  is the very thing we are trying to estimate! If we estimate  $\theta$  with  $\hat{\theta}$  in  $\alpha$  (compared to other proposed methods, this was empirically seen to show the best results by running simulations), we get the following estimator:

$$\hat{\theta}_4 := \frac{\hat{\theta}}{1+K + \frac{\mathbb{V}(\hat{\theta})}{\hat{\theta}^2(1+K)}}.$$

Hence we must know  $\mathbb{V}(\theta)$ , or at least an approximation of it  $\hat{\mathbb{V}}(\hat{\theta})$ , to get a minimum MSE estimator.

Comparing the unbiased estimator  $\hat{\theta}_3$  and the minimum MSE estimator  $\hat{\theta}_4$ , we see that  $\hat{\theta}_3$  is just  $\hat{\theta}_4$  but where  $\mathbb{V}(\hat{\theta})$  is set to be 0. Also, since  $\hat{\theta}_4$  is just  $\hat{\theta}_3$  but with an extra positive term added in the denominator, we see that  $\hat{\theta}_4 \leq \hat{\theta}_3$ , and hence  $\hat{\theta}_4$  should be negatively biased, i.e., tends to under-estimate  $\theta$  on average. Directly calculating  $\text{bias}(\hat{\theta}_4)$ , and going back to the original formula for  $\alpha^*$  by treating the  $\hat{\theta}$  in the denominator of  $\hat{\theta}_4$  (which was initially  $\theta$ , and is now a plug-in estimator for  $\theta$ ) as fixed at  $\theta$ , and setting  $C = \frac{\mathbb{V}(\hat{\theta})}{\theta^2}$ ,

$$\begin{aligned} \text{bias}\left(\frac{\hat{\theta}}{1+K + \frac{C}{1+K}}\right) &= \frac{1}{1+K + \frac{C}{1+K}} \text{bias}(\hat{\theta}) - \left(1 - \frac{1}{1+K + \frac{C}{1+K}}\right)\theta \\ &= \frac{K\theta}{1+K + \frac{C}{1+K}} - \frac{(K + \frac{C}{1+K})\theta}{1+K + \frac{C}{1+K}} \\ &= \frac{-\frac{C}{1+K}}{1+K + \frac{C}{1+K}}\theta \\ &= -\frac{C}{1+2K+K^2+C}\theta \\ &= -\frac{C}{(1+K)^2+C}\theta = -\frac{\theta}{\frac{(1+K)^2}{C}+1}. \end{aligned}$$

Hence we see that  $\text{bias}(\hat{\theta}_4) \leq 0$ . As the value of  $C$  increases, the magnitude of this bias will increase.

## 4 Bias Correction Is Maximum Likelihood

Now what if we instead try to solve the problem using maximum likelihood? We will do this like so: From the equation

$$\text{bias}(\hat{\theta}) = \mathbb{E}(\hat{\theta}) - \theta = K\theta$$

and solving for  $\mathbb{E}(\hat{\theta})$  we get

$$\mathbb{E}(\hat{\theta}) = (1 + K)\theta.$$

And suppose, as in the minimum MSE part, that we know, or at least can approximate,  $\mathbb{V}(\hat{\theta})$ . Write  $\mu = \mathbb{E}(\hat{\theta})$  and  $\sigma^2 = \mathbb{V}(\hat{\theta})$ , and assume that  $\hat{\theta} \sim \mathcal{N}(\mu, \sigma^2)$ . The likelihood function is then

$$\mathcal{L}(\theta) = f_{\hat{\theta}}(\hat{\theta}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\hat{\theta} - \mu)^2}{2\sigma^2}\right).$$

So the log-likelihood is

$$\ell(\theta) = \log \mathcal{L}(\theta) = \log\left(\frac{1}{\sqrt{2\pi}\sigma}\right) - \frac{(\hat{\theta} - \mu)^2}{2\sigma^2}.$$

Differentiate and set the derivative to 0:

$$\frac{d\ell(\theta)}{d\theta} = -\frac{(\hat{\theta} - \mu)(-(1 + K))}{\sigma^2} = \frac{(\hat{\theta} - (1 + K)\theta)(1 + K)}{\sigma^2} = 0$$

Solve for  $\theta$ :

$$\theta = \frac{\hat{\theta}}{1 + K} = \hat{\theta}_3.$$

Let's verify that this is the maximum using the second derivative test:

$$\frac{d^2\ell(\theta)}{d\theta^2} = -\frac{(1 + K)^2}{\sigma^2}$$

which is negative as long as  $K \neq -1$ . So we find that our unbiased estimator is also the MLE assuming a normal distribution.

## 5 Variance and Standard Error

The fact that the unbiased estimator is also the MLE makes this estimator more appealing as MLE is very popular. It is also appealing compared to the minimum-MSE estimator because it does not require an estimation of the variance of  $\hat{\theta}$ . But if we do have an estimation of the variance of  $\hat{\theta}$ , then the variance of  $\hat{\theta}_3$  is

$$\widehat{\mathbb{V}}(\hat{\theta}_3) = \widehat{\mathbb{V}}\left(\frac{\hat{\theta}}{1 + K}\right) = \frac{\widehat{\mathbb{V}}(\hat{\theta})}{(1 + K)^2}.$$

So we can get a standard error:

$$\widehat{\text{se}}(\hat{\theta}_3) = \frac{\widehat{\text{se}}(\hat{\theta})}{|1 + K|}.$$

## 6 Conclusion

In conclusion, we found three “correction” estimators for the value of  $\theta$  if we start with an estimation  $\hat{\theta}$  and know that  $\text{bias}(\hat{\theta}) = K\theta$  for some number  $K$ . Assuming that  $\theta \neq 0$  and  $K \neq -1$ , we find three estimators:

$$\hat{\theta}_2 = (1 - K)\hat{\theta}, \quad \hat{\theta}_3 = \frac{\hat{\theta}}{1 + K}, \quad \hat{\theta}_4 = \frac{\hat{\theta}}{1 + K + \frac{\widehat{\mathbb{V}}(\hat{\theta})}{\hat{\theta}^2(1 + K)}}.$$

$\hat{\theta}_3$  corrects the bias of  $\hat{\theta}$  by multiplying by the factor  $1/(1 + K)$ , and is also the MLE;  $\hat{\theta}_2$  is a naive bias correction that is actually slightly biased but very close to  $\hat{\theta}_3$  for  $K$  near 0; and  $\hat{\theta}_4$  attempts to minimize the mean squared error and requires an estimate of the variance of  $\hat{\theta}$ ,  $\widehat{\mathbb{V}}(\hat{\theta})$ .  $\hat{\theta}_4$  is negatively biased, meaning it on average under-estimates  $\theta$ .  $\hat{\theta}_3$  is  $\hat{\theta}_4$  but with  $\widehat{\mathbb{V}}(\hat{\theta}) = 0$ .