

Properties of Estimators (and the MLE)

Is the MLE any good? More generally, how can we talk in measured terms about the efficacy of a particular estimator?

Let's start by carefully laying out our problem scenario, and what we mean by an “estimator”. Suppose that we observe independent realizations $X_1 = \mathbf{x}_1, \dots, X_N = \mathbf{x}_N$, $\mathbf{x}_n \in \mathbb{R}^D$ of a random vector with distribution

$$X_n \sim f_X(\mathbf{x}; \boldsymbol{\theta}_0),$$

for some fixed $\boldsymbol{\theta}_0 \in \mathbb{R}^P$ (that is, there are P parameters that control the distribution of the X_n). An **estimator** based on these N observations is some fixed mapping $g(\mathbf{x}_1, \dots, \mathbf{x}_N) : \mathbb{R}^D \rightarrow \mathcal{T}$, defined for every N , that takes the observed data and produces¹ an estimate $\hat{\boldsymbol{\theta}}_N$.

With the randomness of the data taken into account, the estimate is a function of a sequence of random variables,

$$\hat{\boldsymbol{\Theta}}_N = g(X_1, \dots, X_N),$$

and so is itself a random variable. We will discuss three properties of $\hat{\boldsymbol{\Theta}}_N$ that might be desirable:

Bias: An estimator is called *unbiased* if it is equal to $\boldsymbol{\theta}_0$ on average.

Consistency: An estimator is called *consistent* if it is guaranteed to return $\boldsymbol{\theta}_0$ as the amount of data we have goes to infinity.

Efficiency: An unbiased estimator estimator is called *efficient* if

¹The MLE, for example, evaluates the log likelihood functions at the N different points for all values of $\boldsymbol{\theta}$, adds these together, then returns the $\boldsymbol{\theta}$ for which this sum is maximized.

it has the minimum mean-squared error of any possible estimator procedure for a particular value of N .

Bias

The bias of an estimator $\hat{\Theta}$ is simply

$$\text{bias}(\hat{\Theta}) = E[\hat{\Theta}] - \theta_0.$$

An estimator is called **unbiased** if the bias is zero for all values of $\theta_0 \in \mathcal{T}$.

The bias plays a role in describing the mean-squared error (MSE) of an estimator. For a scalar $\hat{\Theta}$, the MSE is

$$\begin{aligned} \text{MSE}(\hat{\Theta}) &= E[(\hat{\Theta} - \theta_0)^2] \\ &= E[(\hat{\Theta} - E[\hat{\Theta}])^2] + (E[\hat{\Theta}] - \theta_0)^2 \\ &= \text{var}(\hat{\Theta}) + \text{bias}(\hat{\Theta})^2. \end{aligned}$$

For vector-valued estimators, this extends to

$$\begin{aligned} \text{MSE}(\hat{\Theta}) &= E[\|\hat{\Theta} - \theta_0\|_2^2] \\ &= \text{trace}(\hat{\mathbf{R}}) + \|E[\hat{\Theta}] - \theta_0\|_2^2 \end{aligned}$$

where $\hat{\mathbf{R}}$ is the covariance matrix for the estimator,

$$\hat{\mathbf{R}} = E[(\hat{\Theta} - E[\hat{\Theta}])(\hat{\Theta} - E[\hat{\Theta}])^T]$$

It is often the case that we can decrease the variance but at a cost of increasing the bias or vice versa — this is referred to as the **bias variance trade off**.

An estimator is called **asymptotically unbiased** if

$$\mathbb{E}[\hat{\boldsymbol{\Theta}}_N] \rightarrow \boldsymbol{\theta}_0 \quad \text{as } N \rightarrow \infty,$$

for all $\boldsymbol{\theta}_0 \in \mathcal{T}$.

The MLE is sometimes unbiased and sometimes biased, but it is always asymptotically unbiased. We will not prove this here, but it will come out in our discussion of consistency in the next section.

Consistency

Consistent estimators return estimates that become arbitrarily close to $\boldsymbol{\theta}_0$ as N becomes large. To be a useful notion, this should be true not just for a particular data set, but should be “typical” behavior over large ensembles of data sets. This is made precise in the following way. We say that the estimator is **consistent** if for every $\epsilon > 0$ and $0 < \delta < 1$, there is an N (that of course depends on ϵ, δ) such that

$$\mathbb{P} \left(|\hat{\boldsymbol{\Theta}}_N - \boldsymbol{\theta}_0| > \epsilon \right) \leq \delta.$$

In words, this means that the probability that the estimate is off by more than ϵ can be made arbitrarily small with enough data.

Let’s look at a particular example. Suppose that the X_n are binary-valued random variables,

$$X_n = \begin{cases} 1, & \text{with probability } \theta_0, \\ 0, & \text{with probability } 1 - \theta_0, \end{cases}$$

and let $\hat{\boldsymbol{\Theta}}_N$ be the sample mean (which we have already seen is also

the MLE):

$$\hat{\Theta}_N = \frac{1}{N} \sum_{n=1}^N X_n.$$

It is clear that $\hat{\Theta}_N$ is unbiased, as

$$\mathbb{E}[\hat{\Theta}_N] = \frac{1}{N} \sum_{n=1}^N \mathbb{E}[X_n] = \theta_0,$$

and the variance is

$$\text{var}(\hat{\Theta}_N) = \frac{1}{N^2} \sum_{n=1}^N \text{var}(X_n) = \frac{1}{N^2} \sum_{n=1}^N \theta_0(1 - \theta_0) = \frac{\theta_0(1 - \theta_0)}{N}.$$

The Chebyshev inequality (see the Technical Details section) says that

$$\mathbb{P} \left(|\hat{\Theta}_N - \theta| > \epsilon \right) \leq \frac{\theta_0(1 - \theta_0)}{N\epsilon^2},$$

that is for any $0 < \delta < 1$,

$$\mathbb{P} \left(|\hat{\Theta}_N - \theta| > \epsilon \right) \leq \delta \quad \text{when} \quad N > \frac{\theta_0(1 - \theta_0)}{\epsilon^2\delta}.$$

As $\theta(1 - \theta) \leq 1/4$ for all $\theta \in [0, 1]$, we can ensure the probability above when

$$N > \frac{1}{4\epsilon^2\delta}.$$

So, the sample mean is a consistent estimator for the probability that a binary random variable is equal to 1.

Consistency of the MLE

We can establish the consistency of the MLE from a series of independent samples under fairly general conditions. As above, suppose that

$$X_n \sim f_X(\mathbf{x}; \boldsymbol{\theta}_0), \quad \text{for some } \boldsymbol{\theta}_0 \in \mathcal{T}. \quad (1)$$

The MLE takes observations $X_1 = \mathbf{x}_1, \dots, X_N = \mathbf{x}_N$ and then finds the maximum of the log likelihood:

$$\hat{\boldsymbol{\theta}}_{\text{MLE}} = \arg \max_{\boldsymbol{\theta} \in \mathcal{T}} \ell(\boldsymbol{\theta}; \mathbf{x}_1, \dots, \mathbf{x}_N) = \arg \max_{\boldsymbol{\theta} \in \mathcal{T}} \sum_{n=1}^N \ell(\boldsymbol{\theta}; \mathbf{x}_n),$$

where

$$\ell(\boldsymbol{\theta}; \mathbf{x}_n) = \log f_X(\mathbf{x}_n; \boldsymbol{\theta}).$$

To make this problem well-posed, we will assume that the collection of probability models $\{f_X(\mathbf{x}; \boldsymbol{\theta}), \boldsymbol{\theta} \in \mathcal{T}\}$ is **identifiable**. This means that for each value of $\boldsymbol{\theta} \in \mathcal{T}$, $f_X(\mathbf{x}; \boldsymbol{\theta})$ is a different density function. More precisely, it means that if I draw random a X using the probability law (1), then evaluate² two different densities at this point, then there is at least some chance I get different answers:

$$\mathbb{P}(f_X(X; \boldsymbol{\theta}_1) \neq f_X(X; \boldsymbol{\theta}_2)) > 0 \quad \text{for all } \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathcal{T}, \quad \boldsymbol{\theta}_1 \neq \boldsymbol{\theta}_2.$$

Of course, for different sets of observations $\{\mathbf{x}_n\}$, the log likelihood function $\ell(\boldsymbol{\theta}; \mathbf{x}_1, \dots, \mathbf{x}_N)$ will be different, and hence so will it maximizer $\hat{\boldsymbol{\theta}}_{\text{MLE}}$. Taking the random nature of the observations into

²It is worth re-emphasizing here that the quantities $f_X(X; \boldsymbol{\theta}_1)$ and $f_X(X; \boldsymbol{\theta}_2)$ are random variables — they are created by evaluating densities at a common random point.

account, we can think of

$$\ell(\boldsymbol{\theta}; X_1, \dots, X_N) = \sum_{n=1}^N \ell(\boldsymbol{\theta}; X_n),$$

as a random function of $\boldsymbol{\theta}$ — the draw of the X_n will determine what this function is. This random function has some well-defined expected value at every point $\boldsymbol{\theta} \in \mathcal{T}$; we call this function the

$$\text{expected log likelihood} = \mathbb{E}[\ell(\boldsymbol{\theta}; X_1, \dots, X_N)].$$

While $\ell(\boldsymbol{\theta}; X_1, \dots, X_N)$ is a random function, the expected log likelihood $\mathbb{E}[\ell(\boldsymbol{\theta}; X_1, \dots, X_N)]$ is a deterministic function of $\boldsymbol{\theta}$. Note that since the X_n are iid,

$$\mathbb{E}[\ell(\boldsymbol{\theta}; X_1, \dots, X_N)] = N \mathbb{E}[\ell(\boldsymbol{\theta}; X)],$$

where $X \sim f_X(\mathbf{x}; \boldsymbol{\theta}_0)$.

The key fact is that if the X_n are drawn as in (1), the the expected log likelihood peaks exactly where we need it to:

$$\arg \max_{\boldsymbol{\theta} \in \mathcal{T}} \mathbb{E}[\ell(\boldsymbol{\theta}; X_1, \dots, X_N)] = \boldsymbol{\theta}_0. \quad (2)$$

To prove this, we start by noting that for any $\boldsymbol{\theta}_1 \neq \boldsymbol{\theta}_0$,

$$\mathbb{E} \left[\frac{f_X(X; \boldsymbol{\theta}_1)}{f_X(X; \boldsymbol{\theta}_0)} \right] \leq 1.$$

Indeed, if \mathcal{X}_0 is the support of the density $f(\mathbf{x}; \boldsymbol{\theta}_0)$,

$$\mathcal{X}_0 = \{\mathbf{x} \in \mathbb{R}^D : f_X(\mathbf{x}; \boldsymbol{\theta}_0) > 0\},$$

then

$$\begin{aligned} \mathbb{E} \left[\frac{f_X(X; \boldsymbol{\theta}_1)}{f_X(X; \boldsymbol{\theta}_0)} \right] &= \int_{\mathbf{x} \in \mathcal{X}_0} \frac{f_X(\mathbf{x}; \boldsymbol{\theta}_1)}{f_X(\mathbf{x}; \boldsymbol{\theta}_0)} f_X(\mathbf{x}; \boldsymbol{\theta}_0) \, d\mathbf{x} \\ &= \int_{\mathbf{x} \in \mathcal{X}_0} f_X(\mathbf{x}; \boldsymbol{\theta}_1) \, d\mathbf{x} \\ &\leq 1. \end{aligned}$$

Now consider any $\boldsymbol{\theta}_1 \neq \boldsymbol{\theta}_0$. We have

$$\begin{aligned} \mathbb{E}[\ell(\boldsymbol{\theta}_1; X)] - \mathbb{E}[\ell(\boldsymbol{\theta}_0; X)] &= \mathbb{E} \left[\log \left(\frac{f_X(X; \boldsymbol{\theta}_1)}{f_X(X; \boldsymbol{\theta}_0)} \right) \right] \\ &\leq \mathbb{E} \left[\frac{f_X(X; \boldsymbol{\theta}_1)}{f_X(X; \boldsymbol{\theta}_0)} - 1 \right], \quad \text{since } \log t \leq t - 1, \\ &\leq 0. \end{aligned}$$

Since $\log(t) = t - 1$ only at $t = 1$, the only way the first \leq above holds with equality is if $\mathbb{P}(f_X(X; \boldsymbol{\theta}_1) = f_X(X; \boldsymbol{\theta}_0)) = 1$, which our identifiability condition forbids. Thus

$$\mathbb{E}[\ell(\boldsymbol{\theta}_1; X)] < \mathbb{E}[\ell(\boldsymbol{\theta}_0; X)] \quad \text{for all } \boldsymbol{\theta}_1 \in \mathcal{T}, \boldsymbol{\theta}_1 \neq \boldsymbol{\theta}_0,$$

which establishes (2).

So if we could maximize the expected log likelihood, we would get exactly the right answer every time. Of course, we don't have access to this function; the MLE instead maximizes the sample mean:

$$\hat{\boldsymbol{\theta}}_{\text{MLE}} = \arg \max_{\boldsymbol{\theta} \in \mathcal{T}} \frac{1}{N} \sum_{n=1}^N \ell(\boldsymbol{\theta}; \mathbf{x}_n).$$

By the weak law of large numbers,

$$\frac{1}{N} \sum_{n=1}^N \ell(\boldsymbol{\theta}; \mathbf{x}_n) \rightarrow \mathbb{E}[\ell(\boldsymbol{\theta}; X)], \quad \text{for all } \boldsymbol{\theta} \in \mathcal{T},$$

and so the maxima of the left- and right-hand sides will also eventually coincide.

We easily can make this perfectly precise when \mathcal{T} is finite. In this case, we know that there is some $\gamma > 0$ such that

$$\mathbb{E}[\ell(\boldsymbol{\theta}; X)] \leq \mathbb{E}[\ell(\boldsymbol{\theta}_0; X)] - \gamma, \quad \text{for all } \boldsymbol{\theta} \neq \boldsymbol{\theta}_0.$$

Under the very mild assumption that $\text{var}(\ell(\boldsymbol{\theta}; X)) < \infty$ for all $\boldsymbol{\theta}$, the Chebyshev inequality tells us that for every $\epsilon > 0$ there exists an N such that

$$\mathbb{P} \left(\left| \frac{1}{N} \sum_{n=1}^N \ell(\boldsymbol{\theta}; \mathbf{x}_n) - \mathbb{E}[\ell(\boldsymbol{\theta}; X)] \right| > \gamma/2 \right) \leq \frac{4 \text{var}(\ell(\boldsymbol{\theta}; X))}{N\gamma^2},$$

and so

$$\begin{aligned} & \mathbb{P} \left(\max_{\boldsymbol{\theta} \in \mathcal{T}} \left| \frac{1}{N} \sum_{n=1}^N \ell(\boldsymbol{\theta}; \mathbf{x}_n) - \mathbb{E}[\ell(\boldsymbol{\theta}; X)] \right| > \gamma/2 \right) \\ & \leq \sum_{\boldsymbol{\theta} \in \mathcal{T}} \mathbb{P} \left(\left| \frac{1}{N} \sum_{n=1}^N \ell(\boldsymbol{\theta}; \mathbf{x}_n) - \mathbb{E}[\ell(\boldsymbol{\theta}; X)] \right| > \gamma/2 \right) \\ & \leq \sum_{\boldsymbol{\theta} \in \mathcal{T}} \frac{4 \text{var}(\ell(\boldsymbol{\theta}; X))}{N\gamma^2} \\ & \leq \frac{4m|\mathcal{T}|}{N\gamma^2}, \end{aligned}$$

where $m = \max_{\boldsymbol{\theta} \in \mathcal{T}} \text{var}(\ell(\boldsymbol{\theta}; X))$ and $|\mathcal{T}|$ is the number of elements in the set \mathcal{T} . The point is that this probability goes to 0 as N gets large, so that probability that

$$\frac{1}{N} \sum_{n=1}^N \ell(\boldsymbol{\theta}_0; \mathbf{x}_n) > \mathbb{E}[\ell(\boldsymbol{\theta}_0; X)] - \frac{\gamma}{2}$$

and

$$\frac{1}{N} \sum_{n=1}^N \ell(\boldsymbol{\theta}; \mathbf{x}_n) < \mathbb{E}[\ell(\boldsymbol{\theta}; X)] + \frac{\gamma}{2}, \quad \text{for all } \boldsymbol{\theta} \neq \boldsymbol{\theta}_0,$$

goes to 1. Thus for N large enough,

$$\frac{1}{N} \sum_{n=1}^N \ell(\boldsymbol{\theta}; \mathbf{x}_n) < \frac{1}{N} \sum_{n=1}^N \ell(\boldsymbol{\theta}_0; \mathbf{x}_n), \quad \text{for all } \boldsymbol{\theta} \neq \boldsymbol{\theta}_0,$$

and so

$$\mathbb{P} \left(\arg \max_{\boldsymbol{\theta} \in \mathcal{T}} \frac{1}{N} \sum_{n=1}^N \ell(\boldsymbol{\theta}; \mathbf{x}_n) = \arg \max_{\boldsymbol{\theta} \in \mathcal{T}} \mathbb{E}[\ell(\boldsymbol{\theta}; X)] = \boldsymbol{\theta}_0 \right) \rightarrow 1,$$

as $N \rightarrow \infty$. This establishes the consistency of the MLE.

When \mathcal{T} is an infinite set, we need additional assumptions on the density functions $f_X(\mathbf{x}; \boldsymbol{\theta})$. Using arguments not too dissimilar from the above, you can show that if

1. the $f_X(\mathbf{x}; \boldsymbol{\theta})$ have the same support for all $\boldsymbol{\theta} \in \mathcal{T}$,
2. the $\mathbb{E}[\ell(\boldsymbol{\theta}; X)]$ are all differentiable, and
3. $\boldsymbol{\theta}_0$ is away from the boundary of \mathcal{T} ,

then the MLE is consistent.

Efficiency

For unbiased estimators, where $E[\hat{\boldsymbol{\theta}}] = \boldsymbol{\theta}_0$ and so

$$\text{MSE}(\hat{\boldsymbol{\Theta}}) = \text{var}(\hat{\boldsymbol{\Theta}}),$$

there is a **fundamental limit** on the MSE that no (unbiased) estimator can surpass. This is called the **Cramer-Rao lower bound**, and it is a massive result from classical statistics. Estimators that achieve this limit as the data size gets large are called **efficient**.

The key quantity in deriving this fundamental limit is the **score function**:

$$\mathbf{s}(\boldsymbol{\theta}; \mathbf{x}) = \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathbf{x}).$$

As you can see, this is a vector-valued function of $\boldsymbol{\theta}$ that returns the gradient of the log likelihood at a fixed data point \mathbf{x} . By definition, the MLE looks for a value of $\boldsymbol{\theta}$ such that $\mathbf{s}(\boldsymbol{\theta}; \mathbf{x}) = \mathbf{0}$.

What is the behavior of this score function around the true value of the parameters $\boldsymbol{\theta} = \boldsymbol{\theta}_0$? As we would hope, the expectation of the score at $\boldsymbol{\theta}_0$ with respect to a random draw of $X \sim f_X(\mathbf{x}; \boldsymbol{\theta}_0)$ is indeed $\mathbf{0}$. To see this, we compute

$$\begin{aligned} E[\mathbf{s}(\boldsymbol{\theta}; X)] &= \int \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; X) f_X(\mathbf{x}; \boldsymbol{\theta}_0) \, d\mathbf{x} \\ &= \int \nabla_{\boldsymbol{\theta}} f_X(\mathbf{x}; \boldsymbol{\theta}) \frac{f_X(\mathbf{x}; \boldsymbol{\theta}_0)}{f_X(\mathbf{x}; \boldsymbol{\theta})} \, d\mathbf{x} \quad (\text{chain rule}), \end{aligned}$$

so³

$$\begin{aligned} \mathbb{E}[\mathbf{s}(\boldsymbol{\theta}_0; X)] &= \int \nabla_{\boldsymbol{\theta}} f_X(\mathbf{x}; \boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \, d\mathbf{x} \\ &= \nabla_{\boldsymbol{\theta}} \left(\int f_X(\mathbf{x}; \boldsymbol{\theta}) \, d\mathbf{x} \right) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \\ &= \nabla_{\boldsymbol{\theta}} 1|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \\ &= \mathbf{0}. \end{aligned}$$

So no matter what $\boldsymbol{\theta}_0$ is, in expectation the score function crosses $\mathbf{0}$ at the right place.

Intuitively, how effectively we are able to estimate $\boldsymbol{\theta}_0$ will depend on the *variation* of the score function around $\boldsymbol{\theta} = \boldsymbol{\theta}_0$. If $\mathbf{s}(\boldsymbol{\theta}; X)$ behaves very predictably around $\boldsymbol{\theta}_0$, the score function will pass through $\mathbf{0}$ in almost the same place for typical draws of the data, meaning that the data carry a lot of information about what $\boldsymbol{\theta}_0$ actually is. If $\mathbf{s}(\boldsymbol{\theta}; X)$ has wide variation around $\boldsymbol{\theta}_0$, then we will have a harder time estimating $\boldsymbol{\theta}_0$ from the data.

Mathematically, this variation is quantified with the covariance matrix of the score function:

$$\begin{aligned} \mathbf{J}(\boldsymbol{\theta}_0) &= \mathbb{E}[\mathbf{s}(\boldsymbol{\theta}_0; X)\mathbf{s}(\boldsymbol{\theta}_0; X)^T] \\ &= \mathbb{E}[\nabla_{\boldsymbol{\theta}} \log f_X(X; \boldsymbol{\theta}_0) \nabla_{\boldsymbol{\theta}} \log f_X(X; \boldsymbol{\theta}_0)^T] \end{aligned}$$

This is called the **Fisher information matrix** at $\boldsymbol{\theta}_0$, and as we see below it is indeed what determines our fundamental limit.

³Switching the integral and the gradient here actually requires some very mild smoothness conditions on how the densities change as a function of $\boldsymbol{\theta}$.

Let $\hat{\Theta}$ be any unbiased estimator, $E[\hat{\Theta}] = \theta_0$, that is formed from a *single sample*:

$$\hat{\Theta} = g(X_1),$$

where $g : \mathbb{R}^D \rightarrow \mathbb{R}^P$ is the mapping that actually implements the estimation algorithm. After deriving a bound on the performance of such single-sample estimators, we will quickly be able to generalize to multiple samples.

It is a fact that the cross-correlation between the error $\hat{\Theta} - \theta_0$ and the score vector at θ_0 is the identity:

$$E[\mathbf{s}(\theta_0; X)(\hat{\Theta} - \theta_0)^T] = \mathbf{I}.$$

To see this, we write down what it means for $\hat{\Theta}$ to be unbiased⁴,

$$E[(\hat{\Theta} - \theta_0)^T] = \int f_X(\mathbf{x}; \theta_0)(\hat{\theta} - \theta_0)^T d\mathbf{x} = \mathbf{0}^T,$$

meaning

$$\int f_X(\mathbf{x}; \theta_0)\hat{\theta}^T d\mathbf{x} = \int f_X(\mathbf{x}; \theta_0)\theta_0^T d\mathbf{x}.$$

Taking the gradient of both sides yields

$$\nabla_{\theta} \left(\int f_X(\mathbf{x}; \theta)\hat{\theta}^T d\mathbf{x} \right) \Big|_{\theta=\theta_0} = \nabla_{\theta} \left(\int f_X(\mathbf{x}; \theta)\theta^T d\mathbf{x} \right) \Big|_{\theta=\theta_0}. \quad (3)$$

⁴Note that $\hat{\Theta}$ is a function of the data — inside the integral here, this is a deterministic function of \mathbf{x} , so we write it lowercase $\hat{\theta}$. We are not explicitly denoting the dependence of $\hat{\theta}$ on \mathbf{x} , but it is there.

Switching⁵ the gradient and the integral and applying the chain rule to the right-hand side above yields

$$\begin{aligned}
& \nabla_{\theta} \left(\int f_X(\mathbf{x}; \boldsymbol{\theta}) \boldsymbol{\theta}^T \, d\mathbf{x} \right) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \\
&= \int \nabla_{\theta} \left(f_X(\mathbf{x}; \boldsymbol{\theta}) \boldsymbol{\theta}^T \right) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \, d\mathbf{x} \\
&= \int \nabla_{\theta} (f_X(\mathbf{x}; \boldsymbol{\theta})) \boldsymbol{\theta}^T \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \, d\mathbf{x} + \int f_X(\mathbf{x}; \boldsymbol{\theta}) \nabla_{\theta} (\boldsymbol{\theta}^T) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \, d\mathbf{x} \\
&= \int \nabla_{\theta} (f_X(\mathbf{x}; \boldsymbol{\theta}_0)) \boldsymbol{\theta}_0^T \, d\mathbf{x} + \int f_X(\mathbf{x}; \boldsymbol{\theta}) \mathbf{I} \, d\mathbf{x} \\
&= \int \nabla_{\theta} (f_X(\mathbf{x}; \boldsymbol{\theta}_0)) \boldsymbol{\theta}_0^T \, d\mathbf{x} + \mathbf{I}.
\end{aligned}$$

Returning to (3), pulling the gradient inside the integral on the left hand side and combining with the above yields

$$\int \nabla_{\theta} f_X(\mathbf{x}; \boldsymbol{\theta}_0) \left(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \right)^T \, d\mathbf{x} = \mathbb{E}[\mathbf{s}(\boldsymbol{\theta}_0; X)(\hat{\boldsymbol{\Theta}} - \boldsymbol{\theta}_0)^T] = \mathbf{I}.$$

The pieces are now in place for our lower bound. Again, let $\hat{\boldsymbol{\Theta}}$ be any unbiased estimator, and denote its covariance matrix as

$$\hat{\mathbf{R}} = \mathbb{E}[(\hat{\boldsymbol{\Theta}} - \boldsymbol{\theta}_0)(\hat{\boldsymbol{\Theta}} - \boldsymbol{\theta}_0)^T].$$

Recall that the MSE of $\hat{\boldsymbol{\Theta}}$ is

$$\text{MSE}(\hat{\boldsymbol{\Theta}}) = \mathbb{E}[\|\hat{\boldsymbol{\Theta}} - \boldsymbol{\theta}_0\|_2^2] = \text{trace}(\hat{\mathbf{R}}).$$

Now let Q be the random vector

$$Q = \begin{bmatrix} \hat{\boldsymbol{\Theta}} - \boldsymbol{\theta}_0 \\ \mathbf{s}(\boldsymbol{\theta}_0; X) \end{bmatrix}.$$

⁵Again, we need regularity conditions on the densities to make this completely valid.

The covariance for Q is then

$$\mathbb{E}[QQ^T] = \begin{bmatrix} \hat{\mathbf{R}} & \mathbf{I} \\ \mathbf{I} & \mathbf{J}(\boldsymbol{\theta}_0) \end{bmatrix}$$

Since $\mathbb{E}[QQ^T]$ is symmetric positive semi-definite, so is

$$\begin{bmatrix} \mathbf{I} & -\mathbf{J}(\boldsymbol{\theta}_0)^{-1} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{R}} & \mathbf{I} \\ \mathbf{I} & \mathbf{J}(\boldsymbol{\theta}_0) \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\mathbf{J}(\boldsymbol{\theta}_0)^{-1} & \mathbf{I} \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{R}} - \mathbf{J}(\boldsymbol{\theta}_0)^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{J}(\boldsymbol{\theta}_0) \end{bmatrix}$$

Block diagonal matrices are symmetric positive semi-definite if and only if their blocks are sym+semi-def. Thus

$$\hat{\mathbf{R}} - \mathbf{J}(\boldsymbol{\theta}_0)^{-1} \quad \text{is symmetric positive semi-definite,}$$

which we write as

$$\hat{\mathbf{R}} \succeq \mathbf{J}(\boldsymbol{\theta}_0)^{-1}.$$

In particular,

$$\text{MSE}(\hat{\boldsymbol{\Theta}}) = \text{trace}(\hat{\mathbf{R}}) \geq \text{trace}(\mathbf{J}(\boldsymbol{\theta}_0)^{-1})$$

This result is called the **Cramer-Rao lower bound**.

The analysis for estimates formed from multiple independent samples is not too different. Indeed, we could treat this problem by just drawing a single random vector in \mathbb{R}^{DN} , where the vector can be partitioned into N independent vectors of length D . In the end, the analysis goes precisely as the above, only with the score function

$$\begin{aligned} \mathbf{s}(\boldsymbol{\theta}; X_1, \dots, X_N) &= \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; X_1, \dots, X_N) \\ &= \sum_{n=1}^N \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; X_n). \end{aligned}$$

The the Fisher information matrix becomes

$$\begin{aligned}
\mathbf{J}_N(\boldsymbol{\theta}_0) &= \mathbb{E}[\mathbf{s}(\boldsymbol{\theta}_0; X_1, \dots, X_N) \mathbf{s}(\boldsymbol{\theta}_0; X_1, \dots, X_N)^T] \\
&= \sum_{n=1}^N \sum_{m=1}^N \mathbb{E}[\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}_0; X_n) \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}_0; X_m)^T] \\
&= \sum_{n=1}^N \mathbb{E}[\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}_0; X_n) \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}_0; X_n)^T] \\
&= N \mathbf{J}(\boldsymbol{\theta}_0).
\end{aligned}$$

Thus any unbiased estimator $\hat{\boldsymbol{\Theta}}_N$ that works from N samples must have a covariance matrix $\hat{\mathbf{R}}_N$ that obeys

$$\mathbb{E}[(\hat{\boldsymbol{\Theta}}_N - \boldsymbol{\theta}_0)(\hat{\boldsymbol{\Theta}}_N - \boldsymbol{\theta}_0)^T] = \hat{\mathbf{R}}_N \succeq \frac{1}{N} \mathbf{J}(\boldsymbol{\theta}_0)^{-1}.$$

Cramer-Rao Lower Bound. Let X_1, \dots, X_N be iid random vectors in \mathbb{R}^D with distribution

$$X_n \sim f_X(\mathbf{x}; \boldsymbol{\theta}_0).$$

Let $\hat{\boldsymbol{\Theta}}_N = g(X_1, \dots, X_N) : \mathbb{R}^{D \times N} \rightarrow \mathbb{R}^D$ be any estimator with $\mathbb{E}[\hat{\boldsymbol{\Theta}}_N] = \boldsymbol{\theta}_0$. Then

$$\mathbb{E}[(\hat{\boldsymbol{\Theta}}_N - \boldsymbol{\theta}_0)(\hat{\boldsymbol{\Theta}}_N - \boldsymbol{\theta}_0)^T] = \hat{\mathbf{R}}_N \succeq \frac{1}{N} \mathbf{J}(\boldsymbol{\theta}_0)^{-1},$$

and in particular

$$\text{MSE}(\hat{\boldsymbol{\Theta}}) = \mathbb{E}[\|\hat{\boldsymbol{\Theta}} - \boldsymbol{\theta}_0\|_2^2] \geq \frac{1}{N} \text{trace}(\mathbf{J}(\boldsymbol{\theta}_0)^{-1}),$$

where

$$\mathbf{J}(\boldsymbol{\theta}_0) = \mathbb{E}[\mathbf{s}(\boldsymbol{\theta}_0; X)\mathbf{s}(\boldsymbol{\theta}_0; X)^T]$$

Example: Let's return to our binary random variable example, where

$$X = \begin{cases} 1, & \text{with probability } \theta_0, \\ 0, & \text{with probability } 1 - \theta_0. \end{cases}$$

Recall that with

$$\hat{\Theta}_N = \frac{1}{N} \sum_{n=1}^N X_n$$

we have

$$\text{var}(\Theta_N) = \frac{\theta_0(1 - \theta_0)}{N}.$$

We compute the Fisher information matrix (actually scalar in this case) as follows. We have

$$\ell(x; \theta) = \theta^x (1 - \theta)^{1-x}, \quad \nabla_{\theta} \ell(x; \theta) = \frac{x}{\theta} - \frac{(1-x)}{(1-\theta)},$$

(since x is a scalar in this case, the gradient above is really just a regular derivative). Then

$$\begin{aligned} J(\theta) &= \mathbb{E} \left[\left(\frac{X}{\theta} - \frac{(1-X)}{(1-\theta)} \right)^2 \right] \\ &= \frac{\mathbb{E}[X^2(1-\theta)^2 - 2X(1-X)\theta(1-\theta) + \theta^2(1-X)^2]}{\theta^2(1-\theta)^2} \\ &= \frac{\theta(1-\theta)^2 + \theta^2(1-\theta)}{\theta^2(1-\theta)^2} \\ &= \frac{1}{\theta(1-\theta)}. \end{aligned}$$

Thus⁶

$$\frac{1}{N} \text{trace}(J(\theta_0)^{-1}) = \frac{\theta_0(1-\theta_0)}{N}.$$

This matches the variance, so the sample mean (which is also the MLE) is efficient in this case.

⁶Again, $J(\theta_0)$ is a scalar here, so $\text{trace}(J(\theta_0)^{-1}) = 1/J(\theta_0)$.

Efficiency of the MLE. We will not prove the efficiency of the MLE. Rather we will simply state a more general and very powerful result. Under smoothness conditions on the expected log likelihood $E[\ell(\boldsymbol{\theta}; X)]$ in the neighborhood of $\boldsymbol{\theta}_0$ (which mostly say that $E[\ell(\boldsymbol{\theta}; X)]$ has three well-behaved derivatives), the distribution of the error in the MLE becomes Gaussian:

$$(\hat{\boldsymbol{\Theta}} - \boldsymbol{\theta}_0) \rightarrow Z \sim \text{Normal}(\mathbf{0}, \mathbf{J}(\boldsymbol{\theta}_0)^{-1}/N).$$

This means that the MLE is asymptotically unbiased, and that its MSE approaches the Cramer-Rao lower bound.

If you want more detailed proofs for the asymptotics of the MLE, a good reference is:

L. Wasserman, *All of Statistics*, 2nd ed, Springer, 2010.

<http://www.stat.cmu.edu/~larry/all-of-statistics/>

Technical Details: Weak Law of Large Numbers

In this section, we show that under very mild conditions, the sample mean converges to the true mean. The only condition is that the underlying distribution has finite variance.

We start by stating the main result precisely. Let X be a random variable with pdf $f_X(x)$, mean $E[X] = \mu$, and variance $\text{var}(X) = \sigma^2 < \infty$. We observe *samples* of X labeled X_1, X_2, \dots, X_N . The X_i are independent of one another, and they all have the same distribution as X . We will show that the sample mean formed from a sample of size N :

$$M_N = \frac{1}{N}(X_1 + X_2 + \dots + X_N),$$

obeys

$$P(|M_N - \mu| > \epsilon) \leq \frac{\sigma^2}{N\epsilon^2},$$

where $\epsilon > 0$ is an arbitrarily small number. In the expression above, M_N is the only thing which is random; μ and σ^2 are fixed underlying properties of the distribution, N is the amount of data we see, and ϵ is something we can choose arbitrarily.

Notice that no matter how small ϵ is, the probability on the right hand side above goes to zero as $N \rightarrow \infty$. That is, for any fixed $\epsilon > 0$,

$$\lim_{N \rightarrow \infty} P(|M_N - \mu| > \epsilon) = 0.$$

This result follows from two simple but important tools known as the *Markov* and *Chebyshev* inequalities.

Markov inequality

Let X be a random variable that only takes positive values:

$$f_X(x) = 0, \quad \text{for } x < 0, \quad \text{or} \quad F_X(0) = 0.$$

Then

$$\boxed{P(X \geq a) \leq \frac{E[X]}{a} \quad \text{for all } a > 0.}$$

For example, the probability that X is more than 5 times its mean is $1/5$, 10 times the mean is $1/10$, etc. And this holds for **any distribution**.

The Markov inequality is easy to prove:

$$\begin{aligned} E[X] &= \int_0^\infty x f_X(x) \, dx \\ &\geq \int_a^\infty x f_X(x) \, dx \\ &\geq \int_a^\infty a f_X(x) \, dx \\ &= a \cdot P(X \geq a) \end{aligned}$$

and so $P(X \geq a) \leq \frac{E[X]}{a}$.

Again, this is a very general statement in that we have assumed nothing about X other than it is positive. The price for the generality is that the bound is typically very loose, and does not usually capture the behavior of $P(X \geq a)$. We can, however, cleverly apply the Markov inequality to get something slightly more useful.

Chebyshev inequality

The main use of the Markov inequality turns out to be its use in deriving other, more accurate deviation inequalities. Here we will use it to derive the **Chebyshev inequality**, from which the weak law of large numbers will follow immediately.

Chebyshev inequality: If X is a random variable with mean μ and variance σ^2 , then

$$\boxed{\mathrm{P}(|X - \mu| > c) \leq \frac{\sigma^2}{c^2} \quad \text{for all } c > 0.}$$

The Chebyshev inequality follows immediately from the Markov inequality in the following way. No matter what range of values X takes, the quantity $|X - \mu|^2$ is always positive. Thus

$$\mathrm{P}(|X - \mu|^2 > c^2) \leq \frac{\mathrm{E}[|X - \mu|^2]}{c^2} = \frac{\sigma^2}{c^2}.$$

Since squaring $(\cdot)^2$ is monotonic (invertible) over positive numbers,

$$\mathrm{P}(|X - \mu|^2 > c^2) = \mathrm{P}(|X - \mu| > c) \leq \frac{\sigma^2}{c^2}.$$

We now have a bound which depends on the mean and the variance of X ; this leads to a more accurate approximation of the probability.

The weak law of large numbers (WLLN)

We now turn to the behavior of the the sample mean

$$M_N = \frac{X_1 + X_2 + \cdots + X_N}{N},$$

where again the X_i are iid random variables with $E[X_i] = \mu$ and $\text{var } X_i = \sigma^2$. We know that

$$E[M_N] = \frac{E[X_1] + E[X_2] + \cdots + E[X_N]}{N} = \frac{N\mu}{N} = \mu,$$

and since the X_i are independent,

$$\text{var}(M_N) = \frac{\text{var}(X_1) + \text{var}(X_2) + \cdots + \text{var}(X_N)}{N^2} = \frac{N\sigma^2}{N^2} = \frac{\sigma^2}{N}.$$

For any $\epsilon > 0$, a direct application of the Chebyshev inequality tells us that

$$P(|M_N - \mu| > \epsilon) \leq \frac{\sigma^2}{N\epsilon^2}.$$

The point is that this gets arbitrarily small as $N \rightarrow \infty$ no matter what ϵ was chosen to be. We have established, in some sense, that even though $\{M_N\}$ is a sequence of random numbers, it converges to something deterministic, namely μ .

WLLN: Let X_1, X_2, \dots be iid random variables as above. For **every** $\epsilon > 0$, we have

$$P(|M_N - \mu| > \epsilon) = P\left(\left|\frac{X_1 + \dots + X_N}{N} - \mu\right| > \epsilon\right) \longrightarrow 0,$$

as $N \rightarrow \infty$.

One of the philosophical consequences of the WLLN is that it tells us that probabilities can be estimated through **empirical frequencies**. Suppose I want to estimate the probability of an event A occurring related to some probabilistic experiment. We run a series of (independent) experiments, and set $X_i = 1$ if A occurred in experiment i , and $X_i = 0$ otherwise. Then given X_1, \dots, X_N , we estimate the probability of A in a completely reasonable way, by computing the percentage of times it occurred:

$$p_{\text{empirical}} = \frac{X_1 + \dots + X_N}{N}.$$

The WLLN tells us that

$$p_{\text{empirical}} \rightarrow P(A), \quad \text{as } N \rightarrow \infty.$$

This lends some mathematical weight to our interpretation of probabilities as *relative frequencies*.

All of the above of course applies to functions of random variables. That is, if X is a random variable, and $g(X)$ is a function of that random variable with

$$\text{var}(g(X)) = E[(g(X) - E[g(X)])^2] < \infty,$$

then given independent realizations X_1, \dots, X_N , we have

$$\frac{1}{N} \sum_{n=1}^N g(X_n) \rightarrow \mathbb{E}[g(X)]$$

as $N \rightarrow \infty$.