

**Lecture Note 12: Point Estimation: Method of Moments and Maximum Likelihood Estimation (CB 7.1, 7.2.1, 7.2.2)**

**1. Point Estimation Problem**

In statistical inference, we consider a set of possible probability models, and try to use observations to infer which of the set of probability models is the one that generated the data. To formalize this, suppose we have a parametric family of PDFs or PMFs

$$\{f(x; \theta), \theta \in \Theta\}$$

We refer to  $\theta$  as the parameter and  $\Theta$  as the parameter space. You could think of  $\theta$  as a particular theory for some phenomenon.

Suppose we have a random sample of size  $n$  from a distribution with PDF or PMF  $f(x; \theta^*)$ . That is, we observe  $X_1, X_2, \dots, X_n$ , where each  $X_i$  is IID with PDF or PMF  $f(x; \theta^*)$ . Here,  $\theta^* \in \Theta$  is the “true value” of the parameter, but we do not know what it is. Our goal is to use the observations to provide an estimate of  $\theta^*$ .

A statistic is any function of the observations, say  $T(X_1, \dots, X_n)$ .

A point estimator is a statistic used to provide a guess about  $\theta$ .

Often, we will use notation like  $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$  to denote a point estimator. For a given realization of the data, it is just a number, but if we were to take a new sample of data, it would take on a different value. In this sense a statistic is itself a random variable, and we will evaluate it according to the repeated sampling criteria: we look at the behavior (distribution) of the estimator when we repeatedly get new random samples of the same size. We would like the distribution of the estimator to be concentrated around  $\theta^*$ .

Note: Be careful about the notation here. We use  $\theta^*$  here to denote the true value of the parameter that generated the data, and use  $\theta$  to denote any element of the parameter space. This distinction will be important as we will evaluate the density function sometimes at the true value of the parameter, and sometimes at arbitrary values.

**2. Method of Moments**

The first approach to systematically find estimators is the method of moments. Consider a set of independent and identically distributed random variables with a PDF/PMF  $f(x; \theta)$ . Suppose  $\theta$  is a scalar. The mean of this distribution is

$$E[X] = g(\theta^*),$$

where the function  $g(\theta)$  (for all values of  $\theta$ , and not just for the true value  $\theta^*$ ) is defined as:

$$g(\theta) = \int x \cdot f(x; \theta) dx.$$

The function  $g(\cdot)$  is clearly a known function given knowledge of  $f(\cdot)$ . Now suppose that we calculate the average of our random sample,  $\bar{X} = \sum X_i/n$ . (For notational simplicity, I

have dropped the “ $n$ ” subscript.) With  $n$  reasonably large this average should be close to  $g(\theta^*)$ . Therefore the value of  $\theta$  that solves

$$g(\theta) = \bar{X},$$

would appear to be a sensible estimator for  $\theta^*$ .

Let us look at some examples. Suppose that  $X_i$  has a Bernoulli distribution with probability  $p^*$ . Then the expectation of  $X$  is  $p^*$ , so the method of moment estimator solves

$$\hat{p} = \bar{X}.$$

This is an unbiased estimator: its expectation is equal to the unknown parameter.

Suppose that  $X$  has an exponential distribution with arrival rate  $\lambda$ . The expectation of  $X$  is  $1/\lambda$ , so the method of moments estimator solves:

$$g(\hat{\lambda}) = 1/\hat{\lambda} = \bar{X},$$

or

$$\hat{\lambda} = 1/\bar{X}.$$

This estimator is not unbiased: the expectation of  $X$  is  $1/\lambda$ , so the expectation of  $1/X$  is greater than  $\lambda$  by Jensen's inequality.

**Result 1** (JENSEN'S INEQUALITY, CB 4.7.7)

*For any random variable  $X$ , if  $g(x)$  is a convex function, then*

$$Eg(X) \geq g(EX).$$

*Equality holds if and only if, for every line  $a + bx$  that is tangent to  $g(x)$  at  $x = EX$ ,  $P(g(X) = a + bX) = 1$ .*

**Proof:** Let  $l(x)$  be a tangent line to  $g(x)$  at the point  $g(EX)$ . Write  $l(x) = a + bx$  for some  $a$  and  $b$ . By convexity of  $g$ , we have  $g(x) \geq a + bx$ . Therefore

$$\begin{aligned} Eg(X) &\geq E(a + bX) \\ &= a + bE(X) \\ &= l(EX) \\ &= g(EX) \end{aligned}$$

(Proof of the last part of the result left to the reader.)  $\square$

Now suppose we have two parameters, that is,  $\theta$  is a two dimensional vector. We could calculate the first two moments:

$$\begin{aligned} g_1(\theta) &= E[X] = \int x \cdot f_X(x; \theta) dx, \\ g_2(\theta) &= E[X^2] = \int x^2 \cdot f_X(x; \theta) dx, \end{aligned}$$

and equate them with the corresponding sample moments:

$$g_1(\hat{\theta}) = \bar{X},$$

$$g_2(\hat{\theta}) = \overline{X^2}.$$

For example, suppose that  $X$  has a normal distribution with mean  $\mu$  and variance  $\sigma^2$ . The  $g(\cdot)$  functions are:

$$g_1(\mu, \sigma^2) = \int x f_X(x; \mu, \sigma^2) dx = \mu,$$

$$g_2(\mu, \sigma^2) = \int x^2 f_X(x; \mu, \sigma^2) dx = \mu^2 + \sigma^2.$$

Hence the method of moments estimators are

$$\hat{\mu} = \bar{X},$$

$$\hat{\sigma}^2 = \overline{X^2} - (\bar{X})^2.$$

Again these appear reasonable and intuitive estimators. In other cases the estimators are somewhat less obvious and less attractive. Consider the case where the distribution of  $X$  is multinomial with parameters  $k$  and  $p$  both unknown (a relatively unusual problem in practice). The first two moments are:

$$g_1(k, p) = \sum x f_X(x; k, p) = k \cdot p,$$

$$g_2(k, p) = \sum x^2 f_X(x; k, p) = k \cdot p \cdot (1 - p) + k^2 \cdot p^2.$$

Thus we solve

$$\bar{X} = \hat{k} \cdot \hat{p},$$

$$\overline{X^2} = \hat{k} \cdot \hat{p} \cdot (1 - \hat{p}) + \hat{k}^2 \cdot \hat{p}^2.$$

The solutions are:

$$\hat{p} = \bar{X} / \hat{k},$$

and

$$\hat{k} = \frac{\bar{X} \cdot \overline{X^2}}{\bar{X} \cdot \bar{X} + \bar{X} - \overline{X^2} + 3}.$$

A problem is that these estimates can in fact be negative if the variance is larger than the mean, which is obviously not such a good estimate for a probability and a number of trials.

### 3. Maximum Likelihood

The second general approach to estimation we consider is maximum likelihood estimation. The likelihood function is the density function viewed as a function of the unknown parameters, rather than as a function of the random variable. Let  $X$  be a random variable with PDF/PMF

$$f_X(x; \theta^*),$$

where  $\theta^*$  is the unknown true value of the parameter  $\theta$ . The likelihood function is then:

$$\mathcal{L}(\theta) = f_X(X; \theta).$$

If we have more than one random variable, say  $X_1, X_2, \dots, X_n$ , the likelihood function is based on the joint probability density/mass function:

$$\mathcal{L}(\theta) = f_{X_1, X_2, \dots, X_n}(X_1, X_2, \dots, X_n; \theta).$$

If the random variables are independent and identically distributed, with common density function  $f_X(x; \theta)$ , the likelihood function obviously simplifies:

$$\mathcal{L}(\theta) = f_{X_1, X_2, \dots, X_n}(X_1, X_2, \dots, X_n; \theta) = \prod_{i=1}^n f_X(X_i; \theta).$$

Often we prefer to work with the logarithm of the likelihood function, the log likelihood function, e.g., in the case of  $n$  independent and identically distributed random variables:

$$L(\theta) = \ln \mathcal{L}(\theta) = \ln f_{X_1, X_2, \dots, X_n}(X_1, X_2, \dots, X_n; \theta) = \sum_{i=1}^n \ln f_X(X_i; \theta).$$

The maximum likelihood estimator or mle for  $\theta^*$  is the value of  $\theta$  that maximizes the likelihood function, or equivalently, the log likelihood function (because the logarithm is a one to one, strictly monotone, transformation the maximand of one is the maximand of the other.)

Let us first look at a simple example. Suppose  $X$  has a normal distribution with unknown mean  $\mu$  and known variance 1. The probability density function is

$$f_X(x; \mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x - \mu)^2\right).$$

The likelihood function is

$$\mathcal{L}(\mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(X - \mu)^2\right).$$

The log likelihood function is

$$L(\mu) = \ln \mathcal{L}(\mu) = -\frac{1}{2} \ln(2\pi) - \frac{1}{2}(X - \mu)^2.$$

The value of  $\mu$  that maximizes the log likelihood function is  $\hat{\mu}_{mle} = X$ .

Now suppose we have  $n$  random variables independent and all normally distributed with mean  $\mu$  and variance 1. In that case the likelihood function is

$$\mathcal{L}(\mu) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(X_i - \mu)^2\right),$$

and the log likelihood function is

$$L(\mu) = -n\frac{1}{2} \ln(2\pi) - \sum_{i=1}^n \frac{1}{2}(X_i - \mu)^2.$$

The maximum likelihood estimator is  $\hat{\mu}_{mle} = \bar{X}$ , the sample average.

In both these cases the resulting estimator appears very reasonable, similar to the method of moments estimators for this case. In general, however, it is not obvious why maximizing the likelihood function is a sensible strategy. One motivation, which we will discuss later, is that the method of maximum likelihood is related to another method, the Bayesian method, that can be justified on decision-theoretic grounds. Another motivation is more direct and relies on large sample approximations.

Consider the random variable  $Y$  defined as the ratio of the density function at some arbitrary value of  $\theta$  to the density function at  $\theta^*$ , both evaluated at the random variable  $X$ :

$$Y = f_X(X; \theta) / f_X(X; \theta^*).$$

Take  $g(\cdot)$  to be minus the logarithmic function:  $g(a) = -\ln(a)$ , so  $g'(a) = -1/a$ , and  $g''(a) = 1/a^2 > 0$  and  $g(\cdot)$  is convex. Then, by Jensen's inequality

$$E[-\ln Y] \geq -\ln E[Y],$$

implying

$$E\left[-\ln\left(\frac{f_X(X; \theta)}{f_X(X; \theta^*)}\right)\right] \geq -\ln\left(E\left[\frac{f_X(X; \theta)}{f_X(X; \theta^*)}\right]\right),$$

where the expectation is over the distribution of  $X$ , that is the density  $f_X(x; \theta^*)$ . The expectation on the right therefore simplifies:

$$E\left[\frac{f_X(X; \theta)}{f_X(X; \theta^*)}\right] = \int \frac{f_X(x; \theta)}{f_X(x; \theta^*)} \cdot f_X(x; \theta^*) dx = \int f_X(x; \theta) dx = 1,$$

for all values of  $\theta$ , so that after taking the log, the righthand side is equal to zero, and thus

$$E\left[-\ln\left(\frac{f_X(X; \theta)}{f_X(X; \theta^*)}\right)\right] \geq 0$$

implying

$$-E\left[\ln f_X(X; \theta)\right] + E\left[\ln f_X(X; \theta^*)\right] \geq 0,$$

and thus

$$E\left[\ln f_X(X; \theta^*)\right] \geq E\left[\ln f_X(X; \theta)\right],$$

for all  $\theta$ . This implies that the expected value of the log likelihood is maximized at the true value of  $\theta$ , and therefore there is some hope that the actual log likelihood function is maximized at a value close to  $\theta^*$ , and therefore at a value that is a good estimate of  $\theta^*$ .

Note that we always take the expectation over the true distribution,  $f(x; \theta^*)$ , even if we are taking expectations of functions of  $\theta$ , evaluated at all possible values of  $\theta$ , and in particular evaluated at values other than the true value  $\theta^*$ .

Let us look at an example that illustrates this argument. Let  $X_1, X_2, \dots, X_n$  be Bernoulli trials with unknown probability  $p^*$ . The joint density is

$$f_{X_1, X_2, \dots, X_n}(x_1, \dots, x_n; p) = p^{\sum x_i} (1-p)^{n-\sum x_i}.$$

The log likelihood function is

$$L(p) = \sum X_i \ln p + (n - \sum X_i) \cdot \ln(1-p).$$

The value that maximizes the log likelihood function is the sample average  $\hat{p}_{mle} = \bar{X}$ . Now consider the expected value of the log likelihood function:

$$\begin{aligned} E[L(p)] &= E \sum X_i \ln p + (n - \sum X_i) \cdot \ln(1 - p), \\ &= n \cdot p^* \cdot \ln(p) + n \cdot (1 - p^*) \cdot \ln(1 - p). \end{aligned}$$

This expected log likelihood function is maximized at  $p = p^*$ , which is the true value. So, while the actual log likelihood function is not necessarily maximized at  $p^*$ , its expectation is. If we multiply the log likelihood function by  $1/n$  it is a sample average and as such subject to law of large numbers. (multiplying the log likelihood function by a constant does not affect the location of the maximum.) Hence its expectation may be reasonably close to the average, and that is one motivation for considering maximum likelihood estimators.

Let us consider two more examples. Suppose that  $X_1, X_2, \dots, X_n$  are iid with uniform distribution on the interval zero to  $\theta^*$ . The likelihood function is

$$\mathcal{L}(\theta) = (1/\theta^n) \cdot 1_{\{\max\{x_i\} \leq \theta\}}.$$

The likelihood function is strictly decreasing for values where it differs from zero. The maximum likelihood estimator is equal to  $\max_i\{X_i\}$ .

Suppose  $X_1, \dots, X_n$  are normal with unknown mean  $\mu$  and unknown variance  $\sigma^2$ . The likelihood function is

$$\mathcal{L}(\mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(X_i - \mu)^2\right),$$

with logarithm

$$L(\mu, \sigma^2) = \sum_{i=1}^n -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2}(X_i - \mu)^2.$$

To maximize this over  $\mu$  and  $\sigma^2$  jointly we try to solve the first order conditions:

$$\frac{\partial L}{\partial \mu}(\mu, \sigma^2) = \sum_{i=1}^n \frac{1}{\sigma^2}(X_i - \mu),$$

$$\frac{\partial L}{\partial \sigma^2}(\mu, \sigma^2) = \sum_{i=1}^n -\frac{1}{2\sigma^2} + \frac{1}{2\sigma^4}(X_i - \mu)^2.$$

Setting both derivatives equal to zero leads to

$$\hat{\sigma}_{mle}^2 = \sum_{i=1}^n (X_i - \bar{X})^2/n,$$

and, as before

$$\hat{\mu}_{mle} = \bar{X}.$$

In principle we also have to check that the matrix of second derivatives is negative definite to ensure that we have located a maximum and not a minimum or saddle point. In this case the matrix of second derivatives is indeed negative definite and this is the global maximum.

Finally, an important property of maximum likelihood estimators is their invariance under reparametrization: If  $\hat{\theta}_{mle}$  is the maximum likelihood estimator for  $\theta^*$ , then  $\hat{\pi}_{mle} = g(\hat{\theta}_{mle})$  is the maximum likelihood estimator for any one-to-one transformation  $\pi^* = g(\theta^*)$ .

For example of this consider the exponential distribution. Let  $X_1, X_2, \dots, X_n$  be iid with an exponential distribution with arrival rate  $\lambda$ . The log likelihood function is:

$$L(\lambda) = n \cdot \ln \lambda - \sum X_i \cdot \lambda.$$

The mle is  $\hat{\lambda}_{mle} = 1/\bar{X}$ . Now suppose we are interested in the mle for the mean of this distribution, and therefore reparametrize the distribution in terms of  $\mu = 1/\lambda$ . we can go back to the likelihood function:

$$L(\mu) = -n \cdot \ln \mu - \sum X_i/\mu,$$

and we end up with the maximum likelihood estimator.

$$\hat{\mu}_{mle} = \bar{X}.$$

The invariance property implies the same result:  $\hat{\mu}_{mle} = 1/\hat{\lambda}_{mle} = (1/(1/\bar{X})) = \bar{X}$ .