# Large Sample Theory

Large Sample Theory is a name given to the search for approximations to
the behaviour of statistical procedures which are derived by computing limits
as the sample size, $n$, tends to infinity. Suppose we have a data set with a
fairly large sample size, say $n = 100$. We imagine our data set is one in a
sequence of possible data sets — one for each possible value of $n$. If we have
a sequence of statistics which converges to something as $n$ tends to infinity
we approximate the value of a probability, for instance, when $n = 100$ by
the corresponding value when $n = \infty$. In this section we will investigate
approximations of this kind for Maximimum Likelihood Estimation.

   Most large sample theory uses three main technical tools: the Law of
Large Numbers (LLN), the Central Limit Theorem (CLT) and Taylor ex-
pansion. I assume you have heard of all of these but will state versions of
them as we go.

   These tools are generally easier to apply to statistics for which we have
explicit formulas than to statistics, like maximum likelihood estimates, where
we do not usually have explicit formulas. In this situation we study the
equations you solve to find the MLEs instead.

   We therefore will study the approximate behaviour of the MLE, $\hat{\theta}$, by
studying the function $U$. Notice first that $U$ is a sum of independent random
variables. This will allow us to apply both the LLN and the CLT to $U$.

**Theorem 1 (LLN)** *If $Y_1, Y_2, \ldots$ are iid with mean $\mu$ then*

$$\frac{\sum Y_i}{n} \to \mu$$

   This is called the law of large numbers but it comes in two forms: Strong
and Weak.

**Theorem 2 (SLLN)** *If $Y_1, Y_2, \ldots$ are iid with mean $\mu$ then*

$$P\left(\lim_{n\to\infty \frac{\sum Y_i}{n}=\mu)=1.}\right.$$

   The strong law is harder to prove than the weak law of large numbers:

1

**Theorem 3 (WLLN)** *If $Y_1, Y_2, \ldots$ are iid with mean $\mu$ then for each positive $\epsilon$*

$$\lim P(|\frac{\sum Y_i}{n} - \mu| > \epsilon) = 0$$

For iid $Y_i$ the stronger conclusion (the SLLN) holds but for our heuristics we will ignore the differences between these notions.

Now suppose that $\theta_0$ is the true value of $\theta$. Then

$$U(\theta)/n \to \mu(\theta)$$

where

$$\mu(\theta) = E_{\theta_0}\left[\frac{\partial \log f}{\partial \theta}(X_i, \theta)\right]$$
$$= \int \frac{\partial \log f}{\partial \theta}(x, \theta) f(x, \theta_0) dx$$

**Remark**: This convergence is *pointwise* and might not be *uniform*. That is, there might be some $\theta$ values where the convergence takes much longer than others. It could be that for every $n$ there is a $\theta$ where $U(\theta)/n$ and $\mu(\theta)$ are not close together.

Consider as an example the case of $N(\mu, 1)$ data where

$$U(\mu)/n = \sum(X_i - \mu)/n = \bar{X} - \mu$$

If the true mean is $\mu_0$ then $\bar{X} \to \mu_0$ and

$$U(\mu)/n \to \mu_0 - \mu$$

If we think of a $\mu < \mu_0$ we see that the derivative of $\ell(\mu)$ is likely to be positive so that $\ell$ increases as we increase $\mu$. For $\mu$ more than $\mu_0$ the derivative is probably negative and so $\ell$ tends to be decreasing for $\mu > 0$. It follows that $\ell$ is likely to be maximized close to $\mu_0$.

Now we repeat these ideas for a more general case. We study the random variable $\log[f(X_i, \theta)/f(X_i, \theta_0)]$. You know the inequality

$$E(X)^2 \leq E(X^2)$$

(because the difference is $Var(X) \geq 0$. This inequality has the following generalization, called Jensen's inequality. If $g$ is a convex function (non-negative second derivative, roughly) then

$$g(E(x)) \leq E(g(X))$$

The inequality above has $g(x) = x^2$. We use $g(x) = -\log(x)$ which is convex because $g''(x) = x^{-2} > 0$. We get

$$-\log(E_{\theta_0}[f(X_i, \theta)/f(X_i, \theta_0)] \leq E_{\theta_0}[-\log\{f(X_i, \theta)/f(X_i, \theta_0)\}]$$

But

$$
\begin{aligned}
E_{\theta_0}[f(X_i, \theta)/f(X_i, \theta_0)] &= \int \frac{f(x, \theta)}{f(x, \theta_0)} f(x, \theta_0) dx \\
&= \int f(x, \theta) dx \\
&= 1
\end{aligned}
$$

We can reassemble the inequality and this calculation to get

$$E_{\theta_0}[\log\{f(X_i, \theta)/f(X_i, \theta_0)\}] \leq 0$$

It is possible to prove that the inequality is strict unless the $\theta$ and $\theta_0$ densities are actually the same. Let $\mu(\theta) < 0$ be this expected value. Then for each $\theta$ we find

$$
\begin{aligned}
n^{-1}[\ell(\theta) - \ell(\theta_0)] &= n^{-1} \sum \log[f(X_i, \theta)/f(X_i, \theta_0)] \\
&\to \mu(\theta)
\end{aligned}
$$

This proves that the likelihood is probably higher at $\theta_0$ than at any other single $\theta$. This idea can often be stretched to prove that the MLE is **consistent**.

**Definition**: A sequence $\hat{\theta}_n$ of estimators of $\theta$ is consistent if $\hat{\theta}_n$ converges weakly (or strongly) to $\theta$.

**Proto theorem**: In regular problems the MLE $\hat{\theta}$ is consistent.

Now let us study the shape of the log likelihood near the true value of $\hat{\theta}$ under the assumption that $\hat{\theta}$ is a root of the likelihood equations close to $\theta_0$. We use Taylor expansion to write, for a 1 dimensional parameter $\theta$,

$$
\begin{aligned}
U(\hat{\theta}) &= 0 \\
&= U(\theta_0) + U'(\theta_0)(\hat{\theta} - \theta_0) + U''(\tilde{\theta})(\hat{\theta} - \theta_0)^2/2
\end{aligned}
$$

for some $\tilde{\theta}$ between $\theta_0$ and $\hat{\theta}$. (This form of the remainder in Taylor's theorem is not valid for multivariate $\theta$.) The derivatives of $U$ are each sums of $n$ terms and so should be both proportional to $n$ in size. The second derivative is multiplied by the square of the small number $\hat{\theta} - \theta_0$ so should be negligible compared to the first derivative term. If we ignore the second derivative term we get

$$-U'(\theta_0)(\hat{\theta} - \theta_0) \approx U(\theta_0)$$

Now let's look at the terms $U$ and $U'$.

In the normal case

$$U(\theta_0) = \sum (X_i - \mu_0)$$

has a normal distribution with mean 0 and variance $n$ (SD $\sqrt{n}$). The derivative is simply

$$U'(\mu) = -n$$

and the next derivative $U''$ is 0. We will analyze the general case by noticing that both $U$ and $U'$ are sums of iid random variables. Let

$$U_i = \frac{\partial \log f}{\partial \theta}(X_i, \theta_0)$$

and

$$V_i = -\frac{\partial^2 \log f}{\partial \theta^2}(X_i, \theta)$$

In general, $U(\theta_0) = \sum U_i$ has mean 0 and approximately a normal distribution. Here is how we check that:

$$
\begin{aligned}
E_{\theta_0}(U(\theta_0)) &= nE_{\theta_0}(U_1) \\
&= n \int \frac{\partial \log(f(x,\theta))}{\partial \theta}(x, \theta_0) f(x, \theta_0) dx \\
&= n \int \frac{\partial f/\partial \theta(x, \theta_0)}{f(x, \theta_0)} \theta f(x, \theta_0) dx \\
&= n \int \frac{\partial f}{\partial \theta}(x, \theta_0) dx \\
&= n \frac{\partial}{\partial \theta} \int f(x, \theta) dx \bigg|_{\theta = \theta_0} \\
&= n \frac{\partial}{\partial \theta} 1 \\
&= 0
\end{aligned}
$$

4

Notice that I have interchanged the order of differentiation and integration at one point. This step is usually justified by applying the dominated convergence theorem to the definition of the derivative. The same tactic can be applied by differentiating the identity which we just proved

$$\int \frac{\partial \log f}{\partial \theta}(x, \theta) f(x, \theta) dx = 0$$

Taking the derivative of both sides with respect to $\theta$ and pulling the derivative under the integral sign again gives

$$\int \frac{\partial}{\partial \theta} \left[ \frac{\partial \log f}{\partial \theta}(x, \theta) f(x, \theta) \right] dx = 0$$

Do the derivative and get

$$-\int \frac{\partial^2 \log(f)}{\partial \theta^2} f(x, \theta) dx = \int \frac{\partial \log f}{\partial \theta}(x, \theta) \frac{\partial f}{\partial \theta}(x, \theta) dx$$

$$= \int \left[ \frac{\partial \log f}{\partial \theta}(x, \theta) \right]^2 f(x, \theta) dx$$

**Definition**: The **Fisher Information** is

$$I(\theta) = -E_\theta(U'(\theta)) = n E_{\theta_0}(V_1)$$

We refer to $\mathcal{I}(\theta_0) = E_{\theta_0}(V_1)$ as the information in 1 observation.

The idea is that $I$ is a measure of how curved the log likelihood tends to be at the true value of $\theta$. Big curvature means precise estimates. Our identity above is

$$I(\theta) = Var_\theta(U(\theta)) = n\mathcal{I}(\theta)$$

Now we return to our Taylor expansion approximation

$$-U'(\theta_0)(\hat{\theta} - \theta_0) \approx U(\theta_0)$$

and study the two appearances of $U$.

We have shown that $U = \sum U_i$ is a sum of iid mean 0 random variables. The central limit theorem thus proves that

$$n^{-1/2} U(\theta) \Rightarrow N(0, \sigma^2)$$

5

where $\sigma^2 = Var(U_i) = E(V_i) = \mathcal{I}(\theta)$.

Next observe that

$$-U'(\theta) = \sum V_i$$

where again

$$V_i = -\frac{\partial U_i}{\partial \theta}$$

The law of large numbers can be applied to show

$$-U'(\theta_0)/n \to E_{\theta_0}[V_1] = \mathcal{I}(\theta_0)$$

Now manipulate our Taylor expansion as follows

$$n^{1/2}(\hat{\theta} - \theta_0) \approx \left[ \frac{\sum V_i}{n} \right]^{-1} \frac{\sum U_i}{\sqrt{n}}$$

Apply Slutsky's Theorem to conclude that the right hand side of this converges in distribution to $N(0, \sigma^2/\mathcal{I}(\theta)^2)$ which simplifies, because of the identities, to $N(0, 1/\mathcal{I}(\theta))$.

**Summary**

In regular families:

- Under strong regularity conditions Jensen's inequality can be used to demonstrate that $\hat{\theta}$ which maximizes $\ell$ globally is consistent and that this $\hat{\theta}$ is a root of the likelihood equations.

- It is generally easier to study $\ell$ only close to $\theta_0$. For instance define $A$ to be the event that $\ell$ is concave on the set of $\theta$ such that $|\theta - \theta_0| < \delta$ and the likelihood equations have a unique root in that set. Under weaker conditions than the previous case we can prove that there is a $\delta > 0$ such that

$$P(A) \to 1$$

  In that case we can prove that the root $\hat{\theta}$ of the likelihood equations mentioned in the definition of $A$ is consistent.

- Sometimes we can only get an even weaker conclusion. Define $B$ to be the event that $\ell(\theta)$ is concave for $n^{1/2}|\theta - \theta_0| < L$ and there is a unique root of $\ell$ over this range. Again this root is consistent but there might be other consistent roots of the likelihood equations.

6

- Under any of these scenarios there is a consistent root of the likelihood equations which is definitely the closest to the true value $\theta_0$. This root $\hat\theta$ has the property

$$\sqrt{n}(\hat\theta - \theta_0) \Rightarrow N(0, 1/\mathcal{I}(\theta)).$$

We usually simply say that the MLE is consistent and asymptotically normal with an asymptotic variance which is the inverse of the Fisher information. This assertion is actually valid for vector valued $\theta$ where now $I$ is a matrix with $ij$th entry

$$I_i j = -E\left(\frac{\partial^2 \ell}{\partial\theta_i\partial\theta_j}\right)$$

**Estimating Equations**

The same ideas arise in almost any model where estimates are derived by solving some equation. As an example I sketch large sample theory for **Generalized Linear Models**.

Suppose that for $i = 1, \ldots, n$ we have observations of the numbers of cancer cases $Y_i$ in some group of people characterized by values $x_i$ of some covariates. You are supposed to think of $x_i$ as containing variables like age, or a dummy for sex or average income or … A parametric regression model for the $Y_i$ might postulate that $Y_i$ has a Poisson distribution with mean $\mu_i$ where the mean $\mu_i$ depends somehow on the covariate values. Typically we might assume that $g(\mu_i) = \beta 0 + x_i\beta$ where $g$ is a so-called **link** function, often for this case $g(\mu) = \log(\mu)$ and $x_i\beta$ is a matrix product with $x_i$ written as a row vector and $\mu$ a column vector. This is supposed to function as a "linear regression model with Poisson errors". I will do as a special case $\log(\mu_i) = \beta x_i$ where $x_i$ is a scalar.

The log likelihood is simply

$$\ell(\beta) = \sum(Y_i \log(\mu_i) - \mu_i)$$

ignoring irrelevant factorials. The score function is, since $\log(\mu_i) = \beta x_i$,

$$U(\beta) = \sum(Y_i x_i - x_i\mu_i) = \sum x_i(Y_i - \mu_i)$$

(Notice again that the score has mean 0 when you plug in the true parameter value.) The key observation, however, is that it is not necessary to believe

that $Y_i$ has a Poisson distribution to make solving the equation $U = 0$ sensible. Suppose only that $\log(E(Y_i)) = x_i\beta$. Then we have assumed that

$$E_\beta(U(\beta)) = 0$$

This was the key condition in proving that there was a root of the likelihood equations which was consistent and here it is what is needed, roughly, to prove that the equation $U(\beta) = 0$ has a consistent root $\hat{\beta}$. Ignoring higher order terms in a Taylor expansion will give

$$V(\beta)(\hat{\beta} - \beta) \approx U(\beta)$$

where $V = -U'$. In the MLE case we had identities relating the expectation of $V$ to the variance of $U$. In general here we have

$$Var(U) = \sum x_i^2 Var(Y_i)$$

If $Y_i$ is Poisson with mean $\mu_i$ (and so $Var(Y_i) = \mu_i$) this is

$$Var(U) = \sum x_i^2 \mu_i$$

Moreover we have

$$V_i = x_i^2 \mu_i$$

and so

$$V(\beta) = \sum x_i^2 \mu_i$$

The central limit theorem (the Lyapunov kind) will show that $U(\beta)$ has an approximate normal distribution with variance $\sigma_U^2 = \sum x_i^2 Var(Y_i)$ and so

$$\hat{\beta} - \beta \approx N(0, \sigma_U^2/(\sum x_i^2 \mu_i)^2)$$

If $Var(Y_i) = \mu_i$, as it is for the Poisson case, the asymptotic variance simplifies to $1/\sum x_i^2 \mu_i$.

Notice that other estimating equations are possible. People suggest alternatives very often. If $w_i$ is any set of deterministic weights (even possibly depending on $\mu_i$ then we could define

$$U(\beta) = \sum w_i(Y_i - \mu_i)$$

and still conclude that $U = 0$ probably has a consistent root which has an asymptotic normal distribution. This idea is being used all over the place these days: see, for example Zeger and Liang's Generalized estimating equations (GEE) which the econometricians call Generalized Method of Moments.

## Problems with maximum likelihood

1. In problems with many parameters the approximations don't work very well and maximum likelihood estimators can be far from the right answer. See your homework for the Neyman Scott example where the MLE is not consistent.

2. When there are multiple roots of the likelihood equation you must choose the right root. To do so you might start with a different consistent estimator and then apply some iterative scheme like Newton Raphson to the likelihood equations to find the MLE. It turns out not many steps of NR are generally required if the starting point is a reasonable estimate.

## Finding (good) preliminary Point Estimates

### Method of Moments

Basic strategy: set sample moments equal to population moments and solve for the parameters.

**Definition**: The $r^{\text{th}}$ sample moment (about the origin) is

$$\frac{1}{n} \sum_{i=1}^{n} X_i^r$$

The $r^{\text{th}}$ population moment is

$$\mathrm{E}(X^r)$$

**Definition**: **Central** moments are

$$\frac{1}{n} \sum_{i=1}^{n} (X_i - \bar{X})^r$$

and

$$\mathrm{E}\left[(X - \mu)^r\right] .$$

If we have $p$ parameters we can estimate the parameters $\theta_1, \ldots, \theta_p$ by solving the system of $p$ equations:

$$\mu_1 = \bar{X}$$

9

$$\mu_2' = \overline{X^2}$$

and so on to

$$\mu_p' = \overline{X^p}$$

You need to remember that the population moments $\mu_k'$ will be formulas involving the parameters.

**Gamma Example**

The Gamma$(\alpha, \beta)$ density is

$$f(x; \alpha, \beta) = \frac{1}{\beta \Gamma(\alpha)} \left(\frac{x}{\beta}\right)^{\alpha - 1} \exp\left[-\frac{x}{\beta}\right] 1(x > 0)$$

and has

$$\mu_1 = \alpha\beta$$

and

$$\mu_2' = \alpha\beta^2$$

This gives the equations

$$\alpha\beta = \overline{X}$$
$$\alpha\beta^2 = \overline{X^2}$$

Divide the second by the first to find the method of moments estimate of $\beta$ is

$$\tilde{\beta} = \overline{X^2}/\overline{X}$$

Then from the first equation get

$$\tilde{\alpha} = \overline{X}/\tilde{\beta} = (\overline{X})^2/\overline{X^2}$$

The equations are much easier to solve than the likelihood equations which involve the function

$$\psi(\alpha) = \frac{d}{d\alpha} \log(\Gamma(\alpha))$$

called the digamma function.