Advanced Econometrics I

Jürgen Meinecke Lecture 10 of 12

Research School of Economics, Australian National University

Maximum Likelihood (ML) Estimation

Maximum Likelihood Estimator

Invariance Property

Fisher Information, Cramér Rao Bound, Information Equality

Minimum Variance Unbiased Estimators

Asymptotic Efficiency of ML

Let's take a step back into the univariate world

You have a random sample Y_1, \ldots, Y_N that is generated from either a

- probability mass function (discrete) or
- probability density function (continuous)

For simplicity I will only write pdf (and omit pmf), and denote it by $f(y|\theta)$

Crucial here

- \cdot f is known
- + $\boldsymbol{\theta}$ is unknown and we want to estimate it

Example: f is the pdf of a normal distribution with variance 1 and we are after the expected value θ

This foreshadows one of the biggest drawbacks of ML estimation: you need to know \boldsymbol{f}

Definition (Likelihood Functions)

Given a random sample Y_1, \ldots, Y_N where each Y_i has $pdf f(y|\theta)$

 \cdot likelihood function

 $L(\boldsymbol{\theta}) := f_{Y_1,\dots,Y_N}(y_1,\dots,y_N|\boldsymbol{\theta}) = \prod_{i=1}^N f(y_i|\boldsymbol{\theta})$

• log likelihood function $\ln L(\theta) := \ln(L(\theta)) = \sum_{i=1}^{N} \ln f(y_i|\theta)$

We view both as function in the parameter θ ; given a random sample the N different values y_i are known

Most generally, the likelihood function is the joint pdf, but here it reduces to the product of pdf's because of iid

The log likelihood function therefore is a sum

Definition (Maximum Likelihood (ML) Estimator)

Given a random sample $Y_1, ..., Y_N$ where each Y_i has $pdf f(y|\theta)$ the **maximum likelihood estimator** of θ is defined by

 $\hat{\theta}^{\mathsf{ML}} := \operatorname{argmax}_{\theta} L(\theta).$

Corollary $\operatorname{argmax}_{\theta} L(\theta) = \operatorname{argmax}_{\theta} \ln L(\theta).$

Because the log is strictly monotone

It is usually easier to maximize $\ln L$ analytically and numerically

Example: exponential distribution

Example: Y_i is exponentially distributed

Recall, pdf of
$$Y_i$$
 is $f(y|\mu) = (1/\mu) \exp(-y/\mu)$

 $1/\mu$ is called the arrival rate

You may remember that $EY_i = \mu$ and $Var Y_i = \mu^2$

For a sample of size N = 1, the likelihood functions are $L(\mu) = (1/\mu) \exp(-y_1/\mu)$ $\ln L(\mu) = -(\ln \mu + y_1/\mu)$

The ML estimator is $\hat{\mu}^{ML} := \operatorname{argmax}_{\mu} - (\ln \mu + y_1/\mu)$ First derivative is $-((1/\mu) - y_1/\mu^2)$

Setting equal to zero results in $\hat{\mu}^{\rm ML} = y_1$

(The second order condition could also be checked to confirm a local maximum)

What if your sample size is N > 1?

The likelihood functions are $L(\mu) = \mu^{-N} \exp(-\sum_{i=1}^{N} y_i/\mu)$ $\ln L(\mu) = -N \ln \mu - \sum_{i=1}^{N} y_i/\mu = -N \left(\ln \mu + \bar{Y}/\mu\right)$

The ML estimator is $\hat{\mu}^{ML} := \operatorname{argmax}_{\mu} - N\left(\ln \mu + \bar{Y}/\mu\right)$ First derivative is $-N\left((1/\mu) - \bar{Y}/\mu^2\right)$

Setting equal to zero results in $\hat{\mu}^{\rm ML} = \bar{Y}$

This may convince you that ML estimation results in sensible estimators

But why should it be a sensible estimation approach?

Consider a random variable Y generated according to $f(y|\theta)$ Take any $\tilde{\theta} \neq \theta$ and look at the likelihood ratio $LR = f(Y|\tilde{\theta})/f(Y|\theta)$ Here θ is the *true* parameter and $\tilde{\theta}$ is some other value By Jensen's inequality $E(-\ln LR) > -\ln E(LR)$

(because $-\ln(\cdot)$ is strictly convex)

Substituting back

$$\mathsf{E}\left(-\ln\frac{f(Y|\tilde{\theta})}{f(Y|\theta)}\right) > -\ln\mathsf{E}\left(\frac{f(Y|\tilde{\theta})}{f(Y|\theta)}\right)$$

Notice the expectations are taken over the true $pdf f(y|\theta)$ (and not $f(y|\tilde{\theta})$)

The expectation on the rhs simplifies

$$\mathsf{E}\left(\frac{f(Y|\tilde{\theta})}{f(Y|\theta)}\right) = \int \frac{f(y|\tilde{\theta})}{f(y|\theta)} f(y|\theta) dy = \int f(y|\tilde{\theta}) dy = 1$$

Therefore

$$E\left(-\ln\frac{f(Y|\tilde{\theta})}{f(Y|\theta)}\right) > -\ln E\left(\frac{f(Y|\tilde{\theta})}{f(Y|\theta)}\right) = -\ln 1 = 0$$

Implying

 $\mathsf{E}\left(\ln f(Y|\theta)\right) > \mathsf{E}\left(\ln f(Y|\tilde{\theta})\right)$

This means that the expected value of the log likelihood is maximized at the *true* value of the parameter

So there's some hope that the sample analog will mimic this feature

Writing the log likelihood using the generic placeholder $\tilde{\mu}$ $\ln L(\tilde{\mu}) = -N \left(\ln \tilde{\mu} + \bar{Y}/\tilde{\mu} \right)$ Consider its expected value $E \left(\ln L(\tilde{\mu}) \right) = -N \left(\ln \tilde{\mu} + \mu/\tilde{\mu} \right)$

By the analogy principle, this is maximized at $\tilde{\mu} = \mu$

Maximum Likelihood (ML) Estimation

Maximum Likelihood Estimator

Invariance Property

Fisher Information, Cramér Rao Bound, Information Equality Minimum Variance Unbiased Estimators Asymptotic Efficiency of ML Here's a simple looking result that is useful

Proposition (Invariance Property of ML estimators)

If $\hat{\theta}^{ML}$ is the ML estimator of θ , then for any function g, the ML estimator of $g(\theta)$ is $g(\hat{\theta}^{ML})$.

Note: I didn't specify the domain and range of g because that arises from the context of $\boldsymbol{\theta}$

Also, g doesn't have to be one-to-one

It's easiest to explain by example...

Let $Y_1, ..., Y_N$ be iid with $pdf f(y|\mu) = (1/\mu) \exp(-y/\mu)$ Let's take the case N > 1We saw earlier that $\hat{\mu}^{ML} = \bar{Y}$ Now, define $\lambda := 1/\mu$, which is called the *arrival rate* Rewrite the pdf: $f(y|\lambda) = \lambda \exp(-\lambda y)$ How would we estimate λ using ML? Easy, the invariance property suggests: $\hat{\lambda}^{ML} = 1/\bar{Y}$

Maximum Likelihood (ML) Estimation

Maximum Likelihood Estimator

Invariance Property

Fisher Information, Cramér Rao Bound, Information Equality

Minimum Variance Unbiased Estimators

Asymptotic Efficiency of ML

Definition (Score Function)

Let Y be a random variable with $pdf f(y|\theta)$. The score function is defined by

$$S(y|\theta) := \frac{\partial \ln f}{\partial \theta}(y|\theta)$$

The score plays an important role in ML theory

Obviously, by definition $(1/N) \sum_{i=1}^{N} S(y_i | \hat{\theta}^{ML}) = 0$ (necessary condition for a maximum)

This is interesting bc it is true that $E(S(Y_i|\theta)) = 0$ (see below)

The ML estimator therefore can be motivated as an analog estimator: it is the value for θ that makes the average score equal zero

Going back to the example of the exponential distribution $S(y|\mu) := -((1/\mu) - y/\mu^2)$ In the case N = 1 we get $S(y_1|\hat{\theta}^{ML}) = 0$ In the case N > 1 we get $(1/N) \sum_{i=1}^{N} S(y_i|\hat{\theta}^{ML}) = 0$ Showing that the expected value of the score is zero:

Proof.

For all θ ,

 $1 = \int f(y|\theta) dy$ and therefore $0 = \frac{\partial}{\partial \theta} \int f(y|\theta) dy$

Swapping differentiation and integration, we get

$$0 = \int \frac{\partial}{\partial \theta} f(y|\theta) dy$$

= $\int \frac{\partial \ln f}{\partial \theta} (y|\theta) \cdot f(y|\theta) dy$
= $\int S(y|\theta) \cdot f(y|\theta) dy$
= $E(S(Y|\theta))$

How can we use this result?

Definition (Fisher Information)

Let Y be a random variable with pdf $f(y|\theta)$. The **Fisher Information** is defined as $I(\theta) := E(S(Y|\theta)^2)$.

Notice that the Fisher Information is also equal to the variance of the score bc Var $S(Y|\theta) = E(S(Y|\theta)^2) - E(S(Y|\theta))^2 = E(S(Y|\theta)^2)$

Here an important result for unbiased estimators $T(Y_1, ..., Y_N)$

Proposition (Cramér Rao Bound (CRB))

Let $Y_1, ..., Y_N$ be iid with $pdff(y|\theta)$ and let $T(Y_1, ..., Y_N)$ be an unbiased estimator for θ . Then $Var T(Y_1, ..., Y_N) \ge \frac{1}{N \cdot I(\theta)}$

Proof.

Let's only check the case N = 1. First notice that $Cov(S(Y|\theta), T(Y)) = E(S(Y|\theta) \cdot T(Y))$

$$\begin{split} &= \int S(y|\theta) \cdot T(y) \cdot f(y|\theta) dy \\ &= \int T(y) \frac{\partial \ln f}{\partial \theta}(y|\theta) f(y|\theta) dy \\ &= \int T(y) \cdot \frac{\partial f}{\partial \theta}(y|\theta) dy \\ &= \frac{\partial}{\partial \theta} \int T(y) \cdot f(y|\theta) dy \\ &= \frac{\partial}{\partial \theta} \mathbb{E} \left(T(Y) \right) = \frac{\partial}{\partial \theta} \theta = 1 \end{split}$$

Use Cauchy Schwarz inequality: $Cov(A, B)^2 \leq Var A \cdot Var B$

$$\operatorname{Var} T(Y) \ge \frac{\operatorname{Cov} \left(S(Y|\theta), T(Y) \right)^2}{\operatorname{Var} S(Y|\theta)} = \frac{1}{\operatorname{Var} S(Y|\theta)} = \frac{1}{I(\theta)}$$

What's the point of the CRB?

It gives us a lower bound for the variance of any unbiased estimator (not only ML estimators)

Aside:

Finding a lower bound for the variance is, in principle, not hard at all: zero is always a lower bound

The nice thing about the CRB, as we will see, is that it can actually be attained by some estimators

And ML estimators have a special relationship with the CRB as we will see

Again going back to the example of the exponential distribution The score was $S(y|\mu) := -((1/\mu) - y/\mu^2)$ Notice that $E(S(Y|u)) = -1/u + EY/u^2 = -1/u + u/u^2 = 0$ (This is a quick cross check that we haven't made any mistakes) Therefore $E(S(Y|\mu)^2) = Var S(Y|\mu) = VY/\mu^4 = \mu^{-2}$ Let $T(Y_1, \dots, Y_N)$ be an unbiased estimator for μ It follows that Var $T(Y_1, \dots, Y_N) \ge \mu^2/N$ The CRB for the variance of $T(Y_1, ..., Y_N)$ is μ^2/N

Another important result in ML theory

Proposition (Information Equality)

Let Y be a random variable with $pdff(y|\theta)$. Then $I(\theta) = -E\left(\frac{\partial S}{\partial \theta}(Y,\theta)\right).$

This gives us 3 ways to obtain the Fisher information; via the:

- \cdot variance of the score
- $\cdot\,$ second moment of the score
- first derivative of the score (which is the second derivative of the log likelihood)

Writing the information equality in terms of the original pdf

$$\mathsf{E}\left(\left(\frac{\partial \ln f}{\partial \theta}(Y|\theta)\right)^2\right) = -\mathsf{E}\left(\frac{\partial^2 \ln f}{\partial \theta^2}(Y|\theta)\right)$$

Proof.

When we derived the expected value of the score, we obtained

$$\int \frac{\partial \ln f}{\partial \theta}(y|\theta) \cdot f(y|\theta) dy = 0$$

Differentiating both sides

$$\int \frac{\partial^2 \ln f}{\partial \theta^2} (y|\theta) \cdot f(y|\theta) dy + \int \frac{\partial \ln f}{\partial \theta} (y|\theta) \cdot \frac{\partial f}{\partial \theta} (y|\theta) dy$$

= $E\left(\frac{\partial^2 \ln f}{\partial \theta^2} (Y|\theta)\right) + \int \frac{\partial \ln f}{\partial \theta} (y|\theta) \frac{\partial \ln f}{\partial \theta} (y|\theta) f(y|\theta) dy$
= $E\left(\frac{\partial^2 \ln f}{\partial \theta^2} (Y|\theta)\right) + \int \left(\frac{\partial \ln f}{\partial \theta} (y|\theta)\right)^2 \cdot f(y|\theta) dy$
= $E\left(\frac{\partial^2 \ln f}{\partial \theta^2} (Y|\theta)\right) + E\left(\left(\frac{\partial \ln f}{\partial \theta} (Y|\theta)\right)^2\right) = 0$

Yet again using the example of the exponential distribution The score was $S(y|\mu) := -((1/\mu) - y/\mu^2)$ We saw earlier that $E(S(Y|\mu)^2) = \mu^{-2}$ Let's look at the derivative of the score

$$\frac{\partial S}{\partial \mu}(y|\mu) = -(-1/\mu^2 + 2y/\mu^3)$$

Recalling EY = μ , you see quickly that $E\left(\frac{\partial S}{\partial \mu}(y|\mu)\right) = -\mu^{-2}$ This confirms the information equality Another thing that's interesting here is that we do know the actual variance of $\hat{\mu}^{\rm ML}$ because it is easy to derive

Recall that $\hat{\mu}^{ML} = \bar{Y}$ (the case in which N > 1)

Therefore Var $\hat{\mu}^{\text{ML}} = \text{Var } \bar{Y} = \mu^2 / N$

So when Y_i has an exponential distribution with parameter μ then the ML estimator $\hat{\mu}^{\rm ML}$ has a variance that attains the CRB

We've also seen that the ML estimator is unbiased

So $\hat{\mu}^{\rm ML}$ is the unbiased estimator with minimum variance

Is this true for all maximum likelihood estimators?

Maximum Likelihood (ML) Estimation

Maximum Likelihood Estimator

Invariance Property

Fisher Information, Cramér Rao Bound, Information Equality

Minimum Variance Unbiased Estimators

Asymptotic Efficiency of ML

What have we discovered so far?

The CRB gives a lower bound for the variance of unbiased estimators

Let's say you have obtained an unbiased estimator

You calculate its variance and compare that to the CRB

The CRB is most useful when we can show that our estimator actually attains it

If its variance is equal to the CRB then we know it must be the minimum variance unbiased estimator (as was the case for $\hat{\mu}^{\text{ML}}$ in exponential distribution example)

If its variance is not equal to the CRB then there are two possibilities: the estimator

- is not a minimum variance unbiased estimator (there is another unbiased estimator with smaller variance)
- is a minimum variance unbiased estimator with a variance greater than the CRB (the bound cannot be attained)

Under what conditions does a minimum variance unbiased estimator with variance equal to the CRB exist?

And how do we find such an estimator?

This is where ML estimation is so powerful:

Theorem

Let Y_1, \ldots, Y_N be iid with $pdff(y|\theta)$. An unbiased estimator $T(Y_1, \ldots, Y_N)$ with variance equal to the Cramér Rao bound exists if and only if the score function can be written $\frac{1}{N}\sum_{i=1}^N S(y_i|\theta) = a(\theta) \cdot (T(Y_1, \ldots, Y_N) - \theta),$

for some function $a(\theta)$.

The minimum variance unbiased estimator is then equal to the maximum likelihood estimator. Its variance is $1/a(\theta)$.

Proof.

Let's only check the case N = 1.

(only if)

Suppose T(Y) is an unbiased estimator with variance equal to the CRB. Then it must be the case that

$$\operatorname{Var} T(Y) = \frac{1}{I(\theta)} = \frac{1}{\operatorname{Var} S(Y|\theta)} = \frac{\operatorname{Cov} \left(S(Y|\theta), T(Y)\right)^2}{\operatorname{Var} S(Y|\theta)}$$

which implies $|Corr(T(Y), S(Y|\theta))| = 1$.

It follows that the score is a linear function of T(Y), with coefficients possibly depending on θ :

 $S(Y|\theta) = b(\theta) + a(\theta) \cdot T(Y)$

Since T(Y) is unbiased, and because $E(S(Y|\theta)) = 0$, it follows that $b(\theta) = -a(\theta) \cdot \theta$, therefore $S(Y|\theta) = a(\theta) \cdot (T(Y) - \theta)$

Proof.

(if)

Suppose the score can be written $S(Y|\theta) = a(\theta) \cdot (T(Y) - \theta)$.

```
Because ES(Y|\theta) = 0, ET(Y) = \theta, so it is unbiased. Furthermore
Var S(Y|\theta) = a(\theta)^2 Var T(Y)
```

which implies

$$\operatorname{Var} T(Y) = a(\theta)^{-2} \operatorname{Var} S(Y|\theta) = a(\theta)^{-2} \cdot I(\theta)$$

At the same time

$$\mathsf{E}\left(\frac{\partial S}{\partial \theta}(Y|\theta)\right) = \mathsf{E}\left(\frac{\partial a}{\partial \theta}(\theta)\left(T(Y) - \theta\right) - a(\theta)\right) = -a(\theta)$$

So $I(\theta) = a(\theta)$ by the information equality, and therefore Var $T(Y) = a(\theta)^{-1} = I(\theta)^{-1}$.

Lastly, that estimator would be found by $S(y|\hat{\theta}) = 0$ which is the ML estimator via the analogy principle

We already know that $\hat{\mu}^{\mathsf{ML}}$ attains the CRB

Let's show that the score can be written like in the theorem

Also let's understand the role of $a(\theta)$ in that theorem

We know from earlier that

$$S(y|\mu) = -\left((1/\mu) - y/\mu^2\right)$$
$$= \frac{y-\mu}{\mu^2}$$

Therefore

$$(1/N) \sum_{i=1}^{N} S(y_i|\mu) = \frac{\bar{Y} - \mu}{\mu^2} = \frac{\hat{\mu}^{\mathsf{ML}} - \mu}{\mu^2} = \mu^{-2} \left(\hat{\mu}^{\mathsf{ML}} - \mu \right),$$

implying $a(\mu) = \mu^{-2}$ So $a(\mu)$ is the Fisher Information!

Maximum Likelihood (ML) Estimation

Maximum Likelihood Estimator Invariance Property Fisher Information, Cramér Rao Bound, Information Equality Minimum Variance Unbiased Estimators Asymptotic Efficiency of ML The example of the exponential distribution again is instructive

- We have Y_1, \dots, Y_N iid with $pdf f(y|\mu) = (1/\mu) \exp(-y/\mu)$
- Now, define $\lambda := 1/\mu$ (recall the *arrival rate*)

Rewrite the pdf:
$$f(y|\lambda) = \lambda \exp(-\lambda y)$$

How would we estimate λ using ML?

Easy, by the invariance property: $\hat{\lambda}^{\rm ML} = 1/ar{Y}$

Interestingly, $\hat{\lambda}^{\mathrm{ML}}$ is biased

It can be shown that $\mathbf{E}\hat{\lambda}^{\mathrm{ML}} = \frac{N}{N-1}\lambda \neq \lambda$

So $\hat{\lambda}^{\mathrm{ML}}$ cannot be a minimum variance unbiased estimator

Luckily, we have the following result

Theorem (Asymptotic Normality of ML Estimators) Let $Y_1, ..., Y_N$ be iid with $pdff(y|\theta)$ Then

 $\sqrt{N}\left(\hat{\theta}^{\mathsf{ML}}-\theta\right) \xrightarrow{d} N\left(0, I(\theta)^{-1}\right)$

Notice that $I(\theta)^{-1}$ is the CRB

Loosely, we take this to mean that $\hat{\theta}^{\text{ML}} \stackrel{approx.}{\sim} N\left(\theta, \frac{1}{N \cdot I(\theta)}\right)$

Corollary (Consistency of ML Estimators) $\hat{Q}^{ML} = 0 + 0$ (1)

 $\hat{\theta}^{ML} = \theta + o_p(1)$

That is, maximum likelihood estimators are consistent and asymptotically normal distributed and attain the CRB

Asymptotically, ML estimators are minimum variance unbiased estimators

Sketch of proof is instructive

Start with the log likelihood $\ln L(\theta) = \sum_{i=1}^{N} \ln f(y_i|\theta)$

By definition, the derivative of $\ln L$ evaluated at $\hat{\theta}^{ML}$ is zero $\frac{\partial \ln L}{\partial \theta}(\hat{\theta}^{ML}) = 0$

Apply the mean value theorem around θ (the "true" parameter) $0 = \frac{\partial \ln L}{\partial \theta} (\hat{\theta}^{\text{ML}}) = \frac{\partial \ln L}{\partial \theta} (\theta) + \frac{\partial^2 \ln L}{\partial \theta^2} (\tilde{\theta}) \cdot (\hat{\theta}^{\text{ML}} - \theta),$

for some $\tilde{\theta}$ between θ and $\hat{\theta}^{\rm ML}$

Multiplying by \sqrt{N} and rearranging results in

$$\begin{split} \sqrt{N} \left(\hat{\theta}^{\mathsf{ML}} - \theta \right) &= \left(-\frac{1}{N} \frac{\partial^2 \ln L}{\partial \theta^2} (\tilde{\theta}) \right)^{-1} \cdot \left(\frac{1}{\sqrt{N}} \frac{\partial \ln L}{\partial \theta} (\theta) \right) \\ &\stackrel{\text{d}}{\to} \left(I(\theta)^{-1} + \mathsf{o}_p(1) \right) \cdot \mathsf{N} \left(0, I(\theta) \right) \\ &= \mathsf{N} \left(0, I(\theta)^{-1} \right) \end{split}$$

For $\hat{\lambda}^{ML} = 1/\bar{Y}$ we have $\hat{\lambda}^{ML} = \lambda + o_p(1)$ Asymptotically, the bias goes away: $E\hat{\lambda}^{ML} = \frac{N}{N-1}\lambda \rightarrow \lambda$ $V\hat{\lambda}^{ML} \simeq \lambda^2/N$ Life is good