

# Elements of Likelihood Theory

We will now examine elements of the likelihood theory, regarding properties of the resulting maximum likelihood estimator (MLE), and testing procedures based on the likelihood functions like the likelihood ratio test (LR), and the Score test.

The likelihood function is definable in parametric statistical models  
Even though the empirical likelihood

function, briefly considered in the notes regarding Gull, is a non-parametric adaptation).

Hence in this framework, the sample  $Z_n$  is paired with a parametric model,  $\{P_\theta, \theta \in \Theta\}$ , where  $P_\theta$  is a distribution related to the sample out scenario  $\theta \in \Theta$ . The mapping  $\theta \rightarrow P_\theta$  is 1-1, hence the parameter value  $\theta$  completely specifies the distribution in the model that indexes:

E.g. 1

$Z_n = (Z_1, Z_2, \dots, Z_m)$  is iid

and  $Z_{(1)} \sim \text{Exp}(\lambda_0)$  distribution  
with  $\lambda_0$  unknown.

Remember that the  $\text{Exp}(\lambda)$  distribution -  
 $\lambda > 0$  - is the one with density

$$f(z) = \begin{cases} 0, & z < 0 \\ \lambda e^{-\lambda z}, & z \geq 0 \end{cases}$$

A natural parametric model emerging  
from the above information is:

$$\{\text{Exp}(\lambda), \lambda > 0\}$$

There  $\Theta = \lambda$ ,  $\Theta = (0, +\infty)$ ,  $k=1$ , and each  
value of  $\lambda$  uniquely determines the

relevant exponential distribution.  $\mathbb{R}$

E.g. 2.

$Z_n = (Z_{(1)}, Z_{(2)}, \dots, Z_{(n)})$  is iid

and  $Z_{(n)} \sim N(\mu, \nu)$  distribution

with  $\theta_0 := (\mu_0, \nu_0)$  unknown. This

information uniquely determines the

following Gaussian model:

$$\left\{ Z_n \sim N\left(\begin{pmatrix} \mu \\ \nu \end{pmatrix}, \nu I_{n \times n}\right), \mu \in \mathbb{R}, \nu > 0 \right\}$$

There  $\theta = (\mu, \nu)$ ,  $\Theta = \mathbb{R} \times \mathbb{R}_{++}$ ,  $k=2$  and each

value of  $\theta$  uniquely determines the relevant

Gaussian distribution.  $\square$

E.g. 3 [Gaussian LS]

$$z_n = [y_n, x_n],$$

$$y_n / \phi(x_n) \sim \mathcal{N}(x_n \theta_0, \mathbf{I}_{x_n}),$$

with  $\theta_0 \in \mathbb{R}^k$  unknown. This information uniquely determines the model

$$\left\{ y_n / \phi(x_n) \sim \mathcal{N}(x_n \theta, \mathbf{I}_{x_n}), \theta \in \mathbb{R}^k \right\}$$

where  $\Theta = \mathbb{R}^k$ , etc.  $\square$

[In what follows we will investigate further examples]

**Note:** Remember that we work under the universal assumption of correct specification;  $P_{\theta_0} \in \text{Model}$ . We will provide some misspecification considerations later.

**Question:** Is it possible to construct a criterion  $C_n$  that fully and faithfully represent the information that appears in the model jointly with the sample?

**Brevity Assumption**

We will assume

that for any  $\theta \in \Theta$ ,  $P_{\theta}$  has a density

function (it should be a real function that depends on  $(\theta, z)$ , say  $f(\theta, z)$ ).

The constructions that follow are easily extendable to cases where densities do not exist.  $\square$

Kullback - Liebler Divergence from  $P_{\theta_0}$  and the likelihood function.

Given the Brevity Assumption, for any  $\theta \in \Theta$ , the KL divergence between the parent  $P_{\theta_0}$  and  $P_{\theta}$ , is defined

as:

$$KL(P_{\theta_0}, P_{\theta}) = \int \ln \frac{f(\theta_0, z)}{f(\theta, z)} f(\theta_0, z) dz$$

*this lies in the space where  $z_n$  assumes its values with prob. 1*

It is not difficult to prove that

argmin  $KL(P_{\theta_0}, P_{\theta}) = P_{\theta_0}$ , hence

$P_{\theta} \in \text{Model}$

the latent distribution has a variational characterization. It is however impossible

to recover  $P_{\theta_0}$ , by minimizing  $KL(P_{\theta_0}, P_{\theta})$

since this depends on  $P_{\theta_0}$ . It however

the empirical joint distribution of the

sample  $P_n^*$  (remember  $P_n^*$  is the discrete

distribution that places probability mass at

$z_n$ )

"Approximates",  $P_{\theta_0}$ , and the integral that defines  $KL(P_{\theta_0}, P_{\theta})$  is approximated by an empirical mean w.r.t.  $IP_n$  we obtain

$$\ln \frac{1}{f(\theta, z_n)} =$$

$$= -\ln f(\theta, z_n)$$

$\hookrightarrow$  this can be perceived as an empirical approximation of  $KL(P_{\theta_0}, P_{\theta})$  for any  $\theta \in \Theta$ .

Hence, we may hope that we can approximate  $\theta_0$  (and hence  $P_{\theta_0}$ ), via

minimizing (w.r.t.  $\theta$ ) the empirical approximation.

Notice that minimizing this is equivalent to maximizing its opposite, whereas maximization is covariant to multiplying this by  $\frac{1}{n}$  "for averaging purposes". Hence we arrive at the concept of the (average log-) likelihood function (following the general convention we will denote it with  $l_n$  instead of  $l$ ).

## Definition

Given the model, the brevity assumption, and if  $f(\theta, z_n) > 0$  at the support of  $P_\theta$  (the "smallest," closed subset of the space on which  $z_n$  attains its values of probability one)

the average log-likelihood function

$$\text{is: } \ell_n(\theta) = \frac{1}{n} \ln(f(\theta, z_n))$$

E.g. 1

We have that  $f(\lambda, z_n) =$

$$= \begin{cases} 0, & \exists i: z_i < 0 \\ \prod_{i=1}^n \lambda \exp(-\lambda z_i), & z_i \geq 0 \forall i \end{cases}$$

[the support of  $z_n$  is  $[0, \infty)^n$ ]

$$\text{hence } \ln(\lambda, z_n) = \frac{1}{n} \sum_{i=1}^n \ln(\lambda \exp(-\lambda z_i))$$

$$= \frac{1}{n} \sum_{i=1}^n \left( \ln(\lambda) + \ln(\exp(-\lambda z_i)) \right)$$

$$= \ln(\lambda) - \lambda \frac{1}{n} \sum_{i=1}^n z_i$$

E.g. 2 We have that

$$f(\theta, z_n) = \frac{(2\pi)^{-n/2}}{v^{n/2} \sqrt{\det \Gamma_{nn}}} \exp\left(-\frac{1}{v} (z_n - \begin{pmatrix} \mu \\ \psi \end{pmatrix})' \Gamma_{nn}^{-1} (z_n - \begin{pmatrix} \mu \\ \psi \end{pmatrix})\right)$$

And therefore

$$\begin{aligned} \ln(\theta) &\approx -\frac{1}{2} \ln 2\pi - \frac{1}{2} \ln v - \frac{1}{2v} (z_n - \begin{pmatrix} \mu \\ \psi \end{pmatrix})' (z_n - \begin{pmatrix} \mu \\ \psi \end{pmatrix}) \\ &= -\frac{1}{2} \ln 2\pi - \frac{1}{2} \ln v - \frac{1}{2v} \sum_{i=1}^n (z_i - \psi)^2 \end{aligned}$$

E.g. 3.

We have that

$$f(\theta, y_n | \mathcal{G}(x_n)) = \frac{(2\pi)^{-n/2}}{\sqrt{\det I_{\theta n}}} \exp\left(-\frac{1}{2} (y_n - x_n \theta)' (y_n - x_n \theta)\right)$$

And thereby

$$\ln(\theta, y_n | \mathcal{G}(x_n)) = -\frac{1}{2} \ln 2\pi - \frac{1}{2} \frac{1}{n} (y_n - x_n \theta)' (y_n - x_n \theta)$$

↓  
the LS criterion.

The MLE

Setting  $Q_n := -\ln$  and following

the general theory the OLS Maximum Likelihood Estimator is established:

$$\theta_n \in \arg \max_{\theta \in \Theta} \ell_n(\theta) = \arg \max_{\theta \in \Theta} \ell_n(\theta)$$

MLE

or more generally (allowing for the existence of optimization error)

$$\ell_n(\theta_n) \geq \sup_{\theta \in \Theta} \ell_n(\theta) - \epsilon_n$$

$\epsilon_n$  optimization error

Note: the existence arguments for the QMLE fall into the ones of our general theory: e.g.  $\Theta$  convex, and  $\ell_n(\theta)$  concave, or e.g.  $\Theta$  compact, and  $\ell_n(\theta)$  continuous, suffice for existence.

E.g. 1

Remember that  $\Theta = (0, +\infty)$

$\rightarrow \Theta = \lambda$  here!

$$\text{and } \ln(z) = \ln(\alpha) - \lambda \frac{1}{n} \sum_{i=1}^n z_i u_i$$

$$\ln \text{ is smooth, } \frac{\partial \ln(\alpha)}{\partial \alpha} = \frac{1}{\alpha} - \frac{1}{n} \sum_{i=1}^n z_i u_i$$

$$\text{And } \frac{\partial^2 \ln(\alpha)}{\partial \alpha^2} = -\frac{1}{\alpha^2} < 0 \quad \forall \alpha \in \Theta, \text{ hence}$$

the function is strictly concave with  $\Theta$  convex. The maximizer would be unique and Kuhn-Tucker conditions are not needed to locate it as it is an interior point, as we will see:

$$\text{f.o.c. } \frac{\partial \ln(\alpha)}{\partial \alpha} = 0 \Leftrightarrow \frac{1}{\alpha} - \frac{1}{n} \sum_{i=1}^n z_i u_i = 0$$

→ the unique critical point.

$$\hat{\lambda}_n = \frac{n}{\sum_{i=1}^n z_{(i)}}$$

→ \* Notice that this is almost surely well-defined as  $X_{(i)} > 0$  with Prob. 1  $\forall i$ .

$$\text{s.o.c.} \quad \frac{\partial^2 \ln(\hat{\lambda}_n)}{\partial \lambda^2} = - \left( \frac{\sum_{i=1}^n z_{(i)}}{n} \right)^2 < 0$$

with Prob. 1, hence the MLE is

$$\hat{\lambda}_n = \frac{n}{\sum_{i=1}^n z_{(i)}}$$

→ the reciprocal of the sample mean.

\* It is easily proven that if  $X \sim \text{Exp}(\lambda_0)$  then  $\mathbb{E}(X) = 1/\lambda_0$ . We have that

$$\mathbb{E}(\hat{\lambda}_n) = \mathbb{E}\left(\frac{n}{\sum_{i=1}^n z_{(i)}}\right) = n \mathbb{E}\left(\sum_{i=1}^n z_{(i)}\right)^{-1} \neq$$

$$n \left( \mathbb{E} \left( \sum_{i=1}^n z_{ii} \right) \right)^{-1} = n \left( \sum_{i=1}^n \mathbb{E} (z_{ii}) \right)^{-1}$$

$$\begin{aligned} X_{ii} &\sim \text{Exp}(\lambda_0), \forall i \\ &= n \left( \sum_{i=1}^n \frac{1}{\lambda_0} \right)^{-1} = n n^{-1} \lambda_0 = \lambda_0 \end{aligned}$$

Hence the MLE of  $\lambda_0$  is biased, which implies that the MLE in general is not unbiased (Unbiasedness is a difficult to hold and fragile property!)  $\mathbb{R}$

E.g. 2

Remember that  $\Theta = \mathbb{R} \times \mathbb{R}_{++}$

$$\theta = \begin{pmatrix} \mu \\ \nu \end{pmatrix}, \quad \ln(\theta) = -\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln \nu - \frac{1}{2} \frac{1}{\nu} \sum_{i=1}^n \underbrace{(z_{ii} - \mu)^2}_{\nu}$$

The function is again smooth with

$$\frac{\partial \ln(\theta)}{\partial \theta} = \begin{pmatrix} \partial \ln(\theta) / \partial \mu \\ \partial \ln(\theta) / \partial \nu \end{pmatrix},$$

$$\frac{\partial \ln(\theta)}{\partial \psi} = \frac{1}{\psi} \frac{1}{n} \sum_{i=1}^n (z_{(i)} - \psi) = \frac{1}{\psi} \left( \frac{1}{n} \sum_{i=1}^n z_{(i)} - \psi \right)$$

$$\frac{\partial \ln(\theta)}{\partial \nu} = -\frac{1}{2\nu} + \frac{1}{2\nu^2} \frac{1}{n} \sum_{i=1}^n (z_{(i)} - \psi)^2$$

$$\frac{\partial^2 \ln(\theta)}{\partial \theta \partial \theta'} = \begin{pmatrix} \frac{\partial^2 \ln(\theta)}{\partial \psi^2} & \frac{\partial^2 \ln(\theta)}{\partial \psi \partial \nu} \\ \frac{\partial^2 \ln(\theta)}{\partial \nu \partial \psi} & \frac{\partial^2 \ln(\theta)}{\partial \nu^2} \end{pmatrix}$$

[Second order derivatives are continuous in  $\theta$ , Young's Th. implies that  $\frac{\partial^2 \ln(\theta)}{\partial \nu \partial \psi} = \frac{\partial^2 \ln(\theta)}{\partial \psi \partial \nu}$ , check!]

$$\frac{\partial^2 \ln(\theta)}{\partial \psi^2} = \frac{\partial \left( \frac{1}{2\nu} \left( \frac{1}{n} \sum_{i=1}^n z_{(i)} - \psi \right) \right)}{\partial \psi} = -\frac{1}{\psi}$$

$$\frac{\partial^2 \ln(\theta)}{\partial \nu \partial \psi} = \frac{\partial \left( \frac{1}{2\nu} \left( \frac{1}{n} \sum_{i=1}^n z_{(i)} - \psi \right) \right)}{\partial \nu} = -\frac{1}{2\nu^2} \left( \frac{1}{n} \sum_{i=1}^n z_{(i)} - \psi \right)$$

$$\frac{\partial^2 \ln(\theta)}{\partial v^2} = \frac{\frac{1}{2v} + \frac{1}{2v^2} \cdot \frac{1}{n} \sum_{i=1}^n (z_{(i)} - v)^2}{\partial v} =$$

$$= \frac{1}{2v^2} - \frac{1}{v^3} \cdot \frac{1}{n} \sum_{i=1}^n (z_{(i)} - v)^2.$$

We have that:

$$\text{f.o.c. } \frac{\partial \ln(\theta)}{\partial \theta} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad (=)$$

$$\left. \begin{array}{l} \frac{\partial \ln(\theta)}{\partial v} = 0 \\ \frac{\partial \ln(\theta)}{\partial \mu} = 0 \end{array} \right\} \Rightarrow \left. \begin{array}{l} \frac{1}{n} \left( \frac{1}{n} \sum_{i=1}^n z_{(i)} - v \right) = 0 \\ \frac{1}{2v} + \frac{1}{2v^2} \cdot \frac{1}{n} \sum_{i=1}^n (z_{(i)} - v)^2 = 0 \end{array} \right\} (=)$$

$$\left. \begin{array}{l} \mu_n = \frac{1}{n} \sum_{i=1}^n z_{(i)} \\ v_n = \frac{1}{n} \sum_{i=1}^n (z_{(i)} - \mu_n) \end{array} \right\} \Rightarrow \theta_n = \begin{pmatrix} \mu_n \\ v_n \end{pmatrix}$$

is the critical point

Notice that:

$$\frac{\partial^2 \ln(\theta_n)}{\partial \nu^2} = -\frac{1}{\nu n} < 0 \text{ with Prob } 1 \text{ since}$$

if  $Z_{(i)} \sim N(\mu_0, \nu_0)$ ,  $P(Z_{(i)} = 0) = 0$

the distribution has  
a density

$$\frac{\partial^2 \ln(\theta_n)}{\partial \nu \partial \mu} = -\frac{1}{\nu n} \left( \frac{1}{n} \sum_{i=1}^n Z_{(i)} - \mu_n \right) =$$

$$= -\frac{1}{\nu n} (\mu_n - \mu_n) = 0$$

$$\frac{\partial^2 \ln(\theta_n)}{\partial \nu^2} = \frac{1}{2\nu n^2} - \frac{1}{\nu n^3} \left( \frac{1}{n} \sum_{i=1}^n (Z_{(i)} - \mu_n)^2 \right)$$

$$= \frac{1}{2\nu n^2} - \frac{\nu n}{\nu n^3} = -\frac{1}{\nu n^2} < 0 \text{ with Prob } 1$$

Thereby

$$\frac{\partial^2 h(\theta_n)}{\partial \theta \partial \theta'} = \begin{pmatrix} -1/2v_n & 0 \\ 0 & -1/v_n^2 \end{pmatrix}$$

the Hessian at  $\theta_0$  is diagonal; hence its eigenvalues are its diagonal elements; hence it is negative definite (remember that we dually defined the MLE via Maximization) iff its diagonal elements are negative; they are with Prob 1.

Hence, with Prob 1, the MLE for  $\theta_0$ ,

$$\theta_n = \begin{pmatrix} \mu_n \\ v_n \end{pmatrix} = \begin{pmatrix} 1/n \sum_{i=1}^n z_i \\ 1/n \sum_{i=1}^n (z_i - \mu_n)^2 \end{pmatrix}$$

It is easy to show (do it!) that

$$\mathbb{E}(\mu_n) = \mu_0, \quad \mathbb{E}(v_n) = \frac{n-1}{n} v_0, \quad \text{hence } \mathbb{E}(\theta_n) =$$

$$= \begin{pmatrix} \mathbb{E}(y_n) \\ \mathbb{E}(v_n) \end{pmatrix} = \begin{pmatrix} y_0 \\ \frac{n-1}{n} v_0 \end{pmatrix} \neq \begin{pmatrix} y_0 \\ v_0 \end{pmatrix}, \text{ hence}$$

the MLE is biased also in this example:

Notice however that since  $\lim_{n \rightarrow \infty} \frac{n-1}{n} = 1$ ,

$$\lim_{n \rightarrow \infty} \mathbb{E}(\hat{\theta}_n) = \begin{pmatrix} y_0 \\ y_0 \lim_{n \rightarrow \infty} \frac{n-1}{n} \end{pmatrix} = \begin{pmatrix} y_0 \\ v_0 \end{pmatrix},$$

hence in this example it is asymptotically unbiased.

**Note:** In this example no Kuhn-Tucker type conditions to enforce the restriction  $v > 0$  were needed; the optimizer is an interior point with Prob. 1.

□

E.g. 3

Remember that

$$l_n(\theta) = -\frac{1}{2} \ln 2\pi - \frac{1}{2n} (y_n - x_n \theta)' (y_n - x_n \theta)$$

and due to monotonicity and duality

$$\underset{\theta \in \Theta}{\operatorname{argmax}} l_n(\theta) \stackrel{\text{Mon}}{=} \underset{\theta \in \Theta}{\operatorname{argmax}} \left( -\frac{1}{2n} (y_n - x_n \theta)' (y_n - x_n \theta) \right)$$

$$\stackrel{\text{Mon}}{=} \underset{\theta \in \Theta}{\operatorname{argmax}} \left( -\frac{1}{n} (y_n - x_n \theta)' (y_n - x_n \theta) \right)$$

$$\stackrel{\text{dual}}{=} \underset{\theta \in \Theta}{\operatorname{argmin}} \frac{1}{n} (y_n - x_n \theta)' (y_n - x_n \theta)$$

Hence  $\text{MLE}(\theta_0) = \text{OLS}(\theta_0)$  in this setup. We will try to figure out the reason of this coincidence later.  $\square$

## General Note 1.

Remember that in the construction of the likelihood function

$f(\theta, z_n)$  denotes the density, of the (joint) distribution of  $z_n$ , at

Parameter scenario  $\theta$ . This can be factored as a product of conditional densities (times a marginal) as

$$f(\theta, z_n) = f(\theta, z_{(n)} / z_{(n-1)}, z_{(n-2)}, \dots, z_{(1)}) \\ \times f(\theta, z_{(n-1)} / z_{(n-2)}, z_{(n-3)}, \dots, z_{(1)}) \times \dots \times f(\theta, z_{(1)})$$

where  $f(\theta, z_{(n-k)} / z_{(n-k-1)}, z_{(n-k-2)}, \dots, z_{(n)})$  is the density of the conditional distribution of the sample element  $z_{(n-k)}$  given  $z_{(n-k-1)}, z_{(n-k-2)}, \dots, z_{(n)}$ , and  $f(\theta, z_{(1)})$  is the marginal density of the first sample element, both are parameter scenario  $\theta$ . Hence:

$$\begin{aligned}
 l_n(\theta) &= \frac{1}{n} \ln (f(\theta, z_n)) = \\
 &= \frac{1}{n} \ln \left( \prod_{i=1}^{n-1} f(\theta, z_{(n-i+1)} / z_{(n-i)}, \dots, z_{(n)}) f(\theta, z_{(1)}) \right) \\
 &= \sum_{i=1}^{n-1} \frac{1}{n} \ln f(\theta, z_{(n-i+1)} / z_{(n-i)}, \dots, z_{(n)}) + \frac{1}{n} \ln f(\theta, z_{(1)})
 \end{aligned}$$

Then the part

$$\sum_{i=2}^{n-1} \ln f(\theta, z_{(n-i+1)} / z_{(n-i)}, \dots, z_{(1)})$$

is termed as the **conditional** (part of the) **likelihood**, while the term

$-\ln f(\theta, z_{(n)})$  is termed as the **marginal** part of the likelihood. In several models (see e.g. the GARCH(1,1) below) the marginal part is unspecified, while the conditional part is only used for the construction of the estimator.

It is expected that as  $n \rightarrow \infty$  the marginal part becomes in any case negligible.

- When  $Z_n$  is iid then

$$f(\theta, z_{(n-i+1)} | z_{(n-i)}, \dots, z_{(i)}) = f(\theta, z_{(n-i+1)})$$

$\forall i = 1, \dots, n-1$  and thereby

$$l_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ln f(\theta, z_{(i)})$$

As is essentially the case with Eq. 1  
and Eq. 2 (what about Eq. 3?).

## General Note 2.

The construction of the likelihood function is possible when density functions do not exist; e.g.  $Z_n$  is comprised by

discrete random variables, using the  
at the time probabilistic machinery  
available.

Example.  $Z_n$  is iid, and

$Z_i \sim \text{Ber}(q_0)$  with unknown  $q_0 \in (0,1)$ .

The natural choice for the statistical  
model is  $\{ \text{Ber}^n(q), q \in (0,1) \}$ .

Remember that if  $Z_{(1)} \sim \text{Ber}(q)$  then  
its probability mass function is

$$p_{Z_{(1)}}^{(q)} : \{0,1\} \rightarrow (0,1), \quad p_{Z_{(1)}}^{(q)}(x) = \begin{cases} q, & x=1 \\ 1-q, & x=0 \end{cases}$$

Using an argument based on the definition of the Kullback-Leibler divergence it is provable that in this case the loss functions are usable in place of the non-existent densities. Hence

$$\begin{aligned}
 \hat{I}_n(q) &= \frac{1}{n} \sum_{i=1}^n \ln P_{z^{(i)}}^{(q)}(z^{(i)}) = \\
 &= \frac{1}{n} \sum_{i=1}^n \ln q \cdot \mathbb{1}_{z^{(i)}=L} + \ln(1-q) \mathbb{1}_{z^{(i)}=0} \\
 &= \frac{1}{n} \left[ \ln q \cdot \#\{i=1, \dots, n, z^{(i)}=L\} + \ln(1-q) \#\{i=1, \dots, n, \right. \\
 &\quad \left. z^{(i)}=0\} \right] = \ln q \cdot (\text{sample freq of } L) \\
 &\quad + \ln(1-q) (1 - \text{sample freq of } L).
 \end{aligned}$$