Maximum Likelihood

Sofia Olhede



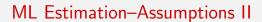
October 12, 2020

Maximum Likelihood





- Provided it exists, the MLE of the natural parameter in a k-parameter natural exponential family with open parameter space Φ is consistent.
- Assuming we can get consistency, we can focus on understanding the sampling distribution of the MLE.
- For simplicity, assume X_1, \ldots, X_n are iid with density/frequency) $f(x; \theta)$ for $\theta \in \Theta$. Write
- Let $\ell(x_i; \theta) = \log f(x_i; \theta)$.
- Let $\ell'(x_i; \theta)$, $\ell''(x_i; \theta)$ and $\ell'''(x_i; \theta)$ denote the partial derivatives wrt θ .
- We need some regularity conditions:
- A1: Θ is an open subset of \mathbb{R} .
- A2: The support of f, supp f is independent of θ .
- A3: f is thrice continuously differentiable w.r.t. θ for all the support of f.





- A4: $\mathbb{E}(\ell'(X_i;\theta)) = 0$ for all θ and \mathbb{V} ar $\{\ell'(X_i;\theta)\} = \mathcal{I}_1(\theta)$.
- A5: $-\mathbb{E}(\ell''(X_i;\theta)) = \mathscr{I}_1(\theta)$.
- A6: $\exists M(x) > 0$ and $\delta > 0$ such that $\mathbb{E} M(x) < \infty$ and

$$|\theta - \theta_0| < \delta \Rightarrow |\ell'''(x;\theta)| < M(x).$$

EPFL

ML Estimation-Explaining Assumptions

- If Θ is an open set then for θ_0 the true parameter, it always makes sense for an estimator $\widehat{\theta}_n$ to have a symmetric distribution around θ_0 (such as the Gaussian).
- Under condition (A2) we have

$$\frac{d}{d\theta} \int_{\text{supp } f} f(x; \theta) \, dx = 0.$$

This means that we are permitted to exchange integration and differentiation. Therefore it follows

$$0 = \int \frac{d}{d\theta} f(x;\theta) dx = \int \ell'(X_i;\theta) f(x;\theta) = \mathbb{E}\{\ell'(X_i;\theta)\}.$$

Thus once A2 and A4 hold, then we can exchange the order of the limits; and it ensures a finite variance.

- The second derivate has a finite moment by A5.
- A2 and A6 are assumptions that allow us to simplify ("linearize") our understanding of the MLE.



ML Estimation-Explaining Assumptions II

 Taking our exchange of limits even further: if we can differentiate twice then

$$0 = \int \frac{d}{d\theta} \{ \ell'(x;\theta) f(x;\theta) \} dx$$

=
$$\int \ell''(x;\theta) f(x;\theta) dx + \int \{ \ell'(x;\theta) \}^2 f(x;\theta) dx.$$

Thus we may deduce that $\mathcal{I}(\theta) = \mathscr{I}(\theta)$



• Theorem: Let X_1, \ldots, X_n be IID random variables with the same density $f(x; \theta)$. Assume that A1–A6 are satisfied. If the MLE $\widehat{\theta}_n$ exists and is unique, and we have consistency then

$$\sqrt{n}\Big\{\widehat{\theta}_n-\theta\Big\}\stackrel{\mathcal{L}}{\to} N\big(0,\mathcal{I}_1(\theta)/\mathscr{I}_1^2(\theta)\big).$$

Furthermore, when we can say that $\mathcal{I}(\theta) = \mathscr{I}(\theta)$ then

$$\sqrt{n}\Big\{\widehat{\theta}_n-\theta\Big\}\stackrel{\mathcal{L}}{\to} N(0,1/\mathcal{I}_1(\theta)).$$

For finite samples we often say

$$\widehat{\theta}_n \stackrel{d}{\approx} N(\theta, 1/\mathcal{I}_n(\theta)).$$

 Thus for large enough samples, the MLE is approximate Gaussian, approximately unbiased, and approximately achieving the Cramér–Rao lower bound.



• Proof: Assuming (A1)-(A3) we can note that

$$\sum_{i=1}^n \ell'(X_i; \hat{\theta}_n) = 0.$$

We now implement a Taylor series to deduce that

$$0 = \sum_{i=1}^{n} \ell'(X_i; \hat{\theta}_n) = \sum_{i=1}^{n} \ell'(X_i; \theta) + (\hat{\theta}_n - \theta) \sum_{i=1}^{n} \ell''(X_i; \hat{\theta}_n)$$
 (1)

$$+\frac{1}{2}(\hat{\theta}_n - \theta)^2 \sum_{i=1}^n \ell'''(X_i; \theta_n^*). \tag{2}$$

We have terminated the Taylor series after 3 terms. To make the equality hold, that means we need an exact form for the remainder. We here use the Lagrane form, other forms include Cauchy. θ_n^* lies between θ and $\widehat{\theta}_n$.



• Dividing the last equation by \sqrt{n} means we obtain

$$0 = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} \ell'(X_i; \theta) + (\hat{\theta}_n - \theta) \frac{1}{\sqrt{n}} \sum_{i=1}^{n} \ell''(X_i; \hat{\theta}_n) + \frac{1}{2\sqrt{n}} (\hat{\theta}_n - \theta)^2 \sum_{i=1}^{n} \ell'''(X_i; \theta_n^*).$$
(3)

Now, we use this equation to re-write

$$\label{eq:theta_n} (\hat{\theta}_n - \theta) = -\frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n \ell'(X_i; \theta)}{\frac{1}{\sqrt{n}} \sum_{i=1}^n \ell''(X_i; \hat{\theta}_n) + \frac{1}{2\sqrt{n}} (\hat{\theta}_n - \theta) \sum_{i=1}^n \ell'''(X_i; \theta_n^*)}.$$

 We can now use the Central Limit Theorem (CLT) that we just derived to deduce

$$\frac{1}{\sqrt{n}}\sum_{i=1}^n \ell'(X_i;\theta) \stackrel{d}{\to} N(0,1/\mathcal{I}_n(\theta)).$$



Furthermore, we can note that (from the Law of Large Numbers):

$$\frac{1}{\sqrt{n}}\sum_{i=1}^n \ell''(X_i; \hat{\theta}_n) \stackrel{p}{\to} -\mathscr{I}(\theta).$$

We define the remainder term to be

$$R_n = \frac{1}{2\sqrt{n}}(\hat{\theta}_n - \theta) \sum_{i=1}^n \ell'''(X_i; \theta_n^*).$$

If we can show that $R_n \stackrel{p}{\to} 0$ then the denominator tends to $\mathscr{I}(\theta)$ and we can use Slutsky's Lemma to deduce the result.

• We have already showed $\hat{\theta}_n - \theta \stackrel{p}{\to} 0$ and $\frac{1}{\sqrt{n}} \sum_{i=1}^n \ell'''(X_i; \theta_n^*)$ does not diverge and so we may deduce that $R_n \stackrel{p}{\to} 0$. Thus the result follows.



- Can we deduce from this result that maximum likelihood estimators are optimal?
- Well this is an asymptotic result, e.g. it holds **eventually** in *n*.
- For a fixed value of *n* things are less clear.
- Also we have assumed that we are comparing unbiased estimatorswhat can be gained by relaxing this assumption?

Shrinkage Estimation



- As simple results by Charles Stein show, things are not as simple as they seem.
- Let Y_1, \ldots, Y_n be independent random variables.
- Assume that each $Y_i \sim N(\mu_i, \sigma^2)$. Each Y_i has a different mean, but the variance is coupled.
- First we might consider the slightly simpler case of $\sigma^2 = 1$.
- Then we wish to estimate μ .
- We use mean square error to assess performance.
- This is not quite a "standard" problem as the number of parameters is growing with the sample size.

Shrinkage Estimation II



The log-likelihood can easily be written up as

$$\ell(\boldsymbol{\mu}) = -\frac{n}{2}\log(2\pi) - \frac{1}{2}\sum_{i=1}^{n}(Y_i - \mu_i)^2.$$

Using differentiation we can easily deduce that

$$\widehat{m{\mu}} = m{Y}$$
.

- This is like having *n* MLE's each of sample size 1. Not great estimation; too few samples.
- The MSE is then n as

$$\mathrm{MSE}(\widehat{\boldsymbol{\mu}}, \boldsymbol{\mu}) = \mathbb{E} \|\widehat{\boldsymbol{\mu}} - \boldsymbol{\mu}\|^2 = n.$$

Is this the best we can do?

Shrinkage Estimation III



- Assume that $\boldsymbol{Y} = (Y_1, \dots Y_n)^T$ such that $\boldsymbol{Y} \sim \mathcal{N}(\boldsymbol{\mu}, \mathsf{I}_n)$ where $\boldsymbol{\mu} \in \mathbb{R}^n$. We call this "Stein's set-up".
- We define a shrinkage estimator to be

$$\widetilde{\mu}_{a} = \left(1 - \frac{a}{\|\mathbf{Y}\|^{2}}\right)\mathbf{Y} = \left(1 - \frac{a}{\|\widehat{\boldsymbol{\mu}}\|^{2}}\right)\widehat{\boldsymbol{\mu}}.$$

This is the shrunken version of the MLE. Assume that $n \ge 3$ then

• for all $a \in (0, 2n - 4)$

$$\mathrm{MSE}(\widetilde{\mu}_a, \mu) \leq \mathrm{MSE}(\widehat{\mu}, \mu),$$

• for a = n - 2

$$\mathrm{MSE}(\widetilde{\boldsymbol{\mu}}_{n-2},0) \leq \mathrm{MSE}(\widehat{\boldsymbol{\mu}},0),$$

• For all $\mu \in \mathbb{R}^n$ and all $a \in (0, 2n - 4)$

$$MSE(\widetilde{\mu}_{n-2}, \mu) \leq MSE(\widetilde{\mu}_a, \mu).$$

Shrinkage Estimation IV



- This was a very surprising result at the time!
- First, the MLE is outperformed.
- The Stein estimator takes the MLE and "shrinks" its magnitude.
- The amount of shrinkage depends on $\|Y\|$.
- This estimate takes the estimate of μ_j into account when estimating μ_i .
- There are no "smoothness" assumptions.
- The performance of the MLE deteriorates as *n* increases.
- To understand this result we need to go in stages.

Shrinkage Estimation V



- Lemma: let $Y \sim \mathcal{N}(\theta, \sigma^2)$ and assume $h: \mathbb{R} \to \mathbb{R}$ is differentiable. If
 - 1. $\mathbb{E}|h(Y)|<\infty$;
 - 2. $\lim_{y\to\pm\infty} \left\{ h(y) \exp\left(-\frac{1}{2\sigma^2}(y-\theta)^2\right) \right\} = 0$ then $\mathbb{E}\{h(Y)(Y-\theta)\} = \sigma^2 \mathbb{E}(h'(Y)).$
- Proof: we note that from the definition of expectation

$$\mathbb{E}\{h(Y)(Y-\theta)\} = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{\mathbb{R}} h(y)(y-\theta) e^{-\frac{1}{2\sigma^2}(y-\theta)^2}$$

$$= \left[-\frac{\sigma}{\sqrt{2\pi}} h(y) e^{-\frac{1}{2\sigma^2}(y-\theta)^2} \right]_{y\to-\infty}^{\infty}$$

$$+ \frac{\sigma}{\sqrt{2\pi}} \int_{-\infty}^{\infty} h'(y) e^{-\frac{1}{2\sigma^2}(y-\theta)^2} dy$$

$$= \sigma^2 \, \mathbb{E}(h'(Y)). \tag{4}$$

This completes the proof of the lemma.





- Now we are ready to prove the Stein theorem.
- We note that from first principles we have that

$$MSE(\widetilde{\boldsymbol{\mu}}_{a}, \boldsymbol{\mu}) = \mathbb{E} \| \left(1 - \frac{a}{\|\mathbf{Y}\|^{2}} \right) \mathbf{Y} - \boldsymbol{\mu} \|^{2}$$

$$= \mathbb{E} \| \mathbf{Y} - \boldsymbol{\mu} - \frac{a}{\|\mathbf{Y}\|^{2}} \mathbf{Y} \|^{2}$$

$$= \mathbb{E} \| \mathbf{Y} - \boldsymbol{\mu} \|^{2} - 2 \mathbb{E} \left(\frac{a \mathbf{Y}^{T} (\mathbf{Y} - \boldsymbol{\mu})}{\|\mathbf{Y}\|^{2}} \right) + \mathbb{E} \left(\frac{a \|\mathbf{Y}\|^{2}}{\|\mathbf{Y}\|^{4}} \right)$$

$$= n - 2a \sum_{i=1}^{n} \mathbb{E} \left[\frac{Y_{i} (Y_{i} - \mu_{i})}{\sum_{j=1}^{n} Y_{j}^{2}} \right] + a^{2} \mathbb{E} \left(\frac{1}{\|\mathbf{Y}\|^{2}} \right). \tag{5}$$

• To understand this quantity we need to study the middle term.

Shrinkage Estimation VII



• We define n differentiable functions $h_i: \mathbb{R}^n \to \mathbb{R}$ by

$$\mathbf{u} = (u_1, \ldots, u_n) \stackrel{h_i}{\mapsto} \frac{u_i}{u_i^2 + \sum_{j \neq i} u_j^2}.$$

ullet We observe that for all $i\in\{1,\ldots,n\}$ and all $\{u_j\}_{j
eq i}\in\mathbb{R}^{n-1}$ then

$$\lim_{u_i\to\pm\infty}\left\{h_i(\boldsymbol{u})\exp\left[-\frac{1}{2\sigma^2}(u_i-\mu_i)^2\right]\right\}=0.$$

- ullet We note that h_i becomes an $\mathbb{R} o \mathbb{R}$ fn once $\{u_j\}_{j \neq i} \in \mathbb{R}^{n-1}$ is fixed.
- We now use the tower property and apply our lemma to re-write:

$$\mathbb{E}\left[\frac{Y_i(Y_i-\mu_i)}{\sum_{j=1}^n Y_j^2}\right].$$





Applying our lemma we get

$$\mathbb{E}\left[\frac{Y_{i}(Y_{i} - \mu_{i})}{\sum_{j=1}^{n} Y_{j}^{2}}\right] = \mathbb{E}\left\{\mathbb{E}\left[\frac{Y_{i}(Y_{i} - \mu_{i})}{\sum_{j=1}^{n} Y_{j}^{2}}|\{Y_{j}\}_{j \neq i}\right]\right\}
= \mathbb{E}\left\{\mathbb{E}[h_{i}(\mathbf{Y})(Y_{i} - \mu_{i})|\{Y_{j}\}_{j \neq i}]\right\}
= \mathbb{E}\left\{\mathbb{E}\left[\frac{\partial}{\partial u_{i}}h_{i}(\mathbf{u})\Big|_{\mathbf{u}=\mathbf{Y}}\Big|\{Y_{j}\}_{j \neq i}\right]\right\}
= \mathbb{E}\left[\frac{\partial}{\partial u_{i}}h_{i}(\mathbf{u})\Big|_{\mathbf{u}=\mathbf{Y}}\right]
= \mathbb{E}\left[\frac{\|\mathbf{Y}\|^{2} - 2Y_{i}^{2}}{\|\mathbf{Y}\|^{4}}\right].$$
(6)

Shrinkage Estimation IX



• It follows that we can re-write the mean square error as:

$$MSE(\widetilde{\boldsymbol{\mu}}_{a}, \boldsymbol{\mu}) = n - 2a \mathbb{E}\left[\frac{n\|\boldsymbol{Y}\|^{2} - 2\|\boldsymbol{Y}\|^{2}}{\|\boldsymbol{Y}\|^{4}}\right] + a^{2} \mathbb{E}\left[\frac{1}{\|\boldsymbol{Y}\|^{2}}\right]$$
$$= n + (a^{2} - 2a(n-2)) \mathbb{E}\left[\frac{1}{\|\boldsymbol{Y}\|^{2}}\right]. \tag{7}$$

- The polynomial $p(a) = a^2 2a(n-2)$ is strictly negative in the range $a \in (0, 2n-4)$. This gives part 1.
- On the same range in a p(a) has a unique minimum at a = n 2. This proves part 3.
- For part (2) we note that if $\mu=0$ then $\|\mathbf{Y}\|^2\sim\chi_n^2$.
- Thus we note that $\mathbb{E}\{\frac{1}{\|\mathbf{Y}\|^2}\}=1/(n-2)$ and recall that $n\geq 3$.
- Therefore it follows $MSE(\widetilde{\mu}_{n-2}, 0) = 2$.

Loss functions I



- We can also replace the mean square error by another convex measure of performance.
- Thus we replace $\|\widehat{\theta} \theta\|$ by a choice deviation measure $\mathcal{L}(\widehat{\theta}, \theta)$ called a loss function.
- The expected loss is the the risk:

$$R(\widehat{\theta}, \theta) = \mathbb{E}\{\mathcal{L}(\widehat{\theta}, \theta)\}.$$

 Selecting the right loss function is crucial and this choice must be made carefully.

Loss functions II



- Example: the exponential distribution. Assume that $Y_1, \ldots Y_n \sim \operatorname{Exponential}(\lambda)$ and that $n \geq 2$.
- The MLE of λ is

$$\widehat{\lambda} = \frac{1}{\overline{Y}}.$$

- Here \overline{Y} is the empirical mean.
- We can calculate

$$\mathbb{E}\{\widehat{\lambda}\} = \frac{n\lambda}{n-1}.$$

• It therefore follows that $\widetilde{\lambda}=(n-1)\widehat{\lambda}/n$ is an unbiased estimator. Observe that

$$MSE(\widetilde{\lambda}) < MSE(\widehat{\lambda}).$$

Thus $\widehat{\lambda}$ is dominated by $\widetilde{\lambda}$.

Loss functions III



- ullet λ here takes a positive value.
- In this case quadratic estimation penalises over estimation more heavily than underestimation.
- The maximum possible under estimation is bounded.
- But could we change the loss function?
- Instead we could use

$$\mathcal{L}(a,b) = a/b - 1 - \log(a/b).$$

- Note that for all fixed a $\lim_{b\to 0} \mathcal{L}(a,b) = \lim_{b\to \infty} \mathcal{L}(a,b) = \infty$.
- Now for n > 1 we calculate the risk

$$R(\lambda, \widetilde{\lambda}) = \mathbb{E}_{\lambda} \left[\frac{n\lambda \overline{Y}}{n-1} - 1 - \log \left(\frac{n\lambda \overline{Y}}{n-1} \right) \right]. \tag{8}$$

Loss functions IV



The risk is

$$R(\lambda, \widetilde{\lambda}) = \mathbb{E}_{\lambda} \left[\frac{n\lambda \overline{Y}}{n-1} - 1 - \log \left(\frac{n\lambda \overline{Y}}{n-1} \right) \right]$$

$$= \mathbb{E}_{\lambda} \left[\lambda \overline{Y} - 1 - \log(\lambda \overline{Y}) \right] + \frac{\mathbb{E}_{\lambda} \left[\lambda \overline{Y} \right]}{n-1} - \log \left(\frac{n}{n-1} \right)$$

$$= \mathbb{E}_{\lambda} \left[\lambda \overline{Y} - 1 - \log(\lambda \overline{Y}) \right] + g(n). \tag{9}$$

- To derive the simplification we write $\bar{Y} = \frac{n-1}{n}\bar{Y} + \frac{1}{n}\bar{Y}$.
- Note that $\mathbb{E}_{\lambda}(\bar{Y}) = \lambda^{-1}$. Thus

$$g(n) = \frac{1}{n-1} - \log\left(\frac{n}{n-1}\right).$$

• We claim that g(n) > 0 once $n \ge 2$.

Loss functions V



• Using that $\log(x) = \int_1^x t^{-1} dt$ this follows if

$$\frac{1}{x} > \log(x+1) - \log(x), \quad x > 1$$

$$\Leftrightarrow \frac{1}{x} > \int_{x}^{x+1} t^{-1} dt, \quad x > 1.$$
(10)

This inequality holds by a rectangle area bound on the integral, as follows:

$$\frac{1}{x} = [(1+x)-x]\frac{1}{x} = \int_{x}^{x+1} \frac{1}{x} dt > \int_{x}^{x+1} \frac{1}{t} dt,$$

when x > 1.

• It therefore follows $R(\widetilde{\lambda}, \lambda) > R(\widehat{\lambda}, \lambda)$ and so $\widetilde{\lambda}$ dominates $\widehat{\lambda}$.

Loss functions VI



- We can push generality even further, and obtain an all encompassing framework.
- Called decision theory, it views inference as a game between nature and the statistician.
- Recall our general framework for statistical inference:
 - 1 Model phenomenon by distribution $F(y_1, ..., y_n; \theta)$ for some $n \ge 1$.
 - **2** Distributional form is known but $\theta \in \Theta$ is not known.
 - 3 Observe realisation of (Y_1, \ldots, Y_n) from this distribution.
 - 4 Use (Y_1, \ldots, Y_n) in order to make assertions concerning the true value of , and quantify the uncertainty associated with these assertions.
- The decision theory framework formalises step (4) to include estimation, testing, and confidence intervals.