Invariance Principle for MLE's

If $\eta = \tau(\Theta)$ and $\hat{\Theta}$ is the MLE of Θ , then $\hat{\eta} = \tau(\hat{\Theta})$ is the MLE of η .

Comments: If T(O) is a 1-1 function, this is a trivial theorem.

If T(0) is <u>not</u> 1-1, this is essentially true by definition of induced likelihood. (see)



The usual parameters $\theta = (\mu, \sigma^2)$ are related to the natural parameters $\eta = (\frac{\mu}{\sigma^2}, \frac{1}{2\sigma^2})$ of the 2pef by a 1-1 function : $\eta = T(\theta)$. The likelihood in terms of θ is $L_1(\theta) = (2\pi\sigma^2)^{-n/2} - n\mu^2/2\sigma^2 + \frac{\mu}{\sigma^2}T_1 - \frac{1}{2\sigma^2}T_2$ where $T_1 = \sum x_2$, $T_2 = \sum x_2^2$

Simple Example X1, X2,..., Xn iid Bernoulli(p) X It is known that MLE of p is $\hat{p} = \overline{X}$. Thus (a) MLE of p^2 is $(\hat{p})^2 = \bar{\chi}^2$ (b) MLE of p(I-p) is $\overline{X}(I-\overline{X})$ The function of p in (a) is 1-1, but not I-1 in (b). Definition of Induced Likelihood:

$$\begin{split} & \text{If } \mathcal{N} = \mathcal{T}(\Theta) \text{, then} \\ L^*(\mathcal{N}) \equiv \sup L(\Theta) \\ & \{\Theta : \mathcal{T}(\Theta) = \mathcal{N}\} \\ & \text{If the } ML \in \hat{\mathcal{N}} \text{ of } \mathcal{N} \text{ is defined to be} \\ & \text{the value which } \max \min zes L^*(\mathcal{N}) \text{, then} \\ & \text{it is easily seen that } \hat{\mathcal{M}} = \mathcal{T}(\hat{\Theta}) \text{.} \end{split}$$

The likelihood in terms of
$$\eta$$
 is

$$L_{2}(\eta) = (-\pi/\eta_{2})^{-n/2} e^{n\eta_{1}^{2}/4\eta_{2}} e^{\eta_{1}T_{1}+\eta_{2}T_{2}}$$
obtained by substituting in $L_{1}(\theta)$
 $\mu = -\eta_{1}/2\eta_{2}$, $\sigma^{2} = -\frac{1}{2}\eta_{2}$,
that is, evaluating L_{1} at
 $\theta = (\mu, \sigma^{2}) = (-\eta_{1}/2\eta_{2}, -\frac{1}{2\eta_{2}}) = \tau^{-1}(\eta)$.
Stated abstractly
 $L_{2}(\eta) = L_{1}(\tau^{-1}(\eta))$
so that L_{2} is maximized when
 $\tau^{-1}(\eta) = \hat{\Theta}$, that is, by $\eta = \tau(\hat{\Theta})$.
The MLE of Θ is known to be γ SS/n
 $\hat{\Theta} = (\hat{\mu}, \hat{\sigma}^{2}) = (\bar{x}, \frac{1}{n} \sum (x_{1} - \bar{x})^{2})$
so the invariance principle says the
MLE of η is
 $\hat{\eta} = \tau(\hat{\Theta}) = (\hat{\mu}, \frac{-1}{2\hat{\sigma}^{2}})$.

Continuation of example: what is the MLE of $\alpha = \mu + \sigma^2 ?$ $\alpha = q(\mu, \sigma^2) = \mu + \sigma^2 \pmod{(not 1 - 1)}$ so that $\hat{\alpha} = g(\hat{\mu}, \hat{\sigma}^2) = \hat{\mu} + \hat{\sigma}^2 = \bar{\chi} + \underline{SS}$ What is the MLE of m? 02? with $g_1(x,y) = x$, $g_2(x,y) = y$ we have $\mu = q_1(\Theta)$ $\sigma^2 = q_2(\Theta)$ so that the MLE's are $\hat{\mu} = q_1(\hat{\Theta}) = \overline{X}$ $\hat{\sigma}_{2} = q_{2}(\hat{\Theta}) = SS/n$. Thus, the invariance principle implies: $(\mu, \sigma^2) = (\hat{\mu}, \hat{\sigma}^2)$ $(\mu, \sigma^2) = (\hat{\mu}, \hat{\sigma}^2)$ $MLE of \sigma^2$ MLEofu MLE of pair

MLE for Exponential Families

The invariance principle for MLEs allows us to work with the natural parameter \mathcal{M} (which is a 1-1 function of Θ).

 $\frac{1 \text{pef}}{f(x|\theta)} = c(\theta)h(x) \exp\{w(\theta)t(x)\}$

Natural param: $M = W(\Theta)$

 $f(x|n) = c(n)h(x)exp \{n(x)\}$ $f(x|n) = c(n)h(x)exp \{n(x)\}$ $f(x) = c(w^{-1}(n))$ $f(x|n) = f(x|w^{-1}(n))$

If $X_{1},...,X_{N}$ iid from f(x|n), then $l(n) = N \log c(n) + \sum_{i=1}^{N} \log h(X_{i}) + \eta \sum_{i=1}^{N} t(X_{i})$

T(X)

$$\begin{split} l^{\prime}(\eta) &= N \underbrace{\partial}_{\eta} \log c(\eta) + \underbrace{\sum_{i=1}^{N} t(X_{i})}_{i=1} (x_{i}) \\ & -Et(X_{i}) \\ & by 3.32(a) \\ &= -E[\underbrace{\sum_{i=1}^{N} t(X_{i})]}_{i=1} + \underbrace{\sum_{i=1}^{N} t(X_{i})}_{i=1} \\ & rv^{3}s \quad observed \\ & Values \\ &= -ET(X) + T(x) \\ & x = (x_{1},...,x_{n}) = observed \\ & data \\ & x = (x_{1},...,x_{n}) = random \\ & odata \\ & x = (x_{1},...,x_{n}) = random \\ & odata \\ & x = (X_{1},...,X_{n}) = random \\ & data \\ & x = (X_{1},...,X_{n}) = random \\ & data \\ & x = (X_{1},...,X_{n}) = random \\ & data \\ & x = (X_{1},...,X_{n}) = random \\ & data \\ & x = (X_{1},...,X_{n}) = random \\ & data \\ & x = (X_{1},...,X_{n}) = random \\ & data \\ & x = (X_{1},...,X_{n}) = random \\ & data \\ & x = (X_{1},...,X_{n}) = random \\ & data \\ & x = (X_{1},...,X_{n}) = random \\ & data \\ & x = (X_{1},...,X_{n}) = random \\ & data \\ & x = (X_{1},...,X_{n}) = random \\ & x =$$



is automatically a global maximum so long as (+)* is convex.

In one dimension ($\Theta \subset IR$), this means Θ^* must be an interval of some sort (can be infinite).

Ignoring the fine print, For a lpef, the log-likelihood will have a unique stationary point which will be the MLE.

kpef $f(x|\theta) = c(\theta)h(x)exp\left\{\sum_{i=1}^{k} w(\theta)t(x)\right\}$ Natural param: $\mathcal{N} = (\mathcal{N}_1, \dots, \mathcal{N}_k) = (W_1(\Theta), \dots, W_k(\Theta)),$ that is, $\mathcal{M}_{i} = \mathcal{W}_{i}(\Theta)$. $f(x|n) = c(n)h(x)exp\left\{\sum_{j=1}^{k} \eta_{j} t_{j}(x)\right\}$ $\int f^{*} \int c^{*}$ If X1,..., XN iid from f(x/7), then $l(n) = N \log c(n) + \sum_{i=1}^{N} \log h(x_i)$ $+\sum_{i=1}^{K} \eta_{i} \left(\sum_{j=1}^{N} t_{j}(x_{i}) \right)$ $\frac{\partial l}{\partial n_i} = N \frac{\partial}{\partial n_i} \log c(n) + \sum_{i=1}^{N} t_i(x_i)$ $= -E t_{i}(X_{i})$ $= -E \sum_{i=1}^{N} t_{i}(X_{i}) + \sum_{i=1}^{N} t_{i}(x_{i})$

 $\frac{\partial^{-\chi}}{\partial n_{j} \partial n_{\ell}} = N\left(\frac{\partial^{2}}{\partial n_{j} \partial n_{\ell}}\log c(n)\right)$

 $= N(-Cov(t_i(X_i), t_g(X_i)))$

proved case j=l in homework

Thus, the equations for a stationary point

$$\frac{\partial \ell}{\partial \eta_j} = 0$$
 for $j = 1, \dots, k$

are equivalent to

$$E_{\eta}T_{j}(X) = T_{j}(x)$$
 for $j = 1, ..., k$ (‡)
where $T_{j}(X) = \sum_{i=1}^{N} t_{j}(X_{i})$ and $T_{j}(x) = \sum_{i=1}^{N} t_{j}(x_{i})$

or in vector notation

$$E_{\eta}T(X) = T(x)$$
 where
 $T(X) = (T_1(X), ..., T_k(X))$ and $T(x) = (T_1(x), ..., T_k(x))$.

The Hessian matrix
$$H(\eta) = \left(rac{\partial^2 \ell}{\partial \eta_i \partial \eta_j}
ight)_{i,j=1}^k$$
 is given by

$$H(\eta) = -N\Sigma(\eta)$$

where $\Sigma(\eta)$ is the $k \times k$ covariance matrix of $(T_1(X_1), T_2(X_1), \ldots, T_k(X_1))$. A covariance matrix will be positive definite (except in degenerate cases), so that $H(\eta)$ will be negative definite for all η .

Conclusion: An interior stationary point (i.e., a solution of (\ddagger)) must be the unique global maximum, and hence the MLE.

This result also holds in the original parameterization with (‡) restated as $E_{\theta}T_j(X) = T_j(x)$, j = 1, ..., k.

Connection with MOM: For a 1pef with t(x) = x, MOM and MLE agree. For a kpef with $t_j(x) = x^j$, MOM and MLE agree. Why? Because then (‡) is equivalent to the equations for the MOM estimator.

Revisiting Gamma Example:

The Gamma family is a 2pef (or a 1pef if α or β is held fixed).

Switching to the natural parameters $\eta_1 = \alpha - 1$ and $\eta_2 = -1/\beta$ (or just making the substitution $\lambda = 1/\beta$) simplifies the second derivatives w.r.t. η_2 (or λ) and makes the sufficient condition for a stationary point to be the global max hold.

The system of equations for the MLE of (α, β) may be easily derived directly from (‡).

MLE's for More General Exponential Families

Proposition: If $X \sim P_{\theta}$, $\theta \in \Theta$ where P_{θ} has a joint pdf (pmf) from an *n*-variate *k* parameter exponential family (nvkpef):

$$f(x \mid \theta) = c(\theta)h(x) \exp\left\{\sum_{j=1}^{k} w_j(\theta)T_j(x)
ight\}$$

for $x \in \mathbb{R}^n$, $\theta \in \Theta \subset \mathbb{R}^k$,

then the MLE of θ based on the observed data x is the solution of the system of equations

$$E_{\theta}T_{j}(X) = T_{j}(x)$$
 for $j = 1, \dots, k$ (Solve for θ)

providing this solution (call it $\hat{\theta}$) satisfies

$$w(\hat{\theta}) \in \text{interior of } \{w(\theta) : \theta \in \Theta\}.$$

Proof: Essentially the same as for the ordinary kpef.

Example: Simple Linear Regression with Known Variance Y1, Y2, ..., Yn independent $Y_{i} \sim N(\beta_{0} + \beta_{1} \chi_{i}, \sigma_{0}^{2}), \Theta = (\beta_{0}, \beta_{1}).$ Joint disth. of $Y = (Y_1, Y_2, ..., Y_n)$ forms exponential family. Natural sufficient statistic is $t(\underline{Y}) = \left(\sum_{i} Y_{i}, \sum_{i} \chi_{i}Y_{i}\right).$ $E_{A}t(Y) = t(Y_{obs})$ has the form $E(\Sigma Y_{i}) = \Sigma Y_{i}$ $E\left(\sum \chi_{i}Y_{i}\right) = \sum \chi_{i}Y_{i}$ Thus MLE $\hat{\Theta} = (\hat{\beta}_0, \hat{\beta}_1)$ is solution of $\sum_{i} (\beta_{0} + \beta_{i} \chi_{i}) = \sum_{i} y_{i}$ $\sum \chi_i(\beta_0 + \beta_1 \chi_i) = \sum \chi_i y_i$

Sufficient Statistics and MLE's

If T = T(X) is a sufficient stat. for Θ , then there is an MLE which is a function of T. (If the MLE is unique, then we can say the MLE is a function of T.) Proof: By FC, $f(x|\theta) = g(T(x), \theta) h(x)$. Assume for convenience the MLE is unique. Then the MLE is $\hat{\Theta}(x) = \arg \max_{A} f(x|\Theta)$ = $\operatorname{argmax}_{\Theta} g(T(x), \Theta)$ which is clearly a function of T(x).



For any 03 (fixed arbitrary value), Maximizing L(B, 02) with respect to B $\iff \operatorname{Minimizing} \underbrace{\sum_{i=1}^{n} (y_i - g(x_i, \beta))^2}_{i=1}$ with respect to B.

MLE and Least squares give same estimates of p parameters.