# Empirical Likelihood (EL)

Art Owen (1988, 1990) introduced.

This is nonparametric methodology for creating likelihood-type inference without specifying a joint distributional form for the data.

$\implies$ we can't misspecify!

EL is going to use the fact that the empirical cdf is a nonparametric MLE to assess how plausible a value of a parameter is to perform inference.

↳ without making distributional assumptions!

# 1 Mean Case

Suppose $\boldsymbol{Y}_1, \ldots, \boldsymbol{Y}_n$ are iid with mean $\boldsymbol{\mu}$ and covariance-variance $\Sigma$. For simplicity, say we are interested in estimating $\boldsymbol{\mu}$.

Imagine assigning probabilities $p_1, \ldots, p_n$ to the data $\underline{Y}_1, \ldots, \underline{Y}_n$ where $0 \leq p_i \leq 1$ and $\sum_{i=1}^{n} p_i = 1$.

$p_i \longmapsto \underline{Y}_i$

$(*)$.

Unlike parametric likelihood, where we assume a functional form for $p_i$'s, only constraints $(*)$.

Define a multinomial likelihood $\prod_{i=1}^{n} p_i$ ( likelihood for $\underline{Y}_1, \ldots, \underline{Y}_n$ using $p_1, \ldots, p_n$).

Recall from class ( likelihood notes pg 10) if you maximize $\prod_{i=1}^{n} p_i$ for $p_1, \ldots, p_n$ the maximizer $p_1 = p_2 = \ldots = p_n = \frac{1}{n}$ ← 1 observation in each "class" and we have also seen the empirical cdf

$$F_n(\underline{y}) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{I}(\underline{Y}_i \leq \underline{y}) \qquad \underline{y} \in \mathbb{R}^q \quad \text{is the MLE (pg 23 likelihood notes).}$$

in other words, given the data the empirical cdf maximizes $\prod_{i=1}^{n} p_i$.

To perform *nonparametric* likelihood inference on $\boldsymbol{\mu}$, we can consider a constrained multinomial likelihood, known as the **Empirical Likelihood function of $\boldsymbol{\mu}$:**

$$L_n(\overset{\text{function of }\underline{\mu}}{\underset{\text{EL function}}{\underbrace{\boldsymbol{\mu}}}}|\boldsymbol{Y}) = \sup\left\{\underset{\text{multinomial likelihood}}{\underbrace{\prod_{i=1}^{n} p_i}} : p_i \mapsto \boldsymbol{Y}_i, \overset{p_i \geq 0}{\underset{}{\sum_{i=1}^{n} p_i = 1}}, \overset{\text{mean of a dsn } (p_1,\dots,p_n) \text{ on } (\underline{Y}_1,\dots,\underline{Y}_n)}{\underset{\text{mean constraint on } (p_1,\dots,p_n)}{\underbrace{\sum_{i=1}^{n} \boldsymbol{Y}_i p_i = \boldsymbol{\mu}}}}\right\}.$$

Given a parameter value $\underline{\mu}$ and $\underline{Y}$, $L_n(\underline{\mu}|\underline{Y})$ assesses how plausible the value of $\underline{\mu}$ is.

$L_n(\underline{\mu}|\underline{Y})$ is the largest multinomial likelihood possible for a probability assignment to the data having mean $\underline{\mu}$.

The largest possible value of $L_n(\boldsymbol{\mu}|\boldsymbol{Y})$ is

$$\underset{\|}{\underbrace{\prod_{i=1}^{n} \frac{1}{n}}} \Rightarrow \underline{\mu} \text{ would be } \sum_{i=1}^{n} Y_i \cdot \frac{1}{n} \Rightarrow \underline{\mu} \text{ would be } \overline{Y}.$$

$$L_n(\overline{\underline{Y}}, \underline{Y}).$$

So $\overline{Y} = \frac{1}{n}\sum_{i=1}^{n} \underline{Y}_i$ is a nonparametric ML estimator of $\underline{\mu}$, i.e. the EL estimator $\hat{\underline{\mu}} = \overline{Y}$ of $\mu$.

# 2 Statistical Inference

We can form an EL ratio for $\boldsymbol{\mu}$

$$R_n(\boldsymbol{\mu}) = \frac{L_n(\boldsymbol{\mu}|\boldsymbol{Y})}{L_n(\hat{\boldsymbol{\mu}}|\boldsymbol{Y})}$$

$$= \frac{L_n(\underline{\mu}|\underline{y})}{\prod\limits_{i=1}^{\hat{n}} \frac{1}{n}} \qquad \hat{\underline{\mu}} = \bar{\underline{y}} \Rightarrow p_i = \frac{1}{n}$$

$$= n^n L_n(\underline{\mu}|\underline{y})$$

$$= \sup\left\{ \prod_{i=1}^{\hat{n}} n p_i : p_i \geq 0, \sum_{i=1}^{\hat{n}} p_i = 1, \underbrace{\sum_{i=1}^{\hat{n}} \underline{y}_i p_i = \mu}_{} \right\}$$

$$\underbrace{\sum_{i=1}^{\hat{n}} (\underline{y}_i - \mu) p_i = 0}_{\text{look familiar??}}$$

**Theorem (Wilk's Theorem):** If $\boldsymbol{Y}_1, \ldots, \boldsymbol{Y}_n \in \mathbb{R}^q$ are iid with mean $\boldsymbol{\mu}_0$ and covariance-variance $\Sigma$ where $\text{rank}(\Sigma) = q$, then

$$-2\log R_n(\boldsymbol{\mu}_0) \xrightarrow{d} \chi^2_q \text{ as } n \to \infty.$$

In other words, for $H_0 : \underline{\mu} = \underline{\mu}_0 \in \mathbb{R}^q$, if $H_0$ is true then $-2\log R_n(\underline{\mu}_0) \xrightarrow{d} \chi^2_q$ as $n \to \infty$.

✳ EL behaves exactly like parametric likelihood for log ratios! ✳

So if $\chi^2_{1-\alpha, q}$ denotes the $1-\alpha$ quantile of $\chi^2_q$, then an approximate $100(1-\alpha)\%$ confidence region for $\underline{\mu}$:

$$CR = \left\{ \underline{\mu} \in \mathbb{R}^q : -2\log R_n(\underline{\mu}) \leq \chi^2_{1-\alpha, q} \right\}.$$

by inverting the EL test

$$P(\underline{\mu}_0 \in CR) = P\left(-2\log R_n(\underline{\mu}_0) \leq \chi^2_{1-\alpha, q}\right) \xrightarrow{as \ n \to \infty} P\left(\chi^2_q \leq \chi^2_{1-\alpha, q}\right) = 1-\alpha \ //.$$

For proof of this theorem, see Owen (1988).

# 3 EL with Estimating Equations

(Qin and Lawless, 1994).

Recall:

For $\underline{Y}_1, \dots, \underline{Y}_n \in \mathbb{R}^q$ iid and $\underline{\theta} \in \mathbb{R}^b$ a parameter of interest

Estimating equations link a data point $\underline{Y}_i$ to parameters through $r \geq b$ functions.

$$\underline{\Psi}(\underline{Y}_i, \underline{\theta}) \quad \text{which satisfy} \quad E\underline{\Psi}(\underline{Y}_i, \underline{\theta}) = \underline{0}_r.$$

For EL inference on $\boldsymbol{\theta} \in \mathbb{R}^b$, we make an EL function

extends mean example to any estimating equation!

$$L_n(\underline{\theta}) = \sup \left\{ \prod_{i=1}^{n} p_i : p_i \geq 0, \sum_{i=1}^{n} p_i = 1, \sum_{i=1}^{n} p_i \underline{\Psi}(\underline{Y}_i, \underline{\theta}) = \underline{0}_r \right\}$$

↑ given value of $\underline{\theta}$

$p_i$'s are placed on $\underline{\Psi}(\underline{Y}_i, \underline{\theta})$ to have expectation zero.

The EL function evaluates the plausibility of a given value of $\underline{\theta}$ based on the data.

Then we can get a point estimate, EL ratio, and corresponding CIs, as well as "profile" EL:

→ regions

point estimate : maximize $L_n(\underline{\theta})$ to obtain maximum EL estimator $\hat{\underline{\theta}}$

EL ratio : $R_n(\underline{\theta}) = \dfrac{L_n(\underline{\theta})}{L_n(\hat{\underline{\theta}})}$ (just like parametric likelihood)

Credible region: $CR = \left\{ \underline{\theta} \in \mathbb{R}^b : -2 \log R_n(\underline{\theta}) \leq \chi^2_{1-\alpha, b} \right\}$ (invert EL ratio).

profile EL: suppose $\underline{\theta} = (\underline{\theta}_1, \underline{\theta}_2)$, $\underline{\theta}_1 \in \mathbb{R}^s$, $\underline{\theta}_2 \in \mathbb{R}^{b-s}$. Given $\underline{\theta}_1$ define $\hat{\underline{\theta}}_{2,\theta_1}$ where

$$L_n\left(\underline{\theta}_1, \hat{\underline{\theta}}_{2,\theta_1}\right) = \sup_{\theta_2} L_n(\underline{\theta}_1, \underline{\theta}_2)$$

Then the profile EL ratio for $\underline{\theta}_1$ is $R_n(\underline{\theta}_1) = \dfrac{L_n\left(\underline{\theta}_1, \hat{\underline{\theta}}_{2,\theta_1}\right)}{L_n(\hat{\underline{\theta}})}$.

*Main EL result*

**Theorem:** Suppose $Y_1, Y_2, \cdots \in \mathbb{R}^q$ are iid with $\mathrm{E}\psi(Y_1, \theta_0) = \mathbf{0}_r$ and $\mathrm{Var}[\psi(Y_1, \theta_0)] = \mathrm{E}\psi(Y_1, \theta_0)\psi(Y_1, \theta_0)^\top$ is positive definite, where $\theta_0$ denotes the true parameter value.

Suppose also that $\partial\psi(y, \theta)/\partial\theta$ and $\partial^2\psi(y, \theta)/\partial\theta\partial\theta^\top$ are continuous in a neighborhood of $\theta_0$ and that, in this neighborhood, $||\psi(Y_1, \theta)||^3$, $||\partial\psi(y, \theta)/\partial\theta||$ and $||\partial^2\psi(y, \theta)/\partial\theta\partial\theta^\top||$ are bounded by an integrable function $\Psi(Y_1)$.

Finally, suppose the $r \times b$ matrix $D_\psi \equiv \mathrm{E}\partial\psi(y, \theta)/\partial\theta$ has full column rank $b$.

Then, as $n \to \infty$,

  i. $\sqrt{n}(\hat{\theta} - \theta_0) \overset{d}{\to} N(\mathbf{0}_b, V)$, where $V = (D_\psi^\top \mathrm{Var}[\psi(Y_1, \theta_0)]D_\psi)^{-1}$.   *EL point estimates are asymptotically Normal*

  ii. If $r > b$, the asymptotic variance $V$ cannot increase if an estimating function is added.   *or decrease if an estimating function is dropped.*

  iii. To test $H_0 : \theta = \theta_0$, we may use $-2\log R_n(\theta_0)$ and when $H_0$ is true,

$$-2\log R_n(\theta_0) \overset{d}{\to} \chi^2_{b}  \quad \text{\# parameters}$$

$$R_n(\theta_0) = \frac{L_n(\theta_0)}{L_n(\hat{\theta})}$$

$\Rightarrow$ *confidence regions:* $CR = \{\theta \in \mathbb{R}^b : -2\log R_n(\theta) \leq \chi^2_{b,1-\alpha}\}.$

  iv. If $r > b$, to test $H_0 : \underline{\mathrm{E}\psi(Y_1, \theta) = \mathbf{0}_r}$ holds for some $\theta$, we may use
  
  *more functions than parameters.*     *moment condition*

$$-2\log \frac{\frac{L_n(\hat{\theta})}{n}}{\prod_{i=1}^{n}(1/n)} = -2\log(n^n L_n(\hat{\theta})).$$

*biggest possible value could ever have for an EL function w/ no moment constraints*

and when $H_0$ is true this quantity converges in distribution to $\chi^2_{r-b}$.

*\# excess estimating functions.*

*Asymptotically, $-2\log R_n(\theta_0)$ and $-2\log n^n L_n(\hat{\theta}))$ are independent.*

  v. To test the profile assumption $H_0 : \theta_1 = \theta_1^0 \in \mathbb{R}^s$, we can use the profile EL ratio

$-2\log R_n(\theta_1^0)$ and , when $H_0$ is true, $-2\log R_n(\theta_1^0) \overset{d}{\to} \chi^2_s$;

*\# parameters in EL function after profiling.*