

Bayesian Empirical Likelihood: a survey

Art B. Owen

Stanford University

These are the slides I presented at BNP 11 in Paris, a very fine meeting. I have added a few interstitial slides like this one to include comments I made while presenting the slides and to give some background thoughts in this area.

The session included a delightful presentation by Nils Lid Hjort on some not yet published theorems where he connects Bayesian empirical likelihood to Dirichlet process Bayesian nonparametrics.

I was pleased to see our session got a positive review on Xian's Og (written by our session chair, Christian Robert).

Also, the conference included two very strong posters by students working on Bayesian empirical likelihood. [Frank Liu](#) (Monash) has a nice way to robustify BETEL. [Laura Turbatu](#) (Geneva) took a close look at high order asymptotics to compare frequentist coverage of Bayes with various kinds of nonparametric priors.

The goal

Science driven prior \times Data driven likelihood

This goal is not everybody's goal. But this goal means you don't have to make up a likelihood which could bring unwanted and consequential restrictions.

If the data deliver the likelihood and the prior is chosen flat or by previously known science, then what is there left for the statistician to do?

The goal

Science driven prior \times Data driven likelihood

The ask

IID data and
an estimand defined by estimating equations.

$$\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n \stackrel{\text{iid}}{\sim} F$$
$$\int m(\mathbf{x}, \theta) dF(\mathbf{x}) = 0, \quad \theta \in \mathbb{R}^p$$

We have IID X and an estimating equation. If you picture the plate diagram, it is one of the simplest possible ones. It covers an important set of problems, but Bayes includes much else as well.

Outline

- 1) Empirical Likelihood
- 2) Bayes EL Lazar (2003)
- 3) Exponential tilting, and ETEL Schennach (2005,2007)
- 4) ABC and EL Mengersen, Pudio, Robert (2013)
- 5) BEL with Hamiltonian MCMC Chaudhuri et al (2017)
- 6) Some other recent works

Notes

- Apologies about necessarily omitted papers.
- Motivation: bioinformatics.

In the bioinformatics motivation, the data are counts. They are assuredly not Poisson. So a common response is to model them as negative binomial. Of course they are not that either. This is work in progress.

SAMSI is having a year on QMC. The original slides touted the opening workshop. However the application deadline will almost surely have passed by the time you're reading this.

Empirical likelihood

For $\mathbf{X}_i \stackrel{\text{iid}}{\sim} F$ observed $\mathbf{X}_i = \mathbf{x}_i \in \mathbb{R}^d$, define

$$L(F) = \prod_{i=1}^n F(\{\mathbf{x}_i\}).$$

The nonparametric MLE (NPMLE) is

$$\hat{F} \equiv \arg \max_F L(F) = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}$$

There are many other NPMLEs.

E.g., [Kaplan-Meier](#) for right censored survival,

Log concave densities: [Walther](#), [Dumbgen](#), [Rufibach](#), [Cule](#), [Samworth](#), [Stewart](#) · · ·

Basic notions

Model: $x_i \stackrel{\text{iid}}{\sim} F$ True value: $F = F_0$ NPMLE: \hat{F}

Estimand $T(F)$

E.g., population mean, median, regression coefficient etc.

True value: $\theta_0 = T(F_0)$ NPMLE: $\hat{\theta} = T(\hat{F})$

Empirical likelihood ratios

	Point est.	Interval / Test
Parametric	$\hat{\theta} = \arg \max_{\theta} R(\theta)$	$-2 \log R(\theta_0) \rightarrow \chi^2$
Non-parametric	$T(\hat{F})$	* * *

The purpose of empirical likelihood ratios is to fill in the * * *,
getting likelihood ratio confidence intervals and tests without parametric assumptions.

I.e., a nonparametric Wilks' theorem

○ (2000) Chapman & Hall, Monograph

Empirical likelihood ratio

Let $w_i = w_i(F) = F(\{\mathbf{x}_i\})$

Then $w_i \geq 0$ and $\sum_{i=1}^n w_i \leq 1$.

$$R(F) = \frac{R(F)}{R(\hat{F})} = \prod_{i=1}^n \frac{w_i}{1/n} = \prod_{i=1}^n (nw_i)$$

Literal EL confidence regions

$$\{T(F) \mid R(F) \geq c\}$$

Profile likelihood

$$\mathcal{R}(\theta) = \max\{R(F) \mid T(F) = \theta\}$$

$$\text{Conf set} = \{\theta \mid \mathcal{R}(\theta) \geq c\}.$$

We still have to choose c

Usual EL

It is usual to force $\sum_i w_i = \sum_i F(\{\mathbf{x}_i\}) = 1$.

$R(F) \geq c$ already implies $1 - F(\{\mathbf{x}_1, \dots, \mathbf{x}_n\}) = O(1/n)$

It is a multinomial likelihood

Supported on $\mathbf{x}_1, \dots, \mathbf{x}_n$

There are other choices

Empirical likelihood for the mean

$$T(F) = \int \mathbf{x} \, dF(\mathbf{x})$$

Profile likelihood

$$\mathcal{R}(\mu) = \max \left\{ \prod_{i=1}^n n w_i \mid \sum_{i=1}^n w_i \mathbf{x}_i = \mu, w_i \geq 0, \sum_{i=1}^n w_i = 1 \right\}$$

Confidence set

$$\{\mu \mid \mathcal{R}(\mu) \geq c\}$$

To pick c

$$-2 \log(\mathcal{R}(\mu_0)) \xrightarrow{d} \chi_{(k)}^2, \quad k = \text{rank}(\text{Var}(\mathbf{x}))$$

$$c = \exp\left(-\frac{1}{2} \chi_{(k)}^{2, 1-\alpha}\right)$$

Dipper, *Cinclus cinclus*

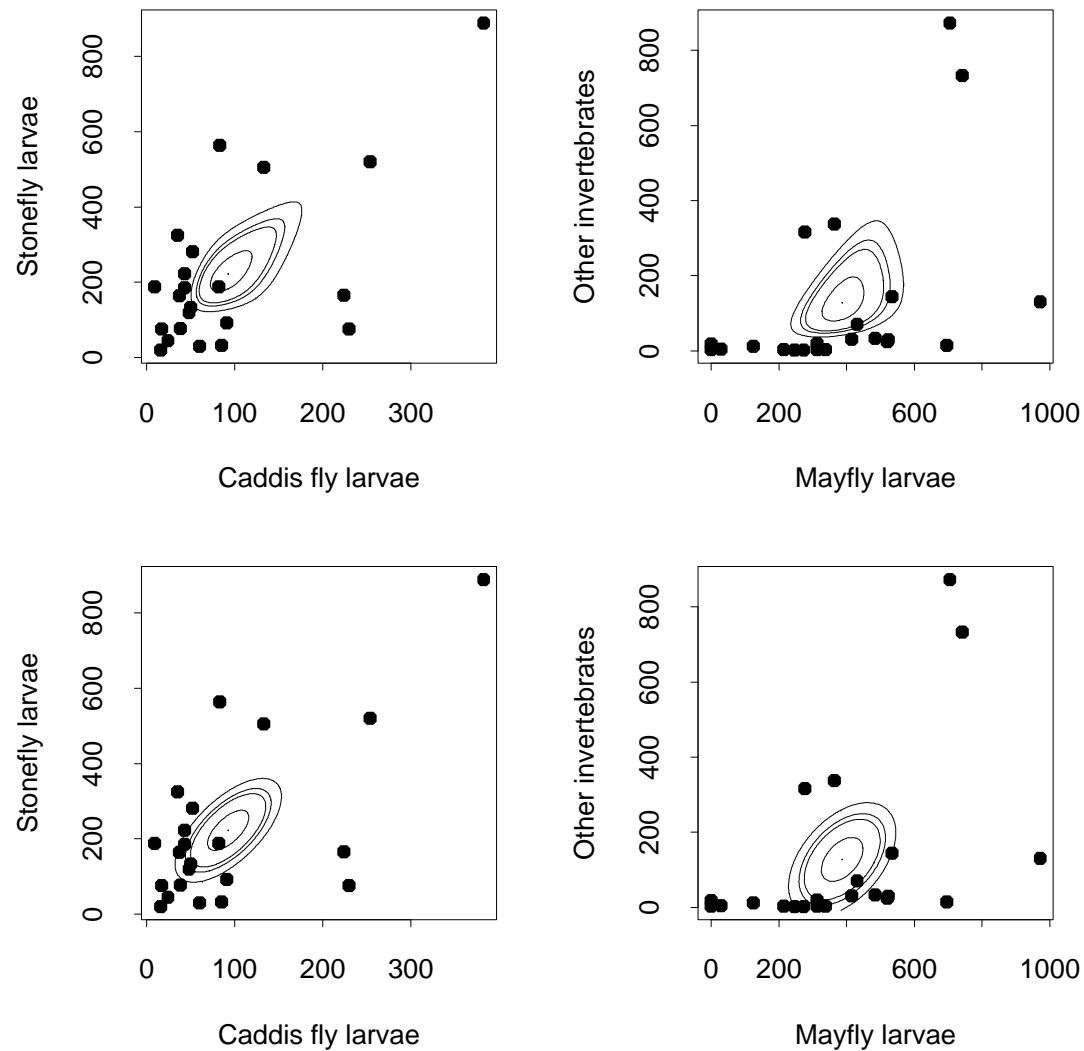


Eats larvae of Mayflies, Stoneflies, Caddis flies, other

Photo By Mark Medcalf. From Wikipedia. Lic under CC by 2.0.

<https://commons.wikimedia.org/w/index.php?curid=15739681>

Dipper diet means



Top row shows EL; bottom Hotelling's T^2 ellipses

Data from Iles (1993) 22 sites in Wales.

Hall (1990) quantifies correctness of shape.

Computing EL for the mean

Easy by convex duality.

The dual is self-concordant

⇒ truncated Newton assured of convergence.

○ (2013) Canadian Journal of Statistics

R code online

Estimating equations

$$\theta \text{ solves } \mathbb{E}(m(\mathbf{x}, \theta)) = 0$$

Examples

$m(\mathbf{x}, \theta)$	Estimand θ
$\mathbf{x} - \theta$	$\mathbb{E}(\mathbf{x})$
$1_{x < \theta} - 0.5$	median(x)
$(\mathbf{x} - \theta) \times 1_{z \in A}$	$\mathbb{E}(\mathbf{x} \mid z \in A)$
$\mathbf{x}(y - \mathbf{x}^\top \beta)$	regression
$\frac{\partial}{\partial \theta} \log f(\mathbf{x}, \theta)$	MLE

Note that some of these estimands can be well defined without making any assumptions about the likelihood. For instance, we could well want a median without wanting to assume a double exponential distribution or any other whose MLE is the median.

Similarly, least squares estimates are interpretable as sample based counterparts to population least squares quantities.

Other cases, like logistic regression, make it harder to interpret what you get if the motivating parametric model does not hold. That said, you might still prefer reliable to unreliable uncertainty quantification for such an estimand.

EL properties

- High power Kitamura (2001), Lazar & Mykland (1998)
Comparable to a true parametric model (when there is one)
- Bartlett correctable DiCiccio, Hall & Romano (1991)
Coverage errors $O(1/n) \rightarrow O(1/n^2)$

Challenges

- Convex hull issue, fixable by pseudo-data,
Chen & Variyath (2008), Tsao & Wu (2013/14), Emerson & O (2009)
- Computing $\max\{\mathcal{R}(\theta_1, \dots, \theta_p) \mid \theta_j = \theta_{j0}\}$ can be hard
O (2000) used expensive NPSOL solver

Possibility

Sampling might work more easily than optimization

Overdetermined estimating equations

$$\theta \in \mathbb{R}^p, \quad m(\mathbf{x}, \theta) \in \mathbb{R}^q, \quad q > p$$

$$q \text{ equations in } p < q \text{ unknowns: } \frac{1}{n} \sum_{i=1}^n m(\mathbf{x}_i, \theta) = 0$$

Popular in econometrics

E.g., generalized method of moments (GMM) Hansen (1982)

E.g., regression through $(0, 0)$

$$\mathbb{E}(y - x\beta) = \mathbb{E}(x(y - x\beta)) = 0$$

Misspecification

Sometimes no $\theta \in \mathbb{R}^p$ gives $\mathbb{E}(m(\mathbf{x}, \theta)) = 0 \in \mathbb{R}^q$

EL for overdetermined

Get n more multinomial parameters w_1, \dots, w_n

NPMLE

$$\max_{\theta, \mathbf{w}} \prod_{i=1}^n w_i \quad \text{subject to} \quad \sum_{i=1}^n w_i m(\mathbf{x}_i, \theta) = 0, \quad \sum_{i=1}^n w_i = 1, \quad w_i > 0$$

Then solve $\sum_i \hat{w}_i m(\mathbf{x}_i, \theta) = 0$ for $\hat{\theta}$

Newey & Smith (2004): same variance & less bias than GMM

Bayesian EL

$$p(\theta | \mathbf{x}) \propto p(\theta) \mathcal{R}(\theta | \mathbf{x})$$

Lazar (2003) justifies by a Bayesian CLT, like Boos & Monahan (1992)

Also Hjort (Today!)

Is that Bayesian enough?

Cressie-Read likelihoods

There are numerous alternative nonparametric likelihoods.

Chang & Mukerjee (2008) give conditions for them to have a very accurate BCLT for very general priors.

EL is uniquely able to do so.

Some others will work with a flat prior.

Also

Poster by Turbatu (yesterday)

Least favorable families

$$\text{Pr}_{\text{EL}}(\mathbf{x}_i) = w_i = \frac{1}{n} \frac{1}{1 + \lambda(\theta)^\top m(\mathbf{x}_i, \theta)}$$

p -dimensional $\theta \implies p$ -dimensional $\lambda(\theta)$.

There is a least favorable p -dimensional family for θ .

Any other family makes the inference artificially easy.

EL family is asymptotically the LFF [DiCiccio & Romano \(1990\)](#)

The connection

$p(\theta)$ induces a prior on the LFF, and

$\mathcal{R}(\theta \mid \mathbf{x}) \doteq$ likelihood on LFF

$\implies p(\theta) \times \mathcal{R}(\theta \mid \mathbf{x}) \doteq$ posterior on LFF

$=$ posterior on θ

To my knowledge nobody has tracked down exactly how the approximations work out for the least favorable family motivation of Bayesian EL.

Exponential tilting (ET)

Schennach (2005)

To get an empirical probability model, for

$$\theta \in \Theta \subset \mathbb{R}^p, \quad \text{partition } \Theta = \bigcup_{\ell=1}^N \Theta_\ell$$

$p(\theta)$ is multinomial on N disjoint cells.

Now let $N \rightarrow \infty$ for universal approximation.

Ultimately

$$p(\theta | \mathbf{x}) \propto p(\theta) \times \prod_{i=1}^n w_i^*(\theta), \quad \text{where}$$

$$w_i^* \quad \text{maximize} \quad - \sum_{i=1}^n w_i \log(w_i)$$

$$\text{s.t.} \quad \sum_{i=1}^n w_i^* \mathbf{x}_i = \theta, \quad \sum_i w_i = 1$$

Exponential tilting EL

Schennach (2007) Form parametric family via exponential tilting

$$w_i^*(\theta) = \frac{e^{\lambda(\theta)^\top m(\mathbf{x}_i, \theta)}}{\sum_{j=1}^n e^{\lambda(\theta)^\top m(\mathbf{x}_j, \theta)}}, \quad \text{where}$$

$$0 = \sum_{i=1}^n e^{\lambda(\theta)^\top m(\mathbf{x}_i, \theta)} m(\mathbf{x}_i, \theta)$$

Likelihood

$$L(\theta) = \prod_{i=1}^n n w_i^*(\theta)$$

ETEL

$$\hat{\theta}_{\text{ETEL}} = \arg \max_{\theta} L(\theta)$$

Works better on misspecified overdetermined models,

i.e., no $\theta \in \mathbb{R}^p$ makes $\mathbb{E}(m(\mathbf{x}, \theta)) = 0 \in \mathbb{R}^q$

Reason: at least there is an estimand.

If we know that no θ can make $\mathbb{E}(m(\mathbf{x}, \theta)) = 0$ then I'm not sure why we still want an estimand. It seems like we should instead change the choice of m .

BETEL

Chib, Shin, Simoni (2016,2017)

Use a slack parameter V_k .

$V_k \neq 0$ for non-truthfulness of moment condition k .

At most $m - p$ nonzero V_k .

Metropolized Independence Sampling from a t distribution.

Asymptotic model selection

Use log marginal likelihood.

If there is a true model then it will be selected over any false one.

Among true models: most equations holding wins.

Robust BETEL

poster by Liu (yesterday)

EL with ABC

Mengersen, Pudlo & Robert (2013)

Basic EL-ABC

Sample $\theta_i \sim p(\theta)$

Compute $w_i = \mathcal{R}(\theta_i | \mathbf{x})$

Normalize the weights

Advanced EL-ABC

Use adaptive multiple importance sampling (AMIS)

Algorithm has 5 loops

Max depth 3

Hamiltonian MCMC

Chaudhuri, Mondal & Yin (2017)

Using $p(\theta)\mathcal{R}(\theta | \boldsymbol{x})$ they report that:

- Gibbs is not applicable
- Random walk Metropolis has problems with irregular support of $\mathcal{R}(\theta | \boldsymbol{x})$

Their solution

Use gradients, Hamiltonian MCMC.

Gradient gets steep just where you need it.

Posterior moments

Vexler, Tao & Hutson (2014)

For $x_i \in \mathbb{R}$

$$\frac{\int_{x_{(1)}}^{x_{(n)}} \theta^k \mathcal{R}(\theta | \mathbf{x}) p(\theta) d\theta}{\int_{x_{(1)}}^{x_{(n)}} \mathcal{R}(\theta | \mathbf{x}) p(\theta) d\theta}$$

Nonparametric James-Stein

For $\mathbf{x}_i = (x_{i1}, \dots, x_{ik}) \in \mathbb{R}^K$, get $\mathcal{R}_j(\theta_j) = \mathcal{R}(\theta_j | x_{1j}, \dots, x_{nj})$

$$\hat{\theta}_j = \frac{\int \theta_j \mathcal{R}_j(\theta) \varphi\left(\frac{\theta_j - \theta_0}{\sigma_*}\right) d\theta_j}{\int \mathcal{R}_j(\theta) \varphi((\theta_j - \theta_0)/\sigma_*) d\theta_j}$$

$$\sigma_*^2 = \arg \max_{\sigma^2} \sum_{j=1}^K \log\left(\frac{1}{\sqrt{2\pi}\sigma} \int \mathcal{R}_j(\theta_j) e^{-\theta_j^2/(2\sigma^2)} d\theta_j\right)$$

Quantiles

Vexler, Yu & Lazar (2017)

EL Bayes factors for quantile $\theta \equiv Q^\alpha = Q_0$ vs $Q^\alpha \neq Q_0$

They overcome a technical challenge handling $\int_{H_A} \mathcal{R}(\theta | \mathbf{x}) p(\theta) d\theta$ because \mathcal{R} has jump discontinuities.

Also two sample Bayes factors.

If scientific or regulatory interest is in a specific quantile then the user does not have to pick a whole parametric family.

Quantile regression

$$\Pr(Y \leq \alpha + \beta x \mid x) = \tau, \quad \text{e.g., } \tau = 0.9$$

Estimating equations

$$0 = \mathbb{E}(1_{y \leq \alpha + \beta x} - \tau)$$

$$0 = \mathbb{E}(x(1_{y \leq \alpha + \beta x} - \tau))$$

Lancaster & Jun (2010)

BETEL for quantiles similar to Jeffrey's prior

Yang & He (2012)

BCLT accounting for non smoothness of estimating equations

Shrinking priors pool β for multiple τ

Both use MCMC

Summary

We can do Bayes without requiring a parametric family for the observations.

It is the “other nonparametric Bayes”.

Thanks

- Luke Bornn, invitation
- Sanjay Chaudhuri, many references
- Scientific and organizing committees
- Especially Judith Rousseau
- NSF DMS-1521145, DMS-1407397