

Non-parametric Empirical Bayes and Compound Bayes Estimation of Independent Normal Means

L. Brown

Statistics Department, Wharton School
University of Pennsylvania
lbrown@wharton.upenn.edu

Joint work with E. Greenshtein

Conference:

Innovation and Inventiveness of John Hartigan

April 17, 2009

Topic Outline

1. Empirical Bayes (Concept)
2. Independent Normal Means (Setting + some theory)
3. The NP-EB Estimator (Heuristics)
4. “A Tale of Two Concepts”
 - Empirical Bayes and Compound Bayes
5. (Somewhat) Sparse Problems
6. Numerical results
7. Theorem **and Proof**
8. The heteroscedastic case – heuristics

Empirical Bayes: General Background

- n Parameters to be estimated: $\theta_1, \dots, \theta_n$. [θ_i real-valued.]
- Observe $X_i \sim f_{\theta_i}$, independent. Let $\mathbf{X} = (X_1, \dots, X_n)$.
- Estimate θ_i by $\delta_i(\mathbf{X})$.
- Component-wise Loss and Overall Risk

$$L(\theta_i, \delta_i) = (\delta_i - \theta_i)^2 \text{ and}$$

$$R(\underline{\theta}, \underline{\delta}) = \frac{1}{n} \sum E_{\theta_i} [L(\theta_i, \delta_i(\mathbf{X}))]$$

Bayes Estimation

- “Pretend” $\{\theta_i\}$ are iid, with a (prior) distribution G_n .
- Under G_n the Bayes Procedure would be

$$\Delta^{G_n} = \left(\delta_1^{G_n}, \dots, \delta_n^{G_n} \right) : \delta_i^{G_n} (X_i) = E(\theta_i | X_i).$$

[*Note:* $\delta_i^{G_n}$ depends only on X_i (and on G_n).]

- It would have Bayes risk

$$B_{[n]}(G_n) = B(G_n, \Delta^{G_n}) = E_{G_n} \left(R(\underline{\theta}, \Delta^{G_n}) \right).$$

[*Note:* Because of the scaling of sum-of-squared-error-loss by $1/n$ it is the case that $B(G_n)$ is also the coordinate-wise Bayes risk, ie,

$$B(G_n) = E_{G_n} \left[E_{\theta_i} \left(\theta_i - \delta_i^{G_n} (X_i) \right)^2 \right].]$$

Empirical Bayes

- Introduced by Robbins:
An empirical Bayes approach to statistics, 3rd Berk Symp, 1956
- “Applicable when the same decision problem presents itself repeatedly and independently with a fixed but unknown a priori distribution of the parameter.” Robbins, *Ann Math Stat*, 1964
- Thus: Fix G . Let $G_n = G$ for all n .
- Try to find a sequence of estimators, $\tilde{\Delta}_n(\mathbf{X}_n)$, that are asymptotically as good as Δ^G .
- *ie*, want

$$B_{[n]}(G, \tilde{\Delta}_n) - B_{[n]}(G) \rightarrow 0.$$

- Much of the subsequent literature emphasized the sequential nature of this problem.

The Fixed-Sample Empirical Goal

- Even earlier Robbins had taken a slightly different perspective. Asymptotically subminimax solutions of compound decision problems, *2nd Berk Symp.*, 1951. See also Zhang (2003). Robbins began,
- “When statistical problems of the same type are considered in large groups...there may exist solutions which are asymptotically ... [desirable]”
- That is, one can benefit even for fixed, but large, n (and even if G_n may change with n).
- To measure the desirability we propose

$$(1) \quad \sup_{G_n \in \mathcal{G}_n} \frac{B_{[n]}(G_n, \tilde{\Delta}_n) - B_{[n]}(G_n)}{B_{[n]}(G_n)} \rightarrow 0.$$

- Here, \mathcal{G}_n is a (very inclusive) subset of priors. [*But not all priors*].

Independent Normal Means

- Observe,

$$X_i \sim N(\theta_i, \sigma^2), i = 1, \dots, n, \text{ indep.}$$

with σ^2 known. Let φ_{σ^2} denote normal density with $\text{Var} = \sigma^2$.

- Assume $\theta_i \sim G_n$, *iid*. Write $G = G_n$, for convenience.
- Consider the i -th coordinate. Write $\theta_i = \theta$, $x_i = x$ for convenience
- The Bayes estimator (for Squared error loss) is

$$\delta^G(x) = E_G(\theta | X = x)$$

- Denote the marginal density of X as

$$g^*(x) = \int \varphi_{\sigma^2}(x - \theta) G(d\theta).$$

- As a general notation, let $\gamma^G(x) = \delta_G(x) - x$

- Note that

$$\gamma^G(x) = \delta^G(x) - x = \frac{\int (\theta - x) \varphi_{\sigma^2}(x - \theta) G(d\theta)}{\int \varphi_{\sigma^2}(x - \theta) G(d\theta)}$$

- Differentiate inside the integral (*always OK*), to write

$$(*) \quad \gamma^G(x) = \sigma^2 \frac{g^{*'}(x)}{g^*(x)}.$$

- Moral of (*): A **really good** estimate of the marginal density $g^*(x)$ should yield a good approximation to $\gamma^G(x)$.

Validity of (*) is proved in Brown (1971).

Proposed Non-Parametric Empirical-Bayes Estimator

- Let h be a bandwidth constant (to be chosen later).
- Estimate $g^*(x)$ by the kernel density estimator

$$\tilde{g}_h^*(x) = \sum \varphi_h(x - X_i) = \frac{1}{n} \sum \frac{1}{h} \varphi_1\left(\frac{x - X_i}{h}\right).$$

- *The normal kernel has some nice properties, to be explained later.*
- Define the NP EB estimator by $\tilde{\Delta} = (\tilde{\delta}_1, \dots, \tilde{\delta}_n)$ with

$$\tilde{\delta}_i(x_i) - x_i = \tilde{\gamma}_i(x_i) = \sigma^2 \frac{\tilde{g}_h^{*'}(x_i)}{\tilde{g}_h^*(x_i)}.$$

- A useful formula is

$$\tilde{g}_h^{*'}(x; \mathbf{X}) = \frac{1}{nh} \sum \frac{X_i - x}{h^2} \varphi\left(\frac{x - X_i}{h}\right).$$

A Key Lemma:

- Let $\hat{G}_n^{\mathbf{X}}$ denote the sample CDF of $\mathbf{X} = (X_1, \dots, X_n)$.
- Let $g_{G,v}^*$ denote the marginal density when $X_i \sim N(\theta_i, v)$.
- Let $\sigma^2 = 1$ and let $v = 1 + h^2$

Lemma 1: $E[\tilde{g}_h^*(x)] = g_{G,v}^*(x)$ and $E[\tilde{g}_h^{*'}(x)] = g_{G,v}^{*'}(x)$.

Proof:

- $\tilde{g}_h^* = \hat{G}_n^{\mathbf{X}} * \varphi_{h^2}$.
- $E[\hat{G}_n^{\mathbf{X}}] = \text{CDF of } X = G * \varphi_1$.
- Hence, $E[\tilde{g}_h^*(x)] = E[\hat{G}_n^{\mathbf{X}} * \varphi_{h^2}] = G * \varphi_1 * \varphi_{h^2} = G * \varphi_v$.

The proof for the derivatives is similar. ☺

Derivation of the Estimator

The expression for the estimator appears **in red** at the beginning and end of the following string of (approximate) equalities.

- $\gamma_1^G = \frac{g_{G,1}^{*'}}{g_{G,1}^*}$ by the fundamental equation (*).
- $\frac{g_{G,1}^{*'}}{g_{G,1}^*} \approx \frac{g_{G,v}^{*'}}{g_{G,v}^*}$ since $v = 1 + h^2 \approx 1$.
- $\frac{g_{G,v}^{*'}}{g_{G,v}^*} \approx \frac{\tilde{g}_h^{*'}}{\tilde{g}_h^*}$ from the Lemma

via plug-in Method-of-Moments in numerator and denominator.
See Jiang and Zhang (2007) for a different scheme based on a Fourier kernel.

A Tale of Two Formulations “Compound” and “Empirical Bayes”

“Compound” Problem:

- Let $\underline{\theta}_{(\cdot)} = \{\theta_{(1)}, \dots, \theta_{(n)}\}$ and $\underline{X}_{(\cdot)} = \{X_{(1)}, \dots, X_{(n)}\}$ denote the order statistics of $\underline{\theta}$ and \underline{X} , resp.
- Consider estimators of the form

$$\underline{\delta} = \{\delta_i\} \ni \delta_i = \delta(x_i; \underline{x}_{(\cdot)})$$

These are called *Simple-Symmetric* est's. **SS** denotes all of them.

- Given $\underline{\theta}_{(\cdot)}$ the optimal **SS** rule is denoted as $\Delta^{\underline{\theta}_{(\cdot)}}$. It satisfies

$$R(\underline{\theta}, \Delta^{\underline{\theta}_{(\cdot)}}) = \inf_{\Delta \in \text{SS}} R(\underline{\theta}, \Delta).$$

- The goal of Compound decision theory is to find rule(s) that do almost as well as $\Delta^{\underline{\theta}_{(\cdot)}}$, as judged by a criterion like (1).

The Link Between the Formulations

EB Implies CO

- Recall that $\hat{G}_n^{\theta_{(\cdot)}}$ denotes the sample CDF of $\theta_{\rightarrow(\cdot)}$.
- Then, $\Delta \in \mathbf{SS}$ implies

$$R(\theta, \Delta) = E \left\{ \frac{1}{n} \sum \left[\theta_i - \delta \left(X_i; \underline{X}_{(\cdot)} \right) \right]^2 \right\} = B \left(\hat{G}_n^{\theta_{(\cdot)}}, \Delta \right).$$

- Consequently: If $\tilde{\Delta}_n$ is EB [in the sense of (1)] then it is also Compound Optimal in the sense of: $\forall \theta_{\rightarrow} \ni \hat{G}_n^{\theta_{\rightarrow}} \in \mathcal{G}_n$.

$$(1') \quad \frac{R_{[n]}(\theta_{\rightarrow n}, \tilde{\Delta}_n) - \inf_{\Delta \in \mathbf{SS}} R_{[n]}(\theta_{\rightarrow n}, \Delta)}{\inf_{\Delta \in \mathbf{SS}} R_{[n]}(\theta_{\rightarrow n}, \Delta)} < \varepsilon_n \rightarrow 0$$

The Converse: CO \Rightarrow EB

- To MOTIVATE the converse, assume $\tilde{\Delta}_n \in \mathbf{SS}$ is CO in that

$$(1') \quad \sup_{\underline{\theta}_n \in \Theta_n} \frac{R_{[n]}(\underline{\theta}_n, \tilde{\Delta}_n) - \inf_{\Delta \in \mathbf{SS}} R_{[n]}(\underline{\theta}_n, \Delta)}{\inf_{\Delta \in \mathbf{SS}} R_{[n]}(\underline{\theta}_n, \Delta)} < \varepsilon_n \rightarrow 0.$$

- Suppose this holds when Θ_n is ALL possible vectors $\underline{\theta}_n$.
- Under a prior G_n the vector $\underline{\theta}$ has iid components, and

$$B_{[n]}(G_n, \Delta) = E \left(B_{[n]}(\hat{G}_n^{\underline{\theta}_{(\cdot)}}, \Delta) \middle| \underline{\theta}_{(\cdot)} \right).$$

- Truth of (1') for all $\underline{\theta}_n$ then implies truth of its Expectation over the distribution of $\underline{\theta}_{(\cdot)}$ under G_n . This yields (1).
- In reality (1') does not hold for all $\underline{\theta}_n$, but only for a very rich subset. Hence the proof requires extra details.

“Sparse” Problems

- An initial motivation for this research was to create a CO - EB method suitable for “Sparse” settings.
- The proto-typical sparse CO setting has

(sparse) $\underline{\theta}_{(\cdot)} = (\vartheta_0, \dots, \vartheta_0, \vartheta_1, \dots, \vartheta_1)$

with $\approx (1 - \alpha)n$ values ϑ_0 and only αn values ϑ_1 .

- Here, α is near 0, but ϑ_0, ϑ_1 may be either known or unknown.
- Situations with $\alpha = O(1/n)$ can be considered *extremely* sparse.
- Situations with, say, $\alpha \approx 1/n^{1-\varepsilon}$, $0 < \varepsilon < 1$ are *moderately* sparse.
- Sparse problems are of interest on their own merits (eg in genetics) – for example, as in Efron (2004+).
- And for their importance in building nonparametric regression estimators – see eg, Johnstone and Silverman (2004).

(Typical) Comparison of NP-EB Estimator with Best Competitor

#Signals	Est'r	Value =3	Value =4	Value =5	Value =7
5	$\tilde{\delta}_{1.15}$	53	49	42	27
5	Best other	34	32	17	7
50	$\tilde{\delta}_{1.15}$	179	136	81	40
50	Best other	201	156	95	52
500	$\tilde{\delta}_{1.15}$	484	302	158	48
500	Best other	829	730	609	505

Table: Total Expected Squared Error (via simulation; to nearest integer) Compound Bayes setup with $n=1000$; ‘most’ means =0 and others =“Value”
“Best Other” is best performing of 18 studied in *J & S* (2004).

$\tilde{\delta}_{1.15}$ is our NP – EB est'r with $v = 1.15$

Statement of EB - CO Theorem

Assumptions:

- $\exists \varepsilon' > 0 \ni \mathcal{G}_n \subseteq \left\{ G_n : B_{[n]}(G_n) > n^{\varepsilon'} \right\}$.

Hence, (only) moderately sparse settings are allowed.

- $\mathcal{G}_n \subseteq \left\{ G_n : G_n \left(\left[-C_n, C_n \right] \right) = 1 \right\} \ni C_n = O(n^\varepsilon) \forall \varepsilon > 0$.

We believe this assumption can be relaxed, but it seems that some sort of uniformly light-tail condition on \mathcal{G}_n is needed.

Theorem: Let $h_n^2 = 1/d_n$ with $d_n/\log(n) \rightarrow \infty$ & $d_n = o(n^\varepsilon) \forall \varepsilon > 0$.

Then (under above assumptions) $\tilde{\Delta}_n$ satisfies (1).

[Note: $n = 1000$ & $d_n = \log(n) \Rightarrow v = 1 + d_n^{-1} \approx 1.15$, as in Table.]

Heteroscedastic Setting

EB formulation:

- $\sigma_1^2, \dots, \sigma_n^2$ known
- Observe $X_i \sim N(\theta_i, \sigma_i^2)$, indep., $i = 1, \dots, n$.
- Assume (EB) $\theta_i \sim G_n$, indep.,
- but G_n unknown, except for $G_n \in \mathcal{G}_n$.
- Loss function, risk function, and optimality target, (1), as before.

Heuristics

- Bayes estimator on i -th coordinate has

$$\gamma_i^G(x_i) = \sigma_i^2 \left(g_{\sigma_i^2}^*{}'(x_i) / g_{\sigma_i^2}^*(x_i) \right).$$

- Previous heuristics suggest approximating $g_{\sigma_i^2}^*(x_i)$ by

$$g_{\sigma_i^2}^*(x_i) \approx g_{\sigma_i^2(1+\ell^2)}^*(x_i)$$

- And then estimating $g_{\sigma_i^2(1+\ell^2)}^*(x_i)$ as the average of

$$g_{\sigma_i^2(1+\ell^2)}^*(x_i) \approx \varphi_{h^2 = \sigma_i^2(1+\ell^2) - \sigma_k^2}(x_i - X_k), k = 1, \dots, n.$$

- To avoid impossibilities, need to use

$$h_{k,i}^2 = \left(\sigma_i^2(1 + \ell^2) - \sigma_k^2 \right)_+.$$

- Resulting estimator has

$$\gamma_i^G(x_i) = \sigma_i^2 \left(\tilde{g}^{*'}(x_i) / \tilde{g}^*(x_i) \right)$$

with $h_{k,i}^2$ as above, and

$$\tilde{g}_i^*(X) = \frac{\sum_k I_{\{k: h_{k,i}^2 > 0\}} (k) \varphi_{h_{k,i}^2}(x_i - X_k)}{\sum_k I_{\{k: h_{k,i}^2 > 0\}}}$$

- (With inessential modifications) this is the estimator used in Brown (AOAS, 2008).