

Quasi-Bayesian inference - pitfalls of incoherence

Jacek Osiewalski (Cracow University of Economics)

Bayesian analysis for a given statistical model:

- probabilistic representation of initial uncertainty about all “unknowns” – not only about observations (available, missing, future) and latent variables, but also classical parameters (unknown constants)
- Bayesian model – joint probability (density) function $p(\mathbf{y}, \omega) = p(\mathbf{y} | \omega) p(\omega)$
- $p(\mathbf{y} | \omega)$ – distribution of available observations given the remaining quantities
- $p(\omega)$ – marginal (multivariate) distribution of all quantities that remain unknown after seeing the data (i.e., after seeing the realization of the vector \mathbf{y} of available observations)
- Bayesian inference is based on simple, general rules of probability calculus

1° conditioning – Bayes formula: $p(\omega | \mathbf{y}) = \frac{p(\mathbf{y} | \omega) p(\omega)}{p(\mathbf{y})} = \frac{p(\mathbf{y} | \omega) p(\omega)}{\int_{\Omega} p(\mathbf{y} | \omega) p(\omega)} \propto p(\mathbf{y} | \omega) p(\omega)$,

2° marginalization – deriving univariate distributions from $p(\omega | \mathbf{y})$

“Coherent inference” – the one following strict rules of probability calculus

Quasi-Bayesian inference:

- **Bayes formula used mechanically, outside the full probabilistic context – incoherence !**
- **$p(y | \omega) = g(y; \omega)$ corresponds to some traditional statistical model**
- **$p(\omega) = f(\omega; y)$ is specified using given y , so it cannot be the marginal distribution !!!**
- **thus $p(\omega | y) \propto g(y; \omega) f(\omega; y)$ IS NOT the posterior in a Bayesian model with initially assumed $p(y | \omega)$, but it can be the posterior in a completely different Bayesian model**
- **question: what are the *true* building blocks (statistical model and prior) corresponding to such $p(\omega | y)$? it would be useful to know *true* assumptions, not only the *declared* ones**
- **fundamental pitfall of incoherence – $p(\omega | y)$ corresponds to some statistical model and prior assumptions to be discovered !**

So-called “Empirical Bayes” (EB) is the most popular quasi-Bayesian approach, advocated in non-Bayesian, sampling-theory texts on inference in hierarchical multi-level statistical models

→ Here we show hidden assumptions behind the EB inference in hierarchical models

SOME SIMPLE EXAMPLE FIRST (Example 1)

Bayesian model: $p(\mathbf{y}, \boldsymbol{\mu}) = p(\mathbf{y} | \boldsymbol{\mu}) p(\boldsymbol{\mu}) = f_N^n(\mathbf{y} | \boldsymbol{\mu} \mathbf{e}_n, c\mathbf{I}_n) f_N^1(\boldsymbol{\mu} | \mathbf{a}, v)$

Decomposition: $p(\mathbf{y}, \boldsymbol{\mu}) = p(\mathbf{y}) p(\boldsymbol{\mu} | \mathbf{y}) = f_N^n(\mathbf{y} | \mathbf{a} \mathbf{e}_n, c\mathbf{I}_n + v \mathbf{e}_n \mathbf{e}_n') f_N^1(\boldsymbol{\mu} | \mathbf{a}_y, v_y)$

where $v_y = \left(\frac{n}{c} + \frac{1}{v}\right)^{-1}$, $\mathbf{a}_y = \left(\frac{n}{c} + \frac{1}{v}\right)^{-1} \left(\frac{n}{c} \bar{\mathbf{y}} + \frac{1}{v} \mathbf{a}\right)$, $\bar{\mathbf{y}} = \frac{1}{n} \mathbf{e}_n' \mathbf{y}$, $\mathbf{e}_n = (\mathbf{1} \ \mathbf{1} \ \dots \ \mathbf{1})'$

Quasi-Bayesian inference: imagine a non-Bayesian statistician who agrees to use Bayes formula

$$p(\boldsymbol{\mu} | \mathbf{y}) \propto p(\mathbf{y} | \boldsymbol{\mu}) p(\boldsymbol{\mu})$$

but disagrees to subjectively specify \mathbf{a} (prior mean); instead he/she puts $\bar{\mathbf{y}}$ (sample average)

and (informally) uses $p^*(\boldsymbol{\mu}) = f_N^1(\boldsymbol{\mu} | \bar{\mathbf{y}}, v)$ and $p^*(\boldsymbol{\mu} | \mathbf{y}) = f_N^1(\boldsymbol{\mu} | \bar{\mathbf{y}}, \left(\frac{n}{c} + \frac{1}{v}\right)^{-1})$

Is there any hidden Bayesian model (sampling + prior) formally justifying such “posterior”?

Consider $\tilde{p}(\mathbf{y}, \boldsymbol{\mu}) = p(\mathbf{y} | \boldsymbol{\mu}) p^*(\boldsymbol{\mu}) = f_N^n(\mathbf{y} - \boldsymbol{\mu} \mathbf{e}_n | \mathbf{0}, c\mathbf{I}_n) f_N^1(\boldsymbol{\mu} - \bar{\mathbf{y}} | \mathbf{0}, v)$

it decomposes into $\tilde{p}(\boldsymbol{\mu} | \mathbf{y}) = p^*(\boldsymbol{\mu} | \mathbf{y})$ and $\tilde{p}(\mathbf{y}) \propto \exp\left(-\frac{1}{2c} \mathbf{y}' M \mathbf{y}\right)$, $M = \mathbf{I}_n - \frac{1}{n} \mathbf{e}_n \mathbf{e}_n'$

or $\tilde{p}(\mathbf{y} | \boldsymbol{\mu}) = f_N^n\left(\mathbf{y} \mid \boldsymbol{\mu} \mathbf{e}_n, c\left(\mathbf{I}_n - \frac{c}{n(c+nv)} \mathbf{e}_n \mathbf{e}_n'\right)\right)$ and $\tilde{p}(\boldsymbol{\mu})$ constant (!!!)

true sampling model assumes dependence (equi-correlation); true prior is flat, improper

MAIN PART: Statistical models with hierarchical structure

conditional distribution of observations: $p(\mathbf{y}|\boldsymbol{\theta}) = g(\mathbf{y}; \boldsymbol{\theta}), \mathbf{y} \in Y, \boldsymbol{\theta} \in \Theta;$

distribution of random parameters (latent variables): $f_0(\boldsymbol{\theta}; \boldsymbol{\alpha}), \boldsymbol{\alpha} \in A \subseteq \mathbb{R}^s;$

joint distribution ($\boldsymbol{\alpha}$ fixed):

$$p(\mathbf{y}|\boldsymbol{\theta}) f_0(\boldsymbol{\theta}; \boldsymbol{\alpha}) = g(\mathbf{y}; \boldsymbol{\theta}) f_0(\boldsymbol{\theta}; \boldsymbol{\alpha}) = f_1(\boldsymbol{\theta}|\mathbf{y}; \boldsymbol{\alpha}) h(\mathbf{y}; \boldsymbol{\alpha}) \quad \leftarrow \text{decomposition}$$

$h(\mathbf{y}; \boldsymbol{\alpha}) \quad \leftarrow$ marginal distribution of \mathbf{y}

$$f_1(\boldsymbol{\theta}|\mathbf{y}; \boldsymbol{\alpha}) = \frac{g(\mathbf{y}; \boldsymbol{\theta}) f_0(\boldsymbol{\theta}; \boldsymbol{\alpha})}{h(\mathbf{y}; \boldsymbol{\alpha})} \propto g(\mathbf{y}; \boldsymbol{\theta}) f_0(\boldsymbol{\theta}; \boldsymbol{\alpha}) \quad \leftarrow \text{conditional distribution of } \boldsymbol{\theta} \text{ (Bayes formula)}$$

SIMPLE EXAMPLE OF A HIERARCHICAL MODEL (Example 2)

θ_i – unobservable characteristic, randomly distributed over n observed units ($i = 1, \dots, n$),

$\theta = (\theta_1 \dots \theta_n)'$, $\theta_i \sim iN(\alpha, d)$, $d > 0$ known;

$x_i = (x_{i1} \dots x_{im})'$, $x_{ij} \sim iN(\theta_i, c_0)$ ($j = 1, \dots, m$) – independent measurements of θ_i (c_0 known)

$y_i = \frac{1}{m} e_m' x_i = \bar{x}_i$ – sufficient statistic (for fixed θ_i); $y_i \sim iN(\theta_i, c)$, $c = \frac{c_0}{m}$, $y = (y_1 \dots y_n)'$

$$p(y|\theta) = f_N^n(y|\theta, cI_n), \quad f_0(\theta; \alpha) = f_N^n(\theta|\alpha e_n, dI_n)$$

Decomposition of the product $p(y|\theta) f_0(\theta; \alpha)$ into $f_1(\theta|y; \alpha) h(y; \alpha)$, where

$$h(y; \alpha) = \int_{\mathbb{R}^n} p(y|\theta) f_0(\theta; \alpha) d\theta = f_N^n(y|\alpha e_n, (c + d)I_n),$$

$$f_1(\theta|y; \alpha) = f_N^n\left(\theta \left| \frac{d^{-1}}{c^{-1}+d^{-1}} \alpha e_n + \frac{c^{-1}}{c^{-1}+d^{-1}} y, \frac{1}{c^{-1}+d^{-1}} I_n\right.\right) \quad (\text{final precision} = \text{sample} + \text{prior})$$

$$E(\theta|y; \alpha) = w \cdot \alpha e_n + (1 - w) \cdot y, \quad w = \frac{d^{-1}}{c^{-1}+d^{-1}} \in (0, 1) \quad (w = \text{prior precision} / \text{final precision})$$

$E(\theta|y; \alpha)$ – point in $\Theta = \mathbb{R}^n$ lying on the line segment between $(\alpha \alpha \dots \alpha)'$ and $(y_1 y_2 \dots y_n)'$

$f_1(\theta|y; \alpha)$ follows Bayes Theorem for any fixed α , so then we have coherence; but how to get α ?

Empirical Bayes (EB)

inference on θ based on the conditional distribution $f_1(\theta|y; \alpha)$ obtained using Bayes Theorem, BUT for some point estimate of unknown $\alpha \in A$, e.g., using so-called type II maximum likelihood:

$$\hat{\alpha} = \hat{\alpha}_{ML} = \arg \max L(\alpha; y) = \arg \max h(y; \alpha), \quad \alpha \in A$$

So EB uses $\hat{p}(\theta|y) = f_1(\theta|y, \hat{\alpha}) \propto p(y|\theta)f_0(\theta; \hat{\alpha})$,

i.e. the “posterior” corresponding to the “prior” with hyper-parameter based on y !!!

EXAMPLE 2 (continued)

$$L(\alpha; y) = h(y; \alpha) = f_N^n(y|\alpha e_n, (c+d)I_n) = (2\pi \cdot \frac{c+d}{n})^{\frac{1}{2}} f_N^1\left(\alpha \middle| \bar{y}, \frac{c+d}{n}\right) f_N^n(My|0, (c+d)I_n),$$

$$\hat{\alpha} = \hat{\alpha}_{ML} = \bar{y} = \frac{1}{n} e_n' y, \quad M = I_n - \frac{1}{n} e_n e_n',$$

$$\hat{p}(\theta|y) = f_1(\theta|y, \hat{\alpha}) = f_N^n(\theta|\hat{\theta}_{EB}, \frac{1}{c^{-1}+d^{-1}}I_n), \quad \hat{\theta}_{EB} = \frac{d^{-1}}{c^{-1}+d^{-1}} \bar{y} e_n + \frac{c^{-1}}{c^{-1}+d^{-1}} y$$

- uncertainty about α not taken into account
- obvious incoherence of inferences on θ

Bayesian hierarchical model (BHM)

$p(\mathbf{y}, \omega) = p(\mathbf{y}, \theta, \alpha) = p(\mathbf{y}|\theta) p(\theta|\alpha) p(\alpha)$, $p(\alpha)$ – the prior for $\alpha \in A$

$\omega = (\theta, \alpha)$, conditional independence: $\mathbf{y} \perp \alpha \mid \theta$ – leads to $p(\mathbf{y}|\omega) = p(\mathbf{y}|\theta)$

$p(\mathbf{y}|\theta) = g(\mathbf{y}; \theta)$, $p(\theta|\alpha) = f_0(\theta; \alpha)$ – the same as in EB

final decomposition of Bayesian model: $p(\mathbf{y}, \theta, \alpha) = p(\mathbf{y}) p(\theta, \alpha|\mathbf{y}) = p(\mathbf{y}) p(\alpha|\mathbf{y}) p(\theta|\mathbf{y}, \alpha)$

$$p(\theta|\mathbf{y}, \alpha) = \frac{p(\mathbf{y}|\theta) p(\theta|\alpha)}{p(\mathbf{y}|\alpha)} = \frac{g(\mathbf{y}; \theta) f_0(\theta; \alpha)}{h(\mathbf{y}; \alpha)} = f_1(\theta|\mathbf{y}; \alpha)$$

$$p(\alpha|\mathbf{y}) = \frac{p(\mathbf{y}|\alpha) p(\alpha)}{p(\mathbf{y})} = \frac{h(\mathbf{y}; \alpha) p(\alpha)}{p(\mathbf{y})}$$

$$p(\mathbf{y}) = \int_A p(\mathbf{y}|\alpha) p(\alpha) d\alpha$$

Remarks:

- $p(\theta|\mathbf{y}) = \int_A f_1(\theta|\mathbf{y}; \alpha) p(\alpha|\mathbf{y}) d\alpha$ – uncertainty about α is formally taken into account
- Bayes Theorem is used twice: for latent variables (given parameters) and for parameters

EXAMPLE 2 (continued) – Bayesian hierarchical model with:

$$p(\mathbf{y}|\boldsymbol{\theta}) = f_N^n(\mathbf{y}|\boldsymbol{\theta}, c\mathbf{I}_n), \quad p(\boldsymbol{\theta}|\alpha) = f_N^n(\boldsymbol{\theta}|\alpha\mathbf{e}_n, d\mathbf{I}_n), \quad p(\alpha) = f_N^1(\alpha|a, v)$$

$$* \quad p(\boldsymbol{\theta}) = \int_{-\infty}^{+\infty} p(\boldsymbol{\theta}|\alpha) p(\alpha) d\alpha = f_N^n(\boldsymbol{\theta}|\mathbf{a}\mathbf{e}_n, d\mathbf{I}_n + v\mathbf{e}_n\mathbf{e}_n')$$

$$* \quad p(\boldsymbol{\theta}|\mathbf{y}) \propto p(\mathbf{y}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) = f_N^n(\mathbf{y}|\boldsymbol{\theta}, c\mathbf{I}_n) p(\boldsymbol{\theta})$$

or, equivalently, $p(\boldsymbol{\theta}|\mathbf{y}) = \int_{-\infty}^{+\infty} p(\boldsymbol{\theta}|\mathbf{y}, \alpha) p(\alpha|\mathbf{y}) d\alpha = \int_{-\infty}^{+\infty} f_1(\boldsymbol{\theta}|\mathbf{y}; \alpha) p(\alpha|\mathbf{y}) d\alpha$

where $p(\alpha|\mathbf{y}) = f_N^1(\alpha | \left(\frac{n}{c+d} + \frac{1}{v}\right)^{-1} \left(\frac{n}{c+d}\bar{\mathbf{y}} + \frac{a}{v}\right), \left(\frac{n}{c+d} + \frac{1}{v}\right)^{-1})$

Finally:

$$p(\boldsymbol{\theta}|\mathbf{y}) = f_N^n(\boldsymbol{\theta} | \frac{c^{-1}}{c^{-1}+d^{-1}}\mathbf{y} + \frac{d^{-1}}{c^{-1}+d^{-1}}\left(\frac{n}{c+d} + \frac{1}{v}\right)^{-1} \left(\frac{n}{c+d}\bar{\mathbf{y}} + \frac{a}{v}\right) \cdot \mathbf{e}_n, \\ \frac{1}{c^{-1}+d^{-1}}\mathbf{I}_n + \left(\frac{n}{c+d} + \frac{1}{v}\right)^{-1} \left(\frac{d^{-1}}{c^{-1}+d^{-1}}\right)^2 \mathbf{e}_n\mathbf{e}_n').$$

If $v^{-1} \approx 0$, then $p(\alpha) \approx \text{const}$, $p(\boldsymbol{\theta}) \propto \exp\left(-\frac{1}{2d} \boldsymbol{\theta}'\mathbf{M}\boldsymbol{\theta}\right)$, $\mathbf{M} = \mathbf{I}_n - \frac{1}{n}\mathbf{e}_n\mathbf{e}_n'$, and

$$p(\boldsymbol{\theta}|\mathbf{y}) \approx f_N^n(\boldsymbol{\theta} | \hat{\boldsymbol{\theta}}_{EB}, \frac{1}{c^{-1}+d^{-1}}\mathbf{I}_n + \frac{c^2}{n(c+d)}\mathbf{e}_n\mathbf{e}_n'); \quad \frac{c^2}{n(c+d)}\mathbf{e}_n\mathbf{e}_n' - \text{reflects uncertainty about } \alpha!$$

If also n is large enough, then $p(\boldsymbol{\theta}|\mathbf{y}) \approx \hat{p}(\boldsymbol{\theta}|\mathbf{y})$; asymptotically, incoherence does not matter

Small-sample interpretation of *Empirical Bayes*

For a given EB form of $\hat{p}(\theta|y)$, we seek for $\tilde{p}(y|\theta)$ and $\tilde{p}(\theta)$ that lead to the Bayesian model $\tilde{p}(y, \theta) = \tilde{p}(y|\theta) \tilde{p}(\theta)$ of the form

$$\tilde{p}(y, \theta) = k(y) p(y|\theta) p(\theta|\alpha = \hat{\alpha}) = k(y) g(y; \theta) f_0(\theta; \hat{\alpha}),$$

resulting in $\hat{p}(\theta|y)$ as the true posterior, i.e.

$$\tilde{p}(\theta|y) = \hat{p}(\theta|y) \propto g(y; \theta) f_0(\theta; \hat{\alpha}).$$

From the form of $\tilde{p}(y, \theta)$ we obtain the (implicit) prior $\tilde{p}(\theta) = \int_Y k(y) g(y; \theta) f_0(\theta; \hat{\alpha}) dy$ and then the (implicit) conditional distribution of observations

$$\tilde{p}(y|\theta) = \frac{\tilde{p}(y, \theta)}{\tilde{p}(\theta)} = k(y) \frac{g(y; \theta) f_0(\theta; \hat{\alpha})}{\tilde{p}(\theta)}.$$

If both f_0 and k are not constant in y , then $\tilde{p}(y|\theta) \neq p(y|\theta) = g(y; \theta)$ and the *true* conditional distribution of observations is different from the *initially assumed* (declared) one.

EXAMPLE 2 (continued)

$$\begin{aligned}\tilde{p}(\mathbf{y}, \boldsymbol{\theta}) &= k(\mathbf{y})g(\mathbf{y}; \boldsymbol{\theta}) f_0(\boldsymbol{\theta}; \hat{\boldsymbol{\alpha}}) = k f_N^n(\mathbf{y}|\boldsymbol{\theta}, c\mathbf{I}_n) f_N^n(\boldsymbol{\theta}|\bar{\mathbf{y}}\mathbf{e}_n, d\mathbf{I}_n) \\ &= f_N^n\left(\mathbf{y}\middle|\boldsymbol{\theta}, \left(\frac{1}{c}\mathbf{I}_n + \frac{1}{dn}\mathbf{e}_n\mathbf{e}_n'\right)^{-1}\right) k(2\pi)^{-\frac{n}{2}} \exp\left(-\frac{1}{2d}\boldsymbol{\theta}'\mathbf{M}\boldsymbol{\theta}\right)\end{aligned}$$

From $\tilde{p}(\mathbf{y}, \boldsymbol{\theta})$ we easily derive:

$\tilde{p}(\boldsymbol{\theta}) = k(2\pi)^{-\frac{n}{2}} \exp\left(-\frac{1}{2d}\boldsymbol{\theta}'\mathbf{M}\boldsymbol{\theta}\right)$ – improper (only σ -finite), but informative (favors equal θ_i)
(for $v^{-1} \approx \mathbf{0}$ we get $p(\boldsymbol{\theta}) \approx \tilde{p}(\boldsymbol{\theta})$, so the *declared* prior coincides with the *true* one)

$$\tilde{p}(\mathbf{y}|\boldsymbol{\theta}) = f_N^n\left(\mathbf{y}\middle|\boldsymbol{\theta}, \left(\frac{1}{c}\mathbf{I}_n + \frac{1}{dn}\mathbf{e}_n\mathbf{e}_n'\right)^{-1}\right)$$

– conditional distribution with equally correlated observations (instead of independent ones!!!)

$$\tilde{\mathbf{V}}(\mathbf{y}|\boldsymbol{\theta}) = c\left(\mathbf{I}_n - \frac{c}{n(c+d)}\mathbf{e}_n\mathbf{e}_n'\right) \Rightarrow \widetilde{\text{Corr}}(\mathbf{y}_i, \mathbf{y}_j|\boldsymbol{\theta}) = \frac{c}{(n-1)c+nd} \quad (i \neq j),$$

true $\tilde{p}(\mathbf{y}|\boldsymbol{\theta})$ is qualitatively different from *declared* $p(\mathbf{y}|\boldsymbol{\theta})$; problem disappears when $n \rightarrow \infty$

Concluding remarks

- From the purely Bayesian perspective, using Bayes formula with “prior” dependent on actual data is completely incoherent.
- Is this, however, of any interest to a non-Bayesian statistician? Perhaps such incoherent quasi-Bayesian approach generates inference tools that are better in terms of sampling-theory properties...
- Remind that, under certain regularity conditions, Bayesian decision functions (estimators) are admissible (cannot be improved – in terms of risk – uniformly in the parameter space) and form complete classes of such decision functions.
- Here it has been shown that incoherent, quasi-Bayesian approaches can be interpreted as Bayesian for other sampling models, not for the *declared* (assumed) ones.
- When a quasi-Bayesian procedure is not Bayesian for the declared sampling model, it may produce inadmissible decision functions (within this sampling model).
- Being coherent (i.e., being Bayesian and obeying rules of probability) prevents from such risks – and it does so for (almost) free...

THANK YOU FOR YOUR ATTENTION!