

An Objective Prior for Hyperparameters in Normal Hierarchical Models

Dongchu Sun, University of Missouri

Jointly with

James Berger, Duke University;

Chengyuan Song, East China Normal University

October 14, 2016

Motivation

- Hierarchical models are the workhorse of much of Bayesian analysis, yet there is uncertainty as to which objective priors to use for hyperparameters (parameters at higher levels of the hierarchical model).
- Formal approaches to objective Bayesian analysis, such as the Jeffreys-rule or reference prior approach, are only implementable in simple hierarchical settings (one-way model).
- It is common to use less formal approaches, i.e., utilize formal priors from non-hierarchical models in hierarchical settings.
- This can be fraught with danger, however. For instance, non-hierarchical Jeffreys-rule priors for variances or covariance matrices result in improper posterior if they are used at higher levels of a hierarchical model.
- Such less formal approaches must be carefully evaluated, and not just from the perspective of posterior propriety.

Normal Hierarchical Model

Consider the hierarchical model,

$$\mathbf{y}_i \mid \boldsymbol{\theta}_i \sim N_k(\boldsymbol{\theta}_i, \boldsymbol{\Sigma}_i); \quad (1)$$

$$\boldsymbol{\theta}_i \mid \boldsymbol{\beta}, \mathbf{V} \sim N_k(\boldsymbol{\beta}, \mathbf{V}), \quad i = 1, \dots, m, \quad (2)$$

where

- \mathbf{y}_i are $k \times 1$ observation vectors,
- $\boldsymbol{\Sigma}_i$ are $k \times k$ known covariance matrices,
- $\boldsymbol{\theta}_i$ are $k \times 1$ unknown mean vectors,
- $\boldsymbol{\beta}$ is a $k \times 1$ unknown "hyper-mean" vector and
- \mathbf{V} is an unknown $k \times k$ "hyper-covariance matrix".
- Here m is the block number and k is the dimension of \mathbf{y}_i .

A Note

- Often the data arise from linear models $\mathbf{y}_i = \mathbf{X}_i \boldsymbol{\theta}_i + \mathbf{e}_i$, where the \mathbf{X}_i are known design matrices of covariates.
- In this situation, one can simply transform to $\mathbf{y}_i^* = (\mathbf{X}_i^t \boldsymbol{\Sigma}_i^{-1} \mathbf{X}_i)^{-1} \mathbf{X}_i^t \boldsymbol{\Sigma}_i^{-1} \mathbf{y}_i$, which will be distributed as in (1), with transformed covariance matrix $\boldsymbol{\Sigma}_i^* = (\mathbf{X}_i^t \boldsymbol{\Sigma}_i^{-1} \mathbf{X}_i)^{-1}$.

Hospital Example, I

- At hospital i , the observations $\mathbf{y}_i = (y_{i1}, \dots, y_{ik})'$ are the sample averages of the costs of k different medical treatments;
- $\boldsymbol{\theta}_i$ is the corresponding unknown vector of true mean costs of the treatments at the hospital;
- $\boldsymbol{\Sigma}_i$ is the covariance matrix associated with these sample averages. If $\boldsymbol{\Sigma}_i$ is not known, it can typically be estimated from the observed data.

Hospital Example, II

- The means θ_i are then modeled, conditional on a $l \times 1$ unknown 'hyper-mean' vector β and $k \times k$ unknown 'hyper-covariance matrix' V , as

$$\theta_i = z_i \beta + e_i^*, \quad e_i^* \sim N_k(\mathbf{0}, V), \quad (3)$$

- z_i are specified $k \times l$ covariate matrices and $l \geq 1$.

Hospital Example, III

- Consider the j^{th} coordinate, θ_{ij} of each θ_i ; it refers to the cost of a certain medical treatment at i^{th} hospital.
- It is natural to model this as depending on p hospital characteristics, such as the number of patients receiving each treatment, the average severity of the condition of the patients for each of the treatments, the average income of the patients, etc.
- Denoting these characteristics for the j^{th} treatment at the i^{th} hospital as $\mathbf{v}_{ij} = (v_{ij1}, \dots, v_{ijp})$;

Hospital Example, IV

- A reasonable model would be a regression model

$$\theta_{ij} = \mathbf{v}_{ij}\boldsymbol{\alpha}_j + e_{ij}, \quad \text{for } i = 1, \dots, m,$$

- where $\boldsymbol{\alpha}_j$ is a column vector of weights determining the effect of hospital characteristics on the cost of treatment j , and e_{ij} is normal error. There would typically be a separate regression of this form for each treatment j , and stacking these regressions vertically leads to equation (3), where $l = kp$,

$$\mathbf{z}_i = \begin{pmatrix} \mathbf{v}_{i1} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{v}_{i2} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{v}_{ik} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \boldsymbol{\alpha}_1 \\ \boldsymbol{\alpha}_2 \\ \vdots \\ \boldsymbol{\alpha}_k \end{pmatrix}. \quad (4)$$

Hospital Example, V

- If it were thought that each of the treatment cost means were independent, V could be taken as diagonal,
- but it is far more likely that there is considerable dependence, so that completely unknown V is reasonable.

Hospital Example, VI

- The above example suggests another possible level to the hierarchical modeling.
- In (3) we stop with the unknown hyperprior mean β , but (4) suggests that it might be sensible to further model this mean, with a third level hierarchy, as

$$\alpha_i \sim N_p(\cdot \mid \xi, \Omega). \quad (5)$$

- Whether or not this is appropriate depends on precise context, but we will see that this analysis can be done with the methodology herein.

The hyperprior distributions: $\pi(\boldsymbol{\beta}, \mathbf{V}) = \pi(\boldsymbol{\beta})\pi(\mathbf{V})$

For \mathbf{V} , satisfies the following condition :

Condition 1. for $0 \leq c \leq 1$,

$$\frac{C_1}{|\mathbf{I} + \mathbf{V}|^{a_2 - a_1} |\mathbf{V}|^{a_1} [\prod_{i < j} (d_i - d_j)]^{1-c}} \leq \pi(\mathbf{V}) \leq \frac{C_2}{|\mathbf{I} + \mathbf{V}|^{a_2 - a_1} |\mathbf{V}|^{a_1} [\prod_{i < j} (d_i - d_j)]^{1-c}}$$

Where C_1 and C_2 are positive constants, $d_1 > d_2 > \dots > d_k > 0$ are the ordered eigenvalues of \mathbf{V} .

Hyperprior of \mathbf{V} : Common Cases

- **Constant prior** ($a_1 = a_2 = 0, c = 1$): $\pi(\mathbf{V}) = 1$;
- **Nonhierarchical indep. Jeffreys priors**
 $(a_1 = a_2 = (k + 1)/2, c = 1)$: $\pi(\mathbf{V}) = |\mathbf{V}|^{-(k+1)/2}$.
- **Hierarchical indep. Jeffreys priors** ($a_1 = 0, a_2 = (k + 1)/2, c = 1$): $\pi(\mathbf{V}) = |\mathbf{I} + \mathbf{V}|^{-(k+1)/2}$.
- **Nonhierarchical indep. reference priors** ($a_1 = a_2 = 1, c = 0$):
 $\pi(\mathbf{V}) = [|\mathbf{V}| \prod_{i < j} (d_i - d_j)]^{-1}$; Yang & Berger (1994).
- **Hierarchical indep. reference priors**
 - 1 $a_1 = 0, a_2 = 1, c = 0$: $\pi(\mathbf{V}) = [|\mathbf{I} + \mathbf{V}| \prod_{i < j} (d_i - d_j)]^{-1}$;
 - 2 $a_1 = a_2 = 1 - \frac{1}{2k}, c = 0$.
 $\pi(\mathbf{V}) = [|\mathbf{V}|^{-(1-1/2k)} \prod_{i < j} (d_i - d_j)]^{-1}$; (**recommended**)

Hyperprior of β

Three commonly considered priors for hyperparameter β are:

- Case 1 **Constant prior**: $\pi(\beta) = 1$.
- Case 2 **Conjugate prior**: $\pi(\beta)$ is $N_k(\beta^0, \mathbf{A})$, where β^0 and \mathbf{A} are subjectively specified.
- Case 3 **Hierarchical prior**:

$$\beta \mid \lambda \sim N_k(\beta^0, \lambda \mathbf{I}), \quad \lambda \sim \pi(\lambda), \lambda > 0,$$

where β^0 is specified, $\pi(\lambda) \propto \lambda^{-1/2} e^{-1/2\lambda}$,

recommended default prior:

$$\pi(\beta) \sim \frac{1}{[1 + \|\beta\|^2]^{(k-1)/2}}.$$

Frequentist Justification in estimating θ

- Frequentist coverage probability of Bayesian credible intervals of a function of θ .
What functions of β ? Each component of θ ?
- Admissibility of Estimation of θ under a loss function.

The definition of admissibility and inadmissibility.

The talk gives conditions under which the hierarchical Bayes estimate $\delta^\pi(\mathbf{y})$ (posterior mean) of $\boldsymbol{\theta}$ is admissible and inadmissible for quadratic loss

$$L(\boldsymbol{\theta}, \boldsymbol{\delta}^\pi) = (\boldsymbol{\theta} - \boldsymbol{\delta}^\pi(\mathbf{y}))^t \mathbf{Q} (\boldsymbol{\theta} - \boldsymbol{\delta}^\pi(\mathbf{y})).$$

where \mathbf{Q} is a known positive-definite matrix.

The performance of an estimator will be evaluated by the usual frequentist risk function

$$R(\boldsymbol{\theta}, \boldsymbol{\delta}) = E_{\boldsymbol{\theta}}^{\mathbf{x}}[L(\boldsymbol{\theta}, \boldsymbol{\delta}^\pi)]$$

The estimator $\boldsymbol{\delta}$ is **inadmissible** if there exists another estimator with risk function nowhere bigger and somewhere smaller. If no such better estimator exists, $\boldsymbol{\delta}$ is **admissible**.

Admissibility

- Berger, Strawderman and Tang (2005) approached the question of choice of hyperpriors in normal hierarchical models by looking at the frequentist notion of admissibility of resulting estimators.
- The motivation was that hyperpriors that are too diffuse result in inadmissible estimators, while hyperpriors that are concentrated enough result in admissible estimators.
- Hyperpriors that are 'on the boundary of admissibility' are sensible choices for objective priors, being as diffuse as possible without resulting in inadmissible procedures.
- The admissibility (and propriety) properties of a number of priors were considered in the paper, but no overall conclusion was reached as to a specific prior to recommend, in part because they were not able to prove admissibility for the leading candidate prior.

Two powerful results from Brown (1971)

Define

$$\overline{m}(r) = \int m(x) d\phi(x), \quad \underline{m}(r) = \int \frac{1}{m(x)} d\phi(x)$$

where $\phi(\cdot)$ is the uniform probability measure on the surface of the sphere of radius $r = \|x\|$.

Result 1. If $\delta^\pi(x) - x$ is **uniformly bounded** and

$$\int_c^\infty [r^{mk-1} \overline{m}(r)]^{-1} dr = \infty \quad (6)$$

for some $c > 0$, then $\delta^\pi(x)$ is **admissible**.

Result 2. If

$$\int_c^\infty r^{1-mk} \underline{m}(r) dr < \infty \quad (7)$$

for some $c > 0$, then $\delta^\pi(x)$ is **inadmissible**.

A Theorem on Admissibility

Theorem 0.1

Consider a class of independent objective priors,

$$\pi_b(\boldsymbol{\beta}) \propto \frac{1}{[1 + \|\boldsymbol{\beta}\|^2]^{(k+2b-2)/2}}, \quad \pi_a(\mathbf{V}) = \frac{1}{|\mathbf{V}|^a \prod_{i < j} (d_i - d_j)}. \quad (8)$$

If $m \geq 2$, $0 \leq b < 1$ and $1 - b/k \leq a < 1$, the posterior mean of $\boldsymbol{\theta}$ is admissible.

- The theorem extends the corresponding result in Berger, Strawderman and Tang (2005).
- It proves admissibility for the boundary $a = 1 - b/k$, that is crucial as it is priors on the boundary of admissibility which are sought.

Computation: Gibbs sampling for β

To compute with the recommended prior, use the representation

$$\beta \mid \lambda \sim N_k(\cdot \mid \mathbf{0}, \lambda \mathbf{I}), \quad \pi(\lambda) \propto \lambda^{-1/2} e^{-1/2\lambda}.$$

- Sample λ from its full conditional, the Inverse Gamma($(k-1)/2, 2/[1 + \|\beta\|^2]$) density; if $k = 1$, this step is not needed as the hyperprior is constant.
- Given λ (and \mathbf{V} and the θ_i), Gibbs sampling of β can be done from its full conditional, (when $k = 1$, set $\lambda = \infty$)

$$N_k\left(\left(\frac{1}{\lambda} \mathbf{I} + \sum_{i=1}^m \mathbf{z}'_i \mathbf{V}^{-1} \mathbf{z}_i\right)^{-1} \sum_{i=1}^m \mathbf{z}'_i \mathbf{V}^{-1} \theta_i, \left(\frac{1}{\lambda} \mathbf{I} + \sum_{i=1}^m \mathbf{z}'_i \mathbf{V}^{-1} \mathbf{z}_i\right)^{-1}\right).$$

Note that the only cost in using the shrinkage prior as opposed to the constant prior is the need to sample λ , but this is an insignificant cost.

Previously suggested sampling methods for \mathbf{V}

The full conditional for \mathbf{V} can be written

$$\pi(\mathbf{V} \mid \boldsymbol{\theta}, \boldsymbol{\beta}) \propto \frac{1}{|\mathbf{V}|^{(\frac{m}{2}+1-\frac{1}{k})} \prod_{i < j} (d_i - d_j)} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{V}^{-1} \mathbf{W}(\boldsymbol{\theta}, \boldsymbol{\beta}))\right),$$

where $\mathbf{W}(\boldsymbol{\theta}, \boldsymbol{\beta}) = \sum_{i=1}^m (\boldsymbol{\theta}_i - \mathbf{z}'_i \boldsymbol{\beta})(\boldsymbol{\theta}_i - \mathbf{z}'_i \boldsymbol{\beta})^t$.

Method 1. The method in Yang and Berger using the log transformation of a covariance matrix. This relies on sampling all of the covariance matrix at one time and hence can only work for modest dimensions (≤ 5).

Previously suggested Metropolis sampling methods for \mathbf{V} .

- *Step 0.* Start with $\mathbf{V}^0 = \mathbf{I}$ or the marginal MLE.
- *Step 1.* At iteration r , generate $\mathbf{V}^* \sim \text{Inv-Wishart}(\mathbf{W}(\boldsymbol{\theta}, \boldsymbol{\beta}), m)$.
- *Step 2.* Set $\mathbf{V}^{r+1} = \begin{cases} \mathbf{V}^* & \text{with probability } \alpha, \\ \mathbf{V}^r & \text{otherwise,} \end{cases}$, where

$$\alpha = \min \left\{ 1, \frac{\prod_{i < j} (d_i^* - d_j^*)}{\prod_{i < j} (d_i^r - d_j^r)} \cdot \frac{|\mathbf{V}^r|^{(k-1+k^{-1})/2}}{|\mathbf{V}^*|^{(k-1+k^{-1})/2}} \right\},$$

the d_i^* and d_i^r being the eigenvalues of \mathbf{V}^* and \mathbf{V}^r , respectively.

- *Step 3.* After iterating Steps 1 and 2 N times, throw away the first M iterations keeping $(\mathbf{V}^{M+1}, \dots, \mathbf{V}^N)$ as the samples from the posterior distribution.

This also samples the entire covariance matrix and hence will only work for modest dimensions.

A New Method

This is a potentially powerful new method that might work for higher dimensions (which the others do not).

Defining $r = (\frac{m}{2} + 1 - \frac{1}{k})$, the full conditional for \mathbf{V} can be written

$$\pi(\mathbf{V} \mid \boldsymbol{\theta}, \boldsymbol{\beta}) \propto \frac{1}{|\mathbf{V}|^r \prod_{i < j} (d_i - d_j)} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{V}^{-1} \mathbf{W})\right).$$

Writing $\mathbf{V} = \mathbf{O}' \mathbf{D} \mathbf{O}$, where \mathbf{O} is orthogonal and \mathbf{D} is the diagonal matrix of ordered eigenvalues, it is shown in Yang and Berger (1994) that the full conditional can be transformed to

$$\begin{aligned} \pi(\mathbf{D}, \mathbf{O} \mid \boldsymbol{\theta}, \boldsymbol{\beta}) &\propto \frac{1}{|\mathbf{D}|^r} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{O} \mathbf{D}^{-1} \mathbf{O}' \mathbf{W})\right) 1_{\{d_1 > d_2, \dots, > d_k\}} d\mathbf{D} d\mathbf{O} \\ &= \frac{1}{|\mathbf{D}|^r} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{D}^{-1} \mathbf{O}' \mathbf{W} \mathbf{O})\right) 1_{\{d_1 > d_2, \dots, > d_k\}} d\mathbf{D} d\mathbf{O} \end{aligned}$$

To sample from the full conditional dist. of \mathbf{O} given \mathbf{D} , write

$$\mathbf{O} = (\mathbf{O}_{12}\mathbf{O}_{13}\cdots\mathbf{O}_{1p})(\mathbf{O}_{23}\cdots\mathbf{O}_{2p})\cdots(\mathbf{O}_{p-1,p})\mathbf{D}_\epsilon,$$

with \mathbf{O}_{ij} being a simple orthogonal matrix such as

$$\mathbf{O}_{ij} = \mathbf{O}_{ij}(o_{ij}) = \begin{matrix} & i & & j & \\ i & \left(\begin{array}{ccccc} I & 0 & 0 & 0 & 0 \\ 0 & \cos o_{ij} & 0 & -\sin o_{ij} & 0 \\ 0 & 0 & I & 0 & 0 \\ 0 & \sin o_{ij} & 0 & \cos o_{ij} & 0 \\ 0 & 0 & 0 & 0 & I \end{array} \right) & & & \\ j & & & & \end{matrix}, \quad (10)$$

where $-\pi/2 < o_{ij} \leq \pi/2$, and \mathbf{D}_ϵ being a diagonal matrix with diagonal elements ± 1 (see Anderson, Olkin, and Underhill, 1987). Insert this in (9), and consider Gibbs sampling of the o_{ij} .

- No need for sampling D_ϵ .
- The full cond. pdf of o_{ij} can be written in the form (for constants a, b, c depending on the other parameters and \mathbf{W}),

$$\pi(o_{ij}) = \pi(o_{ij} \mid \mathbf{D}, \{o_{jk} : jk \neq ij\}) \propto e^{[a \cos^2(o_{ij}) + b \cos(o_{ij}) + c \sin(o_{ij})]}.$$

- A simple rejection sampler to draw from this can be constructed as follows:
 - Find the mle \hat{o}_{ij} . This requires solving a quartic equation.
 - Compute the observed Fisher information \hat{I}_{ij} (i.e., compute the second derivate of $-\log \pi(o_{ij})$ and plug in the mle).
 - Use, as a proposal $p(o_{ij})$, the t-distribution with 4 degrees of freedom and mean and variance \hat{o}_{ij} and \hat{I}_{ij}^{-1} , constrained to the interval $(-\frac{\pi}{2}, \frac{\pi}{2})$.
- Compute $K = \sup_{\{-\pi/2 < o_{ij} < \pi/2\}} \frac{\pi(o_{ij})}{p(o_{ij})}$.
- Do rejection sampling with probability $\pi(o_{ij})/[Kp(o_{ij})]$.

Comments

- We have compare the computational methods. The Gibbs Sampling of the new method converge to the target distribution quite fast.
- We have done the simulation study to see the frequentist performance of the recommend priors with other priors.
- We complete the story and propose a particular objective prior for use in normal hierarchical models, based on considerations of admissibility, ease of implementation (including computational considerations), and performance.
- The method will be extended to three level normal hierarchical models.