

# Bayesian Model Robustness via Disparities

Giles Hooker

Department of Statistical Science

Cornell University

Ithaca, NY 14850

Anand N. Vidyashankar

Department of Statistics

George Mason University

Fairfax, VA, 22030

## **Author's Footnote:**

Giles Hooker is Assistant Professor, Department of Biological Statistics and Computational Biology, Cornell University, Ithaca, NY 14850 (email: giles.hooker@cornell.edu). Anand N. Vidyashankar is Professor, Department of Statistics, George Mason University, VA, 22030 (email: avidyash@gmu.edu). Giles Hooker's research was supported by NSF grant DEB-0813734 and the Cornell University Agricultural Experiment Station federal formula funds Project No. 150446. Anand Vidyashankar's research was supported in part by a grant from NSF DMS 000-03-07057 and also by grants from the NDCHealth Corporation. All computations were performed in R.

## Abstract

This paper develops a methodology for robust Bayesian inference through the use of disparities. Metrics such as Hellinger distance and negative exponential disparity have a long history in robust estimation in frequentist inference. We demonstrate that an equivalent robustification may be made in Bayesian inference by substituting an appropriately scaled disparity for the log likelihood to which standard Monte Carlo Markov Chain methods may be applied. A particularly appealing property of minimum-disparity methods is that while they yield robustness, the resulting parameter estimates are also efficient when the posited probabilistic model is correct. We demonstrate that a similar property holds for disparity-based Bayesian inference. We further show that in the Bayesian setting, it is also possible to extend these methods to robustify regression models, random effects distributions and other hierarchical models. The methods are demonstrated on real world data. Supplementary materials including simulation studies and code are available in an online appendix.

KEYWORDS: Deviance test, Kernel density, Hellinger distance, Negative exponential disparity, MCMC, Bayesian Inference, Posterior, Outliers, and Inliers.

## 1. INTRODUCTION

In this paper we develop a new methodology for providing robust inference in a Bayesian context. When the data at hand are suspected of being contaminated with large outliers it is standard practice to account for these either (1) by postulating a heavy-tailed distribution, (2) by viewing the data as a mixture, with the contamination explicitly occurring as a mixture component or (3) by employing priors that penalize large values of a parameter (see Berger, 1994; Albert, 2009; Andrade and O’Hagan, 2006). In the context of frequentist inference, these issues are investigated using methods such as M-estimation, R-estimation etc. and are part of standard robustness literature (see Hampel et al., 1986; Maronna et al., 2006; Jurečková and Sen, 1996). As is the case for Huberized loss functions in frequentist inference, even though these approaches provide robustness they lead to a loss of precision when contamination is not present when (1) and (2) above hold or to a distortion of prior knowledge when (3) holds. This paper develops an alternative systematic Bayesian approach, based on disparity theory, that is shown to provide robust inference without

loss of efficiency for large samples.

In parametric frequentist inference using independent and identically distributed (i.i.d.) data, several authors (Beran, 1977; Tamura and Boos, 1986; Simpson, 1987, 1989; Cheng and Vidyashankar, 2006) have demonstrated that the dual goal of efficiency and robustness is achievable by using the minimum Hellinger distance estimator (MHDE). In the i.i.d. context, MHDE estimators are defined by minimizing the Hellinger distance between a postulated parametric density  $f_\theta(\cdot)$  and a non-parametric estimate  $g_n(\cdot)$  over the  $p$ -dimensional parameter space  $\Theta$ ; that is,

$$\hat{\theta}_{HD} = \arg \inf_{\theta \in \Theta} \int \left( g_n^{1/2}(x) - f_\theta^{1/2}(x) \right)^2 dx. \quad (1)$$

Typically, for continuous data,  $g_n(\cdot)$  is taken to be a kernel density estimate; if the probability model is supported on discrete values, the empirical distribution is used. More generally, Lindsay (1994) introduced the concept of a minimum disparity procedure; developing a class of divergence measures that have similar properties to minimum Hellinger distance estimates. These have been further developed in Basu et al. (1997) and Park and Basu (2004). Recently, Hooker and Vidyashankar (2011a) have extended these methods to a non-linear regression framework.

A remarkable property of disparity-based estimates is that while they confer robustness, they are also first-order efficient. That is, they obtain the information bound when the postulated density  $f_\theta(\cdot)$  is correct. In this paper we develop robust Bayesian inference using disparities. We show that appropriately scaled disparities approximate  $n$  times the negative log-likelihood near the true parameter values. We use this as a motivation to replace the log likelihood in Bayes rule with a disparity to create what we refer to as the “D-posterior”. We demonstrate that this technique is readily amenable to Markov Chain Monte Carlo (MCMC) estimation methods. Finally, we establish that the expectation of the D-posterior is asymptotically efficient and the resulting credible intervals provide asymptotically accurate coverage, when the proposed parametric model is correct.

Disparity-based robustification in Bayesian inference can be naturally extended to a regression framework through the use of conditional density estimation as discussed in Hooker and Vidyashankar (2011b). We pursue this extension to hierarchical models and replace various terms in the hierarchy with disparities. This creates a novel “plug-in procedure” – allowing the robustification of inference with respect to particular distributional assumptions in complex models. We develop this principle and demonstrate its utility on a number of examples. The use of a disparity

within a Bayesian context imposes an additional computational burden through the estimation of a kernel density estimate and the need to run MCMC methods. Our analysis and simulations demonstrate that while the use of MCMC significantly increases computational costs, the additional cost of the use of disparities is on the order of a factor between 2 and 10, remaining implementable for many applications.

The use of divergence measures for outlier analysis in a Bayesian context has been considered in Dey and Birmiwal (1994) and Peng and Dey (1995). Most of this work is concerned with the use of divergence measures to study Bayesian robustness when the priors are contaminated and to diagnose the effect of outliers. The divergence measures are computed using MCMC techniques. More recently, Zhan and Hettmansperger (2007) and Szpiro and Lumley (2011) have developed analogues of R-estimates and Bayesian Sandwich estimators. These methods can be viewed to be extensions of robust frequentist methods to Bayesian context. By contrast, our paper is based on explicitly replacing the likelihood with a disparity in order to provide a systematic approach to obtain inherently robust and efficient inference.

The remainder of the paper is structured as follows: we provide a formal definition of the disparities in Section 2. Disparity-based Bayesian inference are developed in Section 3. Robustness and efficiency of these estimates are demonstrated theoretically and through a simulation for i.i.d. data in Section 4. The methodology is extended to regression models in Section 5. The plug-in procedure is presented in Section 7 through an application to a one-way random-effects model. Some techniques in dimension reduction for regression problems are given in Section 6. Section 8 is devoted to two real-world data sets where we apply these methods to generalized linear mixed models and a random-slope random-intercept models for longitudinal data. Proofs of technical results and details of simulation studies are relegated to the appendix.

## 2. DISPARITIES AND THEIR NUMERICAL APPROXIMATIONS

In this section we describe a class of disparities and numerical procedures for evaluating them. These disparities compare a proposed parametric family of densities to a non-parametric density estimate. We assume that we have i.i.d. observations  $X_i$  for  $i = 1, \dots, n$  from some density  $h(\cdot)$ .

We let  $g_n$  be the kernel density estimate:

$$g_n(x) = \frac{1}{nc_n} \sum_{i=1}^n K\left(\frac{x - X_i}{c_n}\right) \quad (2)$$

where the kernel  $K$  density and  $c_n$  is a bandwidth for the kernel. If  $c_n \rightarrow 0$  and  $nc_n \rightarrow \infty$  it is known that  $g_n(\cdot)$  is an  $L_1$ -consistent estimator of  $h(\cdot)$  (Devroye and Györfi, 1985). In practice, a number of plug-in bandwidth choices are available for  $c_n$  (e.g. Silverman, 1982; Sheather and Jones, 1991; Engel et al., 1994).

We begin by reviewing the class of disparities described in Lindsay (1994). The definition of disparities involves the residual function,

$$\delta_{\theta,g}(x) = \frac{g(x) - f_{\theta}(x)}{f_{\theta}(x)}, \quad (3)$$

defined on the support of  $f_{\theta}(x)$  and a function  $G : [-1, \infty) \rightarrow \mathcal{R}$ .  $G(\cdot)$  is assumed to be strictly convex and thrice differentiable with  $G(0) = 1$ ,  $G'(0) = 0$  and  $G''(0) = 1$ . The disparity between  $f_{\theta}$  and  $g_n$  is defined to be

$$D(g_n, f_{\theta}) = \int_{\mathcal{R}} G(\delta_{\theta,g_n}(x)) f_{\theta}(x) dx. \quad (4)$$

An estimate of  $\theta$  obtained by minimizing (4) is called a *minimum disparity estimator*. Under differentiability assumptions, this is equivalent to solving the equation

$$\int A(\delta_{\theta}(x)) \nabla_{\theta} f_{\theta}(x) dx = 0,$$

where  $A(\delta) = G(\delta) - (1 + \delta)G'(\delta)$  and  $\nabla_{\theta}$  indicates the derivative with respect to  $\theta$ .

This framework contains Kullback-Leibler divergence as approximation to the likelihood:

$$KL(g_n, f_{\theta}) = \int (\log f_{\theta}(x)) g_n(x) dx \approx \frac{1}{n} \sum_{i=1}^n \log f_{\theta}(x_i)$$

for the choice  $G(\delta) = -(\delta + 1) \log(\delta + 1) + a$  for any constant  $a$ . We note that the choice of  $a$  is arbitrary. In particular, we will assume  $a = 1$  so that  $G(0) = 1$ . The squared Hellinger disparity (HD) corresponds to the choice  $G(x) = [(x + 1)^{1/2} - 1]^2 + 1$ . While robust statistics is typically concerned with the impact of outliers, the alternate problem of *inliers* – defined as nominally-dense regions that lack empirical data and consequently small values of  $\delta_{\theta,g_n}(x)$  – can also cause instability. It has been illustrated in the literature that HD down weighs the effect of large values

of  $\delta_{\theta, g_n}(x)$  (outliers) relative to the likelihood but magnifies the effect of inliers. An alternative, the negative exponential disparity, based on the choice  $G(x) = e^{-x}$  down weighs the effect of both outliers and inliers.

The integrals involved in (4) are not analytically tractable and the use of Monte Carlo integration to approximate the objective function has been suggested in Cheng and Vidyashankar (2006). More specifically, if  $z_1, \dots, z_N$  are i.i.d. random samples generated from  $g_n(\cdot)$ , one can approximate  $D(g_n, f_\theta)$  by

$$\hat{D}(g_n, f_\theta) = \frac{1}{N} \sum_{i=1}^N G(\delta_{\theta, g_n}(z_i)) \frac{f_\theta(z_i)}{g_n(z_i)}. \quad (5)$$

The  $z_i$  can be efficiently generated in the form  $z_i = c_n W_i + X_{N_i}$  for  $W_i$  a random variable generated according to  $K$  and  $N_i$  sampled uniformly from the integers  $1, \dots, N$ . In the specific case of Hellinger distance approximation, the above reduces to

$$\widehat{HD}^2(g_n, f_\theta) = 2 - \frac{2}{N} \sum_{i=1}^N \frac{f_\theta^{1/2}(z_i)}{g_n^{1/2}(z_i)}.$$

The use of a fixed set of Monte Carlo samples from  $g_n(\cdot)$  when optimizing for  $\theta$  provides a stochastic approximation to an objective function that remains a smooth function of  $\theta$  and hence avoids the need for complex stochastic optimization. Similarly, in the present paper, we hold the  $z_i$  constant when applying MCMC methods to generate samples from the posterior distribution in order to improve their mixing properties. If  $f_\theta$  is Gaussian with  $\theta = (\mu, \sigma)$ , Gauss-Hermite quadrature rules can be used to avoid Monte Carlo integration, leading to improved computational efficiency in some circumstances. In this case we have

$$\tilde{D}(g_n, f_\theta) = \sum_{i=1}^M w_i(\theta) G(\delta_{\theta, n}(\xi_i(\theta))), \quad (6)$$

where the  $\xi_i(\theta)$  and  $w_i(\theta)$  are the points and weights for a Gauss-Hermite quadrature scheme for parameters  $\theta = (\mu, \sigma)$ . The choice between (5) and (6) depends on the disparity and the problem under investigation. When  $g_n(\cdot)$  has many local modes, (6) can result in choosing parameters for which some quadrature point coincides with a local modes. However, (5) can be rendered unstable by the factor  $f_\theta(z_i)/g_n(z_i)$  for  $\theta$  far from the optimum. In general, we have found (5) preferable when using Hellinger distance, but that (6) performs better with negative exponential disparity. The relative computational cost of  $\hat{D}(g_n, f_\theta)$  and  $\tilde{D}(g_n, f_\theta)$  in various circumstances is discussed in Section 6.3.

### 3. THE D-POSTERIOR AND MCMC METHODS

We begin this section by a heuristic description of the second-order approximation of  $KL(f_\theta, g_n)$  by  $D(f_\theta, g_n)$ . A Taylor expansion of  $KL(f_\theta, g_n)$  about  $\theta$  has as its first two terms:

$$\nabla_\theta KL(g_n, f_\theta) = \int [\nabla_\theta f_\theta(x)] (\delta_{\theta, g_n}(x) + 1) dx \quad (7)$$

$$\begin{aligned} \nabla_\theta^2 KL(g_n, f_\theta) &= \int \left[ \nabla_\theta^2 f_\theta(x) - \frac{1}{f_\theta(x)} (\nabla_\theta f_\theta(x)) (\nabla_\theta f_\theta(x))^T \right] (\delta_{\theta, g_n}(x) + 1) dx. \\ &= \int \left[ \frac{\nabla_\theta^2 f_\theta(x)}{f_\theta(x)} - \left( \frac{\nabla_\theta f_\theta(x)}{f_\theta(x)} \right) \left( \frac{\nabla_\theta f_\theta(x)}{f_\theta(x)} \right)^T \right] g_n(x) dx \end{aligned} \quad (8)$$

where the second term approximates the observed Fisher Information when the bandwidth is small.

The equivalent terms for  $D(g_n, f_\theta)$  are:

$$\nabla_\theta D(g_n, f_\theta) = \int [\nabla_\theta f_\theta(x)] A(\delta_{\theta, g_n}(x)) dx \quad (9)$$

$$\nabla_\theta^2 D(g_n, f_\theta) = \int \nabla_\theta^2 f_\theta(x) A(\delta_{\theta, g_n}(x)) dx - \int \frac{1}{f_\theta(x)} (\nabla_\theta f_\theta(x)) (\nabla_\theta f_\theta(x))^T (\delta_{\theta, g_n}(x) + 1) A'(\delta_{\theta, g_n}(x)) dx.$$

Now, if  $g_n$  is consistent,  $\delta_{\theta, g_n}(x) \rightarrow 0$  almost surely (a.s.). Observing that  $A(0) = 1$ ,  $A'(0) = -1$  from the conditions on  $G$ , we obtain the equality of (8) and (9). The fact that these heuristics yield efficiency was first noticed by Beran (1977) (eq. 1.1).

In the context of the Bayesian methods examined in this paper, inference is based on the posterior

$$P(\theta|x) = \frac{P(x|\theta)\pi(\theta)}{\int P(x|\theta)\pi(\theta)d\theta}, \quad (10)$$

where  $P(x|\theta) = \exp(\sum_{i=1}^n \log f_\theta(x_i))$ . Following the heuristics above, we propose the simple expedient of replacing the log likelihood,  $\log P(x|\theta)$ , in (10) with a disparity:

$$P_D(\theta|g_n) = \frac{e^{-nD(g_n, f_\theta)} \pi(\theta)}{\int e^{-nD(g_n, f_\theta)} \pi(\theta) d\theta}. \quad (11)$$

In the case of Hellinger distance, the appropriate disparity is  $2HD^2(g_n, f_\theta)$  and we refer to the resulting quantity as the *H-posterior*. When  $D(g_n, f_\theta)$  is based on Negative Exponential disparity, we refer to it as *N-posterior*, and *D-posterior* more generally. These choices are illustrated in Figure 1 where we show the approximation of the log likelihood by Hellinger and negative exponential disparities and the effect of adding an outlier to these in a simple normal-mean example.

Throughout the examples below, we employ a Metropolis algorithm based on a symmetric random walk to draw samples from  $P_D(\theta|g_n)$ . While the cost of evaluating  $D(g_n, f_\theta)$  is greater

than the cost of evaluating the likelihood at each Metropolis step, we have found these algorithms to be computationally feasible and numerically stable. Furthermore, the burn-in period for sampling from  $P_D(\theta|g_n)$  and the posterior are approximately the same, although the acceptance rate of the former is approximately around ten percent higher.

After substituting  $-nD(g_n, f_\theta)$  for the log likelihood, it will be useful to define summary statistics of the  $D$ -posterior in order to demonstrate their asymptotic properties. Since the  $D$ -posterior (11) is a proper probability distribution, the Expected  $D$ -a posteriori (EDAP) estimates exist and are given by

$$\theta_n^* = \int_{\Theta} \theta P_D(\theta|g_n) d\theta.$$

and credible intervals for  $\theta$  can be based on the quantiles of  $P_D(\theta|g_n)$ . These quantities are calculated via Monte Carlo integration using the output from the Metropolis algorithm. We similarly define the Maximum  $D$ -a posteriori (MDAP) estimates by

$$\theta_n^+ = \arg \max_{\theta \in \Theta} P_D(\theta|g_n).$$

In the next section we describe the asymptotic properties of EDAP and MDAP estimators. In particular, we establish the posterior consistency, posterior asymptotic normality and efficiency of these estimators and their robustness properties. Differences between  $P_D(\theta, g_n)$  and the posterior do exist and are described below:

1. The disparities  $D(g_n, f_\theta)$  have strict upper bounds; in the case of Hellinger distance  $0 \leq HD^2(g_n, f_\theta) \leq 2$ , the upper bound for negative exponential disparity is  $e$ . This implies that the likelihood part of the  $D$ -posterior,  $\exp(-nD(g_n, f_\theta))$ , is bounded away from zero. Consequently, a proper prior  $\pi(\theta)$  is required in order to normalize  $P_D(\theta|g_n)$ . In particular, uniform priors on unbounded ranges, along with most reference priors, cannot be employed here. Further, the tails of  $P_D(\theta|g_n)$  are proportional to that of  $\pi(\theta)$ . This leads to a breakdown point of 1 (see below). However, these results do not affect the asymptotic behavior of  $P_D(\theta|g_n)$  since the lower bounds decrease with  $n$ , but they do suggest a potential for alternative disparities that allow  $D(g_n, f_\theta)$  to diverge at a sufficiently slow rate to retain robustness.
2. In Bayesian inference for i.i.d. random variables, the log likelihood is a sum of  $n$  terms. This implies that if new data  $X_{n+1}, \dots, X_{n^*}$  are obtained, the posterior for the combined data



$X_1, \dots, X_{n^*}$  can be obtained by using posterior after  $n$  observations,  $P(\theta|X_1, \dots, X_n)$  as a prior  $\theta$ :

$$P(\theta|X_1, \dots, X_{n^*}) \propto P(X_{n+1}, \dots, X_{n^*}|\theta)P(\theta|X_1, \dots, X_n).$$

By contrast,  $D(g_n, f_\theta)$  is generally not additive in  $g_n$ ; hence  $P_D(\theta|g_n)$  cannot be factored as above. Extending arguments in Park and Basu (2004), we conjecture that no disparity that is additive in  $g_n$  will yield both robust and efficient posteriors.

3. While we have found that the same Metropolis algorithms can be effectively used for the D-posterior as would be used for the posterior, it is not possible to use conjugate priors with disparities. This removes the possibility of using conjugacy to provide efficient sampling methods within a Gibbs sampler, although these could be approximated by combining sampling from a conditional distribution with a rejection step. In that respect, disparity-based methods can incur additional computational cost.

The idea of replacing log likelihood in the posterior with an alternative criterion occurs in other settings. See Sollich (2002), for example, in developing Bayesian methods for support vector machines. However, we replace the log likelihood with an approximation that is explicitly designed to be both robust and efficient, rather than as a convenient sampling tool for a non-probabilistic model.

#### 4. ROBUSTNESS AND EFFICIENCY

In this section, we present theoretical results for i.i.d. data to demonstrate that inference based on the D-posterior is both asymptotically efficient and robust. Results for maximum D-*a posteriori* estimates naturally inherit the properties of minimum disparity estimators and hence we focus on EDAP estimators only.

##### 4.1 Efficiency

We recall that under suitable regularity conditions, Expected *a posteriori* estimators are strongly consistent, asymptotically normal and are statistically efficient; (see Ghosh et al., 2006, Theorems 4.2-4.3). Our results in this section show that this property continues to hold for EDAP estimators under regularity conditions on  $G(\cdot)$  when the model  $\{f_\theta : \theta \in \Theta\}$  contains the true distribution.

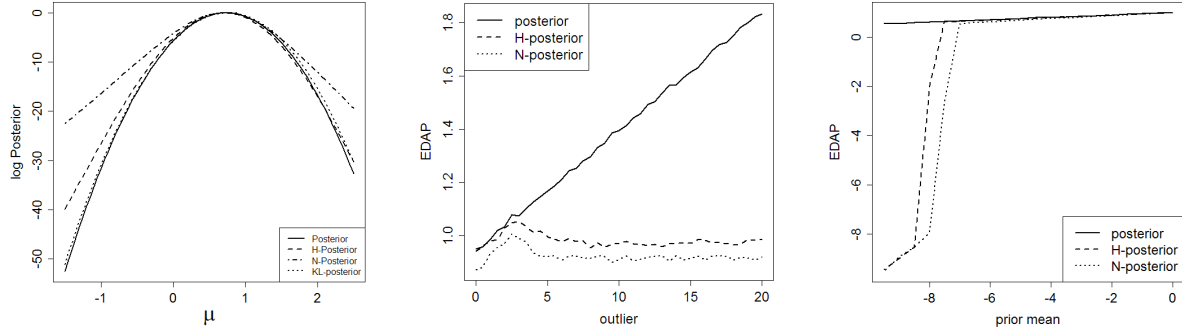


Figure 1: Left: A comparison of log posteriors for  $\mu$  with data generated from  $N(\mu, 1)$  with  $\mu = 1$  using an  $N(0, 1)$  prior for  $\mu$ . Middle: influence of an outlier on expected *D-a posteriori* (EDAP) estimates of  $\mu$  as the value of the outlier is changed from 0 to 20. Right: influence of the prior as the prior mean is changed from 0 to -10.

We define

$$I^D(\theta) = \nabla_{\theta}^2 D(g, f_{\theta}), \text{ and } \hat{I}_n^D(\theta) = \nabla_{\theta}^2 D(g_n, f_{\theta})$$

as the disparity information and  $\theta_g$  the parameter that minimizes  $D(g, f_{\theta})$  (note that  $\theta_g$  here depends on  $g$ ). We note that if  $g = f_{\theta_0}$ ,  $I^D(\theta_g)$  is exactly equal to the Fisher information for  $\theta_0$ .

The proofs of our asymptotic results rely on the assumptions listed below. Among these are that minimum disparity estimators are strongly consistent and efficient; this in turn relies on further assumptions, some of which make those listed below redundant. They are given here to maximize the mathematical clarity of our arguments. We assume that  $X_1, \dots, X_n$  are i.i.d. generated from some distribution  $g(x)$  and that a parametric family,  $f_{\theta}(x)$  has been proposed for  $g(x)$  where  $\theta$  has distribution  $\pi$ . To demonstrate efficiency, we assume

(A1)  $g(x) = f_{\theta_g}(x)$  is a member of the parametric family.

(A2)  $G$  has three continuous derivatives with  $G'(0) = 0$ ,  $G''(0) = 1$  and  $|G'''(0)| \leq \infty$ .

(A3)  $\nabla_{\theta}^2 D(g, f_{\theta})$  is positive definite and continuous in  $\theta$  at  $\theta_g$  and continuous in  $g$  with respect to the  $L_1$  metric.

(A4) For any  $\delta > 0$ , there exists  $\epsilon > 0$  such that

$$\sup_{|\theta - \theta_g| > \delta} (D(g, f_\theta) - D(g, f_{\theta_g})) > \epsilon$$

(A5) The parameter space  $\Theta$  is compact.

(A6) The minimum disparity estimator,  $\hat{\theta}_n$ , satisfies  $\hat{\theta}_n \rightarrow \theta_g$  almost surely and  $\sqrt{n}(\hat{\theta}_n - \theta_g) \xrightarrow{d} N(0, I^D(\theta)^{-1})$ .

Our first result concerns the limit distribution for the posterior density of  $\sqrt{n}(\theta - \hat{\theta}_n)$ , which demonstrates that the D-posterior converges in  $L_1$  to a Gaussian density centered on the minimum disparity estimator  $\hat{\theta}_n$  with variance  $[nI^D(\hat{\theta}_n)]^{-1}$ . This establishes that credible intervals based on either  $P_D(\theta|x_1, \dots, x_n)$  or from  $N(\hat{\theta}_n, I_n^D(\hat{\theta}_n)^{-1})$  will be asymptotically accurate.

**Theorem 1.** *Let  $\hat{\theta}_n$  be the minimum disparity estimator of  $\theta_g$ ,  $\pi(\theta)$  be any prior that is continuous and positive at  $\theta_g$  with  $\int_{\Theta} \|\theta\|_2 \pi(\theta) d\theta < \infty$  where  $\|\cdot\|_2$  is the usual 2-norm, and  $\pi_n^{*D}(t)$  be the D-posterior density of  $t = (t_1, \dots, t_p) = \sqrt{n}(\theta - \hat{\theta}_n)$ . Under conditions (A2)-(A6),*

$$\lim_{n \rightarrow \infty} \int \left| \pi_n^{*D}(t) - \left( \frac{|I^D(\theta_g)|}{2\pi} \right)^{p/2} e^{-\frac{1}{2} t' I^D(\theta_g) t} \right| dt \xrightarrow{a.s.} 0. \quad (12)$$

Furthermore, (12) also holds with  $I^D(\theta_g)$  replaced with  $\hat{I}_n^D(\hat{\theta}_n)$ .

Our next theorem is concerned with the efficiency and asymptotic normality of EDAP estimates.

**Theorem 2.** *Assume conditions (A2)-(A6) and  $\int_{\Theta} \|\theta\|_2 \pi(\theta) d\theta < \infty$ , then  $\sqrt{n}(\theta_n^* - \hat{\theta}_n) \xrightarrow{a.s.} 0$  where  $\theta_n^*$  is the EDAP estimate. Further,  $\sqrt{n}(\theta_n^* - \theta_g) \xrightarrow{d} N(0, I^D(\theta_g))$ .*

The proofs of these theorems are deferred to the Appendix A, but the following remarks concerning the assumptions (A1)-(A6) are in order:

1. Assumption A1 states that  $g$  is a member of the parametric family. When this does not hold, a central limit theorem can be derived for  $\hat{\theta}_n$  but the variance takes a sandwich-type form; see Beran (1977) in the case of Hellinger distance. For brevity, we have followed Basu et al. (1997) and Park and Basu (2004) in restricting to the parametric case.

2. Assumptions A2-A4 are required for the regularity and identifiability of the parametric family  $f_\theta$  in the disparity  $D$ . Specific conditions for (A6) to hold are given in various forms in Beran (1977); Basu et al. (1997); Park and Basu (2004) and Cheng and Vidyashankar (2006), see conditions in Appendix A.
3. We assume that the parameter space,  $\Theta$ , is compact; this result is also used in conditions that guarantee (A6). As noted in Beran (1977), as well as others, the result continues to hold if  $\Theta$  can be appropriately embedded in a compact space. Alternatively, Cheng and Vidyashankar (2006) assume local compactness.
4. The proofs of these results follow the same principles as those given for posterior asymptotic efficiency (see Ghosh et al., 2006, for example). However, here we rely on the second-order convergence of the disparity to the likelihood at appropriate rates and the consequent asymptotic efficiency of minimum-disparity estimators.
5. Since the structure of the proof only requires second-order properties and appropriate rates of convergence, we can replace  $D(g_n, f_\theta)$  for i.i.d. data with an appropriate disparity-based term for more complex models as long as (A6) can be shown hold. In particular, the results in Hooker and Vidyashankar (2011a) and Hooker and Vidyashankar (2011b) suggest that the disparity methods for regression problems detailed in Sections 5 and 6 will also yield efficient estimates.

## 4.2 Robustness

To describe robustness, we view our estimates as functionals  $T_n(h)$  of a density  $h$ . In particular, we examine the EDAP estimate

$$T_n(h) = \frac{\int \theta e^{-nD(h, f_\theta)} \pi(\theta) d\theta}{\int e^{-nD(h, f_\theta)} \pi(\theta) d\theta} \quad (13)$$

and note that in contrast to classical approaches to analyzing robustness the interaction between the disparity and the prior requires us to make the dependence of  $T_n$  on  $n$  explicit. We analyze the behavior of  $T_n(h)$  under the sequence of perturbations  $h_{k,\epsilon}(x) = (1 - \epsilon)g(x) + \epsilon t_k(x)$  for any sequence of densities  $t_k(\cdot)$  and  $0 \leq \epsilon \leq 1$ . We measure robustness via two quantities, namely the

influence function:

$$IF_T(t_k) = \lim_{\epsilon \rightarrow 0} \frac{T_n((1 - \epsilon)g + \epsilon t_k) - T_n(g)}{\epsilon} \quad (14)$$

(Hampel, 1974) and the breakdown point:

$$B(T) = \sup \left\{ \epsilon : \sup_k |T_n((1 - \epsilon)g + \epsilon t_k)| < \infty \right\}, \quad (15)$$

(see Huber (1981)). EDAP estimates turn out to be highly robust. In fact, while the most common robust estimators have breakdown points of 0.5, for most of the commonly-used robust disparities the D-posterior breakdown point is 1. As described previously this is due to the fact that these disparities are bounded above. We point out here that the Kullback-Leibler disparity is not bounded above and is not robust, both in frequentist and in Bayesian settings.

**Theorem 3.** *Let  $D(g, f_\theta)$  be bounded for all  $\theta$  and all densities  $g$  and let  $\int \|\theta\|_2 \pi(\theta) d\theta < \infty$ , then the breakdown point of the EDAP is 1.*

The condition that  $D(g, f_\theta)$  be bounded holds if  $|G(\cdot)|$  is bounded ; this is assumed in Park and Basu (2004) and holds for the negative exponential disparity and Hellinger distance ( $0 \leq 2HD(g, f_\theta) \leq 4$ ). The proof of this theorem is given in Appendix B.

To examine the influence function, we assume that the limit may be taken inside all integrals in (14) and obtain

$$IF(\theta; g, t_k) = E_{P_D(\theta|g)} [\theta C_{nk}(\theta, g)] - [E_{P_D(\theta|g)} \theta] [E_{P_D(\theta|g)} C_{nk}(\theta, g)].$$

where  $E_{P_D(\theta|g)}$  indicates expectation with respect to the D-posterior with density  $g$  and

$$\begin{aligned} C_{nk}(\theta, g) &= \frac{d}{d\epsilon} n \int G \left( \frac{h_{k,\epsilon}(x)}{f_\theta(x)} - 1 \right) f_\theta(x) dx \Big|_{\epsilon=0} \\ &= n \int G' \left( \frac{g(x)}{f_\theta(x)} - 1 \right) (g(x) - t_k(x)) dx. \end{aligned}$$

Thus, if we can establish the posterior integrability of  $C_{nk}(\theta, f)$  and  $\theta C_{nk}(\theta, f)$ , the influence function will be everywhere finite. This is trivially true if  $G'(\cdot)$  is bounded, as is the case for the negative exponential disparity. However,  $G'$  is not bounded at -1 for Hellinger distance.

These results indicate an extreme form of robustness that results from the fact that the disparity approximation to the likelihood is weak in its tails. However, as  $n$  increases the approximation

$T_n(h) - \hat{\theta}(h) = o(n^{-1})$  can be shown to hold where  $\hat{\theta}(h)$  gives the parameter that minimizes  $D(h, f_{\theta})$  (see Appendix B). Following this we also observe that the  $\alpha$ -level influence functions converge at a  $n^{-1}$  rate. This statement can be refined by the asymptotic expansion for a one-parameter family

$$T_n(h) = \hat{\theta}(h) + n^{-1} I^D(h)^{-1} \left( \frac{d^3}{d\theta^2} D(h, f_{\hat{\theta}(h)}) + \frac{\nabla_{\theta} \pi(\hat{\theta}(h))}{\pi(\hat{\theta}(h))} \right) + o_p(n^{-1})$$

where  $a_3$  gives the third derivative of the disparity (see Appendix B). For a location family for which  $\hat{\theta}(h)$  is equivariant, the only non-equivariant term in this expansion is  $\nabla_{\theta} \pi(\hat{\theta}(h)) / \pi(\hat{\theta}(h))$  which also appears in the expansion of the usual posterior (see Ghosh et al., 2006). The additional influence of the prior due to the weak tails of the disparity first appears in terms of  $o_p(n^{-3/2})$ ; although we note that robustness is a finite-sample property and an asymptotic analysis of it should be treated with some caution.

### 4.3 Simulation Studies

To illustrate the small sample performance of D-posteriors, we undertook a simulation study for i.i.d. data from Gaussian and log-Gamma distributions and these are reported in detail in Online Appendix D.1. In both cases, we used the same random-walk Metropolis algorithm to sample from the posterior and the H- and N-posteriors. Here we show very similar performance between disparity-based methods and the log likelihood for uncontaminated data as well as for a Huber M-estimate. When the data are contaminated with outliers disparity-based methods increasingly out-perform the others in terms of bias and coverage as the size of the outlier increases.

The H-posterior, however, demonstrated higher variance for the log-Gamma distribution due to its tendency to create inliers to which Hellinger distance is sensitive. Incorporating outliers into the data strongly biased the posterior for both distributions, but the disparity-based methods were essentially unaffected. The effect of the size of the outlier is investigated in the second plot of Figure 1 where the EDAPs for both disparities smoothly down-weight the outlying point, while the posterior is highly sensitive to it. The influence of the prior is investigated in the right-hand plot of Figure 1 where we observe that the EAP and EDAP estimates are essentially identical until the prior is about 9 standard deviations from the mean of the data: at this point the prior dominates. However, we note that this picture will depend strongly on the prior chosen; a less informative prior will have a smaller range of dominance.

## 5. DISPARITIES BASED ON CONDITIONAL DENSITY FOR REGRESSION MODELS

The discussion above, along with most of the literature on disparity estimation, has focussed on i.i.d. data in which a kernel density estimate may be calculated. The restriction to i.i.d. contexts severely limits the applicability of disparity-based methods. We extend these methods to non-i.i.d. data settings via the use of conditional density estimates. This extension is studied in the frequentist context in the case of minimum-disparity estimates for parameters in non-linear regression in Hooker and Vidyashankar (2011b).

Consider the classical regression framework with data  $(Y_1, X_1), \dots, (Y_n, X_n)$  is a collection of i.i.d. random variables where inference is made conditionally on  $X_i$ . For continuous  $X_i$ , a non-parametric estimate of the conditional density of  $y|x$  is given by Hansen (2004):

$$g_n^c(y|x) = \frac{\frac{1}{nc_{n1}c_{2,n}} \sum_{i=1}^n K\left(\frac{y-Y_i}{c_{n1}}\right) K\left(\frac{\|x-X_i\|}{c_{2,n}}\right)}{\frac{1}{nc_{2,n}} \sum_{i=1}^n K\left(\frac{\|x-X_i\|}{c_{2,n}}\right)}. \quad (16)$$

Under a parametric model  $f_\theta(y|X_i)$  assumed for the conditional distribution of  $Y_i$  given  $X_i$ , we define a disparity between  $g_n^c$  and  $f_\theta$  as follows:

$$D^c(g_n^c, f_\theta) = \sum_{i=1}^n D(g_n^c(\cdot|X_i), f_\theta(\cdot|X_i)). \quad (17)$$

As before, for Bayesian inference we replace the log likelihood by negative of the conditional disparity (17); that is,

$$e^{l(Y|X_i, \theta)} \pi(\theta) \approx e^{-D^c(g_n^c, f_\theta)} \pi(\theta).$$

In the case of simple linear regression,  $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$ ,  $\theta = (\beta_0, \beta_1, \sigma^2)$  and  $f_\theta(\cdot|X_i) = \phi_{\beta_0 + \beta_1 X_i, \sigma^2}(\cdot)$  where  $\phi_{\mu, \sigma^2}$  is Gaussian density with mean  $\mu$  and variance  $\sigma^2$ .

The use of a conditional formulation, involving a density estimate over a multidimensional space, produces an asymptotic bias in MDAP and EDAP estimates similar to that found in Tamura and Boos (1986), who also note that this bias is generally small. Section 6 proposes two alternative formulations that reduce the dimension of the density estimate and also eliminate this bias.

When the  $X_i$  are discrete, (16) reduces to a distinct conditional density for each level of  $X_i$ . For example, in a one-way ANOVA model  $Y_{ij} = X_i + \epsilon_{ij}$ ,  $j = 1, \dots, n_i$ ,  $i = 1, \dots, N$ , this reduces to

$$g_n^c(y|X_i) = \frac{1}{n_i c_n} \sum_{j=1}^{n_i} K\left(\frac{y - Y_{ij}}{c_n}\right).$$

We note that in this case the bias noted above does not appear. However When the  $n_i$  are small, or for high-dimensional covariate spaces the non-parametric estimate  $g_n(y|X_i)$  can become inaccurate. The marginal methods discussed in Section 6 can also be employed in this case.

## 6. DIMENSION REDUCTION METHODS

The conditional disparity formulation outlined above requires the estimation of the density of a response conditioned on a potentially high-dimensional set of covariates; this can result in asymptotic bias and poor performance in small samples. In this section, we explore two methods for reducing the dimension of the conditioning spaces and removing the bias. The first is referred to as the “marginal formulation” and requires only a univariate, unconditional, density estimate. This is a Bayesian extension of the approach suggested in Hooker and Vidyashankar (2011a). It is more stable and computationally efficient than schemes based on nonparametric estimates of conditional densities. However, in a linear-Gaussian model with Gaussian covariates, it requires an external specification of variance parameters for identifiability. For this reason, we propose a two-step Bayesian estimate. The asymptotic analysis for i.i.d. data can be extended to this approach by using the technical ideas in Hooker and Vidyashankar (2011a).

The second method produces a conditional formulation that relies on the structure of a homoscedastic location-scale family  $P(Y_i|X_i, \theta, \sigma) = f_\sigma(y - \eta(X_i, \theta))$  and we refer to it as the “conditional-homoscedastic” approach. This method provides a full conditional estimate by replacing a non-parametric conditional density estimate with a two-step procedure as proposed in Hansen (2004). The method involves first estimating the mean function non-parametrically and then estimating a density from the resulting residuals.

### 6.1 Marginal Formulation

Hooker and Vidyashankar (2011a) proposed basing inference on a marginal estimation of residual density in a nonlinear regression problem. A model of the form

$$Y_i = \eta(X_i, \theta) + \epsilon_i$$

is assumed for independent  $\epsilon_i$  from a scale family with mean zero and variance  $\sigma^2$ .  $\theta$  is an unknown parameter vector of interest. A disparity method was proposed based on a density estimate of the



residuals

$$e_i(\theta) = Y_i - \eta(X_i, \theta)$$

yielding the kernel estimate

$$g_n^m(e, \theta, \sigma) = \frac{1}{nc_n} \sum K\left(\frac{e - e_i(\theta)/\sigma}{c_n}\right) \quad (18)$$

and  $\theta$  was estimated by minimizing  $D(\phi_{0,1}(\cdot), g_n^m(\cdot, \theta, \sigma))$  where  $\phi_{0,1}$  is the postulated density. As described above, in a Bayesian context we replace the log likelihood by  $-nD(\phi_{0,1}(\cdot), g_n^m(\cdot, \theta, \sigma))$ . Here we note that although the estimated of  $g_n^m(\cdot, \theta, \sigma)$  need not have zero mean, it is compared to the centered density  $\phi_{0,1}(\cdot)$  which penalizes parameters for which the residuals are not centered.

This formulation has the advantage of only requiring the estimate of a univariate, unconditional density  $g_n^m(\cdot, \theta, \sigma)$ . This reduces the computational cost considerably as well as providing a density estimate that is more stable in small samples.

Hooker and Vidyashankar (2011a) proposed a two-step procedure to avoid identifiability problems in a frequentist context. This involves replacing  $\sigma$  by a robust external estimate  $\tilde{\sigma}$ . It was observed that estimates of  $\theta$  were insensitive to the choice of  $\tilde{\sigma}$ . After an estimate  $\hat{\theta}$  was obtained by minimizing  $D(\phi_{0,1}(\cdot), g_n^m(\cdot, \theta, \tilde{\sigma}))$ , an efficient estimate of  $\sigma$  was obtained by re-estimating  $\sigma$  based on a disparity for the residuals  $e_i(\hat{\theta})$ . A similar process can be undertaken here.

In a Bayesian context a plug-in estimate for  $\sigma^2$  also allows the use of the marginal formulation: an MCMC scheme is undertaken with the plug-in value of  $\sigma^2$  held fixed. A pseudo-posterior distribution for  $\sigma$  can then be obtained by plugging in an estimate for  $\theta$  to a Disparity-Posterior for  $\sigma$ . More explicitly, the following scheme can be undertaken:

1. Perform an MCMC sampling scheme for  $\theta$  using a plug-in estimate for  $\sigma^2$ .
2. Approximate the posterior distribution for  $\sigma^2$  with an MCMC scheme to sample from the D-posterior  $P_D(\sigma^2|y) = e^{-nD(g_n(\cdot, \hat{\theta}, \sigma), \phi_{0,1}(\cdot))} \pi(\sigma^2)$  where  $\hat{\theta}$  is the EDAP estimate calculated above.

This scheme is not fully Bayesian in the sense that fixed estimators of  $\sigma$  and  $\theta$  are used in each step above. However, the ideas in Hooker and Vidyashankar (2011a) can be employed to demonstrate that under these schemes the two-step procedure will result in statistically efficient estimates and

asymptotically correct credible regions. We note that while we have discussed this formulation with respect to regression problems, it can also be employed with the plug-in procedure for random-effects models and we use this in Section 8.2, below.

The formulation presented here resembles the methods proposed in Pak and Basu (1998) based on a sum of disparities between weighted density estimates of the residuals and their expectations assuming the parametric model. For particular combinations of kernels and densities, these estimates are efficient, and the sum of disparities, appropriately scaled, should also be substitutable for the likelihood in order to achieve an alternative D-posterior.

## 6.2 Nonparametric Conditional Densities for Regression Models in Location-Scale Families

Under a homoscedastic location-scale model (where the errors are assumed to be i.i.d.)  $p(Y_i|X_i, \theta, \sigma) = f_\sigma(Y_i - \eta(X_i, \theta))$  where  $f_\sigma$  is a distribution with zero mean, an alternative density estimate may be used. We first define a non-parametric estimate of the mean function

$$m_n(x) = \frac{\sum Y_i K\left(\frac{x-X_i}{c_{2,n}}\right)}{\sum K\left(\frac{x-X_i}{c_{2,n}}\right)}$$

and then a non-parametric estimate of the residual density

$$g_n^{c2}(e) = \frac{1}{nc_{1,n}} \sum K\left(\frac{e - y_i + m_n(X_i)}{c_{1,n}}\right).$$

We then consider the disparity between the proposed  $f_{\theta,\sigma}$  and  $g_n$ :

$$D^{c2}(g_n, \theta, \sigma) = \sum D(g_n^{c2}(\cdot + m(X_i)), f_\sigma(\cdot + \eta(X_i, \theta))).$$

As before,  $-D^{c2}(g_n, f)$  can be substituted for the log likelihood in an MCMC scheme.

Hansen (2004) remarks that in the case of a homoscedastic conditional density,  $g_n^{c2}$  has smaller bias than  $g_n^c$ . This formulation does not avoid the need to estimate the high-dimensional function  $m_n(x)$ . However, the shift in mean does allow the method to escape the identification problems of the marginal formulation while retaining some of its stabilization.

Online Appendix D.2 gives details of a simulation study of both conditional formulations and the marginal formulation above for a regression problem with a three-dimensional covariate. All disparity-based methods perform similarly to using the posterior with the exception of the conditional form in Section 5 when Hellinger distance is used which demonstrates a substantial increase

in variance. We speculate that this is due to the sparsity of the data in high dimensions creating inliers; negative exponential disparity is less sensitive to this problem (Basu et al., 1997).

### 6.3 Computational Considerations and Implementation

Our experience is that the computational cost of employing Disparity-based methods as proposed above is comparable to employing an MCMC scheme for the equivalent likelihood and generally requires an increase in computation time by a factor of between 2 and 10. Further, the comparative advantage of employing estimates (5) versus (6) depends on the context that is used.

Formally, we assume  $N$  Monte Carlo samples is (5) and  $M$  Gauss-Hermite quadrature points in (6) where typically  $M < N$ . In this case, the cost of evaluating  $g_n(z_i)$  in (5) is  $O(nN)$ , but this may be pre-computed before employing MCMC, and the cost of evaluating (5) for a new value of  $\theta$  is  $O(N)$ . In comparison, the use of (6) requires the evaluation of  $g_n(\xi_i(\theta))$  at each iteration at a  $O(nM)$  each evaluation.

Within the context of conditional disparity metrics, we assume  $N$  Monte Carlo points used for each  $X_i$  in the equivalent version of (5) for (17) and note that in this context  $N$  can be reduced due to the additional averaging over the  $X_i$ . The cost of evaluating  $g_n^c(z_j|X_i)$  from (16) for all  $z_j$  and  $X_i$  is  $O(n^2N)$  for (5) and  $O(n^2M)$  for (6). Here again the computation can be carried out before MCMC is employed for (5), requiring  $O(nN)$  operations. In (6) the denominator of (16) can be pre-computed, reducing the computational cost of each iteration to  $O(nM)$ ; however, in this case we will not necessarily expect  $M < N$ . Similar calculations apply to estimates based on  $g_n^{c2}$ .

For marginal disparities  $g_n^m$  in (18) changes for each  $\theta$ , requiring  $O(nM)$  calculations to evaluate (6). Successful use of (5) would require the  $z_i$  to vary smoothly with  $\theta$  and would also require the re-evaluation of  $g_n^m(z_i)$  at a cost of  $O(nN)$  each iteration. Within the context of hierarchical models above,  $g_n$  varies with latent variables and this the use of (5) will generally be more computationally efficient. The cost of evaluating the likelihood is always  $O(n)$ .

While these calculations provide general guidelines to computational costs, the relative efficiency of (5) and (6) strongly depends on the implementation of the procedure. Our simulations have been carried out in the R programming environment where we have found (5) to be computationally cheaper anywhere it can be employed. However, this will be platform-dependent – changing with what routines are given in pre-compiled code, for example – and will also depend strongly on the

details of our implementation.

## 7. DISPARITY METRICS AND THE PLUG-IN PROCEDURE

The disparity-based techniques developed above can be extended to hierarchical models. In particular, consider the following structure for an observed data vector  $Y$  along with an unobserved latent effect vector  $Z$  of length  $n$ :

$$P(Y, Z, \theta) = P_1(Y|Z, \theta)P_2(Z|\theta)P_3(\theta) \quad (19)$$

where  $P_1$ ,  $P_2$  and  $P_3$  are the conditional distributions of  $Y$  given  $Z$  and  $\theta$  the distribution of  $Z$  given  $\theta$  and the prior distribution of  $\theta$ . Any term in this factorization that can be expressed as the product of densities of i.i.d. random variables can now be replaced by a suitably chosen disparity. This creates a *plug-in procedure* in which particular terms of a complete data log likelihood are replaced by disparities. For example, if the middle term is assumed to be a product:

$$P(Z|\theta) = \prod_{i=1}^n p(Z_i|\theta),$$

inference can be robustified for the distribution of the  $Z_i$  by replacing (19) with

$$P_{D_1}(Y, Z, \theta) = P(Y|Z, \theta)e^{-2D(g_n(\cdot; Z), P_2(\cdot|\theta))}P_3(\theta)$$

where

$$g_n(z; Z) = \frac{1}{nc_n} \sum_{i=1}^n K\left(\frac{z - Z_i}{c_n}\right).$$

In an MCMC scheme, the  $Z_i$  will be imputed at each iteration and the estimate  $g_n(\cdot; Z)$  will change accordingly. If the integral is evaluated using Monte Carlo samples from  $g_n$ , these will also need to be updated. The evaluation of  $D(g_n(\cdot; Z), P_2(\cdot|\theta))$  creates additional computational overhead, but we have found this to remain feasible for moderate  $n$ . A similar substitution may also be made for the first term using the conditional approach suggested above.

To illustrate this principle in a concrete example, consider a one-way random-effects model:

$$Y_{ij} = Z_i + \epsilon_{ij}, \quad i = 1, \dots, n, \quad j = 1, \dots, n_i$$

under the assumptions

$$\epsilon_{ij} \sim N(0, \sigma^2), \quad Z_i \sim N(\mu, \tau^2)$$

where the interest is in the value of  $\mu$ . Let  $\pi(\mu, \sigma^2, \tau^2)$  be the prior for the parameters in the model; an MCMC scheme may be conducted with respect to the probability distribution

$$P(Y, Z, \mu, \sigma^2, \tau^2) = \prod_{i=1}^n \left( \prod_{j=1}^{n_i} \phi_{0, \sigma^2}(Y_{ij} - Z_i) \right) \prod_{i=1}^n \phi_{\mu, \tau^2}(Z_i) \pi(\mu, \sigma^2, \tau^2) \quad (20)$$

where  $\phi_{\mu, \sigma^2}$  is the  $N(\mu, \sigma^2)$  density. There are now two potential sources of distributional errors: either in individual observed  $Y_{ij}$ , or in the unobserved  $Z_i$ . Either (or both) possibilities can be dealt with via the plug-in procedure described above.

If there are concerns that the distributional assumptions on the  $\epsilon_{ij}$  are not correct, we observe that the statistics  $Y_{ij} - Z_i$  are assumed to be i.i.d.  $N(0, \sigma^2)$ . We may then form the conditional kernel density estimate:

$$g_n^c(t|Z_i; Z) = \frac{1}{nc_{1,n}} \sum_{j=1}^{n_i} K \left( \frac{t - (Y_{ij} - Z_i)}{c_{1,n}} \right)$$

and replace (20) with

$$P_{D_2}(Y, Z, \mu, \sigma^2, \tau^2) = e^{-\sum_{i=1}^n n_i D(g_n^c(t|Z_i; Z), \phi_{0, \sigma^2}(\cdot))} \prod_{i=1}^n \phi_{\mu, \tau^2}(Z_i) \pi(\mu, \sigma^2, \tau^2).$$

On the other hand, if the distribution of the  $Z_i$  is misspecified, we form the estimate

$$g_n(z; Z) = \frac{1}{nc_{2,n}} \sum_{i=1}^n K \left( \frac{z - Z_i}{c_{2,n}} \right)$$

and use

$$P_{D_1}(X, Y, \mu, \sigma^2, \tau^2) = \prod_{i=1}^n \left( \prod_{j=1}^{n_i} \phi_{0, \sigma^2}(Y_i - Z_i) \right) e^{-nD(g_n(\cdot; Z), \phi_{\mu, \tau^2}(\cdot))} \pi(\mu, \sigma^2, \tau^2)$$

as the D-posterior. For inference using this posterior, both  $\mu$  and the  $Z_i$  will be included as parameters in every iteration, necessitating the update of  $g_n(\cdot; Z)$  or  $g_n^c(\cdot|z; Z)$ . Naturally, it is also possible to substitute a disparity in both places:

$$P_{D_{12}}(Z, Y, \mu, \sigma^2, \tau^2) = e^{-\sum_{i=1}^n n_i D(g_n^c(\cdot|Z_i; Z), \phi_{0, \sigma^2}(\cdot))} e^{-nD(g_n(\cdot; Z), \phi_{\mu, \tau^2}(\cdot))} \pi(\mu, \sigma^2, \tau^2).$$

A simulation study considering all these approaches with Hellinger distance chosen as the disparity is described in Online Appendix D.3. Our results indicate that all replacements with disparities perform well, although some additional bias is observed in the estimation of variance parameters

which we speculate to be due to the interaction of the small sample size with the kernel bandwidth. Methods that replace the random effect likelihood with a disparity remain largely unaffected by the addition of an outlying random effect while for those that do not the estimation of both the random effect mean and variance is substantially biased.

While a formal analysis of this method is beyond the scope of this paper we remark that the use of density estimates of latent variables requires significant theoretical development in both Bayesian and frequentist contexts. In particular, in the context of using  $P_{D_1}$  appropriate inference on  $\theta$  will require local agreement in the integrated likelihoods

$$\begin{aligned} \int \cdots \int \prod_{i=1}^n \left( \prod_{j=1}^{n_i} \phi_{0,\sigma^2}(Y_i - Z_i) \right) e^{-nD(g_n(\cdot; Z), \phi_{\mu, \tau^2}(\cdot))} dZ_1, \dots, dZ_n \\ \approx \int \cdots \int \prod_{i=1}^n \left( \prod_{j=1}^{n_i} \phi_{0,\sigma^2}(Y_{ij} - Z_i) \right) \prod_{i=1}^n \phi_{\mu, \tau^2}(Z_i) dZ_1, \dots, dZ_n. \end{aligned}$$

This can be demonstrated if the  $n_i \rightarrow \infty$  and hence the conditional variance of the  $Z_i$  is made to shrink at an appropriate rate.

## 8. REAL DATA EXAMPLES

### 8.1 Parasite Data

We begin with a one-way random effect model for binomial data. These data come from one equine farm participating in a parasite control study in Denmark in 2008. Fecal counts of eggs of the Equine Strongyle parasites were taken pre- and post- treatment with the drug Pyrantol; the full study is presented in Nielsen et al. (2010). The data used in this example are reported in Online Appendix E.

For our purposes, we model the post-treatment data from each horse as binomial with probabilities drawn from a logit normal distribution. Specifically, we consider the following model:

$$k_i \sim \text{Bin}(N_i, p_i), \text{ logit}(p_i) \sim N(\mu, \sigma^2), \quad i = 1, \dots, n,$$

where  $N_i$  are the pre-treatment egg counts and  $k_i$  are the post-treatment egg counts. We observe the data  $(k_i, N_i)$  and desire an estimate of  $\mu$  and  $\sigma$ . The likelihood for these data are

$$l(\mu, \sigma | k, N) = - \sum_{i=1}^n [k_i \log p_i + (N_i - k_i) \log(1 - p_i)] - \frac{1}{2\sigma^2} \sum_{i=1}^n (\log(p_i) - \mu)^2.$$

We cannot use conditional disparity methods to account for outlying  $k_i$  since we have only one observation per horse. However, we can consider robustifying the  $p_i$  distribution by use of a negative exponential disparity:

$$g_n(\lambda; p_1, \dots, p_n) = \frac{1}{nc_n} \sum K\left(\frac{\lambda - \text{logit}(p_i)}{c_n}\right)$$

$$l^N(\mu, \sigma | k, N) = - \sum_{i=1}^n [k_i \log p_i + (N_i - k_i) \log(1 - p_i)] - nD(g_n(\cdot; p_1, \dots, p_n), \phi_{\mu, \sigma^2}(\cdot))$$

In order to perform a Bayesian analysis,  $\mu$  was given a  $N(0, 5)$  prior and  $\sigma^2$  an inverse Gamma prior with shape parameter 3 and scale parameter 0.5. These were chosen as conjugates to the assumed Gaussian distribution and are defuse enough to be relatively uninformative while providing reasonable density at the maximum likelihood estimates. A random walk Metropolis algorithm was run for this scheme with parameterization  $(\mu, \log(\sigma), \text{logit}(p_1), \dots, \text{logit}(p_n))$  for 200,000 steps with posterior samples collected every 100 steps in the second half of the chain.  $c_n$  was chosen via the method in Sheather and Jones (1991) treating the empirical probabilities as data.

The resulting posterior distributions, given in Figure 2, indicate a substantial difference between the two posteriors, with the N-posterior having higher mean and smaller variance. This suggests some outlier contamination and a plot of a sample of densities  $g_n$  on the right of Figure 2 suggests a lower-outlier with  $\text{logit}(p_i)$  around -4. In fact, this corresponds to observation 5 which had unusually high efficacy in this horse. Removing the outlier results in good agreement between the posterior and the N-posterior. We note that, as also observed in Stigler (1973), trimming observations in this manner, unless done carefully, may not yield accurate credible intervals.

## 8.2 Class Survey Data

Our second data set are from an in-class survey in an introductory statistics course held at Cornell University in 2009. Students were asked to specify their expected income at ages 35, 45, 55 and 65. Responses from 10 American-born and 10 foreign-born students in the class are used as data in this example; the data are presented and plotted in Online Appendix E. Our object is to examine the expected rate of increase in income and any differences in this rate or in the over-all salary level between American and foreign students. From the plot of these data in Figure 4 in Online Appendix E some potential outliers in both over-all level of expected income and in specific

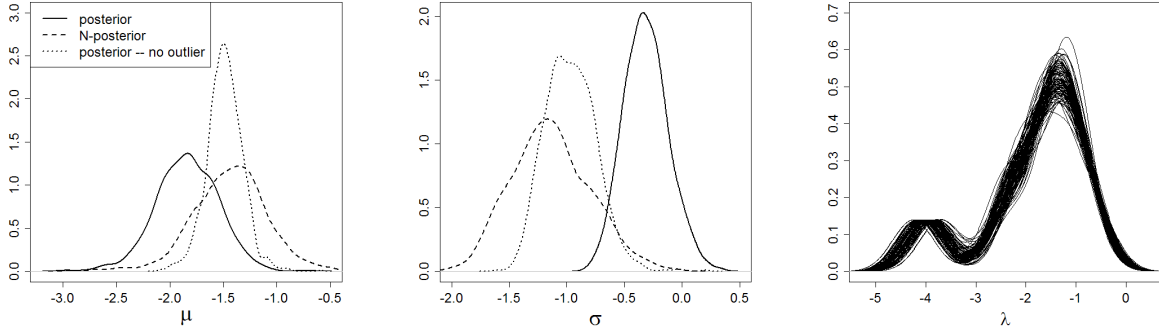


Figure 2: Posterior distributions for the parasite data. Left: posteriors for  $\mu$  with and without an outlier and the N-posterior. Middle: posteriors for  $\sigma$ . Right: samples of  $g_n$  based on draws from the posterior for  $p_1, \dots, p_n$ , demonstrating an outlier at -3.

deviations from income trend are evident.

This framework leads to a longitudinal data model. We begin with a random intercept model

$$Y_{ijk} = b_{0ij} + b_{1j}t_k + \epsilon_{ijk} \quad (21)$$

where  $Y_{ijk}$  is log income for the  $i$ th student in group  $j$  (American ( $a$ ) or foreign ( $f$ )) at age  $t_k$ . We extend to this the distributional assumptions

$$b_{0ij} \sim N(\beta_{0j}, \tau_0^2), \quad \epsilon_{ijk} \sim N(0, \sigma^2)$$

leading to a complete data log likelihood given up to a constant by

$$l(Y, \beta, \sigma^2, \tau_0^2) = - \sum_{i=1}^n \sum_{j \in \{a, f\}} \sum_{k=1}^4 \frac{1}{2\sigma^2} (Y_{ijk} - b_{0ij} - \beta_{1j}t_k)^2 - \sum_{i=1}^n \sum_{j \in \{a, f\}} \frac{1}{2\tau_0^2} (b_{0ij} - \beta_{0j})^2 \quad (22)$$

to which we attach Gaussian priors centered at zero with standard deviations 150 and 0.5 for the  $\beta_{0j}$  and  $\beta_{1j}$  respectively and Gamma priors with shape parameter 3 and scale 0.5 and 0.05 for  $\tau_0^2$  and  $\sigma^2$ . These are chosen to correspond to the approximate orders of magnitude observed in the maximum likelihood estimates of the  $b_{0ij}$ ,  $\beta_{1j}$  and residuals.

As in Section 7 we can robustify this likelihood in two different ways: either against the distributional assumptions on the  $\epsilon_{ijk}$  or on the  $b_{0ij}$ . In the latter case we form the density estimate

$$g_n(b; \beta) = \frac{1}{2nc_n} \sum_{i=1}^n \sum_{j \in \{a, f\}} K\left(\frac{b - b_{0ij} + \beta_{0j}}{c_n}\right)$$



and replace the second term in (22) with  $-2nD(g_n(\cdot; \boldsymbol{\beta}), \phi_{0, \tau_0^2}(\cdot))$ . Here we have used

$$\boldsymbol{\beta} = (\beta_{a0}, \beta_{f0}, \beta_{a1}, \beta_{f1}, b_{0a1}, b_{0f1}, \dots, b_{0an}, b_{0fn})$$

as an argument to  $g_n$  to indicate its dependence on the estimated parameters. We have chosen to combine the  $b_{0ai}$  and the  $b_{0fi}$  together in order to obtain the best estimate of  $g_n$ , rather than using a sum of disparities, one for American and one for foreign students.

To robustify the residual distribution, we observe that we cannot replace the first term with a single disparity based on the density of the combined  $\epsilon_{ijk}$  since the  $b_{0ij}$  cannot be identified marginally. Instead, we estimate a density at each  $ij$ :

$$g_{ij,n}^c(e; \boldsymbol{\beta}) = \frac{1}{4nc_n} \sum_{k=1}^4 K \left( \frac{e - (Y_{ijk} - b_{0ij} - \beta_{1j}t_k)}{c_n} \right)$$

and replace the first term with  $\sum_{i=1}^n \sum_{j \in \{a, f\}} 4D(g_{ij,n}^c(\cdot; \boldsymbol{\beta}), \phi_{0, \sigma^2}(\cdot))$ . This is the conditional form of the disparity. Note that this reduces us to four points for each density estimate; the limit of what could reasonably be employed. Naturally, both replacements can be made.

Throughout our analysis, we used Hellinger distance as a disparity; we also centered the  $t_k$ , resulting in  $b_{0ij}$  representing the expected salary of student  $ij$  at age 50. Bandwidths were fixed within a Metropolis sampling procedures. These were chosen by estimating the  $\hat{b}_{0ij}$  and  $\hat{\beta}_{1j}$  via least squares, and using these to estimate residuals and all other parameters:

$$\hat{\beta}_{0j} = \frac{1}{n} \hat{b}_{0j}, \quad e_{ijk} = Y_{ijk} - \hat{b}_{0ij} - \hat{\beta}_{1j}t_k, \quad \hat{\sigma}^2 = \frac{1}{8n-1} \sum_{ijk} e_{ijk}^2, \quad \hat{\tau}_0^2 = \frac{1}{2n-1} \sum_{ij} (\hat{b}_{0ij} - \hat{\beta}_{0j})^2.$$

The bandwidth selector in Sheather and Jones (1991) was applied to the  $\hat{b}_{0ij} - \hat{b}_{0j}$  to obtain a bandwidth for  $g_n(b; \boldsymbol{\beta})$ . The bandwidth for  $g_{ij,n}^c(e; \boldsymbol{\beta})$  was chosen as the average of the bandwidths selected for the  $e_{ijk}$  for each  $i$  and  $j$ . For each analysis, a Metropolis algorithm was run for 200,000 steps and every 100th sample was taken from the second half of the resulting Markov chain. The results of this analysis can be seen in Figure 3. Here we have plotted only the differences  $\beta_{f0} - \beta_{a0}$  and  $\beta_{f1} - \beta_{a1}$  along with the variance components. We observe that for posteriors that have not robustified the random effect distribution, there appears to be a significant difference in the rate of increase in income ( $P(\beta_{f1} < \beta_{a1}) < 0.02$  for both posterior and replacing the observation likelihood with Hellinger distance), however when the random effect likelihood is replaced with Hellinger

distance, the difference is no longer significant ( $P(\beta_{f1} < \beta_{a1}) > 0.145$  in both cases). We also observe that the estimated observation variance for the model is significantly reduced for posteriors in which the observation likelihood is replaced by Hellinger distance, but that uncertainty in the difference  $\beta_{f0} - \beta_{a0}$  is increased.

Investigating these differences, there were two foreign students who's over-all expected rate of increase is negative and separated from the least-squares slopes for all the other students. Removing these students increased the posterior probability of  $\beta_{a1} > \beta_{f1}$  to 0.11 and decreased the estimate of  $\sigma$  from 0.4 to 0.3. Removing the evident high outlier with a considerable departure from trend at age 45 in Figure 4 in Online Appendix E further reduced the EAP of  $\sigma$  to 0.185, in the same range as those obtained from robustifying the observation distribution.

A further model exploration allows a random slope for each student in addition to the random offset. The model now becomes

$$Y_{ijk} = b_{0ij} + b_{1ij}t_k + \epsilon_{ijk} \quad (23)$$

with additional distributional assumptions

$$b_{1ij} \sim N(\beta_{1j}, \tau_1^2)$$

and an additional term

$$-\frac{1}{2\tau_1^2} \sum_{i=1}^n \sum_{j \in \{a,f\}} (b_{1ij} - \beta_{1j})^2$$

added to (22). Here, this term can be robustified in a similar manner to the robustification of the  $b_{0ij}$ . However, we note that a robustification of the error terms would require the estimation of a conditional density for each  $ij$  – based on only four data points. We viewed this as being too little to achieve reasonable results and therefore employed the marginal formulation described in Section 6. Specifically, we first obtained residuals  $e_{ijk}$  for the random slope model from the maximum likelihood estimates for each subject-specific effect and estimated

$$\hat{\sigma}^2 = \frac{1}{0.674\sqrt{2}} |e_{ijk} - \text{median}(e)|.$$

Following this, we estimated a combined density for all residuals, conditional on the random effects

$$g_n^m(e; \beta) = \frac{1}{8nc_n} \sum_{i=1}^n \sum_{j \in \{a,f\}} \sum_{k=1}^4 K \left( \frac{e - (Y_{ijk} - b_{0ij} - b_{1ij}t_k)}{c_n} \right)$$

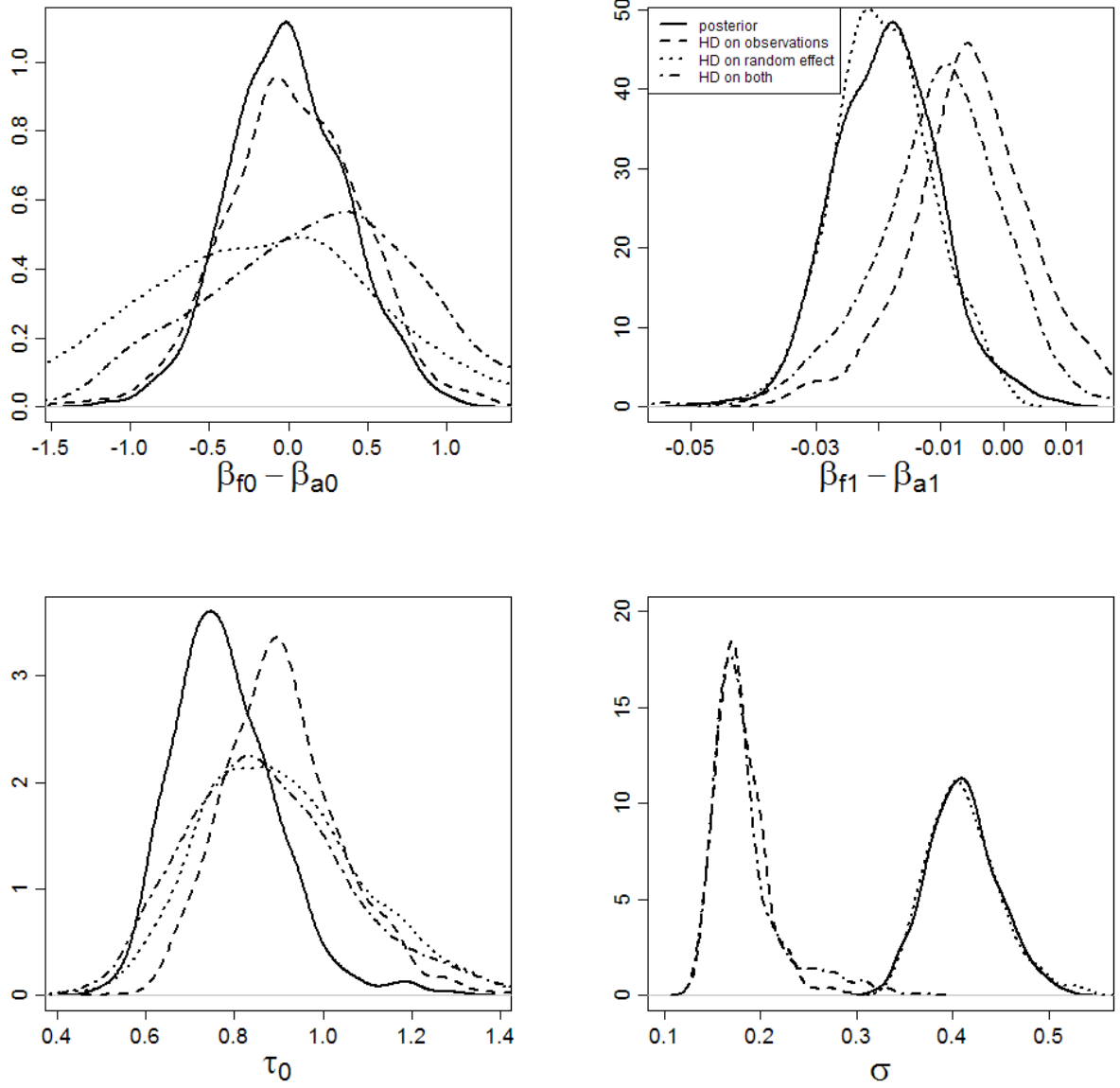


Figure 3: Analysis of the class survey data using a random intercept model with Hellinger distance replacing the observation likelihood, the random effect likelihood or both. Top: differences in intercepts between foreign and American students (left) and differences in slopes (right). Bottom: random effect variance (left) and observation variance (right). Models robustifying the random effect distribution do not show a significant difference in the slope parameters. Those robustifying the observation distribution estimate a significantly smaller observation variance.

and replaced the first term in (22) with  $-8nD(g_n^m(\cdot; \boldsymbol{\beta}), \phi_{0,\sigma^2}(\cdot))$ . Following the estimation of all other parameters, we obtained new residuals  $\tilde{e}_{ijk} = Y_{ijk} - \tilde{b}_{0ij} - b_{1ij}t_k$  where the  $\tilde{b}_{0ij}$  and  $\tilde{b}_{1ij}$  are the EDAP estimators. We then re estimated  $\sigma^2$  based on its H-posterior using the  $\tilde{e}_{ijk}$  as data. In this particular case a large number of outliers from a concentrated peak (see Figure 4) meant that the use of Gauss-Hermite quadrature in the evaluation of

$$HD(g_n^m(\cdot, \tilde{\boldsymbol{\beta}}), \phi_{0,\sigma^2}) = 2 - 2 \int \left( \sqrt{g_n^m(e; \tilde{\boldsymbol{\beta}})} / \sqrt{\phi_{0,\sigma^2}(e)} \right) \phi_{0,\sigma^2}(e) de$$

suffered from large numerical errors and we therefore employed a Monte Carlo integral based on 400 data points drawn from  $g_n^m$  instead, using the estimate (5). To estimate both  $\sigma^2$  and the other parameters we used a Metropolis random walk algorithm which was again run for 200,000 iterations with estimates based on every 100th sample in the second half of the chain.

Some results from this analysis are displayed in Figure 4. The residual distribution of the  $\tilde{e}_{ijk}$  show a very strong peak and a number of isolated outliers. The estimated standard deviation of the residual distribution is therefore very different between those methods that are robust to outliers and those that are not; the mean posterior  $\sigma$  was increased by a factor of four between those methods using a Hellinger disparity and those using the random effect log likelihood. The random slope variance was estimated to be small by all methods – we speculate that the distinction between random effect log likelihoods and Hellinger methods is bias due to bandwidth size – but this was not enough to overcome the differences between the methods concerning the distinction between  $\beta_{f1}$  and  $\beta_{a1}$ .

## 9. CONCLUSIONS

This paper combines disparity methods with Bayesian analysis to provide robust and efficient inference across a broad spectrum of models. In particular, these methods allow the robustification of any portion of a model for which the likelihood may be written as a product of distributions for i.i.d. random variables. This can be done without the need to modify either the assumed data-generating distribution or the prior. In our experience, Metropolis algorithms developed for the parametric model can be used directly to evaluate the D-posterior and generally incur a modest increase in the acceptance rate and computational cost. Our use of Metropolis algorithms in this context is *deliberately naive* in order to demonstrate the immediate applicability of our

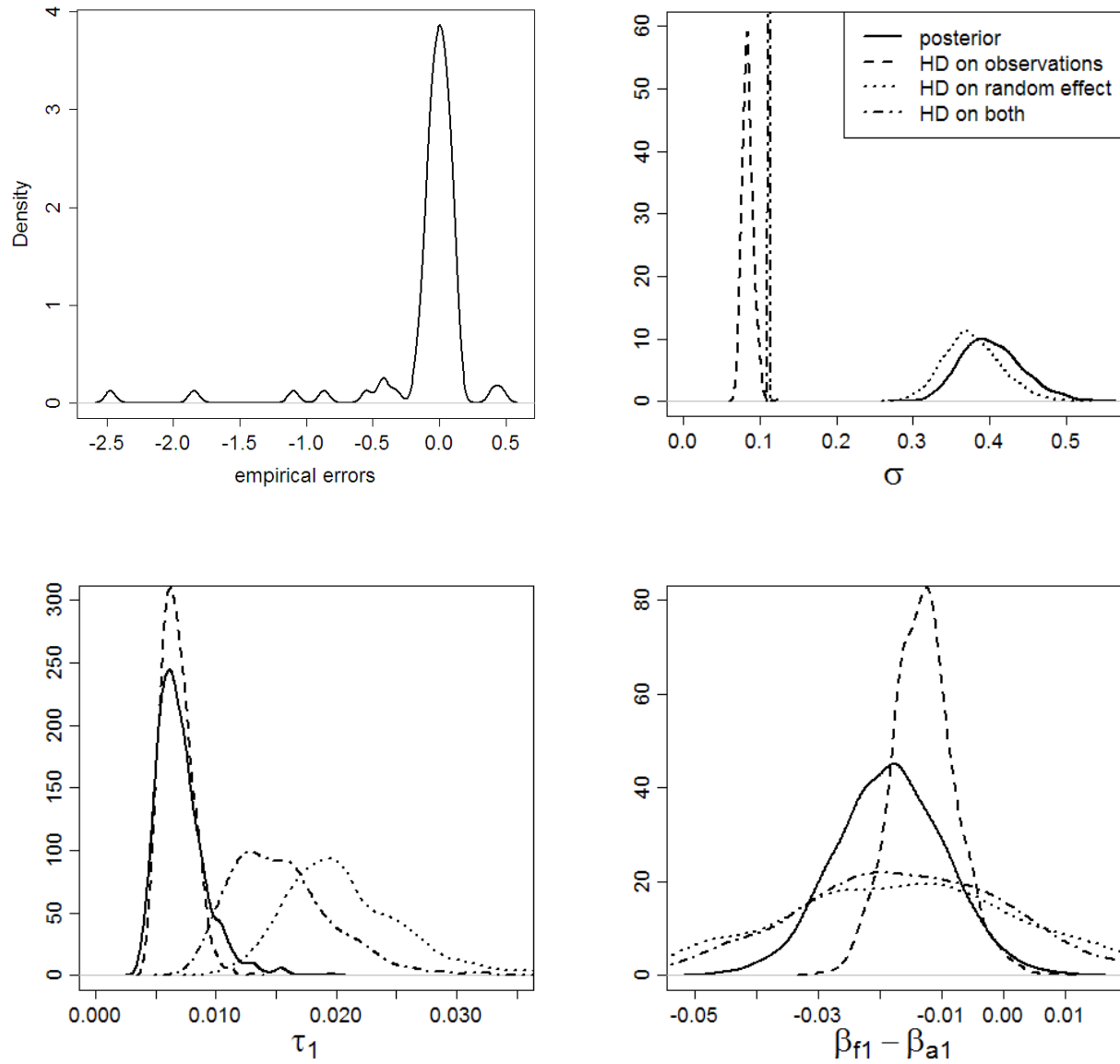


Figure 4: Analysis of a random-slope random-intercept model for the class survey data. Top left: a density estimate of the errors following a two-step procedure with the error variance held constant. This shows numerous isolated outliers than create an ill-conditioned problem for Gauss-Hermite quadrature methods. Top right: estimates of residual standard deviation replacing various terms in the log likelihood with Hellinger distance. The effect of outliers is clearly apparent in producing an over-estimate of variance. Bottom left: estimated variance of the random slope. Bottom right: estimated difference in mean slope between American and foreign students.

methods in combination with existing computational tools. We expect that a more careful study of sampling properties of these methods will yield considerable improvements in both computational and sampling efficiency.

The methods in this paper can be employed as a tool for model diagnostics; differences in results by an application of posterior and D-posterior can indicate problematic components of a hierarchical model. Further, estimated densities can indicate how the current model may be improved. However, the D-posterior can also be used directly to provide robust inference in an automated form.

Our mathematical results are given solely for i.i.d. data; ideas from Hooker and Vidyashankar (2011a) can be used to extend these to the regression framework. Our proposal of hierarchical models remains under mathematical investigation, but we expect that similar results can be established in this case. The methodology can also be applied within a frequentist context to define an alternative marginal likelihood for random effects models, although the numerical estimation of such models is likely to be problematic. Within this context, the choice of bandwidth  $c_n$  can become difficult. We have employed initial least-squares estimates above, but robust estimators could also be used instead. Empirically, we have found our results to be relatively insensitive to the choice of bandwidth.

An opportunity for further development of the proposed methodology lies in removing the boundedness of many disparities in common use. These yield EDAP estimates with breakdown points of 1, indicating hyper-insensitivity to the data. Theoretically, some form of boundedness has been used within proofs of the efficiency of minimum disparity estimators. However these results suggest an investigation of the necessity of this assumption and the development of new disparities which diverge at a rate slow enough to retain robustness.

The use of a kernel density estimate may also be regarded as inconsistent with a Bayesian context and it may therefore be desirable to employ non-parametric Bayesian density estimates as an alternative. Results for disparity estimation are heavily dependent on properties of kernel density estimates and this extension will require significant mathematical development; an initial study of the use of Dirichlet-process priors for density estimates in this context can be found in ?.

There is considerable scope to extend these methods to further problems. Robustification of the

innovation distribution in time-series models, for example, can be readily carried through through disparities and the hierarchical approach will extend this to either the observation or the innovation process in state-space models. The extension to continuous-time models such as stochastic differential equations, however, remains an open and interesting problem. More challenging questions arise in spatial statistics in which dependence decays over some domain and where a collection of i.i.d. random variables may not be available. There are also open questions in the application of these techniques to non-parametric smoothing, and in functional data analysis.

Supplementary Materials

**Background, Simulation and Data:** Appendices C - E detail simulation studies and data.

**R functions for Bayesian Disparity:** provides code to reproduce simulations and data analysis in this article (GNU zipped tar file).

## REFERENCES

- Albert, J. (2009). *Bayesian Computation with R*. New York: Springer.
- Andrade, J. A. A. and A. O’Hagan (2006). Bayesian robustness modeling using regularly varying distributions. *Bayesian Analysis* 1(1), 169–188.
- Basu, A., S. Sarkar, and A. N. Vidyashankar (1997). Minimum negative exponential disparity estimation in parametric models. *Journal of Statistical Planning and Inference* 58, 349–370.
- Beran, R. (1977). Minimum Hellinger distance estimates for parametric models. *Annals of Statistics* 5, 445–463.
- Berger, J. O. (1994). An overview of robust Bayesian analysis. *TEST* 3, 5–124.
- Cheng, A.-L. and A. N. Vidyashankar (2006). Minimum Hellinger distance estimation for randomized play the winner design. *Journal of Statistical Planning and Inference* 136, 1875–1910.
- Devroye, L. and G. Györfi (1985). *Nonparametric Density Estimation: The L1 View*. New York: Wiley.

- Dey, D. K. and L. R. Birmiwai (1994). Robust Bayesian analysis using divergence measures. *Statistics and Probability Letters* 20, 287–294.
- Engel, J., E. Herrmann, and T. Gasser (1994). An iterative bandwidth selector for kernel estimation of densities and their derivatives. *Journal of Nonparametric Statistics* 4, 2134.
- Ghosh, J. K., M. Delampady, and T. Samanta (2006). *An Introduction to Bayesian Analysis*. New York: Springer.
- Hampel, F. R. (1974). The influence curve and its role in robust estimation. *Journal of the american Statistics Association* 69, 383–393.
- Hampel, F. R., E. M. Ronchetti, P. J. Rousseeuw, and W. A. Stahel (1986). *Robust statistics*. Wiley Series in Probability and Mathematical Statistics: Probability and Mathematical Statistics. New York: John Wiley & Sons Inc. The approach based on influence functions.
- Hansen, B. E. (2004). Nonparametric conditional density estimation.
- Hooker, G. and A. N. Vidyashankar (2011a). Minimum disparity methods for nonlinear regression – marginal approach. *in preparation*.
- Hooker, G. and A. N. Vidyashankar (2011b). Minimum disparity methods for nonlinear regression – conditional approach. *in preparation*.
- Huber, P. (1981). *Robust Statistics*. New York: Wiley.
- Jurečková, J. and P. K. Sen (1996). *Robust statistical procedures*. Wiley Series in Probability and Statistics: Applied Probability and Statistics. New York: John Wiley & Sons Inc. Asymptotics and interrelations, A Wiley-Interscience Publication.
- Lindsay, B. G. (1994). Efficiency versus robustness: The case for minimum Hellinger distance and related methods. *Annals of Statistics* 22, 1081–1114.
- Maronna, R. A., R. D. Martin, and V. J. Yohai (2006). *Robust statistics*. Wiley Series in Probability and Statistics. Chichester: John Wiley & Sons Ltd. Theory and methods.



- Nielsen, M., Vidyashankar, A.N., B. Hanlon, S. Petersen, and R. Kaplan (2010). Hierarchical models for evaluating anthelmintic resistance in livestock parasites using observational data from multiple farms. *under review*.
- Pak, R. J. and A. Basu (1998). Minimum disparity estimation in linear regression models: Distribution and efficiency. *Annals of the Institute of Statistical Mathematics* 50, 503–521.
- Park, C. and A. Basu (2004). Minimum disparity estimation: Asymptotic normality and breakdown point results. *Bulletin of Informatics and Cybernetics* 36.
- Peng, F. and D. K. Dey (1995). Bayesian analysis of outlier problems using divergence measures. *Canadian Journal of Statistics* 23, 199–213.
- Sheather, S. J. and M. C. Jones (1991). A reliable data-based bandwidth selection method for kernel density estimation. *Journal of the Royal Statistical Society, Series B* 53, 683–690.
- Silverman, B. W. (1982). *Density Estimation*. Chapman and Hall.
- Simpson, D. G. (1987). Minimum Hellinger distance estimation for the analysis of count data. *Journal of the American statistical Association* 82, 802–807.
- Simpson, D. G. (1989). Hellinger deviance test: efficiency, breakdown points and examples. *Journal of the American Statistical Association* 84, 107–113.
- Sollich, P. (2002). Bayesian methods for support vector machines: Evidence and predictive class probabilities. *Machine Learning* 46, 21–52.
- Stigler, S. M. (1973). The asymptotic distribution of the trimmed mean. *Annals of Statistics* 1, 427–477.
- Szpiro, A.A., R. K. and T. Lumley (2011). Model-robust regression and a bayesian ‘sandwich’ estimator. *Annals of Applied Statistics*, To Appear.
- Tamura, R. N. and D. D. Boos (1986). Minimum Hellinger distances estimation for multivariate location and and covariance. *Journal of the American Statistical Association* 81, 223–229.

Zhan, X. and T. P. Hettmansperger (2007). Bayesian  $R$ -estimates in two-sample location models. *Comput. Statist. Data Anal.* 51(10), 5077–5089.

## A. PROOFS OF EFFICIENCY

### A.1 Proof of Theorem 1

We begin with the following Lemma:

**Lemma 1.** *Let*

$$w_n(t) = \pi(\hat{\theta}_n + t/\sqrt{n})e^{-nD(g_n, f_{\hat{\theta}_n + t/\sqrt{n}}) + nD(g_n, f_{\hat{\theta}_n})} - \pi(\theta_g)e^{-\frac{1}{2}t'I^D(\theta_g)t}$$

then under (A2)-(A6)

$$\int |w_n(t)|dt \xrightarrow{a.s.} 0 \text{ and } \int \|t\|_2 |w_n(t)|dt \xrightarrow{a.s.} 0.$$

*Proof.* We divide the integral into  $A_1 = \{\|t\|_2 > \delta\sqrt{n}\}$  and  $A_2 = \{\|t\|_2 \leq \delta\sqrt{n}\}$ :

$$\int |w_n(t)|dt = \int_{A_1} |w_n(t)|dt + \int_{A_2} |w_n(t)|dt \tag{A.1}$$

and show that each vanishes in turn. First, since  $\sup_{\theta \in \Theta} |D(g_n, f_\theta) - D(g, f_\theta)| \xrightarrow{a.s.} 0$  by Assumption (A4), for some  $\epsilon > 0$  with probability 1 it follows that

$$\exists N : \forall n \geq N, \sup_{\|t\|_2 > \delta} D(g_n, f_{\hat{\theta}_n + t/\sqrt{n}}) - D(g_n, f_{\hat{\theta}_n}) > -\epsilon.$$

This now allows us to demonstrate the convergence of the first term in (A.1):

$$\begin{aligned} \int_{A_1} |w_n(t)|dt &\leq \int_{A_1} \pi(\hat{\theta}_n + t/\sqrt{n})e^{-nD(g_n, f_{\hat{\theta}_n + t/\sqrt{n}}) + nD(g_n, f_{\hat{\theta}_n})} dt \\ &\quad + \int_{A_1} \pi(\theta_g)e^{-\frac{1}{2}t'I^D(\theta_g)t} dt \\ &\leq e^{-n\epsilon} + \pi(\theta_g) \left( \frac{|I^D(\theta_g)|}{2\pi} \right)^{p/2} P(\|Z\|_2 > \sqrt{n}\delta) \end{aligned} \tag{A.2}$$

where  $Z$  is a  $N(0, I^D(\theta_g))$  random variable and (A.2) converges to zero almost surely.

We now deal with the second term in (A.1). Notice that

$$nD(g_n, f_{\hat{\theta}_n + t/\sqrt{n}}) - nD(g_n, f_{\hat{\theta}_n}) = \frac{1}{2}t'I_n^D(\theta'_n)$$

for  $\theta'_n = \hat{\theta}_n + \alpha t/\sqrt{n}$  with  $0 \leq \alpha \leq 1$  and therefore

$$w_n(t) = \pi(\hat{\theta}_n + t/\sqrt{n})e^{-\frac{1}{2}t'I_n^D(\theta'_n)} - \pi(\theta_g)e^{-\frac{1}{2}t'I^D(\theta_g)t} \rightarrow 0$$

for every  $t$ .

By Assumption (A3) we can choose  $\delta$  so that  $I^D(\theta) \succ 2M$  if  $\|\theta - \theta_g\|_2 \leq 2\delta$  for some positive definite matrix  $M$  where  $A \succ B$  indicates  $t'At > t'Bt$  for all  $t$ . Since  $\|\theta'_n - \hat{\theta}_n\| \leq \delta$  with probability 1 for all  $n$  sufficiently large  $\exp\left(-nD\left(g_n, f_{\hat{\theta}_n+t\sqrt{n}}\right) + nD\left(g_n, f_{\hat{\theta}_n}\right)\right) \leq \exp\left(-\frac{1}{2}t'Mt\right)$ . Therefore

$$\int_{A_2} |w_n(t)|dt \leq \int_{A_2} \pi(\hat{\theta}_n + t/\sqrt{n})e^{-\frac{1}{2}t'Mt} + \pi(\theta_g) \int_{A_2} e^{-\frac{1}{2}t'I^D(\theta_g)t}dt < \infty.$$

and the result follows from the pointwise convergence of  $w(t)$  and the dominated convergence theorem.

We can prove  $\int \|t\|_2 |w_n(t)|dt \xrightarrow{a.s.} 0$  in an analogous manner by observing that on  $A_1$

$$\begin{aligned} \int_{A_1} \|t\|_2 |w_n(t)|dt &\leq \int_{A_1} \|t\|_2 \pi(\hat{\theta}_n + t/\sqrt{n})e^{-nD(g_n, f_{\hat{\theta}_n+t\sqrt{n}}) + nD(g_n, f_{\hat{\theta}_n})} dt \\ &\quad + \int_{A_1} \pi(\theta_g) \|t\|_2 e^{-\frac{1}{2}t'I^D(\theta_g)t} dt \end{aligned}$$

and on  $A_2$ ,  $\|t\|_2 |w_n(t)| \xrightarrow{a.s.} 0$  and

$$\int_{A_2} \|t\|_2 |w_n(t)|dt \leq \int_{A_2} \|t\|_2 \pi(\hat{\theta}_n + t/\sqrt{n})e^{-\frac{1}{2}t'Mt} + \pi(\theta_g) \int_{A_2} \|t\|_2 e^{-\frac{1}{2}t'I^D(\theta_g)t}dt < \infty.$$

□

Following this lemma, we prove Theorem 1.

*Proof.* First, from Assumption (A6),  $\sqrt{n}(\hat{\theta}_n - \theta_g) \xrightarrow{d} N(0, I^D(\theta_g))$ , using that  $\int |g_n(t) - f_{\theta_g}(t)|dt \xrightarrow{a.s.} 0$ , the continuity of  $G$  and the compactness of  $\Theta$ , it follows that  $\sup_{\theta \in \Theta} |D(g_n, f_\theta) - D(g, f_\theta)| \xrightarrow{a.s.} 0$  and

$$D(g_n, f_{\hat{\theta}_n}) \xrightarrow{a.s.} D(g, f_{\theta_g}), \quad \nabla_\theta D(g_n, f_{\hat{\theta}_n}) \xrightarrow{a.s.} \nabla_\theta D(g, f_{\theta_g}), \quad \nabla_\theta^2 D(g_n, f_{\hat{\theta}_n}) \xrightarrow{a.s.} \nabla_\theta^2 D(g, f_{\theta_g})$$

Now, we write that  $\pi_n^{*D}(t) = \kappa_n^{-1} \pi(\hat{\theta}_n + t/\sqrt{n}) \exp -nD(g_n, f_{\hat{\theta}_n+t\sqrt{n}}) + nD(g_n, f_{\hat{\theta}_n})$  where  $\kappa_n$  is chosen so that  $\int \pi_n^{*D}(t)dt = 1$ . Let

$$w_n(t) = \pi(\hat{\theta}_n + t/\sqrt{n})e^{-nD(g_n, f_{\hat{\theta}_n+t\sqrt{n}}) + nD(g_n, f_{\hat{\theta}_n})} - \pi(\theta_g)e^{-\frac{1}{2}t'I^D(\theta_g)t}$$

from Lemma 1, we have  $\int |w_n(t)|dt \xrightarrow{a.s.} 0$  from which

$$\kappa_n = \int \pi(\hat{\theta}_n + t/\sqrt{n}) e^{-nD(g_n, f_{\hat{\theta}_n + t/\sqrt{n}}) + nD(g_n, f_{\hat{\theta}_n})} dt \xrightarrow{a.s.} \pi(\theta_g) \int e^{-\frac{1}{2}t' I^D(\theta_g)t} dt = \pi(\theta_g) \left( \frac{2\pi}{|I^D(\theta_g)|} \right)^{p/2}$$

and

$$\begin{aligned} \lim_{n \rightarrow \infty} \int \left| \pi_n^{*D}(t) - \left( \frac{I^D(\theta_g)}{2\pi} \right)^{p/2} e^{-\frac{1}{2}t' I^D(\theta_g)t} \right| dt \\ \leq \kappa_n^{-1} \int |w_n(t)|dt + \left( \frac{2\pi}{|I^D(\theta_g)|} \right)^{p/2} \left| \kappa_n^{-1} \pi(\theta_g) - \left( \frac{|I^D(\theta_g)|}{2\pi} \right)^{p/2} \right| \\ \xrightarrow{a.s.} 0. \end{aligned}$$

That the result holds for  $I^D(\theta_g)$  replaced with  $\hat{I}_n^D(\hat{\theta}_n)$  follows from the almost sure convergence of the latter to the former.  $\square$

## A.2 Proof of Theorem 2

*Proof.* Let  $t = (t_1, \dots, t_p)$ , from Theorem 1

$$\int t_i \pi^{*D}(t|x_1, \dots, x_n) \xrightarrow{a.s.} \left( \frac{2\pi}{|I^D(\theta_g)|} \right)^{p/2} \int t_i e^{-\frac{1}{2}t' I^D(\theta_g)t} dt = 0.$$

Since  $\theta_n^* = E(\hat{\theta}_n + t/\sqrt{n}|X_1, \dots, X_n)$  we have

$$\sqrt{n}(\theta_n^* - \hat{\theta}_n) \xrightarrow{a.s.} \left( \frac{2\pi}{|I^D(\theta_g)|} \right)^{p/2} \int t e^{-\frac{1}{2}t' I^D(\theta_g)t} dt = 0.$$

Since  $\sqrt{n}(\hat{\theta}_n - \theta_g) \xrightarrow{d} N(0, I^D(\theta_g))$ , it follows that  $\sqrt{n}(\theta_n^* - \theta_g) \xrightarrow{d} N(0, I^D(\theta_g))$ ; hence  $\theta_n^*$  is asymptotically normal, efficient as well as robust.  $\square$

## B. PROOFS OF ROBUSTNESS

### B.1 Proof of Theorem 3

*Proof.* Under the assumptions,  $\sup_{\theta, g} D(g, f_\theta) = R < \infty$  and  $\inf_{\theta, g} D(g, f_\theta) = r > -\infty$ . Let  $h_{k, \epsilon} = (1 - \epsilon)g + \epsilon t_k$ , then for all  $\theta$ ,  $e^{-nR} \leq e^{-nD(h_{k, \epsilon}, f_\theta)} < e^{-nr}$ ,  $\forall k \in 1, 2, \dots$ ,  $\forall \epsilon \in [0, 1]$  and therefore

$$e^{n(r-R)} E_{\pi(\theta)} \theta = \frac{\int \theta e^{-nR} \pi(\theta) d\theta}{\int e^{-nr} \pi(\theta) d\theta} \leq E_{P_D(\theta|h_{k, \epsilon})} \theta \leq \frac{\int \theta e^{-nr} \pi(\theta) d\theta}{\int e^{-nR} \pi(\theta) d\theta} = e^{n(R-r)} E_{\pi(\theta)} \theta.$$

$\square$

## B.2 Theorem 4

To study the influence function in a larger class of models for which  $G$  and  $G'$  are unbounded we provide the following Theorem.

**Theorem 4.** *Let  $D(g, f_\theta)$  be bounded and assume that*

$$e_0 = \sup_x \int \left| G' \left( \frac{g(x)}{f_\theta(x)} - 1 \right) \pi(\theta) \right| d\theta < \infty \text{ and } e_1 = \sup_x \int \left| \theta G' \left( \frac{g(x)}{f_\theta(x)} - 1 \right) \pi(\theta) \right| d\theta < \infty \quad (\text{A.3})$$

*then  $|IF(\theta; g, t_k)| < \infty$ .*

In the case of Hellinger distance the conditions of Theorem 4 require the boundedness of  $r(x) = \int (\sqrt{f_\theta(x)}/\sqrt{g(x)})\pi(\theta)d\theta$ .

*Proof.* It is sufficient to show that  $|E_{P_D(\theta|g)} C_{nk}(\theta, g)| < \infty$  and  $|E_{P_D(\theta|g)} [\theta C_{nk}(\theta, g)]| < \infty$ . We will prove the first of these, the second follows analogously.

$$\begin{aligned} |E_{P_D(\theta|g)} C_{nk}(\theta, g)| &\leq e^{n(R-r)} \left| \int C_{nk}(\theta, g) \pi(\theta) d\theta \right| \\ &\leq e^{n(R-r)} \int \left| (g(x) - t_k(x)) \int G' \left( \frac{g(x)}{f_\theta(x)} - 1 \right) \pi(\theta) d\theta \right| dx \\ &\leq e^{n(R-r)} e_0 \int |g(x) - t_k(x)| dx \end{aligned} \quad (\text{A.4})$$

where  $\sup_{\theta, g} D(g, f_\theta) = R < \infty$  and  $\inf_{\theta, g} D(g, f_\theta) = r > -\infty$  and (A.4) follows from the assumption (A.3) and the bound  $\int |g(x) - t_k(x)| dx \leq 2$ .  $\square$

Since  $t_k(x)$  can be made to concentrate on regions where  $r(x)$  is large, we conjecture that the conditions in Theorem 4 are necessary. In fact, this requirement means that the H-posterior influence function will not be bounded for a large collection of parametric families.

## B.3 Convergence of Estimators

**Theorem 5.** *Assume that  $G$  has four continuous derivatives and that  $f_\theta$  is four times continuously differentiable. Define  $T_n(h)$  as in (13) and*

$$\hat{\theta}(h) = \arg \min_{\theta \in \Theta} D(h, f_\theta)$$

*then  $T_n(h) - \hat{\theta}(h) = o_p(n^{-1})$ .*

Before giving the proof we remark that it follows the lines of asymptotic expansions for posterior distributions as outlined in, for example, Ghosh et al. (2006). While we have provided explicit expressions only for the first term of the expansion, further terms can be given analytically.

*Proof.* We begin by taking a Taylor expansion of the log prior

$$\begin{aligned}\pi\left(\hat{\theta}(h) + t/\sqrt{n}\right) &= \pi\left(\hat{\theta}(h)\right) \left[1 + n^{-1/2}t' \frac{\nabla_{\theta}\pi(\hat{\theta}(h))}{\pi(\hat{\theta}(h))} + \frac{1}{2}n^{-1}t' \frac{\nabla_{\theta}^2\pi(\hat{\theta}(h))}{\pi(\hat{\theta}(h))}t\right] + o_p\left(n^{-1}\right) \\ &= \pi\left(\hat{\theta}(h)\right) \left[1 + n^{-1/2}t'b_1 + \frac{1}{2}n^{-1}t'b_2t\right] + o_p\left(n^{-1}\right)\end{aligned}$$

and the corresponding expansion of the disparity

$$\begin{aligned}nD\left(h, f_{\hat{\theta}(h)+t/\sqrt{n}}\right) - nD\left(h, f_{\hat{\theta}(h)}\right) \\ = \frac{1}{2}t'I^D(\hat{\theta}(h))t + \frac{n^{-3/2}}{6} \sum_{i,j,k} t_i t_j t_k a_{3,ijk} + \frac{n^{-1}}{24} \sum_{ijkl} t_i t_j t_k t_l a_{4,ijkl} + o_p\left(n^{-1}\right)\end{aligned}$$

yielding

$$\pi\left(\hat{\theta}(h) + t/\sqrt{n}\right) e^{-nD\left(h, f_{\hat{\theta}(h)+t/\sqrt{n}}\right) + nD\left(h, f_{\hat{\theta}(h)}\right)} = \pi\left(\hat{\theta}(h)\right) e^{-t'I^D(\hat{\theta}(h))t/2} \left[1 + \frac{c_1(t)}{n^{1/2}} + \frac{c_2(t)}{n}\right] + o_p(n^{-1})$$

for  $c_1(t) = \sum_i t_i b_{1,i} + \frac{1}{6} \sum_{ijk} t_i t_j t_k a_{3,ijk}$  and

$$c_2(t) = \sum_{ij} \frac{b_{2,ij}}{2} t_i t_j + \sum_{ijkl} \left( \frac{a_{3,ijk} b_{1,l}}{6} + \frac{a_{4,ijkl}}{24} \right) t^4 + \sum_{ijk} \frac{a_{3,ijk}^2}{72} t_i^2 t_j^2 t_k^2$$

so that in particular we have for the  $i$ th component of the EDAP vector

$$\begin{aligned}T_n(h)_i - \hat{\theta}(h)_i &= \frac{\int \left(\hat{\theta}(h) + t_i/\sqrt{n}\right) e^{-nD\left(h, f_{\hat{\theta}(h)+t/\sqrt{n}}\right) + nD\left(h, f_{\hat{\theta}(h)}\right)} \pi\left(\hat{\theta}(h) + t/\sqrt{n}\right) dt}{\int e^{-nD\left(h, f_{\hat{\theta}(h)+t/\sqrt{n}}\right) + nD\left(h, f_{\hat{\theta}(h)}\right)} \pi\left(\hat{\theta}(h) + t/\sqrt{n}\right) dt} \\ &= \frac{\left(\frac{|I^D(\hat{\theta}(h))|}{2\pi}\right)^{p/2} \left[ \hat{\theta}(h)_i + n^{-1} \left[ I^D(h)^{-1} \right]_{ii} \left( \sum_j \frac{a_{3,jji} [I^D(h)^{-1}]_{jj}}{2} + \frac{\nabla_{\theta}\pi(\hat{\theta}(h))_i}{\pi(\hat{\theta}(h))} \right) \right]}{\left(\frac{|I^D(\hat{\theta}(h))|}{2\pi}\right)^{p/2} + o_p\left(n^{-1}\right)} \\ &= \hat{\theta}(h)_i + n^{-1} \left[ I^D(h)^{-1} \right]_{ii} \left( \sum_j \frac{a_{3,jji} [I^D(h)^{-1}]_{jj}}{2} + \frac{\nabla_{\theta}\pi(\hat{\theta}(h))_i}{\pi(\hat{\theta}(h))} \right) + o_p\left(n^{-1}\right)\end{aligned}$$

□