

Invited Discussion*

Steven N. MacEachern[†] and Juhee Lee[‡]

First, we would like to congratulate the authors for their development of a fast and efficient method of assessing sensitivity to the prior specification of the Dirichlet process mixture (DPM) and related models. Their techniques are widely applicable and will see much use. Their examples give us a taste of what can be done with their method in models that move beyond the DPM. We greatly enjoyed reading the paper.

The past 30-plus years have seen incredible growth in the use of Bayesian methods for data analysis. The initial growth was driven by the development of Markov chain Monte Carlo (MCMC) methods that allowed one to fit models of substantial complexity. This was quickly followed by a realization in the entire statistics community that Bayesian methods simply work better than classical methods in “high information” settings with clean data – settings where one can reliably write down a complete Bayesian model and can clearly specify the inference problem. In practice, this translates to settings where different analysts (or the same analyst on different days) would select similar sampling densities for the data and similar prior distributions for the parameters, including latent structures, to arrive at similar models, and would also choose similar loss functions to formalize the inference problem. It also relies on realism in the model, with choices based on an understanding of the phenomenon under study rather than computational convenience. In these settings, the mathematical theory that links optimal inference to Bayesian methods is borne out. Bayesian methods simply work better than methods that are distant from Bayes.

The success of Bayesian methods is not uniform. In “low information” settings, specification of the sampling density and form of the model are every bit as challenging for the Bayesian as for the classical statistician. For high or infinite dimensional models, specification of the prior distribution remains a challenge. Data of dubious quality have the potential to dramatically impact the final inference. MCMC methods are often slow and sometimes exhibit poor convergence. And analysts commonly adjust their model to make fitting it easier and quicker. These difficulties are not shortcomings of the analyst, but rather features of the low information-complex model setting. They set an agenda for research on Bayesian data analysis.

The centerpiece of this agenda is how to improve Bayesian data analysis. Here, we see three main threads. One focuses on diagnostics through the development of techniques to identify deficiencies in the model (whether sampling density, prior distribution or a combination of the two), to identify cases that do not accord with the model (as in outliers), and to identify cases that have a large impact on inference (influential cases).

*Supported by the NSF under grant numbers SES-1921523, DMS-2015428, and DMS-2015552.

[†]Department of Statistics, The Ohio State University, ORCID 0000-0003-4106-1232, snm@stat.osu.edu

[‡]Department of Statistics, University of California, Santa Cruz, ORCID 0000-0002-9787-3830, juheele@soe.ucsc.edu

The second focuses on computational implementation. The third focuses on strategies to improve model specification and inference, with particular attention to robustness in the low information setting.

Giordano et al.'s delightful paper focuses on the first thread and is informed by the second. The authors' insight into the (in)effective use of Bayesian methods shines sharply through their paper. They consider a high/infinite dimensional setting where the prior distribution is specified through a rule-based strategy and where it would be difficult to place full confidence in any chosen rule. In this same setting, variational methods for fitting the model are far quicker than MCMC methods. Variational Bayes (VB) methods also generate the derivatives needed to examine the local sensitivity of features of the posterior distribution to changes in the prior distribution.

The developments in Giordano et al.'s paper parallel the development of local influence as a diagnostic method in classical statistics. Cook (1977)'s initial development of case influence examined the impact of individual cases on inference in the linear model. The main technique was case deletion. It led to Cook's distance, now a standard diagnostic summary in regression. Cook (1986) subsequently extended measures of influence to classical nonlinear models which, in the early-to-mid 1980's, were subject to the difficulties more recently experienced by MCMC methods. The models were slow to fit (via maximum likelihood) and one needed to be concerned with the numerical accuracy of the fits. With no clean analytical form, these difficulties rendered case deletion methods ineffective for the interactive data analysis that was being developed at the time. Instead, Cook turned to local influence, considering infinitesimal perturbations of case weights, and looked for big directional derivatives.

Cook's methods have been extended to the Bayesian setting, initially with pre-MCMC computation (e.g., Johnson and Geisser (1983)) and then with MCMC (e.g., Weiss (1996), Bradlow and Zaslavsky (1997), MacEachern and Peruggia (2000)). The approaches include both full case deletion and local influence (viz. Thomas et al. (2018)), and they cover a variety of inferences, from impact on the full posterior distribution to impact on marginal summaries of the posterior. While these methods are successful for linear models and low-dimensional nonlinear models the techniques are less effective for high-dimensional problems of the sort considered by Giordano et al.

Our first question for the authors is whether the techniques they develop can be adapted to assess local case influence. If so, is such an adaptation computationally feasible? The extension would provide the analyst with an additional tool to identify cases or sets of cases that have a large impact on inference.

Our second question concerns robust forms of the prior distribution. As described in the paper, DPM models are often used for clustering problems. Inferences on clustering, such as the number of clusters and co-clustering probabilities, are influenced by the prior specification. The parameter of the Dirichlet process (DP) may be split into two parts, the total mass parameter α and a base probability measure, $\mathcal{P}_{\text{base}}$. The distribution $\mathcal{P}_{\text{base}}$ generates cluster specific parameters, i.e., $\beta_k \stackrel{iid}{\sim} \mathcal{P}_{\text{base}}(\beta \mid \xi)$, where the β_k 's are cluster specific parameters and ξ is the hyperparameter vector for $\mathcal{P}_{\text{base}}$. Jointly

with α , $\mathcal{P}_{\text{base}}$ influences inference on the clustering structure. For example, Bush et al. (2010) and Lee et al. (2014) studied the joint impact of α and the dispersion of $\mathcal{P}_{\text{base}}$ on posterior inference. In particular, for fixed α and a given dataset of size n , a DPM model tends to produce fewer clusters with larger sizes under more dispersed $\mathcal{P}_{\text{base}}$. In response to this phenomenon, Lee et al. (2014) developed a local-mass preserving prior distribution for Bayesian nonparametric (BNP) models that produces more stable inference for clustering. The central idea is to define “local mass” as the mass assigned by a measure to a region \mathcal{L} of interest in the parameter space Ω_β prior to analysis and to jointly elicit $\mathcal{P}_{\text{base}}$ and α by holding the mass in \mathcal{L} constant. In addition to providing more stable inference about clustering, this form of prior distribution stabilizes inference for other quantities such as estimation of cluster locations.

It would be of great interest to see whether the authors’ method can be extended to perform a more comprehensive examination of model sensitivity, including sensitivity to the local mass \mathcal{L} . Does the authors’ quick and automated tool to assess sensitivity of inference to the specification of α extend to prior distributions with a local mass structure? If so, does one form of prior distribution show greater robustness than the other?

Our final comment returns to data analysis. The standard Bayesian analysis requires a rigorous determination of three components; (i) the sampling distribution (likelihood function) for the observations, (ii) the prior distribution of the parameters and (iii) loss function for making inference (decision). Bayesian analysis strongly depends on the choice of these components, and it is essential to investigate the sensitivity of the procedure to perturbation of all three components. The loss function has traditionally received less attention than the prior distribution and sampling density, as it often has little impact in a low-dimensional parametric setting. However, the choice of the inference (loss) function is critical for BNP models due to their flexibility. For example, DPMs are good at accommodating local features of the data such as outliers. These cases may be captured as one or more clusters that depart from the general pattern of the data. Thus, an inference function that discounts the impact of the outliers on the overall analysis can be more desirable than traditional inference functions (e.g., the quadratic loss function and the 0-1 loss function) for robust decision making (Lee and MacEachern, 2014). The inference function does not “wash out” as the sample size grows. This is similar to the parameter α not washing out for clustering in DPM models. We see scope for the development of inference functions that target a sensible summary of the posterior (or predictive) distribution and that lead to stable inference as the prior and sampling density are varied. Do the authors have a sense of whether a summary such as “number of clusters exceeding a given size” tends to be more robust than the simple “number of clusters”? Do the authors know of systematic ways to create more robust summaries of clusters?

In keeping with the level innovation in this work, the paper opens a host of questions. Those above seem, to us, to walk the tightrope of computational feasibility that leads from the questions we can answer with our written model to those we would answer with infinite resources. We close our discussion here, congratulating the authors on an interesting paper that develops a technique that will undoubtedly see heavy use.

References

- Bradlow, E. T. and Zaslavsky, A. M. (1997). Case influence analysis in Bayesian inference. *Journal of Computational and Graphical Statistics*, 6(3):314–331. 322
- Bush, C. A., Lee, J., and MacEachern, S. N. (2010). Minimally informative prior distributions for non-parametric Bayesian analysis. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(2):253–268. MR2830767. doi: <https://doi.org/10.1111/j.1467-9868.2009.00735.x>. 323
- Cook, R. D. (1977). Detection of influential observations in linear regression. *Technometrics*, 19(1):15–18. MR0436478. doi: <https://doi.org/10.2307/1268249>. 322
- Cook, R. D. (1986). Assessment of local influence (with discussion). *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 48(2):133–169. MR0867994. 322
- Johnson, W. and Geisser, S. (1983). A predictive view of the detection and characterization of influential observations in regression analysis. *Journal of the American Statistical Association*, 78(381):137–144. MR0696858. doi: <https://doi.org/10.1080/01621459.1983.10477942>. 322
- Lee, J. and MacEachern, S. N. (2014). Inference functions in high dimensional Bayesian inference. *Statistics and Its Interface*, 7(4):477–486. MR3302376. doi: <https://doi.org/10.4310/SII.2014.v7.n4.a5>. 323
- Lee, J., MacEachern, S. N., Lu, Y., and Mills, G. B. (2014). Local-mass preserving prior distributions for nonparametric Bayesian models. *Bayesian Analysis*, 9(2):307–330. MR3216998. doi: <https://doi.org/10.1214/13-BA857>. 323
- MacEachern, S. N. and Peruggia, M. (2000). Importance link function estimation for Markov chain Monte Carlo methods. *Journal of Computational and Graphical Statistics*, 9(1):99–121. MR1819867. doi: <https://doi.org/10.2307/1390615>. 322
- Thomas, Z. M., MacEachern, S. N., and Peruggia, M. (2018). Reconciling curvature and importance sampling based procedures for summarizing case influence in Bayesian models. *Journal of the American Statistical Association*, 113(524):1669–1683. MR3902237. doi: <https://doi.org/10.1080/01621459.2017.1360777>. 322
- Weiss, R. (1996). An approach to Bayesian sensitivity analysis. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 58(4):739–750. MR1410188. 322

robustness measures, but attentive to their ability to answer useful questions in their own modeling contexts.

References

- Averbukh, V. and Smolyanov, O. (1967). “The theory of differentiation in linear topological spaces.” *Russian Mathematical Surveys*, 22(6): 201–258. [MR0223886](#). 356
- Bae, J., Ng, N., Lo, A., Ghassemi, M., and Grosse, R. (2022). “If Influence Functions are the Answer, Then What is the Question?” In *Advances in Neural Information Processing Systems*. 356
- Basu, S. (2000). *Bayesian Robustness and Bayesian Nonparametrics*, 223–240. New York, NY: Springer New York. [MR1795218](#). doi: https://doi.org/10.1007/978-1-4612-1306-2_12. 363
- Basu, S., Jammalamadaka, S. R., and Liu, W. (1996). “Local posterior robustness with parametric priors: Maximum and average sensitivity.” In *Maximum Entropy and Bayesian Methods*, 97–106. Springer. 363
- Billingsley, P. (2008). *Probability and Measure*. John Wiley & Sons. [MR0534323](#). 356
- Carlin, B. and Polson, N. (1991). “An expected utility approach to influence diagnostics.” *Journal of the American Statistical Association*, 86(416): 1013–1021. 358, 363
- Cook, D. (1977). “Detection of influential observation in linear regression.” *Technometrics*, 19(1): 15–18. [MR0436478](#). doi: <https://doi.org/10.2307/1268249>. 357
- Cook, R. (1986). “Assessment of local influence.” *Journal of the Royal Statistical Society: Series B (Methodological)*, 48(2): 133–155. [MR0867994](#). 363
- Diaconis, P. and Freedman, D. (1986). “On the consistency of Bayes estimates.” *The Annals of Statistics*, 1–26. [MR0829555](#). doi: <https://doi.org/10.1214/aos/1176349830>. 356
- Fleming, W. (2012). *Functions of Several Variables*. Springer Science & Business Media. [MR0422527](#). 356
- Gil-Leyva, M. and Mena, R. (2021). “Stick-breaking processes with exchangeable length variables.” *Journal of the American Statistical Association*, 1–14. 357, 358, 359
- Giordano, R., Jordan, M. I., and Broderick, T. (2019a). “A higher-order swiss army infinitesimal jackknife.” *arXiv preprint arXiv:1907.12116*. 356
- Giordano, R., Stephenson, W., Liu, R., Jordan, M. I., and Broderick, T. (2019b). “A Swiss army infinitesimal jackknife.” In *The 22nd International Conference on Artificial Intelligence and Statistics*, 1139–1147. PMLR. 357, 362
- Gustafson, P. (1996). “Local sensitivity of posterior expectations.” *The Annals of Statistics*, 24(1): 174–195. [MR1389886](#). doi: <https://doi.org/10.1214/aos/1033066205>. 356, 363

- Hampel, F. (1974). “The influence curve and its role in robust estimation.” *Journal of the American Statistical Association*, 69(346): 383–393. [MR0362657](#). 356, 357
- Hampel, F., Ronchetti, E., Rousseeuw, P., and Stahel, W. (1986). *Robust Statistics: The Approach Based on Influence Functions*. Wiley-Interscience; New York. [MR0829458](#). 358
- Hattori, S. and Kato, M. (2009). “Approximate subject-deletion influence diagnostics for Inverse Probability of Censoring Weighted (IPCW) method.” *Statistics and Probability Letters*, 79(17): 1833–1838. [MR2749935](#). doi: <https://doi.org/10.1016/j.spl.2009.05.013>. 356
- Huber, P. J. (1964). “Robust estimation of a location parameter.” *The Annals of Mathematical Statistics*, 35(1): 73–101. URL <http://www.jstor.org/stable/2238020> [MR0161415](#). doi: <https://doi.org/10.1214/aoms/1177703732>. 357
- Jacobi, L., Joshi, M., and Zhu, D. (2018). “Automated sensitivity analysis for Bayesian inference via Markov chain Monte Carlo: Applications to Gibbs sampling.” Available at *SSRN 2984054*. 356
- Johnson, W. and Geisser, S. (1983). “A predictive view of the detection and characterization of influential observations in regression analysis.” *Journal of the American Statistical Association*, 78(381): 137–144. [MR0696858](#). 358, 363
- Koh, P. and Liang, P. (2017). “Understanding black-box predictions via influence functions.” In *International Conference on Machine Learning (ICML)*. 357
- Krantz, S. and Parks, H. (2012). *The Implicit Function Theorem: History, Theory, and Applications*. Springer Science & Business Media. [MR2977424](#). doi: <https://doi.org/10.1007/978-1-4614-5981-1>. 356
- Lee, J., James, L., and Choi, S. (2016). “Finite-dimensional BFRY priors and variational Bayesian inference for power law models.” *Advances in Neural Information Processing Systems*. 357
- Lee, J. and MacEachern, S. (2014). “Inference functions in high dimensional Bayesian inference.” *Statistics and its Interface*, 7(4): 477–486. [MR3302376](#). doi: <https://doi.org/10.4310/SII.2014.v7.n4.a5>. 361
- MacEachern, S. and Peruggia, M. (2000). “Importance link function estimation for Markov chain Monte Carlo methods.” *Journal of Computational and Graphical Statistics*, 9(1): 99–121. [MR1819867](#). doi: <https://doi.org/10.2307/1390615>. 363
- Maclaurin, D. (2016). “Modeling, Inference and Optimization With Composable Differentiable Procedures.” [MR3706076](#). 356
- McAuliffe, J., Blei, D., and Jordan, M. I. (2006). “Nonparametric empirical Bayes for the Dirichlet process mixture model.” *Statistics and Computing*, 16: 5–14. [MR2224185](#). doi: <https://doi.org/10.1007/s11222-006-5196-2>. 357, 360, 361
- McCulloch, R. (1989). “Local model influence.” *Journal of the American Statistical Association*, 84(406): 473–478. 358

- McInerney, A. (2013). *First Steps in Differential Geometry*. Springer. MR3098248. doi: <https://doi.org/10.1007/978-1-4614-7732-7>. 356
- Mohamed, S., Rosca, M., Figurnov, M., and Mnih, A. (2020). “Monte Carlo gradient estimation in machine learning.” *Journal of Machine Learning Research*, 21(132): 1–62. MR4138116. 356
- Murray, M. and Rice, J. (1993). *Differential Geometry and Statistics*, volume 48. CRC Press. MR1293124. doi: <https://doi.org/10.1007/978-1-4899-3306-5>. 356
- Ruggeri, F. and Wasserman, L. (1993). “Infinitesimal sensitivity of posterior distributions.” *Canadian Journal of Statistics*, 21(2): 195–203. MR1234761. doi: <https://doi.org/10.2307/3315811>. 356, 363
- Serfling, R. (1980). *Approximation Theorems of Mathematical Statistics*, volume 162. John Wiley & Sons. MR0595165. 357
- Shao, J. and Tu, D. (2012). *The Jackknife and Bootstrap*. Springer Science & Business Media. MR1351010. doi: <https://doi.org/10.1007/978-1-4612-0795-5>. 357
- Shi, L., Lu, J., Zhao, J., and Chen, G. (2016). “Case deletion diagnostics for GMM estimation.” *Computational Statistics & Data Analysis*, 95: 176–191. MR3425947. doi: <https://doi.org/10.1016/j.csda.2015.10.003>. 356
- Stigler, S. (2010). “The changing history of robustness.” *The American Statistician*, 64(4): 277–281. MR2758558. doi: <https://doi.org/10.1198/tast.2010.10159>. 357
- Thomas, W. and Cook, D. (1989). “Assessing influence on regression coefficients in generalized linear models.” *Biometrika*, 76(4): 741–749. MR1041419. doi: <https://doi.org/10.1093/biomet/76.4.741>. 356
- Thomas, Z., MacEachern, S., and Peruggia, M. (2018). “Reconciling curvature and importance sampling based procedures for summarizing case influence in Bayesian models.” *Journal of the American Statistical Association*, 113(524): 1669–1683. MR3902237. doi: <https://doi.org/10.1080/01621459.2017.1360777>. 358
- van der Vaart, A. and Wellner, J. (1996). *Empirical Processes and Weak Convergence*. Springer, New York. MR1385671. doi: <https://doi.org/10.1007/978-1-4757-2545-2>. 357
- von Mises, R. (1947). “On the asymptotic distribution of differentiable statistical functions.” *The Annals of Mathematical Statistics*, 18(3): 309–348. MR0022330. doi: <https://doi.org/10.1214/aoms/1177730385>. 357
- Walker, S. (2007). “Sampling the Dirichlet mixture model with slices.” *Communications in Statistics — Simulation and Computation*, 36(1): 45–54. MR2370888. doi: <https://doi.org/10.1080/03610910601096262>. 359
- Zeidler, E. (1986). *Nonlinear Functional Analysis and its Applications I: Fixed-point Theorems*. Springer-Verlag. MR0816732. doi: <https://doi.org/10.1007/978-1-4612-4838-5>. 356, 361