



Variational Inference: A Review for Statisticians

David M. Blei, Alp Kucukelbir & Jon D. McAuliffe

To cite this article: David M. Blei, Alp Kucukelbir & Jon D. McAuliffe (2017) Variational Inference: A Review for Statisticians, *Journal of the American Statistical Association*, 112:518, 859-877, DOI: [10.1080/01621459.2017.1285773](https://doi.org/10.1080/01621459.2017.1285773)

To link to this article: <https://doi.org/10.1080/01621459.2017.1285773>



[View supplementary material](#) 



Published online: 13 Jul 2017.



[Submit your article to this journal](#) 



Article views: 27190



[View related articles](#) 



[View Crossmark data](#) 



Citing articles: 645 [View citing articles](#) 

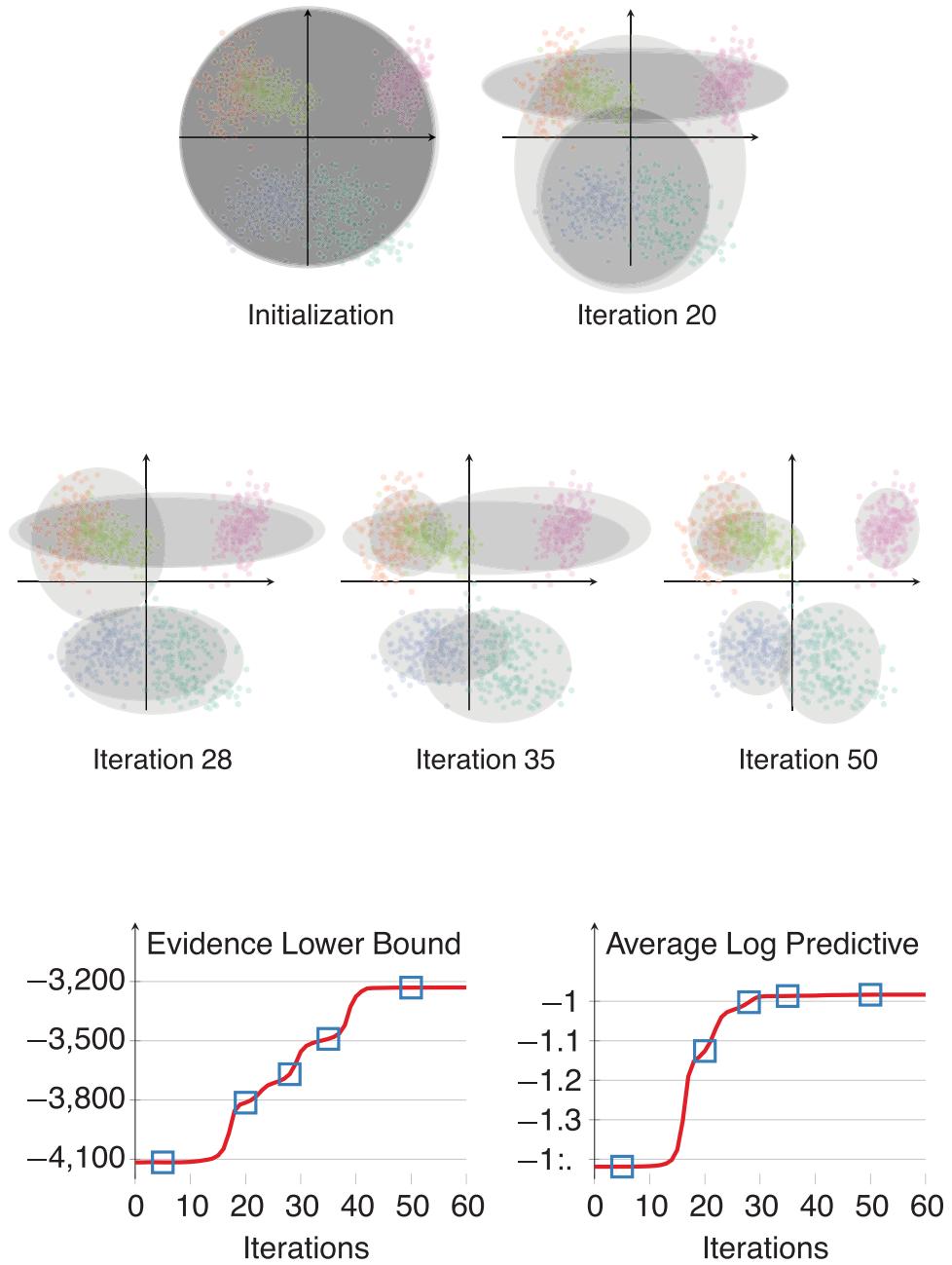


Figure 3. A simulation study of a two-dimensional Gaussian mixture model. The ellipses are 2σ contours of the variational approximating factors.

stochastic variational inference (Hoffman et al. 2013), a stochastic optimization algorithm that scales up variational inference in this setting.

4.1. Complete Conditionals in the Exponential Family

Consider the generic model $p(\mathbf{z}, \mathbf{x})$ of Section 2.1 and suppose each complete conditional is in the exponential family:

$$p(z_j | \mathbf{z}_{-j}, \mathbf{x}) = h(z_j) \exp\{\eta_j(\mathbf{z}_{-j}, \mathbf{x})^\top z_j - a(\eta_j(\mathbf{z}_{-j}, \mathbf{x}))\}, \quad (36)$$

where z_j is its own sufficient statistic, $h(\cdot)$ is a base measure, and $a(\cdot)$ is the log normalizer (Brown 1986). Because this is a

conditional density, the parameter $\eta_j(\mathbf{z}_{-j}, \mathbf{x})$ is a function of the conditioning set.

Consider mean-field variational inference for this class of models, where we fit $q(\mathbf{z}) = \prod_j q_j(z_j)$. The exponential family assumption simplifies the coordinate update of Equation (17),

$$q(z_j) \propto \exp\{\mathbb{E}[\log p(z_j | \mathbf{z}_{-j}, \mathbf{x})]\} \quad (37)$$

$$= \exp\{\log h(z_j) + \mathbb{E}[\eta_j(\mathbf{z}_{-j}, \mathbf{x})]^\top z_j - \mathbb{E}[a(\eta_j(\mathbf{z}_{-j}, \mathbf{x}))]\} \quad (38)$$

$$\propto h(z_j) \exp\{\mathbb{E}[\eta_j(\mathbf{z}_{-j}, \mathbf{x})]^\top z_j\}. \quad (39)$$

This update reveals the parametric form of the optimal variational factors. Each one is in the same exponential family as its corresponding complete conditional. Its parameter has the same

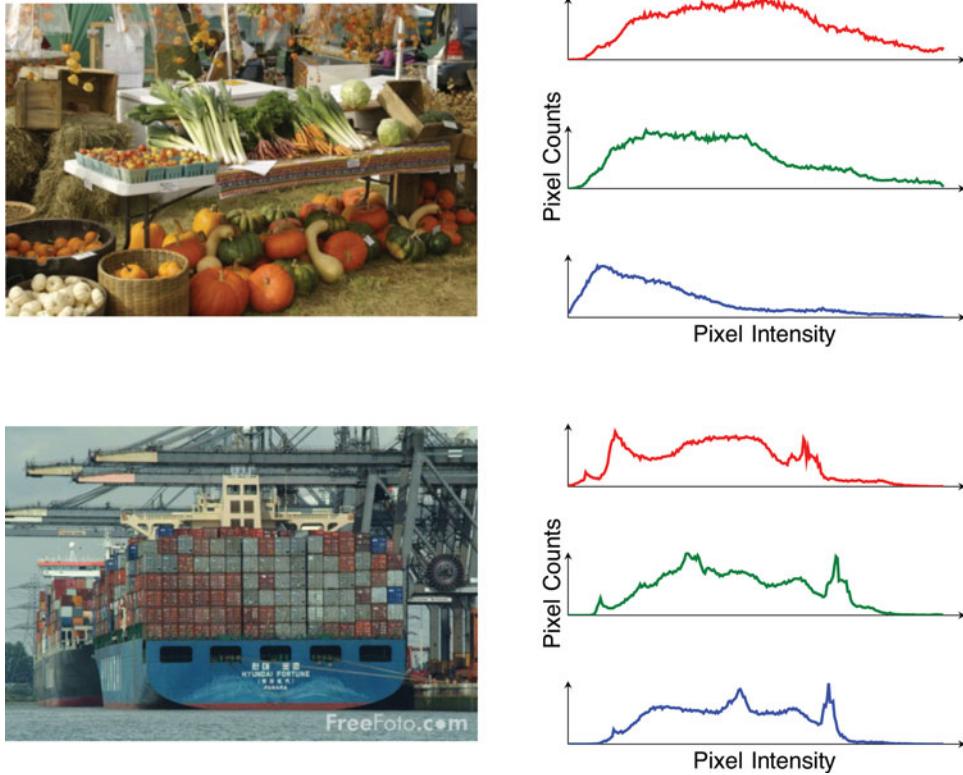


Figure 4. Red, green, and blue channel image histograms for two images from the imageCLEF dataset. The top image lacks blue hues, which is reflected in its blue channel histogram. The bottom image has a few dominant shades of blue and green, as seen in the peaks of its histogram.

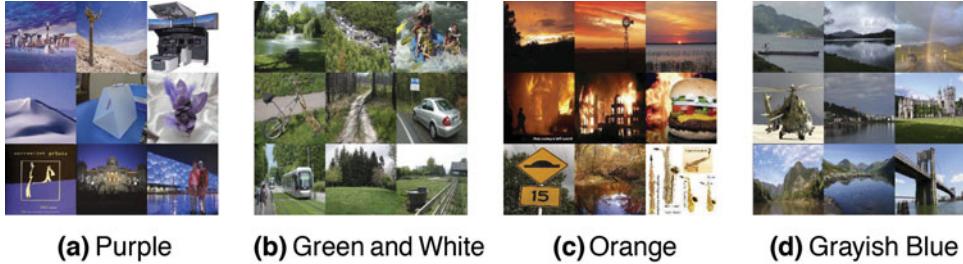


Figure 5. Example clusters from the Gaussian mixture model. We assign each image to its most likely mixture cluster. The subfigures show nine randomly sampled images from four clusters; their namings are subjective.

dimension and it has the same base measure $h(\cdot)$ and log normalizer $a(\cdot)$.

Having established their parametric forms, let ν_j denote the variational parameter for the j th variational factor. When we update each factor, we set its parameter equal to the expected parameter of the complete conditional,

$$\nu_j = \mathbb{E} [\eta_j(\mathbf{z}_{-j}, \mathbf{x})]. \quad (40)$$

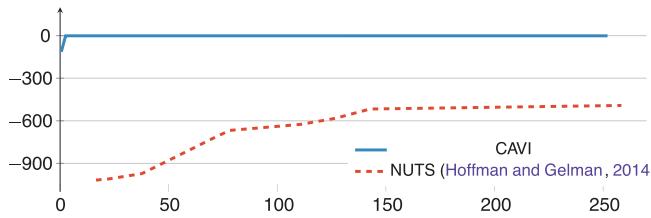


Figure 6. Comparison of CAVI to a Hamiltonian Monte Carlo-based sampling technique. CAVI fits a Gaussian mixture model to 10,000 images in less than a minute.

This expression facilitates deriving CAVI algorithms for many complex models.

4.2. Conditional Conjugacy and Bayesian Models

One important special case of exponential family models are *conditionally conjugate models* with local and global variables. Models like this come up frequently in Bayesian statistics and statistical machine learning, where the global variables are the “parameters” and the local variables are per-data-point latent variables.

Conditionally conjugate models. Let β be a vector of *global latent variables*, which potentially govern any of the data. Let \mathbf{z} be a vector of *local latent variables*, whose i th component only governs data in the i th “context.” The joint density is

$$p(\beta, \mathbf{z}, \mathbf{x}) = p(\beta) \prod_{i=1}^n p(z_i, x_i | \beta). \quad (41)$$

- Wingate, D., and Weber, T. (2013), "Automated Variational Inference in Probabilistic Programming," arXiv preprint, arXiv:1301.1299. Available at <https://arxiv.org/abs/1301.1299> [873]
- Winn, J., and Bishop, C. (2005), "Variational Message Passing," *Journal of Machine Learning Research*, 6, 661–694. [863]
- Wipf, D., and Nagarajan, S. (2009), "A Unified Bayesian Framework for MEG/EEG Source Imaging," *NeuroImage*, 44, 947–966. [871]
- Woolrich, M., Behrens, T., Beckmann, C., Jenkinson, M., and Smith, S. (2004), "Multilevel Linear Modeling for fMRI Group Analysis using Bayesian Inference," *Neuroimage*, 21, 1732–1747. [871]
- Xing, E., Wu, W., Jordan, M. I., and Karp, R. (2004), "Logos: A Modular Bayesian Model for de novo motif Detection," *Journal of Bioinformatics and Computational Biology*, 2, 127–154. [871]
- Yedidia, J. S., Freeman, W. T., and Weiss, Y. (2001), "Generalized Belief Propagation," in *Neural Information Processing Systems*, pp. 689–695. [860]
- Yogatama, D., Wang, C., Routledge, B., Smith, N. A., and Xing, E. (2014), "Dynamic Language Models for Streaming Text," *Transactions of the Association for Computational Linguistics*, 2, 181–192. [871]
- You, C., Ormerod, J., and Muller, S. (2014), "On Variational Bayes Estimation and Variational Information Criteria for Linear Regression Models," *Australian & New Zealand Journal of Statistics*, 56, 73–87. [872]
- Yu, T., and Wu, Y. (2005), "Decentralized Multiple Target Tracking using Netted Collaborative Autonomous Trackers," in *Computer Vision and Pattern Recognition*, pp. 939–946. [871]
- Zumer, J., Attias, H., Sekihara, K., and Nagarajan, S. (2007), "A Probabilistic Algorithm Integrating Source Localization and Noise Suppression for MEG and EEG Data," *NeuroImage*, 37, 102–115. [871]