

---

# A Variational Inference Approach for Locally Private Inference of Poisson Factorization Models

---

**Alexandra Schofield**  
Cornell University  
xanda@cs.cornell.edu

**Aaron Schein**  
University of Massachusetts, Amherst  
aschein@cs.umass.edu

**Zhiwei Steven Wu**  
University of Minnesota  
zsw@umn.edu

**Hanna Wallach**  
Microsoft Research  
hanna@dirichlet.net

## Abstract

Recent work that introduces local privacy to Bayesian Poisson factorization model inference, but its proposed MCMC algorithm for inference suffers from significant performance issues. We derive a coordinate ascent variational inference algorithm to infer factorization models efficiently from private data. Our model relies on several new properties we prove about Bessel distributions. Our method produces a factor of 20 speedup in a synthetic experiment for model inference.

## 1 Introduction

Poisson factorization models are a powerful tool in social science and machine learning for understanding data generated by people. For instance, a corporate email corpus can surface count data such as interactions between parties, counts of words in documents, and logs of types of events and their actors recorded in these messages. However, these data may also encode interactions and content that their human generators wish to keep private, which risks leaking information specific to a few interactions in the final low-rank model.

Most existing methods for private inference of exponential family models operate in the *central model* [15, 3], in which a trusted aggregator collects all the data and perturbs it to ensure privacy. Recent work by [20] provides a privacy-preserving method for Poisson factorization models in the *local model*, in which individual users perturb their own data to ensure privacy without relying on any trusted party. The privacy guarantee uses a generalization of local differential privacy [26] called *limited-precision local privacy*, which allows privacy protection on observations of variable scales (e.g. ranging from an entire document to a single word token). A corresponding MCMC procedure for inference of Poisson factorization models from prior work by [20] relies on the recharacterization of the traditional data generative process to include the addition of perturbation through two-sided geometric noise [12], with the true data characterized a random variable generated using the Skellam and Bessel distributions.

The implementation of this method in prior work applies an MCMC approach to iteratively resample the latent variables, such that after an initial burn-in period of resampling, these samples may be aggregated to estimate the true model parameters. However, this approach to private Bayesian Poisson matrix factorization is prohibitively slow for two primary reasons. First, it requires sampling from the Bessel distribution, which is much slower than sampling for nonprivate Poisson factorization. Second, the addition of private noise produces a denser matrix of observations, in which no observed zero entry in the data is necessarily a true zero in the original data. This addition converts inference from a sparse problem, where zero counts could be ignored, to a dense one, where every entry in the observed data requires an expensive sampling procedure.

We propose a new mean-field variational inference (VI) procedure to infer these models. To do this, we first prove that the Bessel distribution is in the exponential family if its parameter  $\nu$  is fixed, and that it is the optimal Q-distribution for the latent count variables of the true data. However, traditional mean-field inference fails, as the resulting expectation disrupts downstream inference due to the non-integer value of the expectation of the Bessel. We overcome this obstacle by defining the Q-distribution over the count to be a delta-spike at the *mode* of the optimal Bessel distribution (which is integral, by definition). To show that this mode approximation is close to the value obtained using the mean, we prove that the absolute difference between the mean and mode of a Bessel distribution is bounded by one. Finally, we show that in a synthetic experiment, this method converges to reasonable estimates of the true model parameters 20 times faster than the MCMC approach.

## 2 Background

**Bayesian inference** describes a family of methods to infer the posterior distribution  $P(Z | Y)$  of latent variables  $Z$  given data  $Y$ . This is most often performed using Markov chain Monte Carlo (MCMC) or mean-field coordinate ascent variational inference (CAVI). An MCMC inference algorithm iteratively re-samples each latent variable from its *complete conditional*  $z_n \sim P(z_n | Z_{\setminus n}, Y)$  where  $Z_{\setminus n}$  denotes the set of all latent variables except  $z_n$ . CAVI inference algorithms specify a factorized variational distribution over latent variables  $Q(Z) = \prod_n Q(z_n)$  and then optimize its parameters to minimize the KL-divergence from it and to the exact posterior. A general result is that the optimal factorized  $Q(Z)$  consists of  $Q(z_n)$  factors that are proportional to the geometric expectation of the corresponding complete conditional—i.e.,  $Q^*(z_n) \propto \mathbb{G}_{Q_{\setminus z_n}} [P(z_n | Z_{\setminus n}, Y)]$  where the geometric expectation is  $\mathbb{G}[\cdot] = \exp(\mathbb{E}[\ln \cdot])$ . In practice, CAVI requires orders of magnitude fewer iterations than MCMC to converge and is easier to adapt to large-scale streaming data settings.

**Poisson factorization** [23, 8, 30, 13, 14] is a broad class of models for learning latent structure from discrete data, containing many of the most widely used models in the social sciences, including topic models for text corpora [4, 6, 7], population models for genetic data [16], stochastic block models for social networks [2, 13, 29], and tensor factorization for dyadic data [27, 10, 22, 18, 21]. It further includes deep hierarchical models [17, 31], dynamic models [9, 1, 19], and many others.

Poisson factorization assumes that each observed count  $y_i$  is a Poisson random variable— $y_i \sim \text{Poisson}(\mu_i)$ —with unknown rate parameter  $\mu_i$  that is a deterministic function of shared model parameters. We use the multi-index notation  $\mathbf{i}$  for generality; in the special case of Poisson *matrix* factorization, each count is indexed by a row and a column—e.g.,  $\mathbf{i} = (d, v)$ —and its latent rate is defined as the dot product of corresponding row and column parameters—i.e.,  $\mu_i = \mu_i = \theta_d^\top \phi_v$ . In many cases, the latent rate is defined to be a *linear* function of shared model parameters—i.e.,  $\mu_i = \sum_{k=1}^K \mu_{ik}$ , where  $\mu_{ik}$  is a function of latent parameters specific to latent component  $k$ . In Poisson matrix factorization, we define  $\mu_{ik}$  as the product of the latent parameters associated with  $\mathbf{i} = (d, v)$ , or  $\mu_{ik} = \theta_{dk} \phi_{kv}$ .

**MCMC inference for Poisson factorization** is a popular iterative Bayesian inference algorithm. In *nonprivate* Poisson factorization, when the latent rate is a linear function, the observed count can be interpreted as the sum of latent sub-counts—i.e.,  $y_i = \sum_{k=1}^K y_{ik}$ —where each sub-count is an independent Poisson random variable  $y_{ik} \sim \text{Poisson}(\mu_{ik})$ . The first step in MCMC inference is to sample the vector of sub-counts  $\vec{y}_i = (y_{ik})_{k=1}^K$  from its complete conditional, which is a multinomial whose priors are determined by  $\mu_i$ :

$$(\vec{y}_i | -) \sim \text{Multinom}(y_i, \mu_i) \quad (1)$$

where  $\mu_i = (\mu_{ik})_{k=1}^K$ . The multinomial distribution is often parameterized with a *probability* vector  $\mathbf{p}_i$ . In this paper, we parameterize it using a *proportionality* vector  $\mu_i$  which leaves implicit the normalization step—i.e.,  $p_{ik} = \frac{\mu_{ik}}{\sum_{k'=1}^K \mu_{ik'}}$ . Conditioned on samples of the sub-counts, updates to the other latent variables are model-specific—i.e., depend on the particular structure of the rate function and the priors over the parameters. However, in what follows, we derive a general approach that applies to all models with a linear rate  $\mu_i$ .

**Coordinate Ascent Variational Inference (CAVI) for Poisson factorization** is an alternate approach to MCMC. It relies on the approximation of the true posterior distribution of the factorization model with a set of *variational distributions* whose distributions are independent. By iteratively

re-optimizing the values of the arguments, or *variational parameters*, of these approximating distributions, it is possible to iteratively converge to an estimate of the true distribution parameters without sampling.

The optimal variational family for  $\vec{y}_i$  is proportional to the geometric expectation of its complete conditional (given in Equation 1) and equal to:

$$Q^*(\vec{y}_i) = \text{Multinom} \left( \vec{y}_i; y_i, \left( \mathbb{G}_Q [\mu_{ik}] \right)_{k=1}^K \right) \quad (2)$$

The geometric expectations  $\mathbb{G}_Q [\mu_{ik}] = e^{\mathbb{E}_Q [\ln \mu_{ik}]}$  can be understood as messages from the factors  $Q(\mu_{ik})$  in the context of a message-passing algorithm between variational parameters. In standard cases, these expectations are available in closed-form. The factors  $Q(\mu_{ik})$  depend in turn on messages from  $Q(\vec{y}_i)$  expressed via the expectation of  $y_i$ :

$$\mathbb{E}_Q [y_{ik}] = y_i \frac{\mathbb{G}_Q [\mu_{ik}]}{\sum_{k'=1}^K \mathbb{G}_Q [\mu_{ik'}]}. \quad (3)$$

**Locally private MCMC for Poisson factorization** does not permit us to directly observe the count  $y_i$ . Instead, we observe a noised version of it  $\tilde{y}_i^{(\pm)} = \tau_i + y_i$  where  $\tau_i \sim 2\text{Geo}(\alpha_i)$  is a two-sided geometric random variable generated based on privacy parameter  $\alpha_i \in [0, 1]$ . [20] show that such perturbation mechanism satisfies  $(N, N \ln(1/\alpha_i))$ -limited-precision local privacy.<sup>1</sup> MCMC in this case proceeds by treating the true data  $y_i$  itself as a latent variable and re-sampling it from its complete conditional. To do this, one can augment standard Poisson factorization with a set of auxiliary variables for each data point  $\mathcal{A}_i = \left\{ \lambda_i^{(+)}, \lambda_i^{(-)}, g_i^{(+)}, g_i^{(-)}, \tilde{y}_i^{(+)}, m_i \right\}$ , related to a reinterpretation of the generative process for two-sided geometric noise as given in Appendix A. The key feature of this auxiliary variable scheme is the two-step sampling procedure of a value of the sensitive count  $y_i$  as a latent variable:

$$(m_i | -) \sim \text{Bessel} \left( |\tilde{y}_i^{(\pm)}|, 2\sqrt{\lambda_i^{(-)} (\lambda_i^{(+)} + \mu_i)} \right) \quad (4)$$

$$(\tilde{y}_i^{(+)} | -) \sim \text{Multinom} \left( \tilde{y}_i^{(+)}, (\lambda_i^{(+)}, \mu_{i1}, \dots, \mu_{iK}) \right) \quad (5)$$

Here,  $\tilde{y}_i^{(+)}$  is defined deterministically as the maximum of  $m_i$  and  $m_i + \tilde{y}_i^{(\pm)}$ . The noisy observation  $\tilde{y}_i^{(\pm)}$  and  $\tilde{y}_i^{(+)} = (g_i^{(+)}, y_{i1}, \dots, y_{iK})$  is a vector of latent sub-counts that sum to  $\tilde{y}_i^{(+)}$ : the first of these sub-counts  $g_i^{(+)}$  represents the positive noise added to the sensitive data  $y_i = \sum_{k=1}^K y_{ik}$ , while the latter sub-counts constitute the vector of sub-counts corresponding to each latent component in non-private Poisson factorization (see Equation 1).

### 3 Locally private variational inference for Poisson factorization

In this section we describe the challenges to deriving CAVI for locally private Poisson factorization and sketch our solutions. In CAVI, we look to impose a factorized variational distribution over all latent variables which, in this case, includes the set of auxiliary variables  $\mathcal{A}_i$  for each data point mentioned in the last section. The main technical challenge lies in finding a variational families for  $m_i$  and  $\tilde{y}_i^{(+)}$  that yield a good approximation and closed-form messages to the other factors. Towards this goal, we prove the following two related propositions:

**Proposition 1.** (Proven in Appendix B.) *The Bessel distribution  $m \sim \text{Bessel}(\nu, a)$  for fixed  $\nu$  is an exponential family with sufficient statistic  $T_\nu(m) = 2m + \nu$ , natural parameter  $\eta_\nu(m) = \log(\frac{a}{2})$ , and base measure  $h_\nu(m) = \frac{1}{m! \Gamma(m + \nu + 1)}$ .*

**Proposition 2.** *The optimal variational distribution for  $m_i$  is a Bessel distribution:*

$$Q^*(m_i) = \text{Bessel} \left( m_i; |\tilde{y}_i^{(\pm)}|, 2\sqrt{\mathbb{G}_Q [\lambda_i^{(-)}] \mathbb{G}_Q [\lambda_i^{(+)} + \mu_i]} \right). \quad (6)$$

<sup>1</sup>A local randomizer  $R$  satisfies  $(N, \epsilon)$ -limited-precision local privacy if for any two observations  $y, y' \in \mathbb{N}^d$  such that  $\|y - y'\|_1 \leq N$ ,  $R$  satisfies  $\Pr[R(y) = r] \leq \exp(\epsilon) \Pr[R(y') = r]$  for any  $r$  in the output range.

Unfortunately, selecting the Bessel distribution as the family for  $Q(m_{\mathbf{i}})$  prevents a closed-form solution to  $Q^*(\tilde{\mathbf{y}}_{\mathbf{i}}^{(+)})$ : specifically, it results in the expression of the integer number of samples in a binomial distribution downstream in the generative process as a non-integer expectation of  $m_{\mathbf{i}}$ . To overcome this, we instead select the variational family for  $Q(m_{\mathbf{i}})$  to be a Dirac delta function at the *mode* of the optimal family (in Equation 6):

$$Q(m_{\mathbf{i}}) = \delta [m_{\mathbf{i}} = \text{mode}(Q^*(m_{\mathbf{i}}))] \quad (7)$$

This choice allows a closed-form solution for the optimal variational distribution over  $\tilde{\mathbf{y}}_{\mathbf{i}}^{(+)}$ :

$$Q^*(\tilde{\mathbf{y}}_{\mathbf{i}}^{(+)}) = \text{Multinom} \left( \tilde{\mathbf{y}}_{\mathbf{i}}^{(+)}; \mathbb{E}_Q [\tilde{\mathbf{y}}_{\mathbf{i}}^{(+)}], \left( \mathbb{G}_Q [\lambda_{\mathbf{i}}^{(+)}], \mathbb{G}_Q [\mu_{i1}], \dots, \mathbb{G}_Q [\mu_{iK}] \right) \right) \quad (8)$$

Furthermore, we can prove that this degenerate choice variational family is in fact nearly optimal in practice. Since the other factors depend on  $Q(m_{\mathbf{i}})$  only through the message  $\mathbb{E}_Q [m_{\mathbf{i}}]$ , we need only show that the difference between the optimal message and ours is small. Under the optimal  $Q^*(m_{\mathbf{i}})$ , this message equals the *expected value* under the Bessel distribution in Equation 6 while in our solution, it equals the *mode* of that same distribution. To bound the error between these two, we prove the following proposition:

**Proposition 3.** (Proven in Appendix C.) *The absolute difference between the mean and mode of the Bessel distribution is bounded by 1:*

$$| \mathbb{E}_{\text{Bessel}(m; \nu, a)} [m] - \text{mode}(\text{Bessel}(m; \nu, a)) | \leq 1.$$

Equations 6 and 8 provide solutions to the main challenge of deriving a tractable and nearly optimal variational family. The closed-form solutions to the optimal variational families for all other auxiliary variables are more straightforward to derive and we give them in the supplementary material (see Appendix D). There is one final challenge: the log term inside the expectation  $\mathbb{G}_Q [\lambda_{\mathbf{i}}^{(+)} + \mu_{\mathbf{i}}] = e^{\mathbb{E}_Q [\ln \lambda_{\mathbf{i}}^{(+)} + \mu_{\mathbf{i}}]}$  in Equation 6 prevents the derivation of a closed form. However, we can derive a first-order Taylor approximation via the delta method [24].

## 4 Discussion

We present a new CAVI algorithm for inference of Bayesian Poisson factorization models under local privacy. Our method relies on two key theoretical insights about the Bessel distribution: first, that with fixed  $\nu$ , it belongs to the exponential family, and second, that its mode is an integer neighbor of its mean. Using these, we implement a tool with significant performance improvements over even an optimized version of the prior MCMC algorithm. On a synthetic test case of inferring a rank-50 factorization of an 1000 by 1000 data observation matrix, we obtain a 20x speedup in model inference. We detail the construction of this synthetic performance experiment in Appendix E.

## Acknowledgments

We would like to thank David Mimno, and collaborators at Microsoft Research and in the Cornell Natural Language Processing group for their support and feedback. We would also like to thank the reviewers for their helpful comments.

## References

- [1] A. Acharya, J. Ghosh, and M. Zhou. Nonparametric Bayesian factor analysis for dynamic count matrices. arXiv:1512.08996, 2015.
- [2] B. Ball, B. Karrer, and M. E. J. Newman. Efficient and principled method for detecting communities in networks. *Physical Review E*, 84(3):036103, 2011.
- [3] G. Bernstein and D. Sheldon. Differentially private bayesian inference for exponential families. In *Advances in Neural Information Processing Systems 31*, 2018.
- [4] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.

- [5] M. Braun and J. McAuliffe. Variational inference for large-scale models of discrete choice. *Journal of the American Statistical Association*, 105(489):324–335, 2010.
- [6] W. Buntine and A. Jakulin. Applying discrete PCA in data analysis. In *Proceedings of the 20th Conference on Uncertainty in Artificial Intelligence*, pages 59–66, 2004.
- [7] J. Canny. GaP: a factor model for discrete data. In *Proceedings of the 27th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 122–129, 2004.
- [8] A. T. Cemgil. Bayesian inference for nonnegative matrix factorisation models. *Computational Intelligence and Neuroscience*, 2009.
- [9] L. Charlin, R. Ranganath, J. McInerney, and D. M. Blei. Dynamic Poisson factorization. In *Proceedings of the 9th ACM Conference on Recommender Systems*, pages 155–162, 2015.
- [10] E. C. Chi and T. G. Kolda. On tensors, sparsity, and nonnegative factorizations. *SIAM Journal on Matrix Analysis and Applications*, 33(4):1272–1299, 2012.
- [11] L. Devroye. Simulating Bessel random variables. *Statistics & Probability Letters*, 57(3):249–257, 2002.
- [12] A. Ghosh, T. Roughgarden, and M. Sundararajan. Universally utility-maximizing privacy mechanisms. *SIAM Journal on Computing*, 41(6):1673–1693, 2012.
- [13] P. K. Gopalan and D. M. Blei. Efficient discovery of overlapping communities in massive networks. *Proceedings of the National Academy of Sciences*, 110(36):14534–14539, 2013.
- [14] J. Paisley, D. M. Blei, and M. I. Jordan. Bayesian nonnegative matrix factorization with stochastic variational inference. In E. M. Airoldi, D. M. Blei, E. A. Erosheva, and S. E. Fienberg, editors, *Handbook of Mixed Membership Models and Their Applications*, pages 203–222. 2014.
- [15] M. Park, J. Foulds, K. Chaudhuri, and M. Welling. Private topic modeling. *arXiv:1609.04120*, 2016.
- [16] J. K. Pritchard, M. Stephens, and P. Donnelly. Inference of population structure using multilocus genotype data. *Genetics*, 155(2):945–959, 2000.
- [17] R. Ranganath, L. Tang, L. Charlin, and D. Blei. Deep exponential families. In *Proceedings of the 18th International Conference on Artificial Intelligence and Statistics*, pages 762–771, 2015.
- [18] A. Schein, J. Paisley, D. M. Blei, and H. Wallach. Bayesian Poisson tensor factorization for inferring multilateral relations from sparse dyadic event counts. In *Proceedings of the 21st ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1045–1054, 2015.
- [19] A. Schein, H. Wallach, and M. Zhou. Poisson–gamma dynamical systems. In *Advances in Neural Information Processing Systems 29*, pages 5005–5013, 2016.
- [20] A. Schein, Z. S. Wu, A. Schofield, M. Zhou, and H. Wallach. Locally private bayesian inference for count models. *arXiv preprint arXiv:1803.08471*, 2018.
- [21] A. Schein, M. Zhou, D. M. Blei, and H. Wallach. Bayesian Poisson Tucker decomposition for learning the structure of international relations. In *Proceedings of the 33rd International Conference on Machine Learning*, 2016.
- [22] M. N. Schmidt and M. Morup. Nonparametric Bayesian modeling of complex networks: an introduction. *IEEE Signal Processing Magazine*, 30(3):110–128, 2013.
- [23] M. K. Titsias. The infinite gamma–Poisson feature model. In *Advances in Neural Information Processing Systems 21*, pages 1513–1520, 2008.
- [24] J. M. Ver Hoef. Who invented the delta method? *The American Statistician*, 66(2):124–127, 2012.

- [25] C. Wang and D. M. Blei. Variational inference in nonconjugate models. *Journal of Machine Learning Research*, 14(Apr):1005–1031, 2013.
- [26] S. L. Warner. Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 60(309):63–69, 1965.
- [27] M. Welling and M. Weber. Positive tensor factorization. *Pattern Recognition Letters*, 22(12):1255–1261, 2001.
- [28] L. Yuan and J. D. Kalbfleisch. On the Bessel distribution and related problems. *Annals of the Institute of Statistical Mathematics*, 52(3):438–447, 2000.
- [29] M. Zhou. Infinite edge partition models for overlapping community detection and link prediction. In *Proceedings of the 18th Conference on Artificial Intelligence and Statistics*, pages 1135–1143, 2015.
- [30] M. Zhou and L. Carin. Augment-and-conquer negative binomial processes. In *Advances in Neural Information Processing Systems 25*, pages 2546–2554, 2012.
- [31] M. Zhou, Y. Cong, and B. Chen. The Poisson gamma belief network. In *Advances in Neural Information Processing Systems 28*, pages 3043–3051, 2015.

## A Different factorizations of the joint posterior

To find a variational inference algorithm, we leverage existing equivalent characterizations of the generative process for Poisson factorization models under differential privacy.<sup>2</sup> We consider in this case the variable  $y_{\mathbf{i}}$  to be a true observed data point, while  $\tilde{y}_{\mathbf{i}}^{(\pm)}$  is the observed version with random noise. We treat two-sided geometric noise as being the difference of two Poisson variables  $g_{\mathbf{i}}^{(+)}$  and  $g_{\mathbf{i}}^{(-)}$ , generated with Gamma-distributed priors  $\lambda_{\mathbf{i}}^{(+)}$  and  $\lambda_{\mathbf{i}}^{(-)}$ , respectively. We define  $\tilde{y}_{\mathbf{i}}^{(+)} = y_{\mathbf{i}} + g_{\mathbf{i}}^{(+)}$  and  $\tilde{y}_{\mathbf{i}}^{(\pm)} = y_{\mathbf{i}} + g_{\mathbf{i}}^{(+)} - g_{\mathbf{i}}^{(-)}$ . We can determine whether  $\tilde{y}_{\mathbf{i}}^{(\pm)}$  is positive or negative by comparing the values of  $\tilde{y}_{\mathbf{i}}^{(+)}$  and  $g_{\mathbf{i}}^{(-)}$ ; we refer to the minimum of these two as  $m_{\mathbf{i}}$ .

$$P\left(g_{\mathbf{i}}^{(+)}, g_{\mathbf{i}}^{(-)}, (y_{ik})_{k=1}^K, y_{\mathbf{i}}, \tilde{y}_{\mathbf{i}}^{(+)}, \tilde{y}_{\mathbf{i}}^{(\pm)}, m_{\mathbf{i}}\right). \quad (9)$$

The most straightforward factorization of this joint first generates all the Poisson random variables, then computes the remaining variables given their deterministic relationships to the underlying Poissons:

$$\begin{aligned} & P\left(g_{\mathbf{i}}^{(+)}, g_{\mathbf{i}}^{(-)}, (y_{ik})_{k=1}^K, y_{\mathbf{i}}, \tilde{y}_{\mathbf{i}}^{(+)}, \tilde{y}_{\mathbf{i}}^{(\pm)}, m_{\mathbf{i}}\right) \\ &= \text{Poisson}\left(g_{\mathbf{i}}^{(+)}; \lambda_{\mathbf{i}}^{(+)}\right) \text{Poisson}\left(g_{\mathbf{i}}^{(-)}; \lambda_{\mathbf{i}}^{(-)}\right) \left(\prod_{k=1}^K \text{Poisson}(y_{ik}; \theta_{dk} \phi_{kv})\right) \mathbb{1}\left(y_{\mathbf{i}} = \sum_{k=1}^K y_{ik}\right) \\ & \quad \mathbb{1}\left(\tilde{y}_{\mathbf{i}}^{(+)} = y_{\mathbf{i}} + g_{\mathbf{i}}^{(+)}\right) \mathbb{1}\left(\tilde{y}_{\mathbf{i}}^{(\pm)} = \tilde{y}_{\mathbf{i}}^{(+)} - g_{\mathbf{i}}^{(-)}\right) \mathbb{1}\left(m_{\mathbf{i}} = \min\{\tilde{y}_{\mathbf{i}}^{(+)}, g_{\mathbf{i}}^{(-)}\}\right). \end{aligned} \quad (10)$$

We can equivalently first generate the sums of Poissons and then thin them using multinomial and binomial draws. In the following equation, the delta functions are implicitly present in the multinomial and binomial PMFs. Note that we write the probability parameters in the multinomial and binomial PMFs as unnormalized vectors. Also note that  $\mu_{\mathbf{i}} = \sum_k \prod_d \theta_{i_d, k}^{(i)}$ , where  $i_d$  is the index into the  $d$ th

<sup>2</sup>Parts of this and Appendix D have been excerpted from supplementary material in a prior preprint version of (author?) [20].

dimension of index vector  $\mathbf{i}$ .

$$\begin{aligned}
& P\left(g_{\mathbf{i}}^{(+)}, g_{\mathbf{i}}^{(-)}, (y_{\mathbf{i}k})_{k=1}^K, y_{\mathbf{i}}, \tilde{y}_{\mathbf{i}}^{(+)}, \tilde{y}_{\mathbf{i}}^{(\pm)}, m_{\mathbf{i}}\right) \\
&= \text{Poisson}\left(g_{\mathbf{i}}^{(-)}; \lambda_{\mathbf{i}}^{(-)}\right) \text{Poisson}\left(\tilde{y}_{\mathbf{i}}^{(+)}; \lambda_{\mathbf{i}}^{(+)} + \mu_{\mathbf{i}}\right) \\
& \quad \text{Binom}\left((y_{\mathbf{i}}, g_{\mathbf{i}}^{(+)}); \tilde{y}_{\mathbf{i}}^{(+)}, (\mu_{\mathbf{i}}, \lambda_{\mathbf{i}}^{(+)})\right) \text{Mult}\left((y_{\mathbf{i}k})_{k=1}^K; y_{\mathbf{i}}, \left(\prod_d \theta_{i_d, k}^{(i)}\right)_{k=1}^K\right) \\
& \quad \mathbb{1}\left(\tilde{y}_{\mathbf{i}}^{(\pm)} = \tilde{y}_{\mathbf{i}}^{(+)} - g_{\mathbf{i}}^{(-)}\right) \mathbb{1}\left(m_{\mathbf{i}} = \min\{\tilde{y}_{\mathbf{i}}^{(+)}, g_{\mathbf{i}}^{(-)}\}\right). \tag{11}
\end{aligned}$$

We can equivalently first generate the difference  $\tilde{y}_{\mathbf{i}}^{(\pm)}$  and minimum  $m_{\mathbf{i}}$  as Skellam and Bessel random variables. Conditioned on these variables, we can then compute  $\tilde{y}_{\mathbf{i}}^{(+)}$  and  $g_{\mathbf{i}}^{(-)}$  via their deterministic relationship and, finally, thin  $\tilde{y}_{\mathbf{i}}^{(+)}$  using binomial and multinomial draws:

$$\begin{aligned}
& P\left(g_{\mathbf{i}}^{(+)}, g_{\mathbf{i}}^{(-)}, (y_{\mathbf{i}k})_{k=1}^K, y_{\mathbf{i}}, \tilde{y}_{\mathbf{i}}^{(+)}, \tilde{y}_{\mathbf{i}}^{(\pm)}, m_{\mathbf{i}}\right) \\
&= \text{Skel}\left(\tilde{y}_{\mathbf{i}}^{(\pm)}; \lambda_{\mathbf{i}}^{(+)} + \mu_{\mathbf{i}}, \lambda_{\mathbf{i}}^{(-)}\right) \text{Bes}\left(m_{\mathbf{i}}; |\tilde{y}_{\mathbf{i}}^{(\pm)}|, 2\sqrt{\lambda_{\mathbf{i}}^{(-)}(\lambda_{\mathbf{i}}^{(+)} + \mu_{\mathbf{i}})}\right) \\
& \quad \mathbb{1}\left(\tilde{y}_{\mathbf{i}}^{(+)} = m_{\mathbf{i}}\right) \mathbb{1}^{(\tilde{y}_{\mathbf{i}}^{(\pm)} \leq 0)} \mathbb{1}\left(g_{\mathbf{i}}^{(-)} = m_{\mathbf{i}}\right) \mathbb{1}^{(\tilde{y}_{\mathbf{i}}^{(\pm)} > 0)} \mathbb{1}\left(\tilde{y}_{\mathbf{i}}^{(\pm)} = \tilde{y}_{\mathbf{i}}^{(+)} - g_{\mathbf{i}}^{(-)}\right) \\
& \quad \text{Binom}\left((y_{\mathbf{i}}, g_{\mathbf{i}}^{(+)}); \tilde{y}_{\mathbf{i}}^{(+)}, (\mu_{\mathbf{i}}, \lambda_{\mathbf{i}}^{(+)})\right) \text{Mult}\left((y_{\mathbf{i}k})_{k=1}^K; y_{\mathbf{i}}, \left(\prod_d \theta_{i_d, k}^{(i)}\right)_{k=1}^K\right). \tag{12}
\end{aligned}$$

It is this last factorization that enables us to derive the variational distribution.

## B The Bessel distribution as exponential family

**Theorem 1** *The Bessel distribution  $m \sim \text{Bes}(\nu, a)$  for fixed  $\nu$  is an exponential family with sufficient statistic  $T_{\nu}(m) = 2m + \nu$ , natural parameter  $\eta_{\nu}(m) = \log\left(\frac{a}{2}\right)$ , and base measure  $h_{\nu}(m) = \frac{1}{m! \Gamma(m + \nu + 1)}$ .*

*Proof.* The Bessel distribution [28] is a two-parameter distribution over the non-negative integers:

$$f(n; a, \nu) = \frac{\left(\frac{a}{2}\right)^{2n+\nu}}{n! \Gamma(n+\nu+1) I_{\nu}(a)}, \tag{13}$$

where the normalizing constant  $I_{\nu}(a)$  is a modified Bessel function of the first kind—i.e.,

$$I_{\nu}(a) = \sum_{n=0}^{\infty} \frac{\left(\frac{a}{2}\right)^{2n+\nu}}{n! \Gamma(n+\nu+1)}. \tag{14}$$

For fixed and known  $\nu$ , we can rewrite the Bessel PMF as

$$f(n; a, \nu) = \frac{1}{n! \Gamma(n+\nu+1)} \exp\left((2n+\nu) \log\left(\frac{a}{2}\right) - \log I_{\nu}(a)\right). \tag{15}$$

We can then define the following functions:

$$h_{\nu}(n) = \frac{1}{n! \Gamma(n+\nu+1)} \tag{16}$$

$$T_{\nu}(n) = 2n + \nu \tag{17}$$

$$\eta_{\nu}(a) = \log\left(\frac{a}{2}\right) \tag{18}$$

$$A_{\nu}(a) = \log I_{\nu}(a). \tag{19}$$

Finally, we can rewrite the Bessel PMF in the exponential-family form:

$$f(n; a, \nu) = h_{\nu}(n) \exp(\eta_{\nu}(a) \cdot T_{\nu}(n) - A_{\nu}(a)). \tag{20}$$

□

## C The mode and mean of the Bessel distribution

**Theorem 2** *The mode of a Bessel distribution for parameters  $a, \nu$  can be a constant-bounded approximation of the mean:*

$$|\mathbb{E}_{\text{Bessel}(m;\nu,a)}[m] - \text{mode}(\text{Bessel}(m; a, \nu))| \leq 1.$$

The intuitive meaning of this is that the mode of the Bessel distribution is guaranteed to be one of the two integers closest to the mean of the distribution. Given one of these two integers will always be the best integer approximation of this number, we can also say that there cannot exist an integer approximation of the mean of the Bessel that is strictly between the mode and the mean of the Bessel distribution.

*Proof.* A Bessel distribution takes two arguments, which we refer to as its order,  $\nu$  and coordinate,  $a$ . The distribution is defined as:

$$p(x = n | x \sim \text{Bessel}(\nu, a)) = \frac{1}{I_\nu(a)n!\Gamma(n + \nu + 1)} \left(\frac{a}{2}\right)^{2n+\nu},$$

where  $I_\nu(a)$  is a modified Bessel function of the first kind. The arithmetic mean of the distribution is

$$\mathbb{E}_{\text{Bessel}(m;\nu,a)}[m] = \frac{a}{2}R_\nu(a),$$

with  $R_\nu(a)$  referring to the ratio of two Bessel functions:

$$\frac{I_{\nu+1}(a)}{I_\nu(a)}.$$

The Bessel distribution has one or two neighboring integer modes. The mode can be computed directly from the parameters of the distribution without Bessel functions:

$$\text{mode}(\text{Bessel}(\nu, a)) = \left\lfloor \frac{\sqrt{a^2 + \nu^2} - \nu}{2} \right\rfloor.$$

Unlike the mean, which can take arbitrary non-negative real values, the mode is guaranteed by the floor function to be a non-negative integer.

We use the following bound on the mean of a Bessel ratio from **(author?)** [11]:

$$\frac{a}{\nu + 1 + \sqrt{a^2 + (\nu + 1)^2}} \leq R_\nu(a) \leq \frac{a}{\nu + \sqrt{a^2 + \nu^2}}. \quad (21)$$

We multiply through by  $\frac{a}{2}$  to bound the mean of the Bessel distribution:

$$\frac{a^2}{2(\nu + 1 + \sqrt{a^2 + (\nu + 1)^2})} \leq \mathbb{E}_{\text{Bessel}(m;\nu,a)}[m] \leq \frac{a^2}{2(\nu + \sqrt{a^2 + \nu^2})}. \quad (22)$$

We can rewrite these bounds using a difference-of-squares:

$$\begin{aligned} \frac{a^2}{2(\nu + \sqrt{a^2 + \nu^2})} &= \frac{a^2(\sqrt{a^2 + \nu^2} - \nu)}{2(\nu + \sqrt{a^2 + \nu^2})(\sqrt{a^2 + \nu^2} - \nu)} \\ &= \frac{a^2(\sqrt{a^2 + \nu^2} - \nu)}{2a^2 + \nu^2 - \nu^2} \\ &= \frac{\sqrt{a^2 + \nu^2} - \nu}{2}. \end{aligned}$$



This upper bound coincides with the unrounded formulation of the mode. Because the mode is the floor of this quantity, we know it is less than or equal to this upper bound with a difference of less than 1.

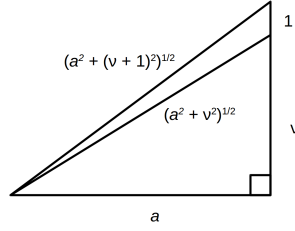
We can convert the lower bound in the same way:

$$\frac{a^2}{2((\nu+1) + \sqrt{a^2 + (\nu^2+1)})} = \frac{\sqrt{a^2 + (\nu+1)^2} - (\nu+1)}{2}.$$

We are interested in bounding the difference between the upper and lower bounds:

$$\frac{\sqrt{a^2 + \nu^2} - \nu}{2} - \frac{\sqrt{a^2 + (\nu+1)^2} - (\nu+1)}{2} = \frac{\sqrt{a^2 + \nu^2} + 1 - \sqrt{a^2 + (\nu+1)^2}}{2}. \quad (23)$$

Knowing that  $\nu$  is positive, we can state that  $\sqrt{a^2 + \nu^2} < \sqrt{a^2 + (\nu+1)^2}$ , or  $\sqrt{a^2 + (\nu+1)^2} - \sqrt{a^2 + \nu^2} > 0$ . Based on the upper slice of the triangle in the figure below, the triangle inequality also gives us another bound, that  $\sqrt{a^2 + \nu^2} + 1 > \sqrt{a^2 + (\nu+1)^2}$ , or  $1 > \sqrt{a^2 + (\nu+1)^2} - \sqrt{a^2 + \nu^2}$ .



Together, these imply that

$$0 \leq \sqrt{a^2 + \nu^2} + 1 - \sqrt{a^2 + (\nu+1)^2} \leq 1.$$

Substituting this in to our computation of the distance between the upper and lower bounds of the mean, we find that

$$\frac{\sqrt{a^2 + \nu^2} + 1 - \sqrt{a^2 + (\nu+1)^2}}{2} < \frac{1}{2}, \quad (24)$$

or that the inferred upper and lower bounds for the mean of the Bessel distribution produce an interval no larger than  $\frac{1}{2}$ . The upper bound of this  $\frac{1}{2}$  interval is the same as the upper bound of the length-1 open interval of the mode of the Bessel. In the most extreme case when the mode is at the bottom end of its range and the mean is at the top, we have

$$\mathbb{E}_{\text{Bessel}(m;\nu,a)}[m] - \text{mode}(\text{Bessel}(\nu, a)) < 1.$$

In the opposite case, we have

$$\text{mode}(\text{Bessel}(\nu, a)) - \mathbb{E}_{\text{Bessel}(m;\nu,a)}[m] \leq \frac{1}{2} < 1.$$

□

## D Deriving the CAVI updates for the variational distribution

Recall that the observed data consists of the difference variables  $\tilde{y}_i^{(\pm)}$ .

$$Q\left(m_i, \tilde{y}_i^{(+)}, g_i^{(+)}, g_i^{(-)}, y_i, (y_{ik})_{k=1}^K\right) \propto \mathbb{G}_Q\left[P(\tilde{y}_i^{(\pm)}, m_i, \tilde{y}_i^{(+)}, g_i^{(+)}, g_i^{(-)}, y_i, (y_{ik})_{k=1}^K \mid -)\right]. \quad (25)$$

Because the likelihood (i.e., the Skellam term) in equation 12 does not depend on any of these latent variables, it disappears entirely. We can then rewrite the right-hand side of equation 25 as:

$$\begin{aligned}
& \mathbb{G}_Q \left[ P(m_{\mathbf{i}}, \tilde{y}_{\mathbf{i}}^{(+)}, g_{\mathbf{i}}^{(+)}, g_{\mathbf{i}}^{(-)}, y_{\mathbf{i}}, (y_{ik})_{k=1}^K \mid \tilde{y}_{\mathbf{i}}^{(\pm)} -) \right] \\
&= \mathbb{G}_Q \left[ \text{Bes} \left( m_{\mathbf{i}}; |\tilde{y}_{\mathbf{i}}^{(\pm)}|, 2\sqrt{\lambda_{\mathbf{i}}^{(-)}(\lambda_{\mathbf{i}}^{(+)} + \mu_{\mathbf{i}})} \right) \right] \\
& \quad \mathbb{1} \left( \tilde{y}_{\mathbf{i}}^{(+)} = m_{\mathbf{i}} \right)^{\mathbb{1}(\tilde{y}_{\mathbf{i}}^{(\pm)} \leq 0)} \mathbb{1} \left( g_{\mathbf{i}}^{(-)} = m_{\mathbf{i}} \right)^{\mathbb{1}(\tilde{y}_{\mathbf{i}}^{(\pm)} > 0)} \mathbb{1} \left( \tilde{y}_{\mathbf{i}}^{(\pm)} = \tilde{y}_{\mathbf{i}}^{(+)} - g_{\mathbf{i}}^{(-)} \right) \\
& \quad \mathbb{G}_Q \left[ \text{Binom} \left( (y_{\mathbf{i}}, g_{\mathbf{i}}^{(+)}); \tilde{y}_{\mathbf{i}}^{(+)}, (\mu_{\mathbf{i}}, \lambda_{\mathbf{i}}^{(+)}) \right) \text{Mult} \left( (y_{ik})_{k=1}^K; y_{\mathbf{i}}, \left( \prod_d \theta_{i_d,k}^{(i)} \right)_{k=1}^K \right) \right]. \quad (26)
\end{aligned}$$

Theorem 6.1 states that the Bessel distribution for fixed first parameter is an exponential family. We can therefore use standard results to push in the geometric expectations:

$$\begin{aligned}
& \mathbb{G}_Q \left[ P(m_{\mathbf{i}}, \tilde{y}_{\mathbf{i}}^{(+)}, g_{\mathbf{i}}^{(+)}, g_{\mathbf{i}}^{(-)}, y_{\mathbf{i}}, (y_{ik})_{k=1}^K \mid \tilde{y}_{\mathbf{i}}^{(\pm)} -) \right] \\
&= \text{Bes} \left( m_{\mathbf{i}}; |\tilde{y}_{\mathbf{i}}^{(\pm)}|, 2\sqrt{\mathbb{G}_Q \left[ \lambda_{\mathbf{i}}^{(-)} \right] \left( \mathbb{G}_Q \left[ \lambda_{\mathbf{i}}^{(+)} + \mu_{\mathbf{i}} \right] \right)} \right) \\
& \quad \mathbb{1} \left( \tilde{y}_{\mathbf{i}}^{(+)} = m_{\mathbf{i}} \right)^{\mathbb{1}(\tilde{y}_{\mathbf{i}}^{(\pm)} \leq 0)} \mathbb{1} \left( g_{\mathbf{i}}^{(-)} = m_{\mathbf{i}} \right)^{\mathbb{1}(\tilde{y}_{\mathbf{i}}^{(\pm)} > 0)} \mathbb{1} \left( \tilde{y}_{\mathbf{i}}^{(\pm)} = \tilde{y}_{\mathbf{i}}^{(+)} - g_{\mathbf{i}}^{(-)} \right) \\
& \quad \text{Binom} \left( (y_{\mathbf{i}}, g_{\mathbf{i}}^{(+)}); \mathbb{E}_Q \left[ \tilde{y}_{\mathbf{i}}^{(+)} \right], \left( \mathbb{G}_Q \left[ \mu_{\mathbf{i}} \right], \mathbb{G}_Q \left[ \lambda_{\mathbf{i}}^{(+)} \right] \right) \right) \\
& \quad \text{Mult} \left( (y_{ik})_{k=1}^K; \mathbb{E}_Q \left[ y_{\mathbf{i}} \right], \left( \mathbb{G}_Q \left[ \prod_d \theta_{i_d,k}^{(i)} \right] \right)_{k=1}^K \right). \quad (27)
\end{aligned}$$

There are two expectations that do not have an analytic form:

$$\mathbb{G}_Q \left[ \lambda_{\mathbf{i}}^{(+)} + \mu_{\mathbf{i}} \right] = \exp \left( \mathbb{E}_Q \left[ \ln \left( \lambda_{\mathbf{i}}^{(+)} + \sum_{k=1}^K \prod_d \theta_{i_d,k}^{(i)} \right) \right] \right) \quad (28)$$

and

$$\mathbb{G}_Q \left[ \mu_{\mathbf{i}} \right] = \exp \left( \mathbb{E}_Q \left[ \ln \left( \sum_{k=1}^K \prod_d \theta_{i_d,k}^{(i)} \right) \right] \right); \quad (29)$$

however, both can be very closely approximated using the delta method [24], which has been previously used in variational inference schemes to approximate intractable expectations [5, 25]. In particular, for some variable  $Y = f(X)$ , expectation  $\mathbb{E}[Y]$  is approximately:

$$\mathbb{E}[Y] = \mathbb{E}[f(X)] \approx f(\mathbb{E}[X]) + \frac{1}{2} f''(\mathbb{E}[X]) \mathbb{V}[X]. \quad (30)$$

In our case, we can consider the two-dimensional Poisson matrix factorization case where  $\mathbf{i}$  are represented by row  $d$  and column  $v$ . We therefore have

$$\mathbb{E}_Q \left[ \ln \mu_{dv} \right] = \mathbb{E}_Q \left[ \ln \left( \sum_{k=1}^K \theta_{dk} \phi_{kv} \right) \right] \approx \ln \left( \mathbb{E}_Q \left[ \sum_{k=1}^K \theta_{dk} \phi_{kv} \right] \right) - \frac{\mathbb{V}_Q \left[ \sum_{k=1}^K \theta_{dk} \phi_{kv} \right]}{2 \left( \mathbb{E}_Q \left[ \sum_{k=1}^K \theta_{dk} \phi_{kv} \right] \right)^2} \quad (31)$$

$$= \ln \left( \sum_{k=1}^K \mathbb{E}_Q \left[ \theta_{dk} \right] \mathbb{E}_Q \left[ \phi_{kv} \right] \right) - \frac{\sum_{k=1}^K \mathbb{V}_Q \left[ \theta_{dk} \phi_{kv} \right]}{2 \left( \sum_{k=1}^K \mathbb{E}_Q \left[ \theta_{dk} \right] \mathbb{E}_Q \left[ \phi_{kv} \right] \right)^2}. \quad (32)$$

Finally, because  $\theta_{dk}$  and  $\phi_{kv}$  are independent, we have

$$\mathbb{V}_Q \left[ \theta_{dk} \phi_{kv} \right] = \mathbb{V}_Q \left[ \theta_{dk} \right] \mathbb{V}_Q \left[ \phi_{kv} \right] + \mathbb{V}_Q \left[ \theta_{dk} \right] \left( \mathbb{E}_Q \left[ \phi_{kv} \right] \right)^2 + \mathbb{V}_Q \left[ \phi_{kv} \right] \left( \mathbb{E}_Q \left[ \theta_{dk} \right] \right)^2. \quad (33)$$

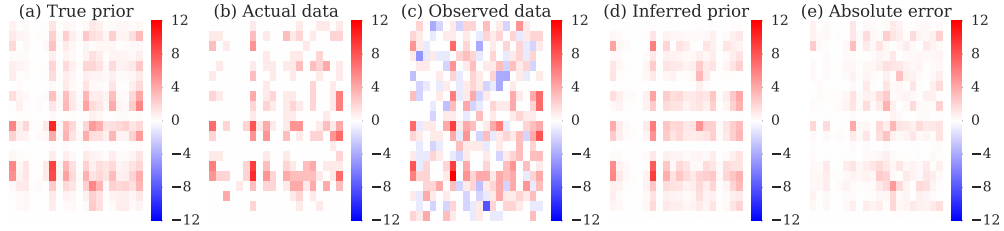


Figure 1: Demonstration of the results of the VI inference process for a 20-word, 20-document synthetic dataset with 3 latent topics. The true data parameters (a) are recovered well by our inference procedure (d), even though the noisy data (d) is much denser than the true data (c). The mean average error (MAE) between the true model parameters and the inferred parameters under privacy is 0.507, with the individual errors pictured in (e).

## E Validation

In order to test this result, we generated synthetic count data using a Poisson matrix factorization formulation analogous to the LDA topic model [4] with  $D$  documents,  $V$  unique terms in a document vocabulary, and  $K$  latent topics. We first used a Gamma prior to generate two matrices of latent parameters,  $\theta$  of dimension  $D \times K$  and  $\phi$  of dimension  $K \times V$ . We then compute the product of these,  $\theta\phi = \mathbf{Y}$ , as the Poisson prior of our data generation process. Finally, we add two-sided geometric noise scaled to  $\epsilon/N = 1$ , a ratio that applies when the privacy budget  $\epsilon$  and the maximum allowed difference between documents  $N$  are equal (e.g., a privacy budget of 2 to privatize 2-word spans). We then test our inference procedure to see how closely it estimates the true parameters of the original model. We find that our model successfully converges to within a reasonable estimate of the true model parameters given the data, as demonstrated in a small example in Figure 1.

Using a larger example of a synthetic 1000-by-1000 matrix of count observations, we test the performance of this algorithm as compared to the performance using inference with MCMC [20]. We observe that a single-threaded version of the code from the MCMC implementation, even with optimizations such as reducing the frequency of parameter resampling for the noise distributions, each iteration of inference takes approximately 0.8 seconds. Using 5000 iterations of burnin and 2500 to collect samples, this takes a combined 1.67 hours to infer, with a final mean average error (MAE) of 0.36. In contrast, our variational inference model takes approximately 2.8 seconds per iteration, and typically requires 100 iterations to converge, resulting in a combined 5 minutes of inference to reach error of 0.52.