

# Variational Inference for the Indian Buffet Process



**Finale Doshi-Velez, Kurt T. Miller, Jurgen Van Gael, Yee Whye Teh**

Computational and Biological Learning Laboratory  
Department of Engineering  
University of Cambridge

Technical Report CBL-2009-001  
May 2009

# Variational Inference for the Indian Buffet Process

Finale Doshi-Velez\*      Kurt T. Miller\*      Jurgen Van Gael\*      Yee Whye Teh  
Cambridge University    University of California, Berkeley    Cambridge University    Gatsby Unit

May 2009

## Abstract

The Indian Buffet Process (IBP) is a nonparametric prior for latent feature models in which observations are influenced by a combination of hidden features. For example, images may be composed of several objects and sounds may consist of several notes. Latent feature models seek to infer these unobserved features from a set of observations; the IBP provides a principled prior in situations where the number of hidden features is unknown. Current inference methods for the IBP have all relied on sampling. While these methods are guaranteed to be accurate in the limit, samplers for the IBP tend to mix slowly in practice. We develop a deterministic variational method for inference in the IBP based on truncating to finite models, provide theoretical bounds on the truncation error, and evaluate our method in several data regimes. This technical report is a longer version of Doshi-Velez et al. (2009).

**Keywords:** Variational Inference, Indian Buffet Process, Nonparametric Bayes, Linear Gaussian, Independent Component Analysis

## 1 Introduction

Many unsupervised learning problems seek to identify a set of unobserved, co-occurring features from a set of observations. For example, given images composed of various objects, we may wish to identify the set of unique objects and determine which images contain which objects. Similarly, we may wish to extract a set of notes or chords from an audio file as well as when each note was played. In scenarios such as these, the number of latent features is often unknown a priori.

Unfortunately, most traditional machine learning approaches require the number of latent features as an input. To apply these traditional approaches to scenarios in which the number of latent features is unknown, we can use model selection to define and manage the trade-off between model complexity and model fit. In contrast, nonparametric Bayesian approaches treat the number of features as a random quantity to be determined as part of the posterior inference procedure.

The most common nonparametric prior for latent feature models is the Indian Buffet Process (IBP) (Griffiths and Ghahramani, 2005). The IBP is a prior on infinite binary matrices that allows us to simultaneously infer which features influence a set of observations and how many features there are. The form of the prior ensures that only a finite number of features will be present in any finite set of observations, but more features may appear as more observations are received. This property is both natural and desirable if we consider, for example, a set of

---

\*Authors contributed equally.

images: any one image contains a finite number of objects, but, as we see more images, we expect to see objects not present in the previous images.

While an attractive model, the combinatorial nature of the IBP makes inference particularly challenging. Even if we limit ourselves to  $K$  features for  $N$  objects, there exist  $2^{NK}$  possible feature assignments. In this combinatorial space, sampling-based inference procedures for the IBP often suffer because they assign specific values to the feature-assignment variables. Hard variable assignments give samplers less flexibility to move between optima, and the samplers may need large amounts of time to escape small optima and find regions with high probability mass. Unfortunately, all current inference procedures for the IBP rely on sampling. These approaches include Gibbs sampling (Griffiths and Ghahramani, 2005), which may be augmented with Metropolis split-merge proposals (Meeds et al., 2007), as well as slice sampling (Teh et al., 2007) and particle filtering (Wood and Griffiths, 2007).

Mean field variational methods, which approximate the true posterior via a simpler distribution, provide a deterministic alternative to sampling-based approaches. Inference involves using optimisation techniques to find a good approximate posterior. Our mean field approximation for the IBP maintains a separate probability for each feature-observation assignment. Optimising these probability values is also fraught with local optima, but using the soft variable assignments—that is, using probabilities instead of sampling hard assignments—gives the variational method a flexibility that samplers lack. In the early stages of the inference, this flexibility can help the variational method avoid small local optima.

Variational approximations have provided benefits for other nonparametric Bayesian models, including Dirichlet Processes (e.g. Blei and Jordan (2004), Kurihara et al. (2007a) and Kurihara et al. (2007b)) and Gaussian Processes (e.g. Winther (2000), Gibbs and MacKay (2000), and Snelson and Ghahramani (2006)). Of all the nonparametric Bayesian models studied so far, however, the IBP is the most combinatorial and is therefore in the most need of a more efficient inference algorithm.

We will show how our variational approximation is better in terms of both computation required and predictive likelihood than the sampling based approaches as we get more data and as the dimensionality of the data grows.

The remainder of this technical report is organised as follows:

- **Background.** Section 2 reviews the likelihood model that we will use throughout the technical report, the Indian Buffet Process (IBP), and the basic Gibbs sampler for the IBP. It also summarises the notation used in later sections. Appendix A details the Gibbs sampling equations used in our tests.
- **Variational Method.** Section 3 reviews the variational inference framework and outlines our specific approach which is based on a truncated representation of the IBP. Sections 4 and 5 contain all of the update equations required to implement the variational inference; full derivations are provided in Appendices C and D. Appendix E describes how similar updates can be derived for the infinite ICA model of Knowles and Ghahramani (2007).
- **Truncation Bound.** Section 6 derives bounds on the expected error due to our use of a truncated representation for the IBP; these bounds can serve as guidelines for what level of truncation may be appropriate. Extended derivations are left to Appendix F.
- **Results.** Section 7 demonstrates how our variational approach scales to high-dimensional data sets while maintaining good predictive performance.

## 2 Background

Section 2.1 describes the likelihood model that we will use throughout this report, Section 2.2 provides background information on the Indian Buffet Process and Section 2.3 summarises the notation used in the remainder of the report.

### 2.1 Data Model

Let  $\mathbf{X}$  be an  $N \times D$  matrix where each of the  $N$  rows contains a  $D$ -dimensional observation. In this report, we consider a likelihood model in which  $\mathbf{X}$  can be approximated by  $\mathbf{Z}\mathbf{A}$  where  $\mathbf{Z}$  is an  $N \times K$  binary matrix and  $\mathbf{A}$  is a  $K \times D$  matrix. Each column of  $\mathbf{Z}$  corresponds to a latent feature where  $z_{nk} \equiv \mathbf{Z}(n, k) = 1$  if feature  $k$  is present in observation  $n$  and 0 otherwise. The parameters for feature  $k$  are stored in row  $k$  of  $\mathbf{A}$ . The observed data  $\mathbf{X}$  is then given by  $\mathbf{Z}\mathbf{A} + \epsilon$ , where  $\epsilon$  is some measurement noise (see Figure 1). We assume that the noise is independent of  $\mathbf{Z}$  and  $\mathbf{A}$  and is uncorrelated across observations.

An important feature of this model is that only the nonzero columns of  $\mathbf{Z}$  affect the likelihood. This will be important when we place a prior on the binary feature matrix  $\mathbf{Z}$  in which  $\mathbf{Z}$  has an infinite number of columns, but only a finite number of entries will be nonzero. Since the all-zero columns do not affect the likelihood, we will only need to reason about the nonzero columns.

$$\begin{matrix} & D \\ \mathbf{X} & \begin{matrix} \square & \square & \square & \square & \square \\ \square & \square & \square & \square & \square \\ \square & \square & \square & \square & \square \\ \square & \square & \square & \square & \square \\ \square & \square & \square & \square & \square \\ \square & \square & \square & \square & \square \end{matrix} \\ N & \end{matrix} = \begin{matrix} & K \\ \mathbf{Z} & \begin{matrix} \blacksquare & \square & \square & \square & \square \\ \square & \blacksquare & \square & \square & \square \\ \square & \square & \blacksquare & \square & \square \\ \square & \square & \square & \blacksquare & \square \\ \square & \square & \square & \square & \blacksquare \\ \square & \square & \square & \square & \square \end{matrix} \\ N & \end{matrix} \cdots \times \begin{matrix} & D \\ \mathbf{A} & \begin{matrix} \square & \square & \square & \square & \square \\ \square & \square & \square & \square & \square \\ \square & \square & \square & \square & \square \\ \square & \square & \square & \square & \square \\ \square & \square & \square & \square & \square \end{matrix} \\ K & \end{matrix} + \epsilon$$

Figure 1: Our likelihood model posits that the data  $\mathbf{X}$  is the product  $\mathbf{Z}\mathbf{A}$  plus some noise.

Given  $\mathbf{X}$  and being in a Bayesian framework, we wish to find the posterior distribution of  $\mathbf{Z}$  and  $\mathbf{A}$ . From Bayes rule,

$$p(\mathbf{Z}, \mathbf{A} | \mathbf{X}) \propto p(\mathbf{X} | \mathbf{Z}, \mathbf{A}) p(\mathbf{Z}) p(\mathbf{A})$$

where we have assumed that  $\mathbf{Z}$  and  $\mathbf{A}$  are independent a priori. The specific application will determine the likelihood function  $p(\mathbf{X} | \mathbf{Z}, \mathbf{A})$  and the feature prior  $p(\mathbf{A})$ . In the linear-Gaussian model which will be the focus of this report, both the noise  $\epsilon$  and the features  $\mathbf{A}$  have Gaussian priors. In Appendix E, we discuss a different likelihood model and show how to apply our variational approximation to this model.

We are left with placing a prior on  $\mathbf{Z}$ . Since we often do not know  $K$ , we desire a prior that allows the number of nonzero columns to be determined at inference time. The Indian Buffet Process is one option for such a prior.

### 2.2 Indian Buffet Process

The IBP places the following prior on  $[\mathbf{Z}]$ , a canonical form of  $\mathbf{Z}$  that is invariant to the ordering of the features (see Griffiths and Ghahramani (2005) for details):

$$p([\mathbf{Z}]) = \frac{\alpha^K}{\prod_{h \in \{0,1\}^N \setminus \mathbf{0}} K_h!} \exp\{-\alpha H_N\} \prod_{k=1}^K \frac{(N - m_k)!(m_k - 1)!}{N!}. \quad (1)$$

Here,  $K$  is the number of nonzero columns in  $\mathbf{Z}$ ,  $m_k$  is the number of ones in column  $k$  of  $\mathbf{Z}$ ,  $H_N$  is the  $N^{\text{th}}$  harmonic number, and  $K_h$  is the number of occurrences of the nonzero binary vector  $h$  among the columns in  $\mathbf{Z}$ . The parameter  $\alpha$  controls the expected number of features present in each observation. For any fixed setting of  $\alpha$ , with probability one, a matrix  $\mathbf{Z}$  drawn from this prior will have a finite number of nonzero entries.

**Restaurant Construction.** The following culinary metaphor is one way to sample a matrix  $\mathbf{Z}$  from the prior described in Equation (1). Imagine the rows of  $\mathbf{Z}$  correspond to customers and the columns correspond to dishes in an infinitely long (Indian) buffet. The customers choose their dishes as follows:

1. The first customer takes the first  $\text{Poisson}(\alpha)$  dishes.
2. The  $i^{\text{th}}$  customer then takes dishes that have been previously sampled with probability  $m_k/i$ , where  $m_k$  is the number of people who have already sampled dish  $k$ . The  $i^{\text{th}}$  customer also takes  $\text{Poisson}(\alpha/i)$  new dishes.

Then,  $z_{nk}$  is one if customer  $n$  tried the  $k^{\text{th}}$  dish and zero otherwise. The resulting matrix is infinitely exchangeable, meaning that the order in which the customers enter the buffet has no impact on the distribution of  $\mathbf{Z}$  (up to permutations of the columns).

The Indian buffet metaphor leads directly to a Gibbs sampler for posterior inference. Bayes' rule states

$$p(z_{nk} | \mathbf{Z}_{-nk}, \mathbf{A}, \mathbf{X}) \propto p(\mathbf{X} | \mathbf{A}, \mathbf{Z}) p(z_{nk} | \mathbf{Z}_{-nk}).$$

The likelihood term  $p(\mathbf{X} | \mathbf{A}, \mathbf{Z})$  is easily computed from the noise model while the prior term  $p(z_{nk} | \mathbf{Z}_{-nk})$  is obtained by imagining that customer  $n$  was the last to enter the restaurant (this assumption is valid due to exchangeability). The prior term  $p(z_{nk} | \mathbf{Z}_{-nk})$  is  $m_k/N$  for active features. New features are sampled by combining the likelihood model with the  $\text{Poisson}(\alpha/N)$  prior on the number of new dishes a customer will try.

If the prior on  $\mathbf{A}$  is conjugate to the likelihood, we can marginalise out  $\mathbf{A}$  from the likelihood  $p(\mathbf{X} | \mathbf{Z}, \mathbf{A})$  and consider  $p(\mathbf{X} | \mathbf{Z})$ . This approach leads to a collapsed Gibbs sampler for the IBP. Marginalising out  $\mathbf{A}$  gives the collapsed Gibbs sampler a level of flexibility that allows it to mix more quickly than an uncollapsed Gibbs sampler.

However, if the likelihood is not conjugate or if the dataset is large and high-dimensional,  $p(\mathbf{X} | \mathbf{Z})$  may be much more expensive to compute than  $p(\mathbf{X} | \mathbf{Z}, \mathbf{A})$ . In these cases, the Gibbs sampler must also sample the feature matrix  $\mathbf{A}$  based on its posterior distribution  $p(\mathbf{A} | \mathbf{X}, \mathbf{Z})$ . For further details and equations on Gibbs samplers for the IBP, please refer to Appendix A.

**Stick-breaking Construction.** The restaurant construction directly lends itself to a Gibbs sampler, but it does not easily lend itself to a variational approach. For the variational approach, we turn to an equivalent alternative construction of the IBP, the stick-breaking construction of Teh et al. (2007). To generate a matrix  $\mathbf{Z}$  using the stick-breaking construction, we begin by assigning a parameter  $\pi_k \in [0, 1]$  to each column of  $\mathbf{Z}$ . Given  $\pi_k$ , each  $z_{nk}$  in column  $k$  is sampled as an independent  $\text{Bernoulli}(\pi_k)$ . Since each 'customer' samples a dish independently of the other customers, this representation makes it clear that the ordering of the customers does not impact the distribution.

The  $\pi_k$  themselves are generated by the following stick-breaking process. We first draw a sequence of independent random variables  $v_1, v_2, \dots$ , each distributed  $\text{Beta}(\alpha, 1)$ . We assign  $\pi_1 = v_1$ . For each subsequent  $k$ , we assign  $\pi_k = v_k \pi_{k-1} = \prod_{i=1}^k v_i$ , resulting in a decreasing sequence of probabilities  $\pi_k$ . In expectation, the probability of seeing feature  $k$  decreases exponentially with  $k$ . The parameter  $\alpha$  affects how quickly these probabilities decrease. Larger values of  $\alpha$

mean that the values  $\pi_k$  decrease more slowly, which means that we expect to see more features in the data. For more details, see (Teh et al., 2007).

### 2.3 Notation

We now summarise the notation which we use throughout the technical report. Vectors or matrices of variables are bold face. A subscript of “ $-i$ ” indicates all components except component  $i$ . A subscript “ $\cdot$ ” indicates all components in a given dimension. For example,  $\mathbf{Z}_{-nk}$  is the full  $\mathbf{Z}$  matrix except the  $(n, k)$  entry, and  $\mathbf{X}_n$  is the entire  $n^{\text{th}}$  row of  $\mathbf{X}$ . For probability distributions, a subscript indicates the parameters used to specify the distribution. For example,  $q_{\boldsymbol{\tau}}(\mathbf{v}) = q(\mathbf{v}; \boldsymbol{\tau})$ .

Commonly recurring variables are:

- $\mathbf{X}$ : The observations are stored in  $\mathbf{X}$ , an  $N \times D$  matrix. The linear-Gaussian model posits  $\mathbf{X} = \mathbf{Z}\mathbf{A} + \boldsymbol{\epsilon}$ , where  $\boldsymbol{\epsilon}$  is a  $N \times D$  matrix of independent elements, each with mean 0 and variance  $\sigma_n^2$ .
- $z_{nk}, \mathbf{Z}$ : Each  $z_{nk}$  indicates whether feature  $k$  is present in observation  $n$ . Here,  $n \in \{1 \dots N\}$  runs over the number of data-points and  $k \in \{1 \dots \infty\}$  runs over the number of features. The matrix  $\mathbf{Z}$  refers to the collection of all  $z_{nk}$ ’s. It has dimensionality  $N \times K$ , where  $K$  is the finite number of nonzero features. All other  $z_{nk}, k > K$  are assumed to be zero. We use  $\alpha$  to denote the concentration parameter of the IBP.
- $\boldsymbol{\pi}$ : The stick lengths (feature probabilities) are  $\pi_k$ .
- $\boldsymbol{\nu}$ : The stick-breaking variables are  $\nu_k$ .
- $\mathbf{A}$ : The collection of Gaussian feature variables, a  $K \times D$  matrix where each feature is represented by the vector  $\mathbf{A}_k$ . In the linear-Gaussian model, the prior states that the elements are  $\mathbf{A}$  are independent with mean 0 and variance  $\sigma_A^2$ .

## 3 Variational Inference

We will focus on variational inference procedures for the linear-Gaussian likelihood model (Griffiths and Ghahramani, 2005), in which  $\mathbf{A}$  and  $\boldsymbol{\epsilon}$  are Gaussian. These updates can be easily adapted to other exponential family likelihood models. As an example, we briefly discuss the variational procedure for the infinite ICA model (Knowles and Ghahramani, 2007) in Appendix E.

We denote the set of hidden variables in the IBP by  $\mathbf{W} = \{\boldsymbol{\pi}, \mathbf{Z}, \mathbf{A}\}$  and the set of parameters by  $\boldsymbol{\theta} = \{\alpha, \sigma_A^2, \sigma_n^2\}$ . Computing the true log posterior

$$\log p(\mathbf{W}|\mathbf{X}, \boldsymbol{\theta}) = \log p(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta}) - \log p(\mathbf{X}|\boldsymbol{\theta})$$

is difficult due to the intractability of computing the log marginal probability  $\log p(\mathbf{X}|\boldsymbol{\theta}) = \log \int p(\mathbf{X}, \mathbf{W}|\boldsymbol{\theta}) d\mathbf{W}$ .

Mean field variational methods approximate the true posterior with a *variational distribution*  $q_{\boldsymbol{\Phi}}(\mathbf{W})$  from some tractable family of distributions  $Q$  (Beal, 2003; Wainwright and Jordan, 2008). Here,  $\boldsymbol{\Phi}$  denotes the set of parameters used to describe the distribution  $q$ . Inference then reduces to performing an optimisation on the parameters  $\boldsymbol{\Phi}$  to find the member  $q \in Q$  that minimises the KL divergence  $D(q_{\boldsymbol{\Phi}}(\mathbf{W})||p(\mathbf{W}|\mathbf{X}, \boldsymbol{\theta}))$ . Since the KL divergence  $D(q||p)$  is nonnegative and equal to zero iff  $p = q$ , the unrestricted solution to our problem is to set

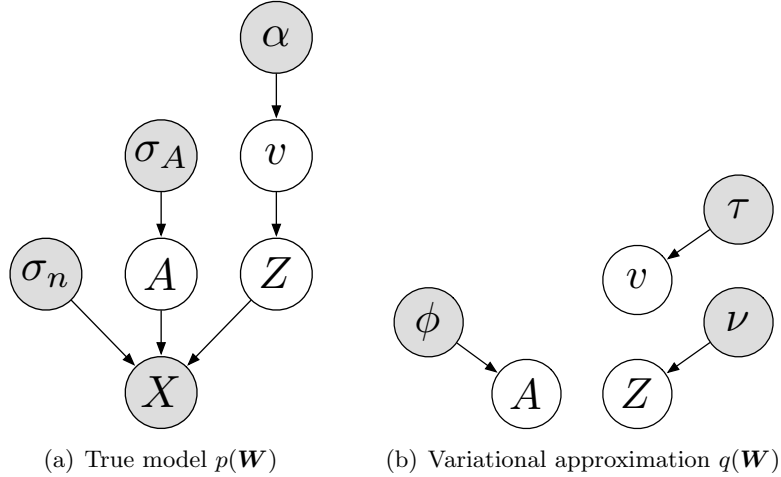


Figure 2: Graphical models for the true posterior distribution of  $\mathbf{W}$  and for our infinite variational approximation.

$q_{\Phi}(\mathbf{W}) = p(\mathbf{W}|\mathbf{X}, \theta)$ . However, this general optimisation problem is intractable. We therefore restrict  $Q$  to a parameterised family of distributions for which this optimisation is tractable.

Specifically, we present two mean field variational approaches with two different families  $Q$ . In both models, we use a truncated model with truncation level  $K$ . A truncation level  $K$  means that  $\mathbf{Z}$  (in our approximating distribution) is nonzero in at most  $K$  columns.

Our first approach minimises the KL-divergence  $D(q||p_K)$  between the variational distribution and a finite approximation  $p_K$  to the IBP described in Section 4; we refer to this approach as the *finite variational* method. In this model, we let  $Q$  be a factorised family

$$q(\mathbf{W}) = q_{\tau}(\boldsymbol{\pi})q_{\phi}(\mathbf{A})q_{\nu}(\mathbf{Z}) \quad (2)$$

where  $\tau$ ,  $\phi$ , and  $\nu$  are optimised to minimise  $D(q||p_K)$ . By minimising the KL-divergence with respect to  $p_K$  and not the true  $p$ , this first approach introduces an additional layer of approximation not present in our second approach. We do not show the graphical model for this approach, but it is equivalent to the one shown in Figure 2 if we replace  $v$  with  $\pi$  and let the true model be the truncated beta-Bernoulli.

Our second approach minimises the KL-divergence to the true IBP posterior  $D(q||p)$ . We call this approach the *infinite variational* method because, while our variational distribution is finite, its updates are based the true IBP posterior (which contains an infinite number of features). In this model, we work directly with the stick-breaking weights  $\mathbf{v}$  instead of directly with  $\boldsymbol{\pi}$ . The family  $Q$  is then the factorised family

$$q(\mathbf{W}) = q_{\tau}(\mathbf{v})q_{\phi}(\mathbf{A})q_{\nu}(\mathbf{Z})$$

where  $\tau$ ,  $\phi$ , and  $\nu$  are the variational parameters. The forms of the distributions  $q$  and the variational updates are specified in Section 5. The graphical model for this variational approximation is shown in Figure 2.

Inference in both approaches consists of optimising the parameters of the approximating distribution to most closely match the true posterior. This optimisation is equivalent to maximising a lower bound on the evidence since

$$\begin{aligned} \log p(\mathbf{X}|\boldsymbol{\theta}) &= \mathbb{E}_q[\log(p(\mathbf{X}, \mathbf{W}|\boldsymbol{\theta}))] + H[q] + D(q||p) \\ &\geq \mathbb{E}_q[\log(p(\mathbf{X}, \mathbf{W}|\boldsymbol{\theta}))] + H[q] \end{aligned} \quad (3)$$

where  $H[q]$  is the entropy of distribution  $q$ , and therefore

$$\arg \min_{\tau, \phi, \nu} D(q||p) = \arg \max_{\tau, \phi, \nu} \mathbb{E}_q[\log(p(\mathbf{X}, \mathbf{W}|\boldsymbol{\theta}))] + H[q]. \quad (4)$$

This optimisation is not convex; in general, we can only hope to find variational parameters that are local optima.

To minimise  $D(q||p)$ , we cycle through each of the variational parameters, and for each one, perform a coordinate ascent that maximises the right side of Equation (4). In doing so, we also improve a lower bound on the log-likelihood of the data.

In Sections 4 and 5, we go over the finite and infinite approaches in detail assuming we have chosen our truncation level  $K$ . Section 6 presents a bound describing how close a truncated stick-breaking approximation of the IBP is to the true IBP, thereby giving a heuristic for how to choose the truncation level  $K$  in both models. This is the first such bound known for the IBP. Appendix B reviews the key concepts for variational inference with exponential family models.

## 4 The Finite Variational Approach

In this section, we introduce our *finite variational approach*, an approximate inference algorithm for an approximation to the IBP. Specifically, we assume that the IBP can be well approximated using the finite beta-Bernoulli model  $p_K$  introduced by (Griffiths and Ghahramani, 2005)

$$\begin{aligned} \pi_k &\sim \text{Beta}(\alpha/K, 1) && \text{for } k \in \{1 \dots K\}, \\ z_{nk} &\sim \text{Bernoulli}(\pi_k) && \text{for } k \in \{1 \dots K\}, n \in \{1 \dots N\}, \\ \mathbf{A}_k &\sim \text{Normal}(0, \sigma_A^2 I) && \text{for } k \in \{1 \dots K\}, \\ \mathbf{X}_n &\sim \text{Normal}(\mathbf{Z}_n \mathbf{A}, \sigma_n^2 I) && \text{for } n \in \{1 \dots N\}, \end{aligned}$$

where  $K$  is some finite (but large) truncation level. Griffiths and Ghahramani (2005) showed that as  $K \rightarrow \infty$ , this finite approximation converges in distribution to the IBP. Under the finite approximation, the joint probability of the data and latent variables is

$$p_K(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta}) = \prod_{k=1}^K \left( p(\pi_k|\alpha) p(\mathbf{A}_k|\sigma_A^2 I) \prod_{n=1}^N p(z_{nk}|\pi_k) \right) \prod_{n=1}^N p(\mathbf{X}_n|\mathbf{Z}_n, \mathbf{A}, \sigma_n^2 I).$$

Even in this simplified finite approximation, working with the log posterior

$$\log p_K(\mathbf{W}|\mathbf{X}, \boldsymbol{\theta}) = \log p_K(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta}) - \log p_K(\mathbf{X}|\boldsymbol{\theta}),$$

is intractable. We will therefore use the following variational distribution as an approximation

$$q(\mathbf{W}) = q_{\tau}(\boldsymbol{\pi}) q_{\phi}(\mathbf{A}) q_{\nu}(\mathbf{Z})$$

where we assume there are at most  $K$  nonzero columns of  $Z$  and

- $q_{\tau_k}(\pi_k) = \text{Beta}(\pi_k; \tau_{k1}, \tau_{k2})$ ,
- $q_{\phi_k}(\mathbf{A}_k) = \text{Normal}(\mathbf{A}_k; \bar{\boldsymbol{\phi}}_k, \boldsymbol{\Phi}_k)$ ,
- $q_{\nu_{nk}}(z_{nk}) = \text{Bernoulli}(z_{nk}; \nu_{nk})$ .



Inference then involves optimising  $\boldsymbol{\tau}$ ,  $\boldsymbol{\phi}$ , and  $\boldsymbol{\nu}$  to either minimise the KL divergence  $D(q||p_K)$  or, equivalently, maximise the following lower bound on  $p_K(\mathbf{X}|\boldsymbol{\theta})$ :

$$\mathbb{E}_q[\log(p_K(\mathbf{X}, \mathbf{W}|\boldsymbol{\theta})) + H[q]].$$

While variational inference with respect to the finite beta-Bernoulli model  $p_K$  is not the same as variational inference with respect to the true IBP posterior, the variational updates are significantly easier and, in the limit of large  $K$ , the finite beta-Bernoulli model is equivalent to the IBP.

#### 4.1 Lower Bound on the Marginal Likelihood

We expand the lower bound in this section, leaving the full set of equations for Appendix C.1. Note that all expectations in this section are taken with respect to the variational distribution  $q$ . We therefore drop the use of  $E_q$  and instead use  $E_{\mathbf{W}}$  to indicate which variables we are taking expectations over. Substituting expressions into Equation (3), our lower bound is

$$\begin{aligned} \log p_K(\mathbf{X}|\boldsymbol{\theta}) &\geq \mathbb{E}_{\mathbf{W}} [\log p_K(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta})] + H[q], \\ &= \sum_{k=1}^K \mathbb{E}_{\boldsymbol{\pi}} [\log p_K(\pi_k|\alpha)] + \sum_{k=1}^K \sum_{n=1}^N \mathbb{E}_{\boldsymbol{\pi}, \mathbf{Z}} [\log p_K(z_{nk}|\pi_k)] \\ &\quad + \sum_{k=1}^K \mathbb{E}_{\mathbf{A}} [\log p_K(\mathbf{A}_{k\cdot}|\sigma_A^2 I)] + \sum_{n=1}^N \mathbb{E}_{\mathbf{Z}, \mathbf{A}} [\log p_K(\mathbf{X}_n|\mathbf{Z}_n, \mathbf{A}, \sigma_n^2 I)] + H[q]. \end{aligned} \quad (5)$$

Evaluating these expectations are all straightforward exponential family calculations. Expanding each of these, the full lower bound is

$$\begin{aligned} \log p_K(\mathbf{X}|\boldsymbol{\theta}) &\geq \sum_{k=1}^K \left[ \log \frac{\alpha}{K} + \left( \frac{\alpha}{K} - 1 \right) (\psi(\tau_{k1}) - \psi(\tau_{k1} + \tau_{k2})) \right] \\ &\quad + \sum_{k=1}^K \sum_{n=1}^N [\nu_{nk} \psi(\tau_{k1}) + (1 - \nu_{nk}) \psi(\tau_{k2}) - \psi(\tau_{k1} + \tau_{k2})] \\ &\quad + \sum_{k=1}^K \left[ \frac{-D}{2} \log(2\pi\sigma_A^2) - \frac{1}{2\sigma_A^2} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_k^T) \right] \\ &\quad + \sum_{n=1}^N \left[ -\frac{D}{2} \log(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} \left( \mathbf{X}_n \mathbf{X}_n^T - 2 \sum_{k=1}^K \nu_{nk} \bar{\boldsymbol{\phi}}_k \mathbf{X}_n^T + 2 \sum_{k < k'} \nu_{nk} \nu_{nk'} \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_{k'}^T + \sum_{k=1}^K \nu_{nk} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_k^T) \right) \right] \\ &\quad + \sum_{k=1}^K \left[ \log \left( \frac{\Gamma(\tau_{k1}) \Gamma(\tau_{k2})}{\Gamma(\tau_{k1} + \tau_{k2})} \right) - (\tau_{k1} - 1) \psi(\tau_{k1}) - (\tau_{k2} - 1) \psi(\tau_{k2}) + (\tau_{k1} + \tau_{k2} - 2) \psi(\tau_{k1} + \tau_{k2}) \right] \\ &\quad + \sum_{k=1}^K \left[ \frac{1}{2} \log((2\pi e)^D |\Phi_k|) \right] + \sum_{k=1}^K \sum_{n=1}^N [-\nu_{nk} \log \nu_{nk} - (1 - \nu_{nk}) \log(1 - \nu_{nk})]. \end{aligned} \quad (6)$$

where  $\psi(\cdot)$  is the digamma function. Full derivations are left to Appendix C.1.

#### 4.2 Parameter Updates

Our approach to variational inference involves cycling through each of the variational parameters and sequentially updating them using standard exponential family variational update

equations. High level details of this approach can be found in Appendix B, Blei and Jordan (2004), and Wainwright and Jordan (2008). Derivations of the following update equations are in Appendix C.2.

1. For  $k = 1, \dots, K$ , we update the  $\bar{\phi}_k$  and  $\Phi_k$  in  $\text{Normal}(\mathbf{A}_k; \bar{\phi}_k, \Phi_k)$  as

$$\begin{aligned}\Phi_k &= \left( \frac{1}{\sigma_A^2} + \frac{\sum_{n=1}^N \nu_{nk}}{\sigma_n^2} \right)^{-1} I \\ \bar{\phi}_k &= \left[ \frac{1}{\sigma_n^2} \sum_{n=1}^N \nu_{nk} \left( \mathbf{X}_{n\cdot} - \left( \sum_{l:l \neq k} \nu_{nl} \bar{\phi}_l \right) \right) \right] \left( \frac{1}{\sigma_A^2} + \frac{\sum_{n=1}^N \nu_{nk}}{\sigma_n^2} \right)^{-1}.\end{aligned}$$

2. For  $k = 1, \dots, K$ ,  $n = 1, \dots, N$ , update  $\nu_{nk}$  in  $\text{Bernoulli}(z_{nk}; \nu_{nk})$  as

$$\nu_{nk} = \frac{1}{1 + e^{-\vartheta}}.$$

where

$$\vartheta = \psi(\tau_{k1}) - \psi(\tau_{k2}) - \frac{1}{2\sigma_n^2} \left( \text{tr}(\Phi_k) + \bar{\phi}_k \bar{\phi}_k^T \right) + \frac{1}{\sigma_n^2} \bar{\phi}_k \left( \mathbf{X}_{n\cdot}^T - \left( \sum_{l:l \neq k} \nu_{nl} \bar{\phi}_l^T \right) \right)$$

3. For  $k = 1, \dots, K$ , we update the  $\tau_{k1}$  and  $\tau_{k2}$  in  $\text{Beta}(\pi_k; \tau_{k1}, \tau_{k2})$  as

$$\begin{aligned}\tau_{k1} &= \frac{\alpha}{K} + \sum_{n=1}^N \nu_{nk}, \\ \tau_{k2} &= N + 1 - \sum_{n=1}^N \nu_{nk}.\end{aligned}$$

## 5 The Infinite Variational Approach

In this section, we introduce the *infinite variational approach*, a method for doing approximate inference for the linear-Gaussian model with respect to a full IBP prior. The model for  $p$  is the full (untruncated) stick-breaking construction for the IBP:

$$\begin{aligned}v_k &\sim \text{Beta}(\alpha, 1) && \text{for } k \in \{1, \dots, \infty\}, \\ \pi_k &= \prod_{i=1}^k v_i && \text{for } k \in \{1 \dots \infty\}, \\ z_{nk} &\sim \text{Bernoulli}(\pi_k) && \text{for } k \in \{1 \dots \infty\}, n \in \{1 \dots N\}, \\ \mathbf{A}_k &\sim \text{Normal}(0, \sigma_A^2 I) && \text{for } k \in \{1 \dots \infty\}, \\ \mathbf{X}_{n\cdot} &\sim \text{Normal}(\mathbf{Z}_n \mathbf{A}, \sigma_n^2 I) && \text{for } n \in \{1 \dots N\}.\end{aligned}$$

The joint probability of the data and variables is

$$p(\mathbf{W}, \mathbf{X} | \boldsymbol{\theta}) = \prod_{k=1}^{\infty} \left( p(\pi_k | \alpha) p(\mathbf{A}_k | \sigma_A^2 I) \prod_{n=1}^N p(z_{nk} | \pi_k) \right) \prod_{n=1}^N p(\mathbf{X}_{n\cdot} | \mathbf{Z}_n, \mathbf{A}, \sigma_n^2 I).$$

Working with the log posterior

$$\log p(\mathbf{W}|\mathbf{X}, \boldsymbol{\theta}) = \log p(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta}) - \log p(\mathbf{X}|\boldsymbol{\theta}).$$

is again intractable, so we use a variational approximation. Similar to the approach used by Blei and Jordan (2004), our variational approach uses a truncated stick-breaking process that truncates the number of nonzero columns to a finite value  $K$ . In the truncated stick-breaking process,  $\pi_k = \prod_{i=1}^k v_i$  for  $k \leq K$  and zero otherwise.

The ordered feature probabilities  $\{\pi_1 \dots \pi_K\}$  are dependent in this model, while the  $\{v_1 \dots v_K\}$  are independent. We therefore use  $\mathbf{v}$  instead of  $\boldsymbol{\pi}$  as our hidden variable because it is simpler to work with and given  $\mathbf{v}$ , it is straightforward to calculate  $\boldsymbol{\pi}$ . Our mean field variational distribution is then:

$$q(\mathbf{W}) = q_{\boldsymbol{\tau}}(\mathbf{v})q_{\boldsymbol{\phi}}(\mathbf{A})q_{\boldsymbol{\nu}}(\mathbf{Z})$$

where

- $q_{\boldsymbol{\tau}_k}(v_k) = \text{Beta}(v_k; \tau_{k1}, \tau_{k2})$ ,
- $q_{\boldsymbol{\phi}_k}(\mathbf{A}_{k\cdot}) = \text{Normal}(\mathbf{A}_{k\cdot}; \bar{\boldsymbol{\phi}}_k, \Phi_k)$ ,
- $q_{\nu_{nk}}(z_{nk}) = \text{Bernoulli}(z_{nk}; \nu_{nk})$ .

As with the finite approach, inference involves optimising  $\boldsymbol{\tau}$ ,  $\boldsymbol{\phi}$ , and  $\boldsymbol{\nu}$  to minimise the KL divergence  $D(q||p)$ , or equivalently to maximise the lower bound on  $p(\mathbf{X}|\boldsymbol{\theta})$

$$\mathbb{E}_q[\log p(\mathbf{X}, \mathbf{W}|\boldsymbol{\theta})] + H[q].$$

Unfortunately, the update equations for this approximation are not as straightforward as in the finite approach.

## 5.1 Lower Bound on the Marginal Likelihood

As in the finite approach, we first derive an expression for the variational lower bound. However, parts of our model are no longer in the exponential family and require nontrivial computations. We expand upon these parts here, leaving the straightforward exponential family calculations to Appendix D.1.

The lower bound on  $p(\mathbf{X}|\boldsymbol{\theta})$  can be decomposed as follows

$$\begin{aligned} \log p(\mathbf{X}|\boldsymbol{\theta}) \geq & \sum_{k=1}^K \mathbb{E}_{\mathbf{v}} [\log p(v_k|\boldsymbol{\alpha})] + \sum_{k=1}^K \sum_{n=1}^N \mathbb{E}_{\mathbf{v}, \mathbf{Z}} [\log p(Z_{nk}|\mathbf{v})] \\ & + \sum_{k=1}^K \mathbb{E}_{\mathbf{A}} [\log p(\mathbf{A}_{k\cdot}|\sigma_A^2)] + \sum_{n=1}^N \mathbb{E}_{\mathbf{Z}, \mathbf{A}} [\log p(\mathbf{X}_{n\cdot}|\mathbf{Z}, \mathbf{A}, \sigma_n^2)] + H[q], \end{aligned} \quad (7)$$

Except for the second term, all of the terms are exponential family calculations; evaluated they

come out to

$$\begin{aligned}
& \log p(\mathbf{X}|\boldsymbol{\theta}) \\
& \geq \sum_{k=1}^K [\log \alpha + (\alpha - 1) (\psi(\tau_{k1}) - \psi(\tau_{k1} + \tau_{k2}))] \\
& \quad + \sum_{k=1}^K \sum_{n=1}^N \left[ \nu_{nk} \left( \sum_{m=1}^k \psi(\tau_{k2}) - \psi(\tau_{k1} + \tau_{k2}) \right) + (1 - \nu_{nk}) \mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right] \right] \\
& \quad + \sum_{k=1}^K \left[ -\frac{D}{2} \log(2\pi\sigma_A^2) - \frac{1}{2\sigma_A^2} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\Phi}}_k \bar{\boldsymbol{\Phi}}_k^T) \right] \\
& \quad + \sum_{n=1}^N \left[ -\frac{D}{2} \log(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} \left( \mathbf{X}_n \cdot \mathbf{X}_n^T - 2 \sum_{k=1}^K \nu_{nk} \bar{\boldsymbol{\Phi}}_k \mathbf{X}_n^T + 2 \sum_{k < k'} \nu_{nk} \nu_{nk'} \bar{\boldsymbol{\Phi}}_k \bar{\boldsymbol{\Phi}}_{k'}^T + \sum_{k=1}^K \nu_{nk} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\Phi}}_k \bar{\boldsymbol{\Phi}}_k^T) \right) \right] \\
& \quad + \sum_{k=1}^K \left[ \log \left( \frac{\Gamma(\tau_{k1}) \Gamma(\tau_{k2})}{\Gamma(\tau_{k1} + \tau_{k2})} \right) - (\tau_{k1} - 1) \psi(\tau_{k1}) - (\tau_{k2} - 1) \psi(\tau_{k2}) + (\tau_{k1} + \tau_{k2} - 2) \psi(\tau_{k1} + \tau_{k2}) \right] \\
& \quad + \sum_{k=1}^K \frac{1}{2} \log((2\pi e)^D |\boldsymbol{\Phi}_k|) + \sum_{k=1}^K \sum_{n=1}^N [-\nu_{nk} \log \nu_{nk} - (1 - \nu_{nk}) \log(1 - \nu_{nk})]
\end{aligned} \tag{8}$$

where  $\psi(\cdot)$  is the digamma function, and we have left  $\mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right]$ , a byproduct of the expectation of  $\mathbb{E}_{\mathbf{v}, \mathbf{Z}} [\log p(Z_{nk}|\mathbf{v})]$ , unevaluated. This expectation has no closed-form solution, so we instead lower bound it (and therefore lower bound the log posterior).

In this section, we present a multinomial approximation which leads to a computationally efficient lower bound and straightforward parameter updates.<sup>1</sup> An approach based on a Taylor series expansion is presented in Appendix D.1. Unlike the multinomial approximation, the Taylor approximation can be made arbitrarily precise; however, empirically we find that the multinomial bound is usually only 2-10% looser than a fifty-term Taylor series expansion—and about thirty times faster to compute. Also, parameter updates under the Taylor approximation do not have a closed form solution and must be numerically optimised. Thus, we recommend using the multinomial approximation and the corresponding parameter updates; the Taylor derivation is provided in the Appendix D.1 largely for reference.

To bound  $\mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right]$  with the multinomial approximation, we introduce an auxiliary  $k$ -multinomial distribution  $q_k = (q_{k1}, q_{k2}, \dots, q_{kk})$  into the expectation and apply Jensen's inequality:

$$\begin{aligned}
\mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right] &= \mathbb{E}_{\mathbf{v}} \left[ \log \left( \sum_{y=1}^k (1 - v_y) \prod_{m=1}^{y-1} v_m \right) \right] \\
&= \mathbb{E}_{\mathbf{v}} \left[ \log \left( \sum_{y=1}^k q_{ky} \frac{(1 - v_y) \prod_{m=1}^{y-1} v_m}{q_{ky}} \right) \right] \\
&\geq \mathbb{E}_y \mathbb{E}_{\mathbf{v}} \left[ \log(1 - v_y) + \sum_{m=1}^{y-1} \log v_m \right] + H[q_k] \\
&= \mathbb{E}_y \left[ \psi(\tau_{y2}) + \left( \sum_{m=1}^{y-1} \psi(\tau_{m1}) \right) - \left( \sum_{m=1}^y \psi(\tau_{m1} + \tau_{m2}) \right) \right] + H[q_k].
\end{aligned}$$

<sup>1</sup>Note that Jensen's inequality cannot be used here; the concavity of the log goes in the wrong direction.

Explicitly writing out  $q_k$ , we get

$$\mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right] \geq \left( \sum_{m=1}^k q_{km} \psi(\tau_{m2}) \right) + \left( \sum_{m=1}^{k-1} \left( \sum_{n=m+1}^k q_{kn} \right) \psi(\tau_{m1}) \right) \quad (9)$$

$$- \left( \sum_{m=1}^k \left( \sum_{n=m}^k q_{kn} \right) \psi(\tau_{m1} + \tau_{m2}) \right) - \sum_{m=1}^k q_{km} \log q_{km}.$$

Since Equation (9) holds for any  $q_{k1}, \dots, q_{kk}$  for all  $1 \leq k \leq K$ , we now optimise  $q_k$  to maximise the lower bound. Taking derivatives with respect to each  $q_{ki}$  and introducing  $\lambda$  as the Lagrangian variable to ensure that  $q_k$  is a distribution, we get

$$0 = \psi(\tau_{i2}) + \sum_{m=1}^{i-1} \psi(\tau_{m1}) - \sum_{m=1}^i \psi(\tau_{m1} + \tau_{m2}) - 1 - \log(q_{ki}) - \lambda$$

Solving for  $q_{ki}$ , we find

$$q_{ki} \propto \exp \left( \psi(\tau_{i2}) + \sum_{m=1}^{i-1} \psi(\tau_{m1}) - \sum_{m=1}^i \psi(\tau_{m1} + \tau_{m2}) \right) \quad (10)$$

where the proportionality ensures that  $q_k$  a valid distribution. If we plug this multinomial lower bound back into Equation (8), we get a lower bound on  $\log p(\mathbf{X}|\boldsymbol{\theta})$ . We then optimise the remaining parameters to maximise the resulting lower bound as described in the next section.

The auxiliary distribution  $q_k$  is largely a computational tool, but it does have the following intuition. Since  $\pi_k = \prod_{i=1}^k v_i$ ; we can imagine the event  $z_{nk} = 1$  is equivalent to the event that a series of variables  $u_i \sim \text{Bernoulli}(v_i)$  all flip to one. If any of the  $u_i$ 's equal zero, then the feature is off. The multinomial distribution  $q_{kj}$  can be thought of as a distribution over the event that the  $j^{\text{th}}$  variable  $u_j$  is the first  $u_i$  to equal 0.

## 5.2 Parameter Updates

The updates for the variational parameters for  $\mathbf{A}$  and  $\mathbf{Z}$  are still in the exponential family. For the parameters of  $\mathbf{A}$ , the updates are identical to those of the finite model. For the parameters of  $\mathbf{Z}$ , the updates are again similar to the finite model, except we must use an approximation for  $\mathbb{E}_{\mathbf{v}}[\log(1 - \prod_{i=1}^k v_i)]$ .

The updates for the parameters for  $\mathbf{v}$ , however, strongly depend on how we approximate the term  $\mathbb{E}_{\mathbf{v}}[\log(1 - \prod_{i=1}^k v_i)]$ . If we use the multinomial lower bound of Section 5.1, the updates have a nice closed form<sup>2</sup>. As in the finite approach, we sequentially update each of the variational parameters in turn. Derivations of the following update equations are in Appendix D.2.

1. For  $k = 1, \dots, K$ , we update the  $\bar{\boldsymbol{\phi}}_k$  and  $\boldsymbol{\Phi}_k$  in  $\text{Normal}(\mathbf{A}_k; \bar{\boldsymbol{\phi}}_k, \boldsymbol{\Phi}_k)$  as

$$\boldsymbol{\Phi}_k = \left( \frac{1}{\sigma_A^2} + \frac{\sum_{n=1}^N \nu_{nk}}{\sigma_n^2} \right)^{-1} \mathbf{I}$$

$$\bar{\boldsymbol{\phi}}_k = \left[ \frac{1}{\sigma_n^2} \sum_{n=1}^N \nu_{nk} \left( \mathbf{X}_n - \left( \sum_{l:l \neq k} \nu_{nl} \bar{\boldsymbol{\phi}}_l \right) \right) \right] \left( \frac{1}{\sigma_A^2} + \frac{\sum_{n=1}^N \nu_{nk}}{\sigma_n^2} \right)^{-1}.$$

<sup>2</sup>Appendix D.2 describes an alternative approach that numerically optimises the variational lower bound based on our Taylor series approximation; however, we found the direct optimisation was less computationally efficient.

2. For  $k = 1, \dots, K$ ,  $n = 1, \dots, N$ , update  $\nu_{nk}$  in  $\text{Bernoulli}(z_{nk}; \nu_{nk})$  as

$$\nu_{nk} = \frac{1}{1 + e^{-\vartheta}}$$

where

$$\begin{aligned} \vartheta = & \sum_{i=1}^k (\psi(\tau_{i1}) - \psi(\tau_{i1} + \tau_{i2})) - \mathbb{E}_{\mathbf{v}}[\log(1 - \prod_{i=1}^k v_i)] \\ & - \frac{1}{2\sigma_n^2} (\text{tr}(\mathbf{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T) + \frac{1}{\sigma_n^2} \bar{\phi}_k \left( \mathbf{X}_n^T - \left( \sum_{l:l \neq k} \nu_{nl} \bar{\phi}_l^T \right) \right). \end{aligned}$$

We leave the term  $\mathbb{E}_{\mathbf{v}}[\log(1 - \prod_{i=1}^k v_i)]$  unevaluated because the choice of how to approximate it does not change the form of the update. In practice, we use the multinomial lower bound for this term.

3. For  $k = 1, \dots, K$ , we must update the  $\tau_{k1}$  and  $\tau_{k2}$  in  $\text{Beta}(v_k; \tau_{k1}, \tau_{k2})$ . If we use the multinomial lower bound for  $\mathbb{E}_{\mathbf{v}}[\log(1 - \prod_{i=1}^k v_i)]$ , then we can first compute  $q_{ki}$  according to Equation (10) (remembering to normalise  $q_{ki}$ ). Then the updates for  $\tau_{k1}$  and  $\tau_{k2}$  have the following closed form

$$\begin{aligned} \tau_{k1} &= \alpha + \sum_{m=k}^K \sum_{n=1}^N \nu_{nm} + \sum_{m=k+1}^K \left( N - \sum_{n=1}^N \nu_{nm} \right) \binom{m}{i=k+1} \\ \tau_{k2} &= 1 + \sum_{m=k}^K \left( N - \sum_{n=1}^N \nu_{nm} \right) q_{mk}. \end{aligned}$$

## 6 Bound for the Infinite Approximation

Both of our variational inference approaches require us to choose a truncation level  $K$  for our variational distribution. Building on results from Thibaux and Jordan (2007) and Teh et al. (2007), we present a bound on how close the marginal distribution of the data  $\mathbf{X}$  using a truncated stick-breaking prior will be to the marginal distribution using the true IBP stick-breaking prior. The bound can serve as a rough guide for choosing  $K$ , though the results do not tell us how good our variational approximations will be.

Our development parallels a bound for the Dirichlet Process by Ishwaran and James (2001) and presents the first such truncation bound for the IBP. Let us denote the marginal distribution of observation  $\mathbf{X}$  by  $m_{\infty}(\mathbf{X})$  when we integrate  $\mathbf{W}$  with respect to the true IBP stick-breaking prior  $p(\mathbf{W}|\boldsymbol{\theta})$ . Let  $m_K(\mathbf{X})$  be the marginal distribution when  $\mathbf{W}$  are integrated out with respect to the truncated stick-breaking prior with truncation level  $K$  as described at the beginning of Section 5. For consistency, we continue to use the notation from the linear-Gaussian model, but the derivation that follows is independent of the likelihood model.

Intuitively, the error in the truncation will depend on the probability that, given  $N$  observations, we observe more than  $K$  features in the data (otherwise the truncation should have no effect). Using the beta process representation for the IBP (Thibaux and Jordan, 2007) and using an analysis similar to the one in (Ishwaran and James, 2001), we can show that the difference

between the marginal distributions of  $\mathbf{X}$  is at most

$$\begin{aligned}
\frac{1}{4} \int |m_K(\mathbf{X}) - m_\infty(\mathbf{X})| d\mathbf{X} &\leq \Pr(\exists k > K, n \text{ with } z_{nk} = 1) \\
&= 1 - \Pr(\text{all } z_{ik} = 0, i \in \{1, \dots, N\}, k > K) \\
&= 1 - \mathbb{E} \left[ \left( \prod_{i=K+1}^{\infty} (1 - \pi_i) \right)^N \right] \\
&\leq 1 - \left( \mathbb{E} \left[ \prod_{i=K+1}^{\infty} (1 - \pi_i) \right] \right)^N. \tag{11}
\end{aligned}$$

We begin the derivation of the formal truncation bound by noting that beta process construction for the IBP (Thibaux and Jordan, 2007) implies that the sequence of  $\pi_1, \pi_2, \dots$  may be modeled as a Poisson process on the unit interval  $[0, 1]$  with rate  $\mu(x) = \alpha x^{-1} dx$ . It follows that the sequence of  $\pi_{K+1}, \pi_{K+2}, \dots$  may be modeled as a Poisson process on the interval  $[0, \pi_K]$  with the same rate. The Levy-Khintchine formula (Applebaum, 2004) states that the moment generating function of a Poisson process  $X$  with rate  $\mu$  can be written as

$$\mathbb{E}[\exp(tf(X))] = \exp \left( \int (\exp(tf(y)) - 1) \mu(y) dy \right).$$

where we use  $f(X)$  to denote  $\sum_{x \in X} f(x)$ .

Returning to Equation (11), if we rewrite the final expectation as

$$\mathbb{E} \left[ \left( \prod_{i=K+1}^{\infty} (1 - \pi_i) \right)^N \right] = \mathbb{E} \left[ \exp \left( \sum_{i=K+1}^{\infty} \log(1 - \pi_i) \right)^N \right],$$

then we can apply the Levy-Khintchine formula to get

$$\begin{aligned}
\mathbb{E} \left[ \exp \left( \sum_{i=K+1}^{\infty} \log(1 - \pi_i) \right)^N \right] &= \mathbb{E}_{\pi_K} \left[ \exp \left( \int_0^{\pi_K} (\exp(\log(1 - x)) - 1) \mu(x) dx \right)^N \right] \\
&= \mathbb{E}_{\pi_K} [\exp(-\alpha \pi_K)].
\end{aligned}$$

Finally, we apply Jensen's inequality, using the fact that  $\pi_K$  is the product of independent Beta( $\alpha, 1$ ) variables:

$$\begin{aligned}
\mathbb{E}_{\pi_K} [\exp(-\alpha \pi_K)] &\geq \exp(\mathbb{E}_{\pi_K} [-\alpha \pi_K]) \\
&= \exp \left( -\alpha \left( \frac{\alpha}{1 + \alpha} \right)^K \right).
\end{aligned}$$

Substituting this expression back into Equation (11) gives us the bound

$$\frac{1}{4} \int |m_K(X) - m_\infty(X)| dX \leq 1 - \exp \left( -N \alpha \left( \frac{\alpha}{1 + \alpha} \right)^K \right). \tag{12}$$

Similar to truncation bound for the Dirichlet Process, the expected error increases as  $N$  and  $\alpha$  – the factors that increase the expected number of features – increase. However, for fixed  $N$  and  $\alpha$ , the bound decreases exponentially quickly as the truncation level  $K$  is increased.

Figure 3 shows our truncation bound and the true  $L_1$  distance based on 1000 Monte Carlo simulations of an IBP matrix with  $N = 30$  observations and  $\alpha = 5$ . As expected, the bound

decreases exponentially fast with the truncation level  $K$ . The bound is loose, however; in practice, we find that a heuristic approximation to the bound using a Taylor series expansion provides tighter estimates of the loss. Appendix F describes both this heuristic bound and other (principled) bounds that can be derived via other applications of Jensen’s inequality. All of these provide guidance for how to choose  $K$  such that the truncated stick-breaking construction is close to the true stick-breaking construction, but do not tell us how close the truncated variational approximation is to the true posterior. For this, we use the bound merely as a heuristic.

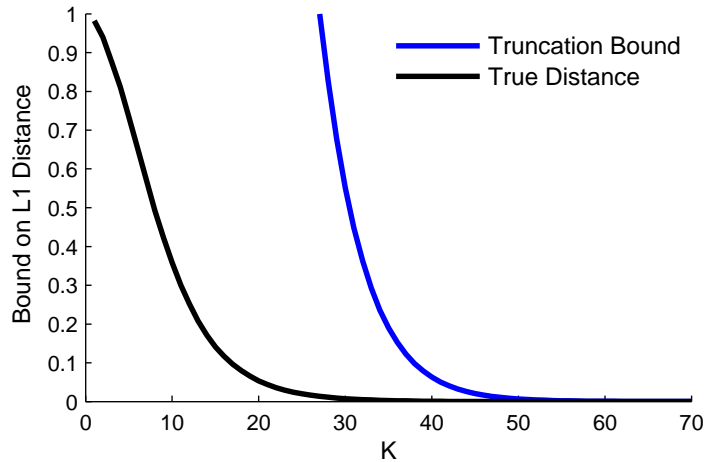


Figure 3: Truncation bound and true  $L_1$  distance.

## 7 Experiments

We compared our variational approaches with both Gibbs sampling (Griffiths and Ghahramani, 2005) and particle filtering (Wood and Griffiths, 2007). As variational algorithms are only guaranteed to converge to a *local* optimum, we applied standard optimisation tricks to avoid small minima. Each run was given a number of random restarts and the hyperparameters for the noise and feature variance were annealed to smooth the posterior. We also experimented with several other techniques such as gradually introducing data and merging correlated features. The latter techniques proved less useful as the size and dimensionality of the datasets increased; they were not included in the final experiments.

We compared against two variants of the Gibbs sampling method introduced in Griffiths and Ghahramani (2005). The first is the collapsed Gibbs sampler described by (Griffiths and Ghahramani, 2005), the second is a partially-uncollapsed alternative in which instantiated features are explicitly represented and new features are integrated out. We compare against these two methods because the time complexity of the collapsed sampler scales cubically, so as we get more data, the cost of this sampler becomes prohibitive. The uncollapsed sampler scales much better in terms of computation per iteration as we get more data, but in general does not explore the posterior as well as the collapsed sampler.

In contrast to the variational methods, the number of features present in the IBP matrix will adaptively grow or shrink in the samplers. To provide a fair comparison with the variational approaches, we also tested finite variants of the collapsed and uncollapsed Gibbs samplers in which the number of potential features was fixed. Details for these samplers are given in Appendix A. We also tested against the particle filter of Wood and Griffiths (2007). All sampling



methods were annealed and given an equal number of restarts as the variational methods.

Both the variational and Gibbs sampling algorithms were heavily optimised for efficient matrix computation so we could evaluate the algorithms both on their running times and the quality of the inference. For the particle filter, we used the implementation provided by Wood and Griffiths (2007). To measure the quality of these methods, we held out one third of the observations on the last half of the dataset. Once the inference was complete, we computed the predictive likelihood of the held out data (averaged over restarts).

## 7.1 Synthetic Data

The synthetic datasets consisted of  $\mathbf{Z}$  and  $\mathbf{A}$  matrices randomly generated from the truncated stick-breaking prior. Figure 4 shows the evolution of the test-likelihood over a thirty minute interval for a dataset with 500 observations of 500 dimensions and with 20 latent features. The error bars indicate the variation over the 5 random starts<sup>3</sup>. The finite uncollapsed Gibbs sampler (dotted green) rises quickly but consistently gets caught in lower optima and has higher variance than the variational approach. Examining the individual runs, we found the higher variance was not due to the Gibbs sampler mixing but due to each run getting stuck in widely varying local optima. The variational methods were slightly slower per iteration but soon found regions of higher predictive likelihoods. The remaining samplers were much slower per iteration, often failing to mix within the allotted interval.

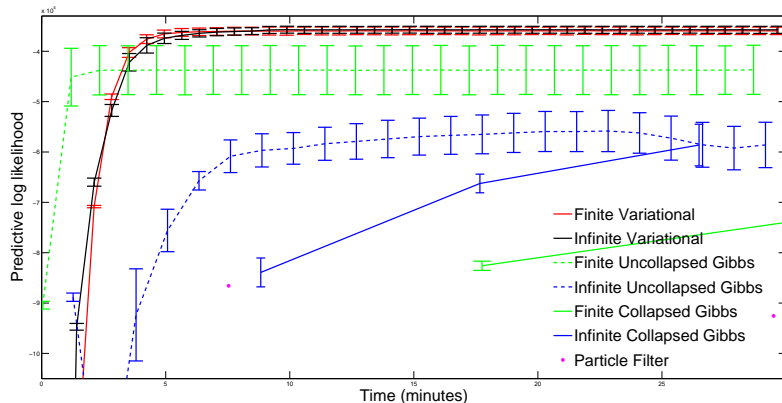


Figure 4: Evolution of test log-likelihoods over a thirty-minute interval for  $N = 500$ ,  $D = 500$ , and  $K = 20$ . The finite uncollapsed Gibbs sampler has the fastest rise but gets caught in worse lower optima than the variational approach.

Figure 5 shows a similar plot for a smaller dataset with  $N = 100$ . Here, the variational approaches do not do as well as the collapsed Gibbs samplers at finding regions of high probability. However, it does as well or better than the uncollapsed samplers. We believe this is because in a smaller dataset, the collapsed Gibbs samplers mix fairly quickly and can therefore find regions of high probability mass. The uncollapsed samplers still get stuck in widely varying local optima as witnessed by the large variance again. The variational approach still converges quickly to optima that on average are better than the uncollapsed samplers, but are not as good the optima the collapsed samplers eventually find.

<sup>3</sup>The particle filter must be run to completion before making prediction, so we cannot test its predictive performance over time. We instead plot the test likelihood only at the end of the inference for particle filters with 10 and 50 particles (the two magenta points in the bottom half of the graph).

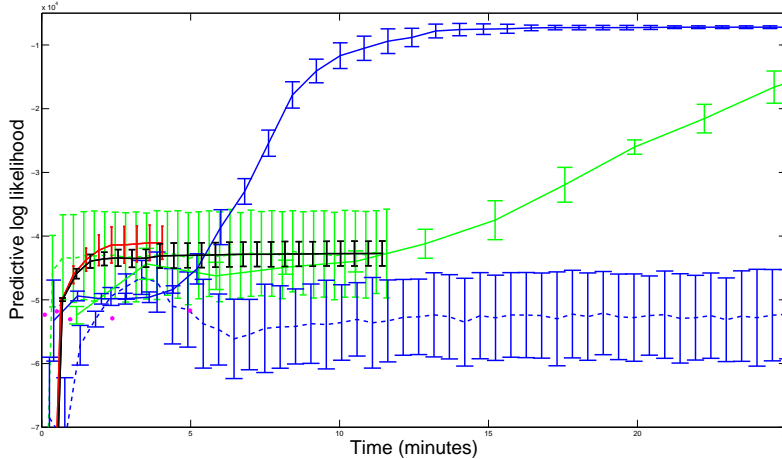


Figure 5: Evolution of test log-likelihoods over a thirty-minute interval for  $N = 100$ ,  $D = 500$ , and  $K = 25$ . For smaller  $N$ , the Gibbs sampler does better at finding an optima of high probability mass.

Figures 6 and 7 show results from a systematic series of tests in which we tested all combinations of observation count  $N = \{5, 10, 50, 100, 500, 1000\}$ , dimensionality  $D = \{5, 10, 50, 100, 500, 1000\}$ , and truncation level  $K = \{5, 10, 15, 20, 25\}$ . Each of the samplers was run for 1000 iterations on three chains and the particle filter was run with 500 particles. For the variational methods, we used a stopping criterion that halted the optimisation when the variational lower bound between the current and previous iterations changed by a multiplicative factor of less than  $10^{-6}$  and the annealing process had completed.

Figure 6 shows how the computation time scales with the truncation level. The variational approaches and the uncollapsed Gibbs sampler are consistently an order of magnitude faster than other algorithms. Figure 7 shows the interplay between dimensionality, computation time, and test log-likelihood for datasets of size  $N = 5$  and  $N = 1000$  respectively. For  $N = 1000$ , the collapsed Gibbs samplers and particle filter did not finish, so they do not appear on the plot. We chose  $K = 20$  as a representative truncation level. Each line represents increasing dimensionality for a particular method (the large dot indicates  $D = 5$ , the subsequent dots correspond to  $D = 10, 50$ , etc.). The nearly vertical lines of the variational methods show that they are quite robust to increasing dimension. Moreover, as dimensionality and dataset size increase, the variational methods become increasingly faster than the samplers. By comparing the lines across the likelihood dimension, we see that for the very small dataset, the variational method often has a lower test log-likelihood than the samplers. In this regime, the samplers mix quickly and find regions of high likelihood.

The short summary of these results is that for smaller data sets (both smaller  $D$  and  $N$ ), the collapsed samplers are preferable. They mix well and find configurations that have high predictive likelihood. However, as the size of the data grows (both  $D$  and  $N$ ), these samplers do not scale well and their cost becomes prohibitive. In these regimes, our variational approximations are preferable over the uncollapsed samplers. Both converge to local optima quickly, but the variational approximations finds better local optima.

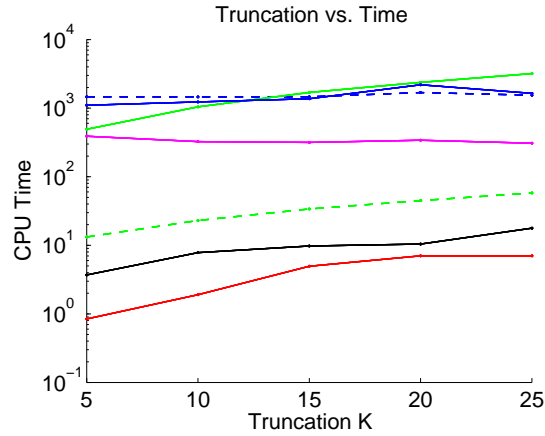


Figure 6: Time versus truncation ( $K$ ). The variational approaches are generally orders of magnitude faster than the samplers (note log scale on the time axis).

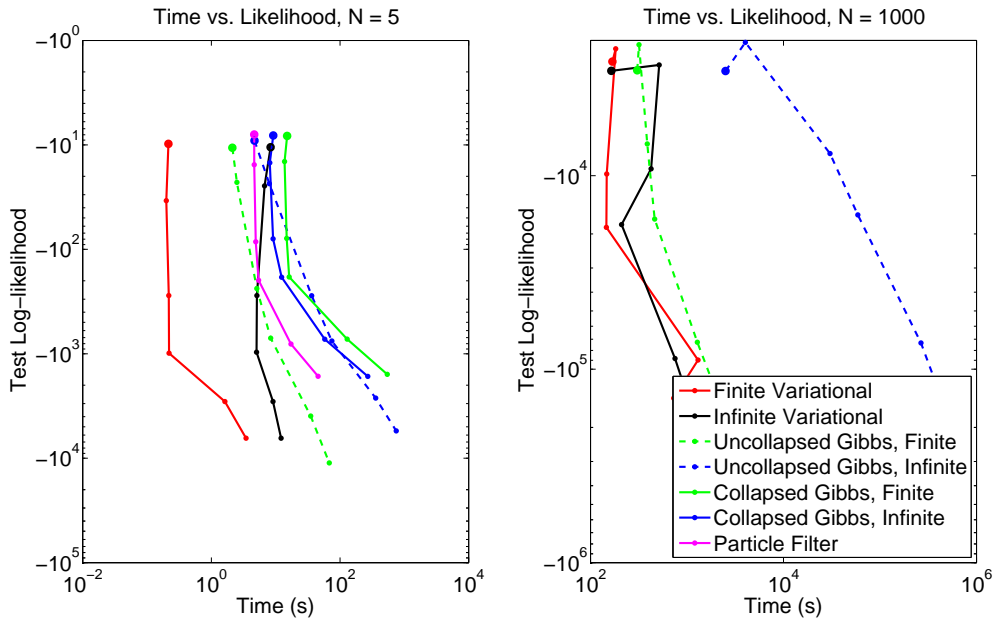


Figure 7: Time versus log-likelihood plot for  $K = 20$ . The larger dots correspond to  $D = 5$  the smaller dots to  $D = 10, 50, 100, 500, 1000$ .

## 7.2 Real Data

To test the conclusions we arrived at based on synthetic data, we applied our variational method to two real-world datasets to see how it would fare with complex, noisy data not drawn from the IBP prior<sup>4</sup>. The first dataset, the Yale Faces (Georghiadis et al., 2001) dataset, consisted of  $N = 721$  32x32 pixel frontal-face images ( $D = 1024$ ) of 14 people with varying expressions and lighting conditions. We set  $\sigma_a$  and  $\sigma_n$  based on the variance of the data. The second dataset, a speech dataset, consisted of  $N = 245$  observations sampled from a 10-microphone audio recording ( $D = 10$ ) of 5 different speakers. We applied the iICA version of our inference algorithm, where a mixing matrix  $S$  modulated the effect of each speaker on the audio signals. The feature and noise variances were taken from an initial run of the Gibbs sampler where  $\sigma_n$  and  $\sigma_a$  were also sampled.

These two datasets were representative of the large  $N, D$  regime (Yale faces dataset) and the small  $N, D$  regime (speech dataset), so based on our synthetic results, we expect the variational approach to give us an advantage on the faces data, but not on the speech data.

Tables 1 and 2 show the results for each of the datasets. All Gibbs samplers were uncollapsed and run for only 200 iterations.<sup>5</sup> In the higher dimensional Yale dataset, the variational methods outperformed the uncollapsed Gibbs sampler. When started from a random position, the uncollapsed Gibbs samplers quickly became stuck in local optima. The variational method was able to find better local optima because it was initially very uncertain about which features were present in which data points; expressing this uncertainty explicitly through the variational parameters (instead of through a sequence of samples with hard assignments) allowed it the flexibility to improve upon its bad initial starting point.

Table 1: Running times in seconds and test log-likelihoods for the Yale Faces dataset.

Algorithm	K	Time	Test Log-Likelihood ( $\times 10^6$ )
Finite Gibbs	5	464.19	-2.250
	10	940.47	-2.246
	25	2973.7	-2.247
Finite Variational	5	163.24	-1.066
	10	767.1	-0.908
	25	10072	-0.746
Infinite Variational	5	176.62	-1.051
	10	632.53	-0.914
	25	19061	-0.750

The story for the speech dataset, however, is quite different. Here, the variational methods were not only slower than the samplers, but they also achieved lower test-likelihoods. The evaluation on the synthetic datasets points to a potential reason for the difference: the speech

<sup>4</sup>Note that our objective was compare inference techniques for the IBP on real data, not to demonstrate state of the art low-rank approximations.

<sup>5</sup>On the Yale dataset, we did not test the collapsed samplers because the finite collapsed Gibbs sampler required one hour per iteration with  $K = 5$  and the infinite collapsed Gibbs sampler required fifty hours per sample. In the iICA model, the collapsed Gibbs sampler could not be run because the features  $\mathbf{A}$  cannot be marginalised.

dataset is much simpler than the Yale dataset, consisting of 10 dimensions (vs. 1032 in the Yale dataset). In this regime, the Gibbs samplers perform well and the approximations made by the variational method become apparent. As the dimensionality grows, the samplers have more trouble mixing, but the variational methods are still able to find regions of high probability mass.

Table 2: Running times in seconds and test log-likelihoods for the speech dataset.

Algorithm	K	Time	Test Log-Likelihood
Finite Gibbs	2	56	-0.7444
	5	120	-0.4220
	9	201	-0.4205
Infinite Gibbs	na	186	-0.4257
Finite Variational	2	2477	-0.8455
	5	8129	-0.5082
	9	8539	-0.4551
Infinite Variational	2	2702	-0.8810
	5	6065	-0.5000
	9	8491	-0.5486

## 8 Summary

The combinatorial nature of the Indian Buffet Process poses specific challenges for sampling-based inference procedures. In this report, we derived a mean field variational inference procedure for the IBP. Whereas sampling methods work in the discrete space of binary matrices, the variational method allows for soft assignments of features because it approaches the inference problem as a continuous optimisation. We showed experimentally that for high dimensional problems, the soft assignments allow the variational methods to find much better optima than sampling-based approaches. On the other hand, we have shown that sampling based approaches are preferable over our variational approximation.

### Acknowledgments

FD was supported by a Marshall scholarship. KTM was supported by contract DE-AC52-07NA27344 from the U.S. Department of Energy through Lawrence Livermore National Laboratory. JVG was supported by a Microsoft Research scholarship.

## References

- D. Applebaum. *Lévy Processes and Stochastic Calculus*. Cambridge University Press, 2004.
- M. J. Beal. *Variational Algorithms for Approximate Bayesian Inference*. PhD thesis, Gatsby Computational Neuroscience Unit, UCL, 2003.
- D. Blei and M. Jordan. Variational methods for the Dirichlet process. In *Proceedings of the 21st International Conference on Machine Learning*, 2004.

- F. Doshi-Velez, K. T. Miller, J. Van Gael, and Y. W. Teh. Variational inference for the indian buffet process. In *Proc. of the Conference on Artificial Intelligence and Statistics*, 2009.
- A.S. Georghiades, P.N. Belhumeur, and D.J. Kriegman. From few to many: Illumination cone models for face recognition under variable lighting and pose. *IEEE Trans. Pattern Anal. Mach. Intelligence*, 23(6), 2001.
- M. N. Gibbs and D. J. C. MacKay. Variational gaussian process classifiers. *IEEE-NN*, 11(6):1458, November 2000.
- T. Griffiths and Z. Ghahramani. Infinite latent feature models and the Indian buffet process. In *TR 2005-001, Gatsby Computational Neuroscience Unit*, 2005.
- Hemant Ishwaran and Lancelot F. James. Gibbs sampling methods for stick breaking priors. *Journal of the American Statistical Association*, 96(453):161–173, 2001.
- D. Knowles and Z. Ghahramani. Infinite Sparse Factor Analysis and Infinite Independent Components Analysis. *Lecture Notes in Computer Science*, 4666:381, 2007.
- K. Kurihara, M. Welling, and Y. W. Teh. Collapsed variational Dirichlet process mixture models. In *Proceedings of the International Joint Conference on Artificial Intelligence*, volume 20, 2007a.
- Kenichi Kurihara, Max Welling, and Nikos Vlassis. Accelerated variational dirichlet process mixtures. In *Advances in Neural Information Processing Systems 19*. 2007b.
- Edward Meeds, Zoubin Ghahramani, Radford M. Neal, and Sam T. Roweis. Modeling dyadic data with binary latent factors. In *Advances in Neural Information Processing Systems 19*, 2007.
- Edward Snelson and Zoubin Ghahramani. Sparse gaussian processes using pseudo-inputs. In *Advances in Neural Information Processing Systems 18*. MIT press, 2006.
- Y. W. Teh, D. Gorur, and Z. Ghahramani. Stick-breaking construction for the Indian buffet process. In *Proceedings of the 11th Conference on Artificial Intelligence and Statistics*, 2007.
- R. Thibaux and M. Jordan. Hierarchical beta processes and the indian buffet process. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2007.
- M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1(1-2):1–305, 2008.
- Ole Winther. Gaussian processes for classification: Mean field algorithms. *Neural Computation*, 12, 2000.
- Frank Wood and Thomas L. Griffiths. Particle filtering for nonparametric Bayesian matrix factorization. In *Advances in Neural Information Processing Systems 19*. 2007.

## A Gibbs Sampling in the IBP

In this appendix, we briefly review the two Gibbs sampling methods that we compared against. Both of these samplers are for the linear-Gaussian likelihood model described in Section 2.1.

### A.1 Collapsed Gibbs Sampler

The collapsed Gibbs sampler maintains samples over  $\mathbf{Z}$  and integrates out  $\mathbf{A}$ , so we only need to specify the posterior distribution for  $\mathbf{Z}$ . If  $\mathbf{A}$  is needed, then given  $\mathbf{Z}$ , it is easy to compute the posterior for  $\mathbf{A}$  as described by Griffiths and Ghahramani (2005).

**Sampling  $z_{nk}$  for Existing Features** For existing (nonzero) features, the collapsed Gibbs sampler for the IBP resamples each element of the feature assignment matrix  $\mathbf{Z}$  via the equation

$$p(z_{nk} = 1 | \mathbf{Z}_{-nk}, \mathbf{X}) \propto \frac{m_{-n,k}}{N-1} p(\mathbf{X} | \mathbf{Z}) \quad (13)$$

where  $m_{-n,k}$  is the number of observations not including  $z_{nk}$  containing feature  $k$ . The likelihood term  $p(\mathbf{X} | \mathbf{Z})$  is given by

$$p(\mathbf{X} | \mathbf{Z}) = \frac{\exp\left(-\frac{1}{2\sigma_x^2} (\mathbf{X}^T (I - \mathbf{Z}(\mathbf{Z}^T \mathbf{Z} + \frac{\sigma_x^2}{\sigma_a^2} I)^{-1} \mathbf{Z}^T) \mathbf{X})\right)}{(2\pi)^{\frac{ND}{2}} \sigma_{\mathbf{X}}^{(N-K)D} \sigma_a^{KD} |\mathbf{Z}^T \mathbf{Z} + \frac{\sigma_x^2}{\sigma_a^2} I|^{\frac{D}{2}}}. \quad (14)$$

**Sampling New Features** There are an infinite number of remaining columns which contain all zeroes. For any particular  $z_{nk}$ ,  $k > K$ , the probability that  $z_{nk} = 1$  is zero. However, we can sample the number of columns that become nonzero,  $k_{\text{new}}$ , as a batch. (See (Griffiths and Ghahramani, 2005) for details.) The number of new features is sampled according to

$$p(k_{\text{new}}) \propto \text{Poisson}\left(k_{\text{new}}; \frac{\alpha}{N}\right) p(\mathbf{X} | \mathbf{Z}_{\text{new}})$$

where  $\mathbf{Z}_{\text{new}}$  is the feature-assignment matrix with  $k_{\text{new}}$  additional columns set to one for object  $n$  and zero otherwise. The term  $p(\mathbf{X} | \mathbf{Z}_{\text{new}})$  can be computed from Equation (14). We compute these probabilities for  $k_{\text{new}} = 0, \dots, K_{\text{max}}$  for some  $K_{\text{max}}$ , normalise and sample from the resulting multinomial.

**Modifications for the Finite Model** If we are sampling from a finite model with  $K$  features and a finite beta-Bernoulli prior on  $\mathbf{Z}$  as described in Section 4, then Equation (13) becomes

$$p(z_{nk} = 1 | \mathbf{Z}_{-nk}, \mathbf{X}) \propto \frac{m_{-n,k} + \alpha/K}{N-1 + \alpha} p(\mathbf{X} | \mathbf{Z}).$$

We never need to sample the number of new features since  $K$  is fixed.

## A.2 Uncollapsed Gibbs Sampler

The disadvantage of the collapsed Gibbs sampler is that Equation (14) can be expensive to compute. The uncollapsed Gibbs sampler explicitly samples the feature matrix  $\mathbf{A}$  and therefore does not need to evaluate Equation (14). Our samples are therefore over  $\mathbf{Z}$  and  $\mathbf{A}$ .

**Sampling  $z_{nk}$  for Existing Features** The Gibbs sampling Equation for  $z_{nk}$  for existing features is now

$$p(z_{nk} = 1 | \mathbf{Z}_{-nk}, \mathbf{A}, \mathbf{X}) \propto \frac{m_{-n,k}}{N-1} p(\mathbf{X} | \mathbf{Z}, \mathbf{A}) \quad (15)$$

where the likelihood term  $p(\mathbf{X} | \mathbf{Z}, \mathbf{A})$  is given by

$$p(\mathbf{X} | \mathbf{Z}, \mathbf{A}) = \frac{1}{(2\pi\sigma_n^2)^{ND/2}} \exp\left(-\frac{1}{2\sigma_n^2} \text{tr}((\mathbf{X} - \mathbf{Z}\mathbf{A})^\top (\mathbf{X} - \mathbf{Z}\mathbf{A}))\right).$$

**Sampling  $\mathbf{A}$  for Existing Features** The posterior for resampling  $\mathbf{A}$  given  $\mathbf{Z}$  and  $\mathbf{X}$  is

$$p(\mathbf{A} | \mathbf{X}, \mathbf{Z}) \sim \mathcal{N}\left(\left(\mathbf{Z}^\top \mathbf{Z} + \frac{\sigma_n^2}{\sigma_A^2} I\right)^{-1} \mathbf{Z}^\top \mathbf{X}, \sigma_n^2 \left(\mathbf{Z}^\top \mathbf{Z} + \frac{\sigma_n^2}{\sigma_A^2} I\right)^{-1}\right).$$

**Sampling New Features** As in the collapsed sampler, we sample  $k_{\text{new}}$ , the number of new nonzero columns instead of sampling each of the infinite number of all-zero columns independently. As before, the probability of  $k_{\text{new}}$  is

$$p(k_{\text{new}}) \propto \text{Poisson}\left(k_{\text{new}}; \frac{\alpha}{N}\right) p(\mathbf{X}|\mathbf{Z}_{\text{new}}, \mathbf{A}) \quad (16)$$

where  $\mathbf{A}$  represents the initialised features. The likelihood  $p(\mathbf{X}|\mathbf{Z}_{\text{new}}, \mathbf{A})$  is given by the integral  $\int_{\mathbf{A}_{\text{new}}} p(\mathbf{X}|\mathbf{Z}_{\text{new}}, \mathbf{A}, \mathbf{A}_{\text{new}}) p(\mathbf{A}_{\text{new}})$ .

If we want the sampler to be fully uncollapsed, one option for drawing  $k_{\text{new}}$  from the distribution in Equation 16 is to perform a Monte Carlo integration to evaluate  $p(\mathbf{X}|\mathbf{Z}_{\text{new}}, \mathbf{A})$ . For this, we first draw many pairs  $(k_{\text{new}}, \mathbf{A}_{\text{new}})$  from their respective priors, assign a weight to each pair based on the data likelihood  $p(\mathbf{X}|\mathbf{Z}_{\text{new}}, \mathbf{A}, \mathbf{A}_{\text{new}})$ , and then we sample a pair  $(k_{\text{new}}, \mathbf{A}_{\text{new}})$  based on the weights and take the  $k_{\text{new}}$  element of the pair as our  $k_{\text{new}}$ . The advantage of using importance sampling in this way is that the approach remains fully uncollapsed — no integrals need be evaluated. However, since the features of  $\mathbf{A}$  are drawn from the prior, the fully uncollapsed approach is slow to mix.

Another option, if we allow ourselves to have a partially collapsed sampler, is to actually compute the integral in the likelihood in Equation (16) — that is, marginalise over the new features. This option results in a faster mixing sampler, and it was the option used in our tests. The equations below describe how to sample  $k_{\text{new}}$  when  $\mathbf{A}_{\text{new}}$  is marginalised out. For notation, let  $\mathbf{Z}_{\text{old}}$  be the current matrix  $\mathbf{Z}$  and  $\mathbf{A}_{\text{old}}$  be the current matrix  $\mathbf{A}$ . Similarly, let  $\mathbf{Z}_{\text{new}}$  and  $\mathbf{A}_{\text{new}}$  be the parts of  $\mathbf{Z}$  and  $\mathbf{A}$  that correspond to the  $k_{\text{new}}$  new features. Finally, let  $\mathbf{Z}^*$  and  $\mathbf{A}^*$  be the concatenation of the new and old matrices.

Using Bayes rule, we can write

$$p(k_{\text{new}}|X, \mathbf{Z}_{\text{old}}, \mathbf{A}_{\text{old}}) \propto p(X|\mathbf{Z}_{\text{old}}, \mathbf{A}_{\text{old}}, k_{\text{new}}) p(k_{\text{new}}) \quad (17)$$

where  $p(k_{\text{new}})$  is  $\text{Poisson}(\alpha/N)$  and  $p(X|\mathbf{Z}_{\text{old}}, \mathbf{A}_{\text{old}}, k_{\text{new}})$  is the likelihood in which  $\mathbf{A}_{\text{new}}$  has been marginalised out.

We must now specify  $p(\mathbf{X}|\mathbf{Z}_{\text{old}}, \mathbf{A}_{\text{old}}, k_{\text{new}})$ :

$$\begin{aligned} p(\mathbf{X}|\mathbf{Z}_{\text{old}}, \mathbf{A}_{\text{old}}, k_{\text{new}}) &= \int p(\mathbf{X}|\mathbf{Z}_{\text{old}}, \mathbf{A}_{\text{old}}, \mathbf{A}_{\text{new}}, k_{\text{new}}) p(\mathbf{A}_{\text{new}}) d\mathbf{A}_{\text{new}} \\ &= \frac{1}{(2\pi\sigma_n^2)^{ND/2}} \frac{1}{(2\pi\sigma_A^2)^{k_{\text{new}}D/2}} \int \exp\left(-\frac{1}{2}\text{tr}\left(\frac{1}{\sigma_n^2}(\mathbf{X} - \mathbf{Z}^* \mathbf{A}^*)^\top (\mathbf{X} - \mathbf{Z}^* \mathbf{A}^*) + \frac{1}{\sigma_A^2} \mathbf{A}_{\text{new}}^\top \mathbf{A}_{\text{new}}\right)\right) d\mathbf{A}_{\text{new}} \end{aligned}$$

where

$$(\mathbf{X} - \mathbf{Z}^* \mathbf{A}^*)^\top (\mathbf{X} - \mathbf{Z}^* \mathbf{A}^*) = \left(\mathbf{X} - \begin{bmatrix} \mathbf{Z}_{\text{old}} & \mathbf{Z}_{\text{new}} \end{bmatrix} \begin{bmatrix} \mathbf{A}_{\text{old}} \\ \mathbf{A}_{\text{new}} \end{bmatrix}\right)^\top \left(\mathbf{X} - \begin{bmatrix} \mathbf{Z}_{\text{old}} & \mathbf{Z}_{\text{new}} \end{bmatrix} \begin{bmatrix} \mathbf{A}_{\text{old}} \\ \mathbf{A}_{\text{new}} \end{bmatrix}\right)$$

Completing squares to integrate our  $\mathbf{A}_{\text{new}}$ , and dropping terms that do not depend on  $k_{\text{new}}$ , we get

$$\begin{aligned} p(\mathbf{X}|\mathbf{Z}_{\text{old}}, \mathbf{A}_{\text{old}}, k_{\text{new}}) &\propto \frac{(\sigma_n/\sigma_A)^{k_{\text{new}}D}}{|1_{k_{\text{new}} \times k_{\text{new}}} + \frac{\sigma_n^2}{\sigma_A^2} I|^{D/2}} \\ &\times \exp\left\{\frac{1}{2\sigma_n^2} \text{tr}\left((\mathbf{X} - \mathbf{Z}_{\text{old}} \mathbf{A}_{\text{old}})^\top \mathbf{Z}_{\text{new}} \left(1_{k_{\text{new}} \times k_{\text{new}}} + \frac{\sigma_n^2}{\sigma_A^2} I\right)^{-1} \mathbf{Z}_{\text{new}}^\top (\mathbf{X} - \mathbf{Z}_{\text{old}} \mathbf{A}_{\text{old}})\right)\right\}. \end{aligned}$$



We can therefore sample  $k_{\text{new}}$  according to Equation (17) where  $p(\mathbf{X}|\mathbf{Z}_{\text{old}}, \mathbf{A}_{\text{old}}, k_{\text{new}})$  is specified above. Once we have sampled  $k_{\text{new}}$ , we need to sample the newly activated features  $\mathbf{A}_{\text{new}}$ . Based on the same calculations that give us  $p(\mathbf{X}|\mathbf{Z}_{\text{old}}, \mathbf{A}_{\text{old}}, k_{\text{new}})$ , we can sample  $\mathbf{A}_{\text{new}}$  from the distribution

$$\begin{aligned} & p(\mathbf{A}_{\text{new}}|\mathbf{X}, \mathbf{Z}_{\text{new}}, \mathbf{Z}_{\text{old}}, \mathbf{A}_{\text{old}}) \\ & \propto p(\mathbf{X}|\mathbf{Z}_{\text{new}}, \mathbf{Z}_{\text{old}}, \mathbf{A}_{\text{old}}, \mathbf{A}_{\text{new}})p(\mathbf{A}_{\text{new}}) \\ & \sim \mathcal{N}\left(\left(1_{k_{\text{new}} \times k_{\text{new}}} + \frac{\sigma_n^2}{\sigma_A^2}I\right)^{-1} \mathbf{Z}_{\text{new}}^\top (\mathbf{X} - \mathbf{Z}_{\text{old}}\mathbf{A}_{\text{old}}), \sigma_n^2 \left(1_{k_{\text{new}} \times k_{\text{new}}} + \frac{\sigma_n^2}{\sigma_A^2}I\right)^{-1}\right). \end{aligned}$$

**Modifications for the Finite Model** If we are sampling from a finite model with  $K$  features and a beta-Bernoulli prior on  $\mathbf{Z}$ , then Equation (15) becomes

$$p(z_{nk} = 1|\mathbf{Z}_{-nk}, \mathbf{X}) \propto \frac{m_{-n,k} + \alpha/K}{N - 1 + \alpha} p(\mathbf{X}|\mathbf{Z}, \mathbf{A}).$$

In addition, we never need to sample the number of new features since  $K$  is fixed.

## B Variational Inference in Exponential Families

Recall that our goal is to find an approximating distribution  $q \in Q$  with minimum KL divergence  $D(q||p)$  to the true distribution  $p$ . Equation (4) rephrased this optimisation problem in terms of certain expectations and entropies:

$$\arg \min_{\tau, \phi, \nu} D(q||p) = \arg \max_{\tau, \phi, \nu} \mathbb{E}_q[\log(p(\mathbf{X}, \mathbf{W}|\boldsymbol{\theta}))] + H[q]. \quad (18)$$

In general, this optimisation can be quite difficult. However, when the conditional distribution and variational distribution are both in the exponential family, each step in the coordinate ascent has a closed form solution (Beal, 2003; Wainwright and Jordan, 2008). If we are updating the variational parameters  $\xi_i$  that correspond to  $W_i$ , then the optimal  $\xi_i$  are the solution to

$$\log q_{\xi_i}(W_i) = \mathbb{E}_{\mathbf{W}_{-i}}[\log p(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta})] + c \quad (19)$$

where the expectation is taken over all  $\mathbf{W}$  except  $W_i$  according to the variational distribution. In the exponential family, this immediately gives us the updated values of the parameters  $\xi_i$ .

See (Beal, 2003; Wainwright and Jordan, 2008) for more details.

## C Derivations for the Finite Variational Approach

This appendix derives the variational lower bound and the variational updates described in Section 4.

### C.1 Variational Lower Bound

We derive expressions for each expectation in Equation (5):

1. For the feature probabilities, which are beta-distributed,

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\pi}} [\log p(\pi_k|\alpha)] &= \mathbb{E}_{\boldsymbol{\pi}} \left[ \log \left( \frac{\alpha}{K} \pi_k^{\alpha/K-1} \right) \right], \\ &= \log \frac{\alpha}{K} + \left( \frac{\alpha}{K} - 1 \right) \mathbb{E}_{\boldsymbol{\pi}} \log(\pi_k), \\ &= \log \frac{\alpha}{K} + \left( \frac{\alpha}{K} - 1 \right) (\psi(\tau_{k1}) - \psi(\tau_{k1} + \tau_{k2})), \end{aligned}$$

where  $\psi(\cdot)$  is the digamma function.

- For the feature assignments, which are Bernoulli-distributed given the feature probabilities,

$$\begin{aligned}\mathbb{E}_{\boldsymbol{\pi}, \mathbf{Z}} [\log p(z_{nk} | \pi_k)] &= \mathbb{E}_{\boldsymbol{\pi}, \mathbf{Z}} [\log (\pi_k^{z_{nk}} (1 - \pi_k)^{1 - z_{nk}})], \\ &= \mathbb{E}_{\boldsymbol{\pi}, \mathbf{Z}} [z_{nk} \log \pi_k + (1 - z_{nk}) \log(1 - \pi_k)], \\ &= \mathbb{E}_{\mathbf{Z}} [z_{nk}] \mathbb{E}_{\boldsymbol{\pi}} [\log \pi_k] + (1 - \mathbb{E}_{\mathbf{Z}} [z_{nk}]) \mathbb{E}_{\boldsymbol{\pi}} [\log(1 - \pi_k)], \\ &= \nu_{nk} \psi(\tau_{k1}) + (1 - \nu_{nk}) \psi(\tau_{k2}) - \psi(\tau_{k1} + \tau_{k2}).\end{aligned}$$

- For the features, which are Gaussian-distributed,

$$\begin{aligned}\mathbb{E}_{\mathbf{A}} [\log p(\mathbf{A}_k | \sigma_A^2 I)] &= \mathbb{E}_{\mathbf{A}} \left[ \log \left( \frac{1}{(2\pi\sigma_A^2)^{D/2}} \exp \left( -\frac{1}{2\sigma_A^2} \mathbf{A}_k^T \mathbf{A}_k \right) \right) \right], \\ &= \mathbb{E}_{\mathbf{A}} \left[ -\frac{D}{2} \log(2\pi\sigma_A^2) - \frac{1}{2\sigma_A^2} \mathbf{A}_k^T \mathbf{A}_k \right], \\ &= -\frac{D}{2} \log(2\pi\sigma_A^2) - \frac{1}{2\sigma_A^2} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_k^T).\end{aligned}$$

- For the likelihood, which is also Gaussian,

$$\begin{aligned}\mathbb{E}_{\mathbf{Z}, \mathbf{A}} [\log p(\mathbf{X}_n | \mathbf{Z}_n, \mathbf{A}, \sigma_n^2 I)] &= \mathbb{E}_{\mathbf{Z}, \mathbf{A}} \left[ \log \left( \frac{1}{(2\pi\sigma_n^2)^{D/2}} \exp \left( -\frac{1}{2\sigma_n^2} (\mathbf{X}_n - \mathbf{Z}_n \mathbf{A}) (\mathbf{X}_n - \mathbf{Z}_n \mathbf{A})^T \right) \right) \right], \\ &= \mathbb{E}_{\mathbf{Z}, \mathbf{A}} \left[ -\frac{D}{2} \log(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} (\mathbf{X}_n - \mathbf{Z}_n \mathbf{A}) (\mathbf{X}_n - \mathbf{Z}_n \mathbf{A})^T \right], \\ &= -\frac{D}{2} \log(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} (\mathbf{X}_n \mathbf{X}_n^T - 2\mathbb{E}_{\mathbf{Z}} [\mathbf{Z}_n] \mathbb{E}_{\mathbf{A}} [\mathbf{A}] \mathbf{X}_n^T + \mathbb{E}_{\mathbf{Z}, \mathbf{A}} [\mathbf{Z}_n \mathbf{A} \mathbf{A}^T \mathbf{Z}_n^T]), \\ &= -\frac{D}{2} \log(2\pi\sigma_n^2) \\ &\quad - \frac{1}{2\sigma_n^2} \left( \mathbf{X}_n \mathbf{X}_n^T - 2 \sum_{k=1}^K \nu_{nk} \bar{\boldsymbol{\phi}}_k \mathbf{X}_n^T + 2 \sum_{k < k'} \nu_{nk} \nu_{nk'} \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_{k'}^T + \sum_{k=1}^K \nu_{nk} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_k^T) \right),\end{aligned}$$

where the final expectation is derived by

$$\begin{aligned}\mathbb{E}_{\mathbf{Z}, \mathbf{A}} [\mathbf{Z}_n \mathbf{A} \mathbf{A}^T \mathbf{Z}_n^T] &= \mathbb{E}_{\mathbf{Z}, \mathbf{A}} \left[ \left( \sum_{k=1}^K z_{nk} \mathbf{A}_k \right) \left( \sum_{k=1}^K z_{nk} \mathbf{A}_k \right)^T \right], \\ &= \mathbb{E}_{\mathbf{Z}, \mathbf{A}} \left[ \sum_{d=1}^D \left( \sum_{k=1}^K z_{nk} A_{kd}^2 + \sum_{k, k': k' \neq k} z_{nk} z_{nk'} A_{kd} A_{k'd} \right) \right], \\ &= \sum_{k=1}^K \nu_{nk} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_k^T) + 2 \sum_{k < k'} \nu_{nk} \nu_{nk'} \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_{k'}^T.\end{aligned}$$

- Finally, for the entropy,

$$\begin{aligned}H[q] &= -\mathbb{E}_q \log \left[ \prod_{k=1}^K q_{\tau_k}(\pi_k) \prod_{k=1}^K q_{\boldsymbol{\phi}_k}(\mathbf{A}_k) \prod_{k=1}^K \prod_{n=1}^N q_{\nu_{nk}}(z_{nk}) \right], \\ &= \sum_{k=1}^K \mathbb{E}_{\boldsymbol{\pi}} (-\log q_{\tau_k}(\pi_k)) + \sum_{k=1}^K \mathbb{E}_{\mathbf{A}} (-\log q_{\boldsymbol{\phi}_k}(\mathbf{A}_k)) + \sum_{k=1}^K \sum_{n=1}^N \mathbb{E}_{\mathbf{Z}} (-\log q_{\nu_{nk}}(z_{nk})),\end{aligned}$$

where

$$\begin{aligned}
\mathbb{E}_{\boldsymbol{\pi}}(-\log q_{\boldsymbol{\tau}_k}(\boldsymbol{\pi}_k)) &= \log \left( \frac{\Gamma(\tau_{k1})\Gamma(\tau_{k2})}{\Gamma(\tau_{k1} + \tau_{k2})} \right) \\
&\quad - (\tau_{k1} - 1)\psi(\tau_{k1}) - (\tau_{k2} - 1)\psi(\tau_{k2}) + (\tau_{k1} + \tau_{k2} - 2)\psi(\tau_{k1} + \tau_{k2}), \\
\mathbb{E}_{\mathbf{A}}(-\log q_{\boldsymbol{\phi}_k}(\mathbf{A}_{k\cdot})) &= \frac{1}{2} \log \left( (2\pi e)^D |\boldsymbol{\Phi}_k| \right), \\
\mathbb{E}_{\mathbf{Z}}(-\log q_{\nu_{nk}}(z_{nk})) &= -\nu_{nk} \log \nu_{nk} - (1 - \nu_{nk}) \log(1 - \nu_{nk}).
\end{aligned}$$

Putting all the terms together gives us the variational lower bound in Equation (6).

## C.2 Parameter Updates

To optimise the variational parameters, we can directly optimise Equation (6). However, since both our  $p_K$  and our variational approximation are in the exponential family, we can instead use Equation (19) from Appendix B to directly give us the update equations for each parameter given all the rest. We take the latter approach in this section to compute the update equations for the variational parameters in the finite model. Throughout this section, we let  $c$  be a constant independent of the variable of interest that may change from line to line.

1. For the feature distribution at the optimal  $\bar{\boldsymbol{\phi}}_k$  and  $\boldsymbol{\Phi}_k$

$$\begin{aligned}
&\log q_{\boldsymbol{\phi}_k}(\mathbf{A}_{k\cdot}) \\
&= \mathbb{E}_{\mathbf{A}_{-k}, \mathbf{Z}} [\log p_K(\mathbf{W}, \mathbf{X} | \boldsymbol{\theta})] + c, \\
&= \mathbb{E}_{\mathbf{A}_{-k}, \mathbf{Z}} \left[ \log p_K(\mathbf{A}_{k\cdot} | \sigma_A^2) + \sum_{n=1}^N \log p_K(\mathbf{X}_n | \mathbf{Z}_n, \mathbf{A}, \sigma_n^2) \right] + c, \\
&= -\frac{1}{2\sigma_A^2} (\mathbf{A}_{k\cdot} \mathbf{A}_{k\cdot}^T) - \frac{1}{2\sigma_n^2} \sum_{n=1}^N \mathbb{E}_{\mathbf{A}_{-k}, \mathbf{Z}} \left[ (\mathbf{X}_n - \mathbf{Z}_n \mathbf{A}) (\mathbf{X}_n - \mathbf{Z}_n \mathbf{A})^T \right] + c, \\
&= -\frac{1}{2} \left[ \mathbf{A}_{k\cdot} \left( \frac{1}{\sigma_A^2} + \frac{\sum_{n=1}^N \nu_{nk}}{\sigma_n^2} \right) \mathbf{A}_{k\cdot}^T - 2\mathbf{A}_{k\cdot} \left( \frac{1}{\sigma_n^2} \sum_{n=1}^N \nu_{nk} \left( \mathbf{X}_n - \left( \sum_{l:l \neq k} \nu_{nl} \bar{\boldsymbol{\phi}}_l \right) \right) \right)^T \right] + c.
\end{aligned}$$

Completing the squares and using Equation (19) gives us that for the optimal  $\bar{\boldsymbol{\phi}}_k$  and  $\boldsymbol{\Phi}_k$ , we must have

$$\log q_{\boldsymbol{\phi}_k}(\mathbf{A}_{k\cdot}) = -\frac{1}{2} (\mathbf{A}_{k\cdot} \boldsymbol{\Phi}_k^{-1} \mathbf{A}_{k\cdot}^T - 2\mathbf{A}_{k\cdot} \boldsymbol{\Phi}_k^{-1} \bar{\boldsymbol{\phi}}_k^T) + c,$$

which gives us the updates

$$\begin{aligned}
\bar{\boldsymbol{\phi}}_k &= \left[ \frac{1}{\sigma_n^2} \sum_{n=1}^N \nu_{nk} \left( \mathbf{X}_n - \left( \sum_{l:l \neq k} \nu_{nl} \bar{\boldsymbol{\phi}}_l \right) \right) \right] \left( \frac{1}{\sigma_A^2} + \frac{\sum_{n=1}^N \nu_{nk}}{\sigma_n^2} \right)^{-1}, \\
\boldsymbol{\Phi}_k &= \left( \frac{1}{\sigma_A^2} + \frac{\sum_{n=1}^N \nu_{nk}}{\sigma_n^2} \right)^{-1} I.
\end{aligned}$$

2. For the feature state distribution at the optimal  $\nu_{nk}$ ,

$$\begin{aligned}
\log q_{\nu_{nk}}(z_{nk}) &= \mathbb{E}_{\boldsymbol{\pi}, \mathbf{A}, \mathbf{Z}_{-nk}} [\log p_K(\mathbf{W}, \mathbf{X} | \boldsymbol{\theta})] + c, \\
&= \mathbb{E}_{\boldsymbol{\pi}, \mathbf{A}, \mathbf{Z}_{-nk}} [\log p_K(z_{nk} | \boldsymbol{\pi}_k) + \log p_K(\mathbf{X}_n | \mathbf{Z}_n, \mathbf{A}, \sigma_n^2)] + c,
\end{aligned}$$

where

$$\mathbb{E}_{\boldsymbol{\pi}, \mathbf{Z}_{-nk}} [\log p_K(z_{nk}|\pi_k)] = z_{nk} \left[ \psi(\tau_{k1}) - \psi(\tau_{k2}) \right] + \psi(\tau_{k2}) - \psi(\tau_{k1} + \tau_{k2}),$$

and

$$\begin{aligned} & \mathbb{E}_{\mathbf{A}, \mathbf{Z}_{-nk}} [\log p_K(\mathbf{X}_{n\cdot} | \mathbf{Z}_{n\cdot}, \mathbf{A}, \sigma_n^2)] \\ &= \mathbb{E}_{\mathbf{A}, \mathbf{Z}_{-nk}} \left[ -\frac{1}{2\sigma_n^2} (\mathbf{X}_{n\cdot} - \mathbf{Z}_{n\cdot} \mathbf{A}) (\mathbf{X}_{n\cdot} - \mathbf{Z}_{n\cdot} \mathbf{A})^T \right] + c, \\ &= -\frac{1}{2\sigma_n^2} \mathbb{E}_{\mathbf{A}, \mathbf{Z}_{-nk}} [-2\mathbf{Z}_{n\cdot} \mathbf{A} \mathbf{X}_{n\cdot}^T + \mathbf{Z}_{n\cdot} \mathbf{A} \mathbf{A}^T \mathbf{Z}_{n\cdot}^T] + c, \\ &= -\frac{1}{2\sigma_n^2} \left[ -2z_{nk} \bar{\boldsymbol{\phi}}_k \mathbf{X}_{n\cdot}^T + z_{nk} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_k^T) + 2z_{nk} \bar{\boldsymbol{\phi}}_k \left( \sum_{l:l \neq k} \nu_{nl} \bar{\boldsymbol{\phi}}_l^T \right) \right] + c. \end{aligned}$$

Therefore

$$\begin{aligned} & \log q_{\nu_{nk}}(z_{nk}) \\ &= z_{nk} \left[ \psi(\tau_{k1}) - \psi(\tau_{k2}) - \frac{1}{2\sigma_n^2} \left( \text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_k^T - 2\bar{\boldsymbol{\phi}}_k \mathbf{X}_{n\cdot}^T + 2\bar{\boldsymbol{\phi}}_k \left( \sum_{l:l \neq k} \nu_{nl} \bar{\boldsymbol{\phi}}_l^T \right) \right) \right] + c. \end{aligned}$$

From the canonical parameterisation of the Bernoulli distribution, we get that

$$\begin{aligned} \log \frac{\nu_{nk}}{1 - \nu_{nk}} &= \psi(\tau_{k1}) - \psi(\tau_{k2}) - \frac{1}{2\sigma_n^2} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_k^T) + \frac{1}{\sigma_n^2} \bar{\boldsymbol{\phi}}_k \left( \mathbf{X}_{n\cdot}^T - \left( \sum_{l:l \neq k} \nu_{nl} \bar{\boldsymbol{\phi}}_l^T \right) \right), \\ &\equiv \vartheta. \end{aligned}$$

Which gives us the update

$$\nu_{nk} = \frac{1}{1 + e^{-\vartheta}}.$$

3. For the feature probabilities at the optimal  $\tau_{k1}$  and  $\tau_{k2}$ ,

$$\begin{aligned} \log q_{\tau_k}(\pi_k) &= \mathbb{E}_{\mathbf{A}, \mathbf{Z}} [\log p_K(\mathbf{W}, \mathbf{X} | \boldsymbol{\theta})] + c, \\ &= \mathbb{E}_{\mathbf{A}, \mathbf{Z}} \left[ \log p_K(\pi_k | \alpha) + \sum_{n=1}^N \log p_K(z_{nk} | \pi_k) \right] + c, \\ &= \left( \frac{\alpha}{K} - 1 \right) \log \pi_k + \sum_{n=1}^N (\nu_{nk} \log \pi_k + (1 - \nu_{nk}) \log(1 - \pi_k)) + c. \end{aligned}$$

Hence the updates are

$$\begin{aligned} \tau_{k1} &= \frac{\alpha}{K} + \sum_{n=1}^N \nu_{nk}, \\ \tau_{k2} &= 1 + \sum_{n=1}^N (1 - \nu_{nk}). \end{aligned}$$

## D Derivations for the Infinite Variational Approach

This appendix derives the variational lower bound and the variational updates described in Section 5.

### D.1 Variational Lower Bound

We derive expressions for each expectation in Equation (7):

1. Each stick is independent, and substituting the form of the beta prior we get

$$\begin{aligned}\mathbb{E}_{\mathbf{v}} [\log p(v_k|\alpha)] &= \mathbb{E}_{\mathbf{v}} [\log (\alpha v_k^{\alpha-1})], \\ &= \log \alpha + (\alpha - 1) \mathbb{E}_{\mathbf{v}} \log(v_k), \\ &= \log \alpha + (\alpha - 1) (\psi(\tau_{k1}) - \psi(\tau_{k1} + \tau_{k2})),\end{aligned}$$

where  $\psi(\cdot)$  is the digamma function.

2. For the feature assignments, which are Bernoulli-distributed given the feature probabilities, we first break the expectation into the following parts

$$\begin{aligned}\mathbb{E}_{\mathbf{v}, \mathbf{Z}} [\log p(z_{nk}|\mathbf{v})] &= \mathbb{E}_{\mathbf{v}, \mathbf{Z}} [\log p(z_{nk} = 1|\mathbf{v})^{z_{nk}} p(z_{nk} = 0|\mathbf{v})^{1-z_{nk}}] \\ &= \mathbb{E}_{\mathbf{Z}} [z_{nk}] \mathbb{E}_{\mathbf{v}} \left[ \log \prod_{m=1}^k v_m \right] + \mathbb{E}_{\mathbf{Z}} [1 - z_{nk}] \mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right] \\ &= \nu_{nk} \left( \sum_{m=1}^k \psi(\tau_{k2}) - \psi(\tau_{k1} + \tau_{k2}) \right) + (1 - \nu_{nk}) \mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right]\end{aligned}$$

The second line follows from the definition of  $\mathbf{v}$ , while the third line follows from the properties of Bernoulli and beta distributions. In Section 5.1, we discussed how to compute a lower bound for  $\mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right]$  using a multinomial approximation since there is no closed form method to evaluate it. In practice, we use this lower bound for  $\mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right]$ . The end of this subsection discusses an alternate approach that can give a closer lower bound, but is also much more computationally expensive. If we use a heuristic evaluation of this term that is not a strict lower bound, then we risk not having a definite means of evaluating convergence of our algorithm since we are no longer optimising a lower bound.

3. For the feature distribution, we simply apply the properties of expectations of Gaussians to get

$$\begin{aligned}\mathbb{E}_{\mathbf{A}} [\log p(\mathbf{A}_k | \sigma_A^2 I)] &= \mathbb{E}_{\mathbf{A}} \left[ \log \left( \frac{1}{(2\pi\sigma_A^2)^{D/2}} \exp \left( -\frac{1}{2\sigma_A^2} \mathbf{A}_k^T \mathbf{A}_k \right) \right) \right], \\ &= \mathbb{E}_{\mathbf{A}} \left[ \frac{-D}{2} \log(2\pi\sigma_A^2) - \frac{1}{2\sigma_A^2} \mathbf{A}_k^T \mathbf{A}_k \right], \\ &= \frac{-D}{2} \log(2\pi\sigma_A^2) - \frac{1}{2\sigma_A^2} (tr(\mathbf{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T).\end{aligned}$$

4. The likelihood for a particular observation is identical to the finite model, so we again have

$$\begin{aligned} \mathbb{E}_{\mathbf{Z}, \mathbf{A}} [\log p(\mathbf{X}_n | \mathbf{Z}_n, \mathbf{A}, \sigma_n^2 I)] \\ = -\frac{D}{2} \log(2\pi\sigma_n^2) \\ - \frac{1}{2\sigma_n^2} \left( \mathbf{X}_n \cdot \mathbf{X}_n^T - 2 \sum_{k=1}^K \nu_{nk} \bar{\phi}_k \mathbf{X}_n^T + 2 \sum_{k < k'} \nu_{nk} \nu_{nk'} \bar{\phi}_k \bar{\phi}_{k'}^T + \sum_{k=1}^K \nu_{nk} (\text{tr}(\Phi_k) + \bar{\phi}_k \bar{\phi}_k^T) \right). \end{aligned}$$

5. The entropy can also be easily computed, since we have chosen exponential family distributions for our variational approximation:

$$\begin{aligned} H[q] &= -\mathbb{E}_q \log \left[ \prod_{k=1}^K q_{\tau_k}(v_k) \prod_{k=1}^K q_{\phi_k}(\mathbf{A}_{k\cdot}) \prod_{k=1}^K \prod_{n=1}^N q_{\nu_{nk}}(z_{nk}) \right], \\ &= \sum_{k=1}^K \mathbb{E}_{\mathbf{v}}(-\log q_{\tau_k}(v_k)) + \sum_{k=1}^K \mathbb{E}_{\mathbf{A}}(-\log q_{\phi_k}(\mathbf{A}_{k\cdot})) + \sum_{k=1}^K \sum_{n=1}^N \mathbb{E}_{\mathbf{Z}}(-\log q_{\nu_{nk}}(z_{nk})), \end{aligned}$$

where

$$\begin{aligned} \mathbb{E}_{\mathbf{v}}(-\log q_{\tau_k}(v_k)) &= \log \left( \frac{\Gamma(\tau_{k1})\Gamma(\tau_{k2})}{\Gamma(\tau_{k1} + \tau_{k2})} \right) \\ &\quad - (\tau_{k1} - 1)\psi(\tau_{k1}) - (\tau_{k2} - 1)\psi(\tau_{k2}) + (\tau_{k1} + \tau_{k2} - 2)\psi(\tau_{k1} + \tau_{k2}). \\ \mathbb{E}_{\mathbf{A}}(-\log q_{\phi_k}(\mathbf{A}_{k\cdot})) &= \frac{1}{2} \log((2\pi e)^D |\Phi_k|). \\ \mathbb{E}_{\mathbf{Z}}(-\log q_{\nu_{nk}}(z_{nk})) &= -\nu_{nk} \log \nu_{nk} - (1 - \nu_{nk}) \log(1 - \nu_{nk}). \end{aligned}$$

Putting all the terms together gives us the variational lower bound in Equation (8).

**Alternate Evaluation of  $\mathbb{E}_{\mathbf{v}}$**   $\left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right]$  We describe a Taylor series alternative to the multinomial lower bound from Section 5.1. As we noted before, the advantage of the Taylor series approximation is that we can make it arbitrarily accurate by including more terms. However, in practice, the multinomial approximation is nearly as accurate, computationally much faster, and leads to straightforward parameter updates.

Recall the Taylor series  $\log(1 - x) = -\sum_n^{\infty} \frac{x^n}{n}$  and that it converges for  $x \in (-1, 1)$ . In our case,  $x$  corresponds to the product of probabilities, so the sum will converge unless *all* of the  $v_m$ s are one or any of them is zero. Since the distribution over the  $v_m$  are continuous densities on  $[0, 1]$ , the series will converge with probability one.

Applying the Taylor expansion to our desired expectation, we obtain

$$\begin{aligned} \mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right] &= \mathbb{E}_{\mathbf{v}} \left[ -\sum_{n=1}^{\infty} \frac{1}{n} \prod_{m=1}^k v_m^n \right] \\ &= -\sum_{n=1}^{\infty} \frac{1}{n} \prod_{m=1}^k \frac{\Gamma(\tau_{m1} + n)\Gamma(\tau_{m2} + \tau_{m1})}{\Gamma(\tau_{m1})\Gamma(\tau_{m2} + \tau_{m1} + n)} \\ &= -\sum_{n=1}^{\infty} \frac{1}{n} \prod_{m=1}^k \frac{(\tau_{m1}) \cdots (\tau_{m1} + n - 1)}{(\tau_{m2} + \tau_{m1}) \cdots (\tau_{m2} + \tau_{m1} + n - 1)} \end{aligned}$$

where we have used the fact that the moments of  $x \sim \text{Beta}(\alpha, \beta)$  are given by

$$\mathbb{E}[x^n] = \frac{\Gamma(\alpha + n)\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\alpha + \beta + n)}.$$

If we simply wish to approximate the variational lower bound, we could truncate the series after a certain number of terms. However, since all of the terms in the Taylor series are negative, truncating the series will not produce a lower bound. Thus, some extra work is required if we wish to preserve the lower bound.

To preserve the lower bound, we first note that if  $\tau_{m2}$  were an integer, most terms in the numerator and denominator would cancel for  $n > \tau_{m2}$ . Let  $T$  be an integer greater than  $\lceil \max_{m \in \{1, \dots, k\}}(\tau_{m2}) \rceil$ , then we can write a lower bound for the series in the following form:

$$\begin{aligned} & \mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right] \\ & \geq - \sum_{n=1}^T \frac{1}{n} \prod_{m=1}^k \frac{(\tau_{m1}) \cdots (\tau_{m1} + n - 1)}{(\tau_{m1} + \tau_{m2}) \cdots (\tau_{m1} + \tau_{m2} + n - 1)} \\ & \quad - \prod_{m=1}^k ((\tau_{m1}) \cdots (\tau_{m1} + \lfloor \tau_{m2} \rfloor - 1)) \cdot \sum_{n=T+1}^{\infty} \frac{1}{n} \prod_{m=1}^k \frac{1}{(\tau_{m1} + n) \cdots (\tau_{m1} + \lfloor \tau_{m2} \rfloor + n - 1)}, \end{aligned}$$

where we have factored the final term to make it clear that the second term has  $n$  only in the denominator. A fairly trivial upper bound on the second sum (and therefore lower bound on the expectation) is:  $\zeta(1 + \sum_{m=1}^k \lfloor \tau_{m2} \rfloor)$ ; a slightly better bound is  $\zeta_H(1 + \sum_{m=1}^k \lfloor \tau_{m2} \rfloor, T + 1)$ , where  $\zeta_H(\cdot, \cdot)$  is the generalised or Hurwitz zeta function. The quality of the bound depends on the choice of  $T$ . For larger  $T$ , we have to compute more terms in the first summation, but the error introduced by the fact that the denominator of the second term is  $(\tau_{m1} + n)$ , not  $n$ , decreases. Empirically, we find that setting  $T = \lceil 2 \max_{m \in \{1, \dots, k\}}(\tau_{m2}) \rceil$  results in very close approximations.

More precisely, we know that the Taylor series reaches the true value from above (since all of the terms in the series are negative) and that the value of the zeta function is a bound on the error. Thus, we can place the true expectation in an interval

$$\mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right] \in - \sum_{n=1}^T \frac{1}{n} \prod_{m=1}^k \frac{(\tau_{m1}) \cdots (\tau_{m1} + n - 1)}{(\tau_{m1} + \tau_{m2}) \cdots (\tau_{m1} + \tau_{m2} + n - 1)} + [-\epsilon, 0],$$

where

$$\epsilon = \prod_{m=1}^k ((\tau_{m1}) \cdots (\tau_{m1} + \lfloor \tau_{m2} \rfloor - 1)) \zeta_H \left( 1 + \sum_{m=1}^k \lfloor \tau_{m2} \rfloor, T + 1 \right).$$

## D.2 Parameter Updates

To optimise the parameters, we can directly optimise Equation (8). However, as in the finite case, when we are in the exponential family, we can sequentially update each of the parameters using Equation (19) from Appendix B. Regardless of how we compute the lower bound, the conditional updates for the features  $\mathbf{A}$  and feature assignments  $\mathbf{Z}$  remain within the exponential family, so we can use the exponential family updates. If we use the multinomial lower bound discussed in Section 5.1, then the updates for  $\boldsymbol{\tau}$  will also be in the exponential family. However, the Taylor series approximation from Appendix D.1 takes us outside of the exponential family and therefore requires a numerical optimisation to update  $\boldsymbol{\tau}$ .

1. The updates for the features  $\mathbf{A}$  are identical to the finite approximation; see Appendix C.2.
2. The updates for the variational distribution on  $\mathbf{Z}$  are slightly different. For the  $\nu$  parameters,

$$\begin{aligned}\log q_{\nu_{nk}}(z_{nk}) &= \mathbb{E}_{\mathbf{v}, \mathbf{A}, \mathbf{Z}_{-nk}} [\log p(\mathbf{W}, \mathbf{X} | \boldsymbol{\theta})] + c, \\ &= \mathbb{E}_{\mathbf{v}, \mathbf{A}, \mathbf{Z}_{-nk}} [\log p(z_{nk} | \mathbf{v}) + \log p(X_n | \mathbf{Z}_n, \mathbf{A}, \sigma_n^2)] + c,\end{aligned}$$

where

$$\mathbb{E}_{\mathbf{v}, \mathbf{Z}_{-nk}} [\log p(z_{nk} | \mathbf{v})] = z_{nk} \sum_{i=1}^k (\psi(\tau_{i1}) - \psi(\tau_{i1} + \tau_{i2})) + (1 - z_{nk}) \mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{i=1}^k v_i \right) \right]$$

and as in Appendix C.2

$$\begin{aligned}\mathbb{E}_{\mathbf{A}, \mathbf{Z}_{-nk}} [\log p(\mathbf{X}_n | \mathbf{Z}_n, \mathbf{A}, \sigma_n^2)] \\ = -\frac{1}{2\sigma_n^2} \left[ -2z_{nk} \bar{\boldsymbol{\phi}}_k \mathbf{X}_n^T + z_{nk} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_k^T) + 2z_{nk} \bar{\boldsymbol{\phi}}_k \left( \sum_{l:l \neq k} \nu_{nl} \bar{\boldsymbol{\phi}}_l^T \right) \right] + c.\end{aligned}$$

Therefore

$$\begin{aligned}\log q_{\nu_{nk}}(z_{nk}) &= z_{nk} \left[ \sum_{i=1}^k (\psi(\tau_{i1}) - \psi(\tau_{i1} + \tau_{i2})) - \mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{i=1}^k v_i \right) \right] \right. \\ &\quad \left. - \frac{1}{2\sigma_n^2} \left( \text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_k^T - 2\bar{\boldsymbol{\phi}}_k \mathbf{X}_n^T + 2\bar{\boldsymbol{\phi}}_k \left( \sum_{l:l \neq k} \nu_{nl} \bar{\boldsymbol{\phi}}_l^T \right) \right) \right] + c.\end{aligned}$$

From the canonical parameterisation of the Bernoulli distribution, we get that

$$\begin{aligned}\log \frac{\nu_{nk}}{1 - \nu_{nk}} &= \sum_{i=1}^k (\psi(\tau_{i1}) - \psi(\tau_{i1} + \tau_{i2})) - \mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{i=1}^k v_i \right) \right] \\ &\quad - \frac{1}{2\sigma_n^2} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_k^T) + \frac{1}{\sigma_n^2} \bar{\boldsymbol{\phi}}_k \left( \mathbf{X}_n^T - \left( \sum_{l:l \neq k} \nu_{nl} \bar{\boldsymbol{\phi}}_l^T \right) \right) \\ &\equiv \vartheta,\end{aligned}$$

where the remaining expectation can be computed using either the multinomial approximation or the Taylor series. This gives us the update

$$\nu_{nk} = \frac{1}{1 + e^{-\vartheta}}.$$

3. The updates for  $\boldsymbol{\tau}$  depend on how we deal with the term  $\mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right]$ . We first discuss the case of using a multinomial lower bound in which we have a closed form, exponential family update. We then discuss how numerical optimisation must be used for the Taylor series lower bound.



### D.2.1 Multinomial Lower Bound.

When we use the multinomial bound, compute  $q_k$  and then hold  $q_k$  fixed, the terms in Equation (8) that contain  $\tau_k$  are

$$\begin{aligned} \mathcal{L}_{\tau_k} = & \left[ \alpha + \sum_{m=k}^K \sum_{n=1}^N \nu_{nm} + \sum_{m=k+1}^K \left( N - \sum_{n=1}^N \nu_{nm} \right) \left( \sum_{i=k+1}^m q_{mi} \right) - \tau_{k1} \right] (\Psi(\tau_{k1}) - \Psi(\tau_{k1} + \tau_{k2})) \\ & + \left[ 1 + \sum_{m=k}^K \left( N - \sum_{n=1}^N \nu_{nm} \right) q_{mk} - \tau_{k2} \right] (\Psi(\tau_{k2}) - \Psi(\tau_{k1} + \tau_{k2})) + \ln \left( \frac{\Gamma(\tau_{k1})\Gamma(\tau_{k2})}{\Gamma(\tau_{k1} + \tau_{k2})} \right). \end{aligned}$$

We can then optimise this with respect to  $\tau_k$  to find that the optimal values of  $\tau_{k1}$  and  $\tau_{k2}$  are

$$\begin{aligned} \tau_{k1} &= \alpha + \sum_{m=k}^K \sum_{n=1}^N \nu_{nm} + \sum_{m=k+1}^K \left( N - \sum_{n=1}^N \nu_{nm} \right) \left( \sum_{i=k+1}^m q_{mi} \right) \\ \tau_{k2} &= 1 + \sum_{m=k}^K \left( N - \sum_{n=1}^N \nu_{nm} \right) q_{mk}. \end{aligned}$$

These are equivalent to exponential family updates just like in the finite variational approximation.

### D.2.2 Taylor Series Lower Bound.

When we use a Taylor series approximation (which we will then truncate as opposed to using a zeta function lower bound), we find that the terms in Equation (8) that contain  $\tau_k$  are

$$\begin{aligned} \mathcal{L}_{\tau_k} = & (\alpha - 1) (\Psi(\tau_{k1}) - \Psi(\tau_{k1} + \tau_{k2})) \\ & + \sum_{m=k}^K \sum_{n=1}^N \nu_{nm} (\Psi(\tau_{k1}) - \Psi(\tau_{k1} + \tau_{k2})) \\ & - \sum_{m=k}^K \sum_{n=1}^N (1 - \nu_{nm}) \sum_{r=1}^{\infty} \frac{1}{r} \frac{(\tau_{k1}) \dots (\tau_{k1} + r - 1)}{(\tau_{k1} + \tau_{k2}) \dots (\tau_{k1} + \tau_{k2} + r - 1)} \\ & + \ln \left( \frac{\Gamma(\tau_{k1})\Gamma(\tau_{k2})}{\Gamma(\tau_{k1} + \tau_{k2})} \right) - (\tau_{k1} - 1)\Psi(\tau_{k1}) - (\tau_{k2} - 1)\Psi(\tau_{k2}) + (\tau_{k1} + \tau_{k2} - 2)\Psi(\tau_{k1} + \tau_{k2}). \end{aligned}$$

These is not a standard exponential family equation, so we must numerically optimise  $\tau_k$  to increase the lower bound. The derivatives with respect to  $\tau_{k1}$  are given by

$$\begin{aligned} \frac{\partial \mathcal{L}_{\tau_k}}{\partial \tau_{k1}} = & (\alpha - 1) (\Psi'(\tau_{k1}) - \Psi'(\tau_{k1} + \tau_{k2})) + \sum_{m=k}^K \sum_{n=1}^N \nu_{nm} (\Psi'(\tau_{k1}) - \Psi'(\tau_{k1} + \tau_{k2})) \\ & - \sum_{m=k}^K \left( N - \sum_{n=1}^N \nu_{nm} \right) \sum_{r=1}^{\infty} \frac{1}{r} \left( \prod_{i=1}^m \frac{(\tau_{i1}) \dots (\tau_{i1} + r - 1)}{(\tau_{i1} + \tau_{i2}) \dots (\tau_{i1} + \tau_{i2} + r - 1)} \right) \sum_{j=1}^r \frac{\tau_{k2}}{(\tau_{k1} + j - 1)(\tau_{k1} + \tau_{k2} + j - 1)} \\ & - (\tau_{k1} - 1)\Psi'(\tau_{k1}) + (\tau_{k1} + \tau_{k2} - 2)\Psi'(\tau_{k1} + \tau_{k2}). \end{aligned}$$

Similarly, the derivatives with respect to  $\tau_{k2}$  are given by

$$\begin{aligned} \frac{\partial \mathcal{L}_{\tau_k}}{\partial \tau_{k2}} &= (\alpha - 1)(-\Psi'(\tau_{k1} + \tau_{k2})) - \sum_{m=k}^K \sum_{n=1}^N \nu_{nm} \Psi'(\tau_{k1} + \tau_{k2}) \\ &\quad - \sum_{m=k}^K \left( N - \sum_{n=1}^N \nu_{nm} \right) \sum_{r=1}^{\infty} \frac{1}{r} \left( \prod_{i=1}^m \frac{(\tau_{i1}) \dots (\tau_{i1} + r - 1)}{(\tau_{i1} + \tau_{i2}) \dots (\tau_{i1} + \tau_{i2} + r - 1)} \right) \sum_{j=1}^r \frac{-1}{\tau_{k1} + \tau_{k2} + j - 1} \\ &\quad - (\tau_{k2} - 1)\Psi'(\tau_{k2}) + (\tau_{k1} + \tau_{k2} - 2)\Psi'(\tau_{k1} + \tau_{k2}). \end{aligned}$$

These can be computed for any particular parameter choices (and an arbitrary truncation level of the infinite sum). Also note that several computations can be reused across  $k$ , and others can be computed iteratively across  $r$ . We can plug these derivatives into an optimisation routine to get updates for  $\tau_{k1}$  and  $\tau_{k2}$ .

## E Variational Inference for the iICA model

In this section we describe the variational approach to do approximate inference for the infinite Independent Component Analysis model. We refer to Knowles and Ghahramani (2007) for more details regarding the model, but to set the notation, we state the iICA model

$$\begin{aligned} v_k &\sim \text{Beta}(\alpha, 1) && \text{for } k \in \{1, \dots, \infty\}, \\ \pi_k &= \prod_{i=1}^k v_i && \text{for } k \in \{1 \dots \infty\}, \\ z_{nk} &\sim \text{Bernoulli}(\pi_k) && \text{for } k \in \{1 \dots \infty\}, \\ s_{nk} &\sim \text{Laplace}(1) && \text{for } k \in \{1 \dots K\}, n \in \{1 \dots N\}, \\ \mathbf{A}_{k\cdot} &\sim \text{Normal}(0, \sigma_A^2 I) && \text{for } k \in \{1 \dots \infty\}, \\ \mathbf{X}_{n\cdot} &\sim \text{Normal}((\mathbf{Z}_{n\cdot} \odot \mathbf{S}_{n\cdot})\mathbf{A}, \sigma_n^2 I) && \text{for } n \in \{1 \dots N\}, \end{aligned}$$

where  $\odot$  denotes pointwise multiplication between two vectors. In other words, we can write the joint probability of the data and latent variables as

$$p(\mathbf{W}, \mathbf{X} | \boldsymbol{\theta}) = \prod_{k=1}^{\infty} \left( p(\pi_k | \alpha) p(\mathbf{A}_{k\cdot} | \sigma_A^2 I) \prod_{n=1}^N p(z_{nk} | \pi_k) p(s_{nk}) \right) \prod_{n=1}^N p(\mathbf{X}_{n\cdot} | \mathbf{Z}_{n\cdot}, \mathbf{A}, \mathbf{S}, \sigma_n^2 I).$$

As in the linear-Gaussian model and as we will discuss in Section E.2, doing exact posterior inference on the latent variables  $\mathbf{W}$  is intractable. We will propose two different approximation schemes just as in the linear-Gaussian model: first we introduce a finite approximation similar to the derivation in Section C and second we describe an infinite variational approximation to the iICA model analogous to Section D. Similar to our discussion of the linear-Gaussian model, we will refer to the distribution  $p_K(\mathbf{W}, \mathbf{X} | \boldsymbol{\theta})$  as the finite beta-Bernoulli approximation of order  $K$  while  $p(\mathbf{W}, \mathbf{X} | \boldsymbol{\theta})$  refers to the exact iICA distribution defined above.

## E.1 The Finite Variational Approach

A finite beta-Bernoulli approximation to the iICA model can be described as follows

$$\begin{aligned}
\pi_k &\sim \text{Beta}(\alpha/K, 1) && \text{for } k \in \{1 \cdots K\}, \\
z_{nk} &\sim \text{Bernoulli}(\pi_k) && \text{for } k \in \{1 \cdots K\}, n \in \{1 \cdots N\}, \\
s_{nk} &\sim \text{Laplace}(1) && \text{for } k \in \{1 \cdots K\}, n \in \{1 \cdots N\}, \\
\mathbf{A}_{k\cdot} &\sim \text{Normal}(0, \sigma_A^2 I) && \text{for } k \in \{1 \cdots K\}, \\
\mathbf{X}_{n\cdot} &\sim \text{Normal}((\mathbf{Z}_{n\cdot} \odot \mathbf{S}_{n\cdot})\mathbf{A}, \sigma_n^2 I) && \text{for } n \in \{1 \cdots N\}.
\end{aligned}$$

where  $K$  is some finite (but large) truncation level. We refer to the set of hidden variables as  $\mathbf{W} = \{\boldsymbol{\pi}, \mathbf{Z}, \mathbf{A}, \mathbf{S}\}$  and the set of parameters as  $\boldsymbol{\theta} = \{\alpha, \sigma_A^2, \sigma_n^2\}$ . Using this notation we can write the joint probability of the data and latent variables as

$$p_K(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta}) = \prod_{k=1}^K \left( p(\pi_k|\alpha) p(\mathbf{A}_{k\cdot}|\sigma_A^2 I) \prod_{n=1}^N p(z_{nk}|\pi_k) p(s_{nk}) \right) \prod_{n=1}^N p(\mathbf{X}_{n\cdot}|\mathbf{Z}_{n\cdot}, \mathbf{A}, \mathbf{S}, \sigma_n^2 I).$$

We are interested in the posterior, or equivalently the log posterior, of the latent variables

$$\log p_K(\mathbf{W}|\mathbf{X}, \boldsymbol{\theta}) = \log p_K(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta}) - \log p_K(\mathbf{X}|\boldsymbol{\theta}). \quad (20)$$

For similar reasons to the linear-Gaussian model in Section 4, computing this quantity is intractable. Hence we use the following variational distribution as an approximation

$$q(\mathbf{W}) = q_{\boldsymbol{\tau}}(\boldsymbol{\pi}) q_{\boldsymbol{\phi}}(\mathbf{A}) q_{\boldsymbol{\nu}}(\mathbf{Z}) q_{\boldsymbol{\mu}, \boldsymbol{\eta}}(\mathbf{S}).$$

where

- $q_{\boldsymbol{\tau}_k}(\pi_k) = \text{Beta}(\pi_k; \tau_{k1}, \tau_{k2})$ ,
- $q_{\boldsymbol{\phi}_k}(\mathbf{A}_{k\cdot}) = \text{Normal}(\mathbf{A}_{k\cdot}; \bar{\boldsymbol{\phi}}_k, \bar{\boldsymbol{\Phi}}_k)$ ,
- $q_{\nu_{nk}}(z_{nk}) = \text{Bernoulli}(z_{nk}; \nu_{nk})$ ,
- $q_{\mu_{nk}, \eta_{nk}}(s_{nk}) = \text{Laplace}(s_{nk}; \mu_{nk}, \eta_{nk})$  where  $\mu_{nk}$  is the mean and  $\eta_{nk}$  is the scale parameter of the Laplace distribution.

In contrast to the linear-Gaussian model, we now need to optimise the parameters  $\boldsymbol{\mu}, \boldsymbol{\eta}$  in addition to  $\boldsymbol{\tau}, \boldsymbol{\phi}, \boldsymbol{\nu}$  with the goal of minimizing KL divergence  $D(q||p_K)$  or equivalently, maximise the following lower bound on  $p_K(\mathbf{X}|\boldsymbol{\theta})$ :

$$\mathbb{E}_q[\log(p_K(\mathbf{X}, \mathbf{W}|\boldsymbol{\theta}))] + H[q].$$

As we discussed in the context of the linear-Gaussian model, inference with respect to this beta-Bernoulli model  $p_K$  is not the same as variational inference with respect to the true iICA model. The variational updates are significantly easier though and in the limit of large  $K$ , the finite beta-Bernoulli model is equivalent to the iICA model.

### E.1.1 Variational Lower Bound

We expand the lower bound on  $\log p_K(\mathbf{X}|\boldsymbol{\theta})$  into its components

$$\begin{aligned}
\log p_K(\mathbf{X}|\boldsymbol{\theta}) &\geq \mathbb{E}_{\mathbf{W}} [\log p(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta})] + H[q], \\
&= \sum_{k=1}^K \mathbb{E}_{\boldsymbol{\pi}} [\log p(\pi_k|\alpha)] + \sum_{k=1}^K \sum_{n=1}^N \mathbb{E}_{\boldsymbol{\pi}, \mathbf{Z}} [\log p(z_{nk}|\pi_k)] + \sum_{k=1}^K \sum_{n=1}^N \mathbb{E}_{\mathbf{S}} [\log p(s_{nk})] \\
&\quad + \sum_{k=1}^K \mathbb{E}_{\mathbf{A}} [\log p(\mathbf{A}_k|\sigma_A^2 I)] + \sum_{n=1}^N \mathbb{E}_{\mathbf{Z}, \mathbf{A}, \mathbf{S}} [\log p(\mathbf{X}_n|\mathbf{Z}_n, \mathbf{A}, \mathbf{S}_n, \sigma_n^2 I)] + H[q],
\end{aligned} \tag{21}$$

where the expectation are computed with respect to the variational distribution  $q$ . We derive expressions for each expectation in Equation (21):

1. The feature probabilities,

$$\begin{aligned}
\mathbb{E}_{\boldsymbol{\pi}} [\log p(\pi_k|\alpha)] &= \mathbb{E}_{\boldsymbol{\pi}} \left[ \log \left( \frac{\alpha}{K} \pi_k^{\alpha/K-1} \right) \right], \\
&= \log \frac{\alpha}{K} + \left( \frac{\alpha}{K} - 1 \right) (\psi(\tau_{k1}) - \psi(\tau_{k1} + \tau_{k2})),
\end{aligned}$$

where  $\psi(\cdot)$  is the digamma function.

2. The signal distribution,

$$\begin{aligned}
\mathbb{E}_{\mathbf{S}} [\log p(s_{nk})] &= \mathbb{E}_{\mathbf{S}} \left[ \log \left( \frac{1}{2} \exp(-|s_{nk}|) \right) \right], \\
&= -\log 2 - \left( |\mu_{nk}| + \eta_{nk} \exp \left( -\frac{|\mu_{nk}|}{\eta_{nk}} \right) \right).
\end{aligned}$$

3. The feature state distribution,

$$\begin{aligned}
\mathbb{E}_{\boldsymbol{\pi}, \mathbf{Z}} [\log p(z_{nk}|\pi_k)] &= \mathbb{E}_{\boldsymbol{\pi}, \mathbf{Z}} [\log (\pi_k^{z_{nk}} (1 - \pi_k)^{1-z_{nk}})], \\
&= \nu_{nk} \psi(\tau_{k1}) + (1 - \nu_{nk}) \psi(\tau_{k2}) - \psi(\tau_{k1} + \tau_{k2}).
\end{aligned}$$

4. The feature distribution,

$$\begin{aligned}
\mathbb{E}_{\mathbf{A}} [\log p(\mathbf{A}_k|\sigma_A^2 I)] &= \mathbb{E}_{\mathbf{A}} \left[ \log \left( \frac{1}{(2\pi\sigma_A^2)^{D/2}} \exp \left( -\frac{1}{2\sigma_A^2} \mathbf{A}_k^T \mathbf{A}_k \right) \right) \right], \\
&= \frac{-D}{2} \log(2\pi\sigma_A^2) - \frac{1}{2\sigma_A^2} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_k^T).
\end{aligned}$$

5. The likelihood,

$$\begin{aligned}
&\mathbb{E}_{\mathbf{Z}, \mathbf{A}, \mathbf{S}} [\log p(\mathbf{X}_n|\mathbf{Z}_n, \mathbf{A}, \mathbf{S}, \sigma_n^2 I)] \\
&= \mathbb{E}_{\mathbf{Z}, \mathbf{A}, \mathbf{S}} \left[ -\frac{D}{2} \log(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} (\mathbf{X}_n - (\mathbf{Z}_n \odot \mathbf{S}_n) \mathbf{A}) (\mathbf{X}_n - (\mathbf{Z}_n \odot \mathbf{S}_n) \mathbf{A})^T \right], \\
&= -\frac{D}{2} \log(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} \left( \mathbf{X}_n \mathbf{X}_n^T - 2 \sum_{k=1}^K \nu_{nk} \mu_{nk} \bar{\boldsymbol{\phi}}_k \mathbf{X}_n^T \right. \\
&\quad \left. + 2 \sum_{k < k'} \nu_{nk} \mu_{nk} \nu_{nk'} \mu_{nk'} \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_{k'}^T + \sum_{k=1}^K \nu_{nk} (2\eta_{nk}^2 + \mu_{nk}^2) (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_k^T) \right),
\end{aligned}$$

where we use the fact that

$$\begin{aligned}
& \mathbb{E}_{\mathbf{Z}, \mathbf{S}, \mathbf{A}}[(\mathbf{Z}_{n\cdot} \odot \mathbf{S}_{n\cdot}) \mathbf{A} \mathbf{A}^T (\mathbf{Z}_{n\cdot} \odot \mathbf{S}_{n\cdot})^T] \\
&= \mathbb{E}_{\mathbf{Z}, \mathbf{A}, \mathbf{S}} \left[ \left( \sum_{k=1}^K (z_{nk} s_{nk}) \mathbf{A}_{k\cdot} \right) \left( \sum_{k=1}^K (z_{nk} s_{nk}) \mathbf{A}_{k\cdot} \right)^T \right], \\
&= \mathbb{E}_{\mathbf{Z}, \mathbf{A}, \mathbf{S}} \left[ \sum_{d=1}^D \left( \sum_{k=1}^K z_{nk} s_{nk}^2 A_{kd}^2 + \sum_{k, k': k' \neq k} (z_{nk} s_{nk}) (z_{nk'} s_{nk'}) \mathbf{A}_{kd} \mathbf{A}_{k'd} \right) \right], \\
&= \sum_{k=1}^K \nu_{nk} (2\eta_{nk}^2 + \mu_{nk}^2) (\text{tr}(\Phi_k) + \bar{\phi}_k \bar{\phi}_k^T) + 2 \sum_{k < k'} \nu_{nk} \mu_{nk} \nu_{nk'} \mu_{nk'} \bar{\phi}_k \bar{\phi}_{k'}^T.
\end{aligned}$$

6. Finally, for the entropy,

$$\begin{aligned}
H[q] &= -\mathbb{E}_q \log \left[ \prod_{k=1}^K q_{\tau_k}(\pi_k) \prod_{k=1}^K q_{\phi_k}(\mathbf{A}_{k\cdot}) \prod_{k=1}^K \prod_{n=1}^N q_{\nu_{nk}}(z_{nk}) \prod_{k=1}^K \prod_{n=1}^N q_{\mu_{nk}, \eta_{nk}}(s_{nk}) \right], \\
&= \sum_{k=1}^K \mathbb{E}_{\boldsymbol{\pi}}(-\log q_{\tau_k}(\pi_k)) + \sum_{k=1}^K \mathbb{E}_{\mathbf{A}}(-\log q_{\phi_k}(\mathbf{A}_{k\cdot})) \\
&\quad + \sum_{k=1}^K \sum_{n=1}^N \mathbb{E}_{\mathbf{Z}}(-\log q_{\nu_{nk}}(z_{nk})) + \sum_{k=1}^K \sum_{n=1}^N \mathbb{E}_{\mathbf{S}}(-\log q_{\mu_{nk}, \eta_{nk}}(s_{nk})),
\end{aligned}$$

where

$$\begin{aligned}
\mathbb{E}_{\boldsymbol{\pi}}(-\log q_{\tau_k}(\pi_k)) &= \log \left( \frac{\Gamma(\tau_{k1}) \Gamma(\tau_{k2})}{\Gamma(\tau_{k1} + \tau_{k2})} \right) - (\tau_{k1} - 1) \psi(\tau_{k1}) \\
&\quad - (\tau_{k2} - 1) \psi(\tau_{k2}) + (\tau_{k1} + \tau_{k2} - 2) \psi(\tau_{k1} + \tau_{k2}), \\
\mathbb{E}_{\mathbf{A}}(-\log q_{\phi_k}(\mathbf{A}_{k\cdot})) &= \frac{1}{2} \log((2\pi e)^D |\Phi_k|), \\
\mathbb{E}_{\mathbf{Z}}(-\log q_{\nu_{nk}}(z_{nk})) &= -\nu_{nk} \log \nu_{nk} - (1 - \nu_{nk}) \log(1 - \nu_{nk}). \\
\mathbb{E}_{\mathbf{S}}(-\log q_{\mu_{nk}, \eta_{nk}}(s_{nk})) &= \log(2e\eta_{nk}).
\end{aligned}$$

Collecting all the above computations together in Equation (21) gives us the variational lower bound on  $\log p_K(\mathbf{X}|\boldsymbol{\theta})$ :

$$\begin{aligned}
& \log p_K(\mathbf{X}|\boldsymbol{\theta}) \\
& \geq \sum_{k=1}^K \left[ \log \frac{\alpha}{K} + \left( \frac{\alpha}{K} - 1 \right) (\psi(\tau_{k1}) - \psi(\tau_{k1} + \tau_{k2})) \right] \\
& \quad - KN \log 2 - \sum_{k=1}^K \sum_{n=1}^N \left( |\mu_{nk}| + \eta_{nk} \exp \left( -\frac{|\mu_{nk}|}{\eta_{nk}} \right) \right) \\
& \quad + \sum_{k=1}^K \sum_{n=1}^N [\nu_{nk} \psi(\tau_{k1}) + (1 - \nu_{nk}) \psi(\tau_{k2}) - \psi(\tau_{k1} + \tau_{k2})] \\
& \quad + \sum_{k=1}^K \left[ \frac{-D}{2} \log(2\pi\sigma_A^2) - \frac{1}{2\sigma_A^2} (\text{tr}(\Phi_k) + \bar{\Phi}_k \bar{\Phi}_k^T) \right] \\
& \quad + \sum_{n=1}^N \left[ -\frac{D}{2} \log(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} \left( \mathbf{X}_n \cdot \mathbf{X}_n^T - 2 \sum_{k=1}^K \nu_{nk} \mu_{nk} \bar{\Phi}_k \mathbf{X}_n^T + 2 \sum_{k < k'} \nu_{nk} \mu_{nk} \nu_{nk'} \mu_{nk'} \bar{\Phi}_k \bar{\Phi}_{k'}^T \right. \right. \\
& \quad \quad \quad \left. \left. + \sum_{k=1}^K \nu_{nk} (2\eta_{nk}^2 + \mu_{nk}^2) (\text{tr}(\Phi_k) + \bar{\Phi}_k \bar{\Phi}_k^T) \right) \right] \\
& \quad + \sum_{k=1}^K \left[ \log \left( \frac{\Gamma(\tau_{k1}) \Gamma(\tau_{k2})}{\Gamma(\tau_{k1} + \tau_{k2})} \right) - (\tau_{k1} - 1) \psi(\tau_{k1}) - (\tau_{k2} - 1) \psi(\tau_{k2}) + (\tau_{k1} + \tau_{k2} - 2) \psi(\tau_{k1} + \tau_{k2}) \right] \\
& \quad + \sum_{k=1}^K \left[ \frac{1}{2} \log((2\pi e)^D |\Phi_k|) \right] + \sum_{k=1}^K \sum_{n=1}^N [-\nu_{nk} \log \nu_{nk} - (1 - \nu_{nk}) \log(1 - \nu_{nk}) + \log(2e\eta_{nk})].
\end{aligned}$$

### E.1.2 Parameter Updates

When we optimise the lower bound on  $\log p_K(\mathbf{X}|\boldsymbol{\theta})$  we perform coordinate-wise gradient ascent by cycling through the variational parameters and update them in turn. For most parameter we will be able to use the standard exponential family variational update from Equation (19). For some parameters we will need to compute the gradient and perform a local gradient ascent step. Throughout this section, we let  $c$  be a constant independent of the variable of interest that may change from line to line.

1. For the feature distribution at the optimal  $\bar{\Phi}_k$  and  $\Phi_k$ ,

$$\begin{aligned}
& \log q_{\Phi_k}(\mathbf{A}_k \cdot) \\
& = \mathbb{E}_{\mathbf{A}_{-k}, \mathbf{Z}, \mathbf{S}} [\log p(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta})] + c, \\
& = \mathbb{E}_{\mathbf{A}_{-k}, \mathbf{Z}, \mathbf{S}} \left[ \log p(\mathbf{A}_k | \sigma_A^2) + \sum_{n=1}^N p(\mathbf{X}_n | \mathbf{Z}_n \cdot, \mathbf{S}_n \cdot, \mathbf{A}, \sigma_n^2) \right] + c, \\
& = -\frac{1}{2\sigma_A^2} (\mathbf{A}_k \cdot \mathbf{A}_k^T) - \frac{1}{2\sigma_n^2} \sum_{n=1}^N \mathbb{E}_{\mathbf{A}_{-k}, \mathbf{Z}, \mathbf{S}} \left[ (\mathbf{X}_n \cdot - (\mathbf{Z}_n \cdot \odot \mathbf{S}_n \cdot) \mathbf{A}) (\mathbf{X}_n \cdot - (\mathbf{Z}_n \cdot \odot \mathbf{S}_n \cdot) \mathbf{A})^T \right] + c, \\
& = -\frac{1}{2} \left[ \mathbf{A}_k \cdot \left( \frac{1}{\sigma_A^2} + \frac{\sum_{n=1}^N \nu_{nk} (2\eta_{nk}^2 + \mu_{nk}^2)}{\sigma_n^2} \right) \mathbf{A}_k^T - 2\mathbf{A}_k \cdot \left( \frac{1}{\sigma_n^2} \sum_{n=1}^N \nu_{nk} \mu_{nk} \left( \mathbf{X}_n \cdot - \left( \sum_{l:l \neq k} \nu_{nl} \mu_{nl} \bar{\Phi}_l \right) \right) \right)^T \right] + c
\end{aligned}$$

Completing the square and using Equation (19) gives us that for the optimal parameter settings we must have

$$\begin{aligned}\log q_{\phi_k}(\mathbf{A}_k) &= -\frac{1}{2} (\mathbf{A}_k \cdot \mathbf{\Phi}_k^{-1} \mathbf{A}_k^T - 2\mathbf{A}_k \cdot \mathbf{\Phi}_k^{-1} \bar{\phi}_k^T) + c, \\ &= -\frac{1}{2} (\mathbf{A}_k - \bar{\phi}_k) \mathbf{\Phi}_k^{-1} (\mathbf{A}_k - \bar{\phi}_k)^T + c.\end{aligned}$$

hence the parameter updates are

$$\begin{aligned}\bar{\phi}_k &= \left[ \frac{1}{\sigma_n^2} \sum_{n=1}^N \nu_{nk} \mu_{nk} \left( \mathbf{X}_{n\cdot} - \left( \sum_{l:l \neq k} \nu_{nl} \mu_{nl} \bar{\phi}_l \right) \right) \right] \left( \frac{1}{\sigma_A^2} + \frac{\sum_{n=1}^N \nu_{nk} (2\eta_{nk}^2 + \mu_{nk}^2)}{\sigma_n^2} \right)^{-1}, \\ \mathbf{\Phi}_k &= \left( \frac{1}{\sigma_A^2} + \frac{\sum_{n=1}^N \nu_{nk} (2\eta_{nk}^2 + \mu_{nk}^2)}{\sigma_n^2} \right)^{-1} I.\end{aligned}$$

2. For the feature state distribution at the optimal  $\nu_{nk}$ ,

$$\begin{aligned}\log q_{\nu_{nk}}(z_{nk}) &= \mathbb{E}_{\boldsymbol{\pi}, \mathbf{A}, \mathbf{Z}_{-nk}, \mathbf{S}} [\log p(\mathbf{W}, \mathbf{X} | \boldsymbol{\theta})] + c, \\ &= \mathbb{E}_{\boldsymbol{\pi}, \mathbf{A}, \mathbf{Z}_{-nk}, \mathbf{S}} [\log p(z_{nk} | \pi_k) + \log p(X_n | \mathbf{Z}_n, \mathbf{A}, \mathbf{S}, \sigma_n^2)] + c,\end{aligned}$$

where

$$\mathbb{E}_{\boldsymbol{\pi}} [\log p(z_{nk} | \pi_k)] = z_{nk} [\psi(\tau_{k1}) - \psi(\tau_{k2})] + \psi(\tau_{k2}) - \psi(\tau_{k1} + \tau_{k2}),$$

and

$$\begin{aligned}\mathbb{E}_{\mathbf{A}, \mathbf{Z}_{-nk}, \mathbf{S}} [\log p(\mathbf{X}_n | \mathbf{Z}_n, \mathbf{A}, \sigma_n^2)] &= \mathbb{E}_{\mathbf{A}, \mathbf{Z}_{-nk}, \mathbf{S}} \left[ -\frac{1}{2\sigma_n^2} (\mathbf{X}_n - (\mathbf{Z}_n \odot \mathbf{S}_n) \mathbf{A}) (\mathbf{X}_n - (\mathbf{Z}_n \odot \mathbf{S}_n) \mathbf{A})^T \right] + c, \\ &= -\frac{1}{2\sigma_n^2} \left[ -2z_{nk} \mu_{nk} \bar{\phi}_k \mathbf{X}_n^T + z_{nk} (2\eta_{nk}^2 + \mu_{nk}^2) (\text{tr}(\mathbf{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T) + 2z_{nk} \mu_{nk} \bar{\phi}_k \left( \sum_{l:l \neq k} \nu_{nl} \mu_{nl} \bar{\phi}_l^T \right) \right] + c.\end{aligned}$$

Therefore

$$\begin{aligned}\log q_{\nu_{nk}}(z_{nk}) &= z_{nk} \left[ \psi(\tau_{k1}) - \psi(\tau_{k2}) - \frac{1}{2\sigma_n^2} ((2\eta_{nk}^2 + \mu_{nk}^2) (\text{tr}(\mathbf{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T) \right. \\ &\quad \left. - 2\mu_{nk} \bar{\phi}_k \mathbf{X}_n^T + 2\mu_{nk} \bar{\phi}_k \left( \sum_{l:l \neq k} \nu_{nl} \mu_{nl} \bar{\phi}_l^T \right) \right] + c.\end{aligned}$$

From the canonical parameterisation of the Bernoulli distribution, we get that

$$\begin{aligned}\log \frac{\nu_{nk}}{1 - \nu_{nk}} &= \psi(\tau_{k1}) - \psi(\tau_{k2}) - \frac{2\eta_{nk}^2 + \mu_{nk}^2}{2\sigma_n^2} (\text{tr}(\mathbf{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T) + \frac{\mu_{nk}}{\sigma_n^2} \bar{\phi}_k \left( \mathbf{X}_n^T - \left( \sum_{l:l \neq k} \mu_{nl} \nu_{nl} \bar{\phi}_l^T \right) \right), \\ &\equiv \vartheta.\end{aligned}$$

which gives us the update

$$\nu_{nk} = \frac{1}{1 + e^{-\vartheta}}.$$

3. The updates  $\tau_{k1}$  and  $\tau_{k2}$  are only dependent on  $\mathbf{Z}$  and  $\boldsymbol{\pi}$  and hence are exactly the same as for the linear-Gaussian model. We refer to Section C.2 for details.
4. Finally, because the Laplace prior on the signal matrix  $\mathbf{S}$  is not in the exponential family we cannot use the standard variational Bayes update formula as the posterior will not be in the exponential family anymore. Hence in order to optimise  $\mu_{nk}$  and  $\eta_{nk}$  we perform a gradient ascent step on the variational lower bound. Let us denote the variational lower bound which we derived earlier with  $\mathcal{L}$ . The components of  $\mathcal{L}$  that depend on  $\mu_{nk}$  and  $\eta_{nk}$  are

$$-\frac{1}{2\sigma_n^2} \left( -2\nu_{nk}\mu_{nk}\bar{\boldsymbol{\phi}}_k\mathbf{X}_n^T + 2(\nu_{nk}\mu_{nk}\bar{\boldsymbol{\phi}}_k) \left( \sum_{k':k' \neq k} \nu_{nk'}\mu_{nk'}\bar{\boldsymbol{\phi}}_{k'}^T \right) + \nu_{nk}(2\eta_{nk}^2 + \mu_{nk}^2) (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k\bar{\boldsymbol{\phi}}_k^T) \right) - |\mu_{nk}| - \eta_{nk} \exp\left(-\frac{|\mu_{nk}|}{\eta_{nk}}\right) + \log(2e\eta_{nk}).$$

It is now straightforward to compute the derivative of  $\mathcal{L}$  with respect to the mean parameter of the Laplace distribution

$$\frac{\partial \mathcal{L}}{\partial \mu_{nk}} = \frac{1}{\sigma_n^2} \left( \nu_{nk}\bar{\boldsymbol{\phi}}_k\mathbf{X}_n^T - (\nu_{nk}\bar{\boldsymbol{\phi}}_k) \left( \sum_{k' \neq k} \nu_{nk'}\mu_{nk'}\bar{\boldsymbol{\phi}}_{k'}^T \right) - \nu_{nk}\mu_{nk} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k\bar{\boldsymbol{\phi}}_k^T) \right) - \text{sign}(\mu_{nk}) \left( 1 - \exp\left(-\frac{|\mu_{nk}|}{\eta_{nk}}\right) \right),$$

and the derivative of  $\mathcal{L}$  with respect to the scale of the Laplace distribution

$$\frac{\partial \mathcal{L}}{\partial \eta_{nk}} = -\frac{2\nu_{nk}\eta_{nk}}{\sigma_n^2} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k\bar{\boldsymbol{\phi}}_k^T) - \exp\left(-\frac{|\mu_{nk}|}{\eta_{nk}}\right) - \frac{|\mu_{nk}|}{\eta_{nk}} \exp\left(-\frac{|\mu_{nk}|}{\eta_{nk}}\right) + \frac{1}{\eta_{nk}}.$$

We can then numerically optimise  $\mu_{nk}$  and  $\eta_{nk}$ .

## E.2 The Infinite Variational Approach

We presented the iICA model at the start of Section E. Recall that our goal is to compute the log posterior

$$\log p(\mathbf{W}|\mathbf{X}, \boldsymbol{\theta}) = \log p(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta}) - \log p(\mathbf{X}|\boldsymbol{\theta}),$$

but this is intractable to compute. We use a variational approximation based on the truncated stick-breaking process as described in Section 5 and use the stick-breaking variables  $\mathbf{v}$  instead of  $\boldsymbol{\mu}$  as before. Our mean field variational distribution is now

$$q(\mathbf{W}) = q_{\boldsymbol{\tau}}(\mathbf{v})q_{\boldsymbol{\phi}}(\mathbf{A})q_{\boldsymbol{\nu}}(\mathbf{Z})q_{\boldsymbol{\mu}, \boldsymbol{\eta}}(\mathbf{S}).$$

where

- $q_{\boldsymbol{\tau}}(v_k) = \text{Beta}(v_k; \tau_{k1}, \tau_{k2})$ ,
- $q_{\boldsymbol{\phi}_k}(\mathbf{A}_{k\cdot}) = \text{Normal}(\mathbf{A}_{k\cdot}; \bar{\boldsymbol{\phi}}_k, \boldsymbol{\Phi}_k)$ ,
- $q_{\nu_{nk}}(z_{nk}) = \text{Bernoulli}(z_{nk}; \nu_{nk})$ ,
- $q_{\mu_{nk}, \eta_{nk}}(s_{nk}) = \text{Laplace}(s_{nk}; \mu_{nk}, \eta_{nk})$  where  $\mu_{nk}$  is the mean and  $\eta_{nk}$  is the scale parameter of the Laplace distribution.



As with the finite approach, inference involves optimising  $\boldsymbol{\tau}$ ,  $\boldsymbol{\phi}$ ,  $\boldsymbol{\nu}$ ,  $\boldsymbol{\mu}$  and  $\boldsymbol{\eta}$  to minimise the KL divergence  $D(q||p)$ , or equivalently to maximise the lower bound on  $p(\mathbf{X}|\boldsymbol{\theta})$

$$\mathbb{E}_q[\log(p(\mathbf{X}, \mathbf{W}|\boldsymbol{\theta})) + H[q]].$$

Although the update equations for this approximation are not as straightforward as in the finite approach, we can reuse many of the computations we did for the linear-Gaussian and beta-Bernoulli iICA approximation.

### E.2.1 Variational Lower Bound

As in the finite approach, we first derive an expression for the variational lower bound. However, as in the linear-Gaussian model, parts of our model are no longer in the exponential family and require nontrivial computations. We expand the lower bound on  $\log p(\mathbf{X}|\boldsymbol{\theta})$  into its components

$$\begin{aligned} \log p(\mathbf{X}|\boldsymbol{\theta}) &\geq \mathbb{E}_{\mathbf{W}}[\log p(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta})] + H[q], \\ &= \sum_{k=1}^K \mathbb{E}_{\boldsymbol{\pi}}[\log p(v_k|\alpha)] + \sum_{k=1}^K \sum_{n=1}^N \mathbb{E}_{\boldsymbol{\pi}, \mathbf{Z}}[\log p(z_{nk}|\tau_k)] + \sum_{k=1}^K \sum_{n=1}^N \mathbb{E}_{\mathbf{S}}[\log p(s_{nk})] \\ &\quad + \sum_{k=1}^K \mathbb{E}_{\mathbf{A}}[\log p(\mathbf{A}_k|\sigma_A^2 I)] + \sum_{n=1}^N \mathbb{E}_{\mathbf{Z}, \mathbf{A}, \mathbf{S}}[\log p(\mathbf{X}_n|\mathbf{Z}_n, \mathbf{A}, \mathbf{S}_n, \sigma_n^2 I)] + H[q], \end{aligned} \quad (22)$$

where the expectation are computed with respect to the variational distribution  $q$ . From Appendix E.1.1 we know how to compute all the expectations for the ICA likelihood; together with the theory from Section 5.1 and Appendix D.1 on the lower bound for the infinite variational approximation to the linear-Gaussian model, we can expand the expectation in Equation (21):

$$\begin{aligned} \log p(\mathbf{X}|\boldsymbol{\theta}) &\geq \sum_{k=1}^K [\log \alpha + (\alpha - 1)(\psi(\tau_{k1}) - \psi(\tau_{k1} + \tau_{k2}))] \\ &\quad + \sum_{k=1}^K \sum_{n=1}^N \left[ \nu_{nk} \left( \sum_{m=1}^k \psi(\tau_{k2}) - \psi(\tau_{k1} + \tau_{k2}) \right) + (1 - \nu_{nk}) \mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right] \right] \\ &\quad - KN \log 2 - \sum_{k=1}^K \sum_{n=1}^N \left( |\mu_{nk}| + \eta_{nk} \exp \left( -\frac{|\mu_{nk}|}{\eta_{nk}} \right) \right) \\ &\quad + \sum_{k=1}^K \left[ \frac{-D}{2} \log(2\pi\sigma_A^2) - \frac{1}{2\sigma_A^2} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\Phi}}_k \bar{\boldsymbol{\Phi}}_k^T) \right] \\ &\quad + \sum_{n=1}^N \left[ -\frac{D}{2} \log(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} \left( \mathbf{X}_n \mathbf{X}_n^T - 2 \sum_{k=1}^K \nu_{nk} \mu_{nk} \bar{\boldsymbol{\Phi}}_k \mathbf{X}_n^T + 2 \sum_{k < k'} \nu_{nk} \mu_{nk} \nu_{nk'} \mu_{nk'} \bar{\boldsymbol{\Phi}}_k \bar{\boldsymbol{\Phi}}_{k'}^T \right. \right. \\ &\quad \left. \left. + \sum_{k=1}^K \nu_{nk} (2\eta_{nk}^2 + \mu_{nk}^2) (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\Phi}}_k \bar{\boldsymbol{\Phi}}_k^T) \right) \right] \\ &\quad + \sum_{k=1}^K \left[ \log \left( \frac{\Gamma(\tau_{k1})\Gamma(\tau_{k2})}{\Gamma(\tau_{k1} + \tau_{k2})} \right) - (\tau_{k1} - 1)\psi(\tau_{k1}) - (\tau_{k2} - 1)\psi(\tau_{k2}) + (\tau_{k1} + \tau_{k2} - 2)\psi(\tau_{k1} + \tau_{k2}) \right] \\ &\quad + \sum_{k=1}^K \left[ \frac{1}{2} \log((2\pi e)^D |\boldsymbol{\Phi}_k|) \right] + \sum_{k=1}^K \sum_{n=1}^N [-\nu_{nk} \log \nu_{nk} - (1 - \nu_{nk}) \log(1 - \nu_{nk}) + \log(2e\eta_{nk})]. \end{aligned}$$

The term  $\mathbb{E}_{\mathbf{v}} \left[ \log \left( 1 - \prod_{m=1}^k v_m \right) \right]$  we left unevaluated. It can be evaluated either with the multinomial approach from Section 5.1 (recommended) or the Taylor series approach from Appendix D.1

## E.2.2 Parameter Updates

For the infinite variational approximation to the iICA model, we need to update the parameters  $\boldsymbol{\tau}$ ,  $\boldsymbol{\phi}$ ,  $\boldsymbol{\nu}$  and  $\boldsymbol{\mu}$ . For the parameters  $\boldsymbol{\phi}$  and  $\boldsymbol{\mu}$ , the updates are exactly the same as for the finite iICA approximation in Section C.2. For the parameters of  $\mathbf{Z}$  we update  $\nu_{nk}$  in Bernoulli( $z_{nk}; \nu_{nk}$ ) as

$$\nu_{nk} = \frac{1}{1 + e^{-\vartheta}}$$

where

$$\begin{aligned} \vartheta = & \sum_{i=1}^k (\psi(\tau_{i1}) - \psi(\tau_{i1} + \tau_{i2})) - \mathbb{E}_{\mathbf{v}}[\log(1 - \prod_{i=1}^k v_i)] \\ & - \frac{2\eta_{nk}^2 + \mu_{nk}^2}{2\sigma_n^2} (\text{tr}(\boldsymbol{\Phi}_k) + \bar{\boldsymbol{\phi}}_k \bar{\boldsymbol{\phi}}_k^T) + \frac{\mu_{nk}}{\sigma_n^2} \bar{\boldsymbol{\phi}}_k \left( \mathbf{X}_n^T - \left( \sum_{l:l \neq k} \mu_{nl} \nu_{nl} \bar{\boldsymbol{\phi}}_l^T \right) \right). \end{aligned}$$

We leave the term  $\mathbb{E}_{\mathbf{v}}[\log(1 - \prod_{i=1}^k v_i)]$  unevaluated because the choice of how to approximate it does not change the form of the update.

Finally, to update  $\tau_{k1}$  and  $\tau_{k2}$  in  $\text{Beta}(v_k; \tau_{k1}, \tau_{k2})$  we use the multinomial lower bound for  $\mathbb{E}_{\mathbf{v}}[\log(1 - \prod_{i=1}^k v_i)]$  and compute  $q_{ki}$  according to Equation (10). As in the linear-Gaussian model, the updates for  $\tau_{k1}$  and  $\tau_{k2}$  have the closed form

$$\begin{aligned} \tau_{k1} &= \alpha + \sum_{m=k}^K \sum_{n=1}^N \nu_{nm} + \sum_{m=k+1}^K \left( N - \sum_{n=1}^N \nu_{nm} \right) \binom{m}{i=k+1} q_{mi} \\ \tau_{k2} &= 1 + \sum_{m=k}^K \left( N - \sum_{n=1}^N \nu_{nm} \right) q_{mk}. \end{aligned}$$

## F Alternate Bounds for the Infinite Approximation

Recall that when giving a bound for how close  $m_K(\mathbf{X})$  is to  $m(\mathbf{X})$ , we must bound

$$1 - \mathbb{E} \left( \left[ \prod_{i=K+1}^{\infty} (1 - \pi_i) \right]^N \right).$$

We gave one approach to bounding this in Section 6. This approach applied Jensen's inequality as follows

$$1 - \mathbb{E} \left( \left[ \prod_{i=K+1}^{\infty} (1 - \pi_i) \right]^N \right) \leq 1 - \left( \mathbb{E} \left[ \prod_{i=K+1}^{\infty} (1 - \pi_i) \right] \right)^N$$

and then relied on the beta process representation from Thibaux and Jordan (2007) to bound the inner term.

Another approach to applying Jensen's inequality is to write

$$1 - \mathbb{E} \left( \left[ \prod_{i=K+1}^{\infty} (1 - \pi_i) \right]^N \right) \leq 1 - \exp \left( N \sum_{i=K+1}^{\infty} \mathbb{E} \log \left( 1 - \prod_{j=1}^i v_j \right) \right), \quad (23)$$

where  $v_j$  are the stick-breaking weights. This section derives a heuristic bound and a principled bound both based on this alternate application of Jensen's inequality.

To use Equation (23), we need to evaluate or bound  $\sum_{i=K+1}^{\infty} \mathbb{E} \log \left( 1 - \prod_{j=1}^i v_j \right)$ . We expand the expectation using the Taylor series approximation of Appendix D.1, noting that the expectation  $\mathbb{E}[v^r]$  of a Beta( $\alpha, 1$ ) random variables is  $\frac{\alpha}{\alpha+r}$ :

$$\begin{aligned} \sum_{i=K+1}^{\infty} \mathbb{E} \log \left( 1 - \prod_{j=1}^i v_j \right) &= - \sum_{i=K+1}^{\infty} \sum_{r=1}^{\infty} \frac{1}{r} \prod_{j=1}^i \frac{\alpha}{\alpha+r} \\ &= - \sum_{r=1}^{\infty} \frac{1}{r} \sum_{i=K+1}^{\infty} \left( \frac{\alpha}{\alpha+r} \right)^i \\ &= - \sum_{r=1}^{\infty} \frac{1}{r^2} \frac{\alpha^{K+1}}{(\alpha+r)^K}. \end{aligned}$$

Substituting the expression above into the original bound, we get

$$1 - \mathbb{E} \left( \left[ \prod_{i=K+1}^{\infty} (1 - \pi_i) \right]^N \right) \leq 1 - \exp \left( -N \sum_{r=1}^{\infty} \frac{1}{r^2} \frac{\alpha^{K+1}}{(\alpha+r)^K} \right).$$

For any truncation level of the sum, we do not necessarily have a true bound, but we find empirically that it is very close to the true truncation bound.

To get a strict bound, we can write<sup>6</sup>

$$\begin{aligned} \sum_{r=1}^{\infty} \frac{1}{r^2} \frac{\alpha^{K+1}}{(\alpha+r)^K} &\leq \frac{\alpha^{K+1}}{(\alpha+1)^K} + \int_1^{\infty} \frac{1}{r^2} \frac{\alpha^{K+1}}{(\alpha+r)^K} \\ &= \frac{\alpha^{K+1}}{(\alpha+1)^K} + \int_0^1 \frac{\alpha^{K+1}}{(\alpha + \frac{1}{t})^K} \\ &= \frac{\alpha^{K+1}}{(\alpha+1)^K} + \frac{\alpha^{K+1}}{K+1} F(K, K+1; K+2; -\alpha) \end{aligned} \quad (24)$$

The first line applies the integral inequality, where we have included the first term to ensure that we have an upper bound. The last line substitutes  $t = \frac{1}{r}$  into the integral and evaluates. Next, we apply the reflection law of hypergeometric functions. The reflection law states

$$\frac{1}{(1-z)^a} F \left( a, b; c; \frac{-z}{1-z} \right) = F(a, c-b; c; z).$$

<sup>6</sup>We thank Professor John Lewis (MIT) for his insights in deriving this bound.

Now we can simplify the hypergeometric function in Equation 24 by expanding it into its sum:

$$\begin{aligned}
\sum_{r=1}^{\infty} \frac{1}{r^2} \frac{\alpha^{K+1}}{(\alpha+r)^K} &\leq \frac{\alpha^{K+1}}{(\alpha+1)^K} + \frac{\alpha^{K+1}}{(K+1)(\alpha+1)^{K+1}} F(2, K+1; K+2; \frac{\alpha}{\alpha+1}) \\
&= \frac{\alpha^{K+1}}{(\alpha+1)^K} + \frac{\alpha^{K+1}}{(\alpha+1)^{K+1}} \sum_{j=0}^{\infty} \left(\frac{\alpha}{\alpha+1}\right)^j \frac{1+j}{K+1+j} \\
&\leq \frac{\alpha^{K+1}}{(\alpha+1)^K} + \frac{\alpha^{K+1}}{(\alpha+1)^{K+1}} \sum_{j=0}^{\infty} \left(\frac{\alpha}{\alpha+1}\right)^j \\
&= 2(\alpha+1) \left(\frac{\alpha}{\alpha+1}\right)^{K+1}
\end{aligned}$$

which we can plug into our original expression to get

$$1 - \exp\left(-N \sum_{r=1}^{\infty} \frac{1}{r^2} \frac{\alpha^{K+1}}{(\alpha+r)^K}\right) \leq 1 - \exp\left(-2N(\alpha+1) \left(\frac{\alpha}{\alpha+1}\right)^{K+1}\right).$$

We note this bound is very similar to the bound derived using the Levy-Khintchine approach from Section 6:

$$1 - \exp\left(-N\alpha \left(\frac{\alpha}{1+\alpha}\right)^K\right).$$