

---

# Dependent Multinomial Models Made Easy: Supplementary Material

---

<b>Scott W. Linderman*</b> Harvard University Cambridge, MA 02138 swl@seas.harvard.edu	<b>Matthew J. Johnson*</b> Harvard University Cambridge, MA 02138 mattjj@csail.mit.edu	<b>Ryan P. Adams</b> Twitter & Harvard University Cambridge, MA 02138 rpa@seas.harvard.edu
---	---	---

## A Symmetry breaking in $\pi_{\text{SB}}$

The stick-breaking map  $\pi_{\text{SB}} : \mathbb{R}^K \rightarrow [0, 1]^K$  is asymmetric in the sense that while the logistic map  $\pi_{\text{LN}} : \mathbb{R}^K \rightarrow [0, 1]^K$  can be written as the composition of an coordinate-wise logistic function and a normalization,

$$\pi_{\text{LN}}(\psi) = \left( \pi_{\text{LN}}^{(1)}(\psi), \dots, \pi_{\text{LN}}^{(K)}(\psi) \right) \quad \pi_{\text{LN}}^{(k)}(\psi) = \frac{e^{\psi_k}}{\sum_{j=1}^K e^{\psi_j}}, \quad (1)$$

the stick-breaking map does not have such a coordinate-wise separation:

$$\pi_{\text{SB}}(\psi) = \left( \pi_{\text{LN}}^{(1)}(\psi), \dots, \pi_{\text{SB}}^{(K)}(\psi) \right) \quad \pi_{\text{SB}}^{(k)}(\psi) = \sigma(\psi_k) \left( \sum_{j < k} \sigma(\psi_j) \right). \quad (2)$$

In particular,  $\pi_{\text{SB}}$  does not preserve permutation symmetries in the density  $p(\psi)$ , so that while for any permutation matrix  $P$  we have

$$p(P\psi) = p(\psi) \implies p(P\pi_{\text{LN}}(\psi)) = p(\pi_{\text{LN}}(\psi)) \quad (3)$$

the same does not hold when  $\pi_{\text{LN}}$  is replaced with  $\pi_{\text{SB}}$ . As a result, the stick-breaking model used in this paper (and in Khan et al. [1]) yields priors (and posteriors) that are not invariant to relabeling of the entries of the corresponding multinomial parameter or the multinomial counts themselves. See Figure 1 and compare it to Figure 1 of the main text.

This symmetry breaking may be undesirable in some cases, but in the models we have studied so far (and in those studied in Khan et al. [1]) the effect does not seem detrimental in terms of learning informative correlation structures or in terms of model predictions. For example, in the correlated topic model (CTM) studied in Section 3, the model is unidentifiable up to permutation on the topic labels and therefore breaking this symmetry does not reduce its representational capacity. For models in which the counts from multinomials with correlated parameters are observed directly, such as in the models of Sections 4 and 5, based on the experiments in this paper the loss of symmetry does not seem to impact performance while the inference advantages are significant. See also the discussion in Khan et al. [1, Section 3].

## B Transforming between $p(\psi)$ and $p(\pi)$

Since the mapping between  $\pi$  and  $\psi$  is invertible, we can compute the distribution on  $\pi$  that is implied by a Gaussian distribution on  $\psi$ . Assume  $\psi \sim \mathcal{N}(\mu, \Sigma)$ . Then,

$$p(\pi \mid \mu, \Sigma) = \mathcal{N}(\pi_{\text{SB}}^{-1}(\pi) \mid \mu, \Sigma) \left| \frac{d\psi}{d\pi} \right|$$

---

\*These authors contributed equally.

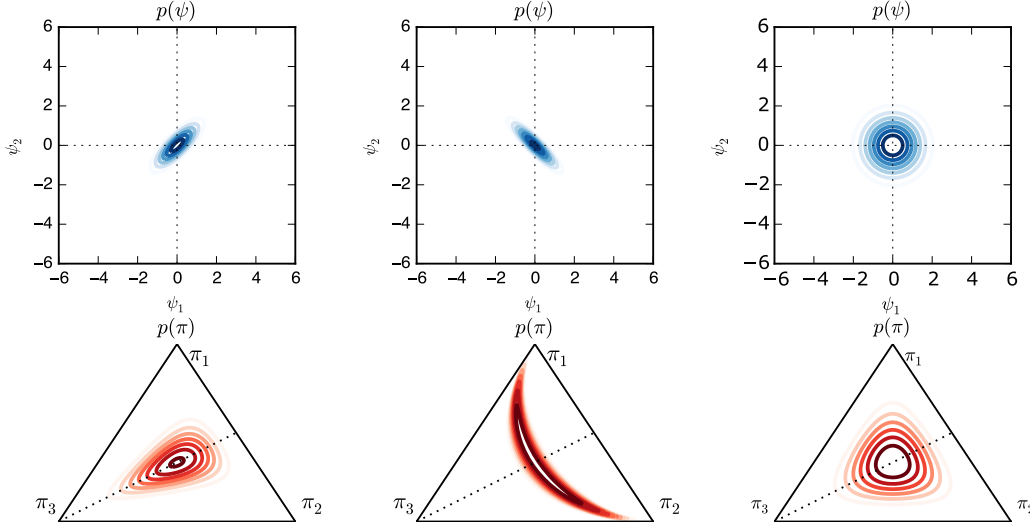


Figure 1: Correlated 2D Gaussian priors on  $\psi$  and their implied densities on  $\pi_{\text{LN}}(\psi)$ . Compare to Figure 1 of the main text, which shows an analogous plot of implied densities on  $\pi_{\text{SB}}(\psi)$ .

From above, we have

$$\psi_1 = \sigma^{-1}(\pi_1), \quad \psi_2 = \sigma^{-1}\left(\frac{\pi_2}{1 - \pi_1}\right), \quad \dots, \quad \psi_k = \sigma^{-1}\left(\frac{\pi_k}{1 - \sum_{j < k} \pi_j}\right).$$

Let

$$g(x) = \left. \frac{d\sigma^{-1}(x)}{dx} \right|_{x=x} = \frac{d}{dx} \log\left(\frac{x}{1-x}\right) = \frac{1}{x} + \frac{1}{1-x} = \frac{1}{x(1-x)}.$$

Then,

$$\frac{\partial \psi_1}{\partial \pi_1} = g(\pi_1), \quad \frac{\partial \psi_k}{\partial \pi_k} = g\left(\frac{\pi_k}{1 - \sum_{j < k} \pi_j}\right) \frac{1}{1 - \sum_{j < k} \pi_j}, \quad \frac{\partial \psi_k}{\partial \pi_{j > k}} = 0.$$

Since the Jacobian of the inverse transformation is lower triangular, its determinant is simply the product of its diagonal entries,

$$\begin{aligned} \left| \frac{d\psi}{d\pi} \right| &= \prod_{k=1}^K \left[ g\left(\frac{\pi_k}{1 - \sum_{j < k} \pi_j}\right) \frac{1}{1 - \sum_{j < k} \pi_j} \right] \\ &= \prod_{k=1}^K \left[ \frac{1 - \sum_{j < k} \pi_j}{\pi_k} \frac{1 - \sum_{j < k} \pi_j}{1 - \sum_{j < k} \pi_j - \pi_k} \frac{1}{1 - \sum_{j < k} \pi_j} \right] \\ &= \prod_{k=1}^K \left[ \frac{1 - \sum_{j=1}^{k-1} \pi_j}{\pi_k (1 - \sum_{j=1}^k \pi_j)} \right] \end{aligned}$$

Thus, the final density is,

$$p(\pi | \mu, \Sigma) = \mathcal{N}(\pi_{\text{SB}}^{-1}(\pi) | \mu, \Sigma) \cdot \prod_{k=1}^K \left[ \frac{1 - \sum_{j=1}^{k-1} \pi_j}{\pi_k (1 - \sum_{j=1}^k \pi_j)} \right].$$

Now, suppose we are given a Dirichlet distribution,  $\pi \sim \text{Dir}(\pi | \alpha)$ , and we wish to compute the density on  $\psi$ . We have,

$$\begin{aligned} p(\psi | \alpha) &= \text{Dir}(\pi_{\text{SB}}(\psi) | \alpha) \cdot \left| \frac{d\pi}{d\psi} \right| \\ &= \text{Dir}(\pi_{\text{SB}}(\psi) | \alpha) \cdot \prod_{k=1}^K \left[ \frac{\pi_k (1 - \sum_{j=1}^k \pi_j)}{1 - \sum_{j=1}^{k-1} \pi_j} \right], \end{aligned}$$

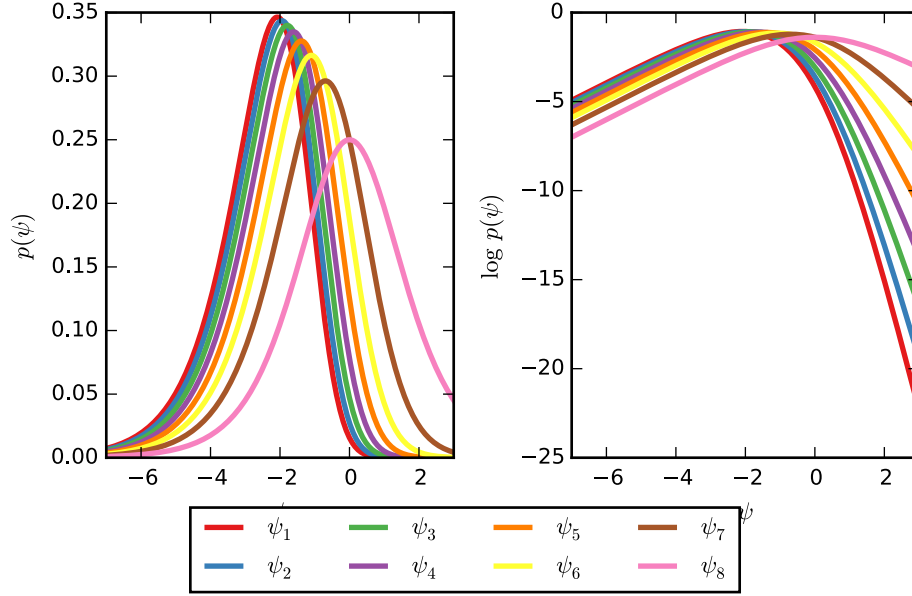


Figure 2: Density and log density of  $p(\psi | \alpha = 1)$ , the density on  $\psi$  implied by a  $K = 9$  dimensional symmetric Dirichlet density on  $\pi$  with parameter  $\alpha = 1$ .

where we have used the fact that the Jacobian of the inverse transformation is simply the inverse of the Jacobian of the forward transformation. We simply need to rewrite the Jacobian in terms of  $\psi$  rather than  $\pi$ . Note that  $1 - \sum_{j < k} \pi_j$  is the length of the remaining stick and  $\sigma(\psi_k)$  is the fraction of the remaining “stick” allocated to  $\pi_k$ . Thus, the remaining stick length is equal to,

$$1 - \sum_{j < k} \pi_j \equiv \prod_{j < k} (1 - \sigma(\psi_j)) \equiv \prod_{j < k} \sigma(-\psi_j).$$

Moreover,  $\pi_k = \sigma(\psi_k)(1 - \sum_{j < k} \pi_j) = \sigma(\psi_k) \prod_{j < k} \sigma(-\psi_j)$ . Thus,

$$\begin{aligned} p(\psi | \alpha) &= \text{Dir}(\pi_{\text{SB}}(\psi) | \alpha) \cdot \prod_{k=1}^K \left[ \frac{(\sigma(\psi_k) \prod_{j < k} \sigma(-\psi_j)) \left( \prod_{j \leq k} \sigma(-\psi_j) \right)}{\prod_{j < k} \sigma(-\psi_j)} \right], \\ &= \text{Dir}(\pi_{\text{SB}}(\psi) | \alpha) \cdot \prod_{k=1}^K \left[ \sigma(\psi_k) \prod_{j \leq k} \sigma(-\psi_j) \right], \end{aligned}$$

Expanding the Dirichlet distribution and substituting  $\psi$  for  $\pi$ , we conclude that,

$$p(\psi | \alpha) = \frac{1}{B(\alpha)} \prod_{k=1}^{K-1} \sigma(\psi_k)^{\alpha_k} \cdot \sigma(-\psi_k)^{\sum_{j=k+1}^K \alpha_j}.$$

This factorized form is unsurprising given that the Dirichlet distribution can be written as a stick-breaking product of beta distributions in the same way that the multinomial can be written as a product of binomials. Each term in the product above corresponds to the transformed beta distribution over  $\tilde{\pi}_k$ .

Figure 2 shows the marginal densities on  $\psi_k$  implied by a  $K = 9$  dimensional symmetric Dirichlet prior on  $\pi$  with  $\alpha = 1$ . The densities of  $\psi_k$  become increasingly skewed for small values of  $k$ , but they are still well approximate by a Gaussian distribution. In order to approximate a uniform distribution, we numerically compute the mean and variance of these densities to set the parameters of a diagonal Gaussian distribution.

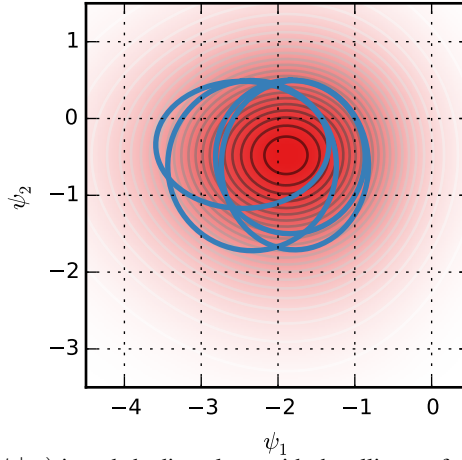


Figure 3: Marginal density,  $p(\psi | x)$  in red shading along with the ellipses of multivariate normal conditional distribution  $p(\psi | x, \omega)$  for 4 steps of the Gibbs sampler. In Gaussian models where we aim to predict  $\psi_{\text{test}}$  on test data, there are substantial gains to be had from making marginal predictions of  $\psi_{\text{test}} | x, \omega$ , integrating out  $\psi_{\text{train}}$ . The key is that the conditional densities overlap substantially with the marginal density.

## C Marginal Predictions with the Augmented Model

One of the primary advantages offered by the Pólya-gamma augmentation is the ability to make marginal predictions about  $\psi_{\text{test}} | x, \omega$ , integrating out the value of  $\psi_{\text{train}}$ . For example, in the GP multinomial regression models described in the main text, the methods were evaluated on the accuracy of their predictions about future name probabilities, which were functions of  $\psi_{\text{test}}$ . When  $p(\psi_{\text{train}})$  and  $p(\psi_{\text{test}} | \psi_{\text{train}})$  are both Gaussian, we can integrate out the latent training variables in order to predict their test values. In a latent Gaussian-multinomial model, the posterior distribution over those latent training variables is non-Gaussian, but after Pólya-gamma augmentation, it is rendered Gaussian.

With the augmentation, we can write

$$p(\psi_{\text{test}} | x) \approx \frac{1}{M} \sum_{m=1}^M \int p(\psi_{\text{test}} | \psi_{\text{train}}) p(\psi_{\text{train}} | x, \omega^{(m)}) d\psi_{\text{train}} \quad \omega^{(m)} \sim p(\omega | x),$$

and perform Monte Carlo integration over  $\omega$  in order to compute the predictive distribution. By contrast, in the standard formulation we must perform Monte Carlo integration over  $\psi$ ,

$$p(\psi_{\text{test}} | x) = \frac{1}{M} \sum_{m=1}^M p(\psi_{\text{test}} | \psi_{\text{train}}^{(m)}) \quad \psi_{\text{train}}^{(m)} \sim p(\psi_{\text{train}} | x).$$

Why does the augmented model confer a predictive advantage? It does not come from performing Monte Carlo integration over a smaller dimension since  $\omega$  and  $\psi_{\text{train}}$  are of the same size. Instead, it comes from the ability of the conjugate Gibbs sampler to efficiently mix over  $\psi$  and  $\omega$ , and from the ability of a single sample of  $\omega$  to render a conditional Gaussian distribution over  $\psi$  that captures much of the volume of the true marginal distribution.

This latter point is illustrated in Figure 3. The red shading shows the true marginal density of  $\psi$  and the blue ellipses show the conditional density for a fixed value of  $\omega$ . Each ellipse capture a significant amount of the marginal distribution, indicating that with a single sample of  $\omega$  we can integrate over a substantial amount of the uncertainty in  $\psi$ . This example is only for a  $K = 3$  dimensional multinomial observation, but this intuition should extend to higher dimensions in which the advantages of analytical integration should be more readily apparent.

## References

- [1] Mohammad E Khan, Shakir Mohamed, Benjamin M Marlin, and Kevin P Murphy. A stick-breaking likelihood for categorical data analysis with latent Gaussian models. In *International Conference on Artificial Intelligence and Statistics*, pages 610–618, 2012.