

**SUPPLEMENT TO:
SHOPPER: A PROBABILISTIC MODEL OF CONSUMER CHOICE
WITH SUBSTITUTES AND COMPLEMENTS**

BY FRANCISCO J. R. RUIZ^{*,†}, SUSAN ATHEY[‡] AND DAVID M. BLEI[†]

University of Cambridge,^{} Columbia University,[†] and Stanford University[‡]*

1. Details on the Inference Algorithm. Here we provide the technical details of the variational inference procedure and a description of the full algorithm.

Recall the notation introduced in Section 5 of the main paper, where $\ell = \{\rho, \alpha, \lambda, \theta, \gamma, \beta, \mu, \delta\}$ denotes the set of all latent variables in the model, $y = \{y_t\}$ is the collection of (unordered) baskets, and $x = x_{1:T}$, where $x_t = (u_t, w_t, r_t)$ are observed covariates of the shopping trips. We use mean-field variational inference, i.e., we approximate the posterior $p(\ell \mid y, \mathbf{x})$ with a fully factorized variational distribution,

$$(1) \quad q(\ell) = \prod_c q(\alpha_c)q(\rho_c)q(\lambda_c)q(\beta_c)q(\mu_c) \times \prod_u q(\theta_u)q(\gamma_u) \times \prod_w q(\delta_w).$$

We set each variational factor in the same family as the prior, i.e., Gaussian variational distributions with diagonal covariance matrices for $q(\rho_i)$, $q(\alpha_i)$, $q(\lambda_i)$, $q(\nu_u)$, $q(\mu_i)$, and $q(\delta_w)$, and independent gamma variational distributions for the price sensitivity terms $q(\gamma_u)$ and $q(\beta_i)$. We parameterize the Gaussian in terms of its mean and standard deviation, and we parameterize the gamma in terms of its shape and mean.

Let v denote the vector containing all the variational parameters. We wish to find the variational parameters v that maximize the evidence lower bound (ELBO),

$$(2) \quad \begin{aligned} \mathcal{L}(v) &= \mathbb{E}_{q(\ell; v)} [\log p(y \mid \mathbf{x}, \ell) + \log p(\ell) - \log q(\ell; v)] \\ &= \mathbb{E}_{q(\ell; v)} \left[\sum_{t=1}^T \log p(y_t \mid x_t, \ell) + \log p(\ell) - \log q(\ell; v) \right] \\ &\leq \log p(y \mid \mathbf{x}). \end{aligned}$$

In this supplement, we show how to apply stochastic optimization to maximize the bound on the log marginal likelihood. More in detail, we first describe how to tackle the intractable expectations using stochastic optimization and the reparameterization trick, and then we show how to leverage stochastic optimization to decrease the computational complexity.

1.1. Intractable expectations: Stochastic optimization and the reparameterization trick. We are interested in maximizing an objective function of the form

$$(3) \quad \widetilde{\mathcal{L}}(v) = \mathbb{E}_{q(\ell; v)} [f(\ell, v)]$$

with respect to the parameters ν . In the particular case that $f(\ell, \nu) = \log p(y, \ell | \mathbf{x}) - \log q(\ell; \nu)$, we recover the ELBO in Eq. 2, but we prefer to keep the notation general because in Section 1.2 we will consider other functions $f(\ell, \nu)$.

The main challenge is that the expectations in Eq. 3 are analytically intractable. Thus, the variational algorithm we develop aims at obtaining and following noisy estimates of the gradient $\nabla_\nu \widetilde{\mathcal{L}}$. We obtain these estimates via stochastic optimization; in particular, we apply the reparameterization trick to form Monte Carlo estimates of the gradient (Kingma and Welling, 2014; Titsias and Lázaro-Gredilla, 2014; Rezende, Mohamed and Wierstra, 2014).

In reparameterization, we first introduce a transformation of the latent variables $\ell = \mathcal{T}(\epsilon; \nu)$ and an auxiliary distribution $\pi(\epsilon; \nu)$, such that we can obtain samples from the variational distribution $q(\ell; \nu)$ following a two-step process:

$$(4) \quad \epsilon \sim \pi(\epsilon; \nu), \quad \ell = \mathcal{T}(\epsilon; \nu).$$

The requirement for $\pi(\epsilon; \nu)$ and $\mathcal{T}(\epsilon; \nu)$ is that this procedure must provide a variable ℓ that is distributed according to $\ell \sim q(\ell; \nu)$. Here we have considered the generalized reparameterization approach (Ruiz, Titsias and Blei, 2016; Naesseth et al., 2017), which allows the auxiliary distribution $\pi(\epsilon; \nu)$ to depend on the variational parameters ν (this is necessary because the gamma random variables are not otherwise reparameterizable).

Once we have introduced the auxiliary variable ϵ , we can rewrite the gradient of the objective in Eq. 3 as an expectation with respect to the auxiliary distribution $\pi(\epsilon; \nu)$,

$$(5) \quad \nabla_\nu \widetilde{\mathcal{L}} = \nabla_\nu \mathbb{E}_{q(\ell; \nu)} [f(\ell, \nu)] = \nabla_\nu \mathbb{E}_{\pi(\epsilon; \nu)} [f(\mathcal{T}(\epsilon; \nu), \nu)].$$

We now push the gradient into the integral¹ and apply the chain rule for derivatives to express the gradient as an expectation,

$$(6) \quad \nabla_\nu \widetilde{\mathcal{L}} = \mathbb{E}_{\pi(\epsilon; \nu)} \left[\nabla_\ell f(\ell, \nu) \Big|_{\ell=\mathcal{T}(\epsilon; \nu)} \nabla_\nu \mathcal{T}(\epsilon; \nu) + f(\mathcal{T}(\epsilon; \nu), \nu) \nabla_\nu \log \pi(\epsilon; \nu) \right].$$

To obtain this expression, we have assumed that $\mathbb{E}_{\pi(\epsilon; \nu)} \left[\nabla_\nu f(\ell, \nu) \Big|_{\ell=\mathcal{T}(\epsilon; \nu)} \right] = 0$ because the only dependence of $f(\ell, \nu)$ on ν is through the term $\log q(\ell; \nu)$, and the expectation of the score function is zero.

We can now obtain a Monte Carlo estimate of the expectation in Eq. 6 (and therefore of the gradient of interest) by drawing a sample from $\pi(\epsilon; \nu)$ and evaluating the argument of the expectation. That is, we form the gradient estimator as

$$(7) \quad \widehat{\nabla}_\nu \widetilde{\mathcal{L}} = \nabla_\ell f(\ell, \nu) \Big|_{\ell=\mathcal{T}(\epsilon; \nu)} \nabla_\nu \mathcal{T}(\epsilon; \nu) + f(\mathcal{T}(\epsilon; \nu), \nu) \nabla_\nu \log \pi(\epsilon; \nu),$$

where $\epsilon \sim \pi(\epsilon; \nu)$. This assumes that we are able to evaluate $f(\ell, \nu)$ and its gradient. (We show in Section 1.2 how to do that efficiently.)

¹In the model that includes thinking-ahead, this step introduces a small bias due to the non-differentiability of the $\max(\cdot)$ operator; see Lee, Yu and Yang (2018).

Algorithm 1: Variational inference algorithm

input : Data y and x , model hyperparameters
output Variational parameters v
:
Initialize v randomly
Initialize iteration number $m \leftarrow 1$
repeat
 Sample $\epsilon \sim \pi(\epsilon; v)$
 Compute $\ell = \mathcal{T}(\epsilon; v)$
 Evaluate $f(\ell, v)$ and $\nabla_{\ell} f(\ell, v)$ (see Algorithm 2)
 Obtain an estimate of the gradient, $\widehat{\nabla}_v \widetilde{\mathcal{L}}$ (Eq. 7)
 Set the step size $\eta^{(m)}$ (e.g., use the schedule proposed in ADVI)
 Take a gradient step, $v \leftarrow v + \eta^{(m)} \odot \widehat{\nabla}_v \widetilde{\mathcal{L}}$
 Increase the iteration number, $m \leftarrow m + 1$
until convergence
return v

We use a transformation $\mathcal{T}(\cdot)$ and auxiliary distribution $\pi(\cdot)$ for each variational factor. For a Gaussian variational factor with mean μ and standard deviation σ , we use the standard reparameterization,

$$(8) \quad \mathcal{T}_{\text{Gauss}}(\epsilon; \mu, \sigma) = \mu + \sigma \epsilon, \quad \pi_{\text{Gauss}}(\epsilon; \mu, \sigma) = \mathcal{N}(0, 1),$$

which makes the last term of Eq. 7 vanish because $\nabla_v \log \pi_{\text{Gauss}}(\epsilon; \mu, \sigma) = 0$. For a gamma variational factor with shape α and mean μ , we use the transformation based on rejection sampling (Marsaglia and Tang, 2000),

$$(9) \quad \mathcal{T}_{\text{Gamma}}(\epsilon; \alpha, \mu) = \frac{\mu}{\alpha} \left(\alpha - \frac{1}{3} \right) \left(1 + \frac{\epsilon}{\sqrt{9\alpha - 3}} \right)^3,$$

and $\pi_{\text{Gamma}}(\epsilon; \alpha, \mu)$ is defined through a rejection sampling procedure. See Naesseth et al. (2017) for further details about the reparameterization trick for gamma random variables.²

Algorithm 1 summarizes the resulting variational inference procedure. At each iteration, we obtain a sample from ℓ via the auxiliary distribution $\pi(\epsilon; v)$ and the transformation $\mathcal{T}(\epsilon; v)$; we evaluate the function $f(\ell, v)$ and its gradient with respect to the latent variables ℓ ; we obtain the gradient estimate in Eq. 7; and we take a gradient step for the variational parameters v . In the stochastic optimization procedure, we adaptively set the step size as proposed in the automatic differentiation variational inference (ADVI) algorithm (Kucukelbir et al., 2017).

1.2. Computational complexity: Stochastic optimization and lower bounds. The algorithm in Section 1.1 requires to evaluate the model log joint (as well as its gradient). There are three issues that make it expensive to evaluate the log joint. First, evaluating the log likelihood is expensive because it involves a summation

²In particular, we also apply the “shape augmentation trick,” which allows us to reparameterize a gamma random variable with shape α in terms of another gamma random variable with shape $\alpha + P$, where P is a positive integer. We use $P = 10$. See Naesseth et al. (2017) for additional details.

over shopping trips. This represents a problem when the dataset is large. Second, evaluating the softmax involves computing its normalization constant, which contains a summation over all items. This becomes an issue when there are thousands of items and we need to evaluate many softmax probabilities. Third, computing the probability over unordered baskets y_t is also expensive, as it involves a summation over all possible permutations.

We address these issues by combining two techniques: data subsampling and lower bounds on the ELBO. We first describe data subsampling for evaluating the log likelihood. Note that the log likelihood involves a summation over many terms,

$$(10) \quad \log p(y | \mathbf{x}, \ell) = \sum_{t=1}^T \log p(y_t | x_t, \ell).$$

We can obtain an unbiased estimate of the log likelihood (and its gradient) by sampling a random subset of shopping trips. Let \mathcal{B}_T be the (randomly chosen) set of trips. The estimator

$$\frac{T}{|\mathcal{B}_T|} \sum_{t \in \mathcal{B}_T} \log p(y_t | x_t, \ell)$$

is unbiased because its expected value is the log likelihood $\log p(y | \mathbf{x}, \ell)$ (Hoffman et al., 2013). Thus, we subsample data terms to obtain unbiased estimates of the log likelihood and its gradient, resulting in a computationally more efficient algorithm.

Second, we describe how to form variational bounds to address the issue of the expensive normalization constant of the softmax. Each softmax log probability is given by

$$(11) \quad \log p(y_{ti} = c | \mathbf{y}_{t,i-1}) = \Psi(c, \mathbf{y}_{t,i-1}) - \log \left(\sum_{c' \notin \mathbf{y}_{t,i-1}} \exp\{\Psi(c', \mathbf{y}_{t,i-1})\} \right).$$

The summation over c' is expensive, and we cannot easily form an unbiased estimator because of the non-linearity introduced by the logarithm. Hence, we apply the one-vs-each bound (Titsias, 2016), which allows us to write

$$(12) \quad \log p(y_{ti} = c | \mathbf{y}_{t,i-1}) \geq \sum_{c' \notin [\mathbf{y}_{t,i-1}, c]} \log \sigma(\Psi(c, \mathbf{y}_{t,i-1}) - \Psi(c', \mathbf{y}_{t,i-1})),$$

where $\sigma(x) = \frac{1}{1+e^{-x}}$ is the sigmoid function. We can form unbiased estimates of the summation via subsampling. More precisely, we randomly sample a set $\mathcal{B}_C^{(t,i)}$ of items (each of them distinct from c and from the other items in the basket). Then, we form the following unbiased estimator:

$$\frac{C-i}{|\mathcal{B}_C^{(t,i)}|} \sum_{c' \in \mathcal{B}_C^{(t,i)}} \log \sigma(\Psi(c, \mathbf{y}_{t,i-1}) - \Psi(c', \mathbf{y}_{t,i-1})).$$

Here, C stands for the total number of items.

Finally, we show how to deal with the issue of unordered baskets. Recall that each log likelihood term involves a summation over all possible permutations of the baskets (holding the checkout item in the last position),

$$(13) \quad \log p(y_t | x_t, \ell) = \log \left(\sum_{\pi_t} p(\mathbf{y}_{t,\pi_t} | x_t, \ell) \right).$$

Following a similar procedure as [Doshi-Velez et al. \(2009\)](#), we introduce an auxiliary distribution $q(\pi_t)$ to rewrite the expression above as an expectation with respect to $q(\pi_t)$, and then we apply Jensen's inequality:

$$(14) \quad \begin{aligned} \log p(y_t | x_t, \ell) &= \log \left(\mathbb{E}_{q(\pi_t)} \left[\frac{p(\mathbf{y}_{t,\pi_t} | x_t, \ell)}{q(\pi_t)} \right] \right) \\ &\geq \mathbb{E}_{q(\pi_t)} [\log p(\mathbf{y}_{t,\pi_t} | x_t, \ell) - \log q(\pi_t)] \\ &= \sum_{\pi_t} q(\pi_t) (\log p(\mathbf{y}_{t,\pi_t} | x_t, \ell) - \log q(\pi_t)). \end{aligned}$$

Since the bound involves a direct summation over permutations, we can subsample terms to alleviate the computational complexity. For simplicity, we set $q(\pi_t)$ to be a uniform distribution over all possible permutations, and thus we do not introduce auxiliary variational parameters that would be too expensive to obtain otherwise. In particular, we form an unbiased estimate of the bound by sampling one random permutation π_t and evaluating the term³

$$\log p(\mathbf{y}_{t,\pi_t} | x_t, \ell).$$

To sum up, we have derived a bound of the ELBO,

$$(15) \quad \tilde{\mathcal{L}}(v) \leq \mathcal{L}(v) \leq \log p(y | x),$$

which can still be written as an expectation with respect to the variational distribution $q(\ell | v)$. More importantly, we can efficiently evaluate the argument $f(\ell, v)$ of such expectation and its gradient with respect to the latent variables ℓ .

Putting all together, the function $f(\ell, v)$ that we use is given by

$$(16) \quad \begin{aligned} f(\ell, v) &= \sum_{t=1}^T \sum_{\pi_t} q(\pi_t) \sum_{i=1}^{n_t} \sum_{c' \notin [\mathbf{y}_{t,i-1}, c]} \log \sigma(\Psi(c, \mathbf{y}_{t,i-1}) - \Psi(c', \mathbf{y}_{t,i-1})) \\ &\quad + \log p(\ell) - \log q(\ell; v). \end{aligned}$$

We obtain an unbiased estimate via subsampling shopping trips, one permutation π_t for each one, and items c' . We use $|\mathcal{B}_T| = 100$ trips and $|\mathcal{B}_C^{(t,i)}| = 50$ items in our experiments.

Algorithm 2 outlines the procedure to obtain an unbiased estimate of $f(\ell, v)$ for a given sample of the latent parameters, $\ell \sim q(\ell; v)$. Differentiation through Algorithm 2 gives the gradient $\nabla_{\ell} f(\ell, v)$, which is also required in the inference procedure (Algorithm 1).

³We ignore the term $\log q(\pi_t)$ because it is a constant.

Algorithm 2: Estimate of $f(\ell, v)$ (Eq. 16)

input : Data y and x , a sample of the latent parameters ℓ , variational distribution $q(\ell; v)$
output An unbiased estimate of $f(\ell, v)$ and its gradient (for the gradient, differentiate through
 :
 the algorithm)
 Initialize $\hat{f} \leftarrow \log p(\ell) - \log q(\ell; v)$
 Sample a set of baskets $\mathcal{B}_T \subseteq \{1, \dots, T\}$
for $t \in \mathcal{B}_T$ **do**
 Sample a permutation π_t of the items in basket t
 Set \mathbf{y}_t to the vector containing the items in basket t ordered according to π_t
 for $i = 1, \dots, n_t$ **do**
 Set c to the i th item in \mathbf{y}_t
 Sample a set of items $\mathcal{B}_C^{(t,i)} \subseteq \{1, \dots, C\} \setminus \{\mathbf{y}_{t,i-1}, c\}$
 for $c' \in \mathcal{B}_C^{(t,i)}$ **do**
 Update $\hat{f} \leftarrow \hat{f} + \frac{T}{|\mathcal{B}_T|} \times \frac{C-i}{|\mathcal{B}_C^{(t,i)}|} \times \log \sigma(\Psi(c, \mathbf{y}_{t,i-1}) - \Psi(c', \mathbf{y}_{t,i-1}))$
 end
 end
end
return \hat{f} and $\nabla_{\ell} \hat{f}$

References.

- DOSHI-VELEZ, F., MILLER, K. T., VAN GAEL, J. and TEH, Y. W. (2009). Variational inference for the Indian buffet process. In *Proceedings of the International Conference on Artificial Intelligence and Statistics* **12**.
- HOFFMAN, M. D., BLEI, D. M., WANG, C. and PAISLEY, J. (2013). Stochastic variational inference. *Journal of Machine Learning Research* **14** 1303–1347.
- KINGMA, D. P. and WELING, M. (2014). Auto-encoding variational Bayes. In *International Conference on Learning Representations*.
- KUCUKELBIR, A., TRAN, D., RANGANATH, R., GELMAN, A. and BLEI, D. M. (2017). Automatic differentiation variational inference. *Journal of Machine Learning Research* **18** 1–45.
- LEE, W., YU, H. and YANG, H. (2018). Reparameterization gradient for non-differentiable models. In *Advances in Neural Information Processing Systems*.
- MARSAGLIA, G. and TANG, W. W. (2000). A simple method for generating gamma variables. *ACM Transactions on Mathematical Software* **26** 363–372.
- NAESSETH, C., RUIZ, F. J. R., LINDERMAN, S. and BLEI, D. M. (2017). Reparameterization gradients through acceptance-rejection methods. In *Artificial Intelligence and Statistics*.
- REZENDE, D. J., MOHAMED, S. and WIERSTRA, D. (2014). Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*.
- RUIZ, F. J. R., TITSIAS, M. K. and BLEI, D. M. (2016). The generalized reparameterization gradient. In *Advances in Neural Information Processing Systems*.
- TITSIAS, M. K. (2016). One-vs-each approximation to softmax for scalable estimation of probabilities. In *Advances in Neural Information Processing Systems*.
- TITSIAS, M. K. and LÁZARO-GREDILLA, M. (2014). Doubly stochastic variational Bayes for non-conjugate inference. In *International Conference on Machine Learning*.