

# Bayesian non-parametric inference with monotonicity restrictions for discrete choice experiments

Remi Daviet<sup>\*</sup> and William McCausland<sup>†</sup>

May 12, 2017

## Abstract

Insert abstract here

**Keywords:** Big Data, Discrete Choice experiments

---

<sup>\*</sup>Corresponding author. Department of economics, University of Toronto; remi.daviet@mail.utoronto.ca

<sup>†</sup>Department of economics, Université de Montréal

# 1 Products attributes and products classes

Every product considered is represented by a point in the chosen attributes space. For instance, a mobile phone can be represented by a set of two attributes, the screen size and the color. Each attribute has a level. In our example, the screen size can have 3 levels (*Small*, *Medium* or *Large*) and the color can have 4 levels (*Black*, *White*, *Blue* or *Other*). Considering this attributes space, two possible phones are phone 1 (*Large*, *Blue*) and phone 2 (*Medium*, *White*).

Products can be regrouped into classes. Classes are set of attributes commons to every product belonging to the class. A class can be very simple, for instance the class ( $\{Small\}, any$ ) which regroups all the phones with a small screen size, no matter the color. Classes can also be more complicated such as ( $\{Small, Large\}, \{White, Other\}$ ) which regroups phones of small or large screen size having a white or "other" color. Note that the ( $\{Small\}, any$ ) class can be written ( $\{Small\}, \{Black, White, Blue, Other\}$ ) and that the keyword "any" is just used for convenience. Finally, a class can specify a unique product in the attribute space, such as the class ( $\{Large\}, \{Blue\}$ ).

More formally, considering a product space with  $L$  attributes, a product  $x_m$  is represented by the vector of attributes  $(x_{m1}, x_{m2}, \dots, x_{mL})$ . Each attribute  $i$  can have  $L_i$  possible levels. For instance, if the attribute (1) "screen size" has 3 possible levels (*Small*, *Medium* or *Large*), we write  $L_1 = 3$ . For the second attribute "Color" which has 4 possible levels, we write  $L_2 = 4$ . The attributes levels attributed to a product  $x_m$  are indicated by their corresponding number. For instance, if *Small* and *Black* are the first levels of their respective attributes, the phone will be represented by the vector  $(x_{m1} = 1, x_{m2} = 1)$ .

A class of product  $\Lambda_k$  is represented by a Cartesian product of  $L$  sets  $(\Lambda_{k1}, \Lambda_{k2}, \dots, \Lambda_{kL})$ . Each  $\Lambda_{ki}$  represents the set of possible levels for the attribute  $i$  for a product in the class  $k$ . The levels belonging to  $\Lambda_{ki}$  are called active levels. For instance, a class regrouping *Black* phones with either *Medium* or *Large* screen size will be represented by  $(\Lambda_{k1} = \{2, 3\}, \Lambda_{k2} = \{1\})$ . When all the levels of an attribute are active, it means that this particular attribute does not matter when defining the class, and the attribute is defined as inactive. On the contrary, when only some levels are active, it means that the attribute is used as a criterion to select if a product is a member of the class, and the attribute is defined as active.

A product  $x_m$  belongs to a class  $\Lambda_k$  if and only if  $x_{m1} \in \Lambda_{k1}, x_{m2} \in \Lambda_{k2}, \dots$  and  $x_{mL} \in \Lambda_{kL}$ . For convenience, we will use the notation  $x_m \in \Lambda_k$  to indicate that the product  $x_m$  is a member of the class  $\Lambda_k$ .

## 2 Individual preferences

We have a population of  $N$  individuals sequentially choosing  $T$  times one product within a set. The choice of an individual  $n$  at a time  $t$  is done by maximizing a function  $u_t(x, \psi_n)$  representing the preferences this individual over the set of available choices  $X_{nt}$ , where  $\psi_n$  is some set of parameters. The choice maximizing the individual preferences is denoted  $x^*$ . The preferences for each individual are built by attributing some value to various classes of products. The value of a class can be positive or negative. For instance, at a time  $t$  an individual  $n$  attributes the value  $\theta_{nkt}$  to the class  $\Lambda_k$ . The value of a given product for an individual is obtained by summing the values of all the classes it belongs to. If the value of a class  $k$  is zero for an individual, we say that the individual does not possess the feature  $k$ . It means

that whether a product belongs to the class  $k$  or not does not change the value of the product for the given individual.

Within the population, when two individuals have a non-zero value  $\theta_{nkt}$  for a class of products  $k$ , we say that they share the common feature  $k$ . It is important to note however that the value of the class does not need to be the same for both individuals sharing the feature. One individual with feature  $k$  may attribute a positive value to the class  $k$  while the other attributes a negative value. Two individuals sharing a feature  $k$  means that they are both sensitive to the corresponding class of products. Some features can be shared by many individuals, while some other features can be rare. Note that there is some redundancy: a same class can appear in two different features. The two classes of products  $\Lambda_k$  and  $\Lambda_{k'}$  can be identical ( $\Lambda_k = \Lambda_{k'}$ ), but to be respectively associated with a  $\theta_{nk}$  and a  $\theta_{nk'}$  following different distributions.

The total number of features in the population is potentially infinite. However, in a given sample only a finite number of features are realized. The features present in our sample are ordered and numbered from 1 to  $K_+$  for convenience.

Formally, the value of an object  $x$  for an individual  $n$  at time  $t$  can be written:

$$u_t(x, \psi_n) = \sum_{k=1}^{K_+} \theta_{nkt} \cdot \mathbb{1}(x \in \Lambda_k), \quad (2.1)$$

where  $\psi_n$  contains the set of  $\theta_{nkt}$  and the set of  $\Lambda_k$ :

$$\psi_n = \{ \{ \theta_{nkt} \}_{t=1}^T, \Lambda_k \}_{k=1}^{K_+}$$

For convenience, the value  $\theta_{nkt}$  may be decomposed into four parts as follow:

$$\theta_{nkt} = z_{nk} \cdot (-1)^{s_{nk}} \cdot (\theta_{nk} + \epsilon_{nkt})$$

where  $z_{nk}$  is a binary variable indicating whether individual  $n$  possesses feature  $k$ ,  $s_{nk}$  is a binary variable that indicates whether the value of the class  $k$  for individual  $n$  is positive or negative,  $\theta_{nk}$  is the stable part of the value in the consumer preferences, and  $\epsilon_{nkt}$  is a tremble part. We can rewrite  $\psi_n$  as:

$$\psi_n = \{ z_{nk}, s_{nk}, \theta_{nk}, \{ \epsilon_{nkt} \}_{t=1}^T, \Lambda_k \}_{k=1}^{K_+}.$$

It is possible that we want to restrict preferences such that they only increase in some attributes. For such attribute, the active levels are restricted to be of the form  $\{l, l+1, \dots, L\}$  if the the class  $\Lambda_k$  is associated with a positive value ( $s_{nk} = 0$ ), and of the form  $\{1, \dots, l-1, l\}$  if the class is associated with a negative value ( $s_{nk} = 1$ ). We denote the first case as a upper active levels, and the second as lower active levels. If a class  $\Lambda_k$  has at least one active monotone attribute, we call it a monotone class. Note that if several monotone attributes are active in the class, they should all be of the same type. If all the active monotone attribute are of the upper active levels type, we say that the class has a positive polarity and denote it  $\rho_k = 0$ . conversely, if the monotone attributes are of the lower active levels type, we say that the class has a negative polarity and denote it  $\rho_k = 1$ .

????????????????QUESTION????????????????

Shall I write  $u_t(x, \psi_n)$  with  $\psi_n = \{ z_{nk}, s_{nk}, \theta_{nk}, \{ \epsilon_{nkt} \}_{t=1}^T, \Lambda_k \}_{k=1}^{K_+}$ ,

Or  $u(x, \psi_{nt})$ , with  $\psi_{nt} = \{z_{nk}, s_{nk}, \theta_{nk}, \epsilon_{nkt}, \Lambda_k\}_{k=1}^{K_+}$  ?

in the second case, the various  $\psi_{nt}$  have common parameters indexed by  $nk$ , which may look weird. I used the first approach where the index  $t$  is on  $u_t$  because the  $t$  is relevant both for the component  $\epsilon_{nkt}$  of  $\psi_n$  to retain, and for the choice set / observation considered.

### 3 Hierarchical model

The Bayesian framework gives us a simple and intuitive way to model our uncertainty about the various parameters. We have to define how the various variables and parameters are distributed in order to be able to obtain a posterior distribution.

Model parameters:

$$\{\{\psi_n\}_{n=1}^N, \{\Lambda_k\}_{k=1}^{K_+}\} = \{\{\{\epsilon_{nkt}\}_{t=1}^T, z_{nk}, s_{nk}, \theta_{nk}\}_{n=1}^N, \Lambda_k\}_{k=1}^{K_+}$$

We add a hierarchical level requiring the addition of the following hyper-parameters:

$$\{\{\pi_k^z, \pi_k^s\}_{k=1}^{K_+}, \alpha^\theta, \sigma^\theta\}$$

The joint distribution of the data and parameters satisfies the conditional independence relationships implied by the following decomposition:

$$P\left(\{\{x_{nt}^*\}_{t=1}^T\}_{n=1}^N, \{\{\{\epsilon_{nkt}\}_{t=1}^T, z_{nk}, s_{nk}, \theta_{nk}\}_{n=1}^N, \Lambda_k, \pi_k^z, \pi_k^s\}_{k=1}^{K_+}, \alpha^\theta, \sigma^\theta\right) \quad (3.2)$$

$$= \prod_t \prod_n P\left(x_{nt}^* | \{\epsilon_{nkt}, z_{nk}, s_{nk}, \theta_{nk}, \Lambda_k\}_{k=1}^{K_+}\right) \quad (3.3)$$

$$\times \prod_k \left[ \prod_n \left( \prod_t P(\epsilon_{nkt} | \alpha^\theta, \sigma^\theta, z_{nk}) \right) P(z_{nk} | \pi_k^z) P(s_{nk} | \pi_k^s, z_{nk}) P(\theta_{nk} | \alpha^\theta, \sigma^\theta, z_{nk}) P(\Lambda_k) \right] \quad (3.4)$$

$$\times \left( \prod_k P(\pi_k^z) P(\pi_k^s) \right) P(\alpha^\theta) P(\sigma^\theta), \quad (3.5)$$

where line (3.3) is the likelihood, line (3.4) is the prior and line (3.5) is the hyper prior.

#### 3.1 Features allocation

Our first object of interest is the binary matrix of features allocation  $Z = \{z_{nk}\}$  where each row represents an individual and each column represents a feature. While the  $Z$  matrix has an infinite number of columns, we are only interested in the columns where there is at least a 1 for our sample. We use the infinite feature model defined by Griffiths and Ghahramani (2011) in their paper about the Indian Buffet Process (IBP). The IBP approach presents a convenient way to perform inference by simulation for models using latent features. Using this model, we assume that the probability that any individual possesses

feature  $k$  is  $\pi_k^z$  and that the features are generated independently. The distribution of the infinite matrix  $Z$  can be obtained by starting with a finite  $N \times K$  matrix of features allocation  $Z_K$  and taking the limit when  $K \rightarrow \infty$ .

In the finite case, the probability of  $Z_K$  given  $\pi^z = \{\pi_1^z, \dots, \pi_K^z\}$  is:

$$P(Z_K|\pi) = \prod_{k=1}^K \prod_{n=1}^N P(z_{nk}|\pi_k^z) = \prod_{k=1}^K (\pi_k^z)^{m_k} (1 - \pi_k^z)^{N-m_k}$$

where  $m_k = \sum_{n=1}^N z_{nk}$  is the number of individuals having feature  $k$ . We can then define a prior on  $\pi^z$  using a *Beta* distribution:

$$\pi_k^z \sim \text{Beta}\left(\frac{\alpha^z \beta^z}{K}, \beta^z\right)$$

In order to take the limit as  $K \rightarrow \infty$ , we need to define a distribution over classes of equivalent matrices. Matrices are defined as equivalent if one can be obtained from the other by permutation of the columns. A class of equivalent matrices is denoted  $[Z]$ . In order to be able to compare matrices, we use permutations on their column in order to put them in left ordered form. Griffiths and Ghahramani (2011) define the left-ordered form as follow: "*lof*( $Z$ ) is obtained by ordering the columns of the binary matrix  $Z$  from left to right by the magnitude of the binary number expressed by that column, taking the first row as the most significant bit". Two equivalent matrices have the same left-ordered form. Taking left-ordered forms let us define a distribution over equivalent classes of matrices. Matrices with left-ordered forms following that distribution can be generated using the IBP.

Using this process, the average number of feature per individual will be  $\alpha^z$ . The  $\beta^z$  parameter lets us set the average overall number of features observed in a sample. The expected overall number of features present in a sample will be  $\alpha^z \sum_{n=1}^N \frac{\beta^z}{\beta^z + n - 1}$ . Depending on  $\beta$ , this number can range from  $\alpha^z$  where everybody share the same features, to  $N\alpha^z$  where no features are shared.

Here are three matrices produced using the described process using various parameter values:

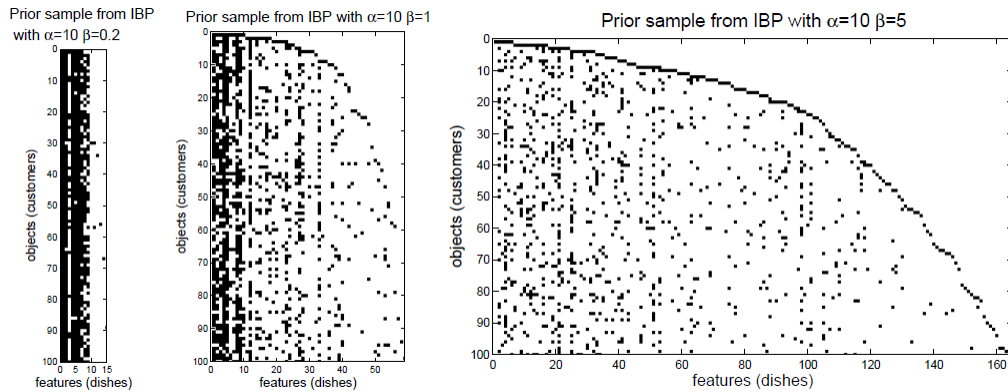


Figure 1: Three matrices produced with the IBP.(Griffiths and Ghahramani, 2011)

### 3.2 Attributes with monotone preferences and $\Lambda_k$

The class  $\Lambda_k$  can be denoted as a set of active attributes  $C_k$  with cardinality  $c_k$ , and for each active attribute  $i \in C_k$ , a set of active levels  $\Lambda_{ki}$  with cardinality  $\lambda_{ki}$ .

The prior distribution of  $\Lambda_k$  can be decomposed as follow:

$$P(\Lambda_k) = P(c_k) \cdot P(\rho_k) \cdot P(C_k|c_k) \cdot \prod_{i \in C_k} P(\Lambda_{ki}|\rho_k) \quad (3.6)$$

The corresponding distributions are detailed below:

- The number of active attributes is uniformly distributed on  $\{1, \dots, L\}$ :

$$c_k \sim U(1, \dots, L)$$

- The polarity has equal probability of being positive or negative:

$$\rho_k \sim \text{Bernoulli}(0.5)$$

- The attributes share the same probability of being active. The probability of observing a set of active attributes  $C_k$  is consequently:

$$P(C_k|c_k) = \binom{L}{c_k}^{-1}$$

- The probability of observing the set of active levels  $\Lambda_{ki}$  depends on several factors:

- If  $\Lambda_{ki}$  is monotone:

$$P(\Lambda_{ki}|\rho_k) = \begin{cases} 1/L_i, & \text{if } \Lambda_{ki} \text{ has upper active levels and } \rho_k = 0 \\ 1/L_i, & \text{if } \Lambda_{ki} \text{ has lower active levels and } \rho_k = 1 \\ 0, & \text{otherwise} \end{cases}$$

- If  $\Lambda_{ki}$  is not monotone, it does not depend on the polarity  $\rho_k$ . The number of active levels  $\lambda_{ki}$  is uniformly distributed on  $\{1, \dots, L_i\}$  and each level has an equal probability of being active.

$$\begin{aligned} P(\Lambda_{ki}|\rho_k) &= P(\Lambda_{ki}|\lambda_{ki}) \cdot P(\lambda_{ki}) \\ &= \binom{L_i}{\lambda_{ki}}^{-1} \cdot \frac{1}{\lambda_{ki}} \end{aligned}$$

### 3.3 Other parameters

The other parameters  $(s_{nk}, \theta_{nk}, \epsilon_{nkt})$  are defined conditional on  $z_{nk}$ . When  $z_{nk} = 0$ , the values of all these parameters are arbitrarily set to 0. We have to define the prior distributions of  $(s_{nk}, \theta_{nk}, \epsilon_{nkt})$  conditional on  $z_{nk} = 1$ :

$$s_{nk} | \pi_k^s, z_{ik} = 1 \sim \begin{cases} 0 & \text{if } m_k = 1 \text{ and } \rho_k = 1 \\ 1 & \text{if } m_k = 1 \text{ and } \rho_k = 0 \\ \text{Bernoulli}(\pi_k^s) & \text{otherwise} \end{cases}$$

$$\theta_{nk} | \sigma^\theta, \alpha^\theta, z_{ik} = 1 \sim \text{Gamma}(\sigma^\theta \alpha^\theta, 1)$$

$$\epsilon_{nkt} | \sigma^\theta, \alpha^\theta, z_{ik} = 1 \sim \text{Gamma}((1 - \sigma^\theta) \alpha^\theta, 1)$$

where  $(s_{nk}, \theta_{nk})$  are independent across  $n$  and  $k$ ; and  $\epsilon_{nkt}$  are independent across  $n, k$  and  $t$ .

Note that the distribution of the value size  $(\theta_{nk} + \epsilon_{nkt})$  is also Gamma with parameters  $\text{Gamma}(\alpha^\theta, 1)$  from the properties of a sum of Gamma distributed variables. The parameter  $\alpha^\theta$  is the shape parameter of the distribution of the value size, while  $\sigma^\theta$  can be interpreted as a stability parameter. For  $\sigma^\theta$  close to one, the stable component  $\theta_{nk}$  is favored and the utility functions have high serial dependence; for  $\sigma^\theta$  close to zero, the variable component  $\epsilon_{nkt}$  is favored and the utility function has low serial dependence.

Finally, we have to define hyperpriors on some of the hyperparameters:

$$\alpha^\theta \sim \text{Gamma}(1, 1)$$

$$\sigma^\theta \sim \text{Gamma}(1, 1)$$

$$\pi_k^s \sim \text{Beta}(1, 1)$$

## 4 Inference

### 4.1 Likelihood and Posetrior

Conditional on the set of parameters  $\psi$ , there is no random term in the utility function. The likelihood becomes a degenerate function that is equal to 1 if all the observed choices are predicted correctly and 0 otherwise. The probability of observing the set of choices  $\{x_{nt}^*\}$  is:

$$P(\{x_{nt}^*\} | \psi) = \mathbb{1}\{\arg \max_{x \in X_{nt}} u(x, \psi) = x_{nt}^*, \forall n, t\}$$

Likelihood-based inference consists in recovering the identified set of parameters such that:

$$P(\{x_{nt}^*\} | \psi) = 1$$

We can see that with the current likelihood function, if at least one choice is not maximizing utility according to the current value of  $\psi$ , the likelihood takes the value of 0.

In the Bayesian framework, the target distribution is the corresponding posterior:

$$P(\psi | \{x_{nt}^*\}) \propto P(\{x_{nt}^*\} | \psi) \cdot P(\psi),$$

where  $P(\psi)$  denotes the prior density of our parameters.

## 4.2 Tempering and Sequence of Distributions

Simulating the posterior can be quite difficult, and we simulate instead a sequence of target distributions converging to this Posterior. The functions used in the sequence are increasingly difficult to simulate. This process is known as tempering, and we achieve it by combining two methods. First, we use Data Tempering and introduce the observations one by one in the likelihood. We start with the first observation of the first individual  $x_{1,1}^*$ , then add the subsequent observations for this individual until all the observations are included in the likelihood. We then add the observations of the following individual, and continue until every observation of every individual are included. For the second tempering approach, we use an instrumental function that does not take a value of zero if the last observed choice is not utility maximizing with the current value of  $\psi$ . We call this approach  $\zeta$ -tempering and substitute the likelihood with the following function:

$$Q(\{x_{nt}^*\}|\psi; \eta, \tau, \zeta) = P(\{\{x_{nt}^*\}_{t=1}^T\}_{n=1}^{\eta-1}|\psi) \cdot P(\{x_{\eta t}^*\}_{t=1}^{\tau-1}|\psi) \cdot \zeta^{u(x_{\eta\tau}^*, \psi) - \max_{x \in X_{\eta\tau}} u(x, \psi)},$$

where  $\zeta > 1$  is a constant chosen by the researcher. We can see that this instrumental function is the product of 3 terms. The first one is the likelihood of observing the choices of individuals 1 to  $\eta - 1$ . The second is the likelihood of observing the  $\tau - 1$  first choices of individual  $\eta$ . The last term requires further attention. If the last observed choice is utility maximizing, this term is equal to 1 as for the likelihood. However, it is equal to a value between 0 and 1 if another choice would be utility maximizing. The smaller the difference between the observed choice and the utility maximizing one, the closer the term gets to 1. It is important to note that  $Q(\{x_{nt}^*\}|\psi; \eta, \tau, \zeta) \rightarrow P(\{x_{\eta t}^*\}_{t=1}^{\tau} \cup \{\{x_{nt}^*\}_{t=1}^T\}_{n=1}^{\eta-1}|\psi)$  as  $\zeta \rightarrow \infty$ .

The corresponding quasi-posterior distribution is defined as

$$Q(\psi|\{x_{nt}^*\}; \eta, \tau, \zeta) \propto Q(\{x_{nt}^*\}|\psi; \eta, \tau, \zeta) \cdot P(\psi).$$

Using these tempering approaches, we can define a sequence of distributions by specifying the values  $(\eta, \tau, \zeta)$  for each distribution in the sequence. We start by the distribution

$$Q(\psi|\{x_{nt}^*\}; \eta = 1, \tau = 1, \zeta = 2) \propto 2^{u(x_{1,1}^*, \psi) - \max_{x \in X_{1,1}} u(x, \psi)} \cdot P(\psi)$$

and progressively increase  $\zeta$  in each subsequent distribution until  $\zeta \rightarrow \infty$  and the target distribution becomes

$$Q(\psi|\{x_{nt}^*\}; \eta = 1, \tau = 1, \zeta \rightarrow \infty) \propto P(\psi|\{x_{1,1}^*\}).$$

We then alternatively increase the value of  $\tau$  by 1 and progressively increase  $\zeta$  from 2 to  $\infty$  until we reach the last observation  $T$  for the current individual. For our application, we choose to increase  $\zeta$  by a factor of 2 with each step in the sequence. We then increase  $\eta$  by one and start again alternatively increasing the value of  $\tau$  by 1 and progressively increasing  $\zeta$  from 2 to  $\infty$  until we reach the last observation of the last individual  $N$ .

The ratio of two consecutive target distributions when increasing  $\zeta$  by a factor of  $a$  is equal to

$$\frac{P(\{\{x_{nt}^*\}_{t=1}^T\}_{n=1}^{\eta-1}|\psi) \cdot P(\{x_{\eta t}^*\}_{t=1}^{\tau-1}|\psi) \cdot (a\zeta)^{u(x_{\eta\tau}^*, \psi) - \max_{x \in X_{\eta\tau}} u(x, \psi)}}{P(\{\{x_{nt}^*\}_{t=1}^T\}_{n=1}^{\eta-1}|\psi) \cdot P(\{x_{\eta t}^*\}_{t=1}^{\tau-1}|\psi) \cdot \zeta^{u(x_{\eta\tau}^*, \psi) - \max_{x \in X_{\eta\tau}} u(x, \psi)}} = a^{u(x_{\eta\tau}^*, \psi) - \max_{x \in X_{\eta\tau}} u(x, \psi)}.$$



The ratio of two consecutive target distributions when increasing  $\tau$  by 1 and resetting  $\zeta$  to 2 is equal to:

$$\frac{P(\{\{x_{nt}^*\}_{t=1}^T\}_{n=1}^{\eta-1}|\psi) \cdot P(\{x_{\eta t}^*\}_{t=1}^{\tau-1}|\psi) \cdot 2^{u(x_{\eta\tau}^*, \psi) - \max_{x \in X_{\eta\tau}} u(x, \psi)}}{P(\{\{x_{nt}^*\}_{t=1}^T\}_{n=1}^{\eta-1}|\psi) \cdot P(\{x_{\eta t}^*\}_{t=1}^{\tau-1}|\psi)} = 2^{u(x_{\eta\tau}^*, \psi) - \max_{x \in X_{\eta\tau}} u(x, \psi)}$$

### 4.3 Quasi-Posterior Simulation

We propose the adaptive Sequential Monte Carlo (SMC) approach (Durham and Geweke, 2013) as a posterior simulation strategy, for two main reasons. First, many utility functions of the form in (2.1) give identical predictions. We cannot expect these functions to be close, in the sense of a practical MCMC chain passing from one to the other in a small number of steps with reasonable probability. Having many particles exploring the parameter space solves that problem. Second, a sequential inference method is required for our tempering approach.

The sequence of target distributions is obtained by setting  $\zeta = 2$  for the first distribution, and by increasing  $\zeta$  by a factor of 2 at each iteration:

$$Q(\psi|\{x_{nt}^*\}, \zeta), \text{ for } \zeta = 2, 4, 8, \dots$$

#### 4.3.1 Algorithm

The algorithm uses a population of  $J$  groups of  $P$  particles, each of them containing a set of realizations of all our parameters and some simulated hyperparameters that do not have a closed form posterior:

$$\psi_{jp} = \left\{ \alpha^\theta, \sigma^\theta, K_+, \{\Lambda_k, \pi_k^z, \{z_{nk}, s_{nk}, \theta_{nk}, \{\epsilon_{nkt}\}_{t=1}^T\}_{n=1}^N\}_{k=1}^{K_+} \right\}.$$

The algorithm starts with an Initialization ( $I$ ) phase and then iterates over Correction ( $C$ ), Selection ( $S$ ) and Mutation ( $M$ ) phases.

The first group of particles will be used to determine when to move from the C phase to the S phase based on Effective Sample Size (ESS). The ESS for the first group of particles will be denoted  $ESS_1$ . The simulated values of our first group of particles will be discarded in our analysis of the simulated posterior.

The general structure of the algorithm is as follow:

- Initialization phase: Initialize  $J \times P$  particles  $\{\{\psi_{jp}\}_{p=1}^P\}_{j=1}^J$
- For  $\eta = 1$  to  $N$  (*Iterate over individuals*)
  - Draw  $\{z_{\eta k}, s_{\eta k}, \theta_{\eta k}\}$  for the newly added individual (*see section 4.3.3*)
  - For  $\tau = 1$  to  $T$  (*Iterate over observations for each individual*)
    - \* Draw  $\epsilon_{\eta k\tau}$  for the newly added observation (*see section 4.3.3*)
    - \* Set  $\zeta = 1$

- \* Repeat until  $\zeta = \infty$ :
  - Correction phase: increase  $\zeta$  until  $ESS_1/P < 0.5$  or  $\zeta = \infty$  (see section 4.3.4)
  - Selection phase: re-sample particles within each group based on weight
  - Mutation phase: move each particle with an Metropolis-Hastings within Gibbs step using the target density  $Q(\psi_{jp}|\{x_{nt}^*\}; \eta, \tau, \zeta)$

The various parts of the algorithm are detailed in the following sections.

### 4.3.2 Initialization phase

Applied to every particle:

- Draw  $\alpha^\theta \sim \text{Gamma}(1, 1)$
- Draw  $\sigma^\theta \sim \text{Gamma}(1, 1)$
- Set  $K_+ = 0$

### 4.3.3 Drawing $\{z_{\eta k}, s_{\eta k}, \theta_{\eta k}\}$ and $\{\epsilon_{\eta k \tau}\}$

To draw the  $z_{\eta k}$  we follow the IBP approach. As detailed in section 3.2, when  $\Lambda_k$  has an active monotone attribute, we set  $m_k = 1$ . When the monotone active attributes have upper active levels (all levels above a specified level are active), we set  $\rho_k = 1$ . When the monotone active attributes have lower active levels (all levels below a specified level are active), we set  $\rho_k = 0$ .

- For  $k = 1$  to  $K_+$ 
  - If  $\sum_{n=1}^{\eta-1} z_{nk} > 0$ :
    - \* Draw  $z_{nk} \sim \text{Bernoulli}(\frac{\sum_{n=1}^{\eta-1} z_{nk}}{\beta^z + \eta - 1})$
    - \* If  $z_{nk} = 1$ :
      - Draw  $s_{\eta k} \sim \begin{cases} 0 & \text{if } m_k = 1 \text{ and } \rho_k = 1 \\ 1 & \text{if } m_k = 1 \text{ and } \rho_k = 0 \\ \text{Bernoulli}(\frac{1 + \sum_{n=1}^{\eta-1} s_{nk}}{2 + \eta - 1}) & \text{otherwise} \end{cases}$
      - Draw  $\theta_{\eta k} \sim \text{Gamma}(\sigma^\theta \alpha^\theta, 1)$
  - Else :
    - \* Set  $z_{nk} = 0$
- If  $K_+ = 0$ :
  - While  $k^* < 1$  : Draw  $k^* \sim \text{Poisson}(\frac{\alpha^z \beta^z}{\beta^z + \eta - 1})$
- Else:

- Draw  $k^* \sim \text{Poisson}(\frac{\alpha^z \beta^z}{\beta^z + \eta - 1})$
- Set  $K_+ = K_+ + k^*$
- For  $k = K_+ - k^*$  to  $K_+$ 
  - Det  $z_{\eta k} = 1$
  - Draw  $\Lambda_k$
  - Draw  $s_{\eta k} \sim \begin{cases} 0 & \text{if } m_k = 1 \text{ and } \rho_k = 1 \\ 1 & \text{if } m_k = 1 \text{ and } \rho_k = 0 \\ \text{Bernoulli}(\frac{1 + \sum_{n=1}^{\eta-1} s_{nk}}{2 + \eta - 1}) & \text{otherwise} \end{cases}$
  - Draw  $\theta_{\eta k} \sim \text{Gamma}(\sigma^\theta \alpha^\theta, 1)$

In order to draw the  $\{\epsilon_{\eta k \tau}\}$  when a new observation is added, we follow the following simple algorithm:

- For  $k = 1$  to  $K_+$ 
  - If  $z_{\eta k} = 1$  : Draw  $\epsilon_{\eta k \tau} \sim \text{Gamma}((1 - \sigma^\theta) \alpha^\theta, 1)$

#### 4.3.4 Correction phase

Applied to every particle:

- Repeat until  $ESS_1/P < 0.5$  or  $\zeta = \infty$ :
  - compute utility differences  $du_{jp} = u(x_{\eta \tau}^*, \psi_{jp}) - \max_{x \in X_{\eta \tau}} u(x, \psi_{jp})$
  - If  $\sum_p du_{1p} = 0$ :
    - \* Set  $\zeta = \infty$
    - \* compute weight  $w_{jp} = \mathbb{1}\{du_{jp} = 0\}$
  - Else:
    - \* Set  $\zeta = 2 \cdot \zeta$
    - \* compute weight  $w_{jp} = w_{jp} \cdot 2^{du_{jp}}$
- compute  $ESS_1/P = \frac{(\sum_p w_{1p})^2}{\sum_p w_{1p}^2}$

#### 4.3.5 Selection phase

The selection phase draws a new set of  $P$  within group  $j$  particles with replacement from the set of old particles using weights  $w_{jp}$ . For each group  $j$  of particle:

- For each group  $j$  of particle:

- For  $p' = 1$  to  $P$ 
  - \* Draw a particle  $\psi^*$  from set  $\{\psi_{jp}\}_{p=1}^P$  using weights  $\{w_{jp}\}_{p=1}^P$
  - \* Set  $\psi_{jp'} = \psi^*$

#### 4.3.6 Mutation phase

The mutation phase do an MCMC step using Gibbs sampling and Metropolis-Hastings (MH) accept-reject algorithm.

- For  $n = 1$  to  $\eta$ 
  - For  $k = 1$  to  $K_+$ 
    - \* If  $\sum_{n' \neq n} z_{n'k} > 0$ :
      - Set  $\psi^* = \psi_{jp}$
      - Draw  $z_{nk}^* \sim \text{Bernoulli}(\frac{\sum_{n' \neq n} z_{n'k}}{\beta^z + \eta - 1})$
      - If  $z_{nk} = 1$ : Draw  $\{s_{nk}^*, \theta_{nk}^*, \{\epsilon_{nkt}^*\}\}$  (*see below*)
      - Store new values into  $\psi^*$
      - Set  $\psi_{jp} = \psi^*$  with probability  $\min\{1, \frac{Q(\{x_{nt}^*\}|\psi^*; \eta, \tau, \zeta)}{Q(\{x_{nt}^*\}|\psi_{jp}; \eta, \tau, \zeta)}\}$
    - \* Else :
      - Set  $z_{nk} = 0$
  - Draw  $k^* \sim \text{Poisson}(\frac{\alpha^z \beta^z}{\beta^z + \eta - 1})$
  - Set  $K_+ = K_+ + k^*$
  - For  $k = K_+ - k^*$  to  $K_+$ 
    - \* Set  $\psi^* = \psi_{jp}$
    - \* Set  $z_{nk}^* = 1$
    - \* Draw  $\Lambda_k^*$
    - \* Draw  $s_{nk}^* \sim \begin{cases} 0 & \text{if } \Lambda_k^* \text{ has active upper monotone attributes} \\ 1 & \text{if } \Lambda_k^* \text{ has active lower monotone attributes} \\ \text{Bernoulli}(\frac{1 + \sum_{n=1}^{\eta-1} s_{nk}}{2 + \eta - 1}) & \text{otherwise} \end{cases}$
    - \* Draw  $\theta_{nk}^* \sim \text{Gamma}(\sigma^\theta \alpha^\theta, 1)$
    - \* Store new values into  $\psi^*$
    - \* Set  $\psi_{jp} = \psi^*$  with probability  $\min\{1, \frac{Q(\{x_{nt}^*\}|\psi^*; \eta, \tau, \zeta)}{Q(\{x_{nt}^*\}|\psi_{jp}; \eta, \tau, \zeta)}\}$
- Draw  $\alpha^* \sim \text{Gamma}(1, 1)$
- Draw  $\sigma^* \sim \text{Gamma}(1, 1)$

- Set  $(\alpha^\theta, \sigma^\theta) = (\alpha^*, \sigma^*)$  with probability

$$\min \left\{ 1, \prod_k^{K_+} \prod_n^\eta \left( \frac{P_\Gamma(\theta_{nk} | \alpha^* \sigma^*, 1) \prod_t^{T_n} P_\Gamma(\epsilon_{nkt} | \alpha^* (1 - \sigma^*), 1)}{P_\Gamma(\theta_{nk} | \alpha^\theta \sigma^\theta, 1) \prod_t^{T_n} P_\Gamma(\epsilon_{nkt} | \alpha^\theta (1 - \sigma^\theta), 1)} \right)^{z_{nk}} \right\}$$

where  $P_\Gamma(\cdot | \alpha, \beta)$  is the PDF of the  $Gamma(\alpha, \beta)$  distribution, and  $T_n = \tau$  if  $n = \eta$  and  $T_n = T$  otherwise.

For this phase, in order to draw  $\{s_{nk}^*, \theta_{nk}^*, \{\epsilon_{nkt}^*\}\}$ , we propose the following process:

- Draw  $s_{nk}^* \sim \begin{cases} 0 & \text{if } m_k^* = 1 \text{ and } \rho_k^* = 1 \\ 1 & \text{if } m_k^* = 1 \text{ and } \rho_k^* = 0 \\ Bernoulli(\frac{1 + \sum_{n' \neq n}^\eta s_{nk}}{2 + \eta - 1}) & \text{otherwise} \end{cases}$
- Draw  $\theta_{nk}^* \sim Gamma(\sigma^\theta \alpha^\theta, 1)$
- If  $n < \eta$ :
  - For  $t = 1$  to  $T$ :
    - \*  $\epsilon_{nkt}^* \sim Gamma((1 - \sigma^\theta) \alpha^\theta, 1)$
- Else:
  - For  $t = 1$  to  $\tau$ :
    - \*  $\epsilon_{nkt}^* \sim Gamma((1 - \sigma^\theta) \alpha^\theta, 1)$

## References

- G. Durham and J. Geweke. Adaptive sequential posterior simulators for massively parallel computing environments. *Available at SSRN 2251635*, 2013.
- T. L. Griffiths and Z. Ghahramani. The indian buffet process: An introduction and review. *The Journal of Machine Learning Research*, 12:1185–1224, 2011.