

# Bayesian Automatic Relevance Determination for Utility Function Specification in Discrete Choice Models

Rodrigues, Filipe; Ortelli, Nicola; Bierlaire, Michel; Pereira, Francisco

*Published in:* I E E E Transactions on Intelligent Transportation Systems

Link to article, DOI: 10.1109/TITS.2020.3031965

Publication date: 2022

Document Version Peer reviewed version

Link back to DTU Orbit

Citation (APA):

Rodrigues, F., Ortelli, N., Bierlaire, M., & Pereira, F. (2022). Bayesian Automatic Relevance Determination for Utility Function Specification in Discrete Choice Models. *I E E Transactions on Intelligent Transportation Systems*, *23*(4), 3126 - 3136. https://doi.org/10.1109/TITS.2020.3031965

### **General rights**

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

• Users may download and print one copy of any publication from the public portal for the purpose of private study or research.

- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

# Bayesian Automatic Relevance Determination for Utility Function Specification in Discrete Choice Models

Filipe Rodrigues<sup>(D)</sup>, Nicola Ortelli, Michel Bierlaire, and Francisco Camara Pereira<sup>(D)</sup>, *Member, IEEE* 

Abstract-Specifying utility functions is a key step towards applying the discrete choice framework for understanding the behaviour processes that govern user choices. However, identifying the utility function specifications that best model and explain the observed choices can be a very challenging and timeconsuming task. This paper seeks to help modellers by leveraging the Bayesian framework and the concept of automatic relevance determination (ARD), in order to automatically determine an optimal utility function specification from an exponentially large set of possible specifications in a purely data-driven manner. Based on recent advances in approximate Bayesian inference, a doubly stochastic variational inference is developed, which allows the proposed MNL-ARD model to scale to very large and high-dimensional datasets. Using semi-artificial choice data, the proposed approach is shown to be able to accurately recover the true utility function specifications that govern the observed choices. Moreover, when applied to real choice data, MNL-ARD is able discover high quality specifications that can outperform previous ones from the literature according to multiple criteria, thereby demonstrating its practical applicability.

*Index Terms*—Discrete choice models, automatic relevance determination, automatic utility specification, doubly stochastic variational inference.

# I. INTRODUCTION

**D** ISCRETE choice models (DCM) provide a powerful framework for understanding user behaviour. By modelling user choices as functions of the alternative-specific attributes and user characteristics, DCMs allow researchers to predict users' future choices given a set of discrete alternatives and understand the behaviour process that governs their choices. Hence, it is without surprise that DCMs have become a widely adopted framework in various domains ranging from psychology to economics, thus making them one of the main work-horses for understanding user travel behaviour, consumer behaviour, and many other kinds of user choices.

In practice, a fundamental part of applying the DCM framework consists in specifying the utility function for each

Manuscript received May 15, 2020; revised August 23, 2020; accepted October 14, 2020. The Associate Editor for this article was J. Xun. (Corresponding author: Filipe Rodrigues.)

Filipe Rodrigues and Francisco Camara Pereira are with the Technical University of Denmark (DTU), 2800 Lyngby, Denmark (e-mail: rodr@dtu.dk).

Nicola Ortelli and Michel Bierlaire are with the École Polytechnique Fédérale de Lausanne (EPFL), 1015 Lausanne, Switzerland.

This article has supplementary downloadable material available at http://ieeexplore.ieee.org, provided by the authors.

Digital Object Identifier 10.1109/TITS.2020.3031965

alternative in the choice set, which are generally assumed to be known a priori. For the sake of interpretability, these utility functions are typically assumed to be functions of a set of explanatory variables that are linear in parameters. Although limiting at first sight, this linear framework can be made rather powerful by exploring variable transformations (e.g. log-transformations, Box-Cox transformations), one-hot encodings, piecewise linear representations, discretizations, interactions between variables, etc. However, all these modelling choices quickly raise the number of possible utility function specifications beyond manageable values for the modeller. On the other hand, given the central role of the utility functions in DCMs, it is essential to determine good specifications, at the risk of obtaining misspecified models and biased parameter estimates [1]. As a consequence, a modeller often spends large portions of time seeking the "best" specification according to different criteria (e.g. convergence, log-likelihood, p-values), typically through a combination of trial-and-error and domain knowledge (e.g. economic theories).

In this paper, we focus on the Multinomial Logit (MNL) model, and we propose leveraging the Bayesian framework in order to automatically determine an optimal utility function specification from an exponentially large set of possible specifications in a purely data-driven manner. Although the proposed approach is not meant to be a replacement for expert intuition and domain knowledge, it is shown to provide key insights about the data that can help the modeller determine the utility function specification that best represents the observed choice data, which can ultimately lead to new understandings about the way people make choices in certain contexts.

Based on the principle of Automatic Relevance Determination (ARD), as developed by [2] in the context of the Relevance Vector Machine and as widely used in the Gaussian Processes literature [3], we propose the use of a hierarchical prior on the preference parameters of each utility function in order to automatically determine their relevance for explaining the observed choice data. The key idea consists in jointly estimating the posterior distribution over the preference parameters, as well as the optimal values for the variances of the Gaussian priors over each possible explanatory variable to be included in each utility function specification. In order to ensure consistency among the selected variables, i.e. that either all or none of the dimensions corresponding to the representation of a given explanatory variable are selected, we propose

1524-9050 © 2020 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information. tying the variance parameters of the Gaussian priors over the parameters that correspond to the same representation of a given choice attribute. Given the estimated optimal values for the variances of the Gaussian priors for a very large set of possible variable representations, a modeller can easily determine the most relevant attributes and corresponding representations for explaining a dataset of observed choices by simply selecting the variables for which the estimated prior variances are non-zero.

Since exact Bayesian inference in the proposed MNL-ARD model is intractable, we propose the use of the variational inference framework. Namely, we develop an efficient approximate inference algorithm using doubly stochastic variational inference (DSVI) [4]. By combining the theory of variational inference with the theory of stochastic optimization, the proposed inference algorithm is able to approximate the true posterior distribution over the preference parameters with a tractable distribution and jointly estimate the optimal Gaussian prior hyper-parameters, while being able to scale to very large datasets with a very high number of dimensions. Although we focus on Multinomial Logit (MNL) models, the proposed approach can be extended to more complex models such as Mixed and Latent Class Logit models.

The validity of the proposed automatic utility function specification framework is empirically demonstrated using both semi-artificial and real choice data. We begin by empirically demonstrating the ability of the proposed approach to discover the correct utility function specifications through an extensive series of experiments on simulated choice data based on the Swissmetro dataset [5]. In particular, we manually specify a series of "artificial" (but realistic) utility function specifications of increasing complexity and, based on the Swissmetro dataset, we sample new artificial choices according to the manuallyspecified utility functions. Our empirical results show that the proposed MNL-ARD model is able to very accurately recover the "true" specifications that were used to generate the artificial choices, even in settings where the number of variables representations and transformations considered for each utility function is in the order of the thousands. Lastly, our empirical results on the real choices from the Swissmetro dataset demonstrate the potential of the proposed framework for discovering novel utility function specifications that can potentially outperform previous ones from the state of the art in terms of explanatory power and generalization to unobserved data.

In summary, the main contributions of this paper are the following:

- We adapt the theory of ARD to the domain of DCMs, making the modifications that are required from a choice modelling perspective (e.g. multiple utility functions with alternative-specific attributes, variable number of dimensions and tied parameters in the hierarchical priors);
- We develop a doubly-stochastic variational inference (DSVI) [4] algorithm for performing fast approximate inference in the proposed MNL-ARD model;
- We empirically show (i) the ability of the the proposed approach to recover the true utility specifications on semiartificial choice data, (ii) that MNL-ARD can discover

new specifications that outperform previous ones from the literature, and (iii) that the developed DSVI algorithm is able to scale to very large datasets and search spaces.

The remainder of this paper is organized as follows. In the next section, we review the relevant literature for this work. Section III presents the proposed MNL-ARD model and derives a scalable doubly-stochastic variational inference algorithm for performing fast approximate Bayesian inference on it. The corresponding experimental results are presented in Section IV. The paper ends with the conclusions (Section V).

### II. LITERATURE REVIEW

The problem of automatically determining the relevant variables for inclusion in a model has been studied to a significant extent in the supervised machine learning literature under the common title of "feature selection". When using feature selection techniques, the main premise is that the considered data contain redundant or irrelevant variables, which can therefore be removed without consequent loss of information [6]. The numerous existing approaches are generally classified as wrapper, filter and embedded methods according to the strategy they employ to search for subsets of variables [7]. Wrappers use the model of interest to score subsets according to the predictive power they allow to achieve. Despite being computationally intensive, wrappers offer a simple way of addressing the problem: a plethora of methods based on simulated annealing [8], tabu search [9], evolutionary algorithms [10] and other combinatorial optimization algorithms have already been applied successfully, both for linear and logistic regressions. In comparison, filter methods are independent of the model under consideration; they use "proxy" measures such as correlation or mutual information [11] to evaluate single features or subsets. While being less computationally intensive than wrappers, filters usually achieve worse results in terms of prediction power. Finally, embedded methods are characterized by the fact that the selection of variables and the estimation of the model are performed simultaneously, in a single process. A good example of such class of methods is the LASSO, initially proposed by [12] and successfully applied both to linear [13] and logisitic [14] regressions. Other existing embedded methods make use of mixed integer optimization [15] or decision trees [16] to effectively incorporate feature selection as part of the training process.

In the field of discrete choice analysis, interest has recently emerged for methods that are able to "mitigate" the need for presumptive structural assumptions. Two main directions of research are explored in the existing literature: the first substitutes DCMs with machine learning classifiers that do not require any prior knowledge concerning the domain [17]–[19], while the second focuses on automatizing the utility specification of DCMs by means of data-driven feature selection algorithms [20]–[22].

A particularly elegant class of methods for performing automatic feature selection in the statistics and machine learning literature relies on the concept of automatic relevance determination (ARD) [2], [23], [24]. The idea behind this class of approaches consists in specifying the a-priori uncertainty and infer a-posteriori uncertainty about regression coefficients explicitly and hierarchically in a Bayesian framework. However, unfortunately, Bayesian inference in such hierarchical models quickly becomes intractable, and effective and scalable methods are required in order to perform approximate inference. To that end, [24] presents a type-II maximum likelihood based on variational inference in a linear regression context, where the hyper-parameters of the hierarchical priors are tuned by maximizing the marginal likelihood of the data. This approach was later extended by [25] to a fully Bayesian approach by further considering a normal inversegamma prior over the hyper-parameters of the hierarchical priors, and then performing variational inference to determine the corresponding posterior distributions. Furthermore, the author also considers ARD in a binary logistic regression context. The difficulty in the latter stems from the nonconjugacy of the sigmoid, which required the authors to consider an additional model-specific parametric lower bound on the sigmoid as proposed by [26], which can raise the computational cost and compromise accuracy. Recently, highly efficient general-purpose black-box variational inference methods have proposed in the literature [4], [27], which allow for approximating the required expectations using inexpensive Monte Carlo approximations. In particular, [4] proposed a doubly stochastic variational inference for performing ARD in binary logistic regression. The approach proposed in this paper builds on the work of [4] to propose an ARD framework for discrete choice models, and to develop a corresponding efficient variational inference algorithm.

### III. APPROACH

#### A. Discrete Choice Models

Following the Random Utility Maximization (RUM) theory, discrete choice models are based on the assumption that each individual  $n \in \{1, ..., N\}$  is a rational decision-maker that aims at maximizing some utility with respect to a choice set  $C_n$ . A key step in discrete choice modeling is then to specify a function  $U_{in}$  that is able to capture the utility of each alternative i for each individual n. The utility function is further assumed to be partitioned into two components,  $U_{in} = V_{in} + \epsilon_{in}$ , where  $V_{in}$  is a systematic (or deterministic) utility and  $\epsilon_{in}$  is an i.i.d. term that captures the uncertainty stemming from the impossibility of  $V_{in}$  to fully capture the choice context. The systematic component  $V_{in}$  is typically assumed to be a function linear in parameters of the observable explanatory variables  $\mathbf{x}_{in} = \{\mathbf{x}_{din}\}_{d=1}^{D_i}$  of the utility of alternative i for each individual n (e.g. alternative attributes, individual's socio-demographic characteristics, etc.):

$$V_{in} = \boldsymbol{\beta}_i^{\mathrm{T}} \mathbf{x}_{in} = \sum_{d=1}^{D_i} \beta_{di} \, x_{din}, \tag{1}$$

where  $\beta_i$  is a vector of alternative-specific preference parameters. This corresponds to the more general setting where preference parameters may vary between different alternatives. Following the same reasoning, our specification further allows for a variable number of explanatory variables  $D_i$  per alternative *i*. Under the standard Multinomial Logit assumption that  $\epsilon_{in} \sim \text{EV}(0, 1)$ , the probability of individual *n* selecting alternative *i* is given by

$$P_n(i) = \frac{e^{V_{in}}}{\sum_{j \in \mathcal{C}_n} e^{V_{jn}}}.$$
(2)

Given a dataset of observed choices and corresponding explanatory variables for a population of size N, the modeler's objective is to determine the preference parameters  $\beta$ , which are typically estimated by maximizing the log-likelihood:

$$\boldsymbol{\beta}^* = \arg \max_{\boldsymbol{\beta}} \sum_{n=1}^{N} \sum_{i \in \mathcal{C}_n} y_{in} \log P_n(i), \qquad (3)$$

where  $y_{in}$  is a one-hot encoding of the observed choice for the  $n^{th}$  individual and **y** and  $\beta$  are used to denote the set of all observed choices and preference parameters, respectively.

Despite the appealing simplicity of maximum likelihood estimation methods, in this paper we shall follow a Bayesian approach. The latter not only allows us to infer full posterior distributions for the preference parameters  $\beta$  that provide a principled way of performing hypotheses testing [28] and uncertainty quantification, but also enable online learning approaches in which the posterior over the parameters is continuously updated as more data becomes available [29]. Moreover, most importantly, it will support the development of the automatic utility function specification approach based on ARD proposed in Section III-B.

We begin by introducing the standard Bayesian framework for the MNL model specified above, which will serve as the starting point for the proposed approach in Section III-B. To enable the Bayesian treatment of model above, we start by placing a prior distribution over the preference parameters for each of the alternatives:

$$\boldsymbol{\beta}_i \sim \mathcal{N}(\boldsymbol{\beta}_i | \mathbf{0}, \lambda \mathbf{I}), \tag{4}$$

where I denotes the identity matrix, thus making  $\lambda I$  a diagonal covariance matrix parametrized by  $\lambda$ . Making use of Bayes' theorem, the posterior distribution over the preference parameters  $\beta$  is

$$p(\boldsymbol{\beta}|\mathbf{y},\lambda) = \frac{p(\boldsymbol{\beta}|\lambda) \prod_{n=1}^{N} \prod_{i \in \mathcal{C}_{n}} (P_{n}(i))^{y_{in}}}{\int p(\boldsymbol{\beta}|\lambda) \prod_{n=1}^{N} \prod_{i \in \mathcal{C}_{n}} (P_{n}(i))^{y_{in}} d\boldsymbol{\beta}}.$$
 (5)

However, the non-conjugacy between the prior (4) and the MNL likelihood in (2) deems the integral in the denominator intractable, thus making exact inference infeasible. Fortunately, over recent years, we have observed very significative improvements in the accuracy and scalability of approximate Bayesian inference, which we shall exploit in Section III-C.

# B. Automatic Utility Function Specification

The main of focus of this paper is on leveraging the Bayesian framework and the concept automatic relevance determination (ARD) [2] to lift the burden of manually searching for an optimal utility function specification for a given discrete choice problem from the modeler. Namely, we wish to automatically determine the relevant variables for the utility function of each alternative *i*, while considering also for different non-linear transformations (e.g. log-transforms, Box-Cox transforms), different continuous variable discretizations, interactions between variables, etc. In order to allow for some of these modeling options and, in particular, variable interactions, let us begin by considering a more flexible parameterization of the utility function. Letting  $s_n$  be a categorical socio-economic variable with *K* categories associated with individual *n* (e.g. age, income, education or profession), we can allow for interactions with the remaining variables by introducing an unknown parameter per category  $\beta_1, \ldots, \beta_K$  and defining the utility function for an alternative *i* as

$$V_{in} = \sum_{d=1}^{D_i} \sum_{k=1}^{K_d} \beta_{kdi} \,\delta_k(s_n) \,h(x_{din}), \tag{6}$$

where  $\delta_k(s_n)$  is an indicator function, which takes the value 1 if the  $n^{th}$  individual belongs to category k and 0 otherwise, and  $h(\cdot)$  is an arbitrary function (e.g. logarithm for a log-transform). Kindly notice that the utility specification in (1) is a special case of (6), when  $K_d = 1$  and  $h(\cdot)$  is the identity function. Similarly, this specification also contains one-hot encodings and discretizations of a variable d as special cases by adapting the functions  $\delta_k(\cdot)$  and  $h(\cdot)$  accordingly.

Based on (6), the problem of automatic utility function specification can then be defined as determining the subset of input dimensions  $S_i \subseteq \{1, \ldots, D_i\}$  that best models the observed choices according to a dataset of observed choices, where  $\{1, \ldots, D_i\}$  is a very large set of possible variable transformations and representations whose usefulness to the model we wish to test. For example, for a cost variable, a modeler may consider including in  $\{1, \ldots, D_i\}$  the variable itself, its log-transformed value, cost interacted with gender, cost interacted with age, cost interacted with both gender and age, a piecewise linear transformation, etc. The goal is then to determine which subset  $S_i$  of these should be included in the utility function specification  $V_i$ .

The starting point for our proposed approach is the concept of automatic relevance determination (ARD), as used for instance in the statistical machine learning literature for the relevance vector machine [2]. The key idea lies in realizing that preference parameters of irrelevant dimensions d should be pushed towards zero. However, the standard prior specification in (4) is too restrictive to allow for some parameters to be pushed arbitrarily close to zero, while others retain their actual values. This restriction stems for the fact in (4), the parameters are assumed to have independent univariate Gaussian priors that share the same prior variance  $\lambda$ . Therefore, we can make progress towards ARD in discrete choice models by constructing a flexible hierarchical prior, in which each parameter is assigned an independent Gaussian prior with its own variance, but parameters belonging to the representation of the same variable share the same variance. Mathematically, this corresponds to

$$\beta_{kdi} \sim \mathcal{N}(\beta_{kdi}|0,\lambda_{di}). \tag{7}$$

Please note that the constraint of sharing the same variance over the index k is crucial in order to ensure that the entire



Fig. 1. Graphical model representation of the proposed model.

group is treated as a whole, i.e. either all k "sub-dimensions" of a variable d are deemed relevant by the model, or none is and their corresponding parameters are all pushed towards zero. The prior over all the preference parameters is then

$$p(\boldsymbol{\beta}|\boldsymbol{\lambda}) = \prod_{i \in \mathcal{C}} \prod_{d=1}^{D_i} \prod_{k=1}^{K_d} \mathcal{N}(\beta_{kdi}|0, \lambda_{di}), \qquad (8)$$

where  $\lambda$  is used to denote the set of all  $\lambda_{di}$ . While one could further place a Gamma prior over the precisions  $\lambda_{di}^{-1}$ , we refrain from doing so because (i) it would introduce a new set of hyper-parameters to specify and (ii), as we shall see in Section III-C, it is possible to optimize over the variance parameters  $\lambda$  analytically. Hence, we shall continue by treating the latter as point parameters rather than random variables in a fully Bayesian setting. The generative process of the proposed model can then be summarized as follows:

- 1) For each alternative i in the entire choice set C
  - a) For each variable  $d \in \{1, \ldots, D_i\}$ 
    - i) Set preference parameter variance  $\lambda_{di}$
    - ii) For each category  $k \in \{1, \ldots, K_d\}$
    - A) Draw pref. param.  $\beta_{kdi} \sim \mathcal{N}(\beta_{kdi}|0, \lambda_{di})$
- 2) For each individual  $n \in \{1, ..., N\}$  (a)

a) Draw observed choice  $y_n \sim \text{Categorical}(y_n|P_n)$ In order to emphasize the hierarchical structure of the proposed model, Figure 1 shows a graphical model representation.

Based on the model specification above, our goal is to be able to jointly infer the preference parameters  $\beta$  and estimate the variance parameters  $\lambda_{di}$  for each explanatory variable, in order to assess which ones should be included in each utility function  $V_i$ . As for the "standard" MNL model in Section III-A, performing exact Bayesian inference in the proposed model is intractable. Therefore, we shall proceed by developing an approximate Bayesian inference algorithm using doubly stochastic variational inference (DSVI) [4].

# C. Doubly Stochastic Variational Inference

The intractability of exact inference for the proposed model stems from the impossibility of obtaining an analytical expression for the marginal likelihood in the denominator of (5), which for the proposed ARD model takes the form

$$p(\mathbf{y}|\boldsymbol{\lambda}) = \int_{\boldsymbol{\beta}} \left( \prod_{i \in \mathcal{C}} \prod_{d=1}^{D_i} \prod_{k=1}^{K_d} \mathcal{N}(\beta_{kdi}|0, \lambda_{di}) \right) \prod_{n=1}^N \prod_{i \in \mathcal{C}_n} P_n(i)^{y_{in}}.$$

In order to obtain an efficient and scalable approximate inference algorithm that is able to cope with large datasets and with very high dimensionalities  $D_i$ , we propose the use of the variational inference framework [30].

Variational inference, or variational Bayes, constructs an approximation to the true posterior distribution  $p(\beta|\mathbf{y})$  by considering a family of tractable distributions  $q(\beta)$ , which can be obtained by relaxing some constraints in the true distribution. In this case, we shall assume the variational distribution  $q(\beta)$  to be a fully-factorized (mean-field) approximation:

$$q(\boldsymbol{\beta}|\boldsymbol{\mu}, \mathbf{c}) = \prod_{i \in \mathcal{C}} \prod_{d=1}^{D_i} \prod_{k=1}^{K_d} \mathcal{N}(\beta_{kdi}|\mu_{kdi}, c_{kdi}), \qquad (9)$$

with variational parameters  $\mu$  and **c**. The goal is then to find the parameters of the variational distribution so that the approximation becomes as close as possible to the true posterior, thereby reducing inference to an optimization problem.

The closeness between the approximate posterior  $q(\boldsymbol{\beta}|\boldsymbol{\mu}, \mathbf{c})$  and the true posterior  $p(\boldsymbol{\beta}|\mathbf{y})$  can be measured by the Kullback-Leibler (KL) divergence [31]:  $\mathbb{KL}(q(\boldsymbol{\beta}|\boldsymbol{\mu}, \mathbf{c})||p(\boldsymbol{\beta}|\mathbf{y}))$ . Unfortunately, this KL divergence cannot be minimized directly. However, we can find a function that we can minimize, which is equal to it up to an additive constant, as follows:

$$\begin{split} \mathbb{K} \mathbb{L}(q(\boldsymbol{\beta}|\boldsymbol{\mu}, \mathbf{c})||p(\boldsymbol{\beta}|\mathbf{y})) \\ &\triangleq \mathbb{E}_q \bigg[ \log \frac{q(\boldsymbol{\beta}|\boldsymbol{\mu}, \mathbf{c})}{p(\boldsymbol{\beta}|\mathbf{y})} \bigg] \\ &= \mathbb{E}_q [\log q(\boldsymbol{\beta}|\boldsymbol{\mu}, \mathbf{c})] - \mathbb{E}_q [\log p(\boldsymbol{\beta}|\mathbf{y})] \\ &= -(\underbrace{\mathbb{E}_q [\log p(\boldsymbol{\beta}, \mathbf{y})] - \mathbb{E}_q [\log q(\boldsymbol{\beta}|\boldsymbol{\mu}, \mathbf{c})]}_{\mathcal{L}(\boldsymbol{\mu}, \mathbf{c}, \lambda)}) + \underbrace{\log p(\mathbf{y})}_{\text{const.}}. \end{split}$$

The log  $p(\mathbf{y})$  term does not depend on the variational parameters and thus can be ignored. Minimizing the KL divergence is then equivalent to maximizing  $\mathcal{L}(\boldsymbol{\mu}, \mathbf{c}, \boldsymbol{\lambda})$ , which is referred to as the evidence lower bound (ELBO). We can further rewrite  $\mathcal{L}(\boldsymbol{\mu}, \mathbf{c}, \boldsymbol{\lambda})$  as a function of simpler terms by exploiting the factorization of the joint and prior distributions, yielding

$$\mathcal{L}(\boldsymbol{\mu}, \mathbf{c}, \boldsymbol{\lambda}) = \sum_{i \in \mathcal{C}} \sum_{d=1}^{D_i} \sum_{k=1}^{K_d} \mathbb{E}_{q(\boldsymbol{\beta})}[\log \mathcal{N}(\boldsymbol{\beta}_{kdi} | 0, \lambda_{di})] + \sum_{n=1}^{N} \sum_{i \in \mathcal{C}_n} y_{in} \mathbb{E}_{q(\boldsymbol{\beta})}[\log P_n(i)] - \sum_{i \in \mathcal{C}} \sum_{d=1}^{D_i} \sum_{k=1}^{K_d} \mathbb{E}_{q(\boldsymbol{\beta})}[\log \mathcal{N}(\boldsymbol{\beta}_{kdi} | \mu_{kdi}, c_{kdi})].$$
(10)

Our goal is then to find the variational parameters { $\mu$ , **c**} and the hyper-parameters  $\lambda$  that maximize  $\mathcal{L}(\mu, \mathbf{c}, \lambda)$ . However, due to the log-sum-exp term resultant from the denominator of the MNL kernel, the expectation  $\mathbb{E}_{q(\beta)}[\log P_n(i)]$  in (10) is still intractable. While some authors proposed the use of complex approximations to further bound this term [32], [33], we shall rely on a more efficient and scalable approximation based on the theory of stochastic optimization. In order to enable it, we begin by reparameterizing our approximate distribution in (9). Consider a random variable  $z \sim \mathcal{N}(z|0, 1)$ . We can change the mean and variance by applying an invertible transformation  $\beta = cz + \mu$  and making use of the change of variables formula for a random vector, which states that for a given function f(x), and given an invertible transformation y = h(x), we have that  $f(y) = f(h(x))|J_{h^{-1}}|$ , where  $|J_{h^{-1}}|$  denotes the determinant of the Jacobian matrix of the inverse transformation  $h^{-1}$ . Hence, given the transformation  $\beta = cz + \mu$  and its inverse  $z = c^{-1}(\beta - \mu)$ , we can rewrite the approximate distribution in (9) as

$$q(\boldsymbol{\beta}|\boldsymbol{\mu}, \mathbf{c}) = \prod_{i \in \mathcal{C}} \prod_{d=1}^{D_i} \prod_{k=1}^{K_d} \frac{1}{|c_{kdi}|} \mathcal{N}(c_{kdi}^{-1}(\beta_{kdi} - \mu_{kdi})|0, 1).$$
(11)

By plugging (11) into (10) and changing variables according to  $z = c^{-1}(\beta - \mu)$ , we can rewrite  $\mathcal{L}(\boldsymbol{\mu}, \mathbf{c}, \boldsymbol{\lambda})$  as follows:

$$\mathcal{L}(\boldsymbol{\mu}, \mathbf{c}, \boldsymbol{\lambda}) = \mathbb{E}_{\mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})}[\log p(\mathbf{y}|\mathbf{c} \circ \mathbf{z} + \boldsymbol{\mu})] + \sum_{i \in \mathcal{C}} \sum_{d=1}^{D_i} \sum_{k=1}^{K_d} \log c_{kdi} - \frac{1}{2} \sum_{i \in \mathcal{C}} \sum_{d=1}^{D_i} K_d \log \lambda_{di} - \frac{1}{2} \sum_{i \in \mathcal{C}} \sum_{d=1}^{D_i} \sum_{k=1}^{K_d} \frac{c_{kdi}^2 + \mu_{kdi}^2}{\lambda_{di}} \quad (12)$$

where  $\circ$  is used to denote the element-wise product and we made use of the factorization of the joint distribution. A detailed derivation of the ELBO is provided as supplementary material.<sup>1</sup> The key insight is that, through the change of variables, the variational parameters have been transferred inside the loglikelihood, thus enabling stochastic optimization by sampling gradients from it and taking optimization steps in the direction pointed by those noisy gradients. We note that this would not be possible if the expectation in (12) was with respect to the variational distribution  $q(\beta | \mu, bc)$ .

Regarding the variance hyper-parameters  $\lambda$ , it is possible to optimize them analytically. This contrasts with other applications of ARD, where the prior variances are estimated using Expectation-Maximization (EM) - a procedure that can exhibit slow convergence due to the strong dependency between the variational parameters { $\mu$ , c} and the hyper-parameters  $\lambda$  [4]. Taking derivatives of  $\mathcal{L}(\mu, c, \lambda)$  w.r.t.  $\lambda_{di}$  and setting them to zero yields the following optimum:

$$\lambda_{di}^* = \frac{1}{K_d} \sum_{k=1}^{K_d} (c_{kdi}^2 + \mu_{kdi}^2).$$
(13)

Substituting back these optimal values in  $\mathcal{L}(\mu, \mathbf{c}, \lambda)$  gives the optimized evidence lower bound

$$\mathcal{L}(\boldsymbol{\mu}, \mathbf{c}) = \mathbb{E}_{\mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})}[\log p(\mathbf{y}|\mathbf{c} \circ \mathbf{z} + \boldsymbol{\mu})] + \sum_{i \in \mathcal{C}} \sum_{d=1}^{D_i} \sum_{k=1}^{K_d} \log c_{kdi}$$
$$-\frac{1}{2} \sum_{i \in \mathcal{C}} \sum_{d=1}^{D_i} K_d \log \sum_{k=1}^{K_d} (c_{kdi}^2 + \mu_{kdi}^2). \quad (14)$$

In order to fit the variational distribution to the true posterior, we must optimize the lower bound in (14) w.r.t.  $\mu$  and **c** by taking derivatives. See the supplementary material for

<sup>1</sup>Supplementary material available at: https://fprodrigues.com/DCM-ARD.

TABLE I

DESCRIPTION OF THE VARIABLES IN THE SWISSMETRO DATASET

Name	Description
TT	Travel time [min]. Based on the car distance.
CO	Train cost [CHF]. If the traveler owns a GA, equal to its price.
HE	Train headway [min].
ga	Travel card ownership. 1 if the traveler owns one, 0 otherwise.
age	It captures the age class of individuals (6 classes).
luggage	0 if none, 1 if one piece, 3 if several pieces.
who	Who pays for the trip. 1 if self, 2 if employer, 3 if half-half.
purpose	Travel purpose (9 categories).

income | Traveler's income per year (4 categories).

a detailed derivation of these gradients. The lower bound  $\mathcal{L}(\mu, \mathbf{c})$  can then be optimized by first sampling a set of preference parameters  $\boldsymbol{\beta} = \mathbf{c} \circ \mathbf{z} + \mu$ ,  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and using the stochastic gradients above to update the all variational parameters  $\boldsymbol{\mu}$  and  $\mathbf{c}$  in parallel:

$$\boldsymbol{\mu}^{(t)} = \boldsymbol{\mu}^{(t-1)} + \rho_t \nabla_{\boldsymbol{\mu}} \mathcal{L}(\boldsymbol{\mu}, \mathbf{c})$$
(15)

$$\mathbf{c}^{(t)} = \mathbf{c}^{(t-1)} + \rho_t \nabla_{\mathbf{c}} \mathcal{L}(\boldsymbol{\mu}, \mathbf{c})$$
(16)

Following the theory of stochastic optimization [34], using a schedule of the learning rates  $\{\rho_t\}$  such that  $\sum \rho_t = \infty$ ,  $\sum \rho_t^2 < \infty$ , iterating these updates will converge to a local maxima of the bound in (14) or to the global maximum when this bound is concave. At convergence, we can assess the relevancy of each explanatory variable *d* in the utility function for alternative *i* by evaluating the magnitude of the estimated variance parameter  $\lambda_{di}$  using (13).

Lastly, we can further scale-up the variational inference algorithm described above by introducing a second type of stochasticity as proposed by [35]. This second type of stochastic stems from using "mini-batches" of data to compute the stochastic gradients rather then the entire dataset at once, hence resulting in a doubly stochastic variational inference algorithm. The final procedure is summarized in Algorithm 1 in the supplementary material.

#### **IV. EXPERIMENTS**

In this section, an empirical evaluation of the proposed MNL-ARD for automatic utility function specification is performed based on both semi-artificial and real choice data. For both sets of experiments, the dataset used is the Swissmetro (SM) dataset described in [5]. This publicly-available dataset consists of survey data collected on the trains between St. Gallen and Geneva, in which the respondents provided information in order to analyze the impact of the construction of the Swissmetro. respondents were asked to state their favorite transportation mode among three alternatives - train, Swissmetro and car (for car owners only) - in nine different hypothetical situations. After discarding respondents for which some variables were not available (e.g. age, purpose), a total of 10692 responses from 1188 individuals were used for the experiments. Table I gives a brief description of the subset of variables used in the context of this study; a more detailed and complete description of the dataset and its collection procedure may be found in [5].

The proposed MNL-ARD model and its corresponding doubly-stochastic variational inference (DSVI) algorithm were

implemented in Matlab. Source code for the implementation and for reproducing all experiments in this paper is publicly available at: https://github.com/fmpr/DCM-ARD.

# A. Semi-Artificial Choice Data

In order to empirically demonstrate the ability of the proposed approach to discover the correct utility function specifications, we began by conducting an extensive series of experiments on semi-artificial choice data based on the Swissmetro dataset. We manually specified a set of "artificial" (but realistic) utility function specifications of varying complexity based on the input and suggestions from domain experts. Then, based on the Swissmetro dataset, we sampled new artificial choices for the respondents according to the manually-specified utility functions. This was done by fitting a standard DCM with the manually-specified utility function to the original data using maximum likelihood estimation and, based on the learned parameters  $\beta^*$ , we then sampled new choices  $y_n \sim \text{Categorical}(y_n|P_n)$ .

We consider two experimental settings for the application of MNL-ARD:

- an experimental setting with a medium-sized utility function search space, in which the number of possible variables to be included in the utility functions is 252; these include the original variables (e.g. alternativespecific constants "ASC", travel-time "TT", cost "CO" and headway "HE"), their log-transformations, and interactions of both the original variables and their logarithms with trip purpose ("pur", 9 categories), respondent age ("age", 5 groups) and annual season ticket availability ("ga", binary). Kindly note that, although this results in 252 variables that can be included in the specification, the dimensionality of the utility function search-space includes all combinations of possible utility functions that can be generated using these variables and therefore grows exponentially with this number. For example, considering just the subset of all utility functions with only 10 variables results in  $\binom{252}{10} = 2.4 \times 10^{17}$  possible utility functions to be considered;
- an experimental setting with a large utility function search space; besides the variables in the medium-sized search space, this search space also considers Box-Cox transformations, variable segmentations based on K-means clustering, and interactions of the original variables with respondent income ("inc", 5 groups), luggage ("lug": none, one piece or multiple pieces) and who pays for the trip ("who": unknown, self, employer or half-half). This results in a total of 602 possible variables to be included in the utility function specifications.

Based on these two search spaces, we manually defined 9 artificial utility function specifications as shown in Table II. Specifications S1-S6 are based on the medium-sized search space, while specifications S7-S9 are based on the large search space. However, in order to verify that MNL-ARD is able to discover the true utility function specification used to generate the choice data regardless of how large the search space RODRIGUES et al.: BAYESIAN ARD FOR UTILITY FUNCTION SPECIFICATION IN DCMs

TABLE II MANUALLY-DEFINED UTILITY FUNCTION SPECIFICATIONS USED TO GENERATE THE SEMI-ARTIFICIAL CHOICE DATA

	Artificial specification						
S#	Variables in V <sub>train</sub>	Variables in V <sub>sm</sub>	Variables in V <sub>car</sub>				
<b>S</b> 1	ASC, TT, CO	ASC, TT, CO	TT, CO				
<u>S2</u>	ASC, TT, TT x age, CO	ASC, TT, CO,	TT, TT x age,				
02		CO x ga	CO				
\$3	ASC, TT, TT x age,	ASC, TT, CO,	TT, TT x age,				
	CO, CO x ga, HE	CO x ga, logHE	СО				
<b>S</b> 4	ASC, ASC x ga, TT, CO	ASC, ASC x ga,	TT, <u>CO</u> ,				
34		TT, CO	CO x purpose				
S5	ASC, logTT, HE	ASC, logTT, <u>HE</u>	TT, CO				
\$6	ASC, logTT,	ASC, logTT	TT, CO				
- 50	logTT x ga, CO						
67	ASC, <u>boxTT</u> ,	ASC, TT	TT, CO				
3/	boxTT x ga, CO						
68	ASC, ASC x ga, TT,	ASC, ASC x ga, TT,	TT, CO,				
30	CO, CO x who	CO, CO x who	CO x luggage				
50	ASC, TT, CO	ASC, TT, TT x age,	TT, CO,				
39	CO x ga	CO, CO x ga	CO x income				

considered is, we also test specifications S1-S3 with the large search space.<sup>2</sup>

Given the semi-artificial choice data generated based on the manually-defined utility function specifications from Table II, our goal is to test the ability MNL-ARD to recover the correct utility function specifications in a purely data-driven way. Table III shows the top-K variables selected by MNL-ARD for the medium-sized search space (i.e. specifications S1-S6) ranked according to their respective learned  $\lambda$  values. In order to simplify the analysis of the results, the variables deemed relevant by MNL-ARD are highlighted in bold. Irrelevant variables are expected to have  $\lambda \approx 0$ . As these results demonstrate, the proposed MNL-ARD is able to discover the true specifications almost perfectly, with all the truly "irrelevant" variables being assigned a  $\lambda$  value of approximately zero. The only minor exceptions can be found in specifications S4 and S5. In the learned utility function for S4, we can observe that cost ("CO") is assigned a  $\lambda$  value of zero for the utility of car despite the fact that it was part of the true specification that was used to generate the semi-artificial data. We believe this to be a consequence of the inclusion of the interaction between "CO" and purpose ("pur") in the true specification for car. Since there is a total of 9 different purposes and some of them have an extremely low number of observations, the effect of "CO" alone can be captured by the baseline and therefore its presence in the specification is essentially not required from a pure data perspective. As for S5, the headway variable ("HE") in the SM utility was assigned a rather low value of  $\lambda$  ( $\lambda = 0.001$ ), despite the fact that it should be clearly identified by MNL-ARD as a relevant variable, since it was part of the true specification of S5.<sup>3</sup>

Let us now consider the large search space. The top of Table IV shows the top-K variables with higher  $\lambda$  value

#### TABLE III

RESULTS FOR MEDIUM-SIZED SEARCH SPACE. \* IS USED TO INDICATE VARIABLES THAT ARE PRESENT IN THE CORRESPONDING TRUE UTILITY FUNCTION SPECIFICATIONS DEFINED IN TABLE II

	Train		Swiss Metro		Car	
S#	Variable	λ	Variable	λ	Variable	λ
	ASC*	1.814	TT*	0.513	TT*	0.744
	TT*	1.174	ASC*	0.126	CO*	0.011
<b>S</b> 1	CO*	0.393	CO*	0.066	logTT x pur1	0.000
	CO x age1	0.000	logHE x age1	0.000	logTT x pur2	0.000
	ASC*	2.353	TT*	0.495	TT*	0.389
	TT x age1*	0.524	CO x ga*	0.195	CO*	0.070
	TT x age2*	0.524	ASC*	0.120	TT x age1*	0.060
	TT x age3*	0.524	CO*	0.030	TT x age2*	0.060
S2	TT x age4*	0.524	ASC x purl	0.000	TT x age3*	0.060
~-	TT*	0.468	ASC x pur2	0.000	TT x age4*	0.060
	CO*	0.416	ASC x pur3	0.000	logTT x purl	0.000
	TT x purl	0.000	ASC x pur4	0.000	logTT x pur2	0.000
	11 A put	0.000	noe a part	0.000	logi i n puiz	0.000
	ASC*	2.536		0.522	 TT*	0.478
	CO x ga*	0.633	CO x ga*	0.426	co*	0.120
	TT x age1*	0.510	ASC*	0.133	TT x age1*	0.061
\$3	TT x age2*	0.510	CO*	0.023	TT x age2*	0.061
	TT x age3*	0.510	logHE*	0.005	TT x age3*	0.061
	TT x age4*	0.510	HE x age1	0.000	TT x age4*	0.061
	TT*	0.300	HE x age?	0.000	logTT x ga	0.000
	co*	0.202	HE x age3	0.000	$\log CO \times pur1$	0.000
	HE*	0.056	HE x aged	0.000	$\log CO \times pur^2$	0.000
	HE x pur1	0.000	$\log CO \times purl$	0.000	logCO x pur3	0.000
		0.000	logeo x puri	0.000	logeo x puis	0.000
	ASC x ga*	6 8 3 6	ASC x ga*	3 401		0.855
		2 323	CO*	1 354	CO x pur1*	0.100
	ASC*	1.338	TT*	0.462	$CO \times pur2*$	0.100
	TT*	0.885	ASC*	0.361	CO x pur <sup>2</sup> *	0.100
	CO x ga	0.001	logHE x age1	0.001	CO x pur4*	0.100
<b>S</b> 4	ASC x pur1	0.000	logHE x age?	0.001	CO x pur5*	0.100
5.	ASC x pur2	0.000	logHE x age3	0.001	CO x pur6*	0.100
	ASC x pur3	0.000	logHE x age4	0.001	CO x pur7*	0.100
	ASC x pur4	0.000	$CO \times pur1$	0.000	CO x pur8*	0.100
	ASC x pur5	0.000	$CO \times pur^2$	0.000	logTT x ga	0.000
	, in the second party of the second s	0.000		0.000		0.000
	ASC*	1.775	logTT*	0.557		0.722
	logTT*	1.405	ASC*	0.087	CO*	0.042
	HE*	0.035	CO	0.002	logTT x purl	0.000
S5	TT x age1	0.000	HE	0.001	logTT x pur?	0.000
	TT x age?	0.000	HE x age1	0.000	logTT x pur3	0.000
		0.000	THE A uger	0.000	log11 x puis	0.000
	ASC*	2.071	 logTT*	0.664	 TT*	0.809
	InoTT v 09*	1 600	ASC*	0 106	CO*	0.042
	logTT*	0.611	logTT v age1	0.000	logTT x purl	0.000
S6	CO*	0.30/	logTT x age?	0.000	logTT x pur?	0.000
	TT x age1	0.000	logTT x age2	0.000	logTT x pur2	0.000
		0.000	Ing II A ages	0.000	log11 x puis	0.000

according to MNL-ARD for S1, S2 and S3. As the obtained results show, MNL-ARD is still able to recover the true specifications that were used to generate the data regardless of the significantly larger search space (602 variables considered, instead of 252 for Table III). However, since the number of variables considered is substantially larger, the execution time of the proposed DSVI algorithm naturally increased from approximately 10 minutes to close to 1 hour on a standard 2.3 GHz dual-core laptop with 16 GB of RAM.

The remaining of Table IV shows the top-K variables deemed relevant by MNL-ARD for inclusion in the utility function specifications for S7, S8 and S9. By comparing these results with the true specifications from Table II, one can again observe that MNL-ARD is able to discover the true specifications almost exactly. The only differences are the fact

<sup>&</sup>lt;sup>2</sup>We further tested other specifications, but omitted their results for conciseness (they lead to similar conclusions). However, they are available at: https://github.com/fmpr/DCM-ARD.

<sup>&</sup>lt;sup>3</sup>Due to space constraints, an analysis of the convergence of the derived DSVI algorithm and the sparsity induced by the hierarchical prior that MNL-ARD uses is provided as supplementary material.

#### IEEE TRANSACTIONS ON INTELLIGENT TRANSPORTATION SYSTEMS

TABLE IV RESULTS OF MNL-ARD FOR LARGE SEARCH SPACE. \* IS USED TO INDICATE VARIABLES THAT ARE PRESENT IN THE CORRESPONDING TRUE UTILITY FUNCTION SPECIFICATIONS DEFINED IN TABLE II

	Train		Swiss Met	tro	Car	
S#	Variable	λ	Variable	λ	Variable	λ
	ASC*	1.879	TT*	0.537	TT*	0.677
	TT*	1.196	ASC*	0.137	CO*	0.011
<b>S</b> 1	CO*	0.513	CO*	0.096	CO x pur1	0.000
	logCO x inc1	0.001	TT x lugg1	0.000	CO x pur2	0.000
	ASC*	2.391	TT*	0.568	TT*	0.411
	TT*	0.606	CO x ga*	0.179	TT x age1*	0.059
	CO*	0.477	ASC*	0.130	TT x age2*	0.059
	TT x age1*	0.352	CO*	0.031	TT x age3*	0.059
S2	TT x age2*	0.352	HE x inc1	0.000	TT x age4*	0.059
	TT x age3*	0.352	HE x inc2	0.000	CO*	0.036
	TT x age4*	0.352	HE x inc3	0.000	TT x lugg1	0.000
	logCO x pur1	0.001	HE x inc4	0.000	TT x lugg2	0.000
	ASC*	2.599	TT*	0.594	TT*	0.446
	CO x ga*	0.717	CO x ga*	0.477	CO*	0.107
	TT*	0.432	ASC*	0.127	TT x age1*	0.074
<b>S</b> 3	TT x age1*	0.364	CO*	0.017	TT x age2*	0.074
	TT x age2*	0.364	logHE*	0.003	TT x age3*	0.074
	TT x age3*	0.364	HE x incl	0.000	TT x age4*	0.074
	TT x age4*	0.364	HE x inc2	0.000	CO x who1	0.001
	CO*	0.194	HE x inc3	0.000	CO x who2	0.001
	HE*	0.057	HE x inc4	0.000	CO x who3	0.001
	ASC x who1	0.002	logHE x incl	0.000	TT x purl	0.000
		2.246	 TPTP *	0.574	···	0.552
	hovTT v go*	2.240		0.574		0.555
	CO*	0.360	CO v purl	0.120	TT v lugg1	0.019
<b>S</b> 7		0.300	$CO \times pur^2$	0.000	$TT \times \log gT$	0.000
	$\log CO \times inc1$	0.001	$CO \times pur2$	0.000	seg(CO 4)	0.000
	logeo x mer	0.001	CO x puis	0.000	305(00,4)	0.000
	ASC x ga*	7.448	ASC x ga*	4.805	 TT*	0.828
	CO*	2.840	CO*	1.695	CO*	0.018
	ASC*	1.611	TT*	0.559	seg(CO.4)	0.000
	TT*	1.120	ASC*	0.336	seg(CO.4)	0.000
<b>S</b> 8	CO x who1*	0.057	CO x who1*	0.025	seg(CO,4)	0.000
	CO x who2*	0.057	CO x who2*	0.025	CO x purl	0.000
	CO x who3*	0.057	CO x who3*	0.025	CO x pur2	0.000
	CO x inc1	0.001	seg(TT,8)	0.000	CO x pur3	0.000
	ASC*	2.255	TT*	1.367	TT*	1.118
	TT*	1.197	CO x ga*	0.501	CO*	0.098
	CO x ga*	0.805	TT x age1*	0.134	CO x inc1*	0.004
	CO*	0.187	TT x age2*	0.134	CO x inc2*	0.004
92	ASC x who1	0.001	TT x age3*	0.134	CO x inc3*	0.004
57	ASC x who2	0.001	TT x age4*	0.134	CO x inc4*	0.004
	ASC x who3	0.001	ASC*	0.110	CO x who1	0.001
	HE x age1	0.000	CO*	0.015	CO x who2	0.001
	HE x age2	0.000	seg(TT,8)	0.000	CO x who3	0.001

that MNL-ARD selected "logTT" instead of "boxTT" in the utility function of train in S7, and the fact that it missed the interaction between "CO" and "luggage" in the utility function of car in S8. While we could not find an obvious explanation for the latter, the former can be easily explained by an analysis of the results of the Box-Cox transform, which uses a maximum likelihood approach to fit the parameters of the transformation. In the particular case of train travel time, we could immediately observe that the transformed values produced by the Box-Cox transformation are almost perfectly correlated with to the ones produced by the log-transformation (correlation coefficient of 0.998), thus leading us to conclude

TABLE V PREDICTION ACCURACY AND LOG-LIKELIHOOD ON HELD-OUT DATA

		MNL		MNL-ARD		MNL-TRUE	
Search Space	Spec	Acc.	LogLik	Acc.	LogLik	Acc.	LogLik
Medium	S1	0.615	-2733.9	0.628	-2569.0	0.627	-2567.4
Medium	S2	0.627	-2697.0	0.638	-2498.2	0.636	-2496.8
Medium	S3	0.639	-2662.5	0.645	-2452.9	0.646	-2450.4
Medium	S4	0.627	-2597.3	0.647	-2454.9	0.648	-2452.7
Medium	S5	0.607	-2788.8	0.623	-2621.9	0.623	-2619.2
Medium	S6	0.624	-2621.3	0.632	-2530.3	0.633	-2527.1
Large	<b>S</b> 1	0.589	-2798.2	0.628	-2569.0	0.627	-2567.4
Large	S2	0.602	-2773.7	0.638	-2498.2	0.636	-2496.8
Large	S3	0.612	-2924.0	0.645	-2452.9	0.646	-2450.4
Large	S7	0.603	-2746.6	0.606	-2675.5	0.617	-2551.1
Large	S8	0.598	-2858.0	0.642	-2489.8	0.646	-2421.7
Large	S9	0.614	-2823.9	0.653	-2466.7	0.660	-2400.3

that both lead to equivalent utility function specifications for the train alternative.

As a further test of scalability and robustness of the proposed approach, we also considered an extremely large search space, which was obtained by expanding the large space space described above with variables that consist of Gaussian random noise, until a total of 1000 variables per alternative was reached (total of 3000 variables). Using the semi-artificial choice data corresponding to specification S2 we were able to verify that, despite the expected increased computational run time (approximately 5 hours), the proposed MNL-ARD was still able to perfectly recover the true specification of S2.

So far we have only been considering the ability of MNL-ARD to infer the correct utility function specifications. However, one can also evaluate MNL-ARD in terms of its prediction accuracy on held-out data. Table V shows the prediction accuracy of MNL-ARD when trained only on 70% of the dataset and tested on 30% held-out data for the different semi-artificial specifications. By comparing these results with the accuracy of a standard MNL model that considers all the variables from the search space as input ("MNL"), one can verify that thanks to the additional flexibility of the proposed hierarchical prior and the sparsity-inducing properties, MNL-ARD is able to generalize better to held-out data, thus resulting in significantly higher prediction accuracies. In fact, is most cases, MNL-ARD achieves almost as good prediction performance as a MNL model estimated using the true specifications that were used to generate the semi-artificial choices ("MNL-TRUE"). On the other hand, a MNL fitted with maximum likelihood estimation with such a high number of input variables is very likely to severely overfit.

# B. Real Choice Data

We will now consider the application of MNL-ARD to perform automatic utility function specification on the real choice data from the Swissmetro dataset. Table VI shows the top-20 variables selected by MNL-ARD for inclusion in the utility functions using the medium-sized search space. Since in this case the correct specification is unknown, we instead evaluate the quality of the MNL models that the specifications inferred by MNL-ARD produce. With this purpose, we developed a series of specifications of increasing complexity based on the results of Table VI. We begin by considering a rather

TABLE VI Results of MNL-ARD for Real SM Data

Train		Swiss Met	ro	Car		
Variable	$\lambda$	Variable	$\lambda$	Variable	$\lambda$	
logTT x ga	9.506	logCO x ga	5.570	logCO	4.479	
ASC	4.002	logCO x pur1	2.251	TT x ga	1.378	
logCO	3.262	logCO x pur2	2.251	logTT x pur1	0.477	
logCO x purl	2.469	logCO	1.184	logTT x pur2	0.477	
logCO x pur2	2.469	logTT	0.506	logCO x age1	0.213	
logCO x ga	1.235	CO	0.349	logCO x age2	0.213	
CO	0.556	ASC x age1	0.250	logCO x age3	0.213	
logCO x age1	0.269	ASC x age2	0.250	logCO x age4	0.213	
logCO x age2	0.269	ASC x age3	0.250	CO x pur1	0.156	
logCO x age3	0.269	ASC x age4	0.250	CO x pur2	0.156	
logCO x age4	0.269	CO x purl	0.236	TT x age1	0.107	
CO x purl	0.228	CO x pur2	0.236	TT x age2	0.107	
CO x pur2	0.228	ASC x ga	0.146	TT x age3	0.107	
logTT	0.175	CO x ga	0.099	TT x age4	0.107	
logHE	0.075	TT x age1	0.027	CO	0.037	
CO x ga	0.068	TT x age2	0.027	logTT	0.000	
CO x age1	0.034	TT x age3	0.027	logTT x ga	0.000	
CO x age2	0.034	TT x age4	0.027	logCO x ga	0.000	
CO x age3	0.034	TT x pur1	0.005	TT	0.000	
CO x age4	0.034	TT x pur2	0.005	TT x pur1	0.000	

#### TABLE VII

UTILITY FUNCTION SPECIFICATIONS FOR TRUE SM DATA

S#Variables in $V_{train}$ Variables in $V_{sm}$ Attrib. in $V_{car}$ R1ASC, TT, COASC, TT, COTT, COR2ASC, logTT,ASC, logTT,TT, logCOlogTT x ga, logCOlogCOlogCOR3ASC, logTT, logTT x ga,ASC, logTT, logCO,It, logCOASC, logTT, logTO x purlogCO x galogCO x gaASC, logTT, logTT x ga,ASC, logTT, logTO x ga,logCO, logCO x ga,logCO, logCO x purlogCO x purlogCO x ga,logCO x purlogCO x purlogCO x purASC, logTT, logTT x ga,ASC, ASC x age,TT, TT x ga,R4logCO, logCO x ga,logCO x purlogCO x purlogCO x purlogCO x purlogCO x purASC, logTT, logTT x ga,ASC, ASC x age,TT, TT x pur,logCO x agelogCO x purlogCO x pur	
R1ASC, TT, COASC, TT, COTT, COR2ASC, logTT, logTT x ga, logCOASC, logTT, logCOTT, logCOR3ASC, logTT, logTT x ga, logCO, logCO x purASC, logTT, logCO, logCO x gaTT, logCOR4logCO, logCO x ga, logCO x purASC, logTT, logTT x ga, logCO, logCO x ga, logCO x purASC, logTT, tTT, TT x ga, logCO x ga, logCO x purR4logCO, logCO x ga, logCO x purIogCO x ga, logCO x purTT, TT x ga, rtt x tr x ga, logCO	
R2ASC, logTT, logTT x ga, logCOASC, logTT, logCOTT, logCOR3ASC, logTT, logTT x ga, logCO, logCO x purASC, logTT, logCO, logCO x gaTT, logCO, TT, logCOR4logCO, logCO x ga, logCO x purlogCO, logCO x ga, logCO x ga, logCO x purTT, TT x ga, logCO x purTT, TT x ga, rt, TT x ga, logCO	
N2     logTT x ga, logCO     logCO       R3     ASC, logTT, logTT x ga, logCO, logCO x pur     ASC, logTT, logCO, logCO x ga     TT, logCO       ASC, logTT, logTT x ga, logCO, logCO x ga, logCO x pur     ASC, logTT, logTT x ga, logCO x ga, logCO x pur     TT, TT x ga, logCO x ga, logCO x ga, logCO x ga, logCO, logCO x ga, logCO, logCO x ga, logCO, logCO x ga, logCO x ga	
R3ASC, logTT, logTT x ga, logCO, logCO x purASC, logTT, logCO, logCO x gaTT, logCOASC, logTT, logTT x ga, logCO, logCO x ga,ASC, logTT, logCO, logCO x ga, logCO x purTT, TT x ga, logCOR4logCO, logCO x ga, logCO x purlogCO, logCO x ga, logCO x ga, logCO x ga, logCO x ga, logCO, logCO x ga, logCO x purTT, TT x ga, logCO tr, TT x pur, logCOR5logCO x age logCO x age logCO x ga, logCO x purlogCO x pur	
NS     logCO, logCO x pur     logCO x ga       ASC, logTT, logTT x ga,     ASC, logTT,     TT, TT x ga,       R4     logCO, logCO x ga,     logCO, logCO x ga,     logCO       logCO x pur     logCO x pur     logCO x ga,     logCO       ASC, logTT, logTT x ga,     ASC, ASC x age,     TT, TT x ga,       logCO, logCO x ga,     logTT, logTT, logTT, logTT, logCO,     TT, TT x pur,       logCO, logCO x ga,     logTT, logCO,     TT x pur,       logCO x age     logCO x pur     logCO	
ASC, logTT, logTT x ga, logCO, logCO x ga, logCO x purASC, logTT, logCO, logCO x ga, logCO x ga, logCO x ga, logCO x ga, logCO x ga, logCO x ga, logCO x ga, logCO, logCO x ga, logCO, logCO x ga, logCO, logCO x ga, logCO, logCO x ga, logCO x ga, logC	
R4     logCO, logCO x ga, logCO x pur     logCO, logCO x ga, logCO x pur     logCO       ASC, logTT, logTT x ga, logCO, logCO x ga, logCO x pur, logCO x ga,     ASC, ASC x age, logTT, logCO, logCO x ga, logCO x ga,     TT, TT x ga, logCO, logCO x pur       No     Index (ASC)     Index (ASC)	
logCO x pur         logCO x pur           ASC, logTT, logTT x ga, logCO, logCO x ga, logCO x pur, logCO x pur, logCO x ga         ASC, ASC x age, logTT, logCO, logCO x ga, logCO x ga, logCO x pur         TT, TT x ga, logCO, logCO x ga, logCO x pur           ASC, logTT, logTT x ga         ASC, ASC x age, logCO x pur         TT, TT x ga, logCO	
ASC, logTT, logTT x ga, logCO, logCO x ga, logCO x pur,     ASC, ASC x age, logTT, logCO, logCO x ga,     TT, TT x ga, logTT, logCO, logCO x ga,       IogCO x ga     logCO x ga, logCO x pur     logCO x ga, logCO x pur	
R5     logCO, logCO x ga, logCO x pur,     logTT, logCO, logCO x ga, logCO x ga,     TT x pur, logCO       IogCO x age     logCO x pur	
IogCO x pur,     logCO x ga,     logCO       logCO x age     logCO x pur	
logCO x age     logCO x pur     ASC logTT logTT x ga     ASC ASC x age     TT TT x ga	
ASC logTT logTT v ga ASC ASC v age TT TT v ga	
<i>Abc, log11, log11 x ga, Abc, Abc x agc, 11, 11 x ga,</i>	
$\mathbf{R}_{6}$ logCO, logCO x ga, logTT, logCO, TT x pur,	
$\log CO x pur,$ $\log CO x ga,$ $\log CO$	
logCO x age, <b>logHE</b> logCO x pur	
ASC, logTT, logTT x ga, ASC, ASC x ga, TT, TT x ga,	
$\mathbf{P7}$ logCO, logCO x ga, ASC x age, logTT, TT x pur, log	CO
$  \log CO x pur,   \log CO x ga,   \log CO x age,$	
logCO x age, logHE logCO x pur logCO x pur	

simplistic specification based only on travel time and cost (R1). We then start adding variables to it according to the results of MNL-ARD in descending order of importance according to the learned values of  $\lambda$ . The complete set of specifications considered is show in Table VII. Kindly note that the last specification (R7), already includes almost all the variables in the top-20 ranking shown in Table VI, and that other additional variables were assigned a  $\lambda$  value of zero (or very close to zero), thus being deemed irrelevant by MNL-ARD. Also, since including both a variable and its log-transform could compromise the interpretability of the MNL models, we decided to include only the version with the higher value of  $\lambda$  in the cases where MNL-ARD selected both variants.<sup>4</sup> Also, due to the fact that the purpose variable

<sup>4</sup>We note that, according to empirical evidence, including both variants does tend to lead to models that fit better the data, including the held-out data.

TABLE VIII Results for Different Specifications on SM Data

	Specification						
	R1	R2	R3	R4	R5	R6	<b>R</b> 7
Loglik full	-8,625	-8,368	-8,064	-7,836	-7,679	-7,645	-7,617
AIC	17,267	16,755	16,152	15,704	15,410	15,345	15,301
BIC	17,326	16,821	16,239	15,820	15,599	15,542	15,549
Pseudo- $R^2$	0.221	0.244	0.272	0.292	0.306	0.309	0.312
Pseudo- $\bar{R}^2$	0.220	0.243	0.271	0.291	0.304	0.307	0.309
Loglik train	-6,032	-5,822	-5,619	-5,429	-5,297	-5,271	-5,247
Loglik test	-2,603	-2,558	-2,457	-2,437	-2,428	-2,421	2,430
Train acc.	0.616	0.636	0.661	0.676	0.689	0.690	0.692
Test acc.	0.615	0.638	0.662	0.670	0.675	0.677	0.679

has 9 categories, with some of them having only a couple of observations, we further grouped the trip purposes into: commuting, shopping and leisure.

Based on the specifications that were generated according to the results of MNL-ARD (Table VI), we then fitted standard MNL models using the PyLogit package [36] in Python. Table VIII shows the results obtained for the different specifications considered. As expected, one can verify that, as we increase the complexity of the specification according to the results of MNL-ARD, the fit of the MNL model improves in terms of log-likelihood. However, the quality of the MNL model also improves in terms of AIC, BIC and pseudo- $\bar{R}^2$ . In order to further assess the quality of the MNL-ARD specifications in terms of generalization ability to held-out data, we also performed a random 70/30% train/test split of the dataset, and computed the likelihood and accuracies in both sets. As the results in Table VIII evidence, as we move towards the full specification inferred by MNL-ARD, the accuracy and held-out data likelihood of the MNL model also improves. Interestingly, it can observed that only when we include essentially all the variables deemed relevant by MNL-ARD we start noticing some signs of overfitting: BIC and testset likelihood do not improve when going from specification R6 to R7. However, indicators such as AIC and pseudo- $\overline{R}^2$ still improve. Furthermore, it should be noted that the variables included from R6 to R7, already consist of variables for which MNL-ARD assigned a relatively low relevance (i.e. low value of  $\lambda$  when compared to the others).

Comparing the results of specifications R6 and R7 with other proposed MNL specifications from the literature for the same dataset (Table IX), it is possible to obtain a better perspective of how good the specifications inferred by MNL-ARD are. For example, the MNL specification proposed in PyLogit for the Swissmetro dataset includes variables such as travel time, cost, headway, seat configuration, luggage and first class. However, it only achieves a loglikelihood of -8,061, a BIC of 16,252 and a pseudo- $\bar{R}^2$  of 0.271. Similarly, the original specification proposed by [5] achieves a loglikelihood of just -8,483, a BIC of 17,050 and a pseudo- $\overline{R}^2$  of 0.233. Moreover, if we consider generalization to held-out data, Table IX also demonstrates that the both R6 and R7 obtain better results than both baseline approaches, thereby highlighting how MNL-ARD can be easily used to automate the search of utility function specifications.

Lastly, Table X shows the estimated coefficients by a MNL model with the specification R6 using PyLogit, and their

	TABLE IX
RESULTS FOR TRUE SM DATA	VS. BASELINE FROM STATE OF THE ART

		Specificatio	n	
	[5]	PyLogit Example	R6	R7
Log-lik full	-8,483	-8,061	-7,645	-7,617
AIC	16,984	16,150	15,345	15,301
BIC	17,050	16,252	15,542	15,549
Pseudo- $R^2$	0.234	0.272	0.309	0.312
Pseudo- $\overline{R}^2$	0.233	0.271	0.307	0.309
Log-lik train	-5,960	-5,633	-5,271	-5,247
Log-lik test	-2,535	-2,450	-2,421	2,430
Train acc.	0.646	0.667	0.690	0.692
Test acc.	0.644	0.650	0.677	0.679

# TABLE X Results for True SM Data, Specification S6

		0.10			10.005	0.0751
	Coef	StdErr	<i>z</i>	p >  z	[0.025	0.975]
ASC (Train)	3.036	0.196	15.478	0.000	2.652	3.421
ASC (SM)	0.900	0.134	6.725	0.000	0.638	1.163
ASC x age1 (SM)	0.575	0.156	3.699	0.000	0.271	0.880
ASC x age2 (SM)	0.784	0.103	7.585	0.000	0.582	0.987
ASC x age3 (SM)	0.704	0.102	6.909	0.000	0.505	0.904
ASC x age4 (SM)	0.479	0.107	4.478	0.000	0.270	0.689
logTT (Train)	-0.964	0.261	-3.697	0.000	-1.477	-0.453
logTT (SM)	-2.570	0.110	-23.465	0.000	-2.785	-2.355
TT (Car)	-0.865	0.218	-3.974	0.000	-1.293	-0.439
logTT x ga (Train)	-2.995	0.275	-10.880	0.000	-3.535	-2.455
TT x ga (Car)	-0.176	0.210	-0.841	0.400	-0.589	0.235
TT x pur1 (Car)	0.273	0.064	4.285	0.000	0.148	0.398
TT x pur2 (Car)	0.289	0.088	3.289	0.001	0.117	0.463
logCO (Train)	-2.637	0.318	-8.297	0.000	-3.261	-2.015
logCO (SM)	-1.984	0.247	-8.023	0.000	-2.470	-1.500
logCO (Car)	-1.875	0.175	-10.714	0.000	-2.218	-1.532
logCO x ga (Train)	-1.997	0.195	-10.248	0.000	-2.379	-1.615
logCO x ga (SM)	-2.249	0.132	-17.024	0.000	-2.509	-1.991
CO x age1 (Train)	-0.317	0.090	-3.539	0.000	-0.493	-0.141
CO x age2 (Train)	-0.578	0.079	-7.336	0.000	-0.733	-0.424
CO x age3 (Train)	-0.647	0.080	-8.134	0.000	-0.804	-0.492
CO x age4 (Train)	-0.525	0.083	-6.301	0.000	-0.690	-0.362
logCO x pur1 (Train)	2.521	0.294	8.574	0.000	1.945	3.098
logCO x pur1 (SM)	1.963	0.231	8.510	0.000	1.511	2.415
logCO x pur2 (Train)	3.282	0.308	10.641	0.000	2.678	3.887
logCO x pur2 (SM)	2.589	0.244	10.606	0.000	2.111	3.068
HE (Train)	-0.948	0.118	-8.059	0.000	-1.179	-0.718

corresponding p-values and other statistics. The full set of results for the other specifications were omitted for brevity but are available at https://github.com/fmpr/DCM-ARD, together with the source code. As the results in Table X demonstrate, the specification learned by MNL-ARD leads to a stable MNL model in which the coefficients for all variables except "TT x ga (Car)", have p-values smaller than 0.001. It should however be noted that, in two cases, the parameter estimates are not entirely behaviourally realistic: for both Train and SM alternatives, the sum of the parameter related to "logCO x pur2" and the corresponding baseline ("logCO (Train)" and "logCO (SM)") is positive, implying that all else being equal, increasing the travel cost of shopping trips improves their attractiveness. Such result is obviously wrong; it indicates that the involved parameters are erroneously capturing or omitting some effects, most probably because the travel cost of the two affected modes is interacted with "ga" and "pur", but not with both simultaneously. However, since such interactions were not considered in the search-space, MNL-ARD is unable to identify them as relevant. Thus, this is a great example that highlights an important limitation of MNL-ARD: its results are dependent of the search-space considered, and it has no

knowledge of behavioural theories. However, we reiterate that its purpose is to assist modellers on specifying utility functions according to data-driven knowledge, rather then serving as a replacement for domain knowledge and modellers' expertise.

# V. CONCLUSION

This paper proposed a Bayesian framework for performing automatically utility function specification in discrete choice models based on the idea of automatic relevance determination (ARD). An efficient doubly stochastic variational inference algorithm was derived in order to perform approximate Bayesian inference in the proposed MNL-ARD model. As our empirical results using both semi-artificial and real choice data showed, the proposed approach is able to automatically discover good utility function specifications in a pure data-driven manner. Even in situations where the number of possible variables considered for inclusion in the utility functions is very large, our proposed approach was shown to be able to recover the "true" utility function specification almost perfectly. The practical advantages and overall realworld feasibility of the proposed approach were demonstrated through an application to the popular Swissmetro dataset [5], where MNL-ARD was shown to be capable of generating specifications that outperform others from the state of the art according to multiple criteria.

Despite the importance of the standard formulation of the multinomial logit in discrete choice theory, it only corresponds to a subset of the models that are used in practice, with modelling approaches like mixed logits and latent class choice models providing important ways of capturing the heterogeneity in preferences among the decision makers. Therefore, our future work focuses on extending the proposed MNL-ARD formulation for these models. Concretely, the idea consists in placing flexible hierarchical priors on the class-specific logit parameters in latent class choice models, and on the population-level taste parameters in the case of mixed logit models. Furthermore, given that a critical factor for the practical applicability of the proposed approach is computational efficiency, future work will also look at recent advances in approximate Bayesian inference, such as amortized variational inference and normalizing flows [37] as proposed in [38], in order to further improve the scalability of the proposed approach for automatically utility function specification.

#### REFERENCES

- C. Torres, N. Hanley, and A. Riera, "How wrong can you be? Implications of incorrect utility function specification for welfare measurement in choice experiments," *J. Environ. Econ. Manage.*, vol. 62, no. 1, pp. 111–121, Jul. 2011.
- [2] M. E. Tipping, "Sparse Bayesian learning and the relevance vector machine," J. Mach. Learn. Res., vol. 1, pp. 211–244, Sep. 2001.
- [3] C. E. Rasmussen, "Gaussian processes in machine learning," in *Summer School on Machine Learning*. Berlin, Germany: Springer, Feb. 2003, pp. 63–71.
- [4] M. Titsias and M. Lázaro-Gredilla, "Doubly stochastic variational Bayes for non-conjugate inference," in *Proc. Int. Conf. Mach. Learn.*, 2014, pp. 1971–1979.
- [5] M. Bierlaire, K. Axhausen, and G. Abay, "The acceptance of modal innovation: The case of swissmetro," in *Proc. 1st Swiss Transp. Res. Conf.*, 2001, pp. 1–17.

RODRIGUES et al.: BAYESIAN ARD FOR UTILITY FUNCTION SPECIFICATION IN DCMs

- [6] M. Dash and H. Liu, "Feature selection for classification," *Intell. Data Anal.*, vol. 1, nos. 1–4, pp. 131–156, 1997.
- [7] I. Guyon and A. Elisseeff, "An introduction to variable and feature selection," J. Mach. Learn. Res., vol. 3, pp. 1157–1182, Jan. 2003.
- [8] S.-W. Lin, Z.-J. Lee, S.-C. Chen, and T.-Y. Tseng, "Parameter determination of support vector machine and feature selection using simulated annealing approach," *Appl. Soft Comput.*, vol. 8, no. 4, pp. 1505–1512, Sep. 2008.
- [9] D. Fouskakis and D. Draper, "Comparing stochastic optimization methods for variable selection in binary outcome prediction, with application to health policy," *J. Amer. Stat. Assoc.*, vol. 103, no. 484, pp. 1367–1381, Dec. 2008.
- [10] N. R. Pal, S. Nandi, and M. K. Kundu, "Self-crossover—A new genetic operator and its application to feature selection," *Int. J. Syst. Sci.*, vol. 29, no. 2, pp. 207–212, Feb. 1998.
- [11] E. P. Xing, M. I. Jordan, and R. M. Karp, "Feature selection for highdimensional genomic microarray data," in *Proc. ICML*, vol. 1, 2001, pp. 601–608.
- [12] R. Tibshirani, "Regression shrinkage and selection via the lasso," J. Roy. Stat. Soc., Ser. B, Methodol., vol. 58, no. 1, pp. 267–288, Jan. 1996.
- [13] C.-H. Zhang and J. Huang, "The sparsity and bias of the lasso selection in high-dimensional linear regression," *Ann. Statist.*, vol. 36, no. 4, pp. 1567–1594, Aug. 2008.
- [14] H. Huttunen, T. Manninen, J.-P. Kauppi, and J. Tohka, "Mind reading with regularized multinomial logistic regression," *Mach. Vis. Appl.*, vol. 24, no. 6, pp. 1311–1325, Aug. 2013.
- [15] T. Sato, Y. Takano, R. Miyashiro, and A. Yoshise, "Feature subset selection for logistic regression via mixed integer optimization," *Comput. Optim. Appl.*, vol. 64, no. 3, pp. 865–880, Jul. 2016.
- [16] H. Deng and G. Runger, "Feature selection via regularized trees," in Proc. Int. Joint Conf. Neural Netw. (IJCNN), Jun. 2012, pp. 1–8.
- [17] T. Brathwaite, A. Vij, and J. L. Walker, "Machine learning meets microeconomics: The case of decision trees and discrete choice," 2017, arXiv:1711.04826. [Online]. Available: https://arxiv.org/abs/1711.04826
- [18] A. Lhéritier, M. Bocamazo, T. Delahaye, and R. Acuna-Agost, "Airline itinerary choice modeling using machine learning," *J. Choice Model.*, vol. 31, pp. 198–209, Jun. 2019.
- [19] B. Sifringer, V. Lurkin, and A. Alahi, "Enhancing discrete choice models with neural networks," in *Proc. 18th Swiss Transp. Res. Conf.*, 2018, pp. 16–18.
- [20] G. Tutz, W. Pößnecker, and L. Uhlmann, "Variable selection in general multinomial logit models," *Comput. Statist. Data Anal.*, vol. 82, pp. 207–222, Feb. 2015.
- [21] A. Paz, C. Arteaga, and C. Cobos, "Specification of mixed logit models assisted by an optimization framework," *J. Choice Model.*, vol. 30, pp. 50–60, Mar. 2019.
- [22] N. Ortelli, F. C. Pereira, F. Rodrigues, and M. Bierlaire, "Assisted utility specification in discrete choice models," *Unpublished Manuscript*, to be published.
- [23] D. J. MacKay, "Bayesian non-linear modeling for the prediction competition," in *Maximum Entropy and Bayesian Methods*. Dordrecht, The Netherlands: Springer, 1996, pp. 221–234.
- [24] C. M. Bishop, Pattern Recognition and Machine Learning. New York, NY, USA: Springer-Verlag, 2006.
- [25] J. Drugowitsch, "Variational Bayesian inference for linear and logistic regression," 2013, arXiv:1310.5438. [Online]. Available: http://arxiv.org/ abs/1310.5438
- [26] T. S. Jaakkola and M. I. Jordan, "Bayesian parameter estimation via variational methods," *Statist. Comput.*, vol. 10, no. 1, pp. 25–37, 2000.
- [27] R. Ranganath, S. Gerrish, and D. Blei, "Black box variational inference," in *Proc. 17th Int. Conf. Artif. Intell. Statist.*, vol. 33. PMLR, 2014, pp. 814–822. [Online]. Available: http://proceedings.mlr. press/v33/ranganath14.html
- [28] Y. Song, F. S. Nathoo, and M. E. J. Masson, "A Bayesian approach to the mixed-effects analysis of accuracy data in repeated-measures designs," *J. Memory Lang.*, vol. 96, pp. 78–92, Oct. 2017.
- [29] M. Danaf, "Personalized recommendations using discrete choice models with inter-and intra-consumer heterogeneity," in *Proc. Int. Choice Modeling Conf.*, Mar. 2017. [Online]. Available: http://www. icmconference.org.uk/index.php/icmc/ICMC2017/paper/view/1350
- [30] M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul, "An introduction to variational methods for graphical models," *Mach. Learn.*, vol. 37, no. 2, pp. 183–233, Nov. 1999.
- [31] D. J. MacKay, Information Theory, Inference and Learning Algorithms. Cambridge, U.K.: Cambridge Univ. Press, 2003.

- [32] D. M. Blei and J. D. Lafferty, "A correlated topic model of science," *Ann. Appl. Statist.*, vol. 1, no. 1, pp. 17–35, Jun. 2007.
- [33] D. A. Knowles and T. Minka, "Non-conjugate variational message passing for multinomial and binary regression," in *Proc. Adv. Neural Inf. Process. Syst.*, 2011, pp. 1701–1709.
- [34] H. Robbins and S. Monro, "A stochastic approximation method," Ann. Math. Statist., vol. 22, no. 3, pp. 400–407, 1951. [Online]. Available: https://projecteuclid.org/euclid.aoms/1177729586#info
- [35] M. D. Hoffman, D. M. Blei, C. Wang, and J. Paisley, "Stochastic variational inference," J. Mach. Learn. Res., vol. 14, no. 1, pp. 1303–1347, 2013.
- [36] T. Brathwaite and J. L. Walker, "Asymmetric, closed-form, finiteparameter models of multinomial choice," *J. Choice Model.*, vol. 29, pp. 78–112, Dec. 2018.
- [37] D. J. Rezende and S. Mohamed, "Variational inference with normalizing flows," 2015, arXiv:1505.05770. [Online]. Available: http://arxiv. org/abs/1505.05770
- [38] F. Rodrigues, "Scaling Bayesian inference of mixed multinomial logit models to very large datasets," 2020, arXiv:2004.05426. [Online]. Available: http://arxiv.org/abs/2004.05426



Filipe Rodrigues is an Associate Professor with the Technical University of Denmark (DTU), where he is working on machine learning models for understanding urban mobility and the behaviour of crowds. Previously, he was a H. C. Ørsted/Marie-Skłodowska Curie Actions (COFUND) Post-Doctoral Fellow, also at DTU, and working on spatio-temporal models of mobility demand with an emphasis on modeling uncertainty and the effects of special events. His research interests include machine learning, probabilistic graphi-

cal models, Bayesian inference, intelligent transportation systems, and urban mobility.



Nicola Ortelli is a Research and Teaching Assistant with the Transport and Mobility Laboratory, École Polytechnique Fédérale de Lausanne (EPFL), Switzerland, where he is currently pursuing the Ph.D. degree under the supervision of Prof. M. Bierlaire.



Michel Bierlaire is a Professor of the École Polytechnique Fédérale de Lausanne (EPFL), Switzerland, where he is the Director of the Transport and Mobility Laboratory. Since 2009, he has been the Director of TraCE, the Transportation Center. Since September 2017, he has been the Head of the Civil Engineering Institute, EPFL. He has been active in demand modeling, operations research, and dynamic traffic management systems. His research interests include the design, development, and applications of models and algorithms for the design, analysis, and management of transportation systems.



Francisco Camara Pereira (Member, IEEE) is a Professor with the Technical University of Denmark (DTU), where he leads the Machine Learning for Smart Mobility (MLSM) Research Group. Previously, he was a Research Scientist at the MIT and an Assistant Professor with the University of Coimbra. His main research interests include applying machine learning and pattern recognition to the context of transportation systems with the purpose of understanding and predicting mobility behavior, and modeling and optimizing the transportation system

as a whole. He was awarded several prestigious prizes, including an IEEE achievements award in 2009, the Singapore GYSS Challenge in 2013, and the Pyke Johnson Award from the Transportation Research Board in 2015.