

A Convex Optimization Approach for Computing Correlated Choice Probabilities with Many Alternatives

Selin Damla Ahipasaoglu, Xiaobo Li, Karthik Natarajan

Abstract—A popular discrete choice model that incorporates correlation information is the Multinomial Probit (MNP) model where the random utilities of the alternatives are chosen from a multivariate normal distribution. Computing the choice probabilities is challenging in the MNP model when the number of alternatives is large. Mishra, Natarajan and Teo (IEEE TAC, 2012) have proposed a semidefinite optimization approach to compute choice probabilities for the distribution of the random utilities that maximizes expected agent utility given only the mean, variance and covariance information. Their model is referred to as the Cross Moment (CMM) model. Computing the choice probabilities with many alternatives is challenging in the CMM model since one needs to solve large scale semidefinite programs. We develop a simpler formulation as a representative agent model by maximizing over the choice probabilities in the unit simplex where the objective function is the sum of the expected utilities and a strongly concave perturbation function. By characterizing the perturbation function for the CMM model and its gradient, we develop a simple first order gradient method with inexact line search to compute choice probabilities. We establish local linear convergence of this algorithm under mild assumptions on the choice probabilities. An implication of our results is that inverting the choice probabilities to compute the mean utilities is straightforward given any positive definite covariance matrix. Numerical experiments show that this method can compute choice probabilities for a large number of alternatives within a reasonable amount of time while explicitly capturing the correlation information. Comparisons with simulation methods for MNP and semidefinite programming methods for CMM indicate the efficacy of the method.

Index Terms—choice probability, optimization, stochastic systems, optimization algorithms

I. INTRODUCTION

Consider a finite and mutually exclusive set of alternatives denoted by $[n] = \{1, 2, \dots, n\}$. For example, the alternatives might represent a set of differentiated products such as cars or laptops from which agents make their choices. Given a set of agents, each of whom chooses their most preferred alternative, a modeler with limited information on the agent preferences is interested in characterizing their discrete choice behavior.

This project was partly funded by the SUTD-MIT International Design Center grant IDG31300105 on ‘Optimization for Complex Discrete Choice’ and the MOE Tier 2 grant MOE2013-T2-2-168 on ‘Distributional Robust Optimization for Consumer Choice in Transportation Systems’.

Engineering Systems and Design, Singapore University of Technology and Design, 8 Somapah Road, SG 487372. Email: ahipasaoglu@sutd.edu.sg

Industrial and Systems Engineering, University of Minnesota, Minneapolis, MN USA 55455. Email: lixx3195@umn.edu

Engineering Systems and Design, Singapore University of Technology and Design, 8 Somapah Road, SG 487372. Email: karthik_natarajan@sutd.edu.sg

A. Discrete Choice Models

Random utility maximization is the most popular approach used to study this choice problem. In the additive random utility maximization model, the utility of alternative i is specified as:

$$\tilde{u}_i = \mu_i + \tilde{\epsilon}_i, \quad \forall i \in [n], \quad (1)$$

where μ_i is the deterministic component of the utility and $\tilde{\epsilon}_i$ is the random component of the utility of alternative i . Let $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)^T$ denote the vector of the deterministic component of the utility which is common across the agents and $\tilde{\boldsymbol{\epsilon}} = (\tilde{\epsilon}_1, \dots, \tilde{\epsilon}_n)^T$ denote the vector of the random components with a joint probability distribution $\theta(\cdot)$. Each realization of the random component vector results in a utility vector for an agent in the population who chooses the alternative with the highest utility. Assuming no ties, the probability that i is the most preferred alternative in this population of agents is given by:

$$p_i = P \left(i = \operatorname{argmax}_{k \in [n]} (\mu_k + \tilde{\epsilon}_k) \right). \quad (2)$$

An alternative derivation of the choice probability that is particularly relevant to this paper is the ‘representative agent’ model from economics (see Anderson, Palma and Thisse [3], [2]). In this model, the aggregate behavior of a population of agents is described through the choices made by a single representative agent who has a preference for diversity and randomizes choice. Consider a representative agent who chooses a probability vector in the unit $(n-1)$ -simplex:

$$\Delta_{n-1} = \left\{ \mathbf{x} \in \mathbb{R}_n^+ \mid \mathbf{e}^T \mathbf{x} = 1 \right\},$$

where \mathbf{e} is a vector of all ones. Given a convex deterministic function $V(\cdot) : \Delta_{n-1} \rightarrow \mathbb{R}$ that characterizes the preference for diversity, the representative agent solves the convex optimization problem (3).

$$\mathbf{p} = \operatorname{argmax} \left\{ \boldsymbol{\mu}^T \mathbf{x} - V(\mathbf{x}) \mid \mathbf{x} \in \Delta_{n-1} \right\}. \quad (3)$$

Hofbauer and Sandholm [21] show that all random utility models (with some regularity conditions) can be derived from the representative agent model where the function $V(\cdot)$ is a continuously differentiable, strictly convex function in the simplex and becomes infinitely steep near the relative boundary of the simplex.

Formulation (3) is interesting in its generality since the choice probability vectors can be computed as the solution

to a convex optimization problem and thus opens up an alternate computational approach to find choice probabilities. Specific choices of perturbation functions that have been studied include the quadratic function $V(\mathbf{x}) = \sum_i x_i^2$, Tsallis entropy $V(\mathbf{x}) = \frac{1}{q(1-q)} \sum_i x_i - x_i^q$ for $q > 0$, Renyi entropy $V(\mathbf{x}) = -\frac{1}{1-q} \ln \sum_i x_i^q$ for $q \in (0, 1)$ and the logarithm function $V(\mathbf{x}) = -\sum_i \ln x_i$. However, it seems difficult to find an explicit computationally tractable representation for the perturbation function $V(\cdot)$ in terms of the distribution function $\theta(\cdot)$, particularly with correlated utilities. To the best of our knowledge, random utility maximization models for which the perturbation function $V(\cdot)$ is explicitly known are the multinomial logit model and the nested logit model (see Verboven [43]). These correspond to instances where the choice probabilities are known in closed form. For example in the MNP model with a general correlation structure among the random terms, no explicit representation of the function $V(\cdot)$ seems to be known. Feng, Li and Wang [13] have recently shown that the representative agent model in (3) is equivalent to a ‘‘semi-parametric’’ choice model proposed by Natarajan, Song and Teo [33], where the choice probability is evaluated for a distribution (or a sequence of distributions) in a prescribed set of distributions Θ that maximize expected utility:

$$Z^* = \max_{\tilde{\epsilon} \sim_{\theta} \Theta} \mathbb{E}_{\theta} \left(\max_{i \in [n]} (\mu_i + \tilde{\epsilon}_i) \right). \quad (4)$$

B. Modeling Correlation

The Multinomial Logit (MNL) model, popularized by Luce [28] and McFadden [29], is perhaps the most popular choice model. The MNL model has the following choice probability formula:

$$p_i^{\text{mnl}} = \frac{e^{\mu_i}}{\sum_{k \in [n]} e^{\mu_k}}, \quad \forall i \in [n].$$

While the MNL choice probability is known in closed form and possesses desirable properties such as concavity of the log-likelihood function, it also suffers from drawbacks. One of the well-known properties of MNL is the Independence of Irrelevant Alternatives (IIA) property which implies that the ratio of the choice probabilities for any two alternatives is independent of the utilities of the other alternatives:

$$\frac{p_i^{\text{mnl}}}{p_j^{\text{mnl}}} = e^{\mu_i - \mu_j}, \quad \forall i \neq j.$$

When the alternatives have correlated utilities, the IIA property of MNL gives rise to misleading choice predictions. Many examples highlighting the limitations of IIA from a behavioral aspect are discussed in the literature, the red bus/blue bus example from McFadden [29] being probably the most famous. A second drawback of the MNL models is the Invariant Proportion of Substitution (IPS) discussed in Steenburgh [38]. Under the MNL model, improving an attribute of a product increases the demand for the product by decreasing the demand for all products proportional to their market shares.

Since similarities between products may exist, IPS property is considered counter-intuitive in many applications.

Correlation information can be captured using the Generalized Extreme Value (GEV) model (see McFadden [30]), which is a generalization of MNL model. The GEV model relaxes the assumption that the perception errors are independent random variables and, therefore, allows modelling meaningful substitution relationships among the alternatives. It also includes all closed form utility maximization formulations based on the extreme value error distribution with equal variance across alternatives. While the choice probabilities in the GEV model still have a closed form expression, the model does not allow for all possible correlation structures among the random terms. Another popular model that overcomes the IIA property is the Mixed Logit (MIX-MNL) model (see McFadden and Train [31]). In the MIX-MNL model, the utility of the alternative i is specified as:

$$\tilde{u}_i = \tilde{\beta}^T z_i + \tilde{\epsilon}_i, \quad \forall i \in [n], \quad (5)$$

where $\tilde{\beta}$ is a random coefficient vector that captures heterogeneity in the tastes of agents. While the agents know their true realization of the taste coefficient vector, from the modeler’s perspective, $\tilde{\beta}$ is a random vector with a density function given by $f(\cdot)$. Under the assumption that the components of $\tilde{\epsilon}$ are still identically and independently distributed Gumbel random variables, the choice probability in the mixed logit model is given by the integral of the logit probabilities as follows:

$$p_i^{\text{mix-mnl}} = \int \frac{e^{\beta^T z_i}}{\sum_{k \in [n]} e^{\beta^T z_k}} f(\beta) d\beta, \quad \forall i \in [n].$$

A common assumption on $\tilde{\beta}$ is that the individual components are independent normal random variables with $\tilde{\beta}_i \sim \text{Normal}(\mu_{\beta,i}, \sigma_{\beta,i}^2)$. However, there is no closed form expression for this integral and simulation is the most commonly used technique to estimate the choice probabilities (see Train [41]).

A model that accounts for any valid correlation matrix is the Multinomial Probit (MNP) model in which $\tilde{\epsilon}$ is assumed to be normally distributed with mean $\mathbf{0}$ and covariance matrix Σ , namely $\tilde{\mathbf{u}} \sim \text{Normal}(\boldsymbol{\mu}, \Sigma)$. The MNP model is flexible in terms of modeling dependence and does not possess the IIA property. However the choice probabilities do not have a closed form expression with Monte Carlo simulation being the most commonly used method to find the choice probabilities. The reader is referred to Hajivassiliou, McFadden and Ruud [18] for an in-depth discussion of simulation techniques used to approximate the choice probabilities in MNP models with the Geweke-Hajivassiliou-Keane (GHK) simulator being the most commonly used technique among them (see Geweke [14], Hajivassiliou and McFadden [17], Keane [25]).

As discussed above, the models proposed to overcome the IIA property of MNL might have their own challenges as well. For example, the GEV model does not allow for general correlation structure among utilities while the MIX-MNL and MNP models are computationally expensive. Mishra, Natarajan and Teo [32] recently proposed a semi-parametric model called the

Cross Moment Model (CMM), which elevates both IIA and IPS properties and allows for general correlation structures. In this choice model, the joint distribution of $\tilde{\epsilon}$ is assumed to be only partially specified to the modeler. Specifically, the available information on the joint distribution is the first two moments of $\tilde{\epsilon}$. Let $\tilde{\epsilon} \sim_{\theta} (\mathbf{0}, \Sigma)$ denote the set of probability distributions for $\tilde{\epsilon}$ that satisfies the following two conditions: $E_{\theta}[\tilde{\epsilon}] = \mathbf{0}$ and $Cov_{\theta}[\tilde{\epsilon}] = \Sigma$. The modeler is then assumed to solve the optimization problem:

$$(CMM) \quad Z_{cmm}^* = \max_{\tilde{\epsilon} \sim_{\theta}(\mathbf{0}, \Sigma)} \mathbb{E}_{\theta} \left(\max_{i \in [n]} (\mu_i + \tilde{\epsilon}_i) \right). \quad (6)$$

The outer optimization in (6) is over all joint distributions of the random components that are consistent with the two moment information. Hence, problem (6) is equivalent to finding a joint distribution for the random components that maximizes the expected agent utility¹. Mishra, Natarajan and Teo [32] solved the moment problem (6) by reformulating it as the following semidefinite program:

$$\begin{aligned} Z_{cmm}^* = \max \quad & \sum_{i \in [n]} e_i^T \mathbf{y}_i \\ \text{s.t.} \quad & \sum_{i \in [n]} \begin{pmatrix} \mathbf{W}_i & \mathbf{y}_i \\ \mathbf{y}_i^T & x_i \end{pmatrix} = \begin{pmatrix} \Sigma + \mu\mu^T & \mu \\ \mu^T & 1 \end{pmatrix} \\ & \begin{pmatrix} \mathbf{W}_i & \mathbf{y}_i \\ \mathbf{y}_i^T & x_i \end{pmatrix} \succeq 0, \quad \forall i \in [n], \end{aligned} \quad (7)$$

where e_i is a vector with 1 in the i th position and 0 otherwise. Let $\{\mathbf{W}_i^*, \mathbf{y}_i^*, x_i^*\}$ for $i \in [n]$ be an optimal solution to the semidefinite program (7). The joint distribution of the random utilities $\tilde{\mathbf{u}}$ that maximizes the expected agent utility is a mixture of multivariate normal distributions given as:

$$\tilde{\mathbf{u}} = \left\{ \text{Normal} \left(\frac{\mathbf{y}_i^*}{x_i^*}, \frac{\mathbf{W}_i^*}{x_i^*} - \frac{\mathbf{y}_i^* \mathbf{y}_i^{*T}}{x_i^*} \right), \text{ w.p. } x_i^*, \quad \forall i \in [n]. \right.$$

More importantly, they showed that the optimal decision variables \mathbf{x}^* in the SDP formulation are the choice probabilities for the mixture of multivariate normal distributions which maximizes the expected agent utility. Mishra, Natarajan and Teo [32] provide applications of this formulation to problems in route choice, random walk theory and product line selection with the number of alternatives up to hundred. Numerical experiments in [32] showed that the CMM model captures correlation information in predicting choices and provides insights often qualitatively similar to MNP.

C. Structure and Contributions

The main contributions and the structure of the paper are summarized next:

- (a) In Section II, we provide a representative agent formulation for the CMM model with an explicit computation-ally tractable representation of the perturbation function

¹The problem of finding the joint distribution of the random components that minimizes the expected agent utility with the first two moment information reduces to Jensen's bound. This is uninteresting from a discrete choice modelling perspective since all the agents then choose the alternative with the highest mean.

$V(\cdot)$. We also prove strong concavity of the objective function for the CMM model under the assumption that the covariance matrix is strictly positive definite. Our results make extensive use of properties of the square root function of matrices from matrix analysis.

- (b) In Section III, we evaluate the gradient of the objective function in the representative agent formulation for the CMM model and provide the optimality conditions. As an application, we show that under the assumption of strict positive definiteness of the covariance matrix, there is a one-to-one mapping between the mean utilities and the choice probabilities (or market shares) in the CMM model. Inferring the mean utilities from the choice probabilities (or inverting the market shares) is easy in the CMM model without having to go through simulation methods.
- (c) In Section IV, we develop a simple first order gradient ascent method with inexact line search to compute choice probabilities. This is particularly appealing since it transforms the solution of a semidefinite program (SDP) for the CMM model to a gradient based method which makes it possible to solve large instances. We prove the algorithm is locally linearly convergent.
- (d) In Section V, we provide computational results for the CMM model. The computational results indicate that the decrease in the running time obtained from the gradient method is potentially of significant value when solving for choice probabilities with many alternatives since solving large scale semidefinite programs remains a challenge. While the choice probabilities in the CMM model are computed using convex optimization, the choice probabilities in the MNP model are computed using simulation. We provide numerical experiments to illustrate the efficacy of the model.

II. A REPRESENTATIVE AGENT FORMULATION FOR THE CROSS MOMENT (CMM) MODEL

We start the section with the result by Natarajan and Teo [34] that further reduces the size of the semidefinite program (7) by reformulating it as:

$$\begin{aligned} Z_{cmm}^* = \max \quad & \text{trace}(\mathbf{Y}) \\ \text{s.t.} \quad & \mathbf{x} \in \Delta_{n-1} \\ & \begin{pmatrix} \Sigma + \mu\mu^T & \mathbf{Y}^T & \mu \\ \mathbf{Y} & \text{Diag}(\mathbf{x}) & \mathbf{x} \\ \mu^T & \mathbf{x}^T & 1 \end{pmatrix} \succeq 0. \end{aligned} \quad (8)$$

In formulation (8) $\text{trace}(\mathbf{Y})$ is the trace of the matrix \mathbf{Y} and $\text{Diag}(\mathbf{x})$ is a diagonal matrix with the entries of \mathbf{x} along the diagonal. In formulation (8), the optimal x_i^* variables represent a lower bound on the choice probability for the distribution θ^* that maximizes the expected agent utility given the first two moments:

$$P_{\theta^*} \left(i = \underset{k \in [n]}{\text{argmax}} (\mu_k + \tilde{\epsilon}_k) \right) \geq x_i^*, \quad \forall i \in [n],$$

where equality holds if there are no ties. Since the multivariate normal distribution is a feasible distribution in the CMM

formulation, Z_{cmm}^* is an upper bound on the expected agent utility in MNP. Computationally these models differ in the way the choice probabilities are computed. In the MNP model, simulation techniques are used to compute the choice probabilities. On the other hand, the CMM model uses convex optimization techniques to solve the semidefinite program.

There has been an increasing interest in the literature on discrete choice models that deal with a large number of alternatives. Examples that have been studied includes the choice of lake recreation sites in the state of Wisconsin involving 589 alternatives (see Parsons and Kealy [36]), choice of car models involving 689 alternatives (see Brownstone, Bunch and Train [9]) and choice of messenger bags involving 3584 alternatives (see Toubia et al. [40]). Models that treat products as bundles of characteristics with an additive error term that accounts for variation in the taste for the products in conjunction with variation in taste for the characteristics of the products results in choices where the number of products (alternatives) is exponential in the number of characteristics. In a recent paper, Ahipasaoglu et. al. [1] used the CMM model as an alternative to MNP for computing choice probabilities in a traffic equilibrium problem and showed that it provides a practical alternative to MNP in estimating traffic flows. The correlation information in their model arises from origin-destination paths (alternatives) sharing common roads (characteristics). The number of paths in such networks might be exponential in the number of roads. In our computational experiments, we have found that solving the semidefinite program (8) using state of art interior point method based solvers such as SDPT3 version 4 (see [42], [39]) in MATLAB R2014 on a laptop with an Intel(R) i7-5600U CPU processor (2.6 GHz) with 4GB RAM works well when the number of alternatives is up to two hundred roughly. Solving large semidefinite programs with matrix size up to a few thousands still remains a computational challenge and is a subject of intense research in the optimization community. One such algorithm that is able to solve large scale SDPs is the research software SDPNAL+ version 0.3 that has been recently developed by Sun, Toh, Yang and Zhao (see [46] and [45])². In contrast to such general purpose codes that solve large scale SDPs, we develop a specialized method based on gradient ascent with inexact line search by deriving a representative agent reformulation of the CMM model. Numerical results in Section V show that such a method is suitable when the number of alternatives is large.

A. Optimization over the Unit Simplex

In this section, we develop a representative agent formulation for the CMM model that transforms the semidefinite program to a nonlinear maximization problem over the unit simplex.

Theorem 1: Assume that $\Sigma \succ 0$. Then the maximum expected agent utility Z_{cmm}^* in the CMM model is the opti-

²This code has kindly been made freely available at <http://www.math.nus.edu.sg/~mattohk/SDPNALplus.html> by the authors and is based on a semismooth Newton-Conjugate gradient augmented Lagrangian method coupled with a alternating direction method of multipliers.

mal objective value to the following nonlinear optimization problem over the unit simplex:

$$Z_{\text{cmm}}^* = \max_{\mathbf{x} \in \Delta_{n-1}} \boldsymbol{\mu}^T \mathbf{x} + \text{trace} \left(\left(\Sigma^{1/2} \mathbf{S}(\mathbf{x}) \Sigma^{1/2} \right)^{1/2} \right) \quad (9)$$

where $\mathbf{S}(\mathbf{x}) = \text{Diag}(\mathbf{x}) - \mathbf{x}\mathbf{x}^T \succeq 0$ and $\mathbf{B} = \mathbf{A}^{1/2}$ is the unique positive semidefinite square root of a matrix $\mathbf{A} \succeq 0$ such that $\mathbf{A} = \mathbf{B}^2$. Furthermore the optimal decision variables \mathbf{x}^* are the choice probabilities for the distribution that maximizes the expected agent utility.

Proof:

Applying Schur's lemma to the positive semidefinite matrix in formulation (8), we obtain the equivalent nonlinear semidefinite program:

$$Z_{\text{cmm}}^* = \max_{\mathbf{x} \in \Delta_{n-1}} \text{trace}(\mathbf{Y}) \quad (10)$$

$$\text{s.t.} \quad \begin{pmatrix} \Sigma & \mathbf{Y}^T - \boldsymbol{\mu}\mathbf{x}^T \\ \mathbf{Y} - \boldsymbol{\mu}\mathbf{x}^T & \text{Diag}(\mathbf{x}) - \mathbf{x}\mathbf{x}^T \end{pmatrix} \succeq 0.$$

Define a transformation of the variables by letting $\hat{\mathbf{Y}} = \mathbf{Y} - \boldsymbol{\mu}\mathbf{x}^T$. Then, $\text{trace}(\hat{\mathbf{Y}}) = \text{trace}(\mathbf{Y}) - \boldsymbol{\mu}^T \mathbf{x}$. This transforms the problem to the equivalent nonlinear semidefinite programming formulation:

$$Z_{\text{cmm}}^* = \max_{\mathbf{x} \in \Delta_{n-1}} \boldsymbol{\mu}^T \mathbf{x} + \text{trace}(\hat{\mathbf{Y}}) \quad (11)$$

$$\text{s.t.} \quad \begin{pmatrix} \Sigma & \hat{\mathbf{Y}}^T \\ \hat{\mathbf{Y}} & \mathbf{S}(\mathbf{x}) \end{pmatrix} \succeq 0,$$

where $\mathbf{S}(\mathbf{x}) = \text{Diag}(\mathbf{x}) - \mathbf{x}\mathbf{x}^T$. The matrix $\mathbf{S}(\mathbf{x})$ is positive semidefinite for all $\mathbf{x} \in \Delta_{n-1}$ since:

$$\mathbf{v}^T \mathbf{S}(\mathbf{x}) \mathbf{v} = \sum_{i \in [n]} v_i^2 x_i - \left(\sum_{i \in [n]} v_i x_i \right)^2, \quad \forall \mathbf{v} \in \mathfrak{R}_n,$$

$$\geq 0,$$

where the last inequality comes from $\mathbb{E}(\tilde{v}^2) \geq \mathbb{E}(\tilde{v})^2$ where the random variable \tilde{v} is defined to take value v_i with probability x_i for $i \in [n]$. The semidefinite program in (11) can be reformulated as a two-stage optimization problem of the form:

$$Z_{\text{cmm}}^* = \max \left\{ \boldsymbol{\mu}^T \mathbf{x} - V(\mathbf{x}) \mid \mathbf{x} \in \Delta_{n-1} \right\}, \quad (12)$$

where $\mathbf{x} \in \Delta_{n-1}$ is the first stage decision vector and $V(\mathbf{x})$ is the optimal value to the following second stage problem where $\hat{\mathbf{Y}}$ is the second stage matrix decision variable:

$$V(\mathbf{x}) = \min -\text{trace}(\hat{\mathbf{Y}}) \quad (13)$$

$$\text{s.t.} \quad \begin{pmatrix} \Sigma & \hat{\mathbf{Y}}^T \\ \hat{\mathbf{Y}} & \mathbf{S}(\mathbf{x}) \end{pmatrix} \succeq 0.$$

The second stage semidefinite program in (13) for a given value of \mathbf{x} has a closed form solution (see Dowson and Landau [11], Olkin and Pukelsheim [35] and Shapiro [37]).

Applying this result since $\text{range}(\mathbf{S}(\mathbf{x})) \subseteq \text{range}(\boldsymbol{\Sigma})$, the optimal second stage solution is given as:

$$\hat{\mathbf{Y}}^{*T} = \boldsymbol{\Sigma} \mathbf{S}(\mathbf{x})^{1/2} \left(\left(\mathbf{S}(\mathbf{x})^{1/2} \boldsymbol{\Sigma} \mathbf{S}(\mathbf{x})^{1/2} \right)^{1/2} \right)^\dagger \mathbf{S}(\mathbf{x})^{1/2}, \quad (14)$$

where $(\cdot)^\dagger$ denotes the Moore-Penrose pseudoinverse of a matrix. Hence, the optimal value of formulation (13) is:

$$\begin{aligned} V(\mathbf{x}) &= -\text{trace} \left(\boldsymbol{\Sigma} \mathbf{S}(\mathbf{x})^{1/2} \left(\left(\mathbf{S}(\mathbf{x})^{1/2} \boldsymbol{\Sigma} \mathbf{S}(\mathbf{x})^{1/2} \right)^{1/2} \right)^\dagger \mathbf{S}(\mathbf{x})^{1/2} \right), \\ &= -\text{trace} \left(\left(\mathbf{S}(\mathbf{x})^{1/2} \boldsymbol{\Sigma} \mathbf{S}(\mathbf{x})^{1/2} \right) \left(\left(\mathbf{S}(\mathbf{x})^{1/2} \boldsymbol{\Sigma} \mathbf{S}(\mathbf{x})^{1/2} \right)^{1/2} \right)^\dagger \right), \\ &= -\text{trace} \left(\left(\mathbf{S}(\mathbf{x})^{1/2} \boldsymbol{\Sigma} \mathbf{S}(\mathbf{x})^{1/2} \right)^{1/2} \right), \\ &= -\text{trace} \left(\left(\boldsymbol{\Sigma}^{1/2} \mathbf{S}(\mathbf{x}) \boldsymbol{\Sigma}^{1/2} \right)^{1/2} \right), \end{aligned} \quad (15)$$

where the second equality comes from the invariance of the trace under cyclic permutations, the third equality comes from the property of the pseudo-inverse that $\mathbf{A}(\mathbf{A}^{1/2})^\dagger = \mathbf{A}^{1/2} \mathbf{A}^{1/2} (\mathbf{A}^{1/2})^\dagger = \mathbf{A}^{1/2}$ and the last equality comes from the observation that for any $n \times n$ real square matrix \mathbf{A} , the matrices $\mathbf{A}\mathbf{A}^T$ and $\mathbf{A}^T\mathbf{A}$ have the same set of eigenvalues (see Horn and Johnson [22]). By substituting into (12), we obtain the desired result. ■

Remark 1: The second stage problem in (13) has been studied in [11], [35], [37] in the following context: Given two n -dimensional random vectors with covariance matrices $\boldsymbol{\Sigma}$ and $\mathbf{S}(\mathbf{x})$, find the cross moment matrix $\hat{\mathbf{Y}}^T$ between the two random vectors that minimizes the expected L_2 distance between the vectors. In the proof of Theorem 1, the two vectors in the second stage corresponds to the random component of the utility vector $\tilde{\mathbf{e}}$ and the random choice vector which chooses e_i (alternative i) with probability x_i . The first stage problem corresponds to finding the best probability vector \mathbf{x} .

B. Strong Concavity of the Objective Function

In this section, we prove strong concavity of the objective function in the representative agent formulation for the CMM model. The result is based on the following definitions of functions of (positive semidefinite) matrices. Consider a symmetric positive semidefinite matrix \mathbf{A} with an eigendecomposition $\mathbf{Q}\text{Diag}(\boldsymbol{\lambda})\mathbf{Q}^T$ where \mathbf{Q} is an orthonormal matrix and $\boldsymbol{\lambda}$ is the vector of nonnegative eigenvalues. Given a function $h(\cdot) : [0, \infty) \rightarrow [0, \infty)$, the matrix function is defined as $h(\mathbf{A}) = \mathbf{Q}\text{Diag}(h(\boldsymbol{\lambda}))\mathbf{Q}^T$ where $h(\boldsymbol{\lambda})$ is a vector whose i^{th} entry stores $h(\lambda_i)$. As is the convention, we use $\mathbf{A} \succeq \mathbf{B}$ to denote $\mathbf{A} - \mathbf{B} \succeq 0$.

Definition 1: Consider a function $h : [0, \infty) \rightarrow [0, \infty)$.

(a) The function h is operator monotone if for all $\mathbf{A}, \mathbf{B} \succeq 0$:

$$\mathbf{A} \succeq \mathbf{B} \implies h(\mathbf{A}) \succeq h(\mathbf{B}).$$

(b) The function h is operator concave if for all $\mathbf{A}, \mathbf{B} \succeq 0$ and $\lambda \in [0, 1]$:

$$h((1-\lambda)\mathbf{A} + \lambda\mathbf{B}) \succeq (1-\lambda)h(\mathbf{A}) + \lambda h(\mathbf{B}).$$

An example of a matrix function that is both operator monotone and operator concave is the square root function.

Theorem 2: (Special case of the Löwner-Heinz Theorem [27], [19])

The function $h(t) = t^{1/2}$ is both operator monotone and operator concave.

Before we introduce the key result of this section, we recall the definition of *strong convexity*.

Definition 2: A function $V(\mathbf{x}) : \mathcal{D} \rightarrow \Re$ where \mathcal{D} is a convex subset of \Re_n is strongly convex if for all $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ and $\lambda \in (0, 1)$, there exists a constant $m > 0$ such that

$$V(\lambda\mathbf{x} + (1-\lambda)\mathbf{y}) \leq \lambda V(\mathbf{x}) + (1-\lambda)V(\mathbf{y}) - \frac{m}{2} \lambda(1-\lambda) \|\mathbf{x} - \mathbf{y}\|^2.$$

This brings us to the following result for the objective function in formulation (9).

Theorem 3: The function $V(\mathbf{x}) = -\text{trace} \left(\left(\boldsymbol{\Sigma}^{1/2} \mathbf{S}(\mathbf{x}) \boldsymbol{\Sigma}^{1/2} \right)^{1/2} \right)$ defined on the unit simplex $\mathbf{x} \in \Delta_{n-1}$ is strongly convex for $\boldsymbol{\Sigma} \succ 0$.

Proof: For all $\mathbf{x}, \mathbf{y} \in \Delta_{n-1}$ with $\mathbf{x} \neq \mathbf{y}$ and $\lambda \in (0, 1)$, we have:

$$\begin{aligned} & \mathbf{S}(\lambda\mathbf{x} + (1-\lambda)\mathbf{y}) \\ &= \text{Diag}(\lambda\mathbf{x} + (1-\lambda)\mathbf{y}) - (\lambda\mathbf{x} + (1-\lambda)\mathbf{y})(\lambda\mathbf{x} + (1-\lambda)\mathbf{y})^T, \\ &= \lambda \text{Diag}(\mathbf{x}) + (1-\lambda) \text{Diag}(\mathbf{y}) - \lambda^2 \mathbf{x}\mathbf{x}^T - (1-\lambda)^2 \mathbf{y}\mathbf{y}^T \\ &\quad - \lambda(1-\lambda)(\mathbf{x}\mathbf{y}^T + \mathbf{y}\mathbf{x}^T), \\ &= \lambda (\text{Diag}(\mathbf{x}) - \mathbf{x}\mathbf{x}^T) + (1-\lambda) (\text{Diag}(\mathbf{y}) - \mathbf{y}\mathbf{y}^T) \\ &\quad + \lambda(1-\lambda)(\mathbf{x} - \mathbf{y})(\mathbf{x} - \mathbf{y})^T, \\ &= \lambda \mathbf{S}(\mathbf{x}) + (1-\lambda) \mathbf{S}(\mathbf{y}) + \lambda(1-\lambda)(\mathbf{x} - \mathbf{y})(\mathbf{x} - \mathbf{y})^T. \end{aligned}$$

Pre-multiplying and post-multiplying by $\boldsymbol{\Sigma}^{1/2}$ implies that:

$$\boldsymbol{\Sigma}^{1/2} \mathbf{S}(\lambda\mathbf{x} + (1-\lambda)\mathbf{y}) \boldsymbol{\Sigma}^{1/2} = \boldsymbol{\Sigma}^{1/2} (\lambda \mathbf{S}(\mathbf{x}) + (1-\lambda) \mathbf{S}(\mathbf{y}) + \lambda(1-\lambda)(\mathbf{x} - \mathbf{y})(\mathbf{x} - \mathbf{y})^T) \boldsymbol{\Sigma}^{1/2}. \quad (16)$$

Now let $\mathbf{A} = \boldsymbol{\Sigma}^{1/2} \mathbf{S}(\lambda\mathbf{x} + (1-\lambda)\mathbf{y}) \boldsymbol{\Sigma}^{1/2}$, $\mathbf{B} = \lambda \boldsymbol{\Sigma}^{1/2} \mathbf{S}(\mathbf{x}) \boldsymbol{\Sigma}^{1/2} + (1-\lambda) \boldsymbol{\Sigma}^{1/2} \mathbf{S}(\mathbf{y}) \boldsymbol{\Sigma}^{1/2}$, $\rho = \lambda(1-\lambda)$ and $\mathbf{w} = \boldsymbol{\Sigma}^{1/2}(\mathbf{x} - \mathbf{y})$. Using this notation, we can rewrite the equation (16) as:

$$\mathbf{A} = \mathbf{B} + \rho \mathbf{w}\mathbf{w}^T.$$

Let $\lambda_1(\mathbf{A}) \leq \lambda_2(\mathbf{A}) \leq \dots \leq \lambda_n(\mathbf{A})$ denote the eigenvalues of \mathbf{A} (and respectively $\lambda_i(\mathbf{B})$ for \mathbf{B}). For $\rho > 0$ with a rank one perturbation, Bunch, Nilsen and Sorensen [10] have shown that:

$$\lambda_i(\mathbf{A}) \geq \lambda_i(\mathbf{B}), \quad \forall i \in [n].$$

Let $\mathbf{a} = \rho \mathbf{w}^T \mathbf{w}$. Then there exists a vector $\boldsymbol{\beta} \geq 0$ with $\sum_i \beta_i = \mathbf{a}$ such that

$$\lambda_i(\mathbf{A}) = \lambda_i(\mathbf{B}) + \beta_i, \quad \forall i \in [n].$$

Hence a lower bound on the sum of the square roots of the eigenvalues of the matrix \mathbf{A} is obtained by solving the optimization problem:

$$\begin{aligned} \sum_{i \in [n]} \lambda_i(\mathbf{A})^{1/2} &\geq \min_{\boldsymbol{\beta}} \sum_{i \in [n]} (\lambda_i(\mathbf{B}) + \beta_i)^{1/2} \\ \text{s.t.} \quad \sum_{i \in [n]} \beta_i &= \mathbf{a}, \\ \beta_i &\geq 0, \quad \forall i \in [n]. \end{aligned}$$

The right hand side of the above inequality corresponds to minimizing a concave function over a simplex, therefore the minimizer must be attained by at least one of the vertices of the simplex. This gives:

$$\begin{aligned} \sum_{i \in [n]} \lambda_i(\mathbf{A})^{1/2} &\geq \min_{j \in [n]} \left\{ \sum_{i \neq j} \lambda_i(\mathbf{B})^{1/2} + (\lambda_j(\mathbf{B}) + a)^{1/2} \right\}, \\ &\geq \min_{j \in [n]} \left\{ \sum_{i \in [n]} \lambda_i(\mathbf{B})^{1/2} + \frac{a}{2\sqrt{\lambda_j(\mathbf{B}) + a}} \right\}, \\ &= \sum_{i \in [n]} \lambda_i(\mathbf{B})^{1/2} + \frac{a}{2\sqrt{\lambda_n(\mathbf{B}) + a}}, \end{aligned}$$

where the second inequality is from the supergradient inequality for the concave square root function. Clearly there exists a positive number M_1 , such that $a = \lambda(1-\lambda)(\mathbf{x} - \mathbf{y})^T \mathbf{\Sigma}(\mathbf{x} - \mathbf{y}) \leq M_1$ for all $\lambda \in (0, 1)$ and $\mathbf{x}, \mathbf{y} \in \Delta_{n-1}$. Similarly, there exists a positive number M_2 , such that $\lambda_n(\mathbf{B}) = \max_i \lambda_i(\mathbf{B}) \leq M_2$, for all $\lambda \in (0, 1)$ and $\mathbf{x}, \mathbf{y} \in \Delta_{n-1}$. Letting $\alpha = \frac{1}{\sqrt{M_1 + M_2}}$, we have

$$\sum_{i \in [n]} \lambda_i(\mathbf{A})^{1/2} \geq \sum_{i \in [n]} \lambda_i(\mathbf{B})^{1/2} + \frac{\alpha}{2} \lambda(1-\lambda)(\mathbf{x} - \mathbf{y})^T \mathbf{\Sigma}(\mathbf{x} - \mathbf{y}). \quad (17)$$

Since $\text{trace}(\mathbf{A}^{1/2}) = \sum_{i \in [n]} \lambda_i(\mathbf{A})^{1/2}$, by (17) we obtain

$$\begin{aligned} &\text{trace}(\mathbf{\Sigma}^{1/2} \mathbf{S}(\lambda \mathbf{x} + (1-\lambda)\mathbf{y}) \mathbf{\Sigma}^{1/2})^{1/2} \\ &\geq \text{trace}(\lambda \mathbf{\Sigma}^{1/2} \mathbf{S}(\mathbf{x}) \mathbf{\Sigma}^{1/2} + (1-\lambda) \mathbf{\Sigma}^{1/2} \mathbf{S}(\mathbf{y}) \mathbf{\Sigma}^{1/2})^{1/2} \\ &\quad + \frac{\alpha}{2} \lambda(1-\lambda)(\mathbf{x} - \mathbf{y})^T \mathbf{\Sigma}(\mathbf{x} - \mathbf{y}), \\ &\geq \lambda \text{trace}(\mathbf{\Sigma}^{1/2} \mathbf{S}(\mathbf{x}) \mathbf{\Sigma}^{1/2})^{1/2} + (1-\lambda) \text{trace}(\mathbf{\Sigma}^{1/2} \mathbf{S}(\mathbf{y}) \mathbf{\Sigma}^{1/2})^{1/2} \\ &\quad + \frac{\alpha}{2} \lambda(1-\lambda)(\mathbf{x} - \mathbf{y})^T \mathbf{\Sigma}(\mathbf{x} - \mathbf{y}), \end{aligned}$$

where the last inequality is from the concavity of the square root function. Let $\lambda_{\min}(\mathbf{\Sigma})$ be the smallest eigenvalue of $\mathbf{\Sigma}$. Since $\mathbf{\Sigma}$ is positive definite, then $\lambda_{\min}(\mathbf{\Sigma}) > 0$ and

$$\|\mathbf{\Sigma}^{1/2}(\mathbf{x} - \mathbf{y})\|^2 \geq \lambda_{\min}(\mathbf{\Sigma}) \|\mathbf{x} - \mathbf{y}\|^2.$$

Then by the definition of the function $V(\cdot)$, we obtain

$$V(\lambda \mathbf{x} + (1-\lambda)\mathbf{y}) \leq \dots$$

$$\lambda V(\mathbf{x}) + (1-\lambda)V(\mathbf{y}) - \frac{\alpha}{2} \lambda_{\min}(\mathbf{\Sigma}) \lambda(1-\lambda) \|\mathbf{x} - \mathbf{y}\|^2,$$

and therefore the function $V(\mathbf{x})$ is strongly convex on its domain for $\mathbf{\Sigma} \succ 0$, where the strong convexity parameter depends on the matrix $\mathbf{\Sigma}$. ■

III. OPTIMALITY CONDITIONS AND ITS IMPLICATIONS

One of the key advantages in developing the representative agent formulation for the CMM model is that it transforms the semidefinite program to a nonlinear strongly concave maximization problem over the unit simplex. In this section, we provide a characterization of the directional derivatives of the objective function and provide optimality conditions for the model. In addition, we show that as we approach the boundary of the feasible region from its interior, the (projected) gradient of the objective function blows to infinity. These results have important implications that we discuss in detail in this section.

A. Projected Gradient of the Objective Function

Consider a vector \mathbf{x} in the relative interior of the simplex and restrict the direction of the perturbation to be in the tangent space of Δ_{n-1} defined as $\bar{\Delta}_{n-1} = \{\mathbf{v} \in \mathfrak{R}_n \mid \mathbf{e}^T \mathbf{v} = 0\}$. Let $\|\mathbf{v}\|_2 = 1$. Then the directional derivative of $V(\mathbf{x})$ in the direction \mathbf{v} is defined as:

$$\nabla_{\mathbf{v}} V(\mathbf{x}) = \lim_{\epsilon \rightarrow 0} \frac{V(\mathbf{x} + \epsilon \mathbf{v}) - V(\mathbf{x})}{\epsilon}.$$

To compute the directional derivative, observe that:

$$\begin{aligned} V(\mathbf{x} + \epsilon \mathbf{v}) &= -\text{trace}(\mathbf{\Sigma}^{1/2} \mathbf{S}(\mathbf{x} + \epsilon \mathbf{v}) \mathbf{\Sigma}^{1/2})^{1/2}, \\ &= -\text{trace}(\mathbf{\Sigma}^{1/2} \mathbf{S}(\mathbf{x}) \mathbf{\Sigma}^{1/2} \\ &\quad + \epsilon \mathbf{\Sigma}^{1/2} (\text{Diag}(\mathbf{v}) - \mathbf{x} \mathbf{v}^T - \mathbf{v} \mathbf{x}^T) \mathbf{\Sigma}^{1/2} \\ &\quad - \epsilon^2 \mathbf{\Sigma}^{1/2} \mathbf{v} \mathbf{v}^T \mathbf{\Sigma}^{1/2})^{1/2}, \\ &= -\text{trace}(\mathbf{T}(\mathbf{x}) + \mathbf{E}_{\mathbf{v}}(\epsilon, \mathbf{x}))^{1/2}, \end{aligned}$$

where:

$$\mathbf{T}(\mathbf{x}) = \mathbf{\Sigma}^{1/2} \mathbf{S}(\mathbf{x}) \mathbf{\Sigma}^{1/2} \text{ and}$$

$$\mathbf{E}_{\mathbf{v}}(\epsilon, \mathbf{x}) = \epsilon \mathbf{\Sigma}^{1/2} (\text{Diag}(\mathbf{v}) - \mathbf{x} \mathbf{v}^T - \mathbf{v} \mathbf{x}^T) \mathbf{\Sigma}^{1/2} - \epsilon^2 \mathbf{\Sigma}^{1/2} \mathbf{v} \mathbf{v}^T \mathbf{\Sigma}^{1/2}.$$

The next lemma provides a characterization of the null space of the matrix $\mathbf{T}(\mathbf{x})$ and its relation to the null space of the matrix $\mathbf{E}_{\mathbf{v}}(\epsilon, \mathbf{x})$. This lemma is needed for the main result of this and the next section.

Lemma 1:

- (a) Let $\mathbf{x} = \begin{pmatrix} \mathbf{0} \\ \underline{\mathbf{x}} \end{pmatrix} \in \Delta_{n-1}$, where $\mathbf{0}$ is a vector of r zeros and $\underline{\mathbf{x}} \in \mathfrak{R}_{n-r}^+$ is a strictly positive vector for some integer r such that $0 \leq r \leq n-1$. Then the null space of the matrix $\mathbf{T}(\mathbf{x})$ is given as:

$$\text{Null}(\mathbf{T}(\mathbf{x})) = \left\{ k \mathbf{\Sigma}^{-1/2} \mathbf{z} \mid \mathbf{z} = \begin{pmatrix} z_1 \\ \mathbf{e} \end{pmatrix}, z_1 \in \mathfrak{R}_r \right\},$$

where \mathbf{e} is a vector of ones of dimension $n-r$ and $\text{rank}(\mathbf{T}(\mathbf{x})) = n-r-1$.

- (b) Particularly, when \mathbf{x} is in the relative interior of the simplex denoted by $\text{rint}(\Delta_{n-1})$, then $\text{rank}(\mathbf{T}(\mathbf{x})) = n-1$ and $\text{Null}(\mathbf{T}(\mathbf{x})) \subseteq \text{Null}(\mathbf{E}_{\mathbf{v}}(\epsilon, \mathbf{x}))$.

Proof:

- (a) Let $\mathbf{A} \circ \mathbf{B}$ denote the Hadamard product of two matrices of the same dimension. Any vector $\mathbf{z} \in \mathfrak{R}_n$ can be expressed as $\mathbf{z} = \mathbf{\Sigma}^{-1/2} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$, for some $z_1 \in \mathfrak{R}_r$ and $z_2 \in \mathfrak{R}_{n-r}$. Then:

$$\begin{aligned} \mathbf{T}(\mathbf{x}) \mathbf{z} &= \mathbf{\Sigma}^{1/2} \left(\text{Diag}(\mathbf{x}) \mathbf{\Sigma}^{1/2} \mathbf{z} - \mathbf{x} \mathbf{x}^T \mathbf{\Sigma}^{1/2} \mathbf{z} \right), \\ &= \mathbf{\Sigma}^{1/2} \mathbf{x} \circ \left(\mathbf{\Sigma}^{1/2} \mathbf{z} - \mathbf{x}^T \mathbf{\Sigma}^{1/2} \mathbf{z} \mathbf{e} \right), \end{aligned}$$

where the last equality comes from the observation that $\text{Diag}(\mathbf{x})(\mathbf{\Sigma}^{1/2} \mathbf{z}) = \mathbf{x} \circ (\mathbf{\Sigma}^{1/2} \mathbf{z})$ and $\mathbf{x}(\mathbf{x}^T \mathbf{\Sigma}^{1/2} \mathbf{z}) = \mathbf{x} \circ (\mathbf{x}^T \mathbf{\Sigma}^{1/2} \mathbf{z}) \mathbf{e}$. This is equivalent to

$$\mathbf{T}(\mathbf{x}) \mathbf{z} = \mathbf{\Sigma}^{1/2} \begin{pmatrix} \mathbf{0} \\ \underline{\mathbf{x}} \end{pmatrix} \circ \left(\begin{pmatrix} z_1 \\ z_2 \end{pmatrix} - \begin{pmatrix} \mathbf{0} \\ \underline{\mathbf{x}} \end{pmatrix}^T \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} \mathbf{e} \right).$$

Since $\underline{\mathbf{x}} > \mathbf{0}$, $\mathbf{T}(\mathbf{x}) \mathbf{z} = \mathbf{0}$ implies that:

$$z_2 - (\underline{\mathbf{x}}^T z_2) \mathbf{e} = \mathbf{0}.$$

Solving this equation gives $\mathbf{z} \in \text{Null}(\mathbf{T}(\mathbf{x})) = \{k\boldsymbol{\Sigma}^{-1/2}\mathbf{z} \mid \mathbf{z} = \begin{pmatrix} \mathbf{z}_1 \\ \mathbf{e} \end{pmatrix}\}$ where $\mathbf{z}_1 \in \mathfrak{R}_r$. Therefore, $\text{rank}(\text{Null}(\mathbf{T}(\mathbf{x}))) = r + 1$ and the rank-nullity theorem implies that $\text{rank}(\mathbf{T}(\mathbf{x})) = n - r - 1$.

- (b) For $\mathbf{x} \in \text{rint}(\boldsymbol{\Delta}_{n-1})$, all the entries are nonzero. To show that $\boldsymbol{\Sigma}^{-1/2}\mathbf{e}$ lies in the null space of the matrix $\mathbf{E}_v(\epsilon, \mathbf{x})$, observe that:

$$\begin{aligned} & \mathbf{E}_v(\epsilon, \mathbf{x})\boldsymbol{\Sigma}^{-1/2}\mathbf{e} \\ = & \epsilon\boldsymbol{\Sigma}^{1/2} \left(\text{Diag}(\mathbf{v}) - \mathbf{x}\mathbf{v}^T - \mathbf{v}\mathbf{x}^T \right) \boldsymbol{\Sigma}^{1/2}\boldsymbol{\Sigma}^{-1/2}\mathbf{e} \\ & - \epsilon^2\boldsymbol{\Sigma}^{1/2}\mathbf{v}\mathbf{v}^T\boldsymbol{\Sigma}^{1/2}\boldsymbol{\Sigma}^{-1/2}\mathbf{e}, \\ = & \epsilon\boldsymbol{\Sigma}^{1/2} \left(\text{Diag}(\mathbf{v})\mathbf{e} - \mathbf{x}\mathbf{v}^T\mathbf{e} - \mathbf{v}\mathbf{x}^T\mathbf{e} \right) - \epsilon^2\boldsymbol{\Sigma}^{1/2}\mathbf{v}\mathbf{v}^T\mathbf{e}, \\ = & \mathbf{0}, \end{aligned}$$

where the final equality comes from the observation that $\mathbf{e}^T\mathbf{x} = 1$ and $\mathbf{v}^T\mathbf{e} = 0$. Hence, $\text{Null}(\mathbf{T}(\mathbf{x})) \subseteq \mathbf{E}_v(\epsilon, \mathbf{x})$. ■

To prove the main theorem of this section, we make use of the Fréchet derivative of a matrix function which is defined as follows.

Definition 3: The Fréchet derivative of a real matrix function $g : \mathfrak{R}_{n \times n} \mapsto \mathfrak{R}_{n \times n}$ at $\mathbf{X} \in \mathfrak{R}_{n \times n}$ is a linear mapping $L_g : \mathfrak{R}_{n \times n} \mapsto \mathfrak{R}_{n \times n}$ such that $g(\mathbf{X} + \mathbf{E}) - g(\mathbf{X}) - L_g(\mathbf{X}, \mathbf{E}) = o(\|\mathbf{E}\|)$ for all $\mathbf{E} \in \mathfrak{R}_{n \times n}$.

The Fréchet derivative, if exists, is known to be unique. The Fréchet derivative for the matrix square root function, which exists when \mathbf{X} is positive definite, is the unique solution to the Sylvester equation (refer to Kenney and Laub [26], Higham [20]):

$$\mathbf{X}^{1/2}L_{1/2}(\mathbf{X}, \mathbf{E}) + L_{1/2}(\mathbf{X}, \mathbf{E})\mathbf{X}^{1/2} = \mathbf{E}. \quad (18)$$

Next, we derive the first order derivative of $V(\cdot)$.

Theorem 4: Define:

$$\mathbf{g}(\mathbf{x}) = -\frac{1}{2}\text{diag} \left(\boldsymbol{\Sigma}^{1/2}(\mathbf{T}^{1/2}(\mathbf{x}))^\dagger \boldsymbol{\Sigma}^{1/2} \right) + \boldsymbol{\Sigma}^{1/2}(\mathbf{T}^{1/2}(\mathbf{x}))^\dagger \boldsymbol{\Sigma}^{1/2}\mathbf{x}, \quad (19)$$

where $\text{diag}(\cdot)$ is the column vector formed by the diagonal elements of the matrix and $\mathbf{T}(\mathbf{x}) = \boldsymbol{\Sigma}^{1/2}\mathbf{S}(\mathbf{x})\boldsymbol{\Sigma}^{1/2}$. The directional derivative of $V(\mathbf{x})$ at $\mathbf{x} \in \text{rint}(\boldsymbol{\Delta}_{n-1})$ in the direction $\mathbf{v} \in \overline{\boldsymbol{\Delta}}_{n-1}$ is $\nabla_{\mathbf{v}}V(\mathbf{x}) = \mathbf{g}(\mathbf{x})^T\mathbf{v}$, and its projected gradient on the tangent space is:

$$\overline{\nabla}V(\mathbf{x}) = \mathbf{g}(\mathbf{x}) - \frac{1}{n}\mathbf{e}^T\mathbf{g}(\mathbf{x})\mathbf{e}. \quad (20)$$

Proof: Lemma 1 implies that for the given symmetric matrices $\mathbf{T}(\mathbf{x})$ and $\mathbf{E}_v(\epsilon, \mathbf{x})$, there exists an orthogonal matrix \mathbf{P} with $\mathbf{P}\mathbf{P}^T = \mathbf{P}^T\mathbf{P} = \mathbf{I}$ such that

$$\begin{aligned} \mathbf{T}(\mathbf{x}) &= \mathbf{P} \begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \overline{\boldsymbol{\Lambda}}(\mathbf{x}) \end{pmatrix} \mathbf{P}^T \text{ and} \\ \mathbf{E}_v(\epsilon, \mathbf{x}) &= \mathbf{P} \begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \overline{\mathbf{E}}_v(\epsilon, \mathbf{x}) \end{pmatrix} \mathbf{P}^T, \end{aligned}$$

where $\overline{\boldsymbol{\Lambda}}(\mathbf{x})$ is a diagonal matrix of size $(n-1) \times (n-1)$ containing the non-zero eigenvalues of $\mathbf{T}(\mathbf{x})$ and \mathbf{P} is the

matrix of eigenvectors of matrix $\mathbf{T}(\mathbf{x})$ with the first eigenvector equal to $\frac{\boldsymbol{\Sigma}^{-1/2}\mathbf{e}}{\|\boldsymbol{\Sigma}^{-1/2}\mathbf{e}\|_2}$. The matrix $\overline{\mathbf{E}}_v(\epsilon, \mathbf{x})$ is however not necessarily diagonal. Thus, we obtain:

$$\begin{aligned} & V(\mathbf{x} + \epsilon\mathbf{v}) - V(\mathbf{x}) \\ = & -\text{trace} \left((\mathbf{T}(\mathbf{x}) + \mathbf{E}_v(\epsilon, \mathbf{x}))^{1/2} \right) + \text{trace} \left(\mathbf{T}(\mathbf{x})^{1/2} \right), \\ = & -\text{trace} \left(\mathbf{P} \begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \overline{\boldsymbol{\Lambda}}(\mathbf{x}) + \overline{\mathbf{E}}_v(\epsilon, \mathbf{x}) \end{pmatrix} \mathbf{P}^T \right)^{1/2} \\ & + \text{trace} \left(\mathbf{P} \begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \overline{\boldsymbol{\Lambda}}(\mathbf{x}) \end{pmatrix} \mathbf{P}^T \right)^{1/2}, \\ = & -\text{trace} \left(\begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \overline{\boldsymbol{\Lambda}}(\mathbf{x}) + \overline{\mathbf{E}}_v(\epsilon, \mathbf{x}) \end{pmatrix} \right)^{1/2} + \text{trace} \left(\begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \overline{\boldsymbol{\Lambda}}(\mathbf{x}) \end{pmatrix} \right)^{1/2}, \\ = & -\text{trace} \left((\overline{\boldsymbol{\Lambda}}(\mathbf{x}) + \overline{\mathbf{E}}_v(\epsilon, \mathbf{x}))^{1/2} - \overline{\boldsymbol{\Lambda}}^{1/2}(\mathbf{x}) \right). \end{aligned}$$

To evaluate the last expression, let $L_{1/2}(\overline{\boldsymbol{\Lambda}}(\mathbf{x}), \overline{\mathbf{E}}_v(\epsilon, \mathbf{x}))$ denote the Fréchet derivative for the matrix square root which is the unique solution to the Sylvester equation:

$$\overline{\boldsymbol{\Lambda}}^{-1/2}(\mathbf{x})L_{1/2}(\overline{\boldsymbol{\Lambda}}(\mathbf{x}), \overline{\mathbf{E}}_v(\epsilon, \mathbf{x})) + L_{1/2}(\overline{\boldsymbol{\Lambda}}(\mathbf{x}), \overline{\mathbf{E}}_v(\epsilon, \mathbf{x}))\overline{\boldsymbol{\Lambda}}^{-1/2}(\mathbf{x}) = \overline{\mathbf{E}}_v(\epsilon, \mathbf{x}).$$

The existence of solution is guaranteed since $\overline{\boldsymbol{\Lambda}}(\mathbf{x}) \succ 0$. The Sylvester equation can then be expressed as:

$$\begin{aligned} L_{1/2}(\overline{\boldsymbol{\Lambda}}(\mathbf{x}), \overline{\mathbf{E}}_v(\epsilon, \mathbf{x})) + \overline{\boldsymbol{\Lambda}}^{-1/2}(\mathbf{x})L_{1/2}(\overline{\boldsymbol{\Lambda}}(\mathbf{x}), \overline{\mathbf{E}}_v(\epsilon, \mathbf{x}))\overline{\boldsymbol{\Lambda}}(\mathbf{x})^{1/2} \\ = \overline{\boldsymbol{\Lambda}}^{-1/2}(\mathbf{x})\overline{\mathbf{E}}_v(\epsilon, \mathbf{x}). \end{aligned}$$

Hence:

$$\text{trace} \left(L_{1/2}(\overline{\boldsymbol{\Lambda}}(\mathbf{x}), \overline{\mathbf{E}}_v(\epsilon, \mathbf{x})) \right) = \frac{1}{2}\text{trace} \left(\overline{\boldsymbol{\Lambda}}^{-1/2}(\mathbf{x})\overline{\mathbf{E}}_v(\epsilon, \mathbf{x}) \right).$$

Using the definition of the Fréchet derivative we have:

$$\begin{aligned} & V(\mathbf{x} + \epsilon\mathbf{v}) - V(\mathbf{x}) \\ = & -\frac{1}{2}\text{trace} \left(\overline{\boldsymbol{\Lambda}}^{-1/2}(\mathbf{x})\overline{\mathbf{E}}_v(\epsilon, \mathbf{x}) \right) + o \left(\|\overline{\mathbf{E}}_v(\epsilon, \mathbf{x})\| \right), \\ = & -\frac{1}{2}\text{trace} \left(\begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \overline{\boldsymbol{\Lambda}}^{-1/2}(\mathbf{x}) \end{pmatrix} \begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \overline{\mathbf{E}}_v(\epsilon, \mathbf{x}) \end{pmatrix} \right) \\ & + o \left(\|\overline{\mathbf{E}}_v(\epsilon, \mathbf{x})\| \right), \\ = & -\frac{1}{2}\text{trace} \left(\begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \overline{\boldsymbol{\Lambda}}^{-1/2}(\mathbf{x}) \end{pmatrix} \mathbf{P}^T \mathbf{E}_v(\epsilon, \mathbf{x}) \mathbf{P} \right) + o \left(\|\overline{\mathbf{E}}_v(\epsilon, \mathbf{x})\| \right), \\ = & -\frac{1}{2}\text{trace} \left(\mathbf{P} \begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \overline{\boldsymbol{\Lambda}}^{-1/2}(\mathbf{x}) \end{pmatrix} \mathbf{P}^T \mathbf{E}_v(\epsilon, \mathbf{x}) \right) + o \left(\|\overline{\mathbf{E}}_v(\epsilon, \mathbf{x})\| \right), \\ = & -\frac{\epsilon}{2}\text{trace} \left((\mathbf{T}^{1/2}(\mathbf{x}))^\dagger \boldsymbol{\Sigma}^{1/2} \left(\text{Diag}(\mathbf{v}) - \mathbf{x}\mathbf{v}^T - \mathbf{v}\mathbf{x}^T \right) \boldsymbol{\Sigma}^{1/2} \right) \\ & + o \left(\|\overline{\mathbf{E}}_v(\epsilon, \mathbf{x})\| \right), \\ = & -\frac{\epsilon}{2} \left(\text{diag} \left(\boldsymbol{\Sigma}^{1/2}(\mathbf{T}^{1/2}(\mathbf{x}))^\dagger \boldsymbol{\Sigma}^{1/2} \right)^T - 2\mathbf{x}^T \boldsymbol{\Sigma}^{1/2}(\mathbf{T}^{1/2}(\mathbf{x}))^\dagger \boldsymbol{\Sigma}^{1/2} \right) \mathbf{v} \\ & + o \left(\|\overline{\mathbf{E}}_v(\epsilon, \mathbf{x})\| \right). \end{aligned}$$

Hence, we obtain the expression for the directional derivative in the direction $\mathbf{v} \in \overline{\boldsymbol{\Delta}}_{n-1}$ as:

$$\begin{aligned} \nabla_{\mathbf{v}}V(\mathbf{x}) &= \lim_{\epsilon \rightarrow 0} \frac{V(\mathbf{x} + \epsilon\mathbf{v}) - V(\mathbf{x})}{\epsilon}, \\ &= -\frac{1}{2} \left(\text{diag} \left(\boldsymbol{\Sigma}^{1/2}(\mathbf{T}^{1/2}(\mathbf{x}))^\dagger \boldsymbol{\Sigma}^{1/2} \right)^T \right. \\ &\quad \left. - 2\mathbf{x}^T \boldsymbol{\Sigma}^{1/2}(\mathbf{T}^{1/2}(\mathbf{x}))^\dagger \boldsymbol{\Sigma}^{1/2} \right) \mathbf{v}, \\ &= \mathbf{g}(\mathbf{x})^T \mathbf{v}, \end{aligned}$$

where

$$\mathbf{g}(\mathbf{x}) = -\frac{1}{2} \left(\text{diag} \left(\boldsymbol{\Sigma}^{1/2}(\mathbf{T}^{1/2}(\mathbf{x}))^\dagger \boldsymbol{\Sigma}^{1/2} \right) - 2\boldsymbol{\Sigma}^{1/2}(\mathbf{T}^{1/2}(\mathbf{x}))^\dagger \boldsymbol{\Sigma}^{1/2} \right) \mathbf{x}.$$

Since this is true for all $\mathbf{v} \in \overline{\Delta}_{n-1}$, we obtain $\overline{\nabla}V(\mathbf{x})$ by projecting $\mathbf{g}(\mathbf{x})$ onto the tangent space $\overline{\Delta}_{n-1}$.³ ■

B. Optimality Conditions

The representative agent formulation for the CMM model is:

$$Z^* = \max \left\{ f(\mathbf{x}) \mid \mathbf{x} \in \Delta_{n-1} \right\} \text{ where } f(\mathbf{x}) = \boldsymbol{\mu}^T \mathbf{x} - V(\mathbf{x}).$$

Since the objective function is strongly concave and given the first-order derivatives of the objective function in Theorem 4, we can now write down the first-order optimality conditions for the CMM model as follows:

$$\begin{aligned} \overline{\nabla}f(\mathbf{x}) &= \left(\boldsymbol{\mu} - \frac{1}{n} \mathbf{e}^T \boldsymbol{\mu} \mathbf{e} \right) - \left(\mathbf{g}(\mathbf{x}) - \frac{1}{n} \mathbf{e}^T \mathbf{g}(\mathbf{x}) \mathbf{e} \right) = 0, \\ \mathbf{x} &\in \Delta_{n-1}, \end{aligned} \quad (21)$$

where $\overline{\nabla}f(\mathbf{x})$ is the projected gradient of f onto the tangent space of the feasible region. Next, we will discuss some of the implications of these optimality conditions.

C. Mapping between Mean Utilities and Choice Probabilities

In this section, we show a one-to-one correspondence between the mean utilities $\boldsymbol{\mu}$ under appropriate normalization and the choice probabilities \mathbf{x} in the relative interior of the simplex in the CMM model. This is important from an modeling viewpoint since it shows that the CMM model is capable of generating all the choice probabilities in the relative interior of the unit simplex. Furthermore, this is important in identification and estimation of demand parameters (see Berry [5]). We show that under mild assumptions on the covariance matrix, inverting the choice probabilities in the CMM model is fairly easy. Towards this, we first prove the following lemma that characterizes the gradient of the objective function in the representative agent formulation of the CMM model near the relative boundary of the simplex.

Theorem 5: Assume that $\boldsymbol{\Sigma} \succ 0$. As \mathbf{x} approaches the relative boundary of the unit simplex, the projected gradient of $V(\cdot)$ blows up to $+\infty$.

Proof: Suppose that the sequence of interior points $\{\mathbf{x}_k\}_{k=1, \dots, \infty}$ approaches a point $\hat{\mathbf{x}}$ on the relative boundary of the unit simplex, along the direction $-\mathbf{z} \in \mathbb{R}_n$. Assume that $\hat{\mathbf{x}}$ has exactly m zeros. Any such \mathbf{z} must satisfy $\mathbf{e}^T \mathbf{z} = 0$ and $\hat{x}_i = 0 \Rightarrow z_i > 0$. We first prove that:

$$\lim_{t \rightarrow 0^+} \frac{V(\hat{\mathbf{x}}) - V(\hat{\mathbf{x}} + t\mathbf{z})}{t} = +\infty.$$

Let $z_0 = \min_{i: \hat{x}_i=0} z_i$ and $\sigma_1 = \lambda_1(\boldsymbol{\Sigma})$. We have:

$$\begin{aligned} \mathbf{T}(\hat{\mathbf{x}} + t\mathbf{z}) &= \boldsymbol{\Sigma}^{1/2} \mathbf{S}(\hat{\mathbf{x}} + t\mathbf{z}) \boldsymbol{\Sigma}^{1/2}, \\ &= \boldsymbol{\Sigma}^{1/2} \text{Diag}(\hat{\mathbf{x}} + t\mathbf{z}) \boldsymbol{\Sigma}^{1/2} \\ &\quad - \boldsymbol{\Sigma}^{1/2} (\hat{\mathbf{x}} + t\mathbf{z}) (\hat{\mathbf{x}} + t\mathbf{z})^T \boldsymbol{\Sigma}^{1/2}. \end{aligned}$$

³We abuse the notation slightly here since the gradient of $V(\mathbf{x})$ does not exist outside the feasible region. For all theoretical and algorithmic purposes, the mathematical quantity that we calculate, i.e., $\overline{\nabla}V(\mathbf{x})$, behaves as the projected gradient. To be mathematically precise, one would embed the function into the affine subspace $\mathbf{e}^T \mathbf{x} = 1$ by substituting for one of the decision variables, but this approach is notationally cumbersome in exposition.

It is clear that $\min_i \{\hat{x}_i + tz_i\} = \min_{i: \hat{x}_i=0} tz_i = tz_0$ if t is sufficiently small. From Lemmas 4 and 5, both given in the Appendix, this implies that

$$\begin{aligned} \lambda_2(\boldsymbol{\Sigma}^{1/2} \mathbf{S}(\hat{\mathbf{x}} + t\mathbf{z}) \boldsymbol{\Sigma}^{1/2}) &\geq \lambda_1(\boldsymbol{\Sigma}^{1/2} \text{Diag}(\hat{\mathbf{x}} + t\mathbf{z}) \boldsymbol{\Sigma}^{1/2}) \\ &\geq t\sigma_1 z_0 \\ &> 0, \end{aligned}$$

since $\boldsymbol{\Sigma}^{1/2} \mathbf{S}(\hat{\mathbf{x}} + t\mathbf{z}) \boldsymbol{\Sigma}^{1/2}$ is a rank one update of $\boldsymbol{\Sigma}^{1/2} \text{Diag}(\hat{\mathbf{x}} + t\mathbf{z}) \boldsymbol{\Sigma}^{1/2}$. Therefore, together with Lemma 1, we have:

$$\begin{aligned} \lambda_1(\mathbf{T}(\hat{\mathbf{x}} + t\mathbf{z})) &= 0, \quad \lambda_2(\mathbf{T}(\hat{\mathbf{x}} + t\mathbf{z})) \geq t\sigma_1 z_0, \text{ and} \\ \lambda_3(\mathbf{T}(\hat{\mathbf{x}} + t\mathbf{z})), \dots, \lambda_n(\mathbf{T}(\hat{\mathbf{x}} + t\mathbf{z})) &> 0. \end{aligned}$$

Recall also that Lemma 1 gives:

$$\begin{aligned} \lambda_1(\mathbf{T}(\hat{\mathbf{x}})) &= \dots = \lambda_{m+1}(\mathbf{T}(\hat{\mathbf{x}})) = 0 \text{ and} \\ \lambda_{m+2}(\mathbf{T}(\hat{\mathbf{x}})), \dots, \lambda_n(\mathbf{T}(\hat{\mathbf{x}})) &> 0. \end{aligned}$$

Furthermore, from standard results for the continuity of the matrix eigenvalue function (see [16]), we know that $\|\lambda_j(\mathbf{T}(\hat{\mathbf{x}} + t\mathbf{z})) - \lambda_j(\mathbf{T}(\hat{\mathbf{x}}))\| \leq O(t)$, for all j . Combining above facts,

$$\begin{aligned} &\lim_{t \rightarrow 0^+} \frac{V(\hat{\mathbf{x}}) - V(\hat{\mathbf{x}} + t\mathbf{z})}{t} \\ &= \lim_{t \rightarrow 0^+} \frac{\sum_{j=1}^n \left(\sqrt{\lambda_j(\mathbf{T}(\hat{\mathbf{x}} + t\mathbf{z}))} - \sqrt{\lambda_j(\mathbf{T}(\hat{\mathbf{x}}))} \right)}{t} \\ &\geq \lim_{t \rightarrow 0^+} \frac{\sqrt{t\sigma_1 z_0} + \sum_{j=m+2}^n \left(\sqrt{\lambda_j(\mathbf{T}(\hat{\mathbf{x}}))} + O(t) - \sqrt{\lambda_j(\mathbf{T}(\hat{\mathbf{x}}))} \right)}{t} \\ &\rightarrow +\infty, \end{aligned}$$

since the ratio \sqrt{t}/t diverges to $+\infty$ as t approaches 0 and $\lim_{t \rightarrow 0^+} O(t)/t = 0$.

We next show that

$$\begin{aligned} &\lim_{t \rightarrow 0^+} |\nabla_{\mathbf{z}} V(\hat{\mathbf{x}} + t\mathbf{z})| \\ &= \lim_{t \rightarrow 0^+} \lim_{s \rightarrow 0^+} \left(\frac{V(\hat{\mathbf{x}} + t\mathbf{z}) - V(\hat{\mathbf{x}} + (t+s)\mathbf{z})}{s} \right) \\ &= +\infty. \end{aligned}$$

Suppose this is not the case. Since $V(\hat{\mathbf{x}})$ is convex, $\nabla_{\mathbf{z}} V(\hat{\mathbf{x}} + t\mathbf{z})$ is monotone in t , which implies that

$$\begin{aligned} &\liminf_{t \rightarrow 0^+} \lim_{s \rightarrow 0^+} \left(\frac{V(\hat{\mathbf{x}} + t\mathbf{z}) - V(\hat{\mathbf{x}} + (t+s)\mathbf{z})}{s} \right) \\ &= \limsup_{t \rightarrow 0^+} \lim_{s \rightarrow 0^+} \left(\frac{V(\hat{\mathbf{x}} + t\mathbf{z}) - V(\hat{\mathbf{x}} + (t+s)\mathbf{z})}{s} \right). \end{aligned}$$

Therefore, there has to be M such that

$$\begin{aligned} &\lim_{s \rightarrow 0^+} \left(\frac{V(\hat{\mathbf{x}} + t\mathbf{z}) - V(\hat{\mathbf{x}} + (t+s)\mathbf{z})}{s} \right) \leq M, \quad \forall t > 0 \quad (22) \\ &\text{with } \hat{\mathbf{x}} + t\mathbf{z} \in \text{rint}(\Delta_{n-1}). \end{aligned}$$

However, since $\lim_{t \rightarrow 0^+} \frac{V(\hat{\mathbf{x}}) - V(\hat{\mathbf{x}} + t\mathbf{z})}{t} = +\infty$, there exists $t_0 > 0$ such that $\frac{V(\hat{\mathbf{x}}) - V(\hat{\mathbf{x}} + t_0\mathbf{z})}{t_0} > 2M$. In addition, since V is continuous on Δ_{n-1} , it is also uniformly continuous on Δ_{n-1} . Therefore, there exists $\delta > 0$ such that $|V(\hat{\mathbf{x}}) - V(\hat{\mathbf{x}} + s\mathbf{z})| <$

$t_0 M$ for all $s \in [0, \delta)$. It is without loss of generality to assume that $\delta < t_0$. Therefore,

$$\begin{aligned} & \lim_{s \rightarrow 0^+} \left(\frac{V(\hat{\mathbf{x}} + \frac{\delta}{2} \mathbf{z}) - V(\hat{\mathbf{x}} + (\frac{\delta}{2} + s) \mathbf{z})}{s} \right) \\ & \geq \left(\frac{V(\hat{\mathbf{x}} + \frac{\delta}{2} \mathbf{z}) - V(\hat{\mathbf{x}} + t_0 \mathbf{z})}{t_0 - \frac{\delta}{2}} \right), \\ & \geq \left(\frac{V(\hat{\mathbf{x}} + \frac{\delta}{2} \mathbf{z}) - V(\hat{\mathbf{x}} + t_0 \mathbf{z})}{t_0} \right), \\ & > \left(\frac{V(\hat{\mathbf{x}}) - t_0 M - V(\hat{\mathbf{x}} + t_0 \mathbf{z})}{t_0} \right) > M. \end{aligned}$$

This is in contradiction with equation (22), which completes the proof. \blacksquare

We are now ready to prove the main result of this section. Hofbauer and Sandholm [21] have shown that given any joint distribution of the noise terms, the mapping from the deterministic components of the utilities $\boldsymbol{\mu}$ (under appropriate normalization) to the set of choice probabilities in the relative interior of the simplex is surjective, namely any vector of choice probabilities can be obtained by selecting suitable mean values. We show in the next theorem that under mild assumptions on the covariance matrix, there is a one-to-one correspondence between the mean utilities under the normalization condition $\mu_1 = 0$ and the choice probabilities in the relative interior of the simplex for the CMM model.

Theorem 6: Assume that $\boldsymbol{\Sigma} \succ 0$. Without loss of generality, set $\mu_1 = 0$. Let $\mathbf{p} = P(\boldsymbol{\mu}) : \{0\} \times \mathfrak{R}_{n-1} \rightarrow \Delta_{n-1}$ be the mapping from the mean utilities to the choice probabilities in the CMM model. Then $P(\cdot)$ is a bijection between $\{0\} \times \mathfrak{R}_{n-1}$ and the relative interior of the simplex Δ_{n-1} , namely there is a one-to-one correspondence between the mean utilities and the choice probabilities.

Proof:

- (a) We first show that every mean vector in $\{0\} \times \mathfrak{R}_{n-1}$ in the CMM model results in a unique vector of choice probabilities in the relative interior of the unit simplex. From the strong concavity of the objective function in the representative agent formulation of the CMM model (Theorems 1 and 3) and the observation that the gradient of the objective function blows up to infinity near the relative boundary of the simplex (Theorem 5), the choice probability vector in the CMM model lies strictly in the relative interior of the simplex and is unique.
- (b) We next show that every choice probability vector in the relative interior of the simplex maps to a unique mean vector in $\{0\} \times \mathfrak{R}_{n-1}$ in the CMM model. From the optimality conditions in (21) and with $\mu_1 = 0$, by multiplying with the vector \mathbf{e}_1 we have:

$$\frac{1}{n} \mathbf{e}^T \boldsymbol{\mu} = \frac{1}{n} \mathbf{e}^T \mathbf{g}(\mathbf{x}) - \mathbf{g}_1(\mathbf{x}).$$

Plugging in back to the optimality conditions, we obtain the mean utilities from the choice probabilities as follows:

$$\boldsymbol{\mu} = \mathbf{g}(\mathbf{x}) - \mathbf{g}_1(\mathbf{x}) \mathbf{e}. \quad (23)$$

Taken together, this implies there is a one-to-one correspondence between the set of deterministic utilities in $\{0\} \times \mathfrak{R}_{n-1}$ and the set of choice probabilities in

the relative interior of the unit simplex. \blacksquare

For the MNL model, the mean utilities are uniquely identified from the following simple formula:

$$\mu_i = \ln(p_i^{\text{mnl}}) - \ln(p_1^{\text{mnl}}), \quad \forall i \in [n].$$

A similar result exists for identifying the mean utilities from the nested logit model (see Berry [5]). For the CMM model, the mean utilities are uniquely identified from the simple calculation in (23) where $\mathbf{g}(\mathbf{x}) = -\frac{1}{2} \left(\text{diag} \left(\boldsymbol{\Sigma}^{1/2} (\mathbf{T}^{1/2}(\mathbf{x}))^\dagger \boldsymbol{\Sigma}^{1/2} \right) - 2 \boldsymbol{\Sigma}^{1/2} (\mathbf{T}^{1/2}(\mathbf{x}))^\dagger \boldsymbol{\Sigma}^{1/2} \mathbf{x} \right)$. To the best of our knowledge, no such easily computable formula is available for the MNP model.

IV. CALCULATING THE CHOICE PROBABILITIES IN THE CMM MODEL

A. The Gradient Ascent Algorithm

In this section we present a projected gradient ascent method⁴ to calculate the choice probabilities in the CMM model. The algorithm is given in Algorithm 1 in which stepsizes are chosen according to the well-known Armijo's line search rule.

Input: $\boldsymbol{\mu}, \boldsymbol{\Sigma}$, starting point \mathbf{x}_0 , initial step size $\alpha_0 \in (0, 1]$, $\beta \in (0, 1)$, $\tau \in (0, 1)$, tolerance $\epsilon > 0$.

Output: Optimal solution \mathbf{x}^* .

Initialize stopping criteria: $\text{criteria} \leftarrow \epsilon + 1$;

while $\text{criteria} > \epsilon$ **do**

$\alpha \leftarrow \alpha_0$

$\mathbf{x} \leftarrow \mathbf{x}_0 + \alpha \bar{\nabla} f(\mathbf{x}_0)$,

while $\mathbf{x} \notin \text{rint}(\Delta_{n-1})$ **or**

$f(\mathbf{x}) < f(\mathbf{x}_0) + \tau \alpha \|\bar{\nabla} f(\mathbf{x}_0)\|^2$ **do**

$\alpha \leftarrow \beta \alpha$

$\mathbf{x} \leftarrow \mathbf{x}_0 + \alpha \bar{\nabla} f(\mathbf{x}_0)$,

end

$\mathbf{x}_0 \leftarrow \mathbf{x}$

$\text{criteria} \leftarrow \|\mathbf{x} - \mathbf{x}_0\|$.

end

Algorithm 1: Projected gradient ascent algorithm with Armijo search

From Theorem 5, we know that the optimal solution lies in $\text{rint}(\Delta_{n-1})$. The algorithm presented in Algorithm 1 converges to the optimal solution (see [23]). While the objective function has a nice curvature (it is strongly concave), it does not have a Lipschitz continuous gradient near the relative boundary. In fact, the function itself does not satisfy the Lipschitz continuity condition near the relative boundary (see Theorem 5). In the next section, we show that if \mathbf{x} is sufficiently far away from the relative boundary of the feasible region, then the algorithm converges linearly for appropriately chosen parameters within a local neighborhood. As we show numerically in Section 5, this helps explain the good behavior

⁴Note that the presentation here is slightly different than the one in classical references such as Section 2.3 of Bertsekas [6], but the algorithm is the same since the projection is onto an affine subspace, i.e., $\text{Proj}_{\mathbf{Ax}=\mathbf{b}}(\hat{\mathbf{x}} + \nabla f) = \hat{\mathbf{x}} + \text{Proj}_{\mathbf{Ax}=\mathbf{0}}(\nabla f)$ for $\mathbf{Ax} = \mathbf{b}$.

of the algorithm in most cases while for some ill-conditioned problems where for example one of the choice probabilities is very low, the algorithm tends to be slower.

B. Local Linear Convergence of the Algorithm

We first show that with $\tau \in [0.5, 1)$, the distance between the solution at successive iterations and the optimal solution is non-increasing.

Lemma 2: Let $d_k = \|\mathbf{x}^k - \mathbf{x}^*\|$, where \mathbf{x}^* is the optimal solution and \mathbf{x}^k is the k -th iterate. Then $d_k \leq d_{k-1}$ for all $k > 0$ if $\tau \geq 0.5$.

Proof: From the definition of d_k , we have

$$\begin{aligned} d_{k+1}^2 &= \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 \\ &= \|\mathbf{x}^k + \alpha_k \bar{\nabla} f(\mathbf{x}^k) - \mathbf{x}^*\|^2 \\ &= d_k^2 + \alpha_k^2 \|\bar{\nabla} f(\mathbf{x}^k)\|^2 - 2\alpha_k \bar{\nabla} f(\mathbf{x}^k)^T (\mathbf{x}^* - \mathbf{x}^k). \end{aligned} \quad (24)$$

Since $f(\cdot)$ is concave, we have

$$\bar{\nabla} f(\mathbf{x}^k)^T (\mathbf{x}^* - \mathbf{x}^k) \geq f(\mathbf{x}^*) - f(\mathbf{x}^k) = \epsilon_k.$$

Note that $\epsilon_k \geq 0$ by definition. Combined with inequality (24), we have

$$d_{k+1}^2 \leq d_k^2 - \alpha_k (2\epsilon_k - \alpha_k \|\bar{\nabla} f(\mathbf{x}^k)\|^2).$$

From the Armijo's rule, we have

$$f(\mathbf{x}^{k+1}) \geq f(\mathbf{x}^k) + \tau \alpha_k \|\bar{\nabla} f(\mathbf{x}^k)\|^2.$$

Therefore,

$$\alpha_k \|\bar{\nabla} f(\mathbf{x}^k)\|^2 \leq \frac{1}{\tau} (f(\mathbf{x}^{k+1}) - f(\mathbf{x}^k)) = \frac{1}{\tau} (\epsilon_k - \epsilon_{k+1}).$$

Hence, we have

$$\begin{aligned} d_{k+1}^2 &\leq d_k^2 - \alpha_k (2\epsilon_k - \alpha_k \|\bar{\nabla} f(\mathbf{x}^k)\|^2), \\ &\leq d_k^2 - \alpha_k \left(2\epsilon_k - \frac{1}{\tau} (\epsilon_k - \epsilon_{k+1}) \right), \\ &\leq d_k^2 - \alpha_k (2\epsilon_k - 2(\epsilon_k - \epsilon_{k+1})), \\ &\leq d_k^2 - 2\alpha_k \epsilon_{k+1}, \\ &\leq d_k^2, \end{aligned}$$

where the third inequality uses the fact that $\tau \geq 0.5$. \blacksquare

Note that the above proof is very similar to the proof of Proposition 9.1.2 in [4] for the unconstrained convex case.

We are now ready to discuss the rate of convergence of the algorithm by choosing the parameter τ , carefully.

Theorem 7: If there exists a $\gamma \in (0, 1)$ such that $\|\mathbf{x}^0 - \mathbf{x}^*\| \leq \gamma x_{\min}^*$ where $x_{\min}^* = \min_i x_i^*$ and if $\tau = 0.5$, then

$$f(\mathbf{x}^*) - f(\mathbf{x}^k) \leq \theta^k (f(\mathbf{x}^*) - f(\mathbf{x}^0)),$$

where $\theta = 1 - \min\{m, \beta m/L\}$ with m defined as the strong convexity constant in the proof of Theorem 3 and

$$L = \frac{9n}{4(1-\gamma)x_{\min}^*\sigma_1} \|\Sigma^{1/2}\|^4 + \frac{n}{((1-\gamma)x_{\min}^*\sigma_1)^{1/2}} \|\Sigma^{1/2}\|^2,$$

where $\sigma_1 = \lambda_1(\Sigma)$.

Proof: From Lemma 2, we know that $\|\mathbf{x}^k - \mathbf{x}^*\| \leq \gamma x_{\min}^*$ for all $k > 0$. So we can restrict the feasible region to $\mathcal{X} = \Delta_{n-1} \cap \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\| \leq \gamma x_{\min}^*\}$. It is easily

seen that for all $\mathbf{x} \in \mathcal{X}$, $\min_i x_i \geq (1-\gamma)x_{\min}^*$ and, therefore, $\mathcal{X} \subset \text{rint}(\Delta_{n-1})$. Applying Theorem 8 given in the Appendix, $\bar{\nabla} f(\mathbf{x}^k)$ is Lipschitz continuous over \mathcal{X} with Lipschitz constant

$$L = \frac{9n}{4(1-\gamma)x_{\min}^*\sigma_1} \|\Sigma^{1/2}\|^4 + \frac{n}{((1-\gamma)x_{\min}^*\sigma_1)^{1/2}} \|\Sigma^{1/2}\|^2.$$

The result follows from the linear convergence rate result (see Boyd and Vandenberghe [8]) for $\tau \leq 0.5$, with $\theta = 1 - \min\{2m\tau, 2\beta\tau m/L\}$, where m is the strong convexity constant. In our case, $\tau = 0.5$, so we have

$$f(\mathbf{x}^*) - f(\mathbf{x}^k) \leq \theta^k (f(\mathbf{x}^*) - f(\mathbf{x}^0)),$$

where $\theta = 1 - \min\{m, \beta m/L\}$ with m being the constant in the proof of Theorem 3. \blacksquare

Since the algorithm converges globally (see Iusem [23]), Theorem 7 shows that there exists a large enough integer M , such that, after M iterations the algorithm converges linearly. The algorithm is thus locally linearly convergent. The typical behavior of the algorithm is presented in Figure 1. We provide a more detailed computational study regarding the convergence of the algorithm in the next section.

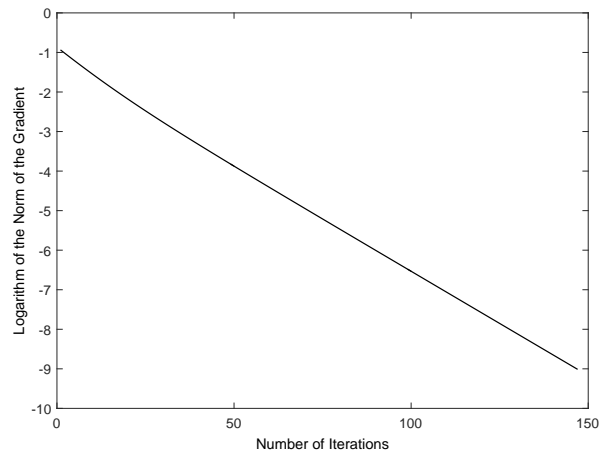


Fig. 1. Local convergence of the algorithm for a random instance with $n = 100$.

We note that while the algorithm we propose in this section is fairly straightforward and commonly used in nonlinear programming, the application of such a method to computing choice probabilities with correlated utilities appears to be new. Furthermore, we find that calculating the gradient of the objective function and proving the local linear convergence behaviour of the algorithm is nontrivial. We also believe that the ‘calculus with matrices’ that we develop here to solve discrete choice problems is novel and interesting on and possibly useful for solving other choice problems in the future as well.

V. COMPUTATIONAL RESULTS

In the first set of numerical experiments, we compare the computational times and accuracy of the gradient ascent

method developed in this paper for the CMM model with an SDP solver that is suitable for solving large scale SDPs. In the second set of experiments, we numerically test the convergence of the gradient method for the CMM model. In the third set of experiments, we compare the choice probabilities from the CMM and MNP model. The results help illustrate the efficacy of the model.

A. Comparison of SDP and the Gradient Method for the CMM Model

The SDP in (8) was solved using the SDPNAL+ version 0.3 (beta) while the code for the gradient ascent method was developed in MATLAB R2014⁵. The computational experiments were run on a laptop with an Intel(R) i7-5600U CPU processor (2.6 GHz) with 4GB RAM.

The number of alternatives n were varied in the set $\{100, 200, \dots, 1000\}$. The mean of the utilities was randomly generated in $[0, 1]^n$. The covariance matrix $\Sigma = \mathbf{V}\text{Diag}(\mathbf{d})\mathbf{V}^T$ was randomly generated by choosing the eigenvalues in the vector \mathbf{d} uniformly in $(0, 1]^n$ and the eigenvectors in \mathbf{V} using an orthogonalization of a random n by n matrix with each entry in $[-1, 1]$. For each size n , 7 instances were randomly generated. In the computational experiments, we used the default settings for SDPNAL+ version 0.3. For the gradient method, the parameters were set as $\alpha_0 = 0.1$, $\tau = 0.5$, $\beta = 0.6$ with $\epsilon = 1e-4$. To compare the accuracy of the methods, we evaluated the error measured in L_2 -distance between the choice probability vectors obtained from the SDP solver and the gradient ascent method:

$$\text{error}_{\text{prob}} = \|\mathbf{x}_{\text{sdp}}^* - \mathbf{x}_{\text{grad}}^*\|_2,$$

where $\mathbf{x}_{\text{sdp}}^*$ and $\mathbf{x}_{\text{grad}}^*$ are the solutions obtained from the SDP solver and the gradient ascent method, respectively. We also evaluated the difference in the optimal objective value as follows:

$$\text{error}_{\text{obj}} = |f(\mathbf{x}_{\text{sdp}}^*) - f(\mathbf{x}_{\text{grad}}^*)|.$$

The results are provided in Table I which clearly indicates that both the methods are very close in terms of choice probabilities and the objective value. In Table II, the computational times for the two methods are provided which illustrates that the gradient method converges much faster for this set of instances in comparison to the SDP solver.

B. Convergence of the Gradient Method for the CMM model

In the second set of numerical experiments, we study the convergence behavior of the algorithm. Theorem 7 shows that the region of linear convergence for the algorithm depends on x_{\min}^* , which is the distance between the optimal solution and the relative boundary. To study the effect of x_{\min}^* , we plot the number of iterations versus the level of accuracy achieved by the algorithm within those iterations, i.e., the tolerance ϵ , in Figure 2. To plot Figure 2, we picked $n = 100$ and randomly generated a covariance matrix Σ . Next, we chose a facet which is the convex combination of 10 randomly

⁵The code for the gradient method and the test instances can be obtained from the webpages of the authors.

n	error	1	2	3	4	5	6	7
100	Prob	1.63e-5	2.46e-5	3.69e-5	0.743e-5	3.08e-5	0.84e-5	2.92e-5
	Obj	0.0056	0.0071	0.0055	0.0020	0.0054	0.0076	0.0067
200	Prob	2.67e-5	2.53e-5	8.76e-5	2.36e-4	2.109e-4	2.672e-4	1.867e-4
	Obj	0.0071	0.0077	0.0077	0.0060	0.0050	0.0067	0.0046
300	Prob	9.40e-5	1.53e-4	0.88e-5	1.14e-4	2.34e-4	1.19e-4	1.04e-4
	Obj	0.0033	0.0021	0.0065	0.0018	0.0016	0.0004	0.0028
400	Prob	1.73e-4	2.89e-4	7.14e-5	0.95e-5	4.17e-5	3.78e-5	5.23e-5
	Obj	0.0031	0.0028	0.0006	0.0055	0.0028	0.0041	0.0026
500	Prob	8.79e-5	5.30e-5	2.95e-5	1.80e-5	3.54e-4	3.01e-5	4.29e-4
	Obj	0.0027	0.0099	0.0026	0.0070	0.0086	0.0065	0.0104
600	Prob	2.14e-5	2.10e-5	2.21e-6	2.82e-5	2.52e-5	2.70e-5	2.13e-5
	Obj	0.0055	0.0042	0.0021	0.0190	0.0044	0.1066	0.0047
700	Prob	7.16e-5	4.61e-5	1.45e-4	3.99e-5	2.31e-5	3.63e-5	1.82e-4
	Obj	0.0015	0.0012	0.0016	0.0010	0.0066	0.0019	0.0051
800	Prob	4.49e-5	1.78e-4	3.01e-5	2.19e-4	4.32e-4	4.51e-5	2.84e-4
	Obj	0.0015	0.0000	0.0010	0.0058	0.0166	0.0017	0.0163
900	Prob	1.12e-4	4.53e-5	4.08e-5	7.41e-4	4.86e-5	3.72e-4	3.06e-5
	Obj	0.0137	0.0108	0.0055	0.0299	0.0120	0.0171	0.0009
1000	Prob	3.88e-5	3.60e-4	3.62e-5	2.53e-4	7.74e-5	3.45e-5	5.56e-5
	Obj	0.0019	0.0088	0.0084	0.0004	0.0045	0.0024	0.0158

TABLE I
COMPARISON OF SDPNAL+ AND THE GRADIENT METHOD IN TERMS OF ACCURACY FOR 7 INSTANCES FOR EACH n .

n	Time (sec)	1	2	3	4	5	6	7
100	Grad	0.40	0.34	0.34	0.34	0.37	0.37	0.34
	SDP	11.82	6.78	8.23	65.49	6.22	7.84	6.35
200	Grad	1.43	1.17	1.23	1.38	1.27	1.35	1.29
	SDP	31.31	30.38	29.53	33.97	34.25	31.63	32.85
300	Grad	2.51	2.68	2.27	2.60	2.24	2.46	2.71
	SDP	114.45	99.74	101.27	114.67	116.59	121.58	113.42
400	Grad	4.35	3.60	3.47	3.86	3.83	3.49	3.65
	SDP	274.00	314.95	282.51	256.59	271.39	266.01	284.29
500	Grad	6.44	8.39	5.75	5.17	4.96	5.05	5.44
	SDP	617.63	548.72	527.94	477.75	585.20	467.12	521.08
600	Grad	13.61	14.24	12.62	12.94	13.35	13.35	13.49
	SDP	715.00	683.00	17903.00	2864.00	755.00	1829.00	655.00
700	Grad	23.16	23.49	20.87	19.92	21.38	20.51	19.84
	SDP	1220.30	1564.10	1692.80	1423.60	1163.10	1484.60	1264.40
800	Grad	33.08	33.83	33.86	38.14	39.65	43.21	39.65
	SDP	2304.20	1880.30	2330.00	2091.30	2812.10	2325.40	2717.60
900	Grad	49.18	52.07	53.41	48.36	52.50	48.95	53.83
	SDP	3665.10	3595.60	3663.40	3809.30	3728.10	3406.50	3571.90
1000	Grad	60.60	71.76	71.93	66.79	63.91	61.58	63.19
	SDP	4846.90	5220.90	5325.90	5398.60	4999.20	4840.20	4954.00

TABLE II
COMPARISON OF SDP SOLVER SDPNAL+ AND THE GRADIENT METHOD IN TERMS OF COMPUTATIONAL TIMES FOR 7 INSTANCES FOR EACH n .

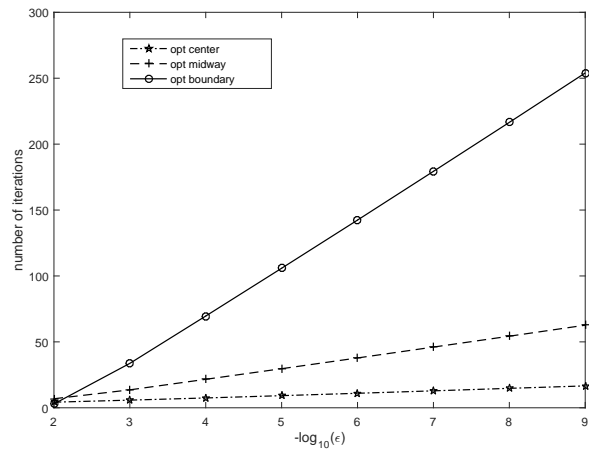


Fig. 2. Average behavior of the algorithm with different x_{\min}^* .

picked extreme points of the unit simplex, and let the center of the facet be \mathbf{x}^{bd} . We considered three scenarios: In the ‘opt center’ scenario, we let $\mathbf{x}^* = 0.9\mathbf{x}^{\text{ct}} + 0.1\mathbf{x}^{\text{bd}}$, where $\mathbf{x}^{\text{ct}} = \{1/100, \dots, 1/100\}$ is the center of the unit simplex; in the ‘opt midway’ scenario, we let $\mathbf{x}^* = 0.5\mathbf{x}^{\text{ct}} + 0.5\mathbf{x}^{\text{bd}}$; in the ‘opt boundary’ scenario, we let $\mathbf{x}^* = 0.1\mathbf{x}^{\text{ct}} + 0.9\mathbf{x}^{\text{bd}}$. Note that we can choose

$$\boldsymbol{\mu} = \mathbf{g}(\mathbf{x}^*)$$

to ensure that the optimal solution is \mathbf{x}^* . Clearly, $x_{\min}^* = 0.009, 0.005$ and 0.001 for these three scenarios. We randomly generate a starting point \mathbf{x}^0 . We then vary ϵ from 10^{-2} to 10^{-9} and record the corresponding number of iterations. Figure 2 is obtained by averaging the results of 20 independent replications. From the plot, we can clearly observe the local linear convergence behavior for all three scenarios. As the optimal solution approaches to the boundary, the slope of the plot increases indicating that the constant in the linear convergence rate result increases as x_{\min}^* decreases. This is also predicted by the theoretical results.

To study the influence of the starting point, we plot the average number of iterations and computation times versus the location of the starting point in the Figure 3. To obtain

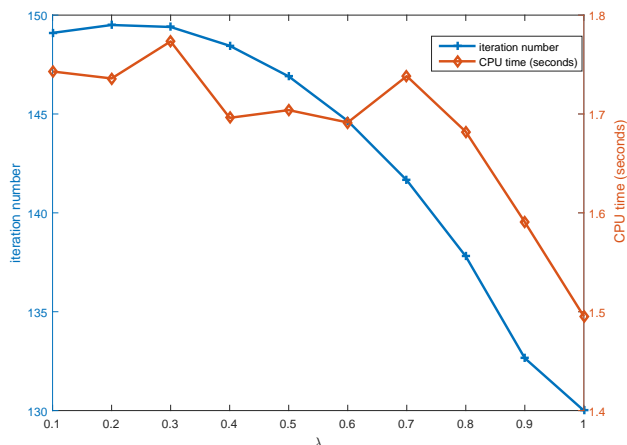


Fig. 3. Average number of iterations and CPU time versus the location of the initial point.

Figure 3, we use the settings as in Figure 2. Instead of varying the optimal solution, we randomly generate μ once to use in all experiments but choose various starting points on a line segment between a boundary point and the center of the simplex by setting $\mathbf{x}^0 = \lambda \mathbf{x}^{\text{ct}} + (1 - \lambda) \mathbf{x}^{\text{bd}}$ and varying the parameter λ . We fix the tolerance to $\epsilon = 10^{-6}$. From the figure, we find that required number of iterations to achieve the fixed tolerance level and the corresponding CPU times do not change much with respect to the starting point. The figure indicates that the location of the optimal solution seems to play a more important role for convergence than the location of the initial starting point.

C. Comparison of the CMM model and MNP model

In the last set of numerical experiments, we compare the choice probabilities for the MNP model obtained with simulation and for the CMM model obtained with the gradient ascent method.

1) *Small size examples from Börsch-Supan and Hajivassiliou [7]:* We first provide a comparison of the MNP and CMM choice probabilities for small size examples taken from Börsch-Supan and Hajivassiliou [7]. A popular alternative to the simple frequency simulator for MNP is the GHK simulator (see Geweke [14], [15], Hajivassiliou [17], and

Keane [24], [25]). The GHK simulator makes use of draws from truncated univariate normal distributions and requires evaluation of univariate integrals. Börsch-Supan and Hajivassiliou [7] have provided four examples with 5 alternatives to show that the GHK simulator produces probability estimates with substantially smaller variances than the simple frequency simulator. The details of the examples are provided in Table III. Example 1 involves mild correlations and has a small choice probability for the first alternative, Example 2 has slightly higher correlations, Example 3 has some large correlation coefficients while Example 4 has a choice probability close to 0.5 with mild correlations. Comparison of the choice probabilities obtained from the GHK simulator for the MNP model and the gradient ascent method for the CMM model are provided in Table III. The results indicate that the choice probability estimate for alternative 1 from the two models are fairly close to each other though developed under different assumptions on the utilities. For examples 1 and 2, where the choice probability of alternative 1 is small, the CMM model gives a higher choice probability for alternative 1 to be the most preferred one in comparison with the MNP model. On the other hand for examples 3 and 4, where the choice probability of alternative 1 is larger. The CMM model gives a slightly lower chance for alternative 1 to be the most preferred one in comparison with the MNP model.

Parameters	MNP	CMM
$\Delta\mu = \begin{pmatrix} -1.00 \\ -0.75 \\ -0.50 \\ -0.20 \end{pmatrix}, \Delta\Sigma = \begin{pmatrix} 1 & 0.2 & 0.3 & 0.1 \\ 0.2 & 1 & 0.4 & 0.3 \\ 0.3 & 0.4 & 1 & 0.5 \\ 0.1 & 0.3 & 0.5 & 1 \end{pmatrix}$	0.02409 (0.00068)	0.05366
$\Delta\mu = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \Delta\Sigma = \begin{pmatrix} 1 & 0.2 & 0.2 & 0.2 \\ 0.2 & 1 & 0.4 & 0.4 \\ 0.2 & 0.4 & 1 & 0.6 \\ 0.2 & 0.4 & 0.6 & 1 \end{pmatrix}$	0.15037 (0.00444)	0.15668
$\Delta\mu = \begin{pmatrix} 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \end{pmatrix}, \Delta\Sigma = \begin{pmatrix} 1 & 0.9 & 0 & 0 \\ 0.9 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0.95 \\ 0 & 0 & 0.95 & 1 \end{pmatrix}$	0.64773 (0.00773)	0.63789
$\Delta\mu = \begin{pmatrix} 1.50 \\ 0.75 \\ 0.50 \\ 0.75 \end{pmatrix}, \Delta\Sigma = \begin{pmatrix} 1 & 0.5 & 0.2 & 0.1 \\ 0.5 & 1 & 0.5 & 0.2 \\ 0.2 & 0.5 & 1 & 0.5 \\ 0.1 & 0.2 & 0.5 & 1 \end{pmatrix}$	0.49716 (0.01394)	0.47787

TABLE III
COMPARISON OF THE CHOICE PROBABILITY FOR ALTERNATIVE 1 BETWEEN THE MNP AND CMM MODEL. $\Delta\mu$ AND $\Delta\Sigma$ ARE THE MEAN AND THE COVARIANCE MATRIX FOR THE UTILITIES $(\tilde{u}_1 - \tilde{u}_2, \tilde{u}_1 - \tilde{u}_3, \tilde{u}_1 - \tilde{u}_4, \tilde{u}_1 - \tilde{u}_5)^T$. THE NUMBER IN PARENTHESIS INDICATES THE STANDARD DEVIATION OF THE ESTIMATOR.

2) *Larger example from Jester rating dataset:* In this example, we compare the choice probabilities from the CMM and the MNP model where data is available regarding the utilities of a large number of alternatives. We use the rating dataset from the Jester Online Joke Recommender System, in particular Dataset 2+⁶. The data consists of more than 2 million continuous ratings for 150 jokes collected from over 50000 individuals. Each individual provides ratings between -10 and 10 for a subset of the jokes, 10 of the jokes have never been rated and therefore excluded from the dataset.

To generate the utility parameter, we capture the data in a matrix of size 50000 by 140, whose the $(i, j)^{\text{th}}$ entry corresponds to the rating of individual i for joke j , if it exists. For the ratings that are incomplete, we use the standard Collaborative Filtering (CF) method, which is widely used

⁶<http://eigentaste.berkeley.edu/dataset/>

in recommendation engines. The user-based version of CF estimates a missing rating from individual i for joke j based on existing ratings for joke j from a set of individuals who are similar to individual i . Alternatively, the item-based version of CF uses the existing ratings of individual i for other items. We use the item-based CF in our application, since it is more suitable in situations where the number of items is significantly smaller than the number of individuals (see Ekstrand, Riedl, and Konstan [12] for a recent exhaustive survey on the topic).

To begin with, let us provide the details of the item-based CF. Let r_j denote the j^{th} column of the data matrix and $r(i, j)$ the existing rating from individual i for joke j . We calculate the *estimated* rating $\hat{r}(i, j)$ for individual i and joke j as follows:

$$\hat{r}(i, j) = \frac{\sum_{k \in J_i} w(j, k) r(i, k)}{\sum_{k \in J_i} |w(j, k)|},$$

where J_i is the set of jokes that have been rated by individual i and $w(j, k)$ is a measure of similarity between jokes j and k . Although there are other similarity measures in the literature, we use the cosine similarity, defined as $w(j, k) = \frac{r_j \cdot r_k}{\|r_j\|_2 \|r_k\|_2}$, for its simplicity, popularity, and good predictive properties. Similarly, the ratings can be estimated with alternative methods as well, nevertheless, the weighted average approach is a popular choice.

We use the *completed* data matrix to estimate the mean ratings $\mu \in \mathbb{R}^{140}$ over all users and the corresponding covariance matrix $\Sigma \in \mathbb{R}^{140 \times 140}$. Using these two parameter values, we calculate the choice probabilities, i.e., the probability that a joke j is the most preferred among all jokes,

- 1) Using the CMM and the gradient ascent algorithm developed in this paper with tolerance level $\epsilon = 10^{-3}$ and the rest of the parameters as in previous section. In our numerical experiments, the computational time for this approach is under 2 seconds.
- 2) Using the MNP model and the GHK simulator described above with 50000 samples. In our numerical experiments, the computational time for this approach is around 10 seconds.

We also calculate a basic in-sample statistic corresponding to the number of times a joke has the highest rating divided by the number of individuals. (Whenever there is a tie between l jokes for the highest rating, the count is incremented by $1/l$ instead of 1.) The choice probabilities from the CMM and MNP models together with the in-sample probabilities are provided in Figure 4. From the figure, we observe that the alternatives with very small choice probabilities in MNP take on higher choice probability values in the CMM model. On the other hand, the alternatives with larger choice probabilities in MNP take on smaller choice probability values in the CMM model. These results mirror the observations from the previous section. A possible explanation for this observation is that the distribution of the random utilities that maximizes expected agent utility in the CMM model is a mixture of multivariate normal distributions. The mixture of normals is a fat-tailed distribution and tends to give higher probabilities to the events that are low probability events in the standard

normal distribution. In terms of the trend, however the results clearly indicate that the alternatives that are more preferred in one model are also more preferred in the other model.

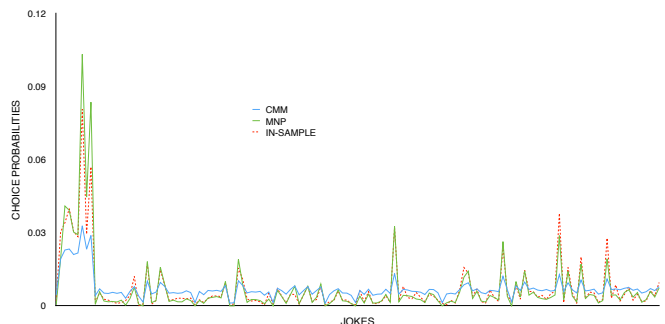


Fig. 4. Choice Probabilities from CMM, MNP, and In-Sample for the Jester Dataset

VI. CONCLUSION

In this paper, we have described a convex optimization approach to compute choice probabilities with correlated utilities. The choice model is derived for the joint distribution of the random utilities that maximizes expected agent utility given only the mean, variance and covariance information. Unlike MNP, the assumption of normality is dropped in this model. In contrast to MNP models where the choice probabilities are evaluated through simulation, we use a simple gradient ascent method to find the choice probabilities. The biggest advantage of the convex optimization approach is that one can compute choice probabilities for many alternatives with correlated utilities in a reasonable amount of time. In this era with consumers having more and more alternative options and increasing amounts of information, this paper proposes a new approach to computing choice probabilities that scales well with size. The next research question is to develop efficient inference techniques for the CMM model.

Acknowledgement

The authors would like to thank Kristin Wood (SUTD), Teo Chung Piau (NUS), Zheng Zhichao (SMU) and Rudabeh Meskarian (SUTD) for valuable discussions on this paper.

APPENDIX

Lemmas for Sections III and IV

Lemma 3: Suppose $\mathbf{x} > \mathbf{0}$ and let $x_{\min} = \min_i x_i$, $x_{\max} = \max_i x_i$. Then:

$$\|L_{1/2}(\text{Diag}(\mathbf{x}), \mathbf{E})\| \leq \frac{\|\mathbf{E}\|}{2x_{\min}^{1/2}},$$

and

$$\|L_{-1/2}(\text{Diag}(\mathbf{x}), \mathbf{E})\| \leq \frac{n\|\mathbf{E}\|}{2x_{\min}^{3/2}}.$$

Proof: The Fréchet derivative for the matrix inverse function is given as: ■

$$L_{-1}(\mathbf{X}, \mathbf{E}) = -\mathbf{X}^{-1}\mathbf{E}\mathbf{X}^{-1}.$$

The Fréchet derivative for the matrix square root function, which exists when \mathbf{X} is positive definite, is the unique solution to the Sylvester equation (refer to Kenney and Laub [26], Higham [20]):

$$\mathbf{X}^{1/2}L_{1/2}(\mathbf{X}, \mathbf{E}) + L_{1/2}(\mathbf{X}, \mathbf{E})\mathbf{X}^{1/2} = \mathbf{E}. \quad (25)$$

Following the chain rule (Theorem 3.4 in Higham [20]), we have

$$L_{-1/2}(\mathbf{X}, \mathbf{E}) = L_{1/2}(\mathbf{X}^{-1}, -\mathbf{X}^{-1}\mathbf{E}\mathbf{X}^{-1}),$$

and therefore,

$$\mathbf{X}^{-1/2}L_{-1/2}(\mathbf{X}, \mathbf{E}) + L_{-1/2}(\mathbf{X}, \mathbf{E})\mathbf{X}^{-1/2} = -\mathbf{X}^{-1}\mathbf{E}\mathbf{X}^{-1}. \quad (26)$$

Combining equations (25) and (26), we have

$$L_{-1/2}(\mathbf{X}, \mathbf{E}) = -\mathbf{X}^{-1/2}L_{1/2}(\mathbf{X}, \mathbf{E})\mathbf{X}^{-1/2}. \quad (27)$$

Define $\mathbf{L} = L_{1/2}(\text{Diag}(\mathbf{x}), \mathbf{E})$. From equation (25), we have:

$$\text{Diag}(\mathbf{x})^{1/2}\mathbf{L} + \mathbf{L}\text{Diag}(\mathbf{x})^{1/2} = \mathbf{E},$$

which implies that

$$|L_{i,j}| = \frac{|E_{i,j}|}{x_i^{1/2} + x_j^{1/2}} \leq \frac{|E_{i,j}|}{2x_{\min}^{1/2}}.$$

Therefore, we have

$$\|\mathbf{L}\| = \sqrt{\sum_{i,j} L_{i,j}^2} \leq \sqrt{\sum_{i,j} \frac{E_{i,j}^2}{4x_{\min}}} = \frac{\|\mathbf{E}\|}{2x_{\min}^{1/2}}.$$

In addition, by equation (27)

$$\begin{aligned} \|L_{-1/2}(\text{Diag}(\mathbf{x}), \mathbf{E})\| &= \|\text{Diag}(\mathbf{x})^{-1/2}\mathbf{L}\text{Diag}(\mathbf{x})^{-1/2}\| \\ &\leq \|\text{Diag}(\mathbf{x})^{-1/2}\|^2\|\mathbf{L}\| \\ &= \frac{n\|\mathbf{E}\|}{2x_{\min}^{3/2}} \end{aligned} \quad \blacksquare$$

Lemma 4: Let $\mathbf{B} = \mathbf{A} - \mathbf{u}\mathbf{u}^T$. Then $\lambda_1(\mathbf{B}) \leq \lambda_1(\mathbf{A})$, and

$$\lambda_{i-1}(\mathbf{A}) \leq \lambda_i(\mathbf{B}) \leq \lambda_i(\mathbf{A}), \quad \forall i = 2, \dots, n.$$

Proof: The proof can be found on page 97-98 of [44].

Lemma 5: Let $\mathbf{D} = \text{Diag}(d_1, \dots, d_n)$ be a diagonal matrix with $d_i > 0, \forall i \in [n]$. Let $\mathbf{\Sigma}$ be a positive definite matrix. Then $\lambda_1(\mathbf{\Sigma}^{1/2}\mathbf{D}\mathbf{\Sigma}^{1/2}) \geq \lambda_1(\mathbf{\Sigma}) \min_i \{d_i\}$.

Proof: Since $\text{eig}(AA^T) = \text{eig}(A^T A)$, we have

$$\lambda_1(\mathbf{\Sigma}^{1/2}\mathbf{D}\mathbf{\Sigma}^{1/2}) = \lambda_1(\mathbf{D}^{1/2}\mathbf{\Sigma}\mathbf{D}^{1/2}).$$

But

$$\mathbf{D}^{1/2}\mathbf{\Sigma}\mathbf{D}^{1/2} \succeq \lambda_1(\mathbf{\Sigma})\mathbf{D}^{1/2}\mathbf{I}\mathbf{D}^{1/2} = \lambda_1(\mathbf{\Sigma})\mathbf{D}.$$

Therefore, for all \mathbf{v} with $\|\mathbf{v}\| = 1$, we have

$$\mathbf{v}^T\mathbf{D}^{1/2}\mathbf{\Sigma}\mathbf{D}^{1/2}\mathbf{v} \geq \mathbf{v}^T\lambda_1(\mathbf{\Sigma})\mathbf{D}\mathbf{v} \geq \lambda_1(\mathbf{\Sigma})\min_i\{d_i\}.$$

Theorem 8: Assume that $\mathbf{\Sigma} \succ 0$. For any feasible direction $\mathbf{v} \in \overline{\Delta}_{n-1}$ with $\|\mathbf{v}\| = 1$ and $\mathbf{x} \in \text{rint}(\Delta_{n-1})$,

$$|f''_{\mathbf{x},\mathbf{v}}(0)| \leq \frac{9n}{4x_{\min}\sigma_1} \|\mathbf{\Sigma}^{1/2}\|^4 + \frac{n}{(x_{\min}\sigma_1)^{1/2}} \|\mathbf{\Sigma}^{1/2}\|^2,$$

where $f_{\mathbf{x},\mathbf{v}}(t) = V(\mathbf{x} + t\mathbf{v}) > 0$, $\sigma_1 = \lambda_1(\mathbf{\Sigma})$ and $x_{\min} = \min_i x_i$.

Proof: Let $\mathbf{g}(\mathbf{x}) = \mathbf{g}_1(\mathbf{x}) + \mathbf{g}_2(\mathbf{x})$ where

$$\mathbf{g}_1(\mathbf{x}) = -\frac{1}{2}\text{diag}\left(\mathbf{\Sigma}^{1/2}(\mathbf{T}^{1/2}(\mathbf{x}))^\dagger\mathbf{\Sigma}^{1/2}\right),$$

$$\mathbf{g}_2(\mathbf{x}) = \mathbf{\Sigma}^{1/2}(\mathbf{T}^{1/2}(\mathbf{x}))^\dagger\mathbf{\Sigma}^{1/2}\mathbf{x}.$$

From the proof of Theorem 4, we know that:

$$\begin{aligned} &\mathbf{g}_1(\mathbf{x} + \epsilon\mathbf{v})^T\mathbf{v} - \mathbf{g}_1(\mathbf{x})^T\mathbf{v} \\ &= -\frac{1}{2}\mathbf{v}^T\text{diag}\left(\mathbf{\Sigma}^{1/2}(\mathbf{T}^{1/2}(\mathbf{x} + \epsilon\mathbf{v}))^\dagger\mathbf{\Sigma}^{1/2}\right) + \frac{1}{2}\mathbf{v}^T\text{diag}\left(\mathbf{\Sigma}^{1/2}(\mathbf{T}^{1/2}(\mathbf{x}))^\dagger\mathbf{\Sigma}^{1/2}\right) \\ &= -\frac{1}{2}\mathbf{v}^T\text{diag}\left(\mathbf{\Sigma}^{1/2}\left((\mathbf{T}^{1/2}(\mathbf{x} + \epsilon\mathbf{v}))^\dagger - (\mathbf{T}^{1/2}(\mathbf{x}))^\dagger\right)\mathbf{\Sigma}^{1/2}\right) \\ &= -\frac{1}{2}\mathbf{v}^T\text{diag}\left(\mathbf{\Sigma}^{1/2}\mathbf{P}\begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & (\overline{\Lambda}(\mathbf{x}) + \overline{\mathbf{E}}_{\mathbf{v}}(\epsilon, \mathbf{x}))^{-1/2} - \overline{\Lambda}(\mathbf{x})^{-1/2} \end{pmatrix}\mathbf{P}^T\mathbf{\Sigma}^{1/2}\right) \\ &= -\frac{1}{2}\mathbf{v}^T\text{diag}\left(\mathbf{\Sigma}^{1/2}\mathbf{P}\begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & L_{-1/2}(\overline{\Lambda}(\mathbf{x}), \overline{\mathbf{E}}_{\mathbf{v}}(\epsilon, \mathbf{x})) \end{pmatrix}\mathbf{P}^T\mathbf{\Sigma}^{1/2}\right) + o(\epsilon). \end{aligned}$$

Therefore, we have:

$$\begin{aligned} &\|\mathbf{g}_1(\mathbf{x} + \epsilon\mathbf{v})^T\mathbf{v} - \mathbf{g}_1(\mathbf{x})^T\mathbf{v}\| \\ &= \left\| -\frac{1}{2}\mathbf{v}^T\text{diag}\left(\mathbf{\Sigma}^{1/2}\mathbf{P}\begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & L_{-1/2}(\overline{\Lambda}(\mathbf{x}), \overline{\mathbf{E}}_{\mathbf{v}}(\epsilon, \mathbf{x})) \end{pmatrix}\mathbf{P}^T\mathbf{\Sigma}^{1/2}\right) + o(\epsilon) \right\| \\ &\leq \frac{1}{2}\|\mathbf{v}\| \cdot \|\mathbf{\Sigma}^{1/2}\|^2 \cdot \left\| \mathbf{P}\begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & L_{-1/2}(\overline{\Lambda}(\mathbf{x}), \overline{\mathbf{E}}_{\mathbf{v}}(\epsilon, \mathbf{x})) \end{pmatrix}\mathbf{P}^T \right\| + o(\epsilon) \\ &\leq \frac{1}{2}\|\mathbf{\Sigma}^{1/2}\|^2 \cdot \left\| \mathbf{P}\begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & L_{-1/2}(\overline{\Lambda}(\mathbf{x}), \overline{\mathbf{E}}_{\mathbf{v}}(\epsilon, \mathbf{x})) \end{pmatrix}\mathbf{P}^T \right\| + o(\epsilon) \\ &\leq \frac{1}{2}\|\mathbf{\Sigma}^{1/2}\|^2 \cdot \|L_{-1/2}(\overline{\Lambda}(\mathbf{x}), \overline{\mathbf{E}}_{\mathbf{v}}(\epsilon, \mathbf{x}))\| + o(\epsilon) \\ &\leq \frac{1}{2}\|\mathbf{\Sigma}^{1/2}\|^2 \cdot \frac{n\|\overline{\mathbf{E}}_{\mathbf{v}}(\epsilon, \mathbf{x})\|}{2\lambda_2^{3/2}(\mathbf{T}(\mathbf{x}))} + o(\epsilon) \\ &= \frac{1}{2}\|\mathbf{\Sigma}^{1/2}\|^2 \cdot \frac{n\|\mathbf{E}_{\mathbf{v}}(\epsilon, \mathbf{x})\|}{2\lambda_2^{3/2}(\mathbf{T}(\mathbf{x}))} + o(\epsilon) \\ &= \frac{1}{2}\|\mathbf{\Sigma}^{1/2}\|^2 \cdot \frac{n\|\epsilon\mathbf{\Sigma}^{1/2}(\text{Diag}(\mathbf{v}) - \mathbf{x}\mathbf{v}^T - \mathbf{v}\mathbf{x}^T)\mathbf{\Sigma}^{1/2}\|}{2\lambda_2^{3/2}(\mathbf{T}(\mathbf{x}))} + o(\epsilon) \\ &\leq \frac{\epsilon}{2}\|\mathbf{\Sigma}^{1/2}\|^4 \cdot \frac{n\|\mathbf{v}\|(1+2\|\mathbf{x}\|)}{2\lambda_2^{3/2}(\mathbf{T}(\mathbf{x}))} + o(\epsilon) \\ &\leq \frac{\epsilon}{2}\|\mathbf{\Sigma}^{1/2}\|^4 \cdot \frac{3n}{2\lambda_2^{3/2}(\mathbf{T}(\mathbf{x}))} + o(\epsilon). \end{aligned}$$

Note that the last inequality holds since $\|\mathbf{x}\| \leq 1$ for all $\mathbf{x} \in$

Δ_{n-1} . On the other hand,

$$\begin{aligned}
& g_2(\mathbf{x} + \epsilon \mathbf{v})^T \mathbf{v} - g_2(\mathbf{x})^T \mathbf{v} \\
&= \mathbf{v}^T \Sigma^{1/2} (\mathbf{T}^{1/2}(\mathbf{x} + \epsilon \mathbf{v}))^\dagger \Sigma^{1/2}(\mathbf{x} + \epsilon \mathbf{v}) \\
&\quad - \mathbf{v}^T \Sigma^{1/2} (\mathbf{T}^{1/2}(\mathbf{x}))^\dagger \Sigma^{1/2} \mathbf{x} \\
&= \mathbf{v}^T \Sigma^{1/2} \mathbf{P} \begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & (\bar{\Lambda}(\mathbf{x}) + \bar{\mathbf{E}}_v(\epsilon, \mathbf{x}))^{-1/2} \end{pmatrix} \mathbf{P}^T \Sigma^{1/2}(\mathbf{x} + \epsilon \mathbf{v}) \\
&\quad - \mathbf{v}^T \Sigma^{1/2} \mathbf{P} \begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \bar{\Lambda}(\mathbf{x})^{-1/2} \end{pmatrix} \mathbf{P}^T \Sigma^{1/2} \mathbf{x} \\
&= \mathbf{v}^T \Sigma^{1/2} \mathbf{P} \begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & L_{-1/2}(\bar{\Lambda}(\mathbf{x}), \bar{\mathbf{E}}_v(\epsilon, \mathbf{x})) \end{pmatrix} \mathbf{P}^T \Sigma^{1/2} \mathbf{x} \\
&\quad + \epsilon \mathbf{v}^T \Sigma^{1/2} \mathbf{P} \begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \bar{\Lambda}(\mathbf{x})^{-1/2} \end{pmatrix} \mathbf{P}^T \Sigma^{1/2} \mathbf{v} + o(\epsilon).
\end{aligned}$$

Therefore, we have

$$\begin{aligned}
& \left\| g_2(\mathbf{x} + \epsilon \mathbf{v})^T \mathbf{v} - g_2(\mathbf{x})^T \mathbf{v} \right\| \\
&\leq \left\| \mathbf{v}^T \Sigma^{1/2} \mathbf{P} \begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & L_{-1/2}(\bar{\Lambda}(\mathbf{x}), \bar{\mathbf{E}}_v(\epsilon, \mathbf{x})) \end{pmatrix} \mathbf{P}^T \Sigma^{1/2} \mathbf{x} \right\| \\
&\quad + \epsilon \left\| \mathbf{v}^T \Sigma^{1/2} \mathbf{P} \begin{pmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \bar{\Lambda}(\mathbf{x})^{-1/2} \end{pmatrix} \mathbf{P}^T \Sigma^{1/2} \mathbf{v} \right\| + o(\epsilon) \\
&\leq \|\mathbf{x}\| \left\| \Sigma^{1/2} \right\|^2 \left\| L_{-1/2}(\bar{\Lambda}(\mathbf{x}), \bar{\mathbf{E}}_v(\epsilon, \mathbf{x})) \right\| \\
&\quad + \epsilon \left\| \Sigma^{1/2} \right\|^2 \left\| \bar{\Lambda}(\mathbf{x})^{-1/2} \right\| + o(\epsilon) \\
&\leq \left\| \Sigma^{1/2} \right\|^2 \cdot \frac{n \|\bar{\mathbf{E}}_v(\epsilon, \mathbf{x})\|}{2\lambda_2^{3/2}(\mathbf{T}(\mathbf{x}))} + \epsilon \left\| \Sigma^{1/2} \right\|^2 \frac{n}{\lambda_2(\mathbf{T}(\mathbf{x}))^{1/2}} + o(\epsilon) \\
&\leq \epsilon \left\| \Sigma^{1/2} \right\|^4 \cdot \frac{3n}{2\lambda_2^{3/2}(\mathbf{T}(\mathbf{x}))} + \epsilon \left\| \Sigma^{1/2} \right\|^2 \frac{n}{\lambda_2(\mathbf{T}(\mathbf{x}))^{1/2}} + o(\epsilon).
\end{aligned}$$

In addition, from Lemma 4 and Lemma 5, we have

$$\begin{aligned}
\lambda_2(\mathbf{T}(\mathbf{x})) &= \lambda_2(\Sigma^{1/2} \mathbf{S}(\mathbf{x}) \Sigma^{1/2}) \\
&\geq \lambda_1(\Sigma^{1/2} \text{Diag}(\mathbf{x}) \Sigma^{1/2}) \\
&\geq x_{\min} \sigma_1 > 0.
\end{aligned}$$

Therefore, we have

$$\begin{aligned}
& \left| f''_{\mathbf{x}, \mathbf{v}}(0) \right| \\
&= \left| \lim_{\epsilon \rightarrow 0^+} \frac{g_1(\mathbf{x} + \epsilon \mathbf{v})^T \mathbf{v} - g_1(\mathbf{x})^T \mathbf{v} + g_2(\mathbf{x} + \epsilon \mathbf{v})^T \mathbf{v} - g_2(\mathbf{x})^T \mathbf{v}}{\epsilon} \right| \\
&\leq \left| \lim_{\epsilon \rightarrow 0^+} \frac{g_1(\mathbf{x} + \epsilon \mathbf{v})^T \mathbf{v} - g_1(\mathbf{x})^T \mathbf{v}}{\epsilon} \right| + \left| \lim_{\epsilon \rightarrow 0^+} \frac{g_2(\mathbf{x} + \epsilon \mathbf{v})^T \mathbf{v} - g_2(\mathbf{x})^T \mathbf{v}}{\epsilon} \right| \\
&\leq \frac{9n}{4\lambda_2^{3/2}(\mathbf{T}(\mathbf{x}))} \left\| \Sigma^{1/2} \right\|^4 + \frac{n}{\lambda_2(\mathbf{T}(\mathbf{x}))^{1/2}} \left\| \Sigma^{1/2} \right\|^2 \\
&\leq \frac{9n}{4(x_{\min} \sigma_1)^{3/2}} \left\| \Sigma^{1/2} \right\|^4 + \frac{n}{(x_{\min} \sigma_1)^{1/2}} \left\| \Sigma^{1/2} \right\|^2.
\end{aligned}$$

Remark 2: Theorem 8 develops an upper bound for the absolute value of the second order derivatives in direction $\mathbf{v} \in \bar{\Delta}_{n-1}$. The significance of the theorem is that the second order derivative of $V(\mathbf{x})$ is bounded for all points $\mathbf{x} \in \text{rint}(\Delta_{n-1})$, with the bound associated with the minimum components of \mathbf{x} .

REFERENCES

- [1] S. D. AHIPASOGLU, R. MESKARIAN, L. MAGNANTI, THOMAS, AND K. NATARAJAN, *Beyond normality: A cross moment-stochastic user equilibrium model*, Transportation Research Part B: Methodological, 81 (2015), pp. 333–354.
- [2] S. ANDERSON, A. PALMA, AND J.-F. THISSE, *Discrete Choice Theory of Product Differentiation*, Cambridge: MIT Press, 1992.
- [3] S. ANDERSON, A. D. PALMA, AND J. THISSE, *A representative consumer theory of the logit model*, International Economic Review, 29 (1988), pp. 461–466.
- [4] A. BEN-TAL AND A. NEMIROVSKI, *Lectures on modern convex optimization: analysis, algorithms, and engineering applications*, vol. 2, Siam, 2001.
- [5] S. T. BERRY, *Estimating discrete-choice models of product differentiation*, RAND Journal of Economics, 25 (1994), pp. 242–262.
- [6] D. P. BERTSEKAS, *Nonlinear programming*, (1999).
- [7] A. BÖRSCH-SUPAN AND V. A. HAJIVASSILIOU, *Smooth unbiased multivariate probability simulators for maximum likelihood estimation of limited dependent variable models*, Journal of Econometrics, 58 (1993), pp. 347–368.
- [8] S. BOYD AND L. VANDENBERGHE, *Convex Optimization*, Cambridge University Press, 2004.
- [9] D. BROWNSTONE, D. S. BUNCH, AND K. TRAIN, *Joint mixed logit models of stated and revealed preferences for alternative-fuel vehicles*, Transportation Research B, 34 (2000), pp. 315–338.
- [10] J. BUNCH, C. NILSEN, AND D. SORENSEN, *Rank-one modification of the symmetric eigenproblem*, Numerische Mathematik, 31 (1978), pp. 31–48.
- [11] D. C. DOWSON AND B. V. LANDAU, *The fréchet distance between multivariate normal distributions*, Journal of Multivariate Analysis, 12 (1982), pp. 450–455.
- [12] M. D. EKSTRAND, J. T. RIEDL, AND J. A. KONSTAN, *Collaborative filtering recommender systems*, Foundations and Trends in Human-Computer Interaction, 4 (2010), pp. 81–173.
- [13] G. FENG, X. LI, AND Z. WANG, *On the relation between several discrete choice models*, Operations Research, (2017).
- [14] J. GEWEKE, *Bayesian inference in econometric models using monte carlo integration*, Econometrica, 57 (1989), pp. 1317–1339.
- [15] J. GEWEKE, *Efficient Simulation from the Multivariate Normal and Student-t Distributions Subject to Linear Constraints*, Defense Technical Information Center, 1992.
- [16] G. H. GOLUB AND C. F. V. LOAN, *Matrix Computations*, John Hopkins University Press, 1996.
- [17] V. A. HAJIVASSILIOU AND D. L. MCFADDEN, *The method of simulated scores for the estimation of ldv models*, Econometrica, 66 (1998), pp. 863–896.
- [18] V. A. HAJIVASSILIOU, D. L. MCFADDEN, AND P. RUUD, *Simulation of multivariate normal rectangle probabilities and their derivatives: Theoretical and computational results*, Journal of Econometrics, 72 (1996), pp. 85–134.
- [19] E. HEINZ, *Beiträge zur Störungstheorie der Spektralzerlegung*, Math. Ann., 123 (1951), pp. 415–438.
- [20] N. J. HIGHAM, *Functions of Matrices: Theory and Computation*, Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2008.
- [21] J. HOFBAUER AND W. H. SANDHOLM, *On the global convergence of stochastic fictitious play*, Econometrica, 70 (2002), pp. 2265–2294.
- [22] R. A. HORN AND C. R. JOHNSON, *Matrix Analysis*, Cambridge University Press, Cambridge, 1985.
- [23] A. IUSEM, *On the convergence properties of the projected gradient method for convex optimization*, Computational & Applied Mathematics, 22 (2003), pp. 37–52.
- [24] M. P. KEANE, *Four Essays in Empirical Macro and Labor Economics*, Brown University, 1990.
- [25] M. P. KEANE, *A computationally practical simulation estimator for panel data*, Econometrica, 62 (1994), pp. 95–116.
- [26] C. KENNEY AND A. J. LAUB, *Condition estimates for matrix functions*, SIAM Journal on Matrix Analysis and Applications, 10 (1989), pp. 191–209.
- [27] K. LÖWNER, *Über monotone matrixfunktionen*, Mathematische Zeitschrift, 38 (1934), pp. 177–216.
- [28] R. D. LUCE, *Individual choice behavior*, John Wiley, New York, 1959.
- [29] D. MCFADDEN, *Conditional logit analysis of qualitative choice behavior*, in *paul zarembka, editor.*, Frontiers in Econometrics, (1974), pp. 105–142.
- [30] D. MCFADDEN, A. KARLQVIST, L. LUNDQVIST, F. SNICKARS, AND E. J. WEIBULL, *Modelling the choice of residential location*, Proc. Spatial Interaction Theory Planning Models, (1978), pp. 75–96.
- [31] D. MCFADDEN AND K. TRAIN, *Mixed mnl models for discrete response*, Journal of Applied Econometrics, 15 (2000), pp. 447–470.
- [32] V. MISHRA, K. NATARAJAN, H. TAO, AND C.-P. TEO, *Choice prediction with semidefinite optimization when utilities are correlated*, IEEE Transactions on Automatic Control, 57 (2012), pp. 2450–2463.

- [33] K. NATARAJAN, M. SONG, AND C.-P. TEO, *Persistency model and its applications in choice modeling*, Management Science, 55 (2009), pp. 453–469.
- [34] K. NATARAJAN AND C.-P. TEO, *On reduced semidefinite programs for second order moment bounds with applications*, To appear in Mathematical Programming, (2016).
- [35] I. OLKIN AND F. PUKELSHEIM, *The distance between two random vectors with given dispersion matrices*, Linear Algebra and its Applications, 48 (1982), pp. 257 – 263.
- [36] G. PARSONS AND M. KEALY, *Randomly drawn opportunity sets in a random utility model of lake recreation*, Land Economics, 68 (1992), pp. 93–206.
- [37] A. SHAPIRO, *Extremal problems on the set of nonnegative definite matrices*, Linear Algebra and its Applications, 67 (1985), pp. 7 – 18.
- [38] T. J. STEENBURGH, *The invariant proportion of substitution property (ips) of discrete-choice models*, Marketing Science, 27 (2008), pp. 300–307.
- [39] K. C. TOH, M. J. TODD, AND R. H. TUTUNCU, *Sdpt3 - a matlab software package for semidefinite programming*, Optimization Methods and Software, 1 (1999), pp. 545–581.
- [40] O. TOUBIA, D. I. SIMESTER, J. R. HAUSER, AND E. DAHAN, *Fast polyhedral adaptive conjoint estimation*, Marketing Science, 22 (2003), pp. 273–303.
- [41] K. TRAIN, *Discrete choice methods with simulation*, (2009).
- [42] R. H. TUTUNCU, K. C. TOH, AND M. J. TODD, *Solving semidefinite-quadratic-linear programs using sdpt3*, Mathematical Programming Series B, 95 (2003), pp. 189–217.
- [43] F. VERBOVEN, *The nested logit model and representative consumer theory*, Economics Letters, 50 (1996), pp. 57–63.
- [44] J. H. WILKINSON, *The algebraic eigenvalue problem*, vol. 87, Clarendon Press Oxford, 1965.
- [45] L. Q. YANG, D. SUN, AND K. C. TOH, *Sdpnal+: a majorized semismooth newton-cg augmented lagrangian method for semidefinite programming with nonnegative constraints*, Mathematical Programming Computation, 7 (2015), pp. 331–366.
- [46] X.-Y. ZHAO, D. SUN, K. C. TOH, AND M. J. TODD, *A newton-cg augmented lagrangian method for semidefinite programming*, SIAM Journal of Optimization, 20 (2010), pp. 1737–1765.

Selin Damla Ahipasaoglu Selin Damla Ahipasaoglu is an Assistant Professor at the Engineering Systems and Design Pillar at the Singapore University of Technology and Design. Her research interests lie in the connections between robust optimization and discrete choice, experimental design and statistical learning and smart grids. She obtained her PhD from Cornell University.

Xiaobo Li Xiaobo Li is currently a PhD student at the Industrial and Systems Engineering Department at the University of Minnesota. His research interests lie in discrete choice models, distributional robust optimization, online learning and the shared economy.

Karthik Natarajan Karthik Natarajan is an Associate Professor at the Engineering Systems and Design Pillar at the Singapore University of Technology and Design. His research interests lie in optimization under uncertainty with a focus on distributionally robust optimization and applications to areas such as marketing, finance, supply chain management and engineering applications. He obtained his PhD from the Singapore-MIT Alliance Program, National University of Singapore.