

Mixed Logit and Pure Characteristics Models*

Jay Lu[†]

Kota Saito[‡]

March 2021

Abstract

Mixed logit or random coefficients logit models are used extensively in empirical work while pure characteristic models feature in much of theoretical work. We provide a theoretical analysis comparing and contrasting the two classes of models. First, we show an approximation theorem that precisely characterizes the extent to which mixed logit models can approximate pure characteristic models. In the process, we introduce a general class of models that corresponds exactly to the closure of logit models. We then present two conditions that highlight behavioral differences between mixed logit and pure characteristic models. Both pertain to choice patterns relating to product differentiation. The first is a substitutability condition that is satisfied by many pure characteristic models (including the Hotelling model of horizontal differentiation) but is violated by almost all mixed logit models. The second is a continuity condition that is satisfied by all pure characteristic models but is violated by all mixed logit models.

* We thank Giovanni Compiani, Yusuke Narita, Matt Shum and Yi Xin for helpful comments. Financial support from the NSF under awards SES-1558757 (Saito), SES-1919263 (Saito) and SES-1919275 (Lu) are gratefully acknowledged.

[†] Department of Economics, UCLA; jay@econ.ucla.edu.

[‡] Division of the Humanities and Social Sciences, Caltech; saito@caltech.edu.

1 Introduction

Mixed logit models, also known as random coefficients logit models, have been widely used in empirical work across different fields (McFadden (1973), Rust (1987) and Berry, Levinsohn and Pakes (1995)). In these models, agents' utilities contain iid extreme-value distributed error terms that generate convenient expressions for choice probabilities which are useful for estimation. On the other hand, much of the theoretical literature in decision theory and industrial organization since Hotelling (1929) have focused on pure characteristic models (Berry and Pakes (2007)). In these models, there are no iid error terms and utilities are purely continuous functions of product characteristics. In this paper, we provide a theoretical analysis comparing and contrasting these two classes of models.

Our main contribution is two-fold. First, we provide an approximation theorem that precisely characterizes the extent to which mixed logit models can approximate pure characteristic models. This sharpens existing approximation results (e.g. McFadden and Train (2000)) and shows that a pure characteristic model can be approximated by mixed logit models if and only if they belong to the same parametric family. Second, we highlight some inherent differences between the two classes of models. We study two patterns of choice behavior related to product differentiation that are natural in pure characteristic models but cannot be easily accommodated by mixed logit models. These results demonstrate that while mixed logit models are flexible enough to accommodate a wide range of behaviors, they also impose certain limitations that impact estimation and counterfactual analysis.¹

In our model, each choice option (i.e. product) corresponds to a vector in \mathbb{R}^k where k is the number of characteristics. We assume there is rich set of products with varying characteristics. Following most empirical work, we focus on parametric families of polynomials up to some degree $d \geq 1$. For example, polynomials of degree 1 correspond to all linear functions $u(x) = \beta \cdot x$. The main approximation result (Theorem 1) states that a pure characteristic model of degree d can and only can be approximated by mixed logits of degree d . For instance, if the pure characteristic has degree d , then it is in general impossible to approximate the model using mixed logits of degree $d' < d$. In practice, this means that specifying the correct specification of the degree of the parametric family of utilities is important for mixed logit approximations.

¹ Our results complement an empirical literature on well-known issues with mixed logit models. See the discussion on related literature.

In the process of showing this result, we also characterize the universal set of all models that can be approximated by logit and mixed logit models. These results may be of independent interest to researchers. The closure of logit is a class of models which we call *lexicographic-logit*; this is a lexicographic choice rule with logit tie-breaking. Lexicographic-logit is an exceptionally rich class of models and includes some models (e.g. lexicographic choice) that cannot be expressed as random utilities.

To see how this is useful for the approximation theorem, consider any pure characteristic model that can be approximated by mixed logits of degree $d = 1$, i.e. linear utilities $u(x) = \beta \cdot x$. This must be a mixture of lexicographic-logits where the lexicographic preferences correspond to linear utilities. Since utilities in a pure characteristic model are continuous functions of product characteristics, logit tie-breaking can never occur as they are discontinuous. This means that only the lexicographic utilities remain and this can be rewritten as a pure characteristic model with linear ($d = 1$) utilities.

While mixed logit models are flexible and can approximate any pure characteristic model, we show that there are inherent differences between the two classes of models. We highlight these differences using two conditions on patterns in choice behavior. The first condition we focus on is called *convex substitutability*. Consider two products x and y and a third product $z = \lambda x + (1 - \lambda)y$ with intermediary characteristics of the other two. Convex substitutability says that the demand for x decreases if we replace y with z . The intuition is that since x is more similar to z than to y , agents will substitute away from x when the more similar product is introduced. This is natural condition that is satisfied by many pure characteristic models, including the classic Hotelling model; however, it is violated by all mixed logit models excepting the special case of uniform choice (Theorem 2).

The second condition is called *continuity in characteristics*. Consider two products x and y and series of products y_n with characteristics that converge to those of y , i.e. $y \rightarrow y_n$ in \mathbb{R}^k . Continuity in characteristics says that the demand for x when both y and y_n are available will eventually converge to the demand for x when only y is available. The intuition is that agents will eventually be unable to distinguish between y and y_n and treat both as the same product. This condition is satisfied by all pure characteristic models but violated by all mixed logit models (Theorem 3).

These conditions show stark differences in choice behavior between mixed logit and pure characteristic models. How significant are these differences for practical purposes? On the one hand, the discrepancies will eventually vanish as mixed logit approximations get ar-

bitrarily close to the pure characteristic model. On the other hand, any mixed logit that eventually emerges from estimation will violate these conditions. Since both conditions pertain to product differentiation, this would complicate counterfactual analysis when product characteristics vary or new products are substituted. The significance and magnitude of these violations will depend on the specific application, but they highlight potential issues to consider when using mixed logit approximations.

We focus on the mixed logit due to its prominence in applied work but our results extend to more generally to a larger class of models. For example, all models with iid error terms would have difficulty satisfying convex substitutability and continuity in characteristics. Importantly, not all classes of models useful for estimation necessarily violate these conditions; we provide an example of a continuous probit model that satisfies both (see Example 7). While iid error terms are useful for modeling unobserved heterogeneity, their convenience imposes restrictions on choice behavior that may be undesirable. When deciding whether to use one class of models versus another, one would need to weigh the importance of adhering to certain choice behaviors with the burden of computational costs.

1.1 Related Literature

Luce (1959) provided an early characterization of multinomial logit. Recent papers in decision theory have considered generalizations of logit. These include mixed logit (Gul, Natenzon, and Pesendorfer (2014), Saito (2018)) and nested logit (Kovach and Tserenjigmid (2020)). Cerreia-Vioglio, Maccheroni, Marinacci and Rustichini (2018a; 2018b) consider the Luce axiom without positivity and obtain a model that is a discrete version of our lexicographic-logit model. Fudenberg and Strzalecki (2015) consider dynamic extensions of logit. Natenzon (2019) studies a Bayesian probit model. Chambers, Cuhadaroglu and Masatlioglu (2020) consider a variation of the logit model in a social setting.

Theoretical work in decision theory has focused on pure characteristic models. These include Gul and Pesendorfer (2006), Ahn and Sarver (2013), Lu (2016), Apestegua, Ballester and Lu (2017), Lu and Saito (2018), Duraj (2018), Frick, Iijima and Strzalecki (2019), Lu (2019) and Lin (2019). Wilcox (2011) and Apestegua and Ballester (2018) discuss issues with respect to comparative statics between logit and pure characteristic models while Frick, Iijima, and Strzalecki (2019) highlight issues associated with assessing option values. These results are similar in spirit to ours highlighting the differences in choice behavior between

the two class of models.

In the empirical literature, logit-based models have been widely applied for discrete choice analysis. These include McFadden (1973), Rust (1987), Hotz and Miller (1993), Berry, Levinsohn and Pakes (1995), Nevo (2001), Hendel and Nevo (2006), Gowrisankaran and Rysman (2012) and Compiani (2019). McFadden and Train (2000) show that that mixed logit models can approximate any pure characteristic model. Our result sharpens their result and captures the precise extent of this approximation. Narita and Saito (2021) consider the case where the set of characteristics is finite. They provide a necessary and sufficient condition for when mixed logit can approximate random utility models. They also provide an algorithm for constructing mixed logit models that can approximate random utility arbitrarily well when the condition is satisfied; when the condition is not satisfied, they find that the size of the approximation error is large.

Other papers comparing logit-based models with pure characteristics models (also known as hedonic models) include Anderson, DePalma and Thisse (1989), Petrin (2002), Bajari and Benkard (2004; 2005), Akerberg and Rysman (2005) and Berry and Pakes (2007). They consider the implications on price elasticities and welfare when new products are introduced. Logit-based models may imply too much “taste for product” while pure characteristics models may imply competition that is too localized. Our results on convex substitutability and continuity in characteristics are similar in spirit and highlight new patterns in choice behavior differentiating the two classes of models.

2 Setup

There are $k \geq 1$ characteristics and we associate each choice option (i.e. product) with a vector $x \in X \subset \mathbb{R}^k$ of characteristics. We assume that there is rich variation in the set of characteristics. Formally, X is full-dimensional, compact and convex. A *menu* $A \subset X$ is a finite set of products and let \mathcal{A} denote the set of all menus. A *stochastic choice* ρ is a mapping on \mathcal{A} such that for any menu $A \in \mathcal{A}$, $\rho_A(\cdot)$ is a probability distribution over elements in A . For binary menus, $A = \{x, y\}$, we use the simpler notation $\rho(x, y) = \rho_A(x)$. The set of all stochastic choice can be thus defined as

$$\mathcal{P} = \prod_{A \in \mathcal{A}} \Delta A$$

where ΔA is the set of all probability distributions over A . We endow \mathcal{P} with the product topology.

We focus on utilities that are continuous in product characteristics. Formally, let $U \subset \mathbb{R}^X$ denote the set of all continuous utility functions $u : X \rightarrow \mathbb{R}$. We say $u \in U$ is a *polynomial of degree $d \geq 1$* if it is a multivariate polynomial where every term has exponents that sum up to at most d , i.e.

$$u(x_1, \dots, x_k) = \sum_{m_1^i + \dots + m_k^i \leq d} \beta_i x_1^{m_1^i} \dots x_k^{m_k^i}$$

for some $\beta \in \mathbb{R}^{m(k,d)}$ where $m(k,d) = \sum_{i=1}^d \frac{(k+i-1)!}{i!(k-1)!}$. For example, if $k = d = 2$, then

$$u(x_1, x_2) = \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_2^2 + \beta_5 x_1 x_2$$

In general, β represents the coefficients for different characteristics including higher-order and interaction terms. Constant terms β_0 are excluded as they have no bearing on any of the subsequent analysis. Given any $x \in X$, let x^* denote the corresponding vector in polynomial space X^* so

$$u(x) = \beta \cdot x^*.$$

Note that if $d = 1$, then $x = x^*$ and $u(x) = \beta \cdot x$ is just a linear function. Let $U_d \subset U$ denote the set of all polynomials of degree d .

Two prominent classes of models are pure characteristic and mixed logit.

Definition 1. (Pure characteristic) A stochastic choice ρ is *pure characteristic* if there exists a distribution μ on U such that

$$\rho_A(x) = \mu(\{u \in U : u(x) \geq u(y) \text{ for all } y \in A\}).$$

It is *pure characteristic of degree d* if $u \in U_d$ a.s.

Definition 2. (Mixed logit) A stochastic choice ρ is *mixed logit* if there exists a distribution ν on U such that

$$\rho_A(x) = \int_U \frac{e^{v(x)}}{\sum_{y \in A} e^{v(y)}} d\nu$$

It is *mixed logit of degree d* if $v \in U_d$ a.s.

A special case of mixed logit is of course when the distribution $\mu = \delta_u$ is degenerate. In this case, we obtain standard logit where v is its *systematic utility*. We summarize this

below.

Definition 3. (Logit) A stochastic choice ρ is *logit* if there exists a $v \in U$ such that

$$\rho_A(x) = \frac{e^{v(x)}}{\sum_{y \in A} e^{v(y)}}$$

It is *logit of degree d* if $v \in U_d$.

Both pure characteristic and mixed logit belong to a more general class of random utility models. To see why, note that we can rewrite the mixed logit model as

$$\rho_A(x) = \mathbb{P}(\{v(x) + \varepsilon(x) \geq v(y) + \varepsilon(y) \text{ for all } y \in A\})$$

where $v(\cdot)$ are distributed according to ν and $\varepsilon(\cdot)$ are extreme-valued distributed and iid across products. On the other hand, the pure characteristic model is a random utility where the utilities $u(\cdot)$ are continuous in product characteristics.

3 Mixed Logit Approximations of Pure Characteristic Models

3.1 An Approximation Theorem

This section provides a precise characterization of the extent to which mixed logit models can be used to approximate any pure characteristic model. McFadden and Train (2000) showed that mixed logit models can be used to approximate any pure characteristic model. We translate their result to our setup. Recall that \mathcal{P} is endowed with the product topology so $\rho^n \rightarrow \rho$ iff $\rho_A^n \rightarrow \rho_A$ for all $A \in \mathcal{A}$.

Proposition 1. *For any pure characteristic ρ , there exists a sequence of mixed logits ρ^n such that $\rho^n \rightarrow \rho$.*

Proof. See Appendix. □

While mixed logit models owe much of their popularity to their computability, the above result also provides some theoretical justification for their use. In practice however, a researcher will usually commit to some class of mixed logit models for approximation. For example, suppose the researcher uses mixed logit of degree $d = 1$, i.e. mixed logit models where the systematic utility $u(x) = \beta \cdot x$ is linear. What is the set of all pure characteristic

models that this can approximate? While this clearly includes pure characteristic models with linear utilities, could it include more? It turns out the answer is no.

The next result provides a precise characterization of the extent to which mixed logit models can be used to approximate pure characteristic models. We say a stochastic choice ρ can be *approximated* by a set of stochastic choice models if it is in the closure of that set.²

Theorem 1. *For any pure characteristic ρ , the following are equivalent:*

- (1) ρ is pure characteristic of degree d
- (2) ρ can be approximated by mixed logit of degree d .

Proof. See Appendix. □

The main implication of this result is that it is important to correctly specify the degree of the pure characteristic model. For example, suppose the pure characteristic model has degree 3 where third-order terms matter. In this case, it would be generally impossible to approximate the pure characteristic model if the researcher only uses mixed logits of degree $d < 3$. A special case when a lower-degree mixed logit would suffice is if the utilities in the pure characteristic model is exactly a monotone transformation of a lower-degree polynomial. For instance, if the utilities in the pure characteristic model satisfy

$$u(x) = \beta_1 x_1^3 + \beta_2 x_1^2 x_2 + \beta_3 x_1 x_2^2 + \beta_4 x_2^3 = (\gamma_1 x_1 + \gamma_2 x_2)^3$$

where $\beta_1 = \gamma_1^3$, $\beta_2 = 3\gamma_1^2\gamma_2$, $\beta_3 = 3\gamma_1\gamma_2^2$ and $\beta_4 = \gamma_2^3$. In this case, mixed logits of degree 1 can be used to approximate this model.

Theorem 1 sharpens the result from McFadden and Train (2000) in two ways. First, while they only provide sufficiency, we provide necessity as well as demonstrated in the example above. Second, our result is more precise about the class of utility functions that can be approximated (i.e. the degree of the polynomial utility). This has a practical importance for empirical analysis given that the majority of empirical work assumes $d = 1$.

3.2 Closure of Logit: Lexicographic-logit

In order to demonstrate the reasoning behind Theorem 1, we first characterize the universal set of all models that can be approximated by (mixed) logit. This characterization may

² Recall that closure here is with respect to the product topology on \mathcal{P} .

be of independent interest to researchers. It shows the full extent in which (mixed) logit models can be used to approximate a rich class of models including some that are not pure characteristic or even random utility. We begin with an example of a model that is the limit of logits.

Example 1 (*Lexicographic choice rule*). Let $X = [0, 1]^2 \subset \mathbb{R}^2$ and $d = 1$. Let ρ^n be logit with $\beta_n = (n^2, n)$ and ρ be the limit of the logits ρ^n . Thus,

$$\rho_A(x) = \lim_n \frac{e^{\beta_n \cdot x}}{\sum_{y \in A} e^{\beta_n \cdot y}} = \left(\sum_{y \in A} e^{\lim_n (n^2(y_1 - x_1) + n(y_2 - x_2))} \right)^{-1}$$

In this case, ρ corresponds to the lexicographic preference \succ on X where $x \succ y$ if $x_1 > y_1$ or $x_1 = y_1$ and $x_2 > y_2$. To see why, note that if $x_1 > y_1$ or $x_1 = y_1$ and $x_2 > y_2$, then

$$\lim_n (n^2(y_1 - x_1) + n(y_2 - x_2)) = -\infty$$

If this is true for all $y \in A$, then $\rho_A(x) = 1$. On the other hand, if $y \succ x$ for some $y \in A$, then $\rho_A(x) = 0$ as desired. Thus, ρ is a lexicographic choice rule.

The above example shows how lexicographic choice rules is one class of models that can be approximated by logit models. Note that in that example, ρ is not pure characteristic or even a random utility. To see why, suppose otherwise and ρ is a random utility model with some distribution over all utility functions $u : X \rightarrow \mathbb{R}$. Thus, for each distinct $x, y \in X$, we have $x \succ y$ if $u(x) > u(y)$ with probability one. If we define the average utility $\bar{u}(x) := \mathbb{E}[u(x)]$, then \bar{u} represents \succ . This yields a contradiction since it is well-known that no utility representation exists for lexicographic preferences. Thus, although every logit is a random utility, the closure of logit includes models that cannot be expressed as a random utility.³

While mixed logit can approximate pure characteristic models and lexicographic choice rules, what is the full class of models that can be approximated? We now characterize that set. Consider a collection of polynomials (u_1, \dots, u_t) where $u_i \in U_d$ for all $i \in \{1, \dots, t\}$. Let $(\beta_1, \dots, \beta_t)$ be their corresponding coefficients so $u_i(x) = \beta_i \cdot x^*$ for all i . We say the collection is *orthogonal* if $\beta_i \cdot \beta_j = 0$ for all $i, j \in \{1, \dots, t\}$.

³ If we allow for only finite-additive distributions or non-measurable utilities, then one could represent lexicographic choice rules using some “random utility” (see Cohen (1980)). However, given that a lexicographic preference has no utility representation, it would be odd for it to have a random utility representation. Moreover, this would exclude the possibility of integrating utilities (e.g. calculating social surplus).

Given $\omega = (u_1, \dots, u_t)$, let \succeq_ω be its induced lexicographic preference relation on X . In other words, $x \sim_\omega y$ if $u_i(x) = u_i(y)$ for all $i \in \{1, \dots, t\}$ and $x \succ_\omega y$ if $u_i(x) > u_i(y)$ for some $i \leq t$ and $u_j(x) = u_j(y)$ for all $j < i$. Let Ω_d be the set of all orthogonal polynomials $\omega = (u_1, \dots, u_t)$ for some $t \leq m(k, d)$. Under lexicographic-logit, choices follow a lexicographic preference relation where ties are broken according to logit.

Definition 4. ρ is *lexicographic-logit of degree d* if there exist $\omega \in \Omega_d$ and $v \in U_d$ such that

$$\rho_A(x) = 1 \{x \succeq_\omega y \text{ for all } y \in A\} \frac{e^{v(x)}}{\sum_{y \in A, y \sim_\omega x} e^{v(y)}}$$

The following result shows that the closure of logit is exactly lexicographic-logit.

Proposition 2. *The following are equivalent:*

- (1) ρ is *lexicographic-logit of degree d*
- (2) ρ can be approximated by logit of degree d .

Proof. See Appendix. □

How about for mixed logit? We now define the mixed lexicographic-logit model.

Definition 5. ρ is *mixed lexicographic-logit of degree d* if there exists a distribution ν on $\Omega_d \times U_d$ such that

$$\rho_A(x) = \int_{\Omega_d \times U_d} 1 \{x \succeq_\omega y \text{ for all } y \in A\} \frac{e^{v(x)}}{\sum_{y \in A, y \sim_\omega x} e^{v(y)}} d\nu$$

The following result parallels Proposition 2 and shows that mixed lexicographic-logit is exactly the set of all stochastic choices that can be approximated by mixed logit. In fact, it is the smallest set of models containing logit that is closed under mixing and approximations.

Proposition 3. *The following are equivalent:*

- (1) ρ is *mixed lexicographic-logit of degree d*
- (2) ρ can be approximated by mixed logit of degree d .

Proof. See Appendix. □

Mixed lexicographic-logit includes a rich class of models. When ties are universal, i.e. $x \sim_\omega y$ for all $x, y \in X$ a.s., this reduces to mixed logit. When ω only consists of a single polynomial and ties never occur, this reduces to pure characteristic. The following are a few additional special cases:

Example 2 (*Mixed lexicographic*). Let $\omega = (u_1, \dots, u_t)$ for $t > 1$ and suppose ties never occur, i.e. $x \not\sim_\omega y$ for all $x, y \in X$ a.s. This corresponds to a population of agents where each agent chooses according to a lexicographic preference. As special case of course is Example 1 above.

Example 3 (*Mixture of logit and pure characteristic*). Let ν_1 correspond to a mixed logit model, ν_2 correspond to a pure characteristic model and $\nu = \alpha\nu_1 + (1 - \alpha)\nu_2$ for $\alpha \in (0, 1)$. Here, α parametrizes the degree of iid noise in the model. Note that this is neither a mixed logit model nor a pure characteristic model.

Example 4 (*Generalized nested-logit*). Let $A = \{x, y, z\}$ and consider $u_1, u_2 \in U_d$ such that $u_1(x) = u_1(y) > u_1(z)$ and $u_2(x) < u_2(y) = u_2(z)$. Suppose ν is such that $(\omega, \nu) = (u_1, \nu_1)$ with probability α and $(\omega, \nu) = (u_2, \nu_2)$ with probability $1 - \alpha$. In this case,

$$\rho_A(x) = \alpha \frac{e^{v_1(y)}}{e^{v_1(y)} + e^{v_1(x)}} + (1 - \alpha) \frac{e^{v_2(y)}}{e^{v_2(y)} + e^{v_2(z)}}$$

This corresponds to an agent who either picks the “nest” $\{x, y\}$ in the first-stage followed by logit with ν_1 in the second stage or picks the “nest” $\{y, z\}$ followed by logit ν_2 . If we interpret the nests as consideration sets, then u_1 and u_2 correspond to salience measures.

Finally, we end this section with a brief outline of how Proposition 3 can be used to prove Theorem 1. Suppose ρ is a pure characteristic model that can also be approximated by mixed logits of degree $d = 1$. Proposition 3 implies that ρ is mixed lexicographic-logit of degree $d = 1$. We first show that $x \sim_\omega y$ can never occur with positive probability so logit tie-breaking never occurs. Let $y_n = \left(1 - \frac{1}{n}\right)y + \frac{1}{n}x$ and note that since $\omega \in \Omega_1$, $x \succeq_\omega y$ iff $x \succeq_\omega y_n$. Let $A_n = \{x, y, y_n\}$ and taking the limit as $n \rightarrow \infty$, we have

$$\begin{aligned} \lim_n \rho_{A_n}(x) &= \lim_n \int_{\Omega_1 \times \mathbb{R}^k} 1\{x \succeq_\omega y\} \frac{e^{\beta \cdot x}}{e^{\beta \cdot x} + 1\{y \sim_\omega x\}(e^{\beta \cdot y} + e^{\beta \cdot y_n})} d\nu \\ &= \int_{\Omega_1 \times \mathbb{R}^k} 1\{x \succeq_\omega y\} \frac{e^{\beta \cdot x}}{e^{\beta \cdot x} + 1\{y \sim_\omega x\} 2e^{\beta \cdot y}} d\nu \\ &\leq \int_{\Omega_1 \times \mathbb{R}^k} 1\{x \succeq_\omega y\} \frac{e^{\beta \cdot x}}{e^{\beta \cdot x} + 1\{y \sim_\omega x\} e^{\beta \cdot y}} d\nu = \rho(x, y) \end{aligned}$$

Since ρ is also pure characteristic, $\lim_n \rho_{A_n}(x) = \rho(x, y)$ (see Theorem 3.1 below) so it must be that $y \sim_\omega x$ with measure zero. We can thus write

$$\rho_A(x) = \nu(\{\omega \in \Omega_1 : x \succeq_\omega y \text{ for all } y \in A\})$$

It is straightforward to show that this satisfies the random linear utility axioms of Gul and Pesendorfer (2006) so ρ is pure characteristic of degree $d = 1$. The full proof of Theorem 1 extends this argument for $d > 1$.

4 Contrasting Mixed Logit and Pure Characteristic Models

This previous section focuses on the extent to which mixed logit models can be used to approximate any pure characteristic model. In this section, we focus on two behavioral differences between the two classes of models. Both pertain to patterns in choice behavior relating to product differentiation. Section 4.1 studies a substitutability condition that is satisfied by many pure characteristic models but is violated by almost all mixed logit models. Section 4.2 studies a continuity condition that is satisfied by all pure characteristic models but is violated by all mixed logit models. The main purpose of this section is to illustrate patterns of choice behavior that clearly separate mixed logit models from pure characteristic models.

4.1 Convex Substitutability

This section introduces a substitutability condition that is natural in many pure characteristic models but cannot be accommodated by almost all mixed logit models. To illustrate, consider an example from the classic Hotelling (1929) model of horizontal differentiation.

Example 5 (*Hotelling*). *Each choice alternative corresponds to a product $x = (\theta, p) \in \mathbb{R}^2$ where θ measures quality and p is its price. Each agent i has utility*

$$u_i(\theta, p) = \alpha_i - \lambda_i(\theta - \beta_i)^2 - p$$

This is a pure characteristic model where $\rho(x, y)$ is the demand of product x over y . Given $x = (\theta, p)$ and $y = (\theta', p')$, let $z = \frac{1}{2}x + \frac{1}{2}y$ denote an intermediary product that is a convex mixture of the characteristic and price of the two products. If an agent prefers x to z , then

he must also prefer x to y . To see why, note that

$$\begin{aligned} u_i(z) &= \alpha_i - \lambda_i \left(\frac{1}{2}\theta + \frac{1}{2}\theta' - \beta_i \right)^2 - \frac{1}{2}p - \frac{1}{2}p' \\ &\geq \alpha_i - \frac{1}{2}\lambda_i (\theta - \beta_i)^2 - \frac{1}{2}\lambda_i (\theta' - \beta_i)^2 - \frac{1}{2}p - \frac{1}{2}p' \\ &\geq \frac{1}{2}u_i(x) + \frac{1}{2}u_i(y) \end{aligned}$$

so $u_i(x) \geq u_i(z)$ implies $u_i(x) \geq u_i(y)$. This means that the demand of x will increase if we replace z with y , i.e.

$$\rho\left(x, \frac{1}{2}x + \frac{1}{2}y\right) \leq \rho(x, y).$$

Intuitively, the product z is more “similar” (in a convex sense) to x than y is to x and demand for x increases if we replace z with y .

We can generalize this substitutability condition as follows.

Definition 6. ρ satisfies *convex substitutability* if $\rho_{A \cup \{y\}}(x) \geq \rho_{A \cup \{\lambda x + (1-\lambda)y\}}(x)$ for all $\lambda \in (0, 1)$.

Convex substitutability says that demand for x increases when the other alternatives become less similar. Intuitively, this is because alternatives that are similar serve as substitutes. Note that similarity here is measured in terms of convexity in the space of characteristics \mathbb{R}^k .

In the Hotelling model, convex substitutability is satisfied because all utilities are concave. In fact, convex substitutability is satisfied as long as all the utilities in the pure characteristic model are quasiconcave.⁴ We say a pure characteristic model is *quasiconcave* if its utility functions $u : X \rightarrow R$ are quasiconcave a.s.

Theorem 2.1. *Any quasiconcave pure characteristic ρ satisfies convex substitutability.*

Proof. Consider distinct $x, y \in X$ and let $z = \lambda x + (1 - \lambda)y$. Since u is quasiconcave, if $u(x) > u(z)$, then $u(x) \geq u(y)$. Since ties occur with measure zero, this means that

$$\begin{aligned} \rho_{A \cup \{z\}}(x) &= \mu(\{u \in U : u(x) > u(w) \text{ for all } w \in A \cup \{z\}\}) \\ &\leq \mu(\{u \in U : u(x) > u(w) \text{ for all } w \in A \cup \{y\}\}) = \rho_{A \cup \{y\}}(x) \end{aligned}$$

as desired. □

⁴ Convex substitutability is related to convexity conditions in Apesteguia, Ballester and Lu (2017) and Lu (2019).

Note that the proof solely relies on utilities being quasiconcave and holds even if utilities are not continuous. A special case of quasiconcave utility is of course linear utility, for example, the case of expected utility for choice under risk.

Example 6 (*Random expected utility*). Each alternative corresponds to a lottery $x \in \mathbb{R}_+^k$ over k prizes where $\sum_j x_j = 1$. Each agent i is an expected utility maximizer with utility

$$u_i(x) = \sum_j x_j \tilde{u}_{ij}$$

where \tilde{u}_{ij} is agent i 's Bernoulli utility of prize $j \in \{1, \dots, k\}$. This is the random expected utility model of Gul and Pesendorfer (2006). Since $u_i(x) \geq u_i(\lambda x + (1 - \lambda)y)$ iff $u_i(x) \geq u_i(y)$, convex substitutability is in fact satisfied with equality.

While convex substitutability is satisfied by many commonly used pure characteristic models, we now show that it cannot be accommodated by almost all mixed logit models. The one exception is when the stochastic choice is *uniform*, i.e. $\rho_A(x) = 1/|A|$,

Theorem 2.2. *Any non-uniform mixed logit ρ violates convex substitutability.*

Proof. Note that by applying convex substitutability repeatedly, we obtain that

$$\rho_A(x) \geq \rho_{\lambda x + (1-\lambda)A}(x)$$

for all $\lambda \in (0, 1)$. Since ρ is mixed logit, this means that

$$\int_U \frac{e^{v(x)}}{\sum_{y \in A} e^{v(y)}} d\nu = \rho_A(x) \geq \rho_{\lambda x + (1-\lambda)A}(x) = \int_U \frac{e^{v(x)}}{\sum_{y \in A} e^{v(\lambda x + (1-\lambda)y)}} d\nu$$

Now, taking the limit as $\lambda \rightarrow 1$, we have

$$\int_U \frac{e^{v(x)}}{\sum_{y \in A} e^{v(y)}} d\nu \geq \int_U \frac{e^{v(x)}}{n e^{v(x)}} d\nu = \frac{1}{n}$$

and this is true for any $x \in A$. Now, suppose the inequality is strict for some $x \in A$. Then

$$1 = \sum_{x \in A} \rho_A(x) = \sum_{x \in A} \left(\int_U \frac{e^{v(x)}}{\sum_{y \in A} e^{v(y)}} d\nu \right) > \sum_{x \in A} \frac{1}{n} = 1$$

yielding a contradiction. Thus, ρ must be uniform choice. □

Convex substitutability can be attractive for various reasons. Normatively, if one believes that the Hotelling model for instance is the true model, then agents should exhibit such

behavior. Descriptively, convex substitutability has certain appeal as it captures the intuitive notion that similar products crowd out demand. The above result shows that barring uniform choice, no mixed logit model can accommodate this condition. This is a novel restriction on allowable patterns of choice behavior under mixed logit.

What are the practical implications? Consider a researcher who is using mixed logit models to approximate the Hotelling model. Theorem 2.2 says that all models used for approximation will violate a property that is natural in the Hotelling model. How significant is this violation? On the one hand, this discrepancy will eventually vanish as approximations get arbitrarily close. On the other hand, any mixed logit that eventually emerges from estimation will violate convex substitutability. This would complicate counterfactual about what would happen to demand as alternatives become more or less similar. Ultimately, the significance and magnitude of this violation will depend on the specific application, but Theorem 2.2 highlights a potential issue for mixed logit approximations.

Violations of convex substitutability extend more generally to models beyond mixed logit. In fact, any model with additive iid error terms would have difficulty satisfying the condition. To illustrate, consider a model where $X \subset \mathbb{R}$ and $\varepsilon(\cdot)$ represents iid error terms. The probability of choosing x over y is

$$\begin{aligned} \rho(x, y) &= \mathbb{P}\{v(x) + \varepsilon(x) \geq v(y) + \varepsilon(y)\} \\ &= \mathbb{P}\{Z \leq v(x) - v(y)\} \\ &= F(v(x) - v(y)) \end{aligned}$$

where Z is the distribution of $\varepsilon(y) - \varepsilon(x)$ with cdf F . This is known as a Fechnerian model in the literature (Debreu (1958), Davidson and Marschak (1959)). Assuming differentiability, we have

$$\frac{\partial}{\partial y} \rho(x, y) = -f(v(x) - v(y)) v'(y)$$

which is positive iff $v'(y) \leq 0$. Since this does not depend on whether $x > y$ or $x < y$, convex substitutability cannot be satisfied.

Can convex substitutability be satisfied when error terms are not iid? The following shows a continuous version of probit where error terms are correlated that satisfies convex substitutability.

Example 7 (*Continuous probit*). Consider a continuous version of the probit model from

Hausman and Wise (1978). Let $X \subset \mathbb{R}$ and consider a pure characteristic model where

$$u(x) = -x^2 + W(x)$$

where W is a one-dimensional Brownian motion. For $x < y$, we have $u(x) \geq u(y)$ iff

$$y^2 - x^2 \geq W(y) - W(x)$$

$$y^2 - x^2 \geq (y - x)Z$$

$$x + y \geq Z$$

where Z is the standard Normal distribution. Clearly, as y increases, $\rho(x, y)$ also increases. On the other hand, for $y < x$, $u(x) \geq u(y)$ iff $x + y \leq Z$ so $\rho(x, y)$ decreases as y increases. It is easy to see that convex substitutability is satisfied.

Error terms in the continuous probit model are correlated depending on how similar alternatives are. By decreasing the variance of these error terms, continuous probit models can be used to approximate the Hotelling model. In contrast to mixed logit models, this has the advantage that convex substitutability will be satisfied for all models along the path of approximation and depending on the specific application, this could be a desirable feature.

4.2 Continuity in Characteristics

This section introduces a continuity condition that is satisfied by all pure characteristic models but violated by all mixed logit models. To illustrate, first consider the classic red-bus/blue-bus example (or Debussy versus Beethoven as formulated by Debreu (1960)).

Example 8 (Red-bus/blue-bus). Consider the choice of transportation alternatives and let x correspond to traveling by car while y correspond to traveling by a red bus. Suppose car and red bus both have equal market share so $\rho(x, y) = \frac{1}{2}$. Consider introducing a blue bus y' . Supposing agents are indifferent to color, $\rho(y, y') = \frac{1}{2}$ so Luce's independence of irrelevant alternatives (IIA) condition implies that $\rho_{\{x, y, y'\}}(x) = \frac{1}{3}$. In reality, one would expect the car market share to remain close to 50% in violation of IIA.

The red-bus/blue-bus example illustrates a limitation of logit. Mixed logit models are not bound by IIA and can accommodate such choice patterns. However, we now present a new choice pattern that no mixed logit model can accommodate. Suppose we introduce buses y_n

with colors that are increasingly closer to red. Eventually, y_n will be indistinguishable from y , so the car market share should approach 50%. In other words,

$$\rho_{\{x,y,y_n\}}(x) \rightarrow \rho(x,y)$$

We generalize this continuity condition as follows.

Definition 7. ρ satisfies *continuity in characteristics* if $\rho_{A \cup \{y_n\}}(x) \rightarrow \rho_A(x)$ for all $y_n \rightarrow y \in A \setminus \{x\}$.

Intuitively, as alternative y_n becomes increasingly similar to alternative y , the two alternatives will eventually be indistinguishable. When evaluating choice, one can replace two indistinguishable alternatives with the single alternative y in the limit.

In a pure characteristic model, the probability of choosing x over y is given by the probability that the utility of x is greater than that of y . Since utilities are continuous, the utility of y_n will converge to the utility of y as y_n converges to y . As a result, continuity in characteristics will be satisfied.

Theorem 3.1. *Any pure characteristic ρ satisfies continuity in characteristics.*

Proof. See Appendix. □

This applies to the Hotelling model in Example 5 and the random expected utility model in Example 6. It would also apply in models where the set of characteristics is not convex (e.g. Salop (1979)) as long as utilities are continuous in product characteristics. While continuity in characteristics is satisfied by all pure characteristic models, it cannot be accommodated by any mixed logit model.

Theorem 3.2. *Any mixed logit ρ violates continuity in characteristics.*

Proof. Consider distinct $x, y \in A$ and a sequence $y_n \rightarrow y$. Let $A_n = A \cup \{y_n\}$. Since ρ is a mixed logit,

$$\begin{aligned} \lim_n \rho_{A_n}(x) &= \lim_n \int_U \frac{e^{v(x)}}{\sum_{z \in A} e^{v(z)} + e^{v(y_n)}} d\nu \\ &= \int_U \frac{e^{v(x)}}{\sum_{z \in A} e^{v(z)} + e^{v(y)}} d\nu \\ &< \int_U \frac{e^{v(x)}}{\sum_{z \in A} e^{v(z)}} d\nu = \rho_A(x) \end{aligned}$$

Thus, $\lim_n \rho_{A_n}(x) < \rho_A(x)$ for any such A and $y_n \rightarrow y$. □

The above proof illustrates that any mixed logit model not only violates continuity in characteristics but it violates it in a specific direction. When y_n converges to y , the market share of x in $\{x, y, y_n\}$ converges to a limit that is strictly less than its market share in $\{x, y\}$. The reason is that logit errors force individual market shares for alternatives no matter how similar they are to each other. This is related to well-known limitations of mixed logit models (e.g. Berry and Pakes (2007)) and Theorem 3.2 formalizes such intuition.

The above results imply the following corollary.

Corollary 1. *No mixed logit model is pure characteristic.*

Proof. Follows from Theorem 3. □

Mixed logit and pure characteristic models belong to two very different classes of models; Corollary 1 shows that they have empty intersection. Although Theorem 1 guarantees that a researcher can always approximate pure characteristic models using mixed logit models, this approximation will always be from “outside” the set of pure characteristic models.

What does this mean in practice? Like the results from Section 4.1, the magnitude of these issues depend on the application at hand. Note that the continuous probit model in Example 7 is a pure characteristic model and thus satisfies continuity in characteristics. Theorem 3 illustrate one behavioral condition that separates the two class of models but there may be other differences that would have non-trivial implications for estimation and counterfactual analysis. Ultimately, when deciding whether to use one class of models versus another for estimation, one would need to weigh the importance of these choice patterns versus the burden of computational costs.

Appendix

A Preliminaries

Define the space

$$\mathcal{V} := \prod_{A \in \mathcal{A}} \mathbb{R}^A$$

and note that $\mathcal{P} \subset \mathcal{V}$. We endow \mathcal{V} also with the product topology. Note that \mathcal{P} is compact by Tychonoff's theorem. Since \mathcal{V} is a Hausdorff space (Theorem 19.4 of Munkres (2000)), \mathcal{P} is also closed (Lemma 2.32 of Aliprantis and Border (2006), henceforth AB). Although the space \mathcal{V} is neither metrizable or even first-countable, the next lemma shows that it is locally convex which allows us to use separating hyperplane theorems.

Lemma 1. *\mathcal{V} is a locally convex topological vector space.*

Proof. Since \mathbb{R}^A is a topological vector space for every $A \in \mathcal{A}$, \mathcal{V} is a topological vector space by Theorem 5.2 of AB. Consider the family of seminorms $(r_A)_{A \in \mathcal{A}}$ where $r_A : \mathcal{V} \rightarrow \mathbb{R}$ is such that

$$r_A(\tau) = |\tau_A|$$

where $|\cdot|$ is the Euclidean norm in \mathbb{R}^A . Since this family of seminorms generates the product topology, \mathcal{V} is locally convex. \square

Throughout this appendix, let $m = m(k, d)$, and for every $x \in X$, let x^* denote the corresponding vector in polynomial space $X^* \subset \mathbb{R}^m$. We define the following subsets of \mathcal{P} :

- \mathcal{P}^{pc} is the set of pure characteristic and \mathcal{P}_d^{pc} is the set of pure characteristic of degree d
- \mathcal{P}^{log} is the set of logit and \mathcal{P}_d^{log} is the set of logit of degree d
- \mathcal{P}^{mlog} is the set of mixed logit and \mathcal{P}_d^{mlog} is the set of mixed logit of degree d
- \mathcal{P}_d^{lex} is the set of lexicographic-logit of degree d
- \mathcal{P}_d^{mlex} is the set of mixed lexicographic-logit of degree d

Also let $cl(\mathcal{P}')$ denote the closure of any subset $\mathcal{P}' \subset \mathcal{P}$. Note that since \mathcal{P} is closed, $cl(\mathcal{P}') \subset \mathcal{P}$ is compact.

B Proof of Proposition 1

We first prove the following lemma.

Lemma 2. For $A \in \mathcal{A}$ and $u : X \rightarrow \mathbb{R}$ such that $u(x) \neq u(y)$ for all distinct $x, y \in A$,

$$\lim_n \frac{e^{nu(x)}}{\sum_{y \in A} e^{nu(y)}} = 1 \{u(x) > u(y) \text{ for all } y \in A\}$$

Proof. Since $u(x) \neq u(y)$ for all distinct $x, y \in A$, we can rewrite

$$\frac{e^{nu(x)}}{\sum_{y \in A} e^{nu(y)}} = \frac{1}{1 + \sum_{y \in A \setminus x} e^{n(u(y)-u(x))}}$$

Consider $n \rightarrow \infty$. First, suppose $u(x) > u(y)$ for all $y \in A \setminus x$. In this case, $e^{n(u(y)-u(x))} \rightarrow 0$ for all $y \in A \setminus x$ so the expression above converges to 1. On the other hand, suppose there exists some $y \in A$ such that $u(y) > u(x)$. In this case, $e^{n(u(y)-u(x))} \rightarrow \infty$ so the expression converges to 0. The result follows. \blacksquare

We now prove Proposition 1. Let $\rho \in \mathcal{P}^{pc}$ with distribution μ and define

$$\rho_A^n(x) = \int_U \frac{e^{nu(x)}}{\sum_{y \in A} e^{nu(y)}} d\mu$$

so $\rho^n \in \mathcal{P}^{mlog}$. We will show that $\rho^n \rightarrow \rho$. Fix some $A \in \mathcal{A}$ and define

$$U_A = \{u \in U : u(x) \neq u(y) \text{ for all distinct } x, y \in A\}$$

Since ties never occur for random utility models, $\mu(U_A) = 1$. Now, by Lemma 2 and dominated convergence,

$$\begin{aligned} \lim_n \rho_A^n(x) &= \lim_n \int_U \frac{e^{nu(x)}}{\sum_{y \in A} e^{nu(y)}} d\mu = \lim_n \int_{U_A} \frac{e^{nu(x)}}{\sum_{y \in A} e^{nu(y)}} d\mu \\ &= \int_{U_A} \lim_n \frac{e^{nu(x)}}{\sum_{y \in A} e^{nu(y)}} d\mu = \int_{U_A} 1 \{u(x) > u(y) \text{ for all } y \in A\} d\mu \\ &= \mu(\{u \in U_A : u(x) > u(y) \text{ for all } y \in A\}) \\ &= \mu(\{u \in U : u(x) \geq u(y) \text{ for all } y \in A\}) = \rho_A(x) \end{aligned}$$

as desired.

C Proof of Theorem 3.1

First, we define continuity for stochastic choice. We endow \mathcal{A} with the Hausdorff metric.

Definition 8. ρ satisfies *continuity* if $\rho_{A_k} \rightarrow \rho_A$ for all $A_k \rightarrow A$.

We now prove a stronger version of Theorem 3.1 below.

Theorem 3.1*. *Any pure characteristic ρ satisfies continuity.*

Let $\rho \in \mathcal{P}^{pc}$ with distribution μ . Consider $A_k \rightarrow A$. Note that $u(x) = u(y)$ with μ -measure zero. Now, define

$$I := \bigcup_{\{x,y\} \subset A_k \cup A} \{u \in U : u(x) = u(y)\}$$

which is measurable and $\mu(I) = 0$. Let $U^* := U \setminus I$ so $\mu(U^*) = 1$. Let μ^* be the restriction of μ on U^* .

We will now define random variables $\xi_k : U^* \rightarrow X$ and $\xi : U^* \rightarrow X$ that have distributions ρ_{A_k} and ρ_A respectively. For each A_k , let $\xi_k : U^* \rightarrow X$ be such that

$$\xi_k(u) := \arg \max_{x \in A_k} u(x)$$

and define ξ similarly for A . Note that these are well-defined because there exists a unique maximizer for every $u \in U^*$. Now, for any measurable set $E \subset X$,

$$\begin{aligned} \xi_k^{-1}(E) &= \{u \in U^* : \xi_k(u) \in E \cap A_k\} \\ &= \bigcup_{x \in E \cap A_k} \{u \in U^* : u(x) > u(y) \text{ for all } y \in A_k\} \end{aligned}$$

which is measurable. Hence, ξ_k and ξ are random variables. Note that

$$\begin{aligned} \mu^*(\xi_k^{-1}(E)) &= \sum_{x \in E \cap A_k} \mu^* \{u \in U^* : u(x) > u(y) \text{ for all } y \in A_k\} \\ &= \sum_{x \in E \cap A_k} \mu \{u \in U^* : u(x) > u(y) \text{ for all } y \in A_k\} \\ &= \rho_{A_k}(E \cap A_k) = \rho_{A_k}(E) \end{aligned}$$

so ρ_{A_k} and ρ_A are the distributions of ξ_k and ξ respectively. Since every $u \in U^* \subset U$ is continuous, by the Maximum Theorem, $\xi_k(u) = \arg \max_{x \in A_k} u(x)$ is upper hemicontinuous in A_k and thus continuous as ξ_k is singleton-valued. Since $A_k \rightarrow A$, $\xi_k \rightarrow \xi$ μ^* -a.s. and since

a.s. convergence implies convergence in distribution, $\rho_{A_k} \rightarrow \rho_A$ as desired.

D Proof of Proposition 2

We will establish Proposition 2 by breaking it down into the following two propositions. The first shows that ρ is a sequential limit of logits of degree d if and only if it is lexicographic-logit of degree d .

Proposition 2.1. $\rho \in \mathcal{P}_d^{lex}$ iff there exists a sequence $\rho^n \in \mathcal{P}_d^{log}$ such that $\rho^n \rightarrow \rho$.

Since \mathcal{P} is not metrizable under the product topology, it may not be sequential. However, the second proposition shows that the sequential limit points of logit coincides with the closure of logit.

Proposition 2.2. $\rho \in cl(\mathcal{P}_d^{log})$ iff there exists a sequence $\rho^n \in \mathcal{P}_d^{log}$ such that $\rho^n \rightarrow \rho$.

Together, these two propositions establish Proposition 2.

D.1 Proof of Proposition 2.1

We first prove the following useful lemma.

Lemma 3. Consider any sequence $\beta_n \in \mathbb{R}^m$.

- (1) If $\limsup_n |\beta_n| < \infty$, then there exists a subsequence β_i such that $\beta_i \rightarrow \beta$.
- (2) If $\limsup_n |\beta_n| = \infty$, then there exists a subsequence β_i such that $\frac{\beta_i}{|\beta_i|} \rightarrow \gamma \neq 0$.
Moreover, for any $z \in \mathbb{R}^m$, (i) $\gamma \cdot z > 0$ implies $\lim_i \beta_i \cdot z = \infty$, and (ii) $\gamma \cdot z < 0$ implies $\lim_i \beta_i \cdot z = -\infty$

Proof. First consider (1). Define

$$B := \left\{ \beta \in \mathbb{R}^m : |\beta| \leq \limsup_n |\beta_n| \right\}$$

Since $\beta_n \in B$ and B is compact, there must exist a convergence subsequence β_i such that $\beta_i \rightarrow \beta \in B$ as desired.

Now, consider (2). Note that we can find a subsequence β_j such that $|\beta_j| \rightarrow \infty$. Let $S \subset \mathbb{R}^m$ be the unit sphere and let $\hat{\beta}_j = \frac{\beta_j}{|\beta_j|} \in S$ be the normalized unit vector. Since S

is compact, there must exist a convergent subsequence β_i such that $\hat{\beta}_i \rightarrow \gamma \in S$ as desired. Since $|\beta_i| \rightarrow \infty$, if $\gamma \cdot z > 0$, then

$$\lim_i \beta_i \cdot z = \lim_i |\beta_i| (\hat{\beta}_i \cdot z) = \infty$$

The case for $\gamma \cdot z < 0$ is symmetric. □

We now prove Proposition 2.1. Let $\rho^n \in \mathcal{P}_d^{\log}$ with $\beta_n \in \mathbb{R}^m$ and $\rho^n \rightarrow \rho$. Now,

$$\rho_A(x) = \lim_n \rho_A^n(x) = \lim_n \frac{e^{\beta_n \cdot x^*}}{\sum_{y \in A} e^{\beta_n \cdot y^*}} = \left(\sum_{y \in A} e^{\lim_n \beta_n \cdot (y^* - x^*)} \right)^{-1}$$

Since $\rho(y, x)$ is well-defined, so is $\lim_n \beta_n \cdot (y^* - x^*)$ on $\bar{\mathbb{R}}$, the extended real line. Let

$$Z := \{y^* - x^* : x, y \in X\} \subset \mathbb{R}^m$$

First, suppose $\limsup_n |\beta_n| < \infty$ so by Lemma 3, there exists a convergent subsequence $\beta_i \rightarrow \beta \in \mathbb{R}^m$. Thus,

$$\rho_A(x) = \frac{e^{\beta \cdot x^*}}{\sum_{y \in A} e^{\beta \cdot y^*}} = \frac{e^{u(x)}}{\sum_{y \in A} e^{u(y)}}$$

where $u \in U_d$. This means that $\rho \in \mathcal{P}_d^{\log} \subset \mathcal{P}_d^{\text{lex}}$ as desired.

Now, suppose $\limsup_n |\beta_n| = \infty$ so by Lemma 3, there exists a convergent subsequence $\hat{\beta}_i := \frac{\beta_i}{|\beta_i|} \rightarrow \gamma_1 \in \mathbb{R}^m$. Moreover, for any $z \in Z$, $\lim_n \beta_n \cdot z = \lim_i \beta_i \cdot z = \infty$ if $\gamma_1 \cdot z > 0$ and $\lim_n \beta_n \cdot z = -\infty$ if $\gamma_1 \cdot z < 0$. Let $H_1 \subset \mathbb{R}^m$ denote the $(m-1)$ -dimensional hyperplane such that $\gamma_1 \cdot z = 0$ for all $z \in \mathbb{R}^m$. Thus

$$\rho_A(x) = 1_{\{\gamma_1 \cdot x^* \geq \gamma_1 \cdot y^* \text{ for all } y \in A\}} \left(\sum_{y \in A, y^* - x^* \in H_1} e^{\lim_i \beta_i \cdot (y^* - x^*)} \right)^{-1}$$

Let $T_1 : \mathbb{R}^m \rightarrow H_1$ be the projection mapping onto H_1 , that is

$$T_1(\beta) := \beta - (\beta \cdot \gamma_1) \gamma_1$$

Now, for any $z \in H_1$,

$$\beta \cdot z = (T_1(\beta) + (\beta \cdot \gamma_1) \gamma_1) \cdot z = T_1(\beta) \cdot z$$

We thus have

$$\rho_A(x) = 1 \{ \gamma_1 \cdot x^* \geq \gamma_1 \cdot y^* \text{ for all } y \in A \} \left(\sum_{y \in A, y^* - x^* \in H_1} e^{\lim_i T_1(\beta_i) \cdot (y^* - x^*)} \right)^{-1}$$

Now, for any $y^* - x^* \in H_1$, so we can repeat the same arguments as above. If $\limsup_i |T_1(\beta_i)| < \infty$, then by Lemma 3, we can assume $T_1(\beta_i) \rightarrow \beta \in \mathbb{R}^m$ and

$$\rho_A(x) = 1 \{ \gamma_1 \cdot x^* \geq \gamma_1 \cdot y^* \text{ for all } y \in A \} \frac{e^{\beta \cdot x^*}}{\sum_{y \in A, \gamma_1 \cdot y^* = \gamma_1 \cdot x^*} e^{\beta \cdot y^*}}$$

On the other hand, if $\limsup_i |T_1(\beta_i)| = \infty$, then by Lemma 3, we can assume $\frac{T_1(\beta_i)}{|T_1(\beta_i)|} \rightarrow \gamma_2 \in H_1$. Let $H_2 \subset \mathbb{R}^m$ denote the $(m-2)$ -dimensional hyperplane such that $\gamma_1 \cdot z = \gamma_2 \cdot z = 0$ for all $z \in \mathbb{R}^m$ and note that $\gamma_1 \cdot \gamma_2 = 0$. If we let $T_2 : \mathbb{R}^m \rightarrow H_2$ be the projection mapping onto H_2 , then by the same arguments as above,

$$\rho_A(x) = 1 \{ \gamma_j \cdot x^* \geq \gamma_j \cdot y^* \text{ for all } y \in A, j \in \{1, 2\} \} \left(\sum_{y \in A, y^* - x^* \in H_2} e^{\lim_i T_2(\beta_i) \cdot (y^* - x^*)} \right)^{-1}$$

We can continue this argument by induction, and since m is finite, we can find a sequence $(\gamma_1, \dots, \gamma_t, \beta)$ such that

$$\begin{aligned} \rho_A(x) &= 1 \{ \gamma_j \cdot x^* \geq \gamma_j \cdot y^* \text{ for all } y \in A, j \in \{1, \dots, t\} \} \frac{e^{\beta \cdot x^*}}{\sum_{y \in A, \gamma_j \cdot y^* = \gamma_j \cdot x^*, j \in \{1, \dots, t\}} e^{\beta \cdot y^*}} \\ &= 1 \{ v_j(x) \geq v_j(y) \text{ for all } y \in A, j \in \{1, \dots, t\} \} \frac{e^{u(x)}}{\sum_{y \in A, v_j(y) = v_j(x), j \in \{1, \dots, t\}} e^{u(y)}} \end{aligned}$$

for $v_1, \dots, v_t, u \in U_d$ where v_1, \dots, v_t are all orthogonal. This means that $\rho \in \mathcal{P}_d^{lex}$ as desired.

Now, suppose $\rho \in \mathcal{P}_d^{lex}$ with $\omega = (v_1, \dots, v_t)$ for $t \leq m$ and $u \in U_d$. Let β^* be the polynomial vector corresponding to u and $\gamma_1, \dots, \gamma_t \in \mathbb{R}^m$ be the polynomial vectors corresponding to ω which are orthogonal. Without loss of generality, we can assume that $\gamma_i \neq 0$. By a change of basis, we can also assume without loss that they correspond to the standard basis in \mathbb{R}^m . Now, define

$$\beta_n = \left(n^k, n^{k-1}, \dots, n^2, n\beta_t^*, n\beta_{t+1}^*, \dots, n\beta_m^* \right)$$

Let ρ^n be the logit corresponding to β_n so $\rho^n \in \mathcal{P}_d^{log}$ and it is straightforward to see that $\rho^n \rightarrow \rho$.

D.2 Proof of Proposition 2.2

We first prove a couple of technical lemmas.

Lemma 4. *Consider non-constant $u, v \in U_d$ and suppose $u(x) > u(y)$ implies $v(x) \geq v(y)$ for all $x, y \in D$ where D is dense in X . Then $u(x) > u(y)$ iff $v(x) > v(y)$ for all $x, y \in X$.*

Proof. We first show that $u(x) > u(y)$ implies $v(x) \geq v(y)$ for all $x, y \in X$. Since D is dense in X , we can find $x_\varepsilon, y_\varepsilon \in D$ arbitrarily close to x and y respectively such that $u(x_\varepsilon) > u(y_\varepsilon)$. Thus, $v(x_\varepsilon) \geq v(y_\varepsilon)$ so by continuity $v(x) \geq v(y)$.

Now, suppose $u(x) > u(y)$ for $x, y \in X$ but $v(x) = v(y)$. Define

$$g(\alpha) = v(x) - v((1 - \alpha)x + \alpha y)$$

and note that g is a polynomial. Since v is non-constant, g has a finite number of roots. Label them as $\{x_0, \dots, x_n\}$ where $x_i = (1 - \alpha_i)x + \alpha_i y$ with $\alpha_0 = 0$ and $\alpha_n = 1$. Since g has no roots for $\alpha \in (0, \alpha_1)$, it must be either $g(\alpha) > 0$ for all $\alpha \in (0, \alpha_1)$ or $g(\alpha) < 0$ for all $\alpha \in (0, \alpha_1)$. In either case, we can find $x_\varepsilon, x'_\varepsilon \in D$ arbitrarily close to x_0 and x_1 respectively such that $u(x_\varepsilon) > u(x'_\varepsilon)$. By continuity, $u(x_0) \geq u(x_1)$. By the same argument, we can also find $\tilde{x}_\varepsilon, \tilde{x}'_\varepsilon \in D$ arbitrarily close to x_0 and x_1 respectively such that $u(\tilde{x}_\varepsilon) < u(\tilde{x}'_\varepsilon)$ so $u(x_0) \leq u(x_1)$. This implies that $u(x_0) = u(x_1)$. By induction, we can repeat this argument for all x_i for $i \in \{2, \dots, n\}$ so we have $u(x) = u(y)$ yielding a contradiction. Thus, $u(x) > u(y)$ implies $v(x) > v(y)$ and the conclusion follows by symmetry. \square

Lemma 5. *If $\mathcal{D} \subset \mathcal{A}$ is countable, then for any limit point σ of \mathcal{P}_d^{\log} there exists a sequence $\rho^n \in \mathcal{P}_d^{\log}$ such that $\rho^n_A \rightarrow \sigma_A$ for every $A \in \mathcal{D}$.*

Proof. Since \mathcal{D} is countable, we can find sets \mathcal{D}_n where $|\mathcal{D}_n| = n$ and $\mathcal{D}_n \nearrow \mathcal{D}$. Now, for any $A \in \mathcal{D}_n$, let

$$Q_n(A) := \left\{ q \in \Delta A : |q - \sigma_A| < \frac{1}{n} \right\}$$

which is open in ΔA . Thus,

$$\mathcal{Q}_n = \prod_{A \in \mathcal{A}_n} Q_n(A) \times \prod_{A \in \mathcal{A} \setminus \mathcal{A}_n} \Delta A$$

is open in the product topology. Since $\sigma \in \mathcal{Q}_n$ is a limit point, we can find some $\rho^n \in \mathcal{Q}_n \cap \mathcal{P}_d^{\log}$ for every n .

Fix some $A \in \mathcal{D}$ and note that we can find some N such that $A \in \mathcal{D}_n$ for all $n > N$. Since $\rho^n \in \mathcal{Q}_n$, that means that $|\rho_A^n - \sigma_A| < \frac{1}{n}$ for all $n > N$. Thus, $\rho_A^n \rightarrow \sigma_A$ for every $A \in \mathcal{D}$ as desired. \square

We now prove Proposition 2.2. Let D_1 be a countable dense subset of X and $\mathcal{D}_1 \subset \mathcal{A}$ denote all binary menus $\{x, y\}$ where $x, y \in D_1$. Let

$$\begin{aligned} Z &:= \{y^* - x^* : x, y \in X\} \\ Z_1 &:= \{y^* - x^* : x, y \in D_1\} \end{aligned}$$

Since \mathcal{D}_1 is countable, by Lemma 5, we can find $\rho^n \in \mathcal{P}_d^{\text{log}}$ such that $\rho_A^n \rightarrow \sigma_A$ for every $A \in \mathcal{D}_1$. Thus,

$$\sigma(x, y) = \lim_n \rho^n(x, y) = \left(\sum_{y \in A} e^{\lim_n \beta_n \cdot (y^* - x^*)} \right)^{-1}$$

so $\lim_n \beta_n \cdot z$ exists for every $z \in Z_1$. We now consider two cases: (1) $\limsup_n |\beta_n| < \infty$ and (2) $\limsup_n |\beta_n| = \infty$.

First, consider case (1). By Lemma 3, there exists a convergent subsequence $\beta_j \rightarrow \beta$ for some $\beta \in \mathbb{R}^m$. Consider any $B \in \mathcal{A}$ and so by Lemma 5, we can find $\tau^n \in \mathcal{P}_d^{\text{log}}$ such that $\tau_A^n \rightarrow \sigma_A$ for every $A \in \mathcal{D}_1 \cup \{B\}$. Let β'_n correspond to τ^n so for all $z \in Z_1$,

$$\lim_n \beta'_n \cdot z = \lim_n \beta_n \cdot z = \beta \cdot z$$

We now show that $\limsup_n |\beta'_n| < \infty$. Suppose otherwise, so by Lemma 3, we can find a subsequence $\frac{\beta'_i}{|\beta'_i|} \rightarrow \gamma \in \mathbb{R}^m$. Moreover, if $\gamma \cdot z > 0$, then $\lim_i \beta'_i \cdot z = \infty > \beta \cdot z$ yielding a contradiction. The case for $\gamma \cdot z < 0$ is symmetric so it must be that $\gamma \cdot z = 0$ for all $z \in Z_1$. Since D_1 is dense in X and the latter is a full-dimensional subset of \mathbb{R}^k , this is impossible. Thus, $\limsup_n |\beta'_n| < \infty$ so by continuity (Corollary 6.40 of AB), $\lim_n \beta'_n \cdot z = \lim_n \beta_n \cdot z$ for all $z \in Z$. This means that

$$\begin{aligned} \sigma_B(x) &= \lim_n \tau_B^n(x) = \lim_n \frac{e^{\beta'_n \cdot x^*}}{\sum_{y \in B} e^{\beta'_n \cdot y^*}} = \left(\sum_{y \in B} e^{\lim_n \beta'_n \cdot (y^* - x^*)} \right)^{-1} \\ &= \left(\sum_{y \in B} e^{\lim_n \beta_n \cdot (y^* - x^*)} \right)^{-1} = \lim_n \frac{e^{\beta_n \cdot x}}{\sum_{y \in B} e^{\beta_n \cdot y}} \end{aligned}$$

Since this is true for all $B \in \mathcal{A}$, σ is a sequential limit of logits of degree d .

Now, consider case (2). By Lemma 3, we can find a subsequence $\frac{\beta_i}{|\beta_i|} \rightarrow \gamma_1 \in \mathbb{R}^m$ such

that $\lim_i \beta_i \cdot z = \infty$ if $\gamma_1 \cdot z > 0$ and $\lim_i \beta_i \cdot z = -\infty$ if $\gamma_1 \cdot z < 0$ for all $z \in Z$. Let $H_1 \subset \mathbb{R}^m$ denote the $(m-1)$ -dimensional hyperplane such that $\gamma_1 \cdot z = 0$. We consider two subcases: (i) $Z \cap H_1 = \emptyset$ and (ii) $Z \cap H_1 \neq \emptyset$.

First, consider subcase (i) where $Z \cap H_1 = \emptyset$ so $\gamma_1 \cdot z \neq 0$ for all $z \in Z$. Consider any $B \in \mathcal{A}$ and so by Lemma 5, we can find $\tau^n \in \mathcal{P}_d^{\log}$ such that $\tau_A^n \rightarrow \sigma_A$ for every $A \in \mathcal{D}_1 \cup \{B\}$. Let β'_n correspond to τ^n so $\lim_n \beta'_n \cdot z = \lim_n \beta_n \cdot z$ for all $z \in Z_1$. Note that if $\limsup_n |\beta'_n| < \infty$, then by the same argument as above, we have a contradiction. Thus, it must be that $\limsup_n |\beta'_n| = \infty$, so by Lemma 3 again, we can find a subsequence $\frac{\beta'_i}{|\beta'_i|} \rightarrow \gamma'_1 \in \mathbb{R}^m$ such that $\lim_i \beta'_i \cdot z = \infty$ if $\gamma'_1 \cdot z > 0$ and $\lim_i \beta'_i \cdot z = -\infty$ if $\gamma'_1 \cdot z < 0$ for all $z \in Z$. Let $u, v \in U_d$ correspond to γ_1 and γ'_1 respectively. Thus, for $x, y \in D$, $u(x) > u(y)$ implies $\infty = \lim_n \beta_n \cdot (x^* - y^*) = \lim_n \beta_n \cdot (x^* - y^*)$ so $v(x) \geq v(y)$. By Lemma 4, this means that $\gamma_1 \cdot z > 0$ iff $\gamma'_1 \cdot z > 0$ for all $z \in Z$. Since $\gamma_1 \cdot z \neq 0$ for all $z \in Z$, $\lim_i \beta'_i \cdot z = \lim_i \beta_i \cdot z$ for all $z \in Z$ so

$$\sigma_B(x) = \lim_n \tau_B^n(x) = \left(\sum_{y \in B} e^{\lim_n \beta'_n \cdot (y^* - x^*)} \right)^{-1} = \left(\sum_{y \in B} e^{\lim_n \beta_n \cdot (y^* - x^*)} \right)^{-1} = \lim_n \frac{e^{\beta_n \cdot x}}{\sum_{y \in B} e^{\beta_n \cdot y}}$$

Since this is true for all $B \in \mathcal{A}$, σ is a sequential limit of logits of degree d .

Next, consider subcase (ii) where $Z \cap H_1 \neq \emptyset$. Let $D_2 \subset X$ be a countable set such that

$$Z_2 := \{y^* - x^* : x, y \in D_2\}$$

is dense in $Z \cap H_1$. Let $\mathcal{D}_2 \subset \mathcal{A}$ denote all binary menus $\{x, y\}$ where $x, y \in D_1 \cup D_2$. Since \mathcal{D}_2 is countable, by Lemma 5, we can find $\tilde{\rho}^n \in \mathcal{P}_d^{\log}$ such that $\tilde{\rho}_A^n \rightarrow \sigma_A$ for every $A \in \mathcal{D}_2$. Let $\tilde{\beta}_n$ correspond to $\tilde{\rho}^n$ so $\lim_n \tilde{\beta}_n \cdot z = \lim_n \beta_n \cdot z$ for all $z \in Z_1 \cup Z_2$. Since we are in case (2), by the same argument as above, it must be that $\limsup_n |\tilde{\beta}_n| = \infty$, so by Lemma 3 again, we can find a subsequence $\frac{\tilde{\beta}_i}{|\tilde{\beta}_i|} \rightarrow \tilde{\gamma}_1 \in \mathbb{R}^m$ such that $\lim_i \tilde{\beta}_i \cdot z = \infty$ if $\tilde{\gamma}_1 \cdot z > 0$ and $\lim_i \tilde{\beta}_i \cdot z = -\infty$ if $\tilde{\gamma}_1 \cdot z < 0$ for all $z \in Z$. Applying Lemma 4 as above, it must be that $\gamma_1 \cdot z = 0$ iff $\tilde{\gamma}_1 \cdot z = 0$ for all $z \in Z$.

Let $T_1 : \mathbb{R}^m \rightarrow H_1$ be the projection mapping onto H_1 , that is

$$T_1(\beta) := \beta - (\beta \cdot \gamma_1) \gamma_1$$

Now, for any $z \in Z \cap H_1$,

$$\beta \cdot z = (T_1(\beta) + (\beta \cdot \gamma_1) \gamma_1) \cdot z = T_1(\beta) \cdot z$$

We now consider two cases: (1') $\limsup_n |T_1(\tilde{\beta}_n)| < \infty$ and (2') $\limsup_n |T_1(\tilde{\beta}_n)| = \infty$.

First, consider case (1'). By Lemma 3, there exists a convergent subsequence $T_1(\tilde{\beta}_i) \rightarrow \beta$ for some $\beta \in \mathbb{R}^m$. Consider any $B \in \mathcal{A}$ and so by Lemma 5, we can find $\tau^n \in \mathcal{P}_d^{\log}$ such that $\tau_A^n \rightarrow \sigma_A$ for every $A \in \mathcal{D}_2 \cup \{B\}$. Let β'_n correspond to τ^n so for all $z \in Z_1 \cup Z_2$,

$$\lim_n \tilde{\beta}_n \cdot z = \lim_n \beta'_n \cdot z$$

By the same argument as above, it must be that $\lim_n \tilde{\beta}_n \cdot z = \lim_n \beta'_n \cdot z$ for all $z \in Z \setminus H_1$. Now, for $z \in Z_2$, we have

$$\lim_n T_1(\beta'_n) \cdot z = \lim_n \beta'_n \cdot z = \lim_n \tilde{\beta}_n \cdot z = \lim_i T_1(\tilde{\beta}_i) \cdot z = \beta \cdot z$$

By the same argument as above, this implies that $\limsup_n |T_1(\beta'_n)| < \infty$. Thus, we have $\lim_n \beta'_n \cdot z = \lim_n \beta_n \cdot z$ for all $z \in Z$ so

$$\sigma_B(x) = \lim_n \tau_B^n(x) = \left(\sum_{y \in B} e^{\lim_n \beta'_n \cdot (y^* - x^*)} \right)^{-1} = \left(\sum_{y \in B} e^{\lim_n \tilde{\beta}_n \cdot (y^* - x^*)} \right)^{-1} = \lim_n \frac{e^{\tilde{\beta}_n \cdot x}}{\sum_{y \in B} e^{\tilde{\beta}_n \cdot y}}$$

Since this is true for all $B \in \mathcal{A}$, σ is a sequential limit of logits of degree d .

Finally, for case (2'), we can inductively apply the same reasoning as in case (2) above. Since the dimension is finite, we obtain the conclusion by induction.

E Proof of Proposition 3

Let $\bar{\mathcal{P}}_d^{\log}$ denote the closed convex hull of \mathcal{P}_d^{\log} , that is

$$\bar{\mathcal{P}}_d^{\log} := \text{cl} \left(\text{co} \left(\mathcal{P}_d^{\log} \right) \right) \subset \mathcal{P}$$

which is compact. We first show the following which will imply that \mathcal{P}_d^{mlex} is closed.

Lemma 6. $\bar{\mathcal{P}}_d^{\log} = \mathcal{P}_d^{mlex}$

Proof. Since $\text{co}(\mathcal{P}_d^{\log})$ is convex, $\bar{\mathcal{P}}_d^{\log}$ is also convex (Lemma 5.27 of AB). We first show that $\mathcal{P}_d^{mlex} \subset \bar{\mathcal{P}}_d^{\log}$. Let $\rho \in \mathcal{P}_d^{mlex}$ so there exists a distribution ν on $\Omega_d \times U_d$ such that

$$\rho = \int_{\Omega_d \times U_d} \rho_{(\omega, u)} d\nu$$

where $\rho_{(\omega, u)} \in \mathcal{P}_d^{\log}$ is the lexicographic-logit stochastic choice corresponding to $(\omega, u) \in \Omega_d \times U_d$. Suppose $\rho \notin \bar{\mathcal{P}}_d^{\log}$. Since \mathcal{V} is locally convex (Lemma 1), continuous linear functionals

separates points in \mathcal{V} . Thus, we can apply the strict separating hyperplane theorem (Theorem 3.5 of Rudin (1991)) and find a continuous linear functional Λ such that for all $\tau \in \bar{\mathcal{P}}_d^{lex}$,

$$\Lambda(\rho) = 1 > 0 = \Lambda(\tau)$$

Now,

$$1 = \Lambda(\rho) = \Lambda\left(\int_{\Omega_d \times U_d} \rho(\omega, u) d\nu\right) = \int_{\Omega_d \times U_d} \Lambda(\rho(\omega, u)) d\nu = 0$$

as $\rho(\omega, u) \in \bar{\mathcal{P}}_d^{lex}$. This yields a contradiction so $\mathcal{P}_d^{mlex} \subset \bar{\mathcal{P}}_d^{lex}$.

Next, we show that $\bar{\mathcal{P}}_d^{lex} \subset \mathcal{P}_d^{mlex}$. Fix some $\rho \in \bar{\mathcal{P}}_d^{lex}$. Since $\mathcal{P}_d^{lex} = cl(\mathcal{P}_d^{log})$ by Proposition 2, it is closed and thus compact. By Theorem 3.28 of Rudin (1991), there exists a Borel probability measure π on \mathcal{P}_d^{lex} such that

$$\rho = \int_{\mathcal{P}_d^{lex}} \tau d\pi$$

We now show ρ must be mixed lexicographic-logit. Consider the mapping $\varphi : \Omega_d \times U_d \rightarrow \mathcal{P}_d^{lex}$ such that $\varphi = \rho(\omega, u)$. Let \mathcal{G} be the σ -algebra on $\Omega_d \times U_d$ generated by φ . We can thus define a measure ν on \mathcal{G} such that

$$\pi = \nu \circ \varphi^{-1}$$

Thus, by a change of variables (Theorem 13.46 of AB),

$$\rho = \int_{\mathcal{P}_d^{lex}} \tau d\pi = \int_{\Omega_d \times U_d} \varphi(\omega, u) d\nu = \int_{\Omega_d \times U_d} \rho(\omega, u) d\nu$$

so $\rho \in \mathcal{P}_d^{mlex}$ as desired. \square

We now prove Proposition 3. Since $\mathcal{P}_d^{log} \subset \mathcal{P}_d^{lex}$, we have that $\mathcal{P}_d^{mlog} \subset \mathcal{P}_d^{mlex}$. Thus,

$$cl(\mathcal{P}_d^{mlog}) \subset cl(\mathcal{P}_d^{mlex}) = \mathcal{P}_d^{mlex}$$

where the last equality follows from Lemma 6. Now, consider $\rho \in \mathcal{P}_d^{mlex}$ and suppose $\rho \notin cl(\mathcal{P}_d^{mlog})$. Applying the strict separating hyperplane theorem again, there exists a continuous linear functional Λ such that for all $\tau \in cl(\mathcal{P}_d^{mlog})$,

$$\Lambda(\rho) = 1 > 0 = \Lambda(\tau)$$

Now

$$1 = \Lambda(\rho) = \Lambda\left(\int_{\Omega_d \times U_d} \rho(\omega, u) d\nu\right) = \int_{\Omega_d \times U_d} \Lambda(\rho(\omega, u)) d\nu = 0$$

where the last equality follows from the fact that $\rho_{(\omega,u)} \in \mathcal{P}_d^{lex} = cl(\mathcal{P}_d^{log}) \subset cl(\mathcal{P}_d^{mlog})$. Thus, we have a contradiction. This shows that $cl(\mathcal{P}_d^{mlog}) = \mathcal{P}_d^{mlex}$ as desired.

F Proof of Theorem 1

We first prove a couple of technical lemmas.

Lemma 7. *Let $x_n = (1 - \frac{1}{n})x + \frac{1}{n}y$ for $x, y \in X$. For any $\omega \in \Omega_d$, exactly one of the following holds*

- (1) $\lim_n 1\{y \succ_\omega x_n\} = 1$
- (2) $\lim_n 1\{y \sim_\omega x_n\} = 1$
- (3) $\lim_n 1\{y \prec_\omega x_n\} = 1$

Proof. For any $v \in U_d$, let

$$g(\alpha) = v(y) - v((1 - \alpha)x + \alpha y)$$

and note that g is a polynomial. Since any non-zero polynomial has a finite number of roots, it means that we can find some N such that either (i) $g(\frac{1}{n}) > 0$ for all $n > N$, (ii) $g(\frac{1}{n}) = 0$ for all $n > N$, or (iii) $g(\frac{1}{n}) < 0$ for all $n > N$. Let U_d^\succ , U_d^\sim and U_d^\prec be the partition of U_d corresponding to these three conditions. It is straightforward to see that for $\omega = (v_1, \dots, v_t) \in \Omega_d$,

- (1) If $v_i \in U_d^\sim$ for all $0 \leq i < j \leq t$ and $v_j \in U_d^\succ$, then $\lim_n 1\{y \succ_\omega x_n\} = 1$.
- (2) If $v_i \in U_d^\sim$ for all $1 \leq i \leq t$, then $\lim_n 1\{y \sim_\omega x_n\} = 1$.
- (3) Otherwise, $\lim_n 1\{y \prec_\omega x_n\} = 1$

The result follows. □

We now prove Theorem 1. Let $\rho \in \mathcal{P}^{pc}$ and suppose $\rho \in cl(\mathcal{P}_d^{mlog})$. By Proposition 3, $\rho \in \mathcal{P}_d^{mlex}$ so there exists some distribution ν on $\Omega_d \times U_d$ such that

$$\rho_A(x) = \int_{\Omega_d \times U_d} 1\{x \succeq_\omega y \text{ for all } y \in A\} \frac{e^{u(x)}}{\sum_{y \in A, y \sim_\omega x} e^{u(y)}} d\nu$$

We will show that any distinct $x, y \in X$, $x \sim_\omega y$ with ν -measure zero. Let $x_n = (1 - \frac{1}{n})x + \frac{1}{n}y$ and define the following sets

$$\begin{aligned}\Omega_1 &= \{\omega \in \Omega_d : y \succ_\omega x\} & \Omega'_1 &= \left\{ \omega \in \Omega_d : \lim_n 1 \{y \succ_\omega x_n\} = 1 \right\} \\ \Omega_2 &= \{\omega \in \Omega_d : y \sim_\omega x\} & \Omega'_2 &= \left\{ \omega \in \Omega_d : \lim_n 1 \{y \sim_\omega x_n\} = 1 \right\} \\ \Omega_3 &= \{\omega \in \Omega_d : y \prec_\omega x\} & \Omega'_3 &= \left\{ \omega \in \Omega_d : \lim_n 1 \{y \prec_\omega x_n\} = 1 \right\}\end{aligned}$$

Note that $\{\Omega_1, \Omega_2, \Omega_3\}$ and $\{\Omega'_1, \Omega'_2, \Omega'_3\}$ are both partitions of Ω_d , where the latter follows from Lemma 7. Suppose $\omega = (v_1, \dots, v_t) \in \Omega'_2$, so $v_i(y) = v_i(x_n)$ for sufficiently large n and all $1 \leq i \leq t$. This implies that $v_i(y) = v_i(x)$ for all $1 \leq i \leq t$ or $y \sim_\omega x$. Thus, $\Omega'_2 \subset \Omega_2$. This implies that

$$\Omega_1 \cap \Omega'_2 = \Omega_3 \cap \Omega'_2 = \emptyset \quad (1)$$

Let $A_n = \{y, x, x_n\}$ and note that $A_n \rightarrow \{y, x\}$. Since $\rho \in \mathcal{P}^{pc}$, by Theorem 3.1,

$$\rho(y, x) = \lim_n \rho(y, x_n) = \lim_n \rho_{A_n}(y)$$

For ease of notation, we suppress the dependence on U_d ; for instance, we let $\nu(\Omega_i)$ denote $\nu(\Omega_i \times U_d)$. Now,

$$\begin{aligned}\rho(y, x) &= \nu(\Omega_1) + \int_{\Omega_2} \frac{e^{u(y)}}{e^{u(y)} + e^{u(x)}} d\nu \\ &= \nu(\Omega_1 \cap \Omega'_1) + \nu(\Omega_1 \cap \Omega'_3) + \int_{\Omega_2} \frac{e^{u(y)}}{e^{u(y)} + e^{u(x)}} d\nu\end{aligned} \quad (2)$$

where the second equality follows from equation (1). By dominated convergence

$$\begin{aligned}\lim_n \rho(y, x_n) &= \nu(\Omega'_1) + \int_{\Omega'_2} \frac{e^{u(y)}}{e^{u(y)} + e^{u(x)}} d\nu \\ &= \nu(\Omega_1 \cap \Omega'_1) + \nu(\Omega_2 \cap \Omega'_1) + \nu(\Omega_3 \cap \Omega'_1) + \int_{\Omega_2 \cap \Omega'_2} \frac{e^{u(y)}}{e^{u(y)} + e^{u(x)}} d\nu\end{aligned} \quad (3)$$

where the last equality follows from (1) again. Applying dominated convergence again,

$$\lim_n \rho_{A_n}(y) = \nu(\Omega_1 \cap \Omega'_1) + \int_{\Omega_2 \cap \Omega'_1} \frac{e^{u(y)}}{e^{u(y)} + e^{u(x)}} d\nu + \int_{\Omega_2 \cap \Omega'_2} \frac{e^{u(y)}}{e^{u(y)} + 2e^{u(x)}} d\nu \quad (4)$$

Subtracting equation (4) from (3), we have

$$0 = \int_{\Omega_2 \cap \Omega'_1} \left(1 - \frac{e^{u(y)}}{e^{u(y)} + e^{u(x)}} \right) d\nu + \nu(\Omega_3 \cap \Omega'_1) + \int_{\Omega_2 \cap \Omega'_2} \left(\frac{e^{u(y)}}{e^{u(y)} + e^{u(x)}} - \frac{e^{u(y)}}{e^{u(y)} + 2e^{u(x)}} \right) d\nu$$

This implies that $\Omega_2 \cap \Omega'_1$, $\Omega_3 \cap \Omega'_1$ and $\Omega_2 \cap \Omega'_2$ are all ν -measure zero sets. Combining equations (2) and (3), we have

$$0 = \nu(\Omega_1 \cap \Omega'_3) + \int_{\Omega_2} \frac{e^{u(y)}}{e^{u(y)} + e^{u(x)}} d\nu$$

so Ω_2 must also be a ν -measure zero set as desired.

We thus have

$$\begin{aligned} \rho_A(x) &= \int_{\Omega_d} 1 \{x \succeq_\omega y \text{ for all } y \in A\} d\nu \\ &= \nu(\{\omega \in \Omega_d : x \succeq_\omega y \text{ for all } y \in A\}) \end{aligned}$$

Now, for every $\omega = (v_1, \dots, v_t) \in \Omega_d$, let $\theta_\omega = (\beta_1^\omega, \dots, \beta_t^\omega)$ denote the collection of coefficients such that

$$v_i(x) = \beta_i^\omega \cdot x^*$$

We can thus extend ρ to a stochastic choice ρ^* in \mathbb{R}^m such that $\rho_A(x) = \rho_{A^*}^*(x^*)$ and for any finite $D \subset \mathbb{R}^m$,

$$\rho_D^*(z) = \nu(\{\omega \in \Omega_d : z \succeq_{\theta_\omega} w \text{ for all } w \in D\})$$

Moreover, since for all $x, y \in X$, $x \sim_\omega y$ with ν -measure zero, without loss of generality, we can assume that for all $z, w \in \mathbb{R}^m$, $z \sim_{\theta_\omega} w$ with ν -measure zero as well.

We now show that ρ^* satisfies the Gul and Pesendorfer (2006) axioms. Since \succeq_{θ_ω} satisfies independence, ρ^* satisfies linearity and since $y \sim_{\theta_\omega} x$ with ν -measure zero, ρ also satisfies extremeness. Mixture continuity follows from the same argument as Lemma 3 in the Supplement of Gul and Pesendorfer (2006). This means that we can find some finitely-additive μ^* on \mathbb{R}^m such that

$$\rho_D^*(z) = \mu^*(\{\beta \in \mathbb{R}^m : \beta \cdot z \geq \beta \cdot w \text{ for all } w \in D\})$$

Since $\rho_A(x) = \rho_{A^*}^*(x^*)$,

$$\begin{aligned}\rho_A(x) &= \mu^* (\{\beta \in \mathbb{R}^m : \beta \cdot x^* \geq \beta \cdot y^* \text{ for all } y \in A\}) \\ &= \mu (\{u \in U_d : u(x) \geq u(y) \text{ for all } y \in A\})\end{aligned}$$

Finally, since ρ is continuous, the countable additivity of μ follows from the same argument as Lemma 6 in Gul and Pesendorfer (2006).

References

- ACKERBERG, D., AND M. RYSMAN (2005): “Unobserved Product Differentiation in Discrete-Choice Models: Estimating Price Elasticities and Welfare Effects,” *RAND Journal of Economics*, 36(4), 1–19.
- AHN, D., AND T. SARVER (2013): “Preference for Flexibility and Random Choice,” *Econometrica*, 81(1), 341–361.
- ALIPRANTIS, C., AND K. BORDER (2006): *Infinite Dimensional Analysis*. Springer.
- ANDERSON, S., A. DE PALMA, AND J. THISSE (1989): “Demand for Differentiated Products, Discrete Choice Models, and the Characteristics Approach,” *The Review of Economic Studies*, 56(1), 21–35.
- APESTEGUIA, J., AND M. BALLESTER (2018): “Monotone Stochastic Choice Models: The Case of Risk and Time Preferences,” *Journal of Political Economy*, 126(1), 74–106.
- APESTEGUIA, J., M. BALLESTER, AND J. LU (2017): “Single-Crossing Random Utility Models,” *Econometrica*, 85(2), 661–674.
- BAJARI, P., AND C. BENKARD (2004): “Comparing Hedonic and Random Utility Models of Demand with an Application to PC’s,” Mimeo.
- (2005): “Demand Estimation with Heterogeneous Consumers and Unobserved Product Characteristics: A Hedonic Approach,” *Journal of Political Economy*, 113(6), 1239–1276.
- BERRY, S., J. LEVINSOHN, AND A. PAKES (1995): “Automobile Prices in Market Equilibrium,” *Econometrica*, pp. 841–890.
- BERRY, S., AND A. PAKES (2007): “The Pure Characteristics Demand Model,” *International Economic Review*, 48(4), 1193–1225.
- CERREIA-VIOGLIO, S., F. MACCHERONI, M. MARINACCI, AND A. RUSTICHINI (2018a): “Law of Demand and Stochastic Choice,” Mimeo.
- (2018b): “Multinomial Logit Processes and Preference Discovery: Inside and Outside the Black Box,” Mimeo.

- CHAMBERS, C., T. CUHADAROGLU, AND Y. MASATLIOGLU (2020): “Behavioral Influence,” Mimeo.
- COHEN, M. (1980): “Random Utility Systems - The Infinite Case,” *Journal of Mathematical Psychology*, 22, 1–23.
- COMPIANI, G. (2019): “Market Counterfactuals and the Specification of Multi-Product Demand: a Nonparametric Approach,” Mimeo.
- DAVIDSON, D., AND J. MARSCHAK (1959): “Experimental Tests of Stochastic Decision Theory,” in *Measurement: Definitions and Theories*, ed. by C. W. Churchman. Wiley.
- DEBREU, G. (1958): “Stochastic Choice and Cardinal Utility,” *Econometrica*, 26(3), 440–444.
- (1960): “Individual Choice Behavior: A Theoretical Analysis by R. Duncan Luce (review),” *American Economic Review*, 50(1), 186–188.
- DURAJ, J. (2018): “Dynamic Random Subjective Expected Utility,” Mimeo.
- FRICK, M., R. IJIMA, AND T. STRZALECKI (2019): “Dynamic Random Utility,” *Econometrica*, 87(6), 1941–2002.
- FUDENBERG, D., AND T. STRZALECKI (2015): “Dynamic Logit with Choice Aversion,” *Econometrica*, 83(2), 651–691.
- GOWRISANKARAN, G., AND M. RYSMAN (2012): “Dynamics of Consumer Demand for New Durable Goods,” *Journal of Political Economy*, 120(6), 1173–1219.
- GUL, F., P. NATENZON, AND W. PESENDORFER (2014): “Random Choice as Behavioral Optimization,” *Econometrica*, 82(5), 1873–1912.
- GUL, F., AND W. PESENDORFER (2006): “Random Expected Utility,” *Econometrica*, 74(1), 121–146.
- HAUSMAN, J., AND D. WISE (1978): “A Conditional Probit Model for Qualitative Choice: Discrete Decisions Recognizing Interdependence and Heterogeneous Preferences,” *Econometrica*, 46(2), 403–426.

- HENDEL, I., AND A. NEVO (2006): “Measuring the Implications of Sales and Consumer Inventory Behavior,” *Econometrica*, 74(6), 1637–1673.
- HOTELLING, H. (1929): “Stability in Competition,” *Economic Journal*, 39, 41–57.
- HOTZ, V. J., AND R. MILLER (1993): “Conditional Choice Probabilities and the Estimation of Dynamic Models,” *Review of Economic Studies*, 60, 497–529.
- LIN, Y. (2019): “Random Non-Expected Utility: Non-Uniqueness,” Mimeo.
- LU, J. (2016): “Random Choice and Private Information,” *Econometrica*, 84(6), 1983–2027.
- (2019): “Random Ambiguity,” Mimeo.
- LU, J., AND K. SAITO (2018): “Random Intertemporal Choice,” *Journal of Economic Theory*, 177.
- LUCE, D. (1959): *Individual Choice Behavior*. New York: Wiley.
- MCFADDEN, D. (1973): “Conditional Logit Analysis of Qualitative Choice Behavior,” in *Frontiers in Econometrics*, ed. by P. Zarembka. Academic Press.
- MCFADDEN, D., AND K. TRAIN (2000): “Mixed MNL Models for Discrete Response,” *Journal of Applied Econometrics*, pp. 447–470.
- MUNKRES, J. (2000): *Topology*. Prentice Hall.
- NARITA, Y., AND K. SAITO (2021): “Approximating Choice Data by Discrete Choice Models,” Mimeo.
- NATENZON, P. (2019): “Random Choice and Learning,” *Journal of Political Economy*, 127(1), 419–457.
- NEVO, A. (2001): “Measuring Market Power in the Ready-to-Eat Cereal Industry,” *Econometrica*, 69(2), 307–342.
- PETRIN, A. (2002): “Quantifying the Benefits of New Products: The Case of the Minivan,” *Journal of Political Economy*, 110(4), 705–729.
- RUDIN, W. (1991): *Functional Analysis*. McGraw-Hill.

- RUST, J. (1987): “Optimal replacement of GMC bus engines, an empirical model of Harold Zurcher,” *Econometrica*, 55(5), 999–1033.
- SAITO, K. (2018): “Axiomatizations of the Mixed Logit Model,” Mimeo.
- SALOP, S. (1979): “Monopolistic Competition with Outside Goods,” *The Bell Journal of Economics*, 10(1), 141–156.
- TSERENJIGMID, G., AND M. KOVACH (2020): “Behavioral Foundations of Nested Stochastic Choice and Nested Logit,” Mimeo.
- WILCOX, N. (2011): “Stochastically more risk averse: A contextual theory of stochastic discrete choice under risk,” *Journal of Econometrics*, 162, 89–104.