

This page is intentionally blank. Proper e-companion title page, with INFORMS branding and exact metadata of the main paper, will be produced by the INFORMS office when the issue is being assembled.

Online Appendix

Estimating Large-Scale Tree Logit Models

Srikanth Jagabathula

Leonard N. Stern School of Business, New York University, New York, NY 10012, sjagabat@stern.nyu.edu

Paat Rusmevichientong

Marshall School of Business, University of Southern California, Los Angeles, CA 90089, rusmevic@marshall.usc.edu

Ashwin Venkataraman

Naveen Jindal School of Management, University of Texas at Dallas, Richardson, TX 75080, axv190029@utdallas.edu

Xinyi Zhao

Amazon Advertising, Palo Alto, CA 94301, xz2197@stern.nyu.edu

Appendix A: Proof of Theorem 2.1 (Random Utility Representation)

Recall that a random variable x follows a Gumbel distribution with a location parameter μ and a scaling parameter β if for all $x \in \mathbb{R}$, $\mathbb{P}\{x \leq x\} = e^{-e^{-(x-\mu)/\beta}}$, and we denote this by $x \sim \text{Gumbel}(\mu, \beta)$. We say that x follows a *standard* Gumbel distribution if $x \sim \text{Gumbel}(0, 1)$. The proof of Theorem 2.1 makes use of the following standard results about the Gumbel distribution (Gumbel 2004).

Lemma A.1 (Gumbel Properties) *Suppose $x \sim \text{Gumbel}(\mu, \beta)$, and let x_1, \dots, x_n be independent Gumbel random variables with $x_j \sim \text{Gumbel}(\mu_j, \beta)$ for all j . Then,*

1. $\mathbb{E}[x] = \mu + \beta\gamma$ where $\gamma = 0.57721\dots$ is the Euler-Mascheroni constant and $\text{Var}[x] = \pi^2\beta^2/6$.
2. For any $b > 0$ and $a \in \mathbb{R}$, $bx + a \sim \text{Gumbel}(b\mu + a, b\beta)$.
3. The random variable $x_1 - x_2$ follows a Logistic distribution with a location parameter $\mu_1 - \mu_2$ and scale parameter β ; that is, for all $x \in \mathbb{R}$, $\mathbb{P}\{x_1 - x_2 \leq x\} = \frac{1}{1 + e^{-[x - (\mu_1 - \mu_2)]/\beta}}$.
4. The random variable $\max_{i=1, \dots, n} x_i$ follows Gumbel distribution with location parameter $\beta \ln\left(\sum_{j=1}^n e^{\mu_j/\beta}\right)$ and scale parameter β , and $\mathbb{P}\{x_j > \max_{\ell \neq j} x_\ell\} = \frac{e^{\mu_j/\beta}}{\sum_{\ell=1}^n e^{\mu_\ell/\beta}}$.
5. The random variable $\max\{x_1, x_2\}$ is independent of $\mathbb{1}\{x_1 > x_2\}$.

The next lemma shows that there exists a unique distribution such that when a random variable following this distribution is added to an independent $\text{Gumbel}(0, \beta)$ random variable, with $\beta < 1$, the resulting sum is a standard Gumbel random variable.

Lemma A.2 (Theorem 2.1 in Cardell 1997) *Suppose $x \sim \text{Gumbel}(0, \beta)$ with $0 < \beta < 1$ and there is another random variable y independent of x . Then, $x + y$ has a standard Gumbel distribution if and only if y has a density function $f_\beta(y) = \frac{1}{\beta} \sum_{k=0}^{\infty} \frac{(-1)^k e^{-ky}}{k! \Gamma(-\beta k)}$ for all $y \in \mathbb{R}$.*

As an immediate corollary, there exists an independent random variable that can be added to another Gumbel random variable to obtain a Gumbel distribution with a higher scaling parameter.

Corollary A.3 (Changing scale through addition) *Suppose $x \sim \text{Gumbel}(\mu, \nu)$, with $\mu \in \mathbb{R}$ and $\nu > 0$, and we are given $\lambda > 0$ such that $\lambda \geq \nu$. Then, there exists a random variable y such that y is independent of x and $x + y \sim \text{Gumbel}(\mu, \lambda)$.*

The requirement that $\nu \leq \lambda$ in the above corollary is necessary. Since x and y are independent, we have that $\text{Var}(x + y) = \text{Var}(x) + \text{Var}(y)$, which implies that $\pi^2 \lambda^2 / 6 = \pi^2 \nu^2 / 6 + \text{Var}(y)$. Because the variance of a random variable is always non-negative, it must always be true that $\lambda \geq \nu$.

For each leaf node $\ell \in \mathcal{N}$ and each node j that is an ancestor of ℓ , recall that $\text{path}(j, \ell]$ denotes the nodes on the unique path from j to ℓ , excluding j . Define

$$z_{j,\ell} = v_j + \sum_{k \in \text{path}(j,\ell]} v_k + \mu_\ell,$$

where we define $v_{\text{root}} = 0$. The following lemma describes the distribution of $z_{j,\ell}$.

Lemma A.4 *For each leaf node ℓ and its ancestor j such that $j \neq \text{root}$, $z_{j,\ell} \sim \text{Gumbel}(\mu_\ell, \lambda_{\text{pa}(j)})$.*

Proof: Fix an arbitrary leaf node ℓ . We prove this result by induction on the height¹¹ of the ancestor j . The base case is when the height of j is zero, which means that $j = \ell$. By definition, $z_{j,\ell} = v_\ell + \mu_\ell \sim \text{Gumbel}(\mu_\ell, \lambda_{\text{pa}(j)})$, which is the desired result.

To establish the induction step, we assume that $z_{j,\ell} \sim \text{Gumbel}(\mu_\ell, \lambda_{\text{pa}(j)})$ for all ancestors j of ℓ whose heights are at most H . Consider an ancestor j with height $H + 1$. Let k denote the child of j that is also an ancestor of ℓ ; note that the height of k is at most H . By definition, $z_{j,\ell} = v_j + z_{k,\ell}$. Since k has a height of H , it follows from the induction hypothesis that $z_{k,\ell} \sim \text{Gumbel}(\mu_\ell, \lambda_{\text{pa}(k)})$. As j is the parent of k , we have that $\lambda_{\text{pa}(k)} = \lambda_j$ and $z_{k,\ell} \sim \text{Gumbel}(\mu_\ell, \lambda_j)$. The definition of v_j then implies that $z_{j,\ell} = v_j + z_{k,\ell} \sim \text{Gumbel}(\mu_\ell, \lambda_{\text{pa}(j)})$, completing the induction step. This completes the proof. ■

¹¹ Recall that the height of a node is the number of edges on the longest path between that node and a leaf.

For each $\mathcal{S} \subseteq \mathcal{N}$ and each node $j \in \mathbb{T}[\mathcal{S}]$, let the $Z_j(\mathcal{S})$ denote the maximum of random variables $z_{j,\ell}$ over all the leaf nodes in $\mathbb{T}[\mathcal{S}]$ that are descendants of j ; that is,

$$Z_j(\mathcal{S}) \stackrel{\text{def}}{=} \max \{z_{j,\ell} : \ell \text{ is a leaf node in } \mathbb{T}_j[\mathcal{S}]\},$$

where we define the maximum over an empty set to be minus infinity. The next lemma characterizes the distribution of $Z_j(\mathcal{S})$.

Lemma A.5 *For each subset $\mathcal{S} \subseteq \mathcal{N}$ and each node $j \in \mathbb{T}[\mathcal{S}]$ such that $j \neq \text{root}$, $Z_j(\mathcal{S}) \sim \text{Gumbel}(W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}), \lambda_{\text{pa}(j)})$.*

Proof: Without loss of generality, we prove the result for the case $\mathcal{S} = \mathcal{N}$, so $\mathbb{T}[\mathcal{S}] = \mathbb{T}$. The proof for a general subset \mathcal{S} follows from an identical argument applied on the sub-tree $\mathbb{T}[\mathcal{S}]$. Since we consider the full assortment, we will drop references to \mathcal{S} , and simply write Z_j and \mathbb{T}_j . We will prove the result by induction on the height of node j . For the base case, suppose that the height of j is zero. So, j is a leaf node. Then,

$$Z_j = \max \{z_{j,\ell} : \ell \text{ is a leaf node in } \mathbb{T}_j\} = z_{j,j} = \mu_j + v_j \sim \text{Gumbel}(W_j(\boldsymbol{\mu}, \boldsymbol{\lambda}), \lambda_{\text{pa}(j)}),$$

where we use the fact that if j is a leaf node, then $v_j \sim \text{Gumbel}(0, \lambda_{\text{pa}(j)})$ and $W_j(\boldsymbol{\mu}, \boldsymbol{\lambda}) = \mu_j$. This completes the base case.

To establish the induction step, we assume that the result holds for all nodes with height of at most H . We now consider node j with height of $H + 1$. Since $Z_j = \max \{z_{j,\ell} : \ell \text{ is a leaf node in } \mathbb{T}_j\}$, it follows that that

$$Z_j = v_j + \max_{k \in \text{Children}(j)} \{ \max \{z_{k,\ell} : \ell \text{ is a leaf node in } \mathbb{T}_k\} \} = v_j + \max_{k \in \text{Children}(j)} Z_k.$$

For each $k \in \text{Children}(j)$, the height of k is at most H . So, invoking the induction hypothesis, we obtain that for all $k \in \text{Children}(j)$,

$$Z_k = \max \{z_{k,\ell} : \ell \text{ is a leaf node in } \mathbb{T}_k\} \sim \text{Gumbel}(W_k(\boldsymbol{\mu}, \boldsymbol{\lambda}), \lambda_j)$$

because $\lambda_{\text{pa}(k)} = \lambda_j$. Further, because the vertex sets of \mathbb{T}_k and $\mathbb{T}_{k'}$ are disjoint for any $k \neq k'$ such that $\{k, k'\} \subseteq \text{Children}(j)$, we have that Z_k is independent of $Z_{k'}$. It then follows from the properties of the Gumbel distribution that

$$\max_{k \in \text{Children}(j)} Z_k \sim \text{Gumbel} \left(\lambda_j \log \left(\sum_{k \in \text{Children}(j)} e^{W_k(\boldsymbol{\mu}, \boldsymbol{\lambda})/\lambda_j} \right), \lambda_j \right) = \text{Gumbel}(W_j(\boldsymbol{\mu}, \boldsymbol{\lambda}), \lambda_j),$$

where the equality follows from the definition of $W_j(\boldsymbol{\mu}, \boldsymbol{\lambda})$. It thus follows from the definition of v_j that $Z_j \sim \text{Gumbel}(W_j(\boldsymbol{\mu}, \boldsymbol{\lambda}), \lambda_{\text{pa}(j)})$. This establishes the induction step, completing the proof. ■

The next lemma allows us to simplify the conditional probability involving $Z_j(\mathcal{S})$.

Lemma A.6 For each subset $\mathcal{S} \subseteq \mathcal{N}$ and each node $j \in \mathbb{T}[\mathcal{S}]$ such that $j \neq \text{root}$,

$$\begin{aligned} & \mathbb{P} \left\{ Z_j(\mathcal{S}) > \max_{k \in \text{Sibling}(j) \setminus \{j\}} Z_k(\mathcal{S}) \mid Z_i(\mathcal{S}) > \max_{v \in \text{Sibling}(i) \setminus \{i\}} Z_v(\mathcal{S}) \quad \forall i \in \text{path}(\text{root}, \text{pa}(j)) \right\} \\ &= \mathbb{P} \left\{ Z_j(\mathcal{S}) > \max_{k \in \text{Sibling}(j) \setminus \{j\}} Z_k(\mathcal{S}) \right\}, \end{aligned}$$

where $\text{Sibling}(j) = \{k \in \mathbb{T}[\mathcal{S}] : \text{pa}(k) = \text{pa}(j)\}$ denote the siblings of j in the tree $\mathbb{T}[\mathcal{S}]$.

Proof: Without loss of generality, assume that $\mathcal{S} = \mathcal{N}$. The proof for the general \mathcal{S} is essentially the same. Since we consider the full assortment, we will drop references to \mathcal{S} . Fix an arbitrary node $j \in \mathbb{T}$ such that $j \neq \text{root}$. Let $x_1 = Z_j$ and $x_2 = \max_{k \in \text{Sibling}(j) \setminus \{j\}} Z_k$. Note that x_1 and x_2 are independent of each other. We will first establish the following claim by induction.

Claim: For every node $i \in \text{path}(\text{root}, \text{pa}(j))$, there is a deterministic function f_i such that $Z_i = f_i(\max\{x_1, x_2\}, \mathbf{U}_i)$, where the random vector \mathbf{U}_i is independent of x_1 and x_2 .

We will prove the claim by induction on node i . For the base case, consider $i = \text{pa}(j)$. Then, by definition,

$$Z_{\text{pa}(j)} = v_{\text{pa}(j)} + \max_{k \in \text{Children}(\text{pa}(j))} Z_k = v_{\text{pa}(j)} + \max \left\{ Z_j, \max_{k \in \text{Sibling}(j) \setminus \{j\}} Z_k \right\} = v_{\text{pa}(j)} + \max\{x_1, x_2\},$$

and the result follows because $v_{\text{pa}(j)}$ is independent of x_1 and x_2 . This proves the base case.

For the induction step, assume the result holds for some node $i \in \text{path}(\text{root}, \text{pa}(j))$. We will now prove that it also holds for node $\text{pa}(i)$. By definition,

$$\begin{aligned} Z_{\text{pa}(i)} &= v_{\text{pa}(i)} + \max_{v \in \text{Children}(\text{pa}(i))} Z_v = v_{\text{pa}(i)} + \max \left\{ Z_i, \max_{v \in \text{Sibling}(i) \setminus \{i\}} Z_v \right\} \\ &= v_{\text{pa}(i)} + \max \left\{ f_i(\max\{x_1, x_2\}, \mathbf{U}_i), \max_{v \in \text{Sibling}(i) \setminus \{i\}} Z_v \right\} \end{aligned}$$

and the desired result follows because $v_{\text{pa}(i)}$ is independent of x_1 , x_2 , and \mathbf{U}_i . Moreover, $\max_{v \in \text{Sibling}(i) \setminus \{i\}} Z_v$ is also independent of x_1 and x_2 because for each $v \in \text{Sibling}(i) \setminus \{i\}$, the sub-tree \mathbb{T}_v rooted at v is completely separate from the sub-tree \mathbb{T}_k for all $k \in \text{Sibling}(j)$. This completes the induction, establishing the claim.

It follows from the above claim that

$$\begin{aligned} & \mathbb{P} \left\{ Z_j > \max_{k \in \text{Sibling}(j) \setminus \{j\}} Z_k \mid Z_i > \max_{v \in \text{Sibling}(i) \setminus \{i\}} Z_v \quad \forall i \in \text{path}(\text{root}, \text{pa}(j)) \right\} \\ &= \mathbb{P} \left\{ x_1 > x_2 \mid f_i(\max\{x_1, x_2\}, \mathbf{U}_i) > \max_{v \in \text{Sibling}(i) \setminus \{i\}} Z_v \quad \forall i \in \text{path}(\text{root}, \text{pa}(j)) \right\} \\ &= \mathbb{P} \{x_1 > x_2\}, \end{aligned}$$

where the last equality follows from Lemma A.1, which shows that $\mathbb{1}\{x_1 > x_2\}$ is independent of $\max\{x_1, x_2\}$. Also, note that for all $i \in \text{path}(\text{root}, \text{pa}(j))$, $\max_{v \in \text{Sibling}(i) \setminus \{i\}} Z_v$ is independent of x_1 and x_2 . This completes the induction, proving the desired result.

Finally, here is the proof of Theorem 2.1.

Proof of Theorem 2.1: First note that since $\lambda_j \leq \lambda_{\text{pa}(j)}$, it follows from Corollary A.3 that v_j exists for all non-leaf nodes j . Next, fix an arbitrary subset \mathcal{S} and $\ell \in \mathcal{S}$. Recall that $\text{Sibling}(j) = \{k \in \mathbb{T}[\mathcal{S}] : \text{pa}(k) = \text{pa}(j)\}$ denote the siblings of j in the tree $\mathbb{T}[\mathcal{S}]$. Let $\text{root} \rightarrow j_1 \rightarrow j_2 \rightarrow \dots \rightarrow j_m \rightarrow \ell$ denote the unique path from root to ℓ in $\mathbb{T}[\mathcal{S}]$. Note that

$$\{\ell\} = \mathbb{T}_\ell[\mathcal{S}] \cap \mathcal{S} \subseteq \mathbb{T}_{j_m}[\mathcal{S}] \cap \mathcal{S} \subseteq \mathbb{T}_{j_{m-1}}[\mathcal{S}] \cap \mathcal{S} \dots \subseteq \mathbb{T}_{j_1}[\mathcal{S}] \cap \mathcal{S} \subseteq \mathbb{T}_{\text{root}}[\mathcal{S}] \cap \mathcal{S} = \mathcal{S}$$

Note that the event $\text{utility}_\ell > \max_{k \in \mathcal{S} \setminus \{\ell\}} \text{utility}_k$ happens if and only if for every node $j \in \text{path}(\text{root}, \ell]$, we have $\max_{k \in \mathbb{T}_j[\mathcal{S}] \cap \mathcal{S}} \text{utility}_k > \max_{i \in \text{Sibling}(j) \setminus \{j\}} \max_{k \in \mathbb{T}_i[\mathcal{S}] \cap \mathcal{S}} \text{utility}_k$. Therefore,

$$\begin{aligned} & \mathbb{P} \left\{ \text{utility}_\ell > \max_{k \in \mathcal{S} \setminus \{\ell\}} \text{utility}_k \right\} \\ &= \mathbb{P} \left\{ \max_{k \in \mathbb{T}_j[\mathcal{S}] \cap \mathcal{S}} \text{utility}_k > \max_{i \in \text{Sibling}(j) \setminus \{j\}} \max_{k \in \mathbb{T}_i[\mathcal{S}] \cap \mathcal{S}} \text{utility}_k \quad \forall j \in \text{path}(\text{root}, \ell] \right\} \\ &= \mathbb{P} \left\{ Z_j(\mathcal{S}) > \max_{i \in \text{Sibling}(j) \setminus \{j\}} Z_i(\mathcal{S}) \quad \forall j \in \text{path}(\text{root}, \ell] \right\}, \end{aligned}$$

where the last equality follows because for each $i \in \text{Sibling}(j)$

$$\max_{k \in \mathbb{T}_i[\mathcal{S}] \cap \mathcal{S}} \text{utility}_k = \sum_{v \in \text{path}(\text{root}, \text{pa}(j))} v_v + \max_{k \in \mathbb{T}_i[\mathcal{S}] \cap \mathcal{S}} Z_{i,k} = \sum_{v \in \text{path}(\text{root}, \text{pa}(j))} v_v + Z_i(\mathcal{S}),$$

and the term $\sum_{v \in \text{path}(\text{root}, \text{pa}(j))} v_v$ is common for all nodes $i \in \text{Sibling}(j)$.

For any collection of random variables x_1, \dots, x_n , $\mathbb{P}\{x_1, \dots, x_n\} = \prod_{j=1}^n \mathbb{P}\{x_j \mid x_{j-1}, \dots, x_1\}$. Thus,

$$\begin{aligned} & \mathbb{P} \left\{ Z_j(\mathcal{S}) > \max_{i \in \text{Sibling}(j) \setminus \{j\}} Z_i(\mathcal{S}) \quad \forall j \in \text{path}(\text{root}, \ell] \right\} \\ &= \prod_{j \in \text{path}(\text{root}, \ell]} \mathbb{P} \left\{ Z_j(\mathcal{S}) > \max_{i \in \text{Sibling}(j) \setminus \{j\}} Z_i(\mathcal{S}) \mid Z_k > \max_{v \in \text{Sibling}(k) \setminus \{k\}} Z_k \quad \forall k \in \text{path}(\text{root}, \text{pa}(j)) \right\} \\ &= \prod_{j \in \text{path}(\text{root}, \ell]} \mathbb{P} \left\{ Z_j(\mathcal{S}) > \max_{i \in \text{Sibling}(j) \setminus \{j\}} Z_i(\mathcal{S}) \right\} \end{aligned}$$

where last equality follows from Lemma A.6.

By Lemma A.5, $Z_j(\mathcal{S}) \sim \text{Gumbel}(W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}), \lambda_{\text{pa}(j)})$, and for each $i \in \text{Sibling}(j) \setminus \{j\}$, $Z_i(\mathcal{S}) \sim \text{Gumbel}(W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}), \lambda_{\text{pa}(i)})$. Since $\lambda_{\text{pa}(i)} = \lambda_{\text{pa}(j)}$ for all $i \in \text{Sibling}(j)$,

$$\mathbb{P} \left\{ Z_j(\mathcal{S}) > \max_{i \in \text{Sibling}(j) \setminus \{j\}} Z_i(\mathcal{S}) \right\} = \frac{e^{W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) / \lambda_{\text{pa}(j)}}}{\sum_{i \in \text{Sibling}(j)} e^{W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) / \lambda_{\text{pa}(j)}}},$$

which implies that

$$\mathbb{P} \left\{ utility_\ell > \max_{k \in \mathcal{S} \setminus \{\ell\}} utility_k \right\} = \prod_{j \in \text{path}(\text{root}, \ell]} \frac{e^{W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) / \lambda_{\text{pa}(j)}}}{\sum_{i \in \text{Sibling}(j)} e^{W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) / \lambda_{\text{pa}(j)}}} = \psi_{\text{root} \rightarrow \ell}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = \mathbb{P}_\ell(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}),$$

and this completes the proof. \blacksquare

Finally, we show that the error terms $(\varepsilon_\ell : \ell \in \mathcal{N})$ are identically distributed but are not independent of each other:

Lemma A.7 (Distribution of error terms) *For all $\ell \in \mathcal{N}$, $\varepsilon_\ell \sim \text{Gumbel}(0, \lambda_{\text{root}})$. Moreover, for all $\ell, \ell' \in \mathcal{N}$ such that $\ell \neq \ell'$*

$$\text{Corr}(\varepsilon_\ell, \varepsilon_{\ell'}) = \frac{\text{Cov}(\varepsilon_\ell, \varepsilon_{\ell'})}{\sqrt{\text{Var}(\varepsilon_\ell) \text{Var}(\varepsilon_{\ell'})}} = 1 - \left(\frac{\lambda_j}{\lambda_{\text{root}}} \right)^2,$$

where j is the nearest common ancestor of ℓ and ℓ' in \mathbb{T} .

Proof: It follows from the definition of the random variables $(z_{j,\ell} : j \in \mathbb{T} \setminus \mathcal{N})$ before Lemma A.4 that $\varepsilon_\ell = z_{\text{root},\ell} - \mu_\ell$ for all $\ell \in \mathcal{N}$. Further since $v_{\text{root}} = 0$, it follows from Lemma A.4 that $z_{\text{root},\ell} \sim \text{Gumbel}(\mu_\ell, \lambda_{\text{root}})$. Combining this with property 2 of Lemma A.1, it follows that $\varepsilon_\ell \sim \text{Gumbel}(0, \lambda_{\text{root}})$ for all $\ell \in \mathcal{N}$. This establishes the first part of the lemma.

Next, since $\varepsilon_\ell = \sum_{k \in \text{path}(\text{root}, \ell]} v_k$ for any $\ell \in \mathcal{N}$, and all the v_k 's are independent, it follows that $\text{Cov}(\varepsilon_\ell, \varepsilon_{\ell'}) = \sum_{k \in \text{path}(\text{root}, j]} \text{Cov}(v_k, v_k) = \sum_{k \in \text{path}(\text{root}, j]} \text{Var}(v_k)$, where j is the common ancestor nearest to both ℓ and ℓ' ; that is, $\{\ell, \ell'\} \subseteq \mathbb{T}_j$ but $\{\ell, \ell'\} \not\subseteq \mathbb{T}_k$ for all $k \in \text{Children}(j)$. Note that if $j = \text{root}$, then $\text{Cov}(\varepsilon_\ell, \varepsilon_{\ell'}) = 0$ which is consistent with the expression in the statement of the lemma. So suppose $j \neq \text{root}$. Then, it follows from our definitions that $\text{Var}(v_k) = \pi^2 \cdot (\lambda_{\text{pa}(k)}^2 - \lambda_k^2) / 6$ for each non-leaf node $k \in \mathbb{T} \setminus (\text{root} \cup \mathcal{N})$. Then, we obtain by telescoping that the covariance $\text{Cov}(\varepsilon_\ell, \varepsilon_{\ell'}) = \pi^2 \cdot (\lambda_{\text{root}}^2 - \lambda_j^2) / 6$. Finally, since $\text{Var}(\varepsilon_\ell) = \text{Var}(\varepsilon_{\ell'}) = \pi^2 \lambda_{\text{root}}^2 / 6$, the result follows. \blacksquare

Appendix B: Properties of the Weight Function

The following lemmas establish important properties of the weight function $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$, which we will use repeatedly to establish properties of the negative log-likelihood function. We first state all the lemmas and then provide their proofs. The first lemma establishes positive homogeneity and additivity.

Lemma B.1 (Positive Homogeneity and Additivity) *For each subset $\mathcal{S} \subseteq \mathcal{N}$, node $j \in \mathbb{T}[\mathcal{S}]$, $\alpha > 0$, and $\xi \in \mathbb{R}$, $W_j(\mathcal{S}; \alpha \boldsymbol{\mu}, \alpha \boldsymbol{\lambda}) = \alpha W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ and $W_j(\mathcal{S}; \boldsymbol{\mu} + \xi \mathbf{e}, \boldsymbol{\lambda}) = \xi + W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$, where \mathbf{e} is the vector of all ones.*

The next lemma establishes convexity, and under additional assumptions, strict convexity. Recall that for each node $j \in \mathbb{T}[\mathcal{S}]$, $\mathbb{T}_j[\mathcal{S}]$ denotes the sub-tree of $\mathbb{T}[\mathcal{S}]$ rooted at j ; that is, $\mathbb{T}_j[\mathcal{S}]$ consists of the node j and all of its descendant in $\mathbb{T}[\mathcal{S}]$.

Lemma B.2 (Convexity and Strict Convexity) *For each subset $\mathcal{S} \subseteq \mathcal{N}$ and $j \in \mathbb{T}[\mathcal{S}]$, the function $(\boldsymbol{\mu}, \boldsymbol{\lambda}) \mapsto W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ satisfies the following properties:*

- (a) *It is convex in $(\boldsymbol{\mu}, \boldsymbol{\lambda})$.*
- (b) *For each $\boldsymbol{\lambda}$ and $0 < \theta < 1$,*

$$W_j(\mathcal{S}; \theta \boldsymbol{\mu} + (1 - \theta) \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) = \theta W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) + (1 - \theta) W_j(\mathcal{S}; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda})$$

if and only if there exists $\xi \in \mathbb{R}$ such that $\bar{\mu}_\ell = \mu_\ell + \xi$ for all leaf nodes $\ell \in \mathbb{T}_j[\mathcal{S}] \cap \mathcal{S}$.

- (c) *For each $a \in \mathbb{R}$ and $\boldsymbol{\lambda}$, the function $\boldsymbol{\mu} \mapsto W_{\text{root}}(\mathcal{N}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ is strictly convex on the set $\{\boldsymbol{\mu} : \mu_1 = a\}$.*

The next lemma establishes the monotonicity of the weight function.

Lemma B.3 (Monotonicity and Strict Monotonicity) *For each $\mathcal{S} \subseteq \mathcal{N}$ and $j \in \mathbb{T}[\mathcal{S}]$, the function $(\boldsymbol{\mu}, \boldsymbol{\lambda}) \mapsto W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ satisfies the following properties.*

- (a) *It is increasing in $(\boldsymbol{\mu}, \boldsymbol{\lambda})$; that is, if $(\boldsymbol{\mu}, \boldsymbol{\lambda}) \leq (\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}})$ where the inequality holds componentwise, then $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \leq W_j(\mathcal{S}; \bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}})$.*
- (b) *It is strictly increasing in μ_ℓ for each leaf node $\ell \in \mathbb{T}_j[\mathcal{S}] \cap \mathcal{S}$.*
- (c) *It is strictly increasing in λ_k for each $k \in \mathbb{T}_j[\mathcal{S}]$ such that k has at least two children in $\mathbb{T}_j[\mathcal{S}]$.*

Finally, the following lemma provides an expression for the derivative of the weight function with respect to the model parameters.

Lemma B.4 (Derivatives of the Weight Functions) *For each $\mathcal{S} \subseteq \mathcal{N}$ and $j \in \mathbb{T}[\mathcal{S}]$,*

$$\begin{aligned} \frac{\partial W_j}{\partial \mu_\ell}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) &= \psi_{j \rightarrow \ell}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) & \forall \ell \in \mathcal{N} \\ \frac{\partial W_j}{\partial \lambda_k}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) &= \psi_{j \rightarrow k}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \Delta_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) & \forall k \in \mathbb{T} \setminus \mathcal{N}, k \neq \text{root} . \end{aligned}$$

where for each non-leaf node k such that $k \neq \text{root}$,

$$\begin{aligned} \Delta_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) &= \log \left(\sum_{i \in \text{Children}(k) \cap \mathbb{T}[\mathcal{S}]} e^{W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) / \lambda_k} \right) - \frac{\sum_{i \in \text{Children}(k) \cap \mathbb{T}[\mathcal{S}]} e^{W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) / \lambda_k} \times W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_k \sum_{i \in \text{Children}(k) \cap \mathbb{T}[\mathcal{S}]} e^{W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) / \lambda_k}} \\ &= \frac{W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) - \sum_{i \in \text{Children}(k) \cap \mathbb{T}[\mathcal{S}]} \psi_{k \rightarrow i}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \times W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_k} . \end{aligned}$$

We now present the proofs of these four lemmas.

B.1 Proof of Lemma B.1

Proof: Consider an arbitrary subset \mathcal{S} . We prove both results by induction on the height of node j . For the base case, suppose the height of j is zero. This means that j is a leaf node of $\mathbb{T}[\mathcal{S}]$, so $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = \mu_j$ and the result is trivially true. Suppose that the results holds for all nodes at height at most H . Now consider a vertex $j \in \mathbb{T}[\mathcal{S}]$ of height $H + 1$. By the induction hypothesis, the results are true for all children k of node j because the height of k is at most H . Therefore, by definition of W_j and the inductive hypothesis,

$$W_j(\mathcal{S}; \alpha\boldsymbol{\mu}, \alpha\boldsymbol{\lambda}) = \alpha\lambda_j \log \left(\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{\alpha W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) / (\alpha\lambda_j)} \right) = \alpha W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \text{ and}$$

$$W_j(\mathcal{S}; \boldsymbol{\mu} + \xi \mathbf{e}, \boldsymbol{\lambda}) = \lambda_j \log \left(\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{[\xi + W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})] / \lambda_j} \right) = \xi + W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}).$$

We have thus established the result for nodes at height $H + 1$, completing the induction step. The result of the lemma now follows. \blacksquare

B.2 Proof of Lemma B.2

Proof: Fix an arbitrary subset \mathcal{S} . We prove each of the three parts separately. Throughout this proof, we make use of the following observation: for each $j \in \mathbb{T}[\mathcal{S}]$, the function $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ only depends the parameters associated with nodes in the sub-tree $\mathbb{T}_j[\mathcal{S}]$ rooted at j ; that is, $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = W_j(\mathcal{S}; \boldsymbol{\mu}_{\mathbb{T}_j[\mathcal{S}]}, \boldsymbol{\lambda}_{\mathbb{T}_j[\mathcal{S}]})$, where $\boldsymbol{\mu}_{\mathbb{T}_j[\mathcal{S}]} = (\mu_\ell : \ell \in \mathbb{T}_j[\mathcal{S}])$ and $\boldsymbol{\lambda}_{\mathbb{T}_j[\mathcal{S}]} = (\lambda_k : k \in \mathbb{T}_j[\mathcal{S}])$.

Proof of part (a): We use induction on the height of node j . For the base case, suppose that the height of j is zero, so j is a leaf node. Then, the result is trivially true by definition. Suppose that $W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ is convex in $(\boldsymbol{\mu}, \boldsymbol{\lambda})$ for all nodes k with height at most H . Now, consider an arbitrary non-leaf node j at height $H + 1$. By the induction hypothesis, $W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ is convex in $(\boldsymbol{\mu}, \boldsymbol{\lambda})$ for all children k of node j . Then, the function

$$(\boldsymbol{\mu}, \boldsymbol{\lambda}) \mapsto \log \left(\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})} \right)$$

must also be convex in $(\boldsymbol{\mu}, \boldsymbol{\lambda})$ because it is a composition of the log-sum-exp function, which is increasing and convex, with a collection of convex functions $(\boldsymbol{\mu}, \boldsymbol{\lambda}) \mapsto W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$, for $k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]$. Note that for each $k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]$, the function $W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ is *independent* of λ_j because $j \notin \mathbb{T}_k[\mathcal{S}]$, so

$$\log \left(\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})} \right) = \log \left(\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_k(\mathcal{S}; \boldsymbol{\mu}_{\mathbb{T}_k[\mathcal{S}]}, \boldsymbol{\lambda}_{\mathbb{T}_k[\mathcal{S}]})} \right).$$

Therefore, the perspective of the above function is given by

$$\begin{aligned} (\boldsymbol{\mu}, \boldsymbol{\lambda}) \mapsto \lambda_j \log \left(\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_k \left(\mathcal{S}; \frac{\boldsymbol{\mu}_{\mathbb{T}_k[\mathcal{S}]}}{\lambda_j}, \frac{\boldsymbol{\lambda}_{\mathbb{T}_k[\mathcal{S}]}}{\lambda_j} \right)} \right) &= \lambda_j \log \left(\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_k(\mathcal{S}; \boldsymbol{\mu}_{\mathbb{T}_k[\mathcal{S}]}, \boldsymbol{\lambda}_{\mathbb{T}_k[\mathcal{S}]}) / \lambda_j} \right) \\ &= W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}), \end{aligned}$$

where the first equality follows from the positive homogeneity property in Lemma B.1. Because the perspective of a convex function is also convex (Boyd and Vandenberghe 2004), it follows that $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ is convex in $(\boldsymbol{\mu}, \boldsymbol{\lambda})$, completing the induction step. This proves part (a).

Proof of part (b): We prove part (b) by induction on the height of node j . Consider the base case where j has a height of zero, so $j = \ell$ for some leaf node ℓ in $\mathbb{T}[\mathcal{S}]$. In this case, $\mathbb{T}_\ell[\mathcal{S}] = \{\ell\}$ and $W_\ell(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = \mu_\ell$. Because $\mathbb{T}_\ell[\mathcal{S}] \cap \mathcal{S}$ contains only the node ℓ in it, the condition $\bar{\mu}_\ell = \mu_\ell + \xi$ can be trivially satisfied by choosing $\xi = \bar{\mu}_\ell - \mu_\ell$. Further, the relationship

$$W_\ell(\mathcal{S}; \theta \boldsymbol{\mu} + (1 - \theta) \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) = \theta W_\ell(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) + (1 - \theta) W_\ell(\mathcal{S}; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda})$$

is also trivially satisfied by definition for all $\theta \in (0, 1)$. Therefore, the base case is true.

Suppose that the result holds for all nodes with height at most H . Consider an arbitrary non-leaf node j at height $H + 1$. By the induction hypothesis, we have that for each $k \in \text{Children}(j)$,

$$W_k(\mathcal{S}; \theta \boldsymbol{\mu} + (1 - \theta) \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) = \theta W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) + (1 - \theta) W_k(\mathcal{S}; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda})$$

if and only if there exists $\xi_k \in \mathbb{R}$ such that $\mu_\ell = \bar{\mu}_\ell + \xi_k$ for all $\ell \in \mathbb{T}_k[\mathcal{S}] \cap \mathcal{S}$. We will now prove the result at node j .

We will first prove the sufficiency. Suppose that there exists $\xi \in \mathbb{R}$ such that $\bar{\mu}_\ell = \mu_\ell + \xi$ for all leaf nodes $\ell \in \mathbb{T}_j[\mathcal{S}] \cap \mathcal{S}$, or equivalently $\bar{\boldsymbol{\mu}}_{\mathbb{T}_j[\mathcal{S}] \cap \mathcal{S}} = \boldsymbol{\mu}_{\mathbb{T}_j[\mathcal{S}] \cap \mathcal{S}} + \xi \mathbf{e}$, where \mathbf{e} is a vector of ones of an appropriate dimension. Then,

$$\begin{aligned} W_j(\mathcal{S}; \theta \boldsymbol{\mu} + (1 - \theta) \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) &= W_j(\mathcal{S}; \boldsymbol{\mu}_{\mathbb{T}_j[\mathcal{S}]} + \xi(1 - \theta) \mathbf{e}, \boldsymbol{\lambda}) && [W_j \text{ only depends on } \boldsymbol{\mu}_{\mathbb{T}_j[\mathcal{S}]}] \\ &= \xi(1 - \theta) + W_j(\mathcal{S}; \boldsymbol{\mu}_{\mathbb{T}_j[\mathcal{S}]}, \boldsymbol{\lambda}) && [\text{by Lemma B.1}] \\ &= \theta W_j(\mathcal{S}; \boldsymbol{\mu}_{\mathbb{T}_j[\mathcal{S}]}, \boldsymbol{\lambda}) + (1 - \theta) \left(\xi + W_j(\mathcal{S}; \boldsymbol{\mu}_{\mathbb{T}_j[\mathcal{S}]}, \boldsymbol{\lambda}) \right) \\ &= \theta W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) + (1 - \theta) W_j(\mathcal{S}; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}), \end{aligned}$$

where the last equality follows because $\bar{\boldsymbol{\mu}}_{\mathbb{T}_j[\mathcal{S}] \cap \mathcal{S}} = \boldsymbol{\mu}_{\mathbb{T}_j[\mathcal{S}] \cap \mathcal{S}} + \xi \mathbf{e}$. This gives the desired result.

We will now prove the necessity. Suppose that $\boldsymbol{\mu}$, $\bar{\boldsymbol{\mu}}$, and $\theta \in (0, 1)$ are such that

$$W_j(\mathcal{S}; \theta \boldsymbol{\mu} + (1 - \theta) \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) = \theta W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) + (1 - \theta) W_j(\mathcal{S}; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}). \quad (\text{EC.1})$$

Our goal is to exhibit a $\xi \in \mathbb{R}$ such that $\bar{\boldsymbol{\mu}}_{\mathcal{T}_j[\mathcal{S}] \cap \mathcal{S}} + \xi = \boldsymbol{\mu}_{\mathcal{T}_j[\mathcal{S}] \cap \mathcal{S}}$. If we have that $\bar{\boldsymbol{\mu}}_{\mathcal{T}_j[\mathcal{S}] \cap \mathcal{S}} = \boldsymbol{\mu}_{\mathcal{T}_j[\mathcal{S}] \cap \mathcal{S}}$, then the result is trivially true. Therefore, we assume that $\bar{\boldsymbol{\mu}}_{\mathcal{T}_j[\mathcal{S}] \cap \mathcal{S}} \neq \boldsymbol{\mu}_{\mathcal{T}_j[\mathcal{S}] \cap \mathcal{S}}$. To simplify notation, let $a = 1/\lambda_j$, $m = |\text{Children}(j)|$. Further, let $\text{LSE}: \mathbb{R}^m \rightarrow \mathbb{R}_{++}$ be defined as $\text{LSE}(\mathbf{x}) = \log(\sum_{i=1}^m e^{x_i})$, for each $\mathbf{x} \in \mathbb{R}^m$. Note that $\text{LSE}(\cdot)$ is the standard log-sum-exp function. Also, define the following three vectors in \mathbb{R}_+^m :

$$\begin{aligned} \mathbf{y}^1 &= (W_k(\mathcal{S}; \theta \boldsymbol{\mu} + (1 - \theta) \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}): k \in \text{Children}(j) \cap \mathcal{T}[\mathcal{S}]) \\ \mathbf{y}^2 &= (W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}): k \in \text{Children}(j) \cap \mathcal{T}[\mathcal{S}]) \\ \mathbf{y}^3 &= (W_k(\mathcal{S}; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}): k \in \text{Children}(j) \cap \mathcal{T}[\mathcal{S}]) \end{aligned}$$

By convexity of W_k from part (a), we have that $\mathbf{y}^1 \leq \theta \mathbf{y}^2 + (1 - \theta) \mathbf{y}^3$, where the inequality holds componentwise. By definition,

$$\begin{aligned} \frac{1}{\lambda_j} W_j(\mathcal{S}; \theta \boldsymbol{\mu} + (1 - \theta) \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) &= \text{LSE}(a \mathbf{y}^1) \\ &\leq \text{LSE}(\theta a \mathbf{y}^2 + (1 - \theta) a \mathbf{y}^3) && [\text{LSE is strictly increasing}] \\ &\leq \theta \text{LSE}(a \mathbf{y}^2) + (1 - \theta) \text{LSE}(a \mathbf{y}^3) && [\text{LSE is convex}] \\ &= \frac{\theta}{\lambda_j} W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) + \frac{1 - \theta}{\lambda_j} W_j(\mathcal{S}; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) \\ &= \frac{1}{\lambda_j} W_j(\mathcal{S}; \theta \boldsymbol{\mu} + (1 - \theta) \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) && [\text{by hypothesis (EC.1)}]. \end{aligned}$$

Because the left and right hand side expressions are equal to each other, it must be true that both inequalities hold with equalities. We have thus shown that

$$\text{LSE}(a \mathbf{y}^1) = \text{LSE}(\theta a \mathbf{y}^2 + (1 - \theta) a \mathbf{y}^3) = \theta \text{LSE}(a \mathbf{y}^2) + (1 - \theta) \text{LSE}(a \mathbf{y}^3)$$

We now derive the implications from the above two equalities. Because the LSE function is strictly increasing, the first equality implies that $\mathbf{y}^1 = \theta \mathbf{y}^2 + (1 - \theta) \mathbf{y}^3$. In other words, we have that for each $k \in \text{Children}(j)$,

$$y_k^1 = \theta y_k^2 + (1 - \theta) y_k^3 \iff W_k(\mathcal{S}; \theta \boldsymbol{\mu} + (1 - \theta) \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) = \theta W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) + (1 - \theta) W_k(\mathcal{S}; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}).$$

It now follows from the induction hypothesis, applied to W_k , that for each $k \in \text{Children}(j)$, there exists $\xi_k \in \mathbb{R}$ such that $\bar{\mu}_\ell = \mu_\ell + \xi_k$ for all $\ell \in \mathcal{T}_k[\mathcal{S}] \cap \mathcal{S}$. To establish our result, it is now sufficient to show that $\xi_k = \xi_{k'}$ for all $k \neq k'$ such that $\{k, k'\} \subseteq \text{Children}(j)$. For that, we first note that because $\bar{\boldsymbol{\mu}}_{\mathcal{T}_k[\mathcal{S}] \cap \mathcal{S}} = \xi_k + \boldsymbol{\mu}_{\mathcal{T}_k[\mathcal{S}] \cap \mathcal{S}}$, we have that

$$y_k^3 = W_k(\mathcal{S}; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) = W_k(\mathcal{S}; \bar{\boldsymbol{\mu}}_{\mathcal{T}_k \cap \mathcal{S}}, \boldsymbol{\lambda}) = \xi_k + W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = \xi_k + y_k^2, \quad (\text{EC.2})$$

where third equality follows from Lemma B.1. Now, if $\mathbf{y}^2 = \mathbf{y}^3$, then it follows that $\xi_k = 0$ for all $k \in \text{Children}(j)$, establishing the result. Therefore, we assume that $\mathbf{y}^2 \neq \mathbf{y}^3$.

Then, we focus on the second equality above:

$$\text{LSE}(\theta a\mathbf{y}^2 + (1-\theta)a\mathbf{y}^3) = \theta \text{LSE}(a\mathbf{y}^2) + (1-\theta)\text{LSE}(a\mathbf{y}^3). \quad (\text{EC.3})$$

It is a well-known result that the function $t \mapsto \text{LSE}(\mathbf{x} + t\mathbf{z})$ is strictly convex if and only if $\mathbf{z} \neq \mathbf{e}$. In other words, for some $\mathbf{x}, \mathbf{z} \in \mathbb{R}^m$, and for any $t_1, t_2 \in \mathbb{R}$, we have that

$$\text{LSE}(\mathbf{x} + (\theta t_1 + (1-\theta)t_2)\mathbf{z}) = \theta \text{LSE}(\mathbf{x} + t_1\mathbf{z}) + (1-\theta)\text{LSE}(\mathbf{x} + t_2\mathbf{z}) \quad \text{if and only if} \quad \mathbf{z} = \mathbf{e}.$$

Now choosing $\mathbf{x}, \mathbf{z}, t_1$, and t_2 to satisfy

$$\mathbf{x} + t_1\mathbf{z} = a\mathbf{y}^2 \quad \text{and} \quad \mathbf{x} + t_2\mathbf{z} = a\mathbf{y}^3,$$

we obtain from (EC.3) that $\text{LSE}(\mathbf{x} + (\theta t_1 + (1-\theta)t_2)\mathbf{z}) = \theta \text{LSE}(\mathbf{x} + t_1\mathbf{z}) + (1-\theta)\text{LSE}(\mathbf{x} + t_2\mathbf{z})$. Therefore, we must have that $\mathbf{z} = \mathbf{e}$. Solving for \mathbf{z} , we get that $\mathbf{e} = \mathbf{z} = a(\mathbf{y}^2 - \mathbf{y}^3)/(t_1 - t_2)$; this equality is well defined because $\mathbf{y}^2 \neq \mathbf{y}^3$ ensures that $t_1 \neq t_2$. As a result, for an appropriately defined constant δ , we obtain that $y_k^2 - y_k^3 = \delta$ for all children $k \in \text{Children}(j)$. Because we also have that $y_2^k - y_3^k = \xi_k$ from (EC.2), it must be that $\xi_k = \delta$ for all $k \in \text{Children}(j)$. Consequently, we have shown that $\bar{\mu}_\ell = \mu_\ell + \delta$ for all leaf nodes $\ell \in \cup_{k \in \text{Children}(j)} \mathbb{T}_k[\mathcal{S}] \cap \mathcal{S}$. Because $\mathbb{T}_j[\mathcal{S}] \cap \mathcal{S} = \cup_{k \in \text{Children}(j)} \mathbb{T}_k[\mathcal{S}] \cap \mathcal{S}$, we have shown that $\bar{\mu}_\ell = \mu_\ell + \xi$ for all leaf nodes $\ell \in \mathbb{T}_j[\mathcal{S}] \cap \mathcal{S}$ and $\xi = \delta$, which is the desired result. This completes the necessity part and finishes the induction. Therefore, part (b) holds for all nodes $j \in \mathbb{T}[\mathcal{S}]$.

Proof of part (c): Part (c) follows immediately from parts (a) and (b). ■

B.3 Proof of Lemma B.3

Proof: We will prove each of the three parts separately by induction on the height of node j .

Proof of part (a): The base case where j has a height of zero is trivially true by definition. Suppose the result is true for all nodes of height at most H . Then, consider node j with height $H + 1$. By definition,

$$W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = \lambda_j \log \left(\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\lambda_j} \right)$$

Note that $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ only depends on $(\mu_\ell : \ell \in \mathcal{S})$, $(\lambda_{\mathbb{T}_k[\mathcal{S}]} : k \in \text{Children}(j))$, and λ_j . By induction, for each $k \in \text{Children}(j)$, $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ is increasing in $\boldsymbol{\mu}$ and $\lambda_{\mathbb{T}_k[\mathcal{S}]}$. Moreover, taking the derivative of the above expression with respect to λ_j , we get

$$\frac{\partial W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_j} = \log \left(\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\lambda_j} \right) - \frac{1}{\lambda_j} \cdot \frac{\sum_{\ell \in \text{Children}(j) \cap \mathcal{S}} W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \times e^{W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\lambda_j}}{\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\lambda_j}} \geq 0,$$

where the inequality follows because $\log(\sum_{i=1}^n e^{x_i}) \geq \max_{i=1, \dots, n} x_i$. This shows that $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ is increasing in $(\boldsymbol{\mu}, \boldsymbol{\lambda})$, completing the induction and proving part (a).

Proof of part (b): By using induction on the height on node j , we will establish that $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ is strictly increasing in μ_ℓ for all $\ell \in \mathbb{T}_j[\mathcal{S}] \cap \mathcal{S}$. For the base case, suppose that j has a height of zero, so $j = \ell$ for some leaf node $\ell \in \mathbb{T}[\mathcal{S}] \cap \mathcal{S}$. We have by definition that $W_\ell(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = \mu_\ell$, which is clearly strictly increasing in μ_ℓ . Because $\mathbb{T}_\ell[\mathcal{S}] \cap \mathcal{S} = \{\ell\}$, we have established the base case. Suppose the result is true for all nodes k with height at most H . Consider a non-leaf node j at height $H + 1$. By definition, $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = \lambda_j \log\left(\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\lambda_j}\right)$. Now consider a leaf node $\ell \in \mathbb{T}_j[\mathcal{S}] \cap \mathcal{S}$. There exists some child node k of j in $\mathbb{T}_j[\mathcal{S}]$ such that $\ell \in \mathbb{T}_k[\mathcal{S}] \cap \mathcal{S}$. Since the height of $k \in \text{Children}(j)$ is at most H , it follows from the induction hypothesis that $W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ is strictly increasing in μ_ℓ . Moreover, because the log-sum-exp function is strictly increasing, it follows from the expression for $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ that it is strictly increasing in μ_ℓ . This completes the induction, establishing the monotonicity in $\boldsymbol{\mu}$ for all nodes $j \in \mathbb{T}[\mathcal{S}]$.

Proof of part (c): Now, consider the last result that $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ is strictly increasing in λ_k if node k has at least two children in $\mathbb{T}_j[\mathcal{S}]$. For the base case, consider node j with height one, so its children are leaf nodes. By definition, $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = \lambda_j \log\left(\sum_{\ell \in \text{Children}(j) \cap \mathcal{S}} e^{\mu_\ell/\lambda_j}\right)$. Taking the derivative of the above function with respect to λ_j , we get

$$\frac{\partial W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_j} = \log\left(\sum_{\ell \in \text{Children}(j) \cap \mathcal{S}} e^{\mu_\ell/\lambda_j}\right) - \frac{1}{\lambda_j} \cdot \frac{\sum_{\ell \in \text{Children}(j) \cap \mathcal{S}} \mu_\ell \times e^{\mu_\ell/\lambda_j}}{\sum_{\ell \in \text{Children}(j) \cap \mathcal{S}} e^{\mu_\ell/\lambda_j}}. \quad (\text{EC.4})$$

Because $|\text{Children}(j) \cap \mathcal{S}| \geq 2$, we have that $\log(\sum_{\ell \in \text{Children}(j) \cap \mathcal{S}} e^{\mu_\ell/\lambda_j}) > \max_{\ell \in \text{Children}(j) \cap \mathcal{S}} \mu_\ell/\lambda_j$. Further, since the second term in the expression above is a weighted average of the set of numbers $\{\mu_\ell/\lambda_j : \ell \in \text{Children}(j) \cap \mathcal{S}\}$, it follows that $\partial W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\partial \lambda_j > 0$, which establishes the base case.

Suppose that the result is true for all nodes k with height at most H . Consider a non-leaf node j at height $H + 1$. Consider an arbitrary non-leaf node $i \in \mathbb{T}_j[\mathcal{S}]$ such that i has at least two children. If $i \neq j$, then $i \in \mathbb{T}_k[\mathcal{S}]$ for some child k of j in $\mathbb{T}[\mathcal{S}]$. Because i has at least two children in $\mathbb{T}[\mathcal{S}]$, it follows from the induction hypothesis that $W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ is strictly increasing in λ_i . Moreover, because the log-sum-exp function is strictly increasing, it follows from the expression for $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ that it is also strictly increasing in λ_i .

Now suppose $i = j$. Because $W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ are independent of λ_j for all children k of j in sub-tree $\mathbb{T}[\mathcal{S}]$, we can compute the partial derivative of W_j with respect to λ_j , as done for the base case. Using identical arguments, we can conclude that $\partial W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\partial \lambda_j > 0$, establishing the induction step. This completes the proof. \blacksquare

B.4 Proof of Lemma B.4

Proof: Fix an arbitrary $\mathcal{S} \subseteq \mathcal{N}$. We first consider the derivative with respect to μ_ℓ and prove it using induction on the height of the node. For the base case, we consider a leaf node ℓ of sub-tree $\mathbb{T}[\mathcal{S}]$. We have by definition that $W_\ell(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = \mu_\ell$. Therefore, $\partial W_\ell(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\partial \mu_\ell = 1 = \psi_{\ell \rightarrow \ell}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$. We have thus established the base case.

Now suppose the result is true for all nodes of height at most H . In other words, suppose that $\partial W_v(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\partial \mu_\ell = \psi_{v \rightarrow \ell}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ for all nodes v of height at most H and all leaf nodes ℓ in $\mathbb{T}[\mathcal{S}]$. Now consider a non-leaf node j at height $H + 1$. We have by definition that

$$\frac{\partial W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \mu_\ell} = \frac{\partial}{\partial \mu_\ell} \lambda_j \log \left(\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\lambda_j} \right) = \frac{\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\lambda_j} \times \frac{\partial W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \mu_\ell}}{\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\lambda_j}}$$

Since ℓ is a leaf node in $\mathbb{T}_j[\mathcal{S}]$, there exists exactly one child node $i \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]$ such that ℓ is a leaf node in $\mathbb{T}_i[\mathcal{S}]$. For any other child node $k \neq i$, we must have that $\partial W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\partial \mu_\ell = 0$. Therefore, we obtain

$$\frac{\partial W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \mu_\ell} = \frac{e^{W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\lambda_j}}{\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\lambda_j}} \times \frac{\partial W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \mu_\ell} = \psi_{j \rightarrow i}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \times \psi_{i \rightarrow \ell}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}),$$

where the first term in the second equality above follows from the definition of $\psi_{j \rightarrow i}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ and the second term follows from the induction hypothesis. Because $\psi_{j \rightarrow i}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \times \psi_{i \rightarrow \ell}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = \psi_{j \rightarrow \ell}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ by definition, we have shown that $\partial W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\partial \mu_\ell = \psi_{j \rightarrow \ell}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$, as desired. We have thus established the induction step. The first result now follows by induction.

We now consider the derivative of $W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ with respect to λ_k for some non-leaf nodes j and k in $\mathbb{T}[\mathcal{S}]$. We prove the result by induction on the height of j . For the base case, we start with a non-leaf node j at height one, so all children of j are leaf nodes. As shown in Equation (EC.4) in the proof of Lemma B.3, we have that

$$\frac{\partial W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_j} = \log \left(\sum_{\ell \in \text{Children}(j) \cap \mathcal{S}} e^{\mu_\ell/\lambda_j} \right) - \frac{1}{\lambda_j} \cdot \frac{\sum_{\ell \in \text{Children}(j) \cap \mathcal{S}} \mu_\ell \times e^{\mu_\ell/\lambda_j}}{\sum_{\ell \in \text{Children}(j) \cap \mathcal{S}} e^{\mu_\ell/\lambda_j}} = \psi_{j \rightarrow j}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \times \Delta_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}),$$

where the second equality follows from the definition of $\Delta_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ and the fact that $\psi_{j \rightarrow j}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = 1$. We have thus established the base case.

For the induction step, we suppose that $\partial W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\partial \lambda_k = \psi_{j \rightarrow k}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \times \Delta_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ for all nodes j of height at most H and for all non-leaf nodes $k \in \mathbb{T}_j[\mathcal{S}]$. Now consider a node j of height $H + 1$. First, suppose that $k \neq j$. We then have by definition that

$$\frac{\partial W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_k} = \frac{\partial}{\partial \lambda_k} \lambda_j \log \left(\sum_{i \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\lambda_j} \right) = \frac{\sum_{i \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\lambda_j} \times \frac{\partial W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_k}}{\sum_{i \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})/\lambda_j}}.$$

Because of the property of a tree, the non-leaf node k belongs to the sub-tree $\mathbb{T}_i[\mathcal{S}]$ of exactly one child node i of node j . For any other child node $v \neq i$, we have that $\partial W_v(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) / \partial \lambda_k = 0$.

Therefore, we can now write

$$\frac{\partial W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_k} = \frac{e^{W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) / \lambda_j}}{\sum_{i \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) / \lambda_j}} \times \frac{\partial W_i(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_k} = \psi_{j \rightarrow i}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \times \psi_{i \rightarrow k}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \times \Delta_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}),$$

where the first term in the second equality above follows from the definition on $\psi_{j \rightarrow i}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ and the second and third terms follow from the induction hypothesis. Because $\psi_{j \rightarrow k}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = \psi_{j \rightarrow i}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \times \psi_{i \rightarrow k}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ by definition, we have shown that $\partial W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) / \partial \lambda_k = \psi_{j \rightarrow k}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \times \Delta_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$, as desired. This completes the induction step, completing the lemma. The final expression for $\Delta_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ follows from plugging in the definitions of $W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ and $\psi_{k \rightarrow i}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ for all $i \in \text{Children}(k) \cap \mathbb{T}[\mathcal{S}]$. ■

It follows from the above lemma that $\frac{\partial W_{\text{root}}}{\partial \mu_\ell}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = \psi_{\text{root} \rightarrow \ell}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = \mathbb{P}_\ell(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$, and this is consistent with the result from McFadden (1978), which provides an expression of the choice probability in terms of the derivative of the generating function for Generalized Extreme Value choice models.

Appendix C: Proofs of Theorem 2.2, Theorem 2.4 and Theorem 2.5

C.1 Proof Of Theorem 2.2

The negative log-likelihood value of the data is equal to $\frac{1}{Q} \sum_{q=1}^Q -\log \psi_{\text{root} \rightarrow c^q}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$. The expression in the statement of the lemma is obtained by invoking the expression in Equation (2) for $-\log \psi_{\text{root} \rightarrow c^q}^q$ and then re-arranging the terms. More precisely, we have

$$\begin{aligned} \sum_{q=1}^Q -\log \psi_{\text{root} \rightarrow c^q}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) &= \sum_{q=1}^Q \sum_{j \in \text{path}(\text{root}, c^q)} \frac{W_{\text{pa}(j)}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) - W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(j)}} \\ &= \sum_{q=1}^Q \left\{ W_{\text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) + \sum_{j \in \text{path}(\text{root}, c^q)} W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) - \frac{\mu_{c^q}}{\lambda_{\text{pa}(c^q)}} \right\}, \end{aligned}$$

where the last equality follows from a straightforward re-arrangement of the terms of the inner summation. A slightly different re-arrangement yields the following expression:

$$\begin{aligned} \sum_{q=1}^Q -\log \psi_{\text{root} \rightarrow c^q}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) &= \sum_{q=1}^Q \sum_{k \in \text{path}(\text{root}, c^q)} \frac{W_{\text{pa}(k)}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) - W_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(k)}} \\ &\stackrel{(a)}{=} \sum_{q=1}^Q \sum_{k \in \mathbb{T} \setminus \{\text{root}\}} \mathbb{1}_{\{c^q \in \mathbb{T}_k\}} \frac{W_{\text{pa}(k)}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(k)}} - \sum_{q=1}^Q \sum_{k \in \mathbb{T} \setminus \{\text{root}\}} \mathbb{1}_{\{c^q \in \mathbb{T}_k\}} \frac{W_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(k)}}, \\ &\stackrel{(b)}{=} \sum_{q=1}^Q \sum_{j \in \mathbb{T}} \mathbb{1}_{\{c^q \in \mathbb{T}_j\}} \frac{W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_j} - \sum_{q=1}^Q \sum_{j \in \mathbb{T}} \mathbb{1}_{\{c^q \in \mathbb{T}_j\}} \frac{W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(j)}}, \end{aligned}$$

where the equality (a) follows because $k \in \text{path}(\text{root}, c^q]$ if and only if c^q belongs to the sub-tree rooted at k ; that is, $c^q \in \mathbb{T}_k$. The first term in the last equality (b) follows from the change of variable $j = \text{pa}(k)$. With this change of variable, as k is varied over the set of nodes $\mathbb{T} \setminus \{\text{root}\}$, the variable j varies over the set $\mathbb{T} \setminus \mathcal{N}$. Because $\lambda_\ell = +\infty$ for each leaf node $\ell \in \mathcal{N}$, we can extend the summation to cover the entire set of nodes in \mathbb{T} , resulting in the first term. The second term is obtained by replacing k with j , and including root in the summation, since $\lambda_{\text{pa}(\text{root})} = +\infty$. ■

C.2 Proof of Theorem 2.4

To facilitate the proof of Theorem 2.4, let us introduce the following function $b: \mathcal{D}_1 \rightarrow \mathbb{R}_+$ defined by: for each $\zeta \in \mathcal{D}_1$,

$$b(\zeta) = \max_{q=1, \dots, Q} \max_{\ell \in \mathcal{S}^q \setminus \{c^q\}} (\zeta_\ell - \zeta_{c^q}),$$

where we use the convention that the maximum of an empty set is 0. Further, we let $b^* = \min_{\zeta \in \mathcal{D}_1 \cap \Delta} b(\zeta)$ where $\Delta = \{\mu \mid \|\mu\|_\infty = 1\}$. In other words, $b(\zeta)$ denotes the maximum possible loss in mean utility from the choices made in the data under the mean utility vector ζ . Maximizing the log-likelihood requires us to make this loss negative, *if possible*. If the loss is made negative for some ζ , then scaling all the mean utility values by a constant makes the loss diverge to $-\infty$, resulting in unbounded optimal solutions. As we shown below, the assumption of strong connectedness of the comparison graph prevents this from happening. In particular, strong connectedness ensures that the loss is always positive for any ζ . We formalize this intuition next. The first lemma shows that b^* is positive when the comparison graph Comp is strongly connected. Recall that, given transaction data $\{(\mathcal{S}^q, c^q) : q = 1, \dots, Q\}$, a comparison graph $\text{Comp} = (\mathcal{N}, \mathbb{E})$ is a directed graph whose nodes correspond to the products, and there is a directed edge $(\ell_1, \ell_2) \in \mathbb{E}$ if there exists an offer set \mathcal{S}^q such that $\{\ell_1, \ell_2\} \subseteq \mathcal{S}^q$ and $c^q = \ell_1$.

Lemma C.1 (Positive Gap) *If the comparison graph Comp is strongly connected, then $b^* > 0$.*

Proof: We start with the observation that because $\Delta \cap \mathcal{D}_1$ is compact, the minimum $\min_{\mu \in \Delta \cap \mathcal{D}_1} b(\mu)$ is attained at some $\mu^* \in \Delta \cap \mathcal{D}_1$, so that $b^* = b(\mu^*)$. Therefore, to show that $b^* > 0$, it is sufficient to show that $b(\mu) > 0$ for every $\mu \in \Delta \cap \mathcal{D}_1$. We prove this result by contradiction. Suppose on the contrary that $b(\zeta) \leq 0$ for some $\zeta \in \Delta \cap \mathcal{D}_1$. Because $b(\zeta) \leq 0$, it follows from definition that $\zeta_\ell \leq \zeta_{c^q}$ for all $\ell \in \mathcal{S}^q$ and for all $q = 1, \dots, Q$. Consider an arbitrary directed edge (j, k) in Comp . By our construction, there exists a $q \in \{1, \dots, Q\}$ such that $j = c^q$ and $k \in \mathcal{S}^q \setminus \{c^q\}$. Therefore, we must have that $\zeta_k \leq \zeta_j$. But, Comp is strongly connected, which implies that there is a directed path $k \rightarrow s_1 \rightarrow s_2 \rightarrow \dots \rightarrow s_m \rightarrow j$ in Comp from node k to j . Applying the result that $\zeta_t \leq \zeta_s$ whenever there is a directed edge from s to t , we obtain that $\zeta_j \leq \zeta_{s_m} \leq \dots \leq \zeta_{s_1} \leq \zeta_k$, which

implies that $\zeta_j \leq \zeta_k$. This inequality together with the result that $\zeta_k \leq \zeta_j$, implies that $\zeta_k = \zeta_j$. Because j and k are arbitrary nodes and any two nodes are connected by directed paths, we have shown that $\zeta_1 = \zeta_2 = \dots = \zeta_n$. But $\zeta \in \Delta \cap \mathcal{D}_1$, so $\zeta_1 = 0$, and consequently that $\mathbf{0} \in \Delta$. This is a contradiction! Thus, $b(\boldsymbol{\mu}) > 0$ for all $\boldsymbol{\mu} \in \Delta \cap \mathcal{D}_1$. This completes the proof. \blacksquare

The next lemma gives an upper bound on the selection probability in terms of the mean utilities.

Lemma C.2 For each $\ell_1 \in \mathcal{S}$, $\ell_2 \in \mathcal{S}$, and $(\boldsymbol{\mu}, \boldsymbol{\lambda})$, $-\log \mathbb{P}_{\ell_1}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \geq \mu_{\ell_2} - \mu_{\ell_1}$.

Proof: By Equation (2),

$$\begin{aligned} -\log \mathbb{P}_{\ell_1}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) &= -\log \psi_{\text{root} \rightarrow \ell_1}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = \sum_{j \in \text{path}(\text{root}, \ell_1]} \frac{W_{\text{pa}(j)}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) - W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(j)}} \\ &\geq \sum_{j \in \text{path}(\text{root}, \ell_1]} W_{\text{pa}(j)}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) - W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = W_{\text{root}}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) - W_{\ell_1}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}), \end{aligned}$$

where the inequality follows because $W_{\text{pa}(j)}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \geq W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ and $\lambda_{\text{pa}(j)} \leq 1$ for each node j such that $j \neq \text{root}$. The last equality follows from the telescoping sum.

Now, we claim that $W_{\text{root}}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \geq W_{\ell}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ for any leaf node $\ell \in \mathcal{S}$. To see this, consider the path $\text{root} \rightarrow j_1 \rightarrow \dots \rightarrow j_m \rightarrow \ell$ from the root node to the leaf node ℓ in tree $\mathbb{T}[\mathcal{S}]$. Because $W_{\text{pa}(j)}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \geq W_j(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$ for all nodes j , we have the sequence of inequalities $W_{\text{root}}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \geq W_{j_1}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \geq \dots \geq W_{j_m}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \geq W_{\ell}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda})$, which establishes the claim.

By choosing $\ell = \ell_2$, we get

$$-\log \mathbb{P}_{\ell_1}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \geq W_{\text{root}}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) - W_{\ell_1}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) \geq W_{\ell_2}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) - W_{\ell_1}(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) = \mu_{\ell_2} - \mu_{\ell_1},$$

where the last equality follows because both ℓ_1 and ℓ_2 are in \mathcal{S} . \blacksquare

The next lemma establishes an upper bound on the $\text{NegLog}(\mathbf{0}, \boldsymbol{\lambda})$ for all $\boldsymbol{\lambda}$.

Lemma C.3 For each $\boldsymbol{\lambda} \in \mathcal{D}_2$, $\text{NegLog}(\mathbf{0}, \boldsymbol{\lambda}) \leq |\mathbb{T}| \log |\mathbb{T}|$.

Proof: Consider an arbitrary $\boldsymbol{\lambda} \in \mathcal{D}_2$ and $\mathcal{S} \subseteq \mathcal{N}$. We will first establish the following claim.

Claim: For each node $j \in \mathbb{T}[\mathcal{S}]$, $W_j(\mathcal{S}; \mathbf{0}, \boldsymbol{\lambda})/\lambda_j \leq \log |\mathbb{T}_j[\mathcal{S}]|$.

We will prove the claim by induction on the height of node j . For the base case where node j has a height of zero, then $j = \ell$ for some leaf node $\ell \in \mathcal{S}$. Since $\lambda_{\ell} = +\infty$ by definition, we have $W_{\ell}(\mathcal{S}; \mathbf{0}, \boldsymbol{\lambda})/\lambda_{\ell} = 0 = \log |\mathbb{T}_{\ell}[\mathcal{S}]|$ because $\mathbb{T}_{\ell}[\mathcal{S}] = \{\ell\}$. This establishes the base case. Suppose the claim holds for all nodes with height at most H . Consider an arbitrary node j with height $H+1$. By the inductive hypothesis, the claim holds for all children of j ; that is, for each $k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]$, $W_k(\mathcal{S}; \mathbf{0}, \boldsymbol{\lambda})/\lambda_k \leq \log |\mathbb{T}_k[\mathcal{S}]|$, and thus, $W_k(\mathcal{S}; \mathbf{0}, \boldsymbol{\lambda})/\lambda_j \leq \log |\mathbb{T}_k[\mathcal{S}]|$ because $\lambda_k \leq \lambda_j$. Therefore,

$$\frac{W_j(\mathcal{S}; \mathbf{0}, \boldsymbol{\lambda})}{\lambda_j} = \log \left(\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} e^{W_k(\mathcal{S}; \mathbf{0}, \boldsymbol{\lambda})/\lambda_j} \right) \leq \log \left(\sum_{k \in \text{Children}(j) \cap \mathbb{T}[\mathcal{S}]} |\mathbb{T}_k[\mathcal{S}]| \right) \leq \log |\mathbb{T}_j[\mathcal{S}]|,$$

which completes the induction argument. So, the claim holds for all nodes $j \in \mathsf{T}[\mathcal{S}]$.

To establish the lemma, note that by Equation (2), for each transaction (\mathcal{S}^q, c^q) ,

$$\begin{aligned} -\log \mathbb{P}_{c^q}(\mathcal{S}^q; \mathbf{0}, \boldsymbol{\lambda}) &= -\log \psi_{\text{root} \rightarrow c^q}(\mathcal{S}^q; \mathbf{0}, \boldsymbol{\lambda}) = \sum_{j \in \text{path}(\text{root}, c^q)} \frac{W_{\text{pa}(j)}(\mathcal{S}^q; \mathbf{0}, \boldsymbol{\lambda}) - W_j(\mathcal{S}^q; \mathbf{0}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(j)}} \\ &\leq \sum_{j \in \text{path}(\text{root}, c^q)} \frac{W_{\text{pa}(j)}(\mathcal{S}^q; \mathbf{0}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(j)}} \leq \sum_{j \in \text{path}(\text{root}, c^q)} \log \left| \mathsf{T}_{\text{pa}(j)}^q \right| \leq |\mathsf{T}^q| \log |\mathsf{T}^q|, \end{aligned}$$

where the first inequality follows because $W_j(\cdot)$ is non-negative, and the second inequality follows from the above claim. Therefore, $\text{NegLog}(\mathbf{0}, \boldsymbol{\lambda}) \leq \frac{1}{Q} \sum_{q=1}^Q |\mathsf{T}^q| \log |\mathsf{T}^q| \leq |\mathsf{T}| \log |\mathsf{T}|$, which is the desired result. \blacksquare

Here is the proof of Theorem 2.4.

Proof of Theorem 2.4: Fix an arbitrary $\boldsymbol{\lambda} \in \mathcal{D}_2$. Recall that $\mathcal{D}_1 = \{\boldsymbol{\mu} \in \mathbb{R}^n : \mu_1 = 0\}$.

PROOF OF NECESSITY: We will first prove the necessity, so suppose that the optimization problem $\min_{\boldsymbol{\mu} \in \mathcal{D}_1} \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda})$ admits a unique and bounded optimal solution, and let $\boldsymbol{\mu}^* \in \mathcal{D}_1$ denote the unique optimal solution. We will prove by contradiction that the comparison graph **Comp** must be strongly connected.

Suppose on the contrary that **Comp** is not strongly connected. This means that there exist vertices ℓ_1 and ℓ_2 such that ℓ_2 is not reachable from ℓ_1 . Let A_1 denote the set of vertices that are reachable from ℓ_1 , and let $A_2 = \mathcal{N} \setminus A_1$ denote the remaining vertices in **Comp**. Note that $\ell_1 \in A_1$ and $\ell_2 \in A_2$, so both A_1 and A_2 are nonempty, disjoint, and $A_1 \cup A_2 = \mathcal{N}$. By definition, there is no directed path from a node in A_1 to a node in A_2 . Without loss of generality, assume that $1 \in A_1$; the argument for the case where $1 \in A_2$ is analogous.

Consider an arbitrary $\xi > 0$. Define a $\bar{\boldsymbol{\mu}}$ as follows:

$$\bar{\mu}_\ell = \begin{cases} \mu_\ell^* & \text{if } \ell \in A_1, \\ \mu_\ell^* + \xi & \text{if } \ell \in A_2, \end{cases}$$

It is easy to verify that $\bar{\boldsymbol{\mu}} \in \mathcal{D}_1$. Now, consider an arbitrary transaction (\mathcal{S}^q, c^q) . There are 3 cases to consider: 1) $\mathcal{S}^q \subseteq A_1$, 2) $\mathcal{S}^q \subseteq A_2$, and 3) $\mathcal{S}^q \cap A_1 \neq \emptyset$ and $\mathcal{S}^q \cap A_2 \neq \emptyset$.

Case 1: $\mathcal{S}^q \subseteq A_1$. Then, we have that $W_j(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) = W_j(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda})$ for all $j \in \mathsf{T}[\mathcal{S}^q]$ because the weight function only depends on the utility of the products in the offer set \mathcal{S}^q , and by our construction, we have that $\bar{\mu}_\ell = \mu_\ell^*$ for all $\ell \in \mathcal{S}^q \subseteq A_1$. Therefore, $-\log \psi_{\text{root} \rightarrow c^q}(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) = -\log \psi_{\text{root} \rightarrow c^q}(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda})$.

Case 2: $\mathcal{S}^q \subseteq A_2$. By Equation (2), we have that

$$-\log \psi_{\text{root} \rightarrow c^q}(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) = \sum_{j \in \text{path}(\text{root}, c^q)} \frac{W_{\text{pa}(j)}^q(\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) - W_j^q(\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(j)}}$$

$$\begin{aligned}
&\stackrel{(a)}{=} \frac{W_{\text{root}}(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda})}{\lambda_{\text{root}}} + \sum_{j \in \text{path}(\text{root}, c^q]} W_j(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right), \\
&\stackrel{(b)}{=} \frac{W_{\text{root}}(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda})}{\lambda_{\text{root}}} + \sum_{j \in \text{path}(\text{root}, c^q]} W_j(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda}) \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \\
&\quad + \frac{\xi}{\lambda_{\text{root}}} + \sum_{j \in \text{path}(\text{root}, c^q]} \left(\frac{\xi}{\lambda_j} - \frac{\xi}{\lambda_{\text{pa}(j)}} \right) \\
&\stackrel{(c)}{=} \frac{W_{\text{root}}(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda})}{\lambda_{\text{root}}} + \sum_{j \in \text{path}(\text{root}, c^q]} W_j(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda}) \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(v)}} \right) \\
&= -\log \psi_{\text{root} \rightarrow c^q}(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda}) .
\end{aligned}$$

where the equality (a) follows from collecting terms and (b) follows because $\bar{\mu}_\ell = \mu_\ell^* + \xi$ for all $\ell \in A_2$ and from Lemma B.1 on the invariance under translation by a constant, which implies that $W_j(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) = \xi + W_j(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda})$. The last equality (c) follows by the telescoping sum and the fact that $\lambda_\ell = +\infty$ for all leaf nodes $\ell \in \mathcal{N}$.

Case 3: $\mathcal{S}^q \cap A_1 \neq \emptyset$ and $\mathcal{S}^q \cap A_2 \neq \emptyset$. In this case, $|\mathcal{S}^q| \geq 2$, and thus, by our construction of the comparison graph Comp , there is a directed edge from c^q to all other elements in \mathcal{S}^q . Since there is no directed path from A_1 to A_2 , it must be the case that $c^q \in A_2$. By re-arranging the terms in Equation (2) and using the fact that $\lambda_\ell = +\infty$ for all leaf nodes $\ell \in \mathcal{N}$, we can write

$$\begin{aligned}
-\log \psi_{\text{root} \rightarrow c^q}(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) &= \frac{W_{\text{root}}(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda})}{\lambda_{\text{root}}} + \sum_{j \in \text{path}(\text{root}, c^q]} W_j(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) - \frac{W_{c^q}(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(c^q)}} \\
&\stackrel{(a)}{=} \frac{W_{\text{root}}(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda})}{\lambda_{\text{root}}} + \sum_{j \in \text{path}(\text{root}, c^q]} W_j(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \\
&\quad - \frac{W_{c^q}(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda})}{\lambda_{\text{pa}(c^q)}} - \frac{\xi}{\lambda_{\text{pa}(c^q)}} \\
&\stackrel{(b)}{\leq} \frac{W_{\text{root}}(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda})}{\lambda_{\text{root}}} + \sum_{j \in \text{path}(\text{root}, c^q]} W_j(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda}) \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \\
&\quad + \frac{\xi}{\lambda_{\text{root}}} + \sum_{j \in \text{path}(\text{root}, c^q]} \left(\frac{\xi}{\lambda_j} - \frac{\xi}{\lambda_{\text{pa}(j)}} \right) - \frac{W_{c^q}(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda})}{\lambda_{\text{pa}(c^q)}} - \frac{\xi}{\lambda_{\text{pa}(c^q)}} \\
&\stackrel{(c)}{=} -\log \psi_{\text{root} \rightarrow c^q}(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda}) .
\end{aligned}$$

where (a) follows because $c^q \in A_2$ for all q , so $W_{c^q}(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) = \bar{\mu}_{c^q} = \mu_{c^q}^* + \xi$. The inequality (b) above follows from the fact for all $j \in \text{path}(\text{root}, c^q)$,

$$W_j(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \leq W_j(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda}) \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) + \xi \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right)$$

because $0 < \lambda_j \leq \lambda_{\text{pa}(j)}$ and by Lemma B.3, the weight function is increasing in $\boldsymbol{\mu}$, so for all $j \in \mathbb{T}[\mathcal{S}^q]$, $W_j(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) \leq W_j(\mathcal{S}^q; \boldsymbol{\mu}^* + \xi \mathbf{e}, \boldsymbol{\lambda}) = \xi + W_j(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda})$, where the equality follows from the

translation variance property. The final equality (c) above follows from collecting terms in the telescoping sum.

Therefore, we observe that in all three cases,

$$-\log \psi_{\text{root} \rightarrow c^q}(\mathcal{S}^q; \bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) \leq -\log \psi_{\text{root} \rightarrow c^q}(\mathcal{S}^q; \boldsymbol{\mu}^*, \boldsymbol{\lambda}) .$$

Since the transaction (\mathcal{S}^q, c^q) is arbitrary, it follows that $\text{NegLog}(\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) \leq \text{NegLog}(\boldsymbol{\mu}^*, \boldsymbol{\lambda})$, but this contradicts the fact that $\boldsymbol{\mu}^*$ is the unique optimal solution! Hence, it must be the case that Comp is strongly connected. This completes the proof of the necessity.

PROOF OF SUFFICIENCY: We will now prove the sufficiency, so assume that the comparison graph Comp is strongly connected. Fix an arbitrary $\boldsymbol{\lambda}$. Consider $\boldsymbol{\mu} \in \mathcal{D}_1$ such that $\|\boldsymbol{\mu}\|_\infty > \frac{Q}{b^*} \text{NegLog}(\mathbf{0}, \boldsymbol{\lambda})$. Note that $b^* > 0$ by Lemma C.1. Then, $\boldsymbol{\mu} = \|\boldsymbol{\mu}\|_\infty \boldsymbol{\zeta}$ where $\|\boldsymbol{\zeta}\|_\infty = 1$. By definition, there exists a transaction (\mathcal{S}^q, c^q) such that $b(\boldsymbol{\zeta}) = \zeta_\ell - \zeta_{c^q}$ for some $\ell \in \mathcal{S}^q \setminus \{c^q\}$. Then, by Lemma C.2

$$\text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda}) \geq -\log \mathbb{P}_{c^q}(\mathcal{S}^q; \boldsymbol{\mu}, \boldsymbol{\lambda}) / Q \geq (\mu_\ell - \mu_{c^q}) / Q = (\zeta_\ell - \zeta_{c^q}) \|\boldsymbol{\mu}\|_\infty / Q = b(\boldsymbol{\zeta}) \|\boldsymbol{\mu}\|_\infty / Q,$$

and it follows that $\text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda}) \geq \|\boldsymbol{\mu}\|_\infty b(\boldsymbol{\zeta}) / Q \geq \|\boldsymbol{\mu}\|_\infty b^* / Q > \text{NegLog}(\mathbf{0}, \boldsymbol{\lambda})$, which shows that $\min_{\boldsymbol{\mu} \in \mathcal{D}_1} \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda}) = \min \{ \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda}) \mid \boldsymbol{\mu} \in \mathcal{D}_1, \|\boldsymbol{\mu}\|_\infty \leq \frac{Q}{b^*} \text{NegLog}(\mathbf{0}, \boldsymbol{\lambda}) \}$, so the optimization problem has a bounded solution.

We now show that the optimal solution must be unique by establishing that the mapping $\boldsymbol{\mu} \mapsto \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda})$ is strictly convex over the set \mathcal{D}_1 . For that, we first note from Theorem 2.2 that

$$\text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda}) = \frac{1}{Q} \sum_{q=1}^Q \left\{ W_{\text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) + \sum_{j \in \text{path}(\text{root}, c^q)} W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) - \frac{\mu_{c^q}}{\lambda_{\text{pa}(c^q)}} \right\}$$

Because $W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$ is convex in $\boldsymbol{\mu}$ for each q and $\boldsymbol{\lambda}$ (by Lemma B.2) and $1/\lambda_j - 1/\lambda_{\text{pa}(j)} \geq 0$ for all $j \in \mathcal{T} \setminus \mathcal{N}$, it follows that the function above is convex in $\boldsymbol{\mu}$.

We are left to show that the convexity is strict, for which it is sufficient to show that the following function is strictly convex: $f(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \sum_{q=1}^Q W_{\text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$. We show strict convexity through a direct application of its definition. The function $f(\cdot)$ is strictly convex if for $\boldsymbol{\mu} \neq \bar{\boldsymbol{\mu}}$ such that $\mu_1 = \bar{\mu}_1 = 0$ and $\theta \in (0, 1)$, we have $f(\theta\boldsymbol{\mu} + (1-\theta)\bar{\boldsymbol{\mu}}) < \theta f(\boldsymbol{\mu}) + (1-\theta)f(\bar{\boldsymbol{\mu}})$. To arrive at a contradiction, suppose that for some $\boldsymbol{\mu} \neq \bar{\boldsymbol{\mu}}$ and $\theta \in (0, 1)$, we have that

$$\begin{aligned} f(\theta\boldsymbol{\mu} + (1-\theta)\bar{\boldsymbol{\mu}}) &= \theta f(\boldsymbol{\mu}) + (1-\theta)f(\bar{\boldsymbol{\mu}}) \\ \iff \sum_{q=1}^Q W_{\text{root}}^q(\theta\boldsymbol{\mu} + (1-\theta)\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) &= \sum_{q=1}^Q \theta W_{\text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) + (1-\theta)W_{\text{root}}^q(\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) . \end{aligned} \quad (\text{EC.5})$$

By Lemma B.2, for each $q \in \{1, \dots, Q\}$, $W_{\text{root}}^q(\theta\boldsymbol{\mu} + (1-\theta)\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}) \leq \theta W_{\text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) + (1-\theta)W_{\text{root}}^q(\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda})$, with equality occurring if and only if $\mu_\ell = \bar{\mu}_\ell + \kappa^q$ for all $\ell \in \mathcal{S}^q$, where $\kappa^q \in \mathbb{R}$ is some constant.

Therefore, for Equation (EC.5) to be satisfied, we must have equality occurring for each $q \in \{1, \dots, Q\}$. As a result, there exists $\kappa^1, \kappa^2, \dots, \kappa^Q$ such that for each $q \in \{1, \dots, Q\}$,

$$\mu_\ell = \bar{\mu}_\ell + \kappa^q \text{ for all } \ell \in \mathcal{S}^q.$$

To arrive at a contradiction, we now establish that $\kappa^q = 0$ for all q , which implies that $\boldsymbol{\mu} = \bar{\boldsymbol{\mu}}$ since $\mathcal{N} = \cup_{q \in [Q]} \mathcal{S}^q$. This contradicts our assumption that $\boldsymbol{\mu} \neq \bar{\boldsymbol{\mu}}$!

To show that $\kappa^q = 0$ for all $q = 1, \dots, Q$, we first prove the following claim.

Claim: $\kappa^1 = \kappa^2 = \dots = \kappa^Q$.

To prove the claim, it suffices to show that $\kappa^1 = \kappa^2$. Exactly the same argument applies to show that $\kappa^{q_1} = \kappa^{q_2}$ for all $q_1 \neq q_2$. If $\mathcal{S}^1 \cap \mathcal{S}^2 \neq \emptyset$, the result is trivially true because there exists $\ell \in \mathcal{S}^1 \cap \mathcal{S}^2$, which implies that $\kappa^1 = \mu_\ell - \bar{\mu}_\ell = \kappa^2$, which is the desired result. So, suppose that $\mathcal{S}^1 \cap \mathcal{S}^2 = \emptyset$. Pick $\ell_1 \in \mathcal{S}^1$ and $\ell_2 \in \mathcal{S}^2$. Since **Comp** is strongly connected, there exists a path $\ell_1 = j_0 \rightarrow j_1 \rightarrow \dots \rightarrow j_{m-1} \rightarrow j_m = \ell_2$ in **Comp** from ℓ_1 to ℓ_2 . By our construction of **Comp**, for each edge (j_{t-1}, j_t) , there exists a set $\mathcal{S}^{h_t} \in \{\mathcal{S}^1, \dots, \mathcal{S}^Q\}$ such that $\{j_{t-1}, j_t\} \subseteq \mathcal{S}^{h_t}$. Therefore, for each $t = 1, \dots, m$, $j_t \in \mathcal{S}^{h_t} \cap \mathcal{S}^{h_{t+1}}$, and thus, $\kappa^{h_t} = \mu_{j_t} - \bar{\mu}_{j_t} = \kappa^{h_{t+1}}$. Since this is true for all $t = 1, \dots, m$, it follows that

$$\kappa^{h_1} = \kappa^{h_2} = \dots = \kappa^{h_m}.$$

Since $\ell_1 \in \mathcal{S}^1 \cap \mathcal{S}^{h_1}$ and $\ell_2 \in \mathcal{S}^2 \cap \mathcal{S}^{h_m}$, we have that $\kappa^1 = \kappa^{h_1} = \kappa^{h_m} = \kappa^2$. This is the desired result, proving the claim.

We will now use the above claim to show that $\kappa^q = 0$ for all q . Consider the subset \mathcal{S}^q that contains product 1. Because $\mu_1 = \bar{\mu}_1 = 0$, it follows that $\kappa^q = 0$. Because $\kappa^1 = \dots = \kappa^Q$, it follows that $\kappa^q = 0$ for all $q = 1, \dots, Q$. This completes the proof of sufficiency. ■

C.3 Proof of Theorem 2.5

Proof: We will show that for each non-leaf node $j \in \mathbb{T} \setminus \mathcal{N}$ such that $j \neq \text{root}$, the function **NegLog** function varies with respect to λ_j if and only if there exists a transaction $q \in \{1, \dots, Q\}$ such that node j has at least two children in sub-tree \mathbb{T}^q . Consider an arbitrary non-leaf node $j \neq \text{root}$. The sufficiency follows immediately from Lemma B.3. To establish the necessity, assume that **NegLog** function varies with respect to λ_j . Suppose, on the contrary, that for every transaction q , node j has at most one child in the tree \mathbb{T}^q . Then, by definition, for every q ,

$$W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) = \lambda_j \log \left(\sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} e^{W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}) / \lambda_j} \right) = \sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} W_k(\mathcal{S}; \boldsymbol{\mu}, \boldsymbol{\lambda}),$$

where the last equality follows because there is at most one term in the summand. Therefore, for every q , $W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$ is independent of λ_j , which implies that $-\log \psi_{\text{root} \rightarrow c^q}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$ is independent of λ_j , which implies that NegLog function is also independent of λ_j . Contradiction! Therefore, there must be at least one transaction q such that node j has at least two children in the tree \mathbb{T}^q . This proves the necessity. \blacksquare

Appendix D: Derivatives of the negative log-likelihood function $\text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda})$

The following lemmas show that the partial derivatives $(\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda}) / \partial \mu_\ell : \ell \in \mathcal{N})$ and $(\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda}) / \partial \lambda_j : j \in \mathbb{T} \setminus (\mathcal{N} \cup \{\text{root}\}))$ can be computed efficiently via a recursion that starts at the `root` node and ends at the leaf nodes of the tree \mathbb{T} . To facilitate our exposition, for each transaction q , define the set of values $(D_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) : j \in \mathbb{T})$ recursively starting from `root` as follows: Initialize $D_{\text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) = 1$ and for all $j \in \mathbb{T} \setminus \{\text{root}\}$,

$$D_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \equiv \mathbb{1}_{\{c^q \in \mathbb{T}_j\}} \cdot \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) + \psi_{\text{pa}(j) \rightarrow j}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \cdot D_{\text{pa}(j)}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}).$$

We emphasize that the values $D_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$ can be computed efficiently through a simple recursion starting from the `root` node. We will use the values $D_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$ to compute the derivative of the log-likelihood function; see Lemmas D.2 and D.3. The following lemma gives an equivalent expression for $D_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$.

Lemma D.1 (Equivalent Expression for $D_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$) *For each transaction q and for all $k \in \mathbb{T}$,*

$$D_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) = \sum_{j \in \mathbb{T}} \mathbb{1}_{\{c^q \in \mathbb{T}_j\}} \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \psi_{j \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}).$$

Proof: We prove the result by induction on the height of node k . When $k = \text{root}$, the RHS in the equality above reduces to

$$\begin{aligned} \sum_{j \in \mathbb{T}} \mathbb{1}_{\{c^q \in \mathbb{T}_j\}} \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \psi_{j \rightarrow \text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) &= \left(\frac{1}{\lambda_{\text{root}}} - \frac{1}{\lambda_{\text{pa}(\text{root})}} \right) \cdot \psi_{\text{root} \rightarrow \text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \\ &= \psi_{\text{root} \rightarrow \text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) = 1 = D_{\text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}), \end{aligned}$$

where the first equality follows since $\psi_{j \rightarrow \text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) = 0$ for all $j \neq \text{root}$, and the second follows since $\lambda_{\text{pa}(\text{root})} = +\infty$. This establishes the base case. Now, suppose the claim is true for all nodes k of height H . Consider any node k with height $H - 1$. Using the definition of $D_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$, it follows that

$$\begin{aligned} D_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) &= \mathbb{1}_{\{c^q \in \mathbb{T}_k\}} \cdot \left(\frac{1}{\lambda_k} - \frac{1}{\lambda_{\text{pa}(k)}} \right) + \psi_{\text{pa}(k) \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \cdot D_{\text{pa}(k)}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \\ &= \mathbb{1}_{\{c^q \in \mathbb{T}_k\}} \cdot \left(\frac{1}{\lambda_k} - \frac{1}{\lambda_{\text{pa}(k)}} \right) + \psi_{\text{pa}(k) \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \cdot \left(\sum_{j \in \mathbb{T}} \mathbb{1}_{\{c^q \in \mathbb{T}_j\}} \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \psi_{j \rightarrow \text{pa}(k)}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \right) \end{aligned}$$

$$\begin{aligned}
&= \mathbb{1}_{\{c^q \in \mathcal{T}_k\}} \cdot \left(\frac{1}{\lambda_k} - \frac{1}{\lambda_{\text{pa}(k)}} \right) + \psi_{\text{pa}(k) \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \cdot \left(\sum_{j \in \mathcal{T} \setminus \{k\}} \mathbb{1}_{\{c^q \in \mathcal{T}_j\}} \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \psi_{j \rightarrow \text{pa}(k)}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \right) \\
&= \mathbb{1}_{\{c^q \in \mathcal{T}_k\}} \cdot \left(\frac{1}{\lambda_k} - \frac{1}{\lambda_{\text{pa}(k)}} \right) + \sum_{j \in \mathcal{T} \setminus \{k\}} \mathbb{1}_{\{c^q \in \mathcal{T}_j\}} \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \psi_{j \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \\
&= \sum_{j \in \mathcal{T}} \mathbb{1}_{\{c^q \in \mathcal{T}_j\}} \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \psi_{j \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}),
\end{aligned}$$

where the second equality follows from the induction hypothesis, the third since $\psi_{k \rightarrow \text{pa}(k)}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) = 0$, the fourth since $\psi_{j \rightarrow \text{pa}(k)}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \times \psi_{\text{pa}(k) \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) = \psi_{j \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$ for all $j \in \text{path}[\text{root}, k)$, and the final since $\psi_{k \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) = 1$. This completes the induction. \blacksquare

Here is the derivative of the negative log-likelihood function with respect to $\boldsymbol{\mu}$.

Lemma D.2 (Derivatives with Respect to $\boldsymbol{\mu}$) For all $\ell \in \mathcal{N}$,

$$\frac{\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \mu_\ell} = \frac{1}{Q} \sum_{q=1}^Q D_\ell^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) = \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathcal{T} \setminus \mathcal{N}} \mathbb{1}_{\{c^q \in \mathcal{T}_j\}} \cdot \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) - \frac{\text{sales}_\ell}{\lambda_{\text{pa}(\ell)}}$$

Proof. From Theorem 2.2, it follows that

$$\text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda}) = \frac{1}{Q} \sum_{q=1}^Q \left\{ \sum_{j \in \mathcal{T}} \mathbb{1}_{\{c^q \in \mathcal{T}_j\}} \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \right\}.$$

Now, noting that $\partial W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) / \partial \mu_\ell$ is equal to $\psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$ (from Lemma B.4 in Appendix B), it follows that for all $\ell \in \mathcal{N}$:

$$\frac{\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \mu_\ell} = \frac{1}{Q} \sum_{q=1}^Q \left\{ \sum_{j \in \mathcal{T}} \mathbb{1}_{\{c^q \in \mathcal{T}_j\}} \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \right\} = \frac{1}{Q} \sum_{q=1}^Q D_\ell^q(\boldsymbol{\mu}, \boldsymbol{\lambda}),$$

where the second equality follows from Lemma D.1. Finally, we note that

$$\begin{aligned}
\frac{1}{Q} \sum_{q=1}^Q D_\ell^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) &= \frac{1}{Q} \sum_{q=1}^Q \left\{ \sum_{j \in \mathcal{T}} \mathbb{1}_{\{c^q \in \mathcal{T}_j\}} \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \right\} \\
&= \frac{1}{Q} \sum_{q=1}^Q \left\{ \sum_{j \in \mathcal{T} \setminus \mathcal{N}} \mathbb{1}_{\{c^q \in \mathcal{T}_j\}} \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) + \sum_{j \in \mathcal{N}} \mathbb{1}_{\{c^q \in \mathcal{T}_j\}} \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \right\} \\
&= \frac{1}{Q} \sum_{q=1}^Q \left\{ \sum_{j \in \mathcal{T} \setminus \mathcal{N}} \mathbb{1}_{\{c^q \in \mathcal{T}_j\}} \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) - \frac{\mathbb{1}_{\{c^q = \ell\}}}{\lambda_{\text{pa}(\ell)}} \right\} \\
&= \frac{1}{Q} \sum_{q=1}^Q \left\{ \sum_{j \in \mathcal{T} \setminus \mathcal{N}} \mathbb{1}_{\{c^q \in \mathcal{T}_j\}} \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \right\} - \frac{1}{Q} \sum_{q=1}^Q \frac{\mathbb{1}_{\{c^q = \ell\}}}{\lambda_{\text{pa}(\ell)}} \\
&= \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathcal{T} \setminus \mathcal{N}} \mathbb{1}_{\{c^q \in \mathcal{T}_j\}} \cdot \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) - \frac{\text{sales}_\ell}{\lambda_{\text{pa}(\ell)}},
\end{aligned}$$

where the third equality follows since $\lambda_j = +\infty$ for all $j \in \mathcal{N}$ and $\psi_{j \rightarrow \ell}^q = 0$ for all $j \in \mathcal{N}; j \neq \ell$, and the last follows from the definition of sales_ℓ . \blacksquare

The next lemma gives the derivative with respect to λ .

Lemma D.3 (Derivatives with Respect to λ) For all $k \in \mathbb{T} \setminus (\mathcal{N} \cup \{\text{root}\})$,

$$\frac{\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_k} = \frac{1}{Q} \sum_{q=1}^Q \Delta_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \cdot D_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) + E_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}),$$

where for all q , $\Delta_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) := \Delta_k(S^q; \boldsymbol{\mu}, \boldsymbol{\lambda})$ is as defined in Lemma B.4, and

$$E_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) = \frac{\sum_{i \in \text{Children}(k)} \mathbb{1}_{\{c^q \in \mathbb{T}_i\}} \log(\psi_{k \rightarrow i}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}))}{\lambda_k}.$$

Proof. From Theorem 2.2, it follows that

$$\text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda}) = \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T}} \mathbb{1}_{\{c^q \in \mathbb{T}_j\}} \frac{W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_j} - \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T}} \mathbb{1}_{\{c^q \in \mathbb{T}_j\}} \frac{W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(j)}}$$

From the quotient rule of derivatives, it follows that

$$\begin{aligned} \frac{\partial}{\partial \lambda_k} \left(\frac{W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_j} \right) &= \frac{1}{\lambda_j} \cdot \frac{\partial W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_k} - \mathbb{1}_{\{j=k\}} \frac{W_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_k^2} \\ \frac{\partial}{\partial \lambda_k} \left(\frac{W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(j)}} \right) &= \frac{1}{\lambda_{\text{pa}(j)}} \cdot \frac{\partial W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_k} - \mathbb{1}_{\{\text{pa}(j)=k\}} \frac{W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_k^2} \end{aligned}$$

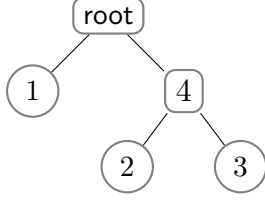
Using the above and the expression for $\frac{\partial W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_k}$ from Lemma B.4, it follows that

$$\begin{aligned} \frac{\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_k} &= \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T}} \mathbb{1}_{\{c^q \in \mathbb{T}_j\}} \cdot \left(\frac{\psi_{j \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \Delta_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_j} - \mathbb{1}_{\{k=j\}} \frac{W_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_k^2} \right) \\ &\quad - \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T}} \mathbb{1}_{\{c^q \in \mathbb{T}_j\}} \cdot \left(\frac{\psi_{j \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \Delta_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(j)}} - \mathbb{1}_{\{k=\text{pa}(j)\}} \frac{W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_k^2} \right) \\ &= \frac{1}{Q} \sum_{q=1}^Q \Delta_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \cdot \left(\sum_{j \in \mathbb{T}} \mathbb{1}_{\{c^q \in \mathbb{T}_j\}} \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \psi_{j \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \right) \\ &\quad - \frac{1}{Q} \sum_{q=1}^Q \left(\mathbb{1}_{\{c^q \in \mathbb{T}_k\}} W_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) - \sum_{i \in \text{Children}(k)} \mathbb{1}_{\{c^q \in \mathbb{T}_i\}} W_i^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \right) \frac{1}{\lambda_k^2} \\ &= \frac{1}{Q} \sum_{q=1}^Q \Delta_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \cdot D_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) - \frac{1}{Q} \sum_{q=1}^Q \sum_{i \in \text{Children}(k)} \mathbb{1}_{\{c^q \in \mathbb{T}_i\}} \cdot (W_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) - W_i^q(\boldsymbol{\mu}, \boldsymbol{\lambda})) \cdot \frac{1}{\lambda_k^2} \\ &= \frac{1}{Q} \sum_{q=1}^Q \Delta_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \cdot D_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) - \frac{1}{Q} \sum_{q=1}^Q \sum_{i \in \text{Children}(k)} \mathbb{1}_{\{c^q \in \mathbb{T}_i\}} \frac{-\log(\psi_{k \rightarrow i}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})) \cdot \lambda_k}{\lambda_k^2} \\ &= \frac{1}{Q} \sum_{q=1}^Q \Delta_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \cdot D_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) + \frac{1}{Q} \sum_{q=1}^Q E_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \end{aligned}$$

where the third equality follows from the expression for $D_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$ in Lemma D.1 and the fact that $\mathbb{1}_{\{c^q \in \mathbb{T}_k\}} = \sum_{i \in \text{Children}(k)} \mathbb{1}_{\{c^q \in \mathbb{T}_i\}}$ for all q , the fourth equality follows from the definition of $\psi_{k \rightarrow i}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$ and the last follows from the definition of $E_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$ in the statement of the lemma.

Appendix E: Proofs of claims and results from Section 3

E.1 Non-convexity of NegLog in λ



Suppose $\mathcal{N} = \{1, 2, 3\}$ and the tree structure is shown in the figure to the left, where **root** has two children: a leaf node 1 and a non-leaf node 4. Node 4 has two children: leaf nodes 2 and 3. The parameters of the model consist of μ_1, μ_2, μ_3 , and λ_4 . Suppose we offer the full assortment, and for each $\ell \in \mathcal{N}$, let $s_\ell \in \mathbb{Z}_{++}$ denote the number of customers who select product ℓ . Then,

$$\begin{aligned} \psi_{\text{root} \rightarrow 1}(\boldsymbol{\mu}, \lambda_4) &= \psi_{\text{root} \rightarrow 1}(\boldsymbol{\mu}, \lambda_4) &&= \frac{e^{\mu_1}}{e^{\mu_1} + [e^{\mu_2/\lambda_4} + e^{\mu_3/\lambda_4}]^{\lambda_4}} \\ \psi_{\text{root} \rightarrow 2}(\boldsymbol{\mu}, \lambda_4) &= \psi_{\text{root} \rightarrow 4}(\boldsymbol{\mu}, \lambda_4) \times \psi_{4 \rightarrow 2}(\boldsymbol{\mu}, \lambda_4) &&= \frac{[e^{\mu_2/\lambda_4} + e^{\mu_3/\lambda_4}]^{\lambda_4}}{e^{\mu_1} + [e^{\mu_2/\lambda_4} + e^{\mu_3/\lambda_4}]^{\lambda_4}} \times \frac{e^{\mu_2/\lambda_4}}{e^{\mu_2/\lambda_4} + e^{\mu_3/\lambda_4}} \\ \psi_{\text{root} \rightarrow 3}(\boldsymbol{\mu}, \lambda_4) &= \psi_{\text{root} \rightarrow 4}(\boldsymbol{\mu}, \lambda_4) \times \psi_{4 \rightarrow 3}(\boldsymbol{\mu}, \lambda_4) &&= \frac{[e^{\mu_2/\lambda_4} + e^{\mu_3/\lambda_4}]^{\lambda_4}}{e^{\mu_1} + [e^{\mu_2/\lambda_4} + e^{\mu_3/\lambda_4}]^{\lambda_4}} \times \frac{e^{\mu_3/\lambda_4}}{e^{\mu_2/\lambda_4} + e^{\mu_3/\lambda_4}}, \end{aligned}$$

which implies that

$$\begin{aligned} \text{NegLog}(\boldsymbol{\mu}, \lambda_4) &= \left((s_1 + s_2 + s_3) \log \left(e^{\mu_1} + [e^{\mu_2/\lambda_4} + e^{\mu_3/\lambda_4}]^{\lambda_4} \right) + (s_2 + s_3)(1 - \lambda_4) \log \left(e^{\mu_2/\lambda_4} + e^{\mu_3/\lambda_4} \right) \right. \\ &\quad \left. - s_1 \mu_1 - s_2 \frac{\mu_2}{\lambda_4} - s_3 \frac{\mu_3}{\lambda_4} \right) / (s_1 + s_2 + s_3). \end{aligned}$$

Suppose we have five customers, with $(s_1, s_2, s_3) = (1, 1, 3)$. For $\boldsymbol{\mu} = (\mu_1, \mu_2, \mu_3) = (0, 1, 1.03)$, we have $\text{NegLog}(\boldsymbol{\mu}, \lambda_4 = 0.1) = 1.0116$ and $\text{NegLog}(\boldsymbol{\mu}, \lambda_4 = 0.3) = 1.0381$, and we have that

$$\text{NegLog}(\boldsymbol{\mu}, \lambda_4 = 0.2) = 1.0317 > \frac{\text{NegLog}(\boldsymbol{\mu}, \lambda_4 = 0.1) + \text{NegLog}(\boldsymbol{\mu}, \lambda_4 = 0.3)}{2} = 1.0249,$$

which shows that $\lambda_4 \mapsto \text{NegLog}(\boldsymbol{\mu}, \lambda_4)$ is not convex in λ_4 .

E.2 Proof of Theorem 3.1

The first part of the theorem follows from the PROOF OF SUFFICIENCY argument in the proof of Theorem 2.4 in Appendix C.2 above.

For the second part, note that from Theorem 2.2, it follows that

$$\begin{aligned} \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda}) &= \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \frac{W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_j} - \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \frac{W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(j)}} \\ &= \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \frac{W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_j} - \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \{\text{root}\}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \frac{W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(j)}}, \end{aligned}$$

where the second equality follows since by definition, $\lambda_\ell = +\infty$ for all $\ell \in \mathcal{N}$ and $\lambda_{\text{pa}(\text{root})} = +\infty$.

Then, substituting $\lambda_j = e^{-\delta_{[j]}}$ for all $j \in \mathbb{T} \setminus \mathcal{N}$, it follows that $\text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\delta}) = F_1(\boldsymbol{\mu}, \boldsymbol{\delta}) - F_2(\boldsymbol{\mu}, \boldsymbol{\delta})$.

Next, we need to show that F_1 and F_2 are strictly convex in $\boldsymbol{\delta}$ for any fixed $\boldsymbol{\mu}$. We first show that the function $F_1(\boldsymbol{\mu}, \boldsymbol{\delta})$ is convex in $\boldsymbol{\delta}$ and then establish later that it is strictly convex. By the homogeneity property of the weight function in Lemma B.1 above, we have

$$F_1(\boldsymbol{\mu}, \boldsymbol{\delta}) = \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda}(\boldsymbol{\delta})) \cdot e^{\delta_{[j]}} = \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} W_j^q(\boldsymbol{\mu} \cdot e^{\delta_{[j]}}, \boldsymbol{\lambda}(\boldsymbol{\delta}) \cdot e^{\delta_{[j]}}),$$

where recall that the vector $\boldsymbol{\lambda}(\boldsymbol{\delta}) = (\exp(-\delta_{[j]}): j \in \mathbb{T} \setminus \mathcal{N})$. Then, for each term in the above summation, the mapping $\boldsymbol{\delta} \mapsto W_j^q(\boldsymbol{\mu} \cdot e^{\delta_{[j]}}, \boldsymbol{\lambda}(\boldsymbol{\delta}) \cdot e^{\delta_{[j]}})$ is convex because it is a composition of an increasing convex function $W_j^q(\cdot, \cdot)$ (by Lemmas B.2 and B.3) with a collection of convex functions, $\mu_\ell \cdot e^{\delta_{[j]}} = \mu_\ell \cdot e^{\sum_{k \in \text{path}[\text{root}, j]} \delta_k}$, and $e^{-\delta_{[k]}} \cdot e^{\delta_{[j]}} = e^{\sum_{h \in \text{path}[\text{root}, j]} \delta_h - \sum_{h \in \text{path}[\text{root}, k]} \delta_h}$. Note that $\mu_\ell \cdot e^{\delta_{[j]}}$ and $e^{-\delta_{[k]}} \cdot e^{\delta_{[j]}}$ are convex in $\boldsymbol{\delta}$ because the expressions in the exponents are linear functions of $\boldsymbol{\delta}$. Therefore, $F_1(\boldsymbol{\mu}, \boldsymbol{\delta})$ is convex in $\boldsymbol{\delta}$.

To establish strict convexity, we show that one of the terms in the summand is strictly convex. Consider the term corresponding to the root node:

$$f(\boldsymbol{\mu}, \boldsymbol{\delta}) = \sum_{q=1}^Q W_{\text{root}}^q(\boldsymbol{\mu} \cdot e^{\delta_{[\text{root}]}}), \boldsymbol{\lambda}(\boldsymbol{\delta}) \cdot e^{\delta_{[\text{root}]}}) = \sum_{q=1}^Q W_{\text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}(\boldsymbol{\delta})) \quad \text{since } \delta_{[\text{root}]} = \delta_{\text{root}} = 0.$$

We now show that $f(\boldsymbol{\mu}, \boldsymbol{\delta})$ is strictly convex in $\boldsymbol{\delta}$. For any $\boldsymbol{\delta} \neq \bar{\boldsymbol{\delta}}$ and scalar $x \in (0, 1)$, we want to show that

$$f(\boldsymbol{\mu}, x \cdot \boldsymbol{\delta} + (1-x) \cdot \bar{\boldsymbol{\delta}}) < x f(\boldsymbol{\mu}, \boldsymbol{\delta}) + (1-x) f(\boldsymbol{\mu}, \bar{\boldsymbol{\delta}}).$$

Since $W_{\text{root}}^{q'}(\boldsymbol{\mu}, \boldsymbol{\lambda}(\boldsymbol{\delta}))$ is convex in $\boldsymbol{\delta}$ for all $q' \in \{1, \dots, Q\}$, it suffices to exhibit one q such that $W_{\text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}(x\boldsymbol{\delta} + (1-x)\bar{\boldsymbol{\delta}})) < x W_{\text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}(\boldsymbol{\delta})) + (1-x) W_{\text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}}))$.

Now, since $\boldsymbol{\delta} \neq \bar{\boldsymbol{\delta}}$, there must exist a non-leaf node j such that $\delta_{[j]} \neq \bar{\delta}_{[j]}$. From the condition in Theorem 2.5, there exists a $q \in \{1, \dots, Q\}$ such that j has at least two children in \mathbb{T}^q . It then follows from Lemmas B.2 and B.3 that $W_{\text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$ is convex and increasing in $\boldsymbol{\lambda}$, and strictly increasing in λ_j . Further, because of the strict convexity of the exponential function, we have that $e^{x \cdot (-\delta_{[j]}) + (1-x) \cdot (-\bar{\delta}_{[j]})} < x e^{-\delta_{[j]}} + (1-x) e^{-\bar{\delta}_{[j]}}$. Together, these facts imply that

$$\begin{aligned} W_{\text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}(x\boldsymbol{\delta} + (1-x)\bar{\boldsymbol{\delta}})) &< W_{\text{root}}^q(\boldsymbol{\mu}, x\boldsymbol{\lambda}(\boldsymbol{\delta}) + (1-x)\boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})) \\ &\leq x W_{\text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}(\boldsymbol{\delta})) + (1-x) W_{\text{root}}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})). \end{aligned}$$

Consequently, F_1 is strictly convex. The proof of the strict convexity of F_2 follows from an almost identical argument, but with the function f re-defined as $f(\boldsymbol{\mu}, \boldsymbol{\delta}) = \sum_{q=1}^Q \sum_{j \in \text{Children}(\text{root}) \cap \mathbb{T}^q} W_j^q(\boldsymbol{\mu}, \boldsymbol{\lambda}(\boldsymbol{\delta}))$. ■

E.3 Step 2 of A-MM algorithm can be solved in closed form

We first show that the partial derivatives $(\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\delta}) / \partial \delta_k : k \in \mathbb{T} \setminus (\mathcal{N} \cup \{\text{root}\}))$ can also be computed efficiently:

Lemma E.1 For all $k \in \mathbb{T} \setminus (\mathcal{N} \cup \{\text{root}\})$, it follows that

$$\frac{\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\delta})}{\partial \delta_k} = - \sum_{j \in \mathbb{T}_k \setminus \mathcal{N}} \lambda_j \cdot \frac{\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_j},$$

Proof. The result follows by invoking the multi-variable chain rule of derivatives:

$$\begin{aligned} \frac{\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\delta})}{\partial \delta_k} &= \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \frac{\partial \lambda_j}{\partial \delta_k} \cdot \frac{\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_j} \\ &= \sum_{j \in \mathbb{T}_k \setminus \mathcal{N}} \frac{\partial \lambda_j}{\partial \delta_k} \cdot \frac{\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_j} \\ &= \sum_{j \in \mathbb{T}_k \setminus \mathcal{N}} -\lambda_j \cdot \frac{\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \lambda_j}, \end{aligned}$$

where the second inequality follows since $\lambda_j = e^{-\sum_{i \in \text{path}[\text{root}, j]} \delta_i}$ and the last since $\frac{\partial \lambda_j}{\partial \delta_k} = -\lambda_j$ for all $j \in \mathbb{T}_k \setminus \mathcal{N}$. The claim then follows. \blacksquare

Next, define the function $J(\alpha) := H^{(s)} \left(\left\{ \boldsymbol{\delta}^{(s)} - \alpha \cdot \nabla_{\boldsymbol{\delta}} \text{NegLog}(\boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)}) \right\}^+ \right)$ for all $\alpha \in \mathbb{R}_+$. The following lemma establishes the piecewise convexity of $J(\cdot)$:

Lemma E.2 For each $j \in \mathbb{T} \setminus (\mathcal{N} \cup \{\text{root}\})$, let $d_j = \frac{\partial \text{NegLog}}{\partial \delta_j} \left(\boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)} \right)$ and define t_j as

$$t_j = \begin{cases} +\infty & \text{if } d_j \leq 0 \\ \delta_j^{(s)} / d_j & \text{otherwise} \end{cases}.$$

Next, let $0 = t_{(0)} < t_{(1)} < \dots < t_{(I)} = +\infty$ denote the sorted values in the set $\{t_j : j \in \mathbb{T} \setminus (\mathcal{N} \cup \{\text{root}\})\} \cup \{0, +\infty\}$.

Then, for each $i \in \{0, 1, \dots, I-1\}$, the function $J(\alpha)$ is either **constant** or **strictly convex** on the interval $[t_{(i)}, t_{(i+1)}]$. In particular, $J(\cdot)$ is piecewise convex on \mathbb{R}_+ with I pieces, where $I \leq |\mathbb{T} \setminus \mathcal{N}|$.

Proof. Define the vector $\mathbf{M}(\alpha) = \left\{ \boldsymbol{\delta}^{(s)} - \alpha \mathbf{d} \right\}^+$ for any $\alpha \geq 0$. Then, for any $\alpha \in [t_{(i)}, t_{(i+1)}]$ and any $j \in \mathbb{T} \setminus (\mathcal{N} \cup \{\text{root}\})$, it can be verified that

$$M_j(\alpha) = \begin{cases} \delta_j^{(s)} - \alpha \cdot d_j & \text{if } t_{(i+1)} \leq t_j \\ 0 & \text{otherwise} \end{cases}. \quad (\text{EC.6})$$

Now, suppose the function $J(\alpha)$ is *not* constant on the interval $[t_{(i)}, t_{(i+1)}]$. In other words, there exists $\bar{\alpha}, \hat{\alpha} \in [t_{(i)}, t_{(i+1)}]$ with $\bar{\alpha} \neq \hat{\alpha}$ such that $J(\bar{\alpha}) \neq J(\hat{\alpha})$. Since $J(\alpha) = H^{(s)}(\mathbf{M}(\alpha))$, we must

have $\mathbf{M}(\bar{\alpha}) \neq \mathbf{M}(\hat{\alpha})$. From (EC.6), this further implies that there exists $j_i \in \mathbb{T} \setminus (\mathcal{N} \cup \{\text{root}\})$ such that $d_{j_i} \neq 0$ and $M_{j_i}(\alpha) = \delta_{j_i}^{(s)} - \alpha \cdot d_{j_i}$ for all $\alpha \in [t_{(i)}, t_{(i+1)}]$.

We show that $J(\alpha)$ must be strictly convex on $[t_{(i)}, t_{(i+1)}]$. For that, consider $\alpha_1, \alpha_2 \in [t_{(i)}, t_{(i+1)}]$ with $\alpha_1 \neq \alpha_2$ and $w \in (0, 1)$. From (EC.6), it can be verified that

$$\mathbf{M}(w \cdot \alpha_1 + (1 - w) \cdot \alpha_2) = w \cdot \mathbf{M}(\alpha_1) + (1 - w) \cdot \mathbf{M}(\alpha_2) \quad (\text{EC.7})$$

Moreover, since $M_{j_i}(\alpha_1) \neq M_{j_i}(\alpha_2)$ it follows that $\mathbf{M}(\alpha_1) \neq \mathbf{M}(\alpha_2)$. Then, it follows that

$$\begin{aligned} J(w \cdot \alpha_1 + (1 - w) \cdot \alpha_2) &= H^{(s)}(\mathbf{M}(w \cdot \alpha_1 + (1 - w) \cdot \alpha_2)) \\ &= H^{(s)}(w \cdot \mathbf{M}(\alpha_1) + (1 - w) \cdot \mathbf{M}(\alpha_2)) \\ &< w \cdot H^{(s)}(\mathbf{M}(\alpha_1)) + (1 - w) \cdot H^{(s)}(\mathbf{M}(\alpha_2)) \\ &= w \cdot J(\alpha_1) + (1 - w) \cdot J(\alpha_2), \end{aligned}$$

where the second equality follows from (EC.7) and the inequality follows since $\mathbf{M}(\alpha_1) \neq \mathbf{M}(\alpha_2)$ and the fact that $H^{(s)}(\cdot)$ is strictly convex as established in Lemma 3.8. This establishes the strict convex of $J(\cdot)$ on $[t_{(i)}, t_{(i+1)}]$. The result then follows from observing that $\cup_{i=0}^{I-1} [t_{(i)}, t_{(i+1)}] = \mathbb{R}_+$.

Finally, note that the number of pieces I satisfies

$$\begin{aligned} I &= |\{t_j : j \in \mathbb{T} \setminus (\mathcal{N} \cup \{\text{root}\})\} \cup \{0, +\infty\}| - 1 \\ &\leq |\{t_j : j \in \mathbb{T} \setminus (\mathcal{N} \cup \{\text{root}\})\}| + 2 - 1 \\ &= |\{t_j : j \in \mathbb{T} \setminus (\mathcal{N} \cup \{\text{root}\})\}| + 1 \\ &\leq |\mathbb{T} \setminus \mathcal{N}| \end{aligned}$$

■

The above lemma can be leveraged to efficiently compute the step size $\alpha^{(s)}$ in Step 2 of the A-MM algorithm as follows: we first obtain $\alpha_0, \alpha_1, \dots, \alpha_{I-1}$ as the following

$$\alpha_i = \arg \min_{\alpha \in [t_{(i)}, t_{(i+1)}]} J(\alpha),$$

which can be done efficiently since $J(\alpha)$ is either constant—in which case any α is optimal—or strictly convex—in which case well-known algorithms such as the golden-section search (Kiefer 1953) can be used—on each interval $[t_{(i)}, t_{(i+1)}]$. Then, we compute $\alpha^{(s)}$ as the following

$$\alpha^{(s)} = \arg \min_{\alpha \in \{\alpha_0, \alpha_1, \dots, \alpha_{I-1}\}} J(\alpha).$$

E.4 Proof of Theorem 3.3

We begin with the following key lemma that will be useful in the proof.

Lemma E.3 *For any $\boldsymbol{\mu} \in \text{Dom}_1$ and $\boldsymbol{\delta} \in \text{Dom}_2$, it follows that*

$$\sum_{\ell \in \mathcal{N}} a_\ell(\boldsymbol{\mu}, \boldsymbol{\delta}) \times \exp(\delta_{[\text{pa}(\ell)]}) = \sum_{\ell \in \mathcal{N}} \text{sales}_\ell \times \exp(\delta_{[\text{pa}(\ell)]}),$$

where $a_\ell(\boldsymbol{\mu}, \boldsymbol{\delta})$ is as defined in (5).

Proof. We leverage the fact that $\text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\delta})$ is shift-invariant in the variable $\boldsymbol{\mu}$. Specifically, for any $x \in \mathbb{R}$, we have that

$$\text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\delta}) = \text{NegLog}(\boldsymbol{\mu} + x \cdot \mathbf{e}, \boldsymbol{\delta}),$$

where $\mathbf{e} \in \mathbb{R}^N$ denotes the vector of all ones. In other words, for any $\boldsymbol{\mu}$ and $\boldsymbol{\delta}$, the function $g: \mathbb{R} \rightarrow \mathbb{R}$ defined as $g(x) = \text{NegLog}(\boldsymbol{\mu} + x \cdot \mathbf{e}, \boldsymbol{\delta})$ is a constant as a function of x . Therefore, the derivative $g'(x)$ of $g(x)$ with respect to x is zero everywhere.

We can compute the derivative of $g(\cdot)$ using the chain rule as follows:

$$g'(x) = \frac{dg}{dx} = \sum_{\ell \in \mathcal{N}} \frac{\partial \text{NegLog}(\boldsymbol{\mu} + x \cdot \mathbf{e}, \boldsymbol{\delta})}{\partial \mu_\ell} \frac{d}{dx}(\mu_\ell + x) = \sum_{\ell \in \mathcal{N}} \frac{\partial \text{NegLog}(\boldsymbol{\mu} + x \cdot \mathbf{e}, \boldsymbol{\delta})}{\partial \mu_\ell}.$$

Equating the above derivative to zero at $x = 0$ yields $\sum_{\ell \in \mathcal{N}} \partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\delta}) / \partial \mu_\ell = 0$. We now use the definition of the constant $a_\ell(\boldsymbol{\mu}, \boldsymbol{\delta})$ for each $\ell \in \mathcal{N}$ to obtain

$$0 = \sum_{\ell \in \mathcal{N}} \frac{\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\delta})}{\partial \mu_\ell} = \sum_{\ell \in \mathcal{N}} (a_\ell(\boldsymbol{\mu}, \boldsymbol{\delta}) - \text{sales}_\ell) \cdot \exp(\delta_{[\text{pa}(\ell)]})$$

The result of the lemma now follows. ■

We are now ready to prove Theorem 3.3. Below, we denote $\boldsymbol{\lambda}^{(s)} = \boldsymbol{\lambda}(\boldsymbol{\delta}^{(s)})$ to simplify certain expressions. In Theorem 3.7, we show that $\sum_{\ell \in \mathcal{N}} G_\ell(\mu_\ell | \boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})$ is a majorizing surrogate [see Definition 3.6] for the mapping $\boldsymbol{\mu} \mapsto \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\delta}^{(s)})$ at $\boldsymbol{\mu}^{(s)}$. From this, it follows that for all $\boldsymbol{\mu} \in \text{Dom}_1$:

$$\begin{aligned} \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\delta}^{(s)}) &\leq \sum_{\ell \in \mathcal{N}} G_\ell(\mu_\ell | \boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \\ &= \sum_{\ell \in \mathcal{N}} C_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) + \sum_{\ell \in \mathcal{N}} a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \exp\left(\frac{\mu_\ell - \mu_\ell^{(s)}}{\lambda_{\text{pa}(\ell)}^{(s)}}\right) - \sum_{\ell \in \mathcal{N}} \text{sales}_\ell \cdot \left(\frac{\mu_\ell - \mu_\ell^{(s)}}{\lambda_{\text{pa}(\ell)}^{(s)}}\right) \end{aligned} \quad (\text{EC.8})$$

and,

$$\text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) = \sum_{\ell \in \mathcal{N}} C_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) + \sum_{\ell \in \mathcal{N}} a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \quad (\text{EC.9})$$

Combining equations (EC.8) and (EC.9), we obtain that for all $\boldsymbol{\mu} \in \text{Dom}_1$:

$$\begin{aligned} & \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\delta}^{(s)}) - \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \\ & \leq \sum_{\ell \in \mathcal{N}} a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \left(\exp\left(\frac{\mu_\ell - \mu_\ell^{(s)}}{\lambda_{\text{pa}(\ell)}^{(s)}}\right) - 1 \right) - \sum_{\ell \in \mathcal{N}} \text{sales}_\ell \cdot \left(\frac{\mu_\ell - \mu_\ell^{(s)}}{\lambda_{\text{pa}(\ell)}^{(s)}}\right) \end{aligned}$$

Now, noting that $\tilde{\mu}_\ell^{(s+1)} = \mu_\ell^{(s)} + \lambda_{\text{pa}(\ell)}^{(s)} \cdot \log\left(\text{sales}_\ell / a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})\right)$ for all $\ell \in \mathcal{N}$, we get (we show in the proof of Theorem 3.7 below that $a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) > 0$ for all s)

$$\begin{aligned} & \text{NegLog}(\tilde{\boldsymbol{\mu}}^{(s+1)}, \boldsymbol{\delta}^{(s)}) - \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \\ & \leq \sum_{\ell \in \mathcal{N}} a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \left(\frac{\text{sales}_\ell}{a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})} - 1 \right) - \sum_{\ell \in \mathcal{N}} \text{sales}_\ell \cdot \log\left(\frac{\text{sales}_\ell}{a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})}\right) \\ & = \sum_{\ell \in \mathcal{N}} \left(\text{sales}_\ell - a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) - \text{sales}_\ell \cdot \log\left(\frac{\text{sales}_\ell}{a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})}\right) \right). \end{aligned} \quad (\text{EC.10})$$

We first establish that in each iteration, the MM update makes a non-positive improvement (recall that we are minimizing the objective function) by showing that the upper bound in equation (EC.10) above is always non-positive. In fact, we establish a stronger result that each term in the upper bound summation is non-positive; that is, $\text{sales}_\ell - a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) - \text{sales}_\ell \cdot \log\left(\text{sales}_\ell / a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})\right) \leq 0$ for all $\ell \in \mathcal{N}$.

To that end, for each $\ell \in \mathcal{N}$, note that $\tilde{\mu}_\ell^{(s+1)}$ is the minimizer of $G_\ell(\cdot | \boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})$ by definition. Therefore, we have that

$$\begin{aligned} & G_\ell(\tilde{\mu}_\ell^{(s+1)} | \boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \leq G_\ell(\mu_\ell^{(s)} | \boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \\ \implies & a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \exp\left(\frac{\tilde{\mu}_\ell^{(s+1)} - \mu_\ell^{(s)}}{\lambda_{\text{pa}(\ell)}^{(s)}}\right) - \text{sales}_\ell \cdot \frac{\tilde{\mu}_\ell^{(s+1)} - \mu_\ell^{(s)}}{\lambda_{\text{pa}(\ell)}^{(s)}} \leq a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \\ \implies & a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \frac{\text{sales}_\ell}{a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})} - \text{sales}_\ell \cdot \log\left(\frac{\text{sales}_\ell}{a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})}\right) \leq a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \\ \implies & \text{sales}_\ell - a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) - \text{sales}_\ell \cdot \log\left(\frac{\text{sales}_\ell}{a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})}\right) \leq 0. \end{aligned}$$

It now follows from the arguments above that

$$\begin{aligned} & \text{NegLog}(\tilde{\boldsymbol{\mu}}^{(s+1)}, \boldsymbol{\delta}^{(s)}) - \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \\ & \leq \sum_{\ell \in \mathcal{N}} \left(\text{sales}_\ell - a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) - \text{sales}_\ell \cdot \log\left(\frac{\text{sales}_\ell}{a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})}\right) \right) \\ & \leq \sum_{\ell \in \mathcal{N}} \frac{\lambda_{\text{lower}}^{(s)}}{\lambda_{\text{pa}(\ell)}^{(s)}} \left[\text{sales}_\ell - a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) - \text{sales}_\ell \cdot \log\left(\frac{\text{sales}_\ell}{a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})}\right) \right] \\ & = \lambda_{\text{lower}} \cdot \left(\sum_{\ell \in \mathcal{N}} \frac{\text{sales}_\ell - a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})}{\lambda_{\text{pa}(\ell)}^{(s)}} \right) - \lambda_{\text{lower}} \cdot \sum_{\ell \in \mathcal{N}} \frac{\text{sales}_\ell}{\lambda_{\text{pa}(\ell)}^{(s)}} \cdot \log\left(\frac{\text{sales}_\ell}{a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})}\right) \end{aligned}$$

$$\begin{aligned}
&= \lambda_{\text{lower}} \cdot \left(\sum_{\ell \in \mathcal{N}} \left\{ \text{sales}_\ell - a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \right\} \times \exp\left(\delta_{[\text{pa}(\ell)]}^{(s)}\right) \right) - \lambda_{\text{lower}} \sum_{\ell \in \mathcal{N}} \text{sales}_\ell \cdot \exp\left(\delta_{[\text{pa}(\ell)]}^{(s)}\right) \cdot \log\left(\frac{\text{sales}_\ell}{a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})}\right) \\
&= -\lambda_{\text{lower}} \sum_{\ell \in \mathcal{N}} \text{sales}_\ell \cdot \exp\left(\delta_{[\text{pa}(\ell)]}^{(s)}\right) \cdot \log\left(\frac{\text{sales}_\ell}{a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})}\right),
\end{aligned}$$

where the first inequality follows because $\lambda_{\text{pa}(\ell)}^{(s)} = e^{-\sum_{k \in \text{path}[\text{root}, \text{pa}(\ell)]} \delta_k^{(s)}} \geq e^{-(\text{height}(\text{root})-1) \times \delta_{\text{upper}}} := \lambda_{\text{lower}}$ for all $\ell \in \mathcal{N}$ and the last equality follows from the result of Lemma E.3 which implies $\sum_{\ell \in \mathcal{N}} \left\{ \text{sales}_\ell - a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \right\} \times \exp\left(\delta_{[\text{pa}(\ell)]}^{(s)}\right) = 0$.

The last expression on the right hand side above has the form of KL-divergence between two probability distributions, except that the corresponding terms in the expressions do not sum to 1. To address that, we normalize $(a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}): \ell \in \mathcal{N})$ and $(\text{sales}_\ell: \ell \in \mathcal{N})$ as follows. We let $T(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})$ denote the sums $\sum_{\ell \in \mathcal{N}} a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \times \exp\left(\delta_{[\text{pa}(\ell)]}^{(s)}\right) = \sum_{\ell \in \mathcal{N}} \text{sales}_\ell \times \exp\left(\delta_{[\text{pa}(\ell)]}^{(s)}\right)$ and define

$$\bar{a}_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) = \frac{a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \cdot \exp\left(\delta_{[\text{pa}(\ell)]}^{(s)}\right)}{T(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})}; \quad \overline{\text{sales}}_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) = \frac{\text{sales}_\ell \cdot \exp\left(\delta_{[\text{pa}(\ell)]}^{(s)}\right)}{T(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})}.$$

It is clear from our definitions that $\bar{\boldsymbol{a}}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) = (\bar{a}_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}): \ell \in \mathcal{N})$ and $\overline{\text{sales}} = (\overline{\text{sales}}_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}): \ell \in \mathcal{N})$ are valid distributions over the set \mathcal{N} . It now follows from these definitions that

$$\begin{aligned}
&\text{NegLog}(\tilde{\boldsymbol{\mu}}^{(s+1)}, \boldsymbol{\delta}^{(s)}) - \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \\
&\leq -T(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \cdot \lambda_{\text{lower}} \sum_{\ell \in \mathcal{N}} \overline{\text{sales}}_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \log\left(\frac{\overline{\text{sales}}_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})}{\bar{a}_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})}\right) \\
&= -T(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \cdot \lambda_{\text{lower}} \cdot D_{KL}\left(\overline{\text{sales}}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \parallel \bar{\boldsymbol{a}}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})\right).
\end{aligned}$$

We can simplify the above expression further by invoking Pinsker's inequality (Csiszár and Körner 2011), which states that $D_{KL}\left(\overline{\text{sales}}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \parallel \bar{\boldsymbol{a}}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})\right) \geq 1/2 \left\| \overline{\text{sales}}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) - \bar{\boldsymbol{a}}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \right\|_1^2$. We now obtain

$$\begin{aligned}
&\text{NegLog}(\tilde{\boldsymbol{\mu}}^{(s+1)}, \boldsymbol{\delta}^{(s)}) - \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \\
&\leq -T(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \cdot \lambda_{\text{lower}} \cdot \frac{1}{2} \left\| \overline{\text{sales}}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) - \bar{\boldsymbol{a}}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \right\|_1^2 \\
&= -\frac{\lambda_{\text{lower}}}{2 \cdot T(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})} \left(\sum_{\ell \in \mathcal{N}} \left| \left\{ \text{sales}_\ell - a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \right\} \times \exp\left(\delta_{[\text{pa}(\ell)]}^{(s)}\right) \right| \right)^2 \\
&= -\frac{\lambda_{\text{lower}}}{2 \cdot T(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})} \left(\sum_{\ell \in \mathcal{N}} \left| \frac{\partial \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})}{\partial \mu_\ell} \right| \right)^2 = -\frac{\lambda_{\text{lower}}}{2 \cdot T(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})} \left\| \nabla_{\boldsymbol{\mu}} \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \right\|_1^2,
\end{aligned}$$

where the second equality follows from the definition of $a_\ell(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)})$.

To complete the proof, we note that

$$T(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) = \sum_{\ell \in \mathcal{N}} \text{sales}_\ell \times \exp\left(\delta_{[\text{pa}(\ell)]}^{(s)}\right) = \sum_{\ell \in \mathcal{N}} \frac{\text{sales}_\ell}{\lambda_{\text{pa}(\ell)}^{(s)}} \leq \sum_{\ell \in \mathcal{N}} \frac{\text{sales}_\ell}{\lambda_{\text{lower}}} = \frac{\sum_{\ell \in \mathcal{N}} \text{sales}_\ell}{\lambda_{\text{lower}}} = \frac{1}{\lambda_{\text{lower}}}.$$

The inequality follows because $\lambda_{\text{pa}(\ell)}^{(s)} \geq \lambda_{\text{lower}}$ for all $\ell \in \mathcal{N}$. The last equality follows from the definition of sales_ℓ , which denotes the fraction of sales of product ℓ in the dataset, so that $\sum_{\ell \in \mathcal{N}} \text{sales}_\ell = 1$.

Putting everything together, we obtain

$$\text{NegLog}(\tilde{\boldsymbol{\mu}}^{(s+1)}, \boldsymbol{\delta}^{(s)}) - \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \leq \frac{-\lambda_{\text{lower}}^2}{2} \left\| \nabla_{\boldsymbol{\mu}} \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \right\|_1^2.$$

The improvement bound for the MM update follows from observing that $\text{NegLog}(\boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)}) = \text{NegLog}(\tilde{\boldsymbol{\mu}}^{(s+1)}, \boldsymbol{\delta}^{(s)})$ since NegLog is shift-invariant w.r.t. $\boldsymbol{\mu}$. \square

E.5 Proof of Theorem 3.4

We begin with the following lemma:

Lemma E.4 *If $\lambda_j \geq \lambda_{\text{lower}}$ for all $j \in \mathbb{T} \setminus \mathcal{N}$, then the mapping $\boldsymbol{\mu} \mapsto \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda})$ is L -smooth with $L \leq 1/\lambda_{\text{lower}}^2$.*

Proof. We use the definition of smoothness stated in (Bubeck 2015, Section 3.2). A continuously differentiable function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is L -smooth if its gradient ∇f is L -Lipschitz continuous, i.e.

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L \|\mathbf{x} - \mathbf{y}\|_2 \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$$

For a twice continuously differentiable function, the above condition is equivalent to $\nabla^2 f(\mathbf{x}) \preceq L \mathbf{I}_n$ for all $\mathbf{x} \in \mathbb{R}^n$, where \mathbf{I}_n is the $n \times n$ identity matrix. In other words, all the eigenvalues of the Hessian are bounded above by L . We will use this definition to derive the smoothness constant L .

In our context, $f(\boldsymbol{\mu}) = \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda})$. We use the fact that the trace (sum of diagonal elements) of a matrix is equal to the sum of its eigenvalues. Consider the sum of the diagonal entries of the Hessian matrix $\nabla^2 f(\boldsymbol{\mu})$:

$$\begin{aligned} & \sum_{\ell \in \mathcal{N}} \frac{\partial^2 f(\boldsymbol{\mu})}{\partial \mu_\ell^2} \\ &= \sum_{\ell \in \mathcal{N}} \frac{\partial}{\partial \mu_\ell} \left(\frac{\partial f(\boldsymbol{\mu})}{\partial \mu_\ell} \right) \\ &\stackrel{(a)}{=} \sum_{\ell \in \mathcal{N}} \frac{\partial}{\partial \mu_\ell} \left(\frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbb{1}_{\{c^q \in \mathbb{T}_j\}} \cdot \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) - \frac{\text{sales}_\ell}{\lambda_{\text{pa}(\ell)}} \right) \\ &= \sum_{\ell \in \mathcal{N}} \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbb{1}_{\{c^q \in \mathbb{T}_j\}} \cdot \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \frac{\partial}{\partial \mu_\ell} \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \end{aligned}$$

$$\begin{aligned}
&\stackrel{(b)}{=} \sum_{\ell \in \mathcal{N}} \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \cdot \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \frac{\partial}{\partial \mu_\ell} \prod_{k \in \text{path}(j, \ell]} \psi_{\text{pa}(k) \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \\
&\stackrel{(c)}{=} \sum_{\ell \in \mathcal{N}} \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \cdot \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \cdot \left(\sum_{k \in \text{path}(j, \ell]} \frac{\partial \psi_{\text{pa}(k) \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \mu_\ell} \right) / \psi_{\text{pa}(k) \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \\
&\stackrel{(d)}{=} \sum_{\ell \in \mathcal{N}} \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \cdot \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \cdot \left(\sum_{k \in \text{path}(j, \ell]} \frac{\psi_{k \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) - \psi_{\text{pa}(k) \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_{\text{pa}(k)}} \right) \\
&\stackrel{(e)}{\leq} \sum_{\ell \in \mathcal{N}} \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \cdot \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \cdot \left(\sum_{k \in \text{path}(j, \ell]} \frac{\psi_{k \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) - \psi_{\text{pa}(k) \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_{\text{lower}}} \right) \\
&= \sum_{\ell \in \mathcal{N}} \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \cdot \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \cdot (\psi_{\ell \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) - \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})) / \lambda_{\text{lower}} \\
&\stackrel{(f)}{\leq} \sum_{\ell \in \mathcal{N}} \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \cdot \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) / \lambda_{\text{lower}} \\
&= \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \cdot \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \cdot \frac{\sum_{\ell \in \mathcal{N}} \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\lambda_{\text{lower}}} \\
&\stackrel{(g)}{=} \frac{1}{Q \cdot \lambda_{\text{lower}}} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \cdot \left(\frac{1}{\lambda_j} - \frac{1}{\lambda_{\text{pa}(j)}} \right) \\
&\stackrel{(h)}{=} \frac{1}{Q \cdot \lambda_{\text{lower}}} \sum_{q=1}^Q \frac{1}{\lambda_{\text{pa}(c^q)}} \\
&\stackrel{(i)}{\leq} \frac{1}{\lambda_{\text{lower}}^2}
\end{aligned}$$

where the justifications for the equalities and inequalities in (a) - (i) are given below.

- (a) The equality follows from the expression for $\frac{\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \mu_\ell}$ from Lemma D.2 with $\boldsymbol{\lambda} = \boldsymbol{\lambda}(\boldsymbol{\delta})$.
- (b) The equality since $\psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) = \prod_{k \in \text{path}(j, \ell]} \psi_{\text{pa}(k) \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$.
- (c) The equality follows from the product rule of derivatives.
- (d) The equality follows since $\psi_{\text{pa}(k) \rightarrow k}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) = e^{-(W_{\text{pa}(k)}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) - W_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda})) / \lambda_{\text{pa}(k)}}$ and $\frac{\partial W_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\partial \mu_\ell} = \psi_{k \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$ for all $k \in \mathbb{T} \setminus \mathcal{N}$ from Lemma B.4.
- (e) The inequality follows since $\lambda_j \geq \lambda_{\text{lower}}$ for all $j \in \mathbb{T} \setminus \mathcal{N}$ by assumption and $\psi_{k \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \geq \psi_{\text{pa}(k) \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda})$ for all $k \in \text{path}(j, \ell]$.
- (f) The equality follows since $\psi_{\ell \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) - \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) \leq 1$.
- (g) The equality follows since $\sum_{\ell \in \mathcal{N}} \psi_{j \rightarrow \ell}^q(\boldsymbol{\mu}, \boldsymbol{\lambda}) = 1$ for any j such that $c^q \in \mathbb{T}_j$.
- (h) The equality follows since $\frac{1}{\lambda_{\text{pa}(\text{root})}} = +\infty$.
- (i) The inequality follows since $\lambda_{\text{pa}(c^q)} \geq \lambda_{\text{lower}}$ by assumption.

Since $f(\boldsymbol{\mu})$ is convex on \mathbb{R}^n (see proof of Theorem 3.1 above), all the eigenvalues of the Hessian $\nabla^2 f(\boldsymbol{\mu})$ are non-negative, for any $\boldsymbol{\mu} \in \mathbb{R}^n$. Then, it follows that:

$$\begin{aligned} \max \text{ eigenvalue of } \nabla^2 f(\boldsymbol{\mu}) &\leq \text{sum of eignvalues of } \nabla^2 f(\boldsymbol{\mu}) \\ &= \text{trace of } \nabla^2 f(\boldsymbol{\mu}) \\ &= \sum_{\ell \in \mathcal{N}} \frac{\partial^2 f(\boldsymbol{\mu})}{\partial \mu_\ell^2} \\ &\leq \frac{1}{\lambda_{\text{lower}}^2} \end{aligned}$$

This establishes the smoothness of $f(\boldsymbol{\mu})$. \square

We are now ready to prove Theorem 3.4. Since $\delta_j^{(s)} \leq \delta_{\text{upper}}$ for all $j \in \mathbb{T} \setminus (\mathcal{N} \cup \{\text{root}\})$ by assumption, it follows that $\lambda_j^{(s)} \geq \lambda_{\text{lower}}$ for all $j \in \mathbb{T} \setminus \mathcal{N}$, where $\boldsymbol{\lambda}^{(s)} = \boldsymbol{\lambda}(\boldsymbol{\delta}^{(s)})$. Then, it follows from Lemma E.4 that $\text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\delta}^{(s)})$ is L -smooth. We then leverage the standard quadratic upper bound property (Bubeck 2015, Lemma 3.4) for smooth functions to obtain:

$$\begin{aligned} &\text{NegLog}(\boldsymbol{\mu}_{\text{GD}}^{(s+1)}, \boldsymbol{\delta}^{(s)}) \\ &\leq \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) + \left\langle \nabla_{\boldsymbol{\mu}} \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}), \boldsymbol{\mu}_{\text{GD}}^{(s+1)} - \boldsymbol{\mu}^{(s)} \right\rangle + \frac{L}{2} \left\| \boldsymbol{\mu}_{\text{GD}}^{(s+1)} - \boldsymbol{\mu}^{(s)} \right\|_2^2 \\ &= \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) - \frac{1}{L} \left\langle \nabla_{\boldsymbol{\mu}} \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}), \nabla_{\boldsymbol{\mu}} \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \right\rangle + \frac{1}{2L} \left\| \nabla_{\boldsymbol{\mu}} \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \right\|_2^2 \\ &\text{(using the definition of } \boldsymbol{\mu}_{\text{GD}}^{(s+1)}) \\ &= \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) - \frac{1}{2L} \left\| \nabla_{\boldsymbol{\mu}} \text{NegLog}(\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) \right\|_2^2 \end{aligned}$$

■

E.5.1 Tightness of improvement bound. We now present an instance where both the upper bound on the smoothness constant L in Lemma E.4 and the ℓ_1 - ℓ_2 norm inequality are tight, so that the improvement bound in Theorem 3.4 is the best possible bound.

For that, we consider the special where each offer set has exactly two products and customers make choices according to the MNL model. Specifically, suppose that we are given the following dataset consisting of $Q = 2 \cdot (n - 1)$ transactions:

$$(\mathcal{S}^q, c^q) = \begin{cases} (\{1, q + 1\}, q + 1) & \text{for } 1 \leq q \leq n - 1 \\ (\{1, q - n + 2\}, 1) & \text{for } n \leq q \leq Q \end{cases} \quad (\text{EC.11})$$

Tightness of upper bound for L . In the special case when each offer set has exactly two products, the MNL model is equivalent to the Bradley-Terry model. For the Bradley-Terry model, Vojnovic et al. (2020, Lemma 3.1) showed that the negative log-likelihood function is L -smooth on \mathbb{R}^n , where $L = \lambda_{\max}(\mathbf{L}_{\mathbf{M}})/4Q$.¹² Here, \mathbf{M} is an $n \times n$ matrix with $M_{i,j}$ equal to the

¹² Vojnovic et al. (2020) consider the un-normalized negative log-likelihood function, and therefore, do not have the normalization by the total number of transactions Q in their result.

number of times offer-set $\{i, j\}$ is observed in the transaction data, and $\lambda_{\max}(\mathbf{L}_{\mathbf{M}})$ is the largest eigenvalue of the Laplacian matrix $\mathbf{L}_{\mathbf{M}} = \mathbf{D}_{\mathbf{M}} - \mathbf{M}$, where $\mathbf{D}_{\mathbf{M}}$ is a diagonal matrix whose diagonal elements are the row sums of \mathbf{M} .

For the above dataset, it is easy to check that the product co-occurrence matrix $\mathbf{M} = 2\mathbf{M}'$ where \mathbf{M}' is the adjacency matrix of the undirected version of the comparison graph Comp defined in Definition 2.3, obtained by replacing the pair of directed edges $\{(1, j), (j, 1)\}$ by the undirected edge $\{1, j\}$ for all $j \in \mathcal{N} \setminus \{1\}$. We refer to this graph as Comp' . Now, since the degree of the node corresponding to product 1 in Comp' is equal to $n - 1$, it follows from Zhang (2011, Theorem 3.19) that $\lambda_{\max}(\mathbf{L}_{\mathbf{M}'}) = n$. Further, since $\mathbf{L}_{\mathbf{M}} = 2\mathbf{L}_{\mathbf{M}'}$, it follows that $L = \lambda_{\max}(\mathbf{L}_{\mathbf{M}})/4Q = 2\lambda_{\max}(\mathbf{L}_{\mathbf{M}'})/4Q = 2n/4Q = \frac{n}{4 \cdot (n-1)}$. For large n , we have that $L \approx 1/4$. Finally, note that $\lambda_{\text{lower}} = 1$ for the MNL model. This shows that the upper bound $1/\lambda_{\text{lower}}^2$ is tight up to constant factors.

Tightness of the ℓ_1 - ℓ_2 norm inequality. To establish this result, we exhibit an initial solution $\boldsymbol{\mu}^{(0)}$ for which the inequality is tight. First, note that for the dataset described in Equation (EC.11), it is easy to check that $\text{sales}_1 = \frac{n-1}{2 \cdot (n-1)} = 1/2$, whereas $\text{sales}_\ell = \frac{1}{2 \cdot (n-1)}$ for all $\ell \in \mathcal{N} \setminus \{1\}$. Further, from Lemma D.2 it follows that for all $\ell \in \mathcal{N}$

$$\frac{\partial \text{NegLog}(\boldsymbol{\mu})}{\partial \mu_\ell} = \frac{1}{Q} \sum_{q=1}^Q \psi_{\text{root} \rightarrow \ell}^q(\boldsymbol{\mu}) - \text{sales}_\ell, \quad (\text{EC.12})$$

where note that we have dropped the $\boldsymbol{\lambda}$ term because we are working with an MNL model. Now, supposing that n is even, consider an initial solution $\boldsymbol{\mu}^{(0)}$ of the following form:

$$\mu_\ell^{(0)} = \begin{cases} 0 & \text{if } \ell = 1 \\ \log(2n - 1) & \text{if } 1 < \ell \leq n/2 \\ -\log(2n - 1) & \text{o.w.} \end{cases}$$

Next, consider any $1 < \ell \leq n/2$. Then, it follows that for the dataset described in in Equation (EC.11):

$$\begin{aligned} \frac{1}{Q} \sum_{q=1}^Q \psi_{\text{root} \rightarrow \ell}^q(\boldsymbol{\mu}^{(0)}) &= \frac{1}{Q} \cdot (\psi_{\text{root} \rightarrow \ell}^{\ell-1}(\boldsymbol{\mu}^{(0)}) + \psi_{\text{root} \rightarrow \ell}^{n+\ell-2}(\boldsymbol{\mu}^{(0)})) \\ &= \frac{2}{Q} \frac{\exp(\mu_\ell^{(0)})}{1 + \exp(\mu_\ell^{(0)})} \\ &= \frac{2 \cdot (2n - 1)}{2 \cdot (n - 1) \cdot 2n} \\ &= \frac{2n - 1}{2n \cdot (n - 1)} \\ &= \frac{1}{2 \cdot (n - 1)} + \frac{1}{2n} \\ &= \text{sales}_\ell + \frac{1}{2n} \end{aligned} \quad (\text{EC.13})$$

In a similar fashion, it can be verified that

$$\frac{1}{Q} \sum_{q=1}^Q \psi_{\text{root} \rightarrow \ell}^q(\boldsymbol{\mu}^{(0)}) = \text{sales}_\ell - \frac{1}{2n} \quad \text{for all } n/2 < \ell \leq n \quad (\text{EC.14})$$

Then, combining (EC.12), (EC.13), and (EC.14), it follows that

$$\frac{\partial \text{NegLog}(\boldsymbol{\mu}^{(0)})}{\partial \mu_\ell} = \begin{cases} \frac{1}{2n} & \text{if } 1 < \ell \leq n/2 \\ -\frac{1}{2n} & \text{if } n/2 < \ell \leq n \end{cases} \quad (\text{EC.15})$$

Further, from the proof of Lemma E.3, it follows that

$$\begin{aligned} \sum_{\ell=1}^n \frac{\partial \text{NegLog}(\boldsymbol{\mu}^{(0)})}{\partial \mu_\ell} = 0 &\implies \frac{\partial \text{NegLog}(\boldsymbol{\mu}^{(0)})}{\partial \mu_1} + (n/2 - 1) \cdot 1/2n - (n/2) \cdot (1/2n) = 0 \\ &\implies \frac{\partial \text{NegLog}(\boldsymbol{\mu}^{(0)})}{\partial \mu_1} = \frac{1}{2n} \end{aligned} \quad (\text{EC.16})$$

Finally, given (EC.16) and (EC.15), it can easily be verified that $\|\nabla_{\boldsymbol{\mu}} \text{NegLog}(\boldsymbol{\mu}^{(0)})\|_2^2 = \frac{1}{n} \|\nabla_{\boldsymbol{\mu}} \text{NegLog}(\boldsymbol{\mu}^{(0)})\|_1^2$. ■

E.6 Proof of Theorem 3.5

We first establish that the constant D is finite, for which it is sufficient to show that the set $\mathcal{U} = \{\boldsymbol{\mu} \in \text{Dom}_1 : \text{NegLog}(\boldsymbol{\mu}) \leq \text{NegLog}(\boldsymbol{\mu}^{(0)})\}$ is bounded. From the PROOF OF SUFFICIENCY argument in the proof of Theorem 2.4 in Appendix C.2 above, it follows that for all $\boldsymbol{\mu} \in \text{Dom}_1$ (since $\text{Dom}_1 = \mathcal{D}_1$):

$$\text{NegLog}(\boldsymbol{\mu}) \geq \|\boldsymbol{\mu}\|_\infty b^*/Q,$$

where b^* is as defined in Appendix C.2 and $b^* > 0$ by Lemma C.1. Therefore, it follows that for any $\boldsymbol{\mu} \in \mathcal{U}$:

$$\text{NegLog}(\boldsymbol{\mu}) \leq \text{NegLog}(\boldsymbol{\mu}^{(0)}) \implies \|\boldsymbol{\mu}\|_\infty \leq \frac{\text{NegLog}(\boldsymbol{\mu}^{(0)}) \cdot Q}{b^*},$$

and therefore the set \mathcal{U} is bounded.

Next, we derive the convergence rates. We begin with the guarantee for the MM algorithm. The proof follows from existing results, see e.g., Allen-Zhu and Orecchia (2014, Fact B.1 in Appendix B). We reproduce the arguments here for completeness. For each $s' \geq 0$, define $\text{NegLogGap}^{(s')} := \text{NegLog}(\boldsymbol{\mu}^{(s')}) - \text{NegLog}(\boldsymbol{\mu}^*)$. Note that $\text{NegLogGap}^{(s')} \geq 0$ for all $s' \geq 0$. Moreover, since the result is trivially true for any $s \geq 1$ such that $\text{NegLogGap}^{(s)} = 0$, we suppose that $\text{NegLogGap}^{(s')} > 0$ for all $0 \leq s' \leq s$.

Now, recall that $\text{NegLog}(\boldsymbol{\mu})$ is convex in $\boldsymbol{\mu}$. Then, it follows that for any $s' \geq 0$:

$$\begin{aligned} \text{NegLogGap}^{(s')} &= \text{NegLog}(\boldsymbol{\mu}^{(s')}) - \text{NegLog}(\boldsymbol{\mu}^*) \\ &\leq \left\langle \nabla \text{NegLog}(\boldsymbol{\mu}^{(s')}), \boldsymbol{\mu}^{(s')} - \boldsymbol{\mu}^* \right\rangle \\ &\leq \left\| \nabla \text{NegLog}(\boldsymbol{\mu}^{(s')}) \right\|_1 \cdot \left\| \boldsymbol{\mu}^{(s')} - \boldsymbol{\mu}^* \right\|_\infty \\ &\leq \left\| \nabla \text{NegLog}(\boldsymbol{\mu}^{(s')}) \right\|_1 \cdot D \end{aligned} \quad (\text{EC.17})$$

where the first inequality follows from the subgradient inequality for convex functions (note that we write the gradient as ∇NegLog instead of $\nabla_{\mu} \text{NegLog}$), the second follows from the Holder's inequality, and the final follows since the MM algorithm guarantees an improving solution in each iteration and using the definition of D in the statement of the theorem.

Next, by plugging in $\lambda_{\text{lower}} = 1$ in the improvement bound from Theorem 3.3, it follows that for all $0 \leq s' < s$:

$$\begin{aligned}
& \text{NegLog}(\boldsymbol{\mu}^{(s'+1)}) - \text{NegLog}(\boldsymbol{\mu}^{(s')}) \leq -\frac{1}{2} \left\| \nabla \text{NegLog}(\boldsymbol{\mu}^{(s')}) \right\|_1^2 \\
\implies & \text{NegLogGap}^{(s'+1)} - \text{NegLogGap}^{(s')} \leq -\frac{1}{2} \left\| \nabla \text{NegLog}(\boldsymbol{\mu}^{(s')}) \right\|_1^2 \\
\implies & \text{NegLogGap}^{(s')} - \text{NegLogGap}^{(s'+1)} \geq \frac{1}{2} \left\| \nabla \text{NegLog}(\boldsymbol{\mu}^{(s')}) \right\|_1^2 \\
\implies & \text{NegLogGap}^{(s')} - \text{NegLogGap}^{(s'+1)} \geq \frac{\left(\text{NegLogGap}^{(s')} \right)^2}{2D^2} \\
\implies & \frac{\text{NegLogGap}^{(s')} - \text{NegLogGap}^{(s'+1)}}{\text{NegLogGap}^{(s')} \cdot \text{NegLogGap}^{(s'+1)}} \geq \frac{\text{NegLogGap}^{(s')}}{2D^2 \cdot \text{NegLogGap}^{(s'+1)}} \\
\implies & \frac{1}{\text{NegLogGap}^{(s'+1)}} - \frac{1}{\text{NegLogGap}^{(s')}} \geq \frac{\text{NegLogGap}^{(s')}}{2D^2 \cdot \text{NegLogGap}^{(s'+1)}} \\
\implies & \frac{1}{\text{NegLogGap}^{(s'+1)}} - \frac{1}{\text{NegLogGap}^{(s')}} \geq \frac{1}{2D^2},
\end{aligned}$$

where the third implication follows from (EC.17), the fourth follows from dividing both sides by $\text{NegLogGap}^{(s')} \cdot \text{NegLogGap}^{(s'+1)}$, and the final follows since $\text{NegLogGap}^{(s')} \geq \text{NegLogGap}^{(s'+1)}$ as the MM algorithm guarantees an improving solution in each iteration.

Summing the last inequality above from $s' = 0$ to $s - 1$ (note that we have $s \geq 1$), it follows that

$$\begin{aligned}
& \sum_{s'=0}^{s-1} \left(\frac{1}{\text{NegLogGap}^{(s'+1)}} - \frac{1}{\text{NegLogGap}^{(s')}} \right) \geq \sum_{s'=0}^{s-1} \frac{1}{2D^2} \\
\implies & \frac{1}{\text{NegLogGap}^{(s)}} - \frac{1}{\text{NegLogGap}^{(0)}} \geq \frac{s}{2D^2} \\
& \implies \frac{1}{\text{NegLogGap}^{(s)}} \geq \frac{s}{2D^2} \\
& \implies \text{NegLogGap}^{(s)} \leq \frac{2D^2}{s} \\
& \implies \text{NegLog}(\boldsymbol{\mu}^{(s)}) - \text{NegLog}(\boldsymbol{\mu}^*) \leq \frac{2D^2}{s}
\end{aligned}$$

The result then follows.

For the GD algorithm, starting from the improvement bound in Theorem 3.4 and using the fact that GD with step size $1/L$ also guarantees an improving solution in each iteration, the above sequence of arguments can be repeated to show that for all $s \geq 1$

$$\text{NegLog}(\boldsymbol{\mu}_{\text{GD}}^{(s)}) - \text{NegLog}(\boldsymbol{\mu}^*) \leq \frac{2LnD^2}{s}$$

The result then follows.

E.7 Proof of Theorem 3.7

The proof of the theorem makes use of the following lemma:

Lemma E.5 *Given any $\bar{\delta} \in \text{Dom}_2$, for each $q = 1, \dots, Q$ and each nonleaf node $j \in \mathbb{T}^q \setminus \mathcal{S}^q$, a majorizing surrogate for the function $\boldsymbol{\mu} \mapsto \frac{\sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} e^{W_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}(\bar{\delta})) \times \exp(\bar{\delta}_{[j]})}}{\sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} e^{W_k^q(\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}(\bar{\delta})) \times \exp(\bar{\delta}_{[j]})}}$ at $\bar{\boldsymbol{\mu}}$ is given by the following separable function:*

$$\boldsymbol{\mu} \mapsto \mathbb{C}(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}}) + \exp(\bar{\delta}_{[j]}) \sum_{\ell \in \mathcal{N}} \exp(-\bar{\delta}_{[\text{pa}(\ell)]}) \times \psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})) \times e^{(\mu_\ell - \bar{\mu}_\ell) \times \exp(\bar{\delta}_{[\text{pa}(\ell)])}},$$

where $\mathbb{C}(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}})$ is a constant depending only on $(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}})$ that is irrelevant for our optimization.

Proof. Fix an arbitrary q and let $\bar{\boldsymbol{\lambda}} = \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})$. We will prove the result by induction on the height of j . For the base case, suppose j has a height of one, that is, it is a parent of some leaf node. In this case, note that for each $\ell \in \text{Children}(j)$, $\exp(\bar{\delta}_{[j]}) \times \exp(-\bar{\delta}_{[\text{pa}(\ell)]}) = 1$. Then,

$$\begin{aligned} & \frac{\sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} e^{W_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})) \times \exp(\bar{\delta}_{[j]})}}{\sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} e^{W_k^q(\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})) \times \exp(\bar{\delta}_{[j]})}} \\ & \stackrel{(a)}{=} \frac{\sum_{\ell \in \text{Children}(j) \cap \mathcal{S}^q} e^{\mu_\ell / \bar{\lambda}_j}}{\sum_{i \in \text{Children}(j) \cap \mathcal{S}^q} e^{\bar{\mu}_i / \bar{\lambda}_j}} \\ & = \sum_{\ell \in \text{Children}(j) \cap \mathcal{S}^q} \frac{e^{\bar{\mu}_\ell / \bar{\lambda}_j}}{\sum_{i \in \text{Children}(j) \cap \mathcal{S}^q} e^{\bar{\mu}_i / \bar{\lambda}_j}} \times e^{(\mu_\ell - \bar{\mu}_\ell) / \bar{\lambda}_{\text{pa}(\ell)}} \\ & = \sum_{\ell \in \text{Children}(j) \cap \mathcal{S}^q} \psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \times e^{(\mu_\ell - \bar{\mu}_\ell) / \bar{\lambda}_{\text{pa}(\ell)}} \\ & \stackrel{(b)}{=} \sum_{\ell \in \mathcal{N}} \psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \times e^{(\mu_\ell - \bar{\mu}_\ell) / \bar{\lambda}_{\text{pa}(\ell)}} \\ & = \sum_{\ell \in \mathcal{N}} \psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})) \times e^{(\mu_\ell - \bar{\mu}_\ell) \times \exp(\bar{\delta}_{[\text{pa}(\ell)])}}, \end{aligned}$$

where (a) follows from the definition of the weight function at leaf nodes and because $\bar{\boldsymbol{\lambda}} = \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})$, and (b) follows because if $\ell \notin \text{Children}(j) \cap \mathcal{S}^q$, then $\ell \notin \mathbb{T}^q$, so $\psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) = 0$. This completes the base case.

To establish the induction step, suppose that the result holds for all nodes k of height at most H . In other words, for every nonleaf node $k \in \mathbb{T}^q$ of height at most H , there is a constant \mathbb{C}_k that depend only on $\bar{\boldsymbol{\mu}}$ and $\bar{\boldsymbol{\delta}}$ such that for all $\boldsymbol{\mu}$,

$$\begin{aligned} & \frac{\sum_{m \in \text{Children}(k) \cap \mathbb{T}^q} e^{W_m^q(\boldsymbol{\mu}, \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})) \times \exp(\bar{\delta}_{[k]})}}{\sum_{m \in \text{Children}(k) \cap \mathbb{T}^q} e^{W_m^q(\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})) \times \exp(\bar{\delta}_{[k]})}} \\ & \leq \mathbb{C}_k + \exp(\bar{\delta}_{[k]}) \sum_{\ell \in \mathcal{N}} \exp(-\bar{\delta}_{[\text{pa}(\ell)]}) \times \psi_{k \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})) \times e^{(\mu_\ell - \bar{\mu}_\ell) \times \exp(\bar{\delta}_{[\text{pa}(\ell)])}}, \end{aligned}$$

and the above inequality is tight at $\boldsymbol{\mu} = \bar{\boldsymbol{\mu}}$. Now, consider a nonleaf node j of height $H + 1$. Letting C denote a constant that depend only on $(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}})$, we have

$$\begin{aligned}
& \frac{\sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} e^{W_k^q(\boldsymbol{\mu}, \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})) \times \exp(\bar{\delta}_{[j]})}}{\sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} e^{W_k^q(\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})) \times \exp(\bar{\delta}_{[j]})}} \\
& \stackrel{(a)}{=} \frac{\sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} e^{W_k^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_j}}{\sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} e^{W_k^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_j}} \\
& \stackrel{(b)}{=} \sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} \frac{\left(\sum_{m \in \text{Children}(k) \cap \mathbb{T}^q} e^{W_m^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_k} \right)^{\bar{\lambda}_k / \bar{\lambda}_j}}{\sum_{s \in \text{Children}(j) \cap \mathbb{T}^q} e^{W_s^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_j}} \\
& \stackrel{(c)}{\leq} C + \sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} \frac{\bar{\lambda}_k}{\bar{\lambda}_j} \times \frac{\left(\sum_{m \in \text{Children}(k) \cap \mathbb{T}^q} e^{W_m^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_k} \right)^{\bar{\lambda}_k / \bar{\lambda}_j}}{\sum_{s \in \text{Children}(j) \cap \mathbb{T}^q} e^{W_s^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_j}} \times \frac{\sum_{m \in \text{Children}(k) \cap \mathbb{T}^q} e^{W_m^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_k}}{\sum_{m \in \text{Children}(k) \cap \mathbb{T}^q} e^{W_m^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_k}} \\
& \stackrel{(d)}{=} C + \sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} \frac{\bar{\lambda}_k}{\bar{\lambda}_j} \times \frac{e^{W_k^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_j}}{\sum_{s \in \text{Children}(j) \cap \mathbb{T}^q} e^{W_s^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_j}} \times \frac{\sum_{m \in \text{Children}(k) \cap \mathbb{T}^q} e^{W_m^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_k}}{\sum_{m \in \text{Children}(k) \cap \mathbb{T}^q} e^{W_m^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_k}} \\
& \stackrel{(e)}{=} C + \sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} \frac{\bar{\lambda}_k}{\bar{\lambda}_j} \times \psi_{j \rightarrow k}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \times \frac{\sum_{m \in \text{Children}(k) \cap \mathbb{T}^q} e^{W_m^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_k}}{\sum_{m \in \text{Children}(k) \cap \mathbb{T}^q} e^{W_m^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_k}} \\
& \stackrel{(f)}{\leq} C + \sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} C_k + \left(\frac{\bar{\lambda}_k}{\bar{\lambda}_j} \times \psi_{j \rightarrow k}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \times \frac{1}{\bar{\lambda}_k} \sum_{\ell \in \mathcal{N}} \bar{\lambda}_{\text{pa}(\ell)} \times \psi_{k \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \times e^{(\mu_\ell - \bar{\mu}_\ell) / \bar{\lambda}_{\text{pa}(\ell)}} \right) \\
& \stackrel{(g)}{=} C_j + \sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} \frac{1}{\bar{\lambda}_j} \times \psi_{j \rightarrow k}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \sum_{\ell \in \mathcal{S}^q \cap \mathbb{T}_k^q} \bar{\lambda}_{\text{pa}(\ell)} \times \psi_{k \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \times e^{(\mu_\ell - \bar{\mu}_\ell) / \bar{\lambda}_{\text{pa}(\ell)}} \\
& \stackrel{(h)}{=} C_j + \frac{1}{\bar{\lambda}_j} \sum_{\ell \in \mathcal{S}^q \cap \mathbb{T}_j^q} \bar{\lambda}_{\text{pa}(\ell)} \times \psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \times e^{(\mu_\ell - \bar{\mu}_\ell) / \bar{\lambda}_{\text{pa}(\ell)}} \\
& \stackrel{(i)}{=} C_j + \frac{1}{\bar{\lambda}_j} \sum_{\ell \in \mathcal{N}} \bar{\lambda}_{\text{pa}(\ell)} \times \psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \times e^{(\mu_\ell - \bar{\mu}_\ell) / \bar{\lambda}_{\text{pa}(\ell)}} \\
& \stackrel{(j)}{=} C_j + \exp(\bar{\delta}_{[j]}) \sum_{\ell \in \mathcal{N}} \exp(-\bar{\delta}_{[\text{pa}(\ell)]}) \times \psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})) \times e^{(\mu_\ell - \bar{\mu}_\ell) \times \exp(\bar{\delta}_{[\text{pa}(\ell)]})},
\end{aligned}$$

where the justifications for the equalities and inequalities in (a) - (j) are given below.

(a) The equality follows because $\bar{\boldsymbol{\lambda}} = \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})$.

(b) The equality follows from the definition that for each $k \in \text{Children}(j) \cap \mathbb{T}^q$, $W_k^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) = \bar{\lambda}_k \log \left(\sum_{m \in \text{Children}(k) \cap \mathbb{T}^q} e^{W_m^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_k} \right)$, so

$$e^{W_k^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_j} = \left(\sum_{m \in \text{Children}(k) \cap \mathbb{T}^q} e^{W_m^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_k} \right)^{\bar{\lambda}_k / \bar{\lambda}_j}.$$

(c) For the inequality (c), we apply the following sub-gradient inequality: for all $\alpha \in (0, 1]$ and $x, y \in \mathbb{R}_+$, $x^\alpha \leq y^\alpha + \alpha y^{\alpha-1}(x - y) = (1 - \alpha)y^\alpha + \alpha y^{\frac{\alpha}{y}}$, with equality if and only if $x = y$. The

inequality (c) follows from the application of the sub-gradient inequality where for each $k \in \text{Children}(k) \cap \mathbb{T}^q$, we let

$$x_k = \sum_{m \in \text{Children}(k) \cap \mathbb{T}^q} e^{W_m^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}})/\bar{\lambda}_k}, \quad y_k = \sum_{m \in \text{Children}(k) \cap \mathbb{T}^q} e^{W_m^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}})/\bar{\lambda}_k}, \quad \text{and} \quad \alpha_k = \frac{\bar{\lambda}_k}{\lambda_j},$$

and note that $\alpha_k \in (0, 1]$ because $\bar{\lambda}_k \leq \bar{\lambda}_j$ since $k \in \text{Children}(j)$. Also, the constant \mathbf{C} corresponds to $\sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} (1 - \alpha_k) y_k^{\alpha_k}$, which only depends on $\bar{\boldsymbol{\mu}}$ and $\bar{\boldsymbol{\lambda}}$.

- (d) This follows from the definition of $e^{W_k^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}})/\bar{\lambda}_j}$.
- (e) This equality follows from the definition of $\psi_{j \rightarrow k}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}})$.
- (f) The inequality (f) follows from an application of the inductive hypothesis to the children of j , and these children have height of at most H , and from collecting terms.
- (g) This follows from defining $\mathbf{C}_j = \mathbf{C} + \sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} \mathbf{C}_k$ and use the fact that $\psi_{k \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) = 0$ for all $\ell \notin \mathcal{S}^q \cap \mathbb{T}_k^q$.
- (h) The equality (h) follows because $\psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) = \psi_{j \rightarrow k}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \times \psi_{k \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}})$.
- (i) The equality (i) follows by noting that for all $\ell \notin \mathcal{S}^q \cap \mathbb{T}_j^q$, $\psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) = 0$.
- (j) Finally, the last equality (j) again follows because $\bar{\boldsymbol{\lambda}} = \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})$.

This completes the induction step. Therefore, the result holds for all nodes j . ■

We are now ready to prove the theorem. For each $\ell \in \mathcal{N}$, define $G_\ell(\cdot | \bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}}) : \mathbb{R} \rightarrow \mathbb{R}$ as

$$G_\ell(x | \bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}}) = e^{(x - \bar{\mu}_\ell) \times \exp(\bar{\delta}_{[\text{pa}(\ell)])} \cdot a_\ell(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}}) - (x - \bar{\mu}_\ell) \cdot \exp(\bar{\delta}_{[\text{pa}(\ell)])} \cdot \text{sales}_\ell \quad (\text{EC.18})$$

Then, letting $\bar{\boldsymbol{\lambda}} = \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})$, it follows from Theorem 2.2 that

$$\text{NegLog}(\boldsymbol{\mu}, \bar{\boldsymbol{\delta}}) = \frac{1}{Q} \sum_{q=1}^Q \left\{ W_{\text{root}}^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) + \sum_{j \in \text{path}(\text{root}, c^q)} W_j^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) \cdot (1/\bar{\lambda}_j - 1/\bar{\lambda}_{\text{pa}(j)}) - \mu_{c^q} / \bar{\lambda}_{\text{pa}(c^q)} \right\}.$$

For each nonleaf node $j \in \mathbb{T} \setminus \mathcal{N}$, let $\bar{\zeta}_j = (1/\bar{\lambda}_j) - (1/\bar{\lambda}_{\text{pa}(j)})$ and we set $\bar{\zeta}_{\text{root}} = 1$. Note that $\bar{\zeta}_j \geq 0$ for all j because $0 < \bar{\lambda}_j \leq \bar{\lambda}_{\text{pa}(j)}$. Then, we have

$$\begin{aligned} \text{NegLog}(\boldsymbol{\mu}, \bar{\boldsymbol{\delta}}) &= \frac{1}{Q} \sum_{q=1}^Q \left\{ W_{\text{root}}^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) + \sum_{j \in \text{path}(\text{root}, c^q)} \bar{\zeta}_j W_j^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) \right\} - \frac{1}{Q} \sum_{q=1}^Q \mu_{c^q} / \bar{\lambda}_{\text{pa}(c^q)} \\ &\stackrel{(a)}{=} \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbb{1}_{\{c^q \in \mathbb{T}_j\}} \bar{\zeta}_j W_j^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) - \frac{1}{Q} \sum_{\ell \in \mathcal{N}} (\mu_\ell / \bar{\lambda}_{\text{pa}(\ell)}) \cdot \left(\sum_{q=1}^Q \mathbb{1}_{\{c^q = \ell\}} \right) \\ &\stackrel{(b)}{=} \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbb{1}_{\{c^q \in \mathbb{T}_j\}} \bar{\zeta}_j W_j^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) - \sum_{\ell \in \mathcal{N}} (\mu_\ell - \bar{\mu}_\ell) \cdot \text{sales}_\ell / \bar{\lambda}_{\text{pa}(\ell)} - \sum_{\ell \in \mathcal{N}} (\bar{\mu}_\ell / \bar{\lambda}_{\text{pa}(\ell)}) \cdot \text{sales}_\ell, \end{aligned}$$

where the equality (a) follows because, for each $j \in \mathbb{T} \setminus \mathcal{N}$, $c^q \in \mathbb{T}_j$ if and only if $j \in \text{path}[\text{root}, c^q]$ and equality (b) follows from the definition of sales_ℓ .

Let $H(\boldsymbol{\mu}) = \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \bar{\zeta}_j W_j^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}})$. We will apply Lemma E.5 to construct a majorizing surrogate function for $H(\boldsymbol{\mu})$. Let \mathbf{C} denote a constant that depend only on $(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}})$. Then,

$$\begin{aligned}
H(\boldsymbol{\mu}) &\stackrel{(c)}{=} \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \bar{\zeta}_j \bar{\lambda}_j \log \left(\sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} e^{W_k^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_j} \right) \\
&\stackrel{(d)}{\leq} \mathbf{C} + \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \bar{\zeta}_j \bar{\lambda}_j \frac{\sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} e^{W_k^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_j}}{\sum_{k \in \text{Children}(j) \cap \mathbb{T}^q} e^{W_k^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_j}} \\
&\stackrel{(e)}{\leq} \mathbf{C} + \frac{1}{Q} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \bar{\zeta}_j \bar{\lambda}_j \times \frac{1}{\bar{\lambda}_j} \sum_{\ell \in \mathcal{N}} \bar{\lambda}_{\text{pa}(\ell)} \times \psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \times e^{(\mu_\ell - \bar{\mu}_\ell) / \bar{\lambda}_{\text{pa}(\ell)}} \\
&\stackrel{(f)}{=} \mathbf{C} + \frac{1}{Q} \sum_{\ell \in \mathcal{N}} e^{(\mu_\ell - \bar{\mu}_\ell) / \bar{\lambda}_{\text{pa}(\ell)}} \times \left[\bar{\lambda}_{\text{pa}(\ell)} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} \bar{\zeta}_j \psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \right] \\
&\stackrel{(g)}{=} \mathbf{C} + \frac{1}{Q} \sum_{\ell \in \mathcal{N}} e^{(\mu_\ell - \bar{\mu}_\ell) / \bar{\lambda}_{\text{pa}(\ell)}} \times \left[\bar{\lambda}_{\text{pa}(\ell)} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} (1/\bar{\lambda}_j - 1/\bar{\lambda}_{\text{pa}(j)}) \psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \right]
\end{aligned}$$

where the justifications for the equalities and inequalities in (c) - (g) are given below.

- (c) The equality follows from the definition of $W_j^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}})$.
- (d) The inequality follows from applying the subgradient inequality for logarithm, which shows that for all $x \in \mathbb{R}_+$ and $y \in \mathbb{R}_+$, $\log(x) \leq \log(y) + \frac{1}{y}(x - y) = \frac{x}{y} + \log(y) - 1$, with equality if and only if $x = y$. Here, for each $q = 1, \dots, Q$ and $j \in \mathbb{T} \setminus \mathcal{N}$,

$$x_j^q = \sum_{k \in \text{Children}(j)} e^{W_k^q(\boldsymbol{\mu}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_j} \quad \text{and} \quad y_j^q = \sum_{k \in \text{Children}(j)} e^{W_k^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) / \bar{\lambda}_j}$$

and the constant $\mathbf{C} = \sum_{q=1}^Q \sum_{j \in \mathbb{T}} (\log(y_j^q) - 1)$, which only depends on $(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}})$.

- (e) The inequality follows from Lemma E.5, and the constant \mathbf{C} is updated accordingly.
- (f) The equality follows from algebra and rearranging terms.
- (g) The equality follows from the definition of $\bar{\zeta}_j$ and the fact that $\bar{\lambda}_{\text{pa}(\text{root})} = +\infty$.

Putting everything together, we have the following majorizing surrogate for $\boldsymbol{\mu} \mapsto \text{NegLog}(\boldsymbol{\mu}, \bar{\boldsymbol{\delta}})$:

$$\begin{aligned}
&\mathbf{C} + \frac{1}{Q} \sum_{\ell \in \mathcal{N}} e^{(\mu_\ell - \bar{\mu}_\ell) / \bar{\lambda}_{\text{pa}(\ell)}} \times \left[\bar{\lambda}_{\text{pa}(\ell)} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} (1/\bar{\lambda}_j - 1/\bar{\lambda}_{\text{pa}(j)}) \psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \right] \\
&\quad - \sum_{\ell \in \mathcal{N}} (\mu_\ell - \bar{\mu}_\ell) \cdot \text{sales}_\ell / \bar{\lambda}_{\text{pa}(\ell)} \\
&= \mathbf{C} + \frac{1}{Q} \sum_{\ell \in \mathcal{N}} e^{(\mu_\ell - \bar{\mu}_\ell) \times \exp(\bar{\delta}_{[\text{pa}(\ell)])}} \times \left[\exp(-\bar{\delta}_{[\text{pa}(\ell)])} \sum_{q=1}^Q \sum_{j \in \mathbb{T} \setminus \mathcal{N}} \mathbf{1}_{\{c^q \in \mathbb{T}_j\}} (\exp(\bar{\delta}_{[j]}) - \exp(\bar{\delta}_{[\text{pa}(j)]})) \psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})) \right] \\
&\quad - \sum_{\ell \in \mathcal{N}} (\mu_\ell - \bar{\mu}_\ell) \times \exp(\bar{\delta}_{[\text{pa}(\ell)])} \times \text{sales}_\ell
\end{aligned}$$

$$\begin{aligned}
&= \mathbf{C} + \sum_{\ell \in \mathcal{N}} e^{(\mu_\ell - \bar{\mu}_\ell) \times \exp(\bar{\delta}_{[\text{pa}(\ell)]})} \times \exp(-\bar{\delta}_{[\text{pa}(\ell)]}) \times \left[\frac{\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\delta})}{\partial \mu_\ell} \Big|_{\boldsymbol{\mu}=\bar{\boldsymbol{\mu}}, \boldsymbol{\delta}=\bar{\boldsymbol{\delta}}} + \text{sales}_\ell \cdot \exp(\bar{\delta}_{[\text{pa}(\ell)]}) \right] \\
&\quad - \sum_{\ell \in \mathcal{N}} (\mu_\ell - \bar{\mu}_\ell) \times \exp(\bar{\delta}_{[\text{pa}(\ell)]}) \times \text{sales}_\ell \\
&= \mathbf{C} + \sum_{\ell \in \mathcal{N}} e^{(\mu_\ell - \bar{\mu}_\ell) \times \exp(\bar{\delta}_{[\text{pa}(\ell)]})} \cdot a_\ell(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}}) - \sum_{\ell \in \mathcal{N}} (\mu_\ell - \bar{\mu}_\ell) \times \exp(\bar{\delta}_{[\text{pa}(\ell)]}) \times \text{sales}_\ell \\
&= \mathbf{C} + \sum_{\ell \in \mathcal{N}} G_\ell(\mu_\ell | \bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}})
\end{aligned}$$

where the first equality follows since $\bar{\boldsymbol{\lambda}} = \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})$, the second follows from Lemma D.2, the third from the definition of $a_\ell(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}})$, and the final from the definition of $G_\ell(\mu_\ell | \bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}})$ in (EC.18).

To show that $G_\ell(\cdot | \bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}})$ is strictly convex, it suffices to show that $a_\ell(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}}) > 0$. Using the expression for $\frac{\partial \text{NegLog}(\boldsymbol{\mu}, \boldsymbol{\lambda}(\boldsymbol{\delta}))}{\partial \mu_\ell}$ from Lemma D.2, it follows that

$$\begin{aligned}
a_\ell(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}}) &= \frac{\exp(-\bar{\delta}_{[\text{pa}(\ell)]})}{Q} \sum_{q=1}^Q \sum_{j \in \mathcal{T} \setminus \mathcal{N}} \mathbb{1}_{\{c^q \in \mathcal{T}_j\}} (\exp(\bar{\delta}_{[j]}) - \exp(\bar{\delta}_{[\text{pa}(j)]})) \psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \boldsymbol{\lambda}(\bar{\boldsymbol{\delta}})) \\
&= \frac{\bar{\lambda}_{\text{pa}(\ell)}}{Q} \sum_{q=1}^Q \sum_{j \in \mathcal{T} \setminus \mathcal{N}} \mathbb{1}_{\{c^q \in \mathcal{T}_j\}} (1/\bar{\lambda}_j - 1/\bar{\lambda}_{\text{pa}(j)}) \psi_{j \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \\
&\geq \frac{\bar{\lambda}_{\text{pa}(\ell)}}{Q} \sum_{q=1}^Q \mathbb{1}_{\{c^q \in \mathcal{T}_{\text{root}}\}} (1/\bar{\lambda}_{\text{root}} - 1/\bar{\lambda}_{\text{pa}(\text{root})}) \psi_{\text{root} \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \\
&= \frac{\bar{\lambda}_{\text{pa}(\ell)}}{Q} \sum_{q=1}^Q \psi_{\text{root} \rightarrow \ell}^q(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) \\
&> 0,
\end{aligned}$$

where the first inequality follows since each term inside the summation is non-negative, the last equality follows because $\mathbb{1}_{\{c^q \in \mathcal{T}_{\text{root}}\}} = 1$, $\bar{\lambda}_{\text{root}} = 1$ and $\bar{\lambda}_{\text{pa}(\text{root})} = +\infty$, and the final inequality follows because $\bar{\lambda}_{\text{pa}(\ell)} > 0$ and the identifiability condition in Theorem 2.4 ensures that $\ell \in \mathcal{S}^{\hat{q}}$ for some \hat{q} so that $\psi_{\text{root} \rightarrow \ell}^{\hat{q}}(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\lambda}}) > 0$.

Finally, the expression for the minimizer follows from setting $\partial G_\ell(\mu_\ell | \bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}}) / \partial \mu_\ell = 0$.

E.8 Proof of Lemma 3.8

Recall from Theorem 3.1 that $\text{NegLog}(\bar{\boldsymbol{\mu}}, \boldsymbol{\delta}) = F_1(\bar{\boldsymbol{\mu}}, \boldsymbol{\delta}) - F_2(\bar{\boldsymbol{\mu}}, \boldsymbol{\delta})$ for all $\boldsymbol{\delta} \in \text{Dom}_2$, where both $F_1(\bar{\boldsymbol{\mu}}, \cdot)$ and $F_2(\bar{\boldsymbol{\mu}}, \cdot)$ are strictly convex on Dom_2 . By applying the sub-gradient inequality to the strictly convex function $F_2(\bar{\boldsymbol{\mu}}, \boldsymbol{\delta})$ at $\boldsymbol{\delta} = \bar{\boldsymbol{\delta}}$, it follows that

$$\begin{aligned}
\text{NegLog}(\bar{\boldsymbol{\mu}}, \boldsymbol{\delta}) &\leq F_1(\bar{\boldsymbol{\mu}}, \boldsymbol{\delta}) - F_2(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}}) - \sum_{j \in \mathcal{T} \setminus (\mathcal{N} \cup \{\text{root}\})} \frac{\partial F_2}{\partial \delta_j}(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\delta}}) \cdot (\delta_j - \bar{\delta}_j) \quad \forall \boldsymbol{\delta} \in \text{Dom}_2 \\
&\quad \text{with equality if and only if } \boldsymbol{\delta} = \bar{\boldsymbol{\delta}}. \quad (\text{EC.19})
\end{aligned}$$

The result then follows from the definition of a majorizing surrogate. ■

E.9 Improvements via MM and PGD updates for δ

As mentioned in the main text, the improvement guarantee for Step 2 of the A-MM algorithm follows from existing results. In particular, it follows from the standard guarantee for the projected gradient descent (PGD) algorithm for constrained smooth problems, which we reproduce here for completeness. Consider the optimization problem

$$\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}), \quad (\text{EC.20})$$

where f is L -smooth on \mathcal{X} and $\mathcal{X} \subset \mathbb{R}^d$ is a convex set. We define the *gradient mapping* (Nesterov 2013, Section 2.2.3) of f over \mathcal{X} , denoted by $\mathcal{G}_{f,\mathcal{X}}: \mathcal{X} \rightarrow \mathbb{R}^d$ as

$$\mathcal{G}_{f,\mathcal{X}}(\mathbf{x}) = L \cdot \left(\mathbf{x} - \text{Proj}_{\mathcal{X}} \left(\mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x}) \right) \right), \quad (\text{EC.21})$$

where $\text{Proj}_{\mathcal{X}}(\cdot)$ is the projection operator onto \mathcal{X} , i.e. $\text{Proj}_{\mathcal{X}}(\bar{\mathbf{x}}) = \arg \min_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x} - \bar{\mathbf{x}}\|_2^2$ for all $\bar{\mathbf{x}} \in \mathbb{R}^d$. The gradient mapping plays the role of the gradient in unconstrained optimization problems since it can be shown that \mathbf{x}^* is a stationary point for problem (EC.20) if and only if $\|\mathcal{G}_{f,\mathcal{X}}(\mathbf{x}^*)\|_2 = 0$.

Now, starting from an initial solution $\mathbf{x}^{(0)} \in \mathcal{X}$, consider solving problem (EC.20) using PGD with step size $1/L$, that is, for each $s \geq 0$:

$$\mathbf{x}^{(s+1)} = \text{Proj}_{\mathcal{X}} \left(\mathbf{x}^{(s)} - \frac{1}{L} \nabla f(\mathbf{x}^{(s)}) \right)$$

Then, using the quadratic upper bound for smooth functions, it follows that

$$f(\mathbf{x}^{(s+1)}) \leq f(\mathbf{x}^{(s)}) + \langle \nabla f(\mathbf{x}^{(s)}), \mathbf{x}^{(s+1)} - \mathbf{x}^{(s)} \rangle + \frac{L}{2} \|\mathbf{x}^{(s+1)} - \mathbf{x}^{(s)}\|_2^2 \quad (\text{EC.22})$$

Next, for any $\mathbf{x} \in \mathcal{X}$, denote $D(\mathbf{x}) = \|\mathbf{x} - (\mathbf{x}^{(s)} - \frac{1}{L} \nabla f(\mathbf{x}^{(s)}))\|_2^2$ and note that $\mathbf{x}^{(s+1)} = \arg \min_{\mathbf{x} \in \mathcal{X}} D(\mathbf{x})$. The optimality of $\mathbf{x}^{(s+1)}$ implies that there is no descent direction from $\mathbf{x}^{(s+1)}$:

$$\langle \nabla D(\mathbf{x}^{(s+1)}), \mathbf{x} - \mathbf{x}^{(s+1)} \rangle \geq 0 \quad \forall \mathbf{x} \in \mathcal{X}$$

By plugging in the gradient of $D(\cdot)$ and choosing $\mathbf{x} = \mathbf{x}^{(s)}$, it follows that

$$\langle \nabla f(\mathbf{x}^{(s)}), \mathbf{x}^{(s+1)} - \mathbf{x}^{(s)} \rangle \leq -L \|\mathbf{x}^{(s+1)} - \mathbf{x}^{(s)}\|_2^2$$

Substituting the above in (EC.22), it follows that

$$\begin{aligned} f(\mathbf{x}^{(s+1)}) &\leq f(\mathbf{x}^{(s)}) - L \|\mathbf{x}^{(s+1)} - \mathbf{x}^{(s)}\|_2^2 + \frac{L}{2} \|\mathbf{x}^{(s+1)} - \mathbf{x}^{(s)}\|_2^2 \\ &= f(\mathbf{x}^{(s)}) - \frac{L}{2} \|\mathbf{x}^{(s+1)} - \mathbf{x}^{(s)}\|_2^2 \\ &= f(\mathbf{x}^{(s)}) - \frac{1}{2L} \|\mathcal{G}_{f,\mathcal{X}}(\mathbf{x}^{(s)})\|_2^2, \end{aligned} \quad (\text{EC.23})$$

where the final equality follows from the definition of $\mathcal{G}_{f,\mathcal{X}}(\mathbf{x}^{(s)})$.

We adapt the above improvement guarantee to our context. Define the domain $\overline{\text{Dom}}_2$ as:

$$\overline{\text{Dom}}_2 = \left\{ \boldsymbol{\delta} \in \mathbb{R}^{|\mathcal{T} \setminus (\mathcal{N} \cup \{\text{root}\})|} : 0 \leq \delta_j \leq \delta_{\text{upper}} \quad \forall j \in \mathcal{T} \setminus (\mathcal{N} \cup \{\text{root}\}) \right\} \quad (\text{EC.24})$$

Then, we can establish the following:

Lemma E.6 (Improvements via MM and PGD updates for $\boldsymbol{\delta}$) *Suppose $\boldsymbol{\delta}^{(s)} \in \overline{\text{Dom}}_2$ and let $\gamma, \gamma_{\text{GD}} > 0$ denote the smoothness constants of the mappings $\boldsymbol{\delta} \mapsto H(\boldsymbol{\delta} | \boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)})$ and $\boldsymbol{\delta} \mapsto \text{NegLog}(\boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta})$, respectively, over the domain $\overline{\text{Dom}}_2$. If $\boldsymbol{\delta}^{(s+1)} \in \overline{\text{Dom}}_2$, then the MM update guarantees the following improvement in the negative log-likelihood objective:*

$$\text{NegLog}(\boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s+1)}) - \text{NegLog}(\boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)}) \leq -\frac{1}{2\gamma} \left\| \mathcal{G}_{H, \overline{\text{Dom}}_2}(\boldsymbol{\delta}^{(s)}) \right\|_2^2$$

Similarly, let $\boldsymbol{\delta}_{\text{GD}}^{(s+1)}$ denote the PGD update with step size $1/\gamma_{\text{GD}}$. If $\boldsymbol{\delta}_{\text{GD}}^{(s+1)} \in \overline{\text{Dom}}_2$, then the PGD update guarantees the following improvement in the negative log-likelihood objective:

$$\text{NegLog}(\boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}_{\text{GD}}^{(s+1)}) - \text{NegLog}(\boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)}) \leq -\frac{1}{2\gamma_{\text{GD}}} \left\| \mathcal{G}_{\text{NegLog}, \overline{\text{Dom}}_2}(\boldsymbol{\delta}^{(s)}) \right\|_2^2$$

Proof. We start with the MM update. First, it can be shown—using an argument analogous to that in the proof of Theorem 3.3 above—that the eigenvalues of the Hessian matrix $\nabla^2 H(\boldsymbol{\delta} | \boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)})$ are bounded above by a constant $\gamma > 0$ that depends on δ_{upper} , for all $\boldsymbol{\delta} \in \overline{\text{Dom}}_2$. This ensures that $H(\boldsymbol{\delta} | \boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)})$ is γ -smooth on $\overline{\text{Dom}}_2$. Now, define $\bar{\boldsymbol{\delta}}^{(s+1)} \in \overline{\text{Dom}}_2$ as follows:

$$\bar{\boldsymbol{\delta}}^{(s+1)} = \text{Proj}_{\overline{\text{Dom}}_2} \left(\boldsymbol{\delta}^{(s)} - \frac{1}{\gamma} \nabla H(\boldsymbol{\delta}^{(s)} | \boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)}) \right),$$

Then, it follows from (EC.23) that

$$\begin{aligned} H(\bar{\boldsymbol{\delta}}^{(s+1)} | \boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)}) &\leq H(\boldsymbol{\delta}^{(s)} | \boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)}) - \frac{1}{2\gamma} \left\| \mathcal{G}_{H, \overline{\text{Dom}}_2}(\boldsymbol{\delta}^{(s)}) \right\|_2^2 \\ &= \text{NegLog}(\boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)}) - \frac{1}{2\gamma} \left\| \mathcal{G}_{H, \overline{\text{Dom}}_2}(\boldsymbol{\delta}^{(s)}) \right\|_2^2, \end{aligned}$$

where the equality follows since $H(\boldsymbol{\delta} | \boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)})$ is a majorizing surrogate. The improvement bound then follows since

$$\text{NegLog}(\boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s+1)}) \leq H(\boldsymbol{\delta}^{(s+1)} | \boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)}) \leq H(\bar{\boldsymbol{\delta}}^{(s+1)} | \boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)}),$$

where the second inequality follows since $\nabla H(\boldsymbol{\delta}^{(s)} | \boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)}) = \nabla_{\boldsymbol{\delta}} \text{NegLog}(\boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s)})$ and the fact that the MM update $\boldsymbol{\delta}^{(s+1)} \in \overline{\text{Dom}}_2$ corresponds to the optimal step size.

The improvement bound for the GD update $\boldsymbol{\delta}_{\text{GD}}^{(s+1)}$ follows in an identical fashion by leveraging the smoothness of $\text{NegLog}(\boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta})$. In particular, we can show that $\text{NegLog}(\boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta})$ is γ_{GD} -smooth on $\overline{\text{Dom}}_2$, for some constant $\gamma_{\text{GD}} > 0$ that depends on δ_{upper} . The result then follows from (EC.23). \blacksquare

E.10 Sublinear convergence rate of the A-MM algorithm

As stated in the main text, the improvements bounds from Theorem 3.3 and Lemma E.6 can be combined to establish a sublinear rate of convergence of the A-MM algorithm, as shown in the following theorem, where we leverage the notations introduced above:

Theorem E.7 (Sublinear convergence of A-MM to a stationary point) *Suppose that $\boldsymbol{\delta}^{(s)} \in \overline{\text{Dom}}_2$ for all $s \geq 0$. Then, the sequence of iterates $\left((\boldsymbol{\mu}^{(s)}, \boldsymbol{\delta}^{(s)}) : s \geq 0 \right)$ converges to a stationary point of the MLE problem. Moreover, we can establish the following guarantee*

$$\min_{0 \leq s' \leq s} \left[\left\| \nabla_{\boldsymbol{\mu}} \text{NegLog} \left(\boldsymbol{\mu}^{(s')}, \boldsymbol{\delta}^{(s')} \right) \right\|_1^2 + \left\| \mathcal{G}_{H, \overline{\text{Dom}}_2}(\boldsymbol{\delta}^{(s')}) \right\|_2^2 \right] \leq \frac{2 \cdot \left(\text{NegLog}(\boldsymbol{\mu}^{(0)}, \boldsymbol{\delta}^{(0)}) - \text{NegLog}^* \right)}{R \cdot (s+1)},$$

where NegLog^* is the optimal objective for the MLE problem, and $R = \min(\lambda_{\text{lower}}^2, 1/\gamma)$.

Proof. The first part of the theorem follows from the fact that the A-MM algorithm guarantees an improving solution as long as the current solution is not a stationary point. For the second part, we leverage the improvements from Theorem 3.3 and Lemma E.6 to obtain, for each $s' \geq 0$:

$$\begin{aligned} & \text{NegLog} \left(\boldsymbol{\mu}^{(s'+1)}, \boldsymbol{\delta}^{(s'+1)} \right) - \text{NegLog} \left(\boldsymbol{\mu}^{(s')}, \boldsymbol{\delta}^{(s')} \right) \\ &= \left[\text{NegLog} \left(\boldsymbol{\mu}^{(s'+1)}, \boldsymbol{\delta}^{(s'+1)} \right) - \text{NegLog} \left(\boldsymbol{\mu}^{(s'+1)}, \boldsymbol{\delta}^{(s')} \right) \right] + \left[\text{NegLog} \left(\boldsymbol{\mu}^{(s'+1)}, \boldsymbol{\delta}^{(s')} \right) - \text{NegLog} \left(\boldsymbol{\mu}^{(s')}, \boldsymbol{\delta}^{(s')} \right) \right] \\ &\leq -\frac{1}{2\gamma} \left\| \mathcal{G}_{H, \overline{\text{Dom}}_2}(\boldsymbol{\delta}^{(s')}) \right\|_2^2 - \frac{\lambda_{\text{lower}}^2}{2} \left\| \nabla_{\boldsymbol{\mu}} \text{NegLog}(\boldsymbol{\mu}^{(s')}, \boldsymbol{\delta}^{(s')}) \right\|_1^2 \\ &\leq -\frac{R}{2} \left[\left\| \mathcal{G}_{H, \overline{\text{Dom}}_2}(\boldsymbol{\delta}^{(s')}) \right\|_2^2 + \left\| \nabla_{\boldsymbol{\mu}} \text{NegLog}(\boldsymbol{\mu}^{(s')}, \boldsymbol{\delta}^{(s')}) \right\|_1^2 \right], \end{aligned}$$

where the second inequality follows from the definition of R . Summing the above from 0 to any iteration $s \geq 0$, we get

$$\sum_{s'=0}^s \left[\text{NegLog} \left(\boldsymbol{\mu}^{(s'+1)}, \boldsymbol{\delta}^{(s'+1)} \right) - \text{NegLog} \left(\boldsymbol{\mu}^{(s')}, \boldsymbol{\delta}^{(s')} \right) \right] \leq -\frac{R}{2} \sum_{s'=0}^s \left[\left\| \mathcal{G}_{H, \overline{\text{Dom}}_2}(\boldsymbol{\delta}^{(s')}) \right\|_2^2 + \left\| \nabla_{\boldsymbol{\mu}} \text{NegLog}(\boldsymbol{\mu}^{(s')}, \boldsymbol{\delta}^{(s')}) \right\|_1^2 \right]$$

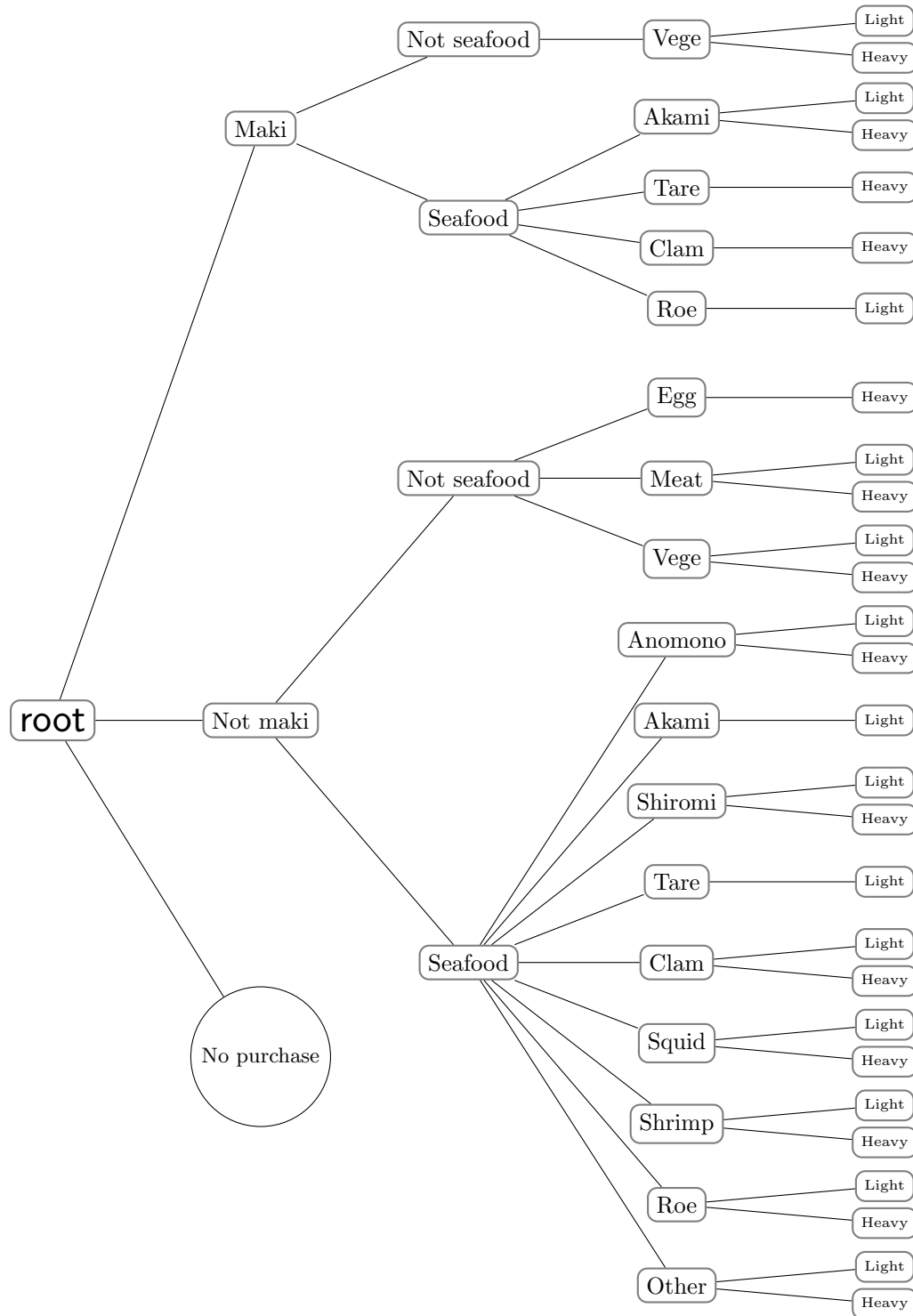
The left hand side of the inequality above is a telescoping sum, and it is equal to $\text{NegLog} \left(\boldsymbol{\mu}^{(s+1)}, \boldsymbol{\delta}^{(s+1)} \right) - \text{NegLog} \left(\boldsymbol{\mu}^{(0)}, \boldsymbol{\delta}^{(0)} \right)$, which is bounded below by $\text{NegLog}^* - \text{NegLog} \left(\boldsymbol{\mu}^{(0)}, \boldsymbol{\delta}^{(0)} \right)$. We thus have

$$\text{NegLog}^* - \text{NegLog} \left(\boldsymbol{\mu}^{(0)}, \boldsymbol{\delta}^{(0)} \right) \leq -\frac{R}{2} \sum_{s'=0}^s \left[\left\| \mathcal{G}_{H, \overline{\text{Dom}}_2}(\boldsymbol{\delta}^{(s')}) \right\|_2^2 + \left\| \nabla_{\boldsymbol{\mu}} \text{NegLog}(\boldsymbol{\mu}^{(s')}, \boldsymbol{\delta}^{(s')}) \right\|_1^2 \right] \quad (\text{EC.25})$$

Finally, using (EC.25) it follows that

$$\begin{aligned} & \min_{0 \leq s' \leq s} \left[\left\| \nabla_{\boldsymbol{\mu}} \text{NegLog} \left(\boldsymbol{\mu}^{(s')}, \boldsymbol{\delta}^{(s')} \right) \right\|_1^2 + \left\| \mathcal{G}_{H, \overline{\text{Dom}}_2}(\boldsymbol{\delta}^{(s')}) \right\|_2^2 \right] \\ &\leq \frac{1}{s+1} \sum_{s'=0}^s \left[\left\| \nabla_{\boldsymbol{\mu}} \text{NegLog} \left(\boldsymbol{\mu}^{(s')}, \boldsymbol{\delta}^{(s')} \right) \right\|_1^2 + \left\| \mathcal{G}_{H, \overline{\text{Dom}}_2}(\boldsymbol{\delta}^{(s')}) \right\|_2^2 \right] \\ &\leq \frac{2}{R} \cdot \frac{1}{s+1} \left[\text{NegLog} \left(\boldsymbol{\mu}^{(0)}, \boldsymbol{\delta}^{(0)} \right) - \text{NegLog}^* \right]. \end{aligned}$$

■

Appendix F: Tree structure used in SUSHI Dataset study

The figure above depicts the tree structure employed in our real-world study on the SUSHI Preference Dataset; since there are 100 products (= the number of sushi varieties), we only show the non-leaf nodes. At the root node, the customer either decides to purchase a sushi variety or leave without a purchase. If the customer decides to make a purchase, she first chooses the *style*

of the sushi type: maki or not maki. Then, for each style, she decides to purchase a sushi variety based on whether or not it contains *seafood*. If the sushi contains seafood, she further chooses from nine minor groups: aomono (blue-skinned fish), akami (red meat fish), shiromi (white-meat fish), tare (something like baste; for eel or sea eel), clam or shell, squid or octopus, shrimp or crab, roe, and other seafood. Otherwise, she chooses from three minor groups: egg, meat other than fish, and vegetables. Lastly, she determines whether to purchase a sushi type with heavy/oily or light taste.

Appendix G: Additional Performance Measure

Here, we compare the PGD, Knitro and A-MM benchmarks on the root mean square error (RMSE) metric. For the simulation study, let $(\mathcal{S}_i^m : i = 1, 2, \dots, 700)$ denote the sampled offer-sets for instance m . Then, we compute the average RMSE value across all the 100 instances for each $\text{algo} \in \{\text{A-MM}, \text{Knitro}, \text{PGD}\}$ as follows:

$$\text{RMSE}^{\text{algo}} = \frac{1}{100} \sum_{m=1}^{100} \sqrt{\frac{1}{60} \sum_{i=1}^{60} \frac{1}{|\mathcal{S}_i^m|} \sum_{\ell \in \mathcal{S}_i^m} (\mathbb{P}_\ell(\mathcal{S}_i^m; \boldsymbol{\mu}^{\text{true},m}, \boldsymbol{\lambda}^{\text{true},m}) - \mathbb{P}_\ell(\mathcal{S}_i^m; \boldsymbol{\mu}^{\text{algo},m}, \boldsymbol{\lambda}^{\text{algo},m}))^2},$$

where for each instance m , $\mathbb{P}_\ell(\mathcal{S}_i^m; \boldsymbol{\mu}^{\text{true},m}, \boldsymbol{\lambda}^{\text{true},m})$ is the ground-truth purchase probability, and $\mathbb{P}_\ell(\mathcal{S}_i^m; \boldsymbol{\mu}^{\text{algo},m}, \boldsymbol{\lambda}^{\text{algo},m})$ is the estimated purchase probability by method algo, of product ℓ in offer-set \mathcal{S}_i^m . We report the RMSE values as well as the percentage improvements $100 \times (\text{RMSE}^{\text{algo}} - \text{RMSE}^{\text{A-MM}}) / \text{RMSE}^{\text{algo}}$ for $\text{algo} \in \{\text{PGD}, \text{Knitro}\}$ in Figure EC.1. Similar to the `NegLogGap` values in the main text, the A-MM method achieves significantly lower RMSE than the benchmarks for all ground-truth problem sizes except the two smallest ones.

For the Sushi dataset study, let $(\mathcal{S}_i^m : i = 1, 2, \dots, 700)$ denote the training offer-sets for instance m , and $(\mathcal{S}_i^m : i = 701, 702, \dots, 1000)$ denote the test offer-sets. Then, we compute the average RMSE value on the training and test offer-sets for each method $\text{algo} \in \{\text{A-MM}, \text{PGD}, \text{Knitro}\}$ as follows:

$$\begin{aligned} \text{RMSE}_{\text{Train}}^{\text{algo}} &= \frac{1}{400} \sum_{m=1}^{400} \sqrt{\frac{1}{700} \sum_{i=1}^{700} \frac{1}{|\mathcal{S}_i^m|} \sum_{\ell \in \mathcal{S}_i^m} (\text{sales}_\ell(\mathcal{S}_i^m) - \mathbb{P}_\ell(\mathcal{S}_i^m; \boldsymbol{\mu}^{\text{algo},m}, \boldsymbol{\lambda}^{\text{algo},m}))^2} \\ \text{RMSE}_{\text{Test}}^{\text{algo}} &= \frac{1}{400} \sum_{m=1}^{400} \sqrt{\frac{1}{300} \sum_{i=701}^{1000} \frac{1}{|\mathcal{S}_i^m|} \sum_{\ell \in \mathcal{S}_i^m} (\text{sales}_\ell(\mathcal{S}_i^m) - \mathbb{P}_\ell(\mathcal{S}_i^m; \boldsymbol{\mu}^{\text{algo},m}, \boldsymbol{\lambda}^{\text{algo},m}))^2} \end{aligned}$$

where $\text{sales}_\ell(\mathcal{S}_i^m)$ is the fraction of observed sales for product ℓ in offer set \mathcal{S}_i^m , and $\mathbb{P}_\ell(\mathcal{S}_i^m; \boldsymbol{\mu}^{\text{algo},m}, \boldsymbol{\lambda}^{\text{algo},m})$ is the probability that a customer purchases product ℓ from the offer set \mathcal{S}_i^m under the parameters estimated by algo in instance m . We report the RMSE numbers as well as the percentage improvements $100 \times (\text{RMSE}_{\text{Train}}^{\text{algo}} - \text{RMSE}_{\text{Train}}^{\text{A-MM}}) / \text{RMSE}_{\text{Train}}^{\text{algo}}$ and $100 \times (\text{RMSE}_{\text{Test}}^{\text{algo}} - \text{RMSE}_{\text{Test}}^{\text{A-MM}}) / \text{RMSE}_{\text{Test}}^{\text{algo}}$ for $\text{algo} \in \{\text{PGD}, \text{Knitro}\}$ in Figure EC.2. Similar to the `NegLogImp` values in the main text, the A-MM method achieves lower RMSE values than the benchmarks on both training and test data, and is robust to initialization.

Degree	Height	# Prods.	# Nodes	λ_{lower}	RMSE ($\times 10^{-5}$)			RMSE Impr.		% better		
					A-MM	PGD	Knitro	over PGD	over Knitro	over PGD	over Knitro	
5	4	625	781	0.50	5.02	6.52	7.21	23.0%	30.4%	72	91	
				0.10	8.60	7.73	9.32	-11.2%	7.8%	40	74	
				0.01	9.75	9.49	9.89	-2.8%	1.4%	50	68	
	5	3,125	3,906	0.50	1.18	2.40	2.12	50.6%	44.1%	100	99	
				0.10	2.15	2.87	3.45	25.1%	37.5%	95	100	
				0.01	2.68	3.53	3.95	24.1%	32.2%	100	100	
	6	4	1,296	1,555	0.50	2.38	5.54	5.89	56.9%	59.5%	100	100
					0.10	4.32	6.45	6.95	33.0%	37.8%	90	100
					0.01	5.04	8.15	7.45	38.2%	32.4%	100	100
5		7,776	9,331	0.50	0.46	1.57	3.19	70.8%	85.6%	100	100	
				0.10	0.87	1.84	6.01	52.5%	85.5%	100	100	
				0.01	1.09	2.32	6.57	53.2%	83.5%	100	100	
7		4	2,401	2,801	0.50	1.26	4.42	2.99	71.6%	58.0%	100	100
					0.10	2.36	5.62	6.22	58.1%	62.1%	90	100
					0.01	2.80	6.39	6.56	56.2%	57.4%	100	100
	5	16,807	19,608	0.50	0.20	0.99	—	79.8%	—	100	—	
				0.10	0.40	1.11	—	63.8%	—	100	—	
				0.01	0.50	1.38	—	63.7%	—	100	—	
	8	4	4,096	4,681	0.50	0.85	3.59	—	76.3%	—	100	—
					0.10	1.37	4.28	—	68.0%	—	100	—
					0.01	1.63	5.07	—	67.9%	—	100	—
5		32,768	37,449	0.50	0.10	0.62	—	83.9%	—	100	—	
				0.10	0.17	0.78	—	78.2%	—	100	—	
				0.01	0.23	0.91	—	74.7%	—	100	—	

Figure EC.1 Comparison of the performances of PGD, Knitro, and our proposed A-MM method in fitting tree logit models to choice data. The columns “Degree”, “Height”, and λ_{lower} report the degree of each non-leaf node, the height of the tree, and the lower bound on the nest dissimilarity parameters, respectively. The columns “# Prods.” and “# Nodes” report the number of products and the number of nodes in the tree, respectively. The columns under A-MM, PGD and Knitro report the average RMSE for each method. The columns under “RMSE Impr.” reports the percentage improvement in the average RMSE value that our A-MM method achieves over the benchmarks. Finally, the columns under “% better” reports the percentage of instances in which the A-MM method obtains a lower RMSE value than the corresponding benchmark. The Knitro benchmark is unable to complete even a single iteration for large problem sizes, and we use “—” to denote such instances. All numbers under the “RMSE Impr.” columns are significantly different from zero at the 1% significance level under a paired samples t-test.

Initialization	Train RMSE ($\times 10^{-4}$)			Test RMSE ($\times 10^{-4}$)			Impr. over PGD		Impr. over Knitro	
	A-MM	Knitro	PGD	A-MM	Knitro	PGD	Train	Test	Train	Test
0/1 start	13.64	13.72	14.35	13.69	13.78	14.40	5.23%	5.20%	0.59%	0.61%
warm start	13.64	13.67	13.79	13.69	13.73	13.85	1.16%	1.15%	0.27%	0.27%

Figure EC.2 Comparison of the performances of PGD, Knitro and our proposed A-MM method in fitting tree logit models to the Sushi Preference Dataset. The first (resp. second) row reports the performance when the methods are initialized using *0/1 start* (resp. *warm start*). The second, third and fourth columns report the RMSE value of A-MM, Knitro and PGD on the training data, while the fifth, sixth and seventh columns report the corresponding RMSE values on the test data. The eighth and tenth columns report the percentage improvement in the average RMSE of the A-MM method over the PGD and Knitro benchmarks, respectively, on the training data. The corresponding improvements on the test data are reported in columns nine and eleven. All numbers under the “Impr.” columns are significantly different from zero at the 1% significance level under a paired sample t-test.

Appendix H: Performance of Alternating GD method

As mentioned in the main text, we evaluate the performance of our method, A-MM, against the alternating gradient descent (A-GD) method to tease apart the effect of the MM approach from effect the variable transformation procedure on the overall performance. Starting from an initial solution $(\boldsymbol{\mu}_{\text{A-GD}}^{(0)}, \boldsymbol{\delta}_{\text{A-GD}}^{(0)})$, the A-GD method performs the following two updates in each iteration $s \geq 0$:

$$\begin{aligned}
 \text{(i)} \quad & \tilde{\boldsymbol{\mu}}_{\text{A-GD}}^{(s+1)} = \boldsymbol{\mu}_{\text{A-GD}}^{(s)} - \beta^{(s)} \cdot \nabla_{\boldsymbol{\mu}} \text{NegLog}(\boldsymbol{\mu}_{\text{A-GD}}^{(s)}, \boldsymbol{\delta}_{\text{A-GD}}^{(s)}) \quad \text{for some step size } \beta^{(s)} > 0; \\
 & \left(\boldsymbol{\mu}_{\text{A-GD}}^{(s+1)} \right)_{\ell} = \left(\tilde{\boldsymbol{\mu}}_{\text{A-GD}}^{(s+1)} \right)_{\ell} - \left(\tilde{\boldsymbol{\mu}}_{\text{A-GD}}^{(s+1)} \right)_1 \quad \forall \ell \in \mathcal{N} \\
 \text{(ii)} \quad & \boldsymbol{\delta}_{\text{A-GD}}^{(s+1)} = \left\{ \boldsymbol{\delta}_{\text{A-GD}}^{(s)} - \gamma^{(s)} \cdot \nabla_{\boldsymbol{\delta}} \text{NegLog}(\boldsymbol{\mu}_{\text{A-GD}}^{(s+1)}, \boldsymbol{\delta}_{\text{A-GD}}^{(s)}) \right\}^+ \quad \text{for some step size } \gamma^{(s)} > 0
 \end{aligned}$$

We use backtracking linesearch to choose the step sizes in each iteration and run 300 iterations of the A-GD method starting from the initial solution $\boldsymbol{\mu}_{\text{A-GD}}^{(0)} = \mathbf{0}$, $\boldsymbol{\delta}_{\text{A-GD}}^{(0)} = \mathbf{0}$, which resulted in comparable runtime to the A-MM method (avg. 6-hr runtime per instance for A-GD compared to 5.9-hr for A-MM for instances with degree 8 and height 5).

The NegLogGap values for both methods are presented in a table shown in Figure EC.3. While the magnitude of improvement is lower compared to the Knitro and the PGD benchmarks in Section 4, our proposed A-MM method still outperforms A-GD across all problem sizes, obtaining significantly lower negative log-likelihood values for the (harder) problem instances with smaller values of the nest dissimilarity parameters.

To further understand the source of these improvements, in Figure EC.4, we plot the NegLogGap values (on a log-scale) for both methods as a function of the iteration number on a sample problem

Degree	Height	# Prods.	# Nodes	λ_{lower}	NegLogGap		NegLog Impr.	% better	
					A-MM	A-GD			
5	4	625	781	0.50	2.6	2.7	0.1	64	
				0.10	7.8	15.5	7.7	100	
				0.01	9.9	27.2	17.3	100	
	5	3,125	3,906	0.50	3.6	4.5	0.9	100	
				0.10	13.2	24.3	11.1	100	
				0.01	21.6	49.0	27.4	100	
	6	4	1,296	1,555	0.50	2.5	2.7	0.2	70
					0.10	8.5	16.1	7.6	100
					0.01	11.1	30.8	19.7	100
5		7,776	9,331	0.50	3.3	3.8	0.5	100	
				0.10	12.8	19.6	6.8	100	
				0.01	20.9	43.0	22.1	100	
7		4	2,401	2,801	0.50	2.4	2.9	0.5	99
					0.10	8.6	14.7	6.1	100
					0.01	11.4	29.6	18.2	100
	5	16,807	19,608	0.50	2.9	3.3	0.4	100	
				0.10	12.1	16.7	4.6	100	
				0.01	19.0	37.8	18.8	100	
	8	4	4,096	4,681	0.50	2.2	2.6	0.4	100
					0.10	8.0	13.0	5.0	100
					0.01	10.5	28.3	17.8	100
5		32,768	37,449	0.50	2.5	2.8	0.3	100	
				0.10	10.7	14.6	3.9	100	
				0.01	17.2	34.2	17.0	100	

Figure EC.3 Comparison of the performances of the A-GD benchmark and our proposed A-MM method in fitting tree logit models to choice data. The columns “Degree”, “Height”, and λ_{lower} report the degree of each non-leaf node, the height of the tree, and the lower bound on the nest dissimilarity parameters, respectively. The columns “# Prods.” and “# Nodes” report the number of products and the number of nodes in the tree, respectively. The columns under A-MM and A-GD report the average NegLogGap for each method. Recall that smaller values for the gaps are preferred. The “NegLog Impr.” column reports the average improvement in the negative log-likelihood value that our A-MM method achieves over the A-GD benchmark, and the “% better” column reports the percentage of instances in which the A-MM method obtains a lower NegLog value than A-GD. All the “NegLog Impr.” numbers are significantly different from zero at the 1% significance level under a paired samples t-test.

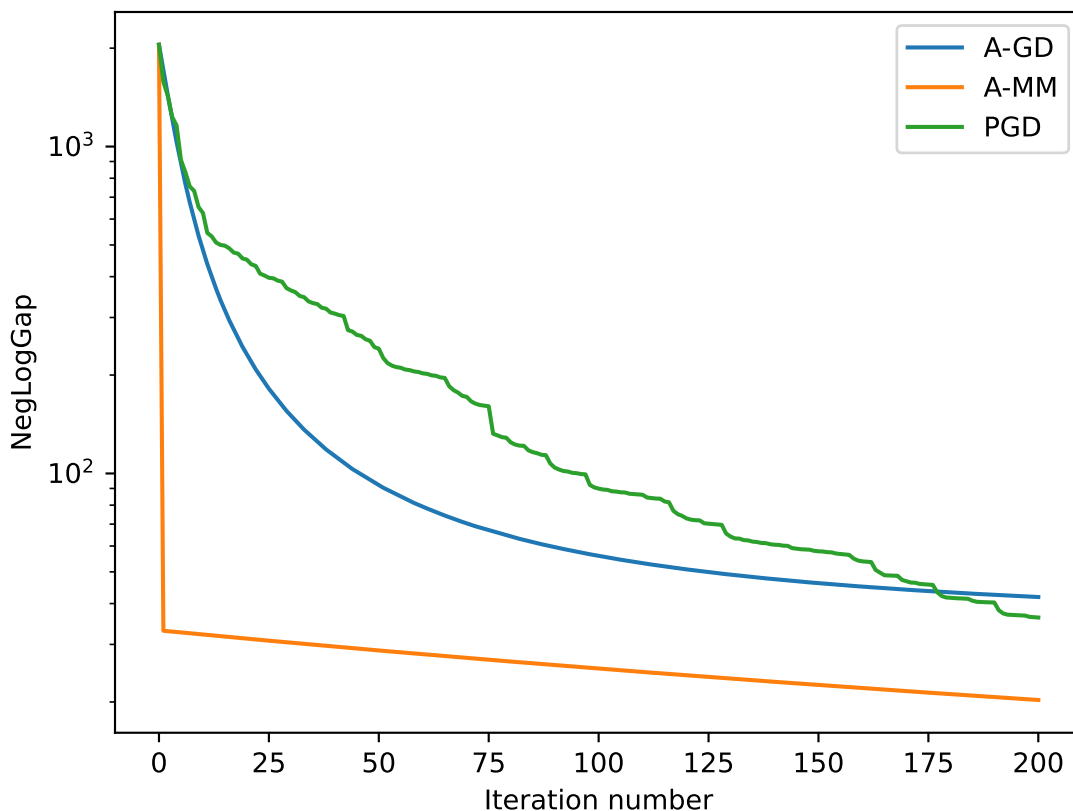


Figure EC.4 NegLogGap as a function of the iteration number for our proposed A-MM algorithm as well as the A-GD and PGD benchmarks on a sample problem instance under ground-truth parameters degree = 8, height = 5, and $\lambda_{\text{lower}} = 0.01$. Note the y-axis is on a log-scale.

instance under ground-truth parameters $r = 8$ (degree), $H = 5$ (height), and $\lambda_{\text{lower}} = 0.01$. The NegLogGap values achieved by the PGD method are also shown for comparison. It can be seen that the A-MM method consistently achieves a lower NegLogGap compared to both benchmarks. In particular, our proposed method makes rapid progress during the first few iterations, which is in line with the theoretical results. In contrast, the two benchmarks make more gradual progress leading to a larger NegLogGap for the same number of iterations. Moreover, an interesting thing to note is that the A-GD method converges faster than PGD, although PGD eventually achieves a lower NegLogGap value for this instance. This could be because PGD utilizes the “full” gradient of the negative log-likelihood, whereas A-GD alternates between gradient descent steps for the mean utility and nest dissimilarity parameters. Future work can investigate this difference in more detail.

Appendix I: Code and Data

The data and source code used in all the numerical studies in the paper can be downloaded at <https://github.com/ashwin90/TreeLogitEstimation>.

References

- Allen-Zhu, Zeyuan, Lorenzo Orecchia. 2014. Linear coupling: An ultimate unification of gradient and mirror descent. *arXiv preprint arXiv:1407.1537* .
- Boyd, S., L. Vandenberghe. 2004. *Convex Optimization*. Cambridge University Press, Cambridge, UK.
- Bubeck, S. 2015. Convex optimization: Algorithms and complexity. *Foundations and Trends in Machine Learning* **8**(3-4) 231–357.
- Cardell, N. S. 1997. Variance components structures for the extreme-value and logistic distributions with applications to models of heterogeneity. *Economic Theory* **13**(2) 185–213.
- Csiszár, Imre, János Körner. 2011. *Information theory: coding theorems for discrete memoryless systems*. Cambridge University Press.
- Gumbel, E. J. 2004. *Statistics of Extremes*. Dover Publications, Meneola, NY.
- Kiefer, J. 1953. Sequential minimax search for a maximum. *Proceedings of the American Mathematical Society* **4**(3) 502–506.
- McFadden, D. 1978. Modeling the choice of residential location. *Transportation Research Record* (672) 72–77.
- Nesterov, Y. 2013. *Introductory Lectures on Convex Optimization: A Basic Course*, vol. 87. Springer Science & Business Media.
- Vojnovic, Milan, Se-Young Yun, Kaifang Zhou. 2020. Convergence rates of gradient descent and mm algorithms for bradley-terry models. *International Conference on Artificial Intelligence and Statistics*. PMLR, 1254–1264.
- Zhang, Xiao-Dong. 2011. The laplacian eigenvalues of graphs: a survey. *arXiv preprint arXiv:1111.2897* .