

LEARNING MIXED MULTINOMIAL LOGITS WITH PROVABLE GUARANTEES

Yiqun Hu   

David Simchi-Levi, MIT

Zhenzhen Yan, Nanyang Technological University

INTRODUCTION

MIXED MULTINOMIAL LOGIT (MMNL)

- Consider a set of $[m] = \{1, \dots, m\}$ alternatives
- Population modeled by K MNL mixtures
- Each mixture k :
 - ▶ shares the same utility $V_{kj}, j \in [m]$
 - ▶ exhibits the same logit model

$$q_{kj} = \frac{\exp(V_{kj})}{\sum_{i \in [m]} \exp(V_{ki})} \quad \forall j \in [m]$$

- ▶ associates with a mixture weight α_k
- ▶ Also known as the softmax function

MIXED MULTINOMIAL LOGIT (MMNL)

- Assume linear utility given observed candidate feature vector $\mathbf{z}_j \in \mathbb{R}^d, \forall j \in [M]$

$$V_{kj} = \boldsymbol{\beta}_k^\top \mathbf{z}_j, \forall k$$

- *Aggregated choice probability:*

$$\mathbf{g} = \sum_{k=1}^K \alpha_k \mathbf{q}_k \in \mathbb{R}^m$$

where $q_{kj} := q_j(\boldsymbol{\beta}_k) = \frac{\exp(\boldsymbol{\beta}_k^\top \mathbf{z}_j)}{\sum_{i \in [m]} \exp(\boldsymbol{\beta}_k^\top \mathbf{z}_i)}$

- Goal: estimate $\alpha_k, \boldsymbol{\beta}_k, k = 1, \dots, K$

Learning MMNL

- Heuristics
 - ▶ EM algorithm (Train 09')
- Learning algorithms
 - ▶ Uniform 2-MNL (Chierichetti et al 18')
 - ▶ Arbitrary 2-MNL (Tang 20')
- Hybrid (convergence only at aggregated level)
 - ▶ Frank-Wolfe (Jagabathula et al 20')

STOCHASTIC SUBREGION FRANK-WOLFE (SSRFW)

PROBLEM FORMULATION

- Data assumption
 - ▶ Population of size N
 - ▶ For each time period $t = 1, \dots, T$:
 - ◆ Observe historical decision for each decision maker i : $Y_i^{(t)} \in \mathbb{R}^M$ with $Y_{ij}^{(t)} = \mathbb{1}_{[i \text{ chose } j \text{ at time } t]}$
 - ◆ Compute **observed share**: $y_j^{(t)} = \frac{1}{N} \sum_{i=1}^N Y_{ij}^{(t)}$
- Learning objective:

$$\min_{\mathbf{g} \in \text{Conv}(\overline{\mathcal{P}})} \mathcal{L}(\mathbf{g}; \mathbf{y}) \quad \equiv \quad \min_{\mathbf{g} \in \text{Conv}(\overline{\mathcal{P}})} \frac{1}{2} \sum_{t=1}^T \left\| \mathbf{y}^{(t)} - \mathbf{g} \right\|^2$$

$$\text{where } \mathbf{g} = \sum_{k=1}^K \alpha_k \mathbf{q}_k(\beta_k)$$

STOCHASTIC SUBREGION FRANK-WOLFE

- Let \mathcal{P} be the set of all logit vectors given $z_j, \forall j$
- Construct a candidate set $\mathcal{Q} \subset \mathcal{P}$ of logit vectors
- Require $\exists \pi : [L] \rightarrow [K], \left\| \mathbf{q}_\ell - \mathbf{q}_{\pi(\ell)}^* \right\| \leq \epsilon, \forall \mathbf{q}_\ell \in \mathcal{Q}$
 L is the number of elements in \mathcal{Q}
- Each of the extreme points of $\text{Conv}(\mathcal{Q})$ is close to some ground truth \mathbf{q}_k

SCORE MATRIX

For decision maker i (of type k)

- Historical data: $Y_i^{(t)} \in \mathbb{R}^M, t = 1, \dots, T$
- Define $X_i^{(t)}, t = 1, \dots, T$ i.i.d random variable with pmf q_k^*
 $X_i^{(t)} = j$ if $Y_{ij}^{(t)} = 1$
- Compute the empirical CDF

$$F_T(x; i) = \frac{1}{T} \sum_{t=1}^T \mathbb{1}_{\{X_i^{(t)} \leq x\}}, x \in [M]$$

Pairwise score (dissimilarity) between i and j

$$s(i, j) = \|F_T(x; i) - F_T(x; j)\|_\infty$$

Q CONSTRUCTION ALGORITHM

$$\mathcal{Q} = \left\{ \mathbf{q}_\ell \mid \mathbf{q}_\ell = \frac{1}{nT} \sum_{i \in I_\ell} \sum_{t=1}^T Y_i^{(t)} \right\}_{\ell=1, \dots, L}$$

Input: score matrix S , number of subsamples L , subsample size n

Initialization: $\mathcal{Q} = \text{set}()$

```
1 for  $\ell \leftarrow 1$  to  $L$  do
2   Choose seed:  $i \sim U(0, N)$ 
3   Initiate:  $I = \text{set}()$ 
4   while  $|I| \neq n$  do
5      $j \leftarrow \text{random\_sample}([N] \setminus I)$ 
6     Generate  $u \sim U(0, 1)$ 
7     if  $u < p_{j|i}$  then
8        $I.\text{add}(j)$ 
9     end
10  end
11  Compute  $\mathbf{q}_\ell = \frac{1}{nT} \sum_{i \in I} \sum_{t=1}^T Y_i^{(t)}$ 
12   $\mathcal{Q}.\text{add}(\mathbf{q}_\ell)$ 
13 end
```

Output: \mathcal{Q}

THE SSRFW ALGORITHM

Input: data \mathbf{y} , \mathcal{Q}

Initialization: $k = 0$; $\alpha^{(0)} = [1]$, a random $\mathbf{g}^{(0)}$

1 **while** *stopping condition not met* **do**

2 $k \leftarrow k + 1$

3 Compute $\mathbf{q} = \arg \min_{\mathbf{v} \in \text{Conv}(\mathcal{Q})} \langle \nabla \mathcal{L}(\mathbf{g}^{(k-1)}; \mathbf{y}), \mathbf{v} - \mathbf{g}^{(k-1)} \rangle$

→ *support finding step*

4 Compute $\alpha^{(k)} = \arg \min_{\alpha \in \Delta_k} \mathcal{L} \left(\alpha_0^{(k)} \mathbf{g}^{(0)} + \sum_{s=1}^k \alpha_s^{(k)} \mathbf{q}^{(s)} \right)$

→ *proportions update step*

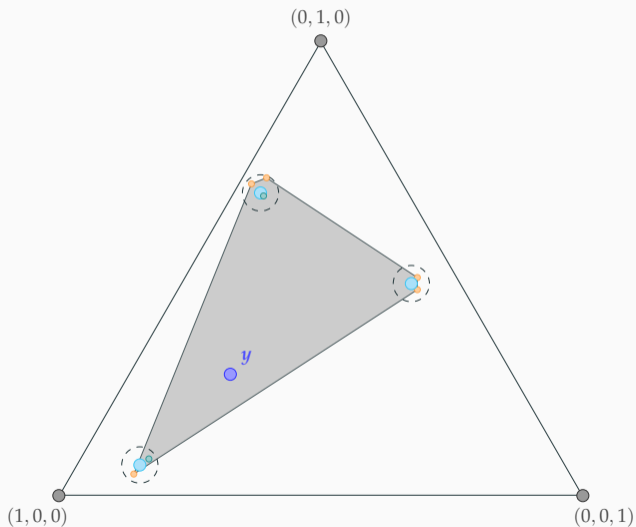
5 Update $\mathbf{g}^{(k)} := \alpha_0^{(k)} \mathbf{g}^{(0)} + \sum_{s=1}^k \alpha_s^{(k)} \mathbf{q}^{(s)}$

6 **end**

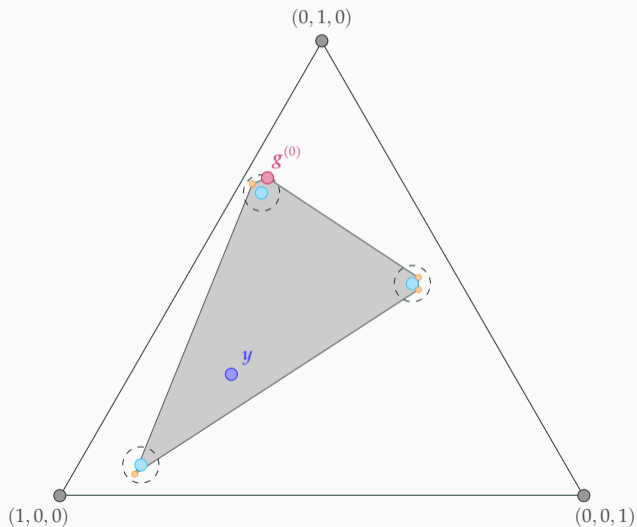
Output: choice prob. $\mathbf{q}^{(0)}, \dots, \mathbf{q}^{(k)}$

mixture weights. $\alpha^{(k)} \in \Delta_k \subset \mathbb{R}^{k+1}$

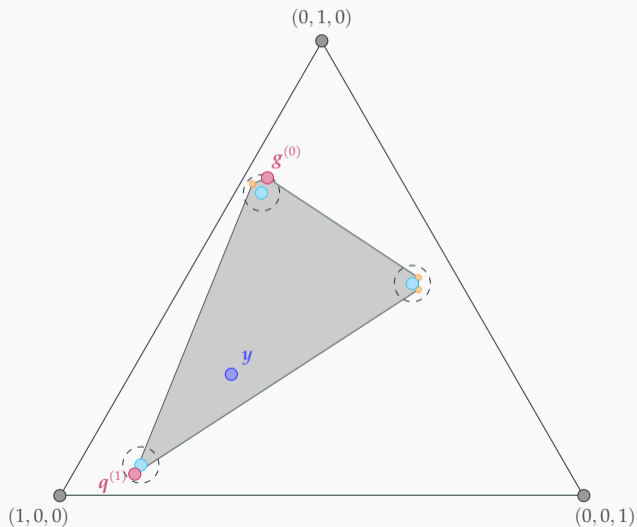
SSRFW ILLUSTRATION



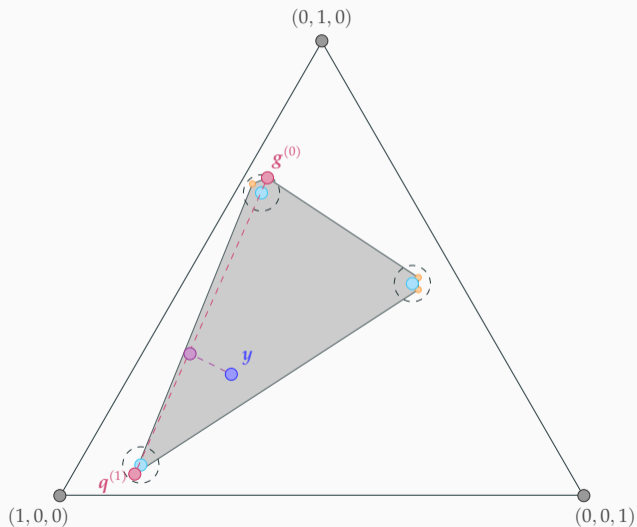
SSRFW ILLUSTRATION



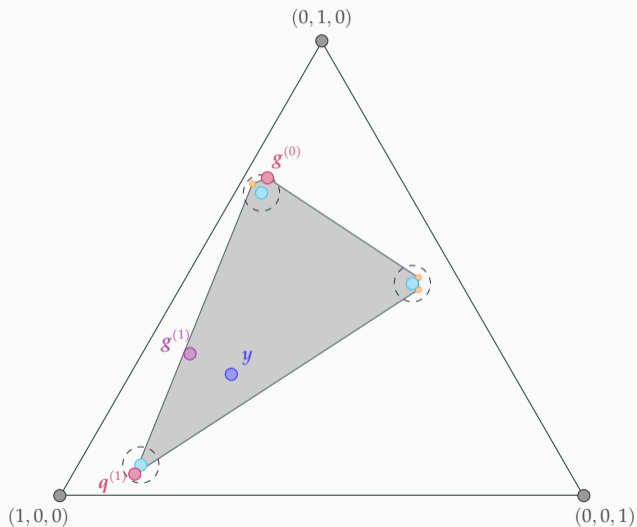
SSRFW ILLUSTRATION



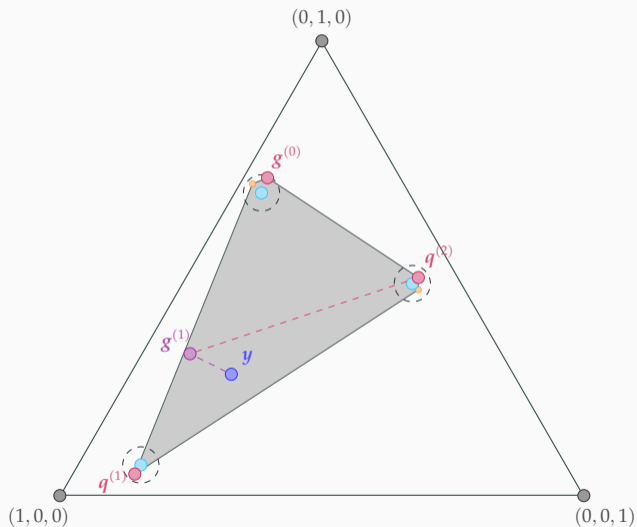
SSRFW ILLUSTRATION



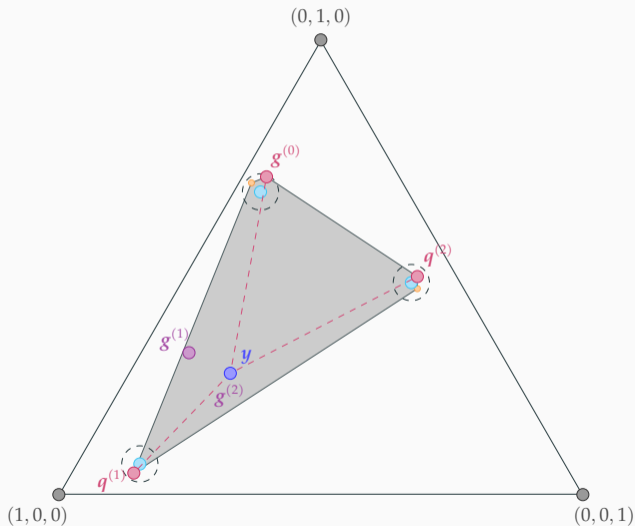
SSRFW ILLUSTRATION



SSRFW ILLUSTRATION



SSRFW ILLUSTRATION



ADVANTAGES

- Utilize individual choice data
- No prior knowledge on the number of mixtures (K) needed - no model misspecification
- Provable convergence on the estimators based on the Q construction algorithm

MAIN RESULT

Theorem 1

Let $\mathbf{g} = \sum_{k=1}^K \alpha_k \mathbf{q}_k$ be a mixed multinomial logit (MMNL) model over a set of M items. Assume $M \geq K$. For any $\epsilon > 0, 0 < \delta < 1$, SSRFW outputs an MMNL $\hat{\mathbf{g}} = \sum_{k=1}^{K'} \hat{\alpha}_k \hat{\mathbf{q}}_k$ where $K' \geq K$ such that, with probability $\geq 1 - \delta$, there exists a many-to-one mapping $\pi : j \mapsto i, j \in [K'], i \in [K]$ such that

$$\left\| \hat{\mathbf{q}}_j - \mathbf{q}_{\pi(j)} \right\| \leq \epsilon, \forall j \text{ and } \left| \sum_{j:\pi(j)=i} \hat{\alpha}_j - \alpha_i \right| \leq \epsilon, \forall i$$

The number of samples required $n(\epsilon, \delta)$ is polynomial in $\frac{1}{\epsilon}$ and $\frac{1}{\delta}$.

NUMERICAL EXPERIMENTS

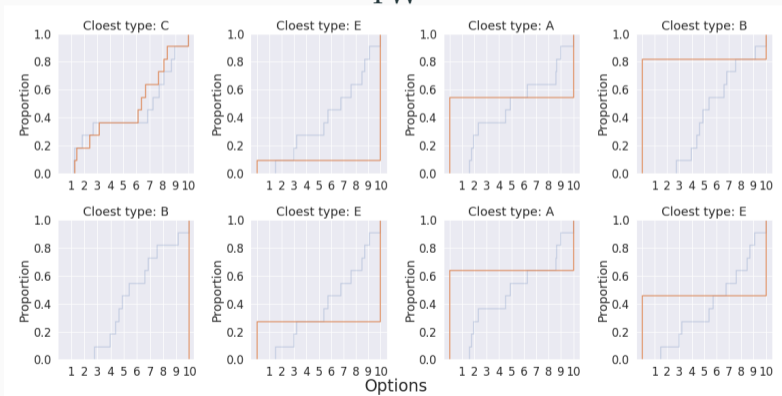
Simulation Studies

Case Study on Nielsen Panel Data

MIXTURE RECOVERY

$q^{(0)}, \dots, q^{(K')}$ output from the algorithm

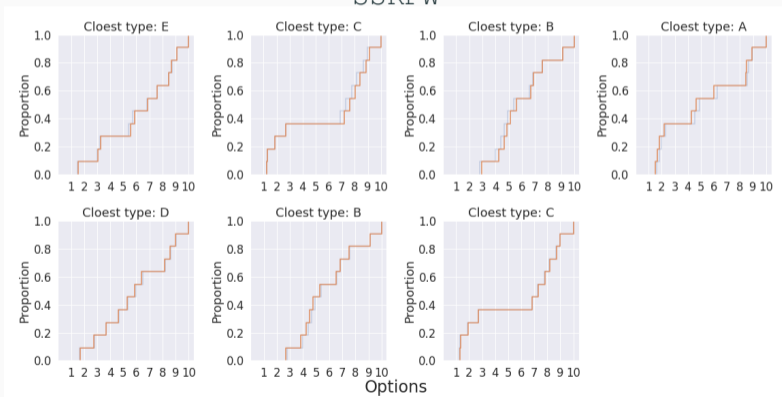
FW



MIXTURE RECOVERY

$q^{(0)}, \dots, q^{(K')}$ output from the algorithm

SSRFW



EXAMPLE RESULTS (NIELSEN PANEL DATA)

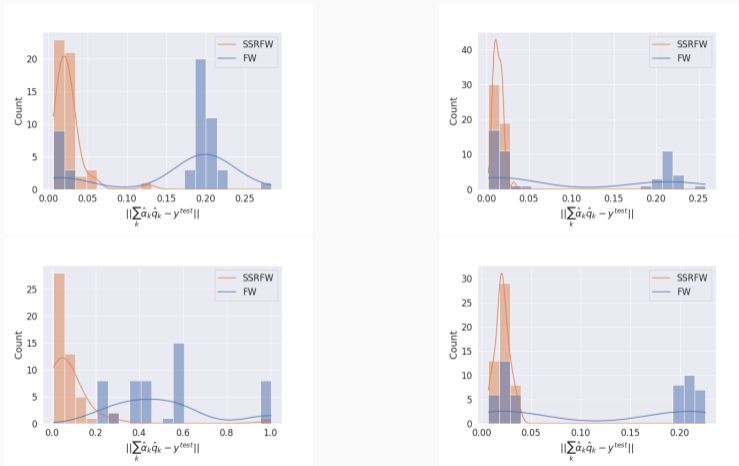


Figure 1: Categories: yogurt, pet food, candy, cereal

CONCLUSION

CONCLUSION

- Developed a new algorithm: SSRFW
 - ▶ based on the FW framework
 - ▶ utilize personal-level choice data
 - ▶ provide theoretical guarantees on the estimators
 - ▶ recover true model parameters
- Conducted various numerical experiments to compare SSRFW and the original FW
 - ▶ Simulation study to compare to the ground truth
 - ▶ Case studies on other Nielsen Consumer Panel Data

THANK YOU