

Stochastic adaptive resampling for the estimation of discrete choice models

15th Workshop on Discrete Choice Models
1–3 June 2023 | EPFL, Switzerland

Nicola Ortelli^{1,2}, Matthieu de Lapparent¹, Michel Bierlaire²

¹ IIDE, HEIG-VD

² TRANSP-OR, EPFL

DCMs in the era of big data

Ever-larger choice datasets

- Exponential growth of collected data:
 - “Wider” data — more variables;
 - “Taller” data — more observations.
- Two distinct problems: specification and estimation.

Model specification is iterative

- More candidates are tested before reaching an adequate model...
- ... and each of them takes longer to estimate.

Estimating DCMs

Maximum likelihood estimation (MLE)

- Suppose a dataset $\mathcal{N} = \{(\mathbf{x}_n, i_n) : n = 1, \dots, N\}$.
- Maximize the joint probability of the observed choices:

$$\max_{\theta} \mathcal{L}(\theta) = \max_{\theta} \sum_{n=1}^N \log P(i_n | \mathbf{x}_n; \theta).$$

- Solved using iterative methods: BFGS, Newton, etc.
- Each iteration is $\mathcal{O}(N)$.

Intuition

Factoring-out redundancy

- If the data contains identical observations:

$$\mathcal{L}(\theta) = \sum_{u=1}^U N_u \times \log P(i_u | \mathbf{x}_u; \theta).$$

- U unique observations ($U \leq N$).
- Each appears N_u times in the original data.
- Evaluating $\mathcal{L}(\theta)$ is faster by a ratio of approximately $\frac{U}{N}$.

⇒ Extend this trick to “nearly identical” observations!

Resampling estimation of DCMs

Procedure

- Group similar observations.
- Sample from each group.
- Associate weights based on group sizes.
- Use MLE on the **weighted** subsample.

Challenges

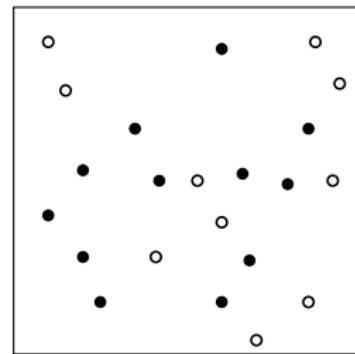
- Minimize information loss.
- Still, clustering must be fast!

⇒ **Use locality-sensitive hashing (LSH).**

LSH-based dataset reduction (LSH-DR)

Toy data

- 2 alternatives.
- 2 explanatory variables.

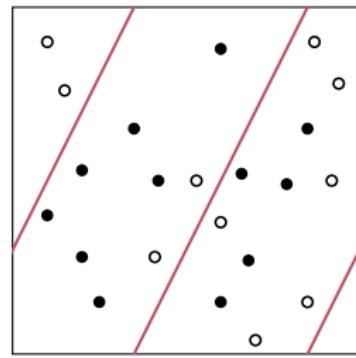


LSH-based dataset reduction (LSH-DR)

LSH function

$$h_{\mathbf{a}, b}(\mathbf{x}_n) = \left\lfloor \frac{\mathbf{a} \cdot \mathbf{x}_n + b}{w} \right\rfloor$$

- $\mathbf{a} \stackrel{iid}{\sim} \mathcal{N}(0, 1)$.
- $b \sim \mathcal{U}(0, w)$.
- w is the **bucket width**.

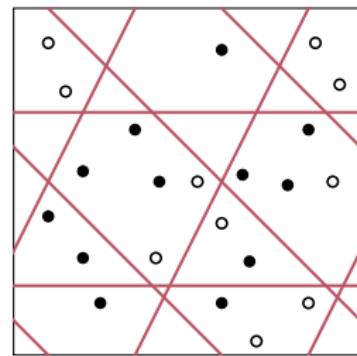


LSH-based dataset reduction (LSH-DR)

Combined LSH functions

$$\begin{aligned} H_{\mathbf{A}, \mathbf{B}}(\mathbf{x}_n) = H_{\mathbf{A}, \mathbf{B}}(\mathbf{x}_p) \Leftrightarrow \\ h_{\mathbf{a}_r, b_r}(\mathbf{x}_n) = h_{\mathbf{a}_r, b_r}(\mathbf{x}_p) \forall r = 1, \dots, R. \end{aligned}$$

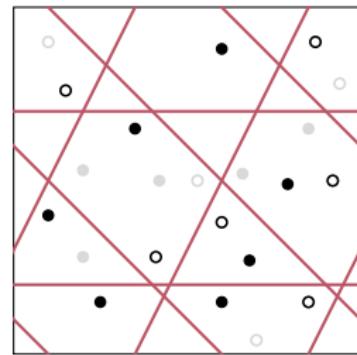
- $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_R)$.
- $\mathbf{B} = (b_1, \dots, b_R)$.



LSH-based dataset reduction (LSH-DR)

Sampling

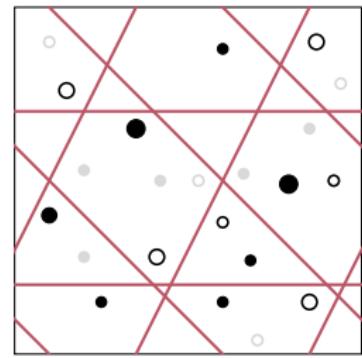
In each bucket, sample one observation (x_g, i_g) per alternative.



LSH-based dataset reduction (LSH-DR)

Associated weights N_g

$$N_g = |\{(x_n, i_n) : H_{A,B}(x_g) = H_{A,B}(x_n), i_g = i_n\}|$$



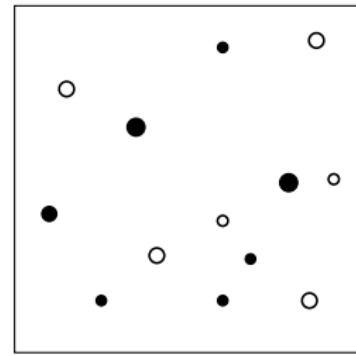
LSH-based dataset reduction (LSH-DR)

Subsample, weighted

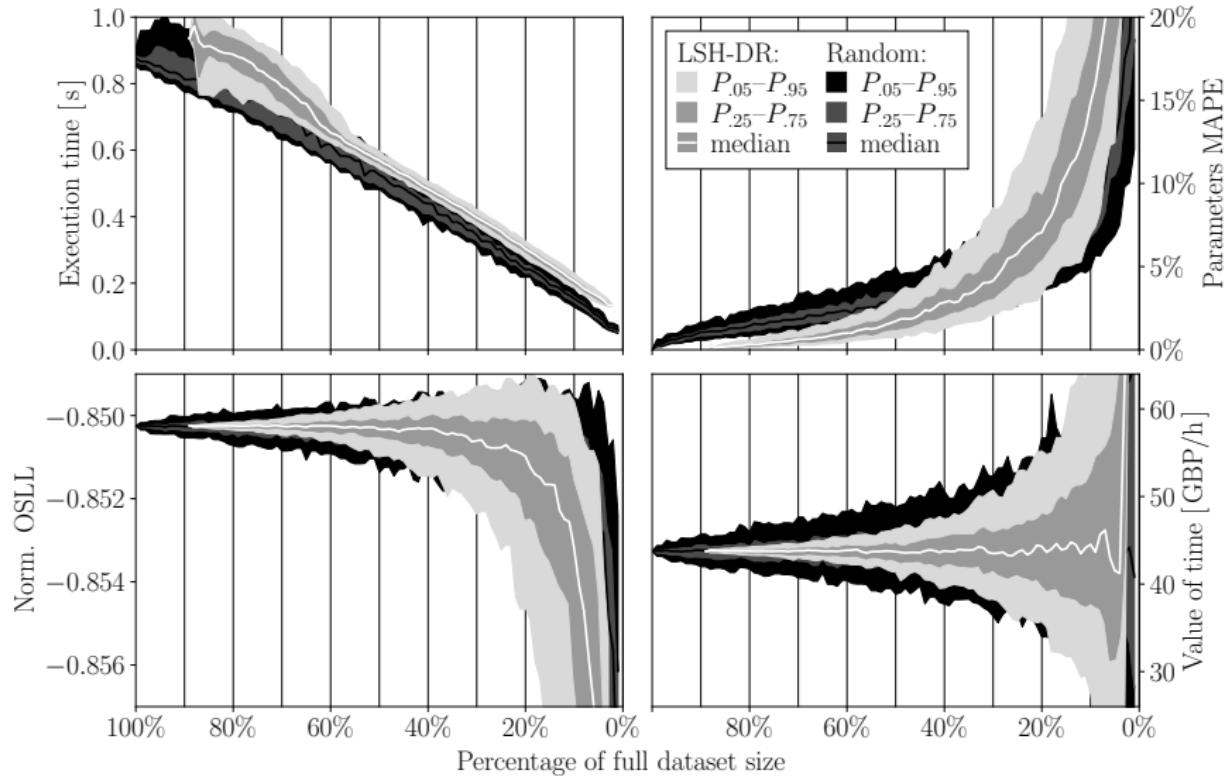
$$\mathcal{N}^* = \{(x_g, i_g, N_g) : g = 1, \dots, G\}$$

Weighted MLE

$$\mathcal{L}^*(\theta) = \sum_{g=1}^G N_g \cdot \log P(i_g | x_g; \theta)$$



Results





Resampling estimation of discrete choice models

Nicola Ortelli *† Matthieu de Lapparent * Michel Bierlaire †

March 30, 2023

Report TRANSP-OR 230330
Transport and Mobility Laboratory
School of Architecture, Civil and Environmental Engineering
Ecole Polytechnique Fédérale de Lausanne
transp-or.epfl.ch

*School of Management and Engineering Vaud, HES-SO // University of Applied Sciences and Arts Western Switzerland, {nicola.ortelli,matthieu.delapparent}@heig-vd.ch
†Transport and Mobility Laboratory, School of Architecture, Civil and Environmental Engineering, Ecole Polytechnique Fédérale de Lausanne (EPFL), {nicola.ortelli,michel.bierlaire}@epfl.ch

Spin-off

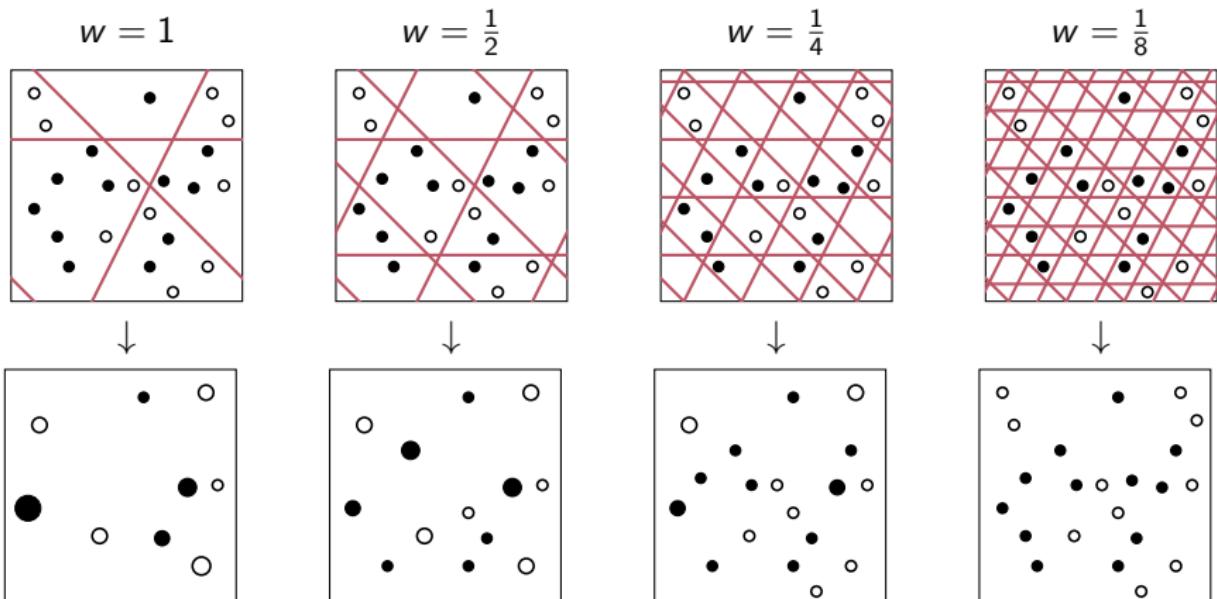
LSH-DR

- Substantial time savings when rough estimates are sufficient.
- What if the full-dataset estimates are needed?

Stochastic adaptive resampling (STAR)

- Embed LSH-DR within the model estimation process.
- Generate **weighted** batches for stochastic optimization. [Lederrey *et al.*, 2021]
- Start small and increase batch size dynamically.

Illustration



Stochastic adaptive resampling (STAR)

Generic algorithm

- Input:
 - \mathcal{N} : full dataset;
 - θ_0 : initial solution;
- Initialization:
 - $k \leftarrow 0$;
- Repeat:
 - ① $\theta_{k+1} \leftarrow \text{newCandidate}(\theta_k, \mathcal{N})$;
 - ② $k \leftarrow k + 1$;
- Until $||\nabla_{\text{rel}} \mathcal{L}(\theta_k)|| < \varepsilon$.

Stochastic adaptive resampling (STAR)

Generic algorithm

- Input:
 - \mathcal{N} : full dataset;
 - θ_0 : initial solution;
- Initialization:
 - $k \leftarrow 0$;
- Repeat:
 - ① $\theta_{k+1} \leftarrow \text{newCandidate}(\theta_k, \mathcal{N})$;
 - ② $k \leftarrow k + 1$;
- Until $||\nabla_{\text{rel}} \mathcal{L}(\theta_k)|| < \varepsilon$.

Relative gradient

$$[\nabla_{\text{rel}} \mathcal{L}(\theta)]_j = [\nabla \mathcal{L}(\theta)]_j \cdot \frac{\theta_j}{\mathcal{L}(\theta)}$$

Stochastic adaptive resampling (STAR)

Generic algorithm + STAR

- Input:
 - \mathcal{N} : full dataset;
 - θ_0 : initial solution;
 - w_0 : initial bucket width.
- Initialization:
 - $k \leftarrow 0$;
- Repeat:
 - ① $\theta_{k+1} \leftarrow \text{newCandidate}(\theta_k, \mathcal{N})$;
 - ② $k \leftarrow k + 1$;
- Until $\|\nabla_{\text{rel}} \mathcal{L}(\theta_k)\| < \varepsilon$.

Relative gradient

$$[\nabla_{\text{rel}} \mathcal{L}(\theta)]_j = [\nabla \mathcal{L}(\theta)]_j \cdot \frac{\theta_j}{\mathcal{L}(\theta)}$$

Stochastic adaptive resampling (STAR)

Generic algorithm + STAR

- Input:
 - \mathcal{N} : full dataset;
 - θ_0 : initial solution;
 - w_0 : initial bucket width.
- Initialization:
 - $k \leftarrow 0$;
- Repeat:
 - ① $\mathcal{N}_k^* \leftarrow \text{LSH-DR}(\mathcal{N}_k, \mathcal{N})$;
 - ② $\theta_{k+1} \leftarrow \text{newCandidate}(\theta_k, \mathcal{N}_k^*)$;
 - ④ $k \leftarrow k + 1$;
- Until $\|\nabla_{\text{rel}} \mathcal{L}(\theta_k)\| < \varepsilon$.

Relative gradient

$$[\nabla_{\text{rel}} \mathcal{L}(\theta)]_j = [\nabla \mathcal{L}(\theta)]_j \cdot \frac{\theta_j}{\mathcal{L}(\theta)}$$

Stochastic adaptive resampling (STAR)

Generic algorithm + STAR

- Input:
 - \mathcal{N} : full dataset;
 - θ_0 : initial solution;
 - w_0 : initial bucket width.
- Initialization:
 - $k \leftarrow 0$;
- Repeat:
 - ① $\mathcal{N}_k^* \leftarrow \text{LSH-DR}(w_k, \mathcal{N})$;
 - ② $\theta_{k+1} \leftarrow \text{newCandidate}(\theta_k, \mathcal{N}_k^*)$;
 - ③ $w_{k+1} \leftarrow \text{updateW}(w_k, \theta_k, \theta_{k+1})$;
 - ④ $k \leftarrow k + 1$;
- Until $\|\nabla_{\text{rel}} \mathcal{L}(\theta_k)\| < \varepsilon$.

Relative gradient

$$[\nabla_{\text{rel}} \mathcal{L}(\theta)]_j = [\nabla \mathcal{L}(\theta)]_j \cdot \frac{\theta_j}{\mathcal{L}(\theta)}$$

Stochastic adaptive resampling (STAR)

Generic algorithm + STAR

- Input:
 - \mathcal{N} : full dataset;
 - θ_0 : initial solution;
 - w_0 : initial bucket width.
- Initialization:
 - $k \leftarrow 0$;
- Repeat:
 - ① $\mathcal{N}_k^* \leftarrow \text{LSH-DR}(w_k, \mathcal{N})$;
 - ② $\theta_{k+1} \leftarrow \text{newCandidate}(\theta_k, \mathcal{N}_k^*)$;
 - ③ $w_{k+1} \leftarrow \text{updateW}(w_k, \theta_k, \theta_{k+1})$;
 - ④ $k \leftarrow k + 1$;
- Until $\|\nabla_{\text{rel}} \mathcal{L}(\theta_k)\| < \varepsilon$.

Relative gradient

$$[\nabla_{\text{rel}} \mathcal{L}(\theta)]_j = [\nabla \mathcal{L}(\theta)]_j \cdot \frac{\theta_j}{\mathcal{L}(\theta)}$$

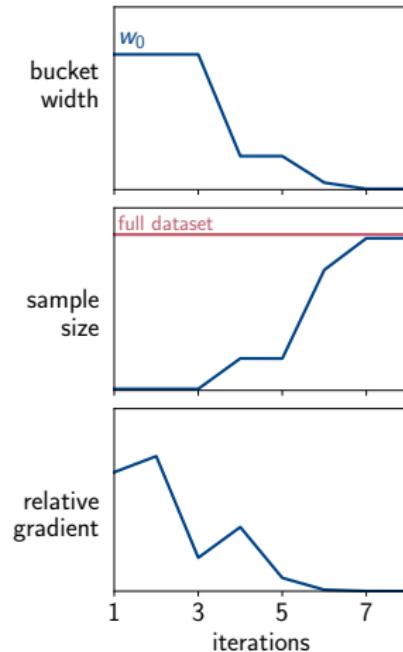
Bucket width update

$$w_{k+1} = w_k \cdot \min \left(1, \frac{\|\nabla_{\text{rel}} \mathcal{L}(\theta_{k+1})\|}{\|\nabla_{\text{rel}} \mathcal{L}(\theta_k)\|} \right)$$

Stochastic adaptive resampling (STAR)

Generic algorithm + STAR

- Input:
 - \mathcal{N} : full dataset;
 - θ_0 : initial solution;
 - w_0 : initial bucket width.
- Initialization:
 - $k \leftarrow 0$;
- Repeat:
 - ① $\mathcal{N}_k^* \leftarrow \text{LSH-DR}(w, \mathcal{N})$;
 - ② $\theta_{k+1} \leftarrow \text{newCandidate}(\theta_k, \mathcal{N}_k^*)$;
 - ③ $w_{k+1} \leftarrow \text{updateW}(w_k, \theta_k, \theta_{k+1})$;
 - ④ $k \leftarrow k + 1$;
- Until $\|\nabla_{\text{rel}} \mathcal{L}(\theta_k)\| < \varepsilon$.



Experimental design

LPMC data [Hillel et al., 2018]

- Mode choice, 4 alternatives.
- 81'086 observations.

Models [Hillel, 2019]

MNL-S

- 10 cont. variables.
- 0 dummies.
- **13 parameters.**

MNL-M

- 11 cont. variables.
- 15 dummies.
- **53 parameters.**

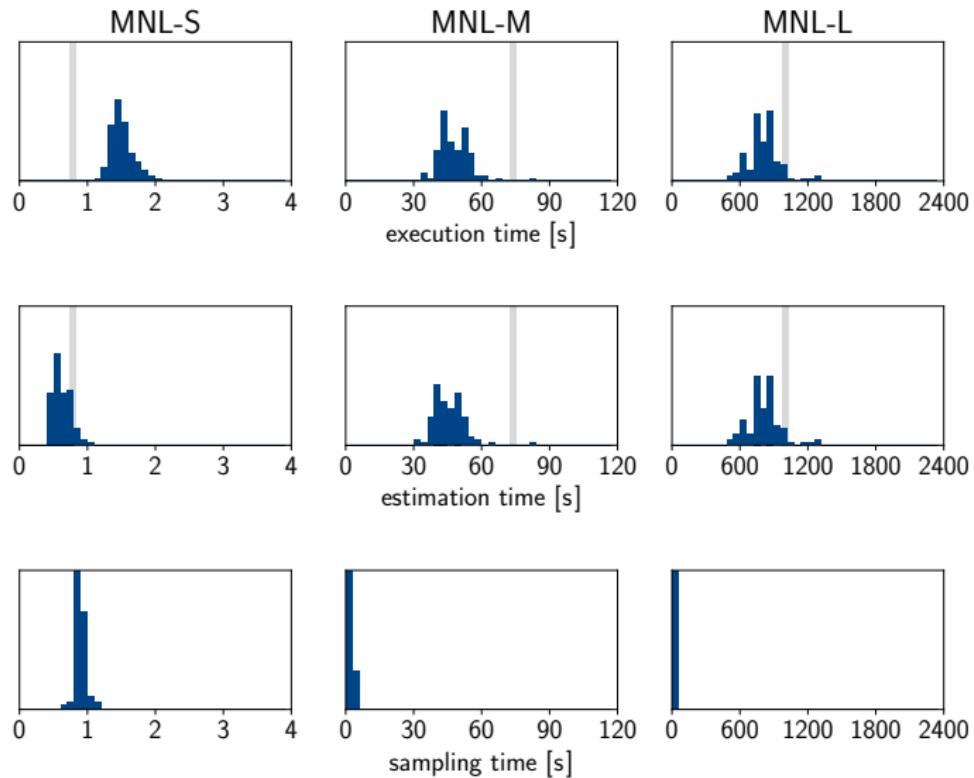
MNL-L

- 13 cont. variables.
- 18 dummies.
- **100 parameters.**

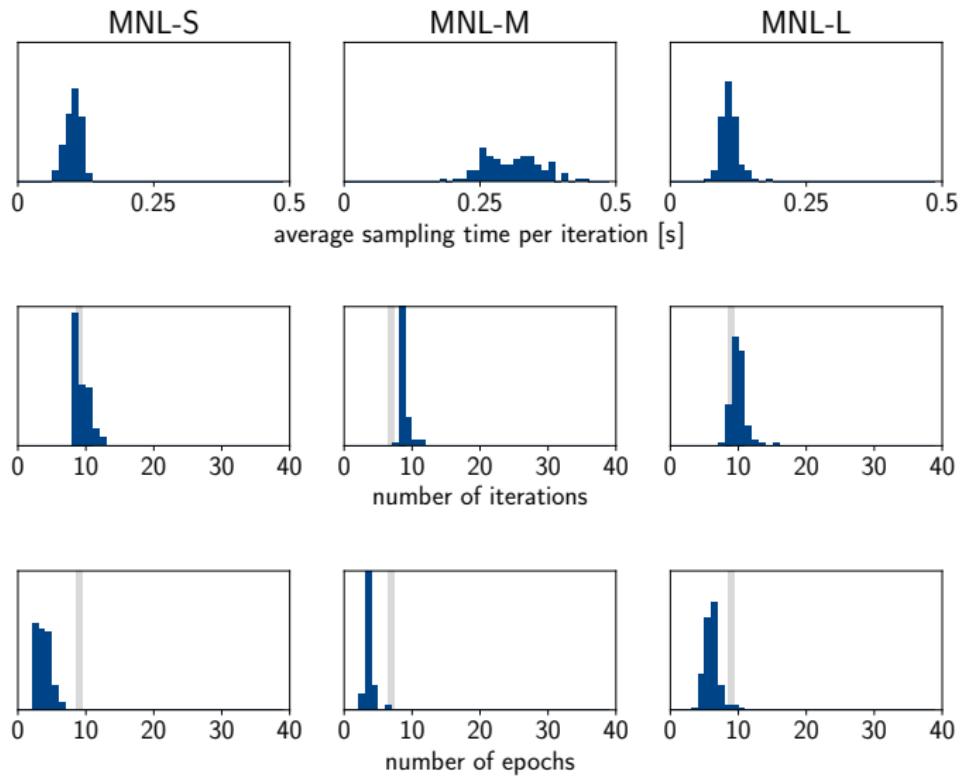
Optimization algorithm

- Newton-TR + STAR (100 repetitions).
- Benchmark: standard Newton-TR.

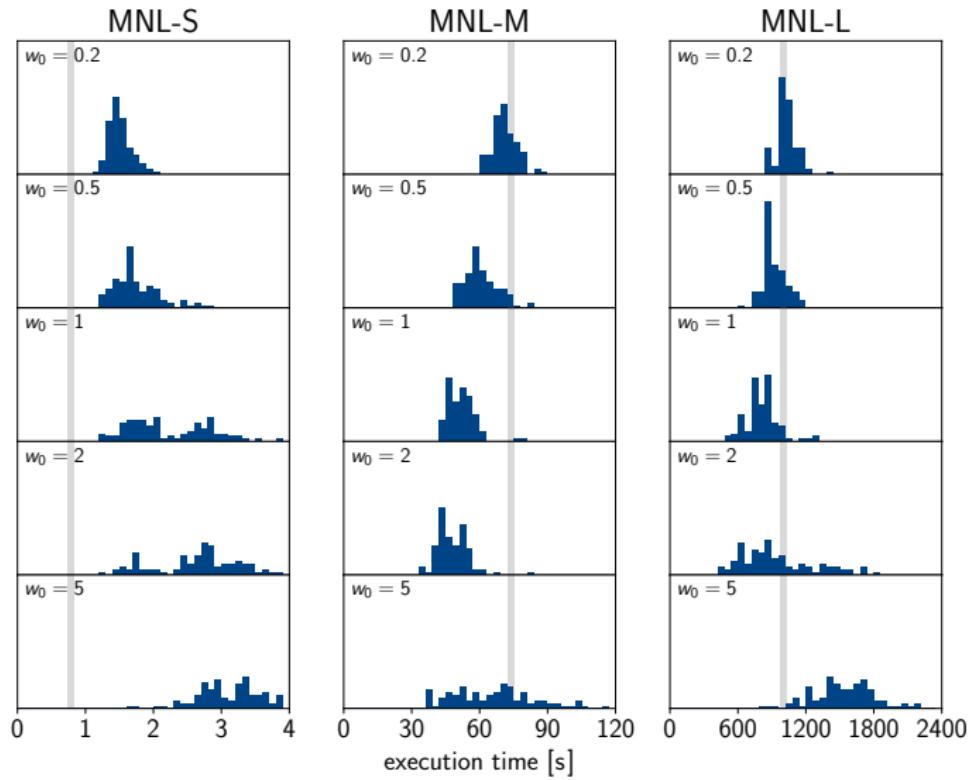
Results: estimation and sampling times



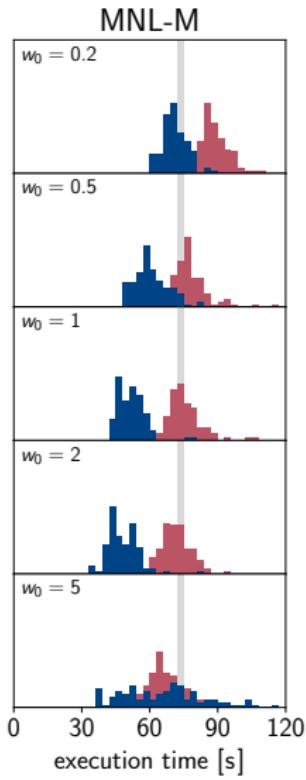
Results: iterations, epochs



Results: effect of w_0



Results: weighted/unweighted subsamples



Conclusion

Summary

- LSH-DR if rough estimates are sufficient.
- STAR when full-dataset estimates are needed.

Next steps

- Bucket width update.
- More advanced DCMs.
- Test with other optimization algorithms.

References

LSH-DR

- Ortelli, N., de Lapparent, M. and Bierlaire, M. (2023). Resampling estimation of discrete choice models, Technical Report, TRANSP-OR 230330. Transport and Mobility Laboratory, ENAC, EPFL.
- Ortelli, N., de Lapparent, M. and Bierlaire, M. (2022). Faster estimation of discrete choice models via dataset reduction, Proceedings of the 23rd Swiss Transportation Research Conference.

Direct precedents

- Lederrey, G., Lurkin, V., Hillel, T. and Bierlaire, M. (2021). Estimation of discrete choice models with hybrid stochastic adaptive batch size algorithms, Journal of choice modelling 38.

Dataset & models

- Hillel, T., Elshafie, M. Z. and Jin, Y. (2018). Recreating passenger mode choice-sets for transport simulation: A case study of London, UK, Proceedings of the Institution of Civil Engineers-Smart Infrastructure and Construction 171(1).
- Hillel, T. (2019). Understanding travel mode choice: A new approach for city scale simulation, PhD thesis, University of Cambridge.

Stochastic adaptive resampling for the estimation of discrete choice models

15th Workshop on Discrete Choice Models
1–3 June 2023 | EPFL, Switzerland

Nicola Ortelli^{1,2}, Matthieu de Lapparent¹, Michel Bierlaire²

¹ IIDE, HEIG-VD

² TRANSP-OR, EPFL