

hEART
2022

1-3 JUNE 2022
LEUVEN, BELGIUM

10th Symposium of the European Association
for Research in Transportation

EPFL

A Benders decomposition for maximum simulated likelihood estimation of advanced discrete choice models

hEART 2022

Tom Häring, Claudia Bongiovanni, Michel Bierlaire

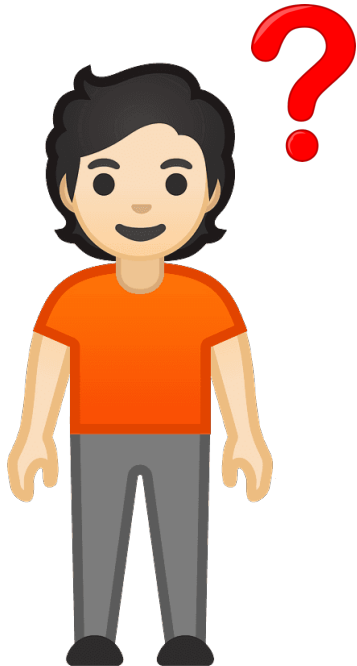
TRANSP-OR Laboratory, EPFL

Contents

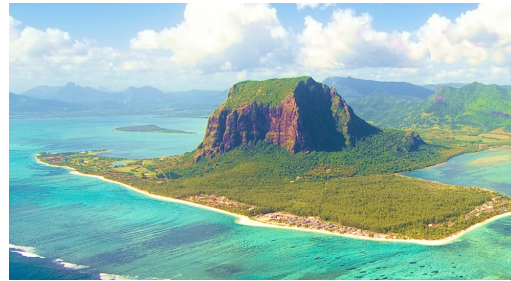
1. Why maximum likelihood estimation (MLE)?
2. Why simulated MLE?
3. Why a mixed integer linear program (MILP)?
4. Simulated MLE as an MILP
5. Why decomposition?
6. The Benders decomposition
7. Results
8. Ideas for future work

Why maximum likelihood estimation (MLE)?

- **MLE** is for example used to estimate the parameters of **discrete choice models**



hEART 2022



TRISTAN XI



ICMC 2022

Why maximum likelihood estimation (MLE)?

- For each **individual** n , every **alternative** i has an associated **utility**:

The diagram shows the utility function $U_{in} = U_{in}(\beta, x, \epsilon_{in})$ with three annotations:

- A green box labeled "parameters to be estimated" points to β .
- A blue box labeled "exogenous attributes" points to x .
- A red box labeled "random error term" points to ϵ_{in} .

- Assumptions:
 - I.) **linear** in parameters
 - II.) we can **draw** from error terms

Why maximum likelihood estimation (MLE)?

- For each **individual** n , every **alternative** i has an associated **utility**:

$$U_{in} = \sum_k \beta_k x_{ink} + \epsilon_{in} = \underbrace{V_{in}}_{\text{deterministic part}} + \underbrace{\epsilon_{in}}_{\text{stochastic part}}$$

- Behavioral assumption: the individual chooses the alternative with **the highest utility**

Why maximum likelihood estimation (MLE)?

- Data: **observed choices** y_{in} (= 1 if ind. n chose alternative i , else = 0)
- Find parameters β_k such that the **likelihood** of this outcome is **maximized**
- **Log-Likelihood function:**

$$\ln \left(\prod_n \prod_i P_n(i)^{y_{in}} \right) = \sum_n \sum_i y_{in} \ln P_n(i)$$

where

$$P_n(i) = \mathbb{P}(V_{in} + \epsilon_{in} \geq V_{jn} + \epsilon_{jn} \forall j \in J)$$

Why simulated MLE?

- DCMs model choices **realistically** [1], but in general lead to **non-convex** probabilities [2]
 - ➔ No global optimality certificates, **danger of local optima**
 - ➔ Non-convex solver \approx **Blackbox**
- **Simulation** mitigates this by giving a **linear** approximation [3] and allows DCMs to be easily **integrated** in **optimization models** [2]

[1] *Bierlaire: Discrete choice models (1998)*

[2] *Pacheco: Integrating advanced discrete choice models in mixed integer linear optimization (2021)*

[3] *Train: Discrete choice methods with simulation (2009)*

Why simulated MLE?

- How:

- **Simulate** R scenarios, utilities become **deterministic**:

$$U_{inr} = V_{in} + \epsilon_{inr} \leftarrow \text{Draw from distribution}$$

- Let ω_{inr} be the **choice variables**

- **Approximated** probabilities: $\hat{P}_n(i) = \frac{1}{R} \sum_{r=0}^{R-1} \omega_{inr}$

Why a mixed integer linear program (MILP)?

- Allow inclusion of **integer variables** in estimation procedure
 - Model **advanced** DCMs, e. g. **latent variables / classes**
 - Additional features, e. g. **automatic / assisted specification**
- Vast literature on efficient **modeling & performance**
- Gives **control** over **optimization process**: information on **bounds, optimality gaps, user-generated cuts**, etc.

Simulated MLE as an MILP

• **Objective:** max Log-Likelihood $\sum_n \sum_i y_{in} \ln P_n(i)$



max sim. Log-Likelihood $\sum_{in} y_{in} \ln \sum_{r=0}^{R-1} \omega_{inr} - y_{in} \ln R$



$$S_{in} = \sum_r \omega_{inr}$$

$$z_{in} \leq L_r - K_r S_{in}$$

$$\max \sum_n \sum_i y_{in} z_{in}$$

Simulated MLE as an MILP

- **Constraints:**

$$\sum_i \omega_{inr} = 1 \quad \forall n, r$$

$$U_{inr} = \sum_k \beta_k x_{ink} + \epsilon_{inr} \quad \forall i, n, r$$

$$U_{nr} \geq U_{inr} \quad \forall i, n, r$$

$$U_{nr} = \sum_i U_{inr} \omega_{inr} \quad \forall n, r$$

$$S_{in} = \sum_r \omega_{inr} \quad \forall i, n$$

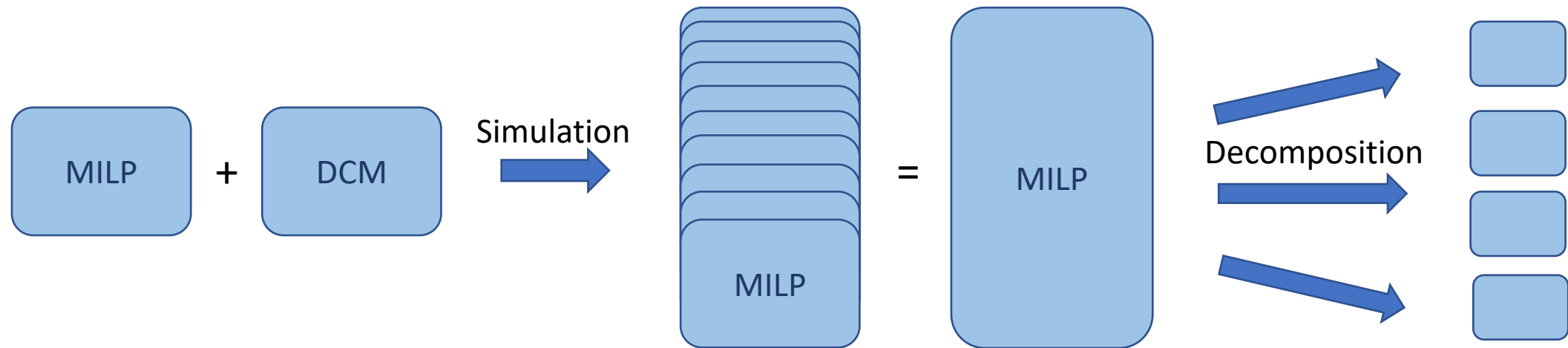
$$z_{in} \leq L_r - K_r S_{in} \quad \forall i, n$$

$$\omega_{inr} \in \{0, 1\}$$

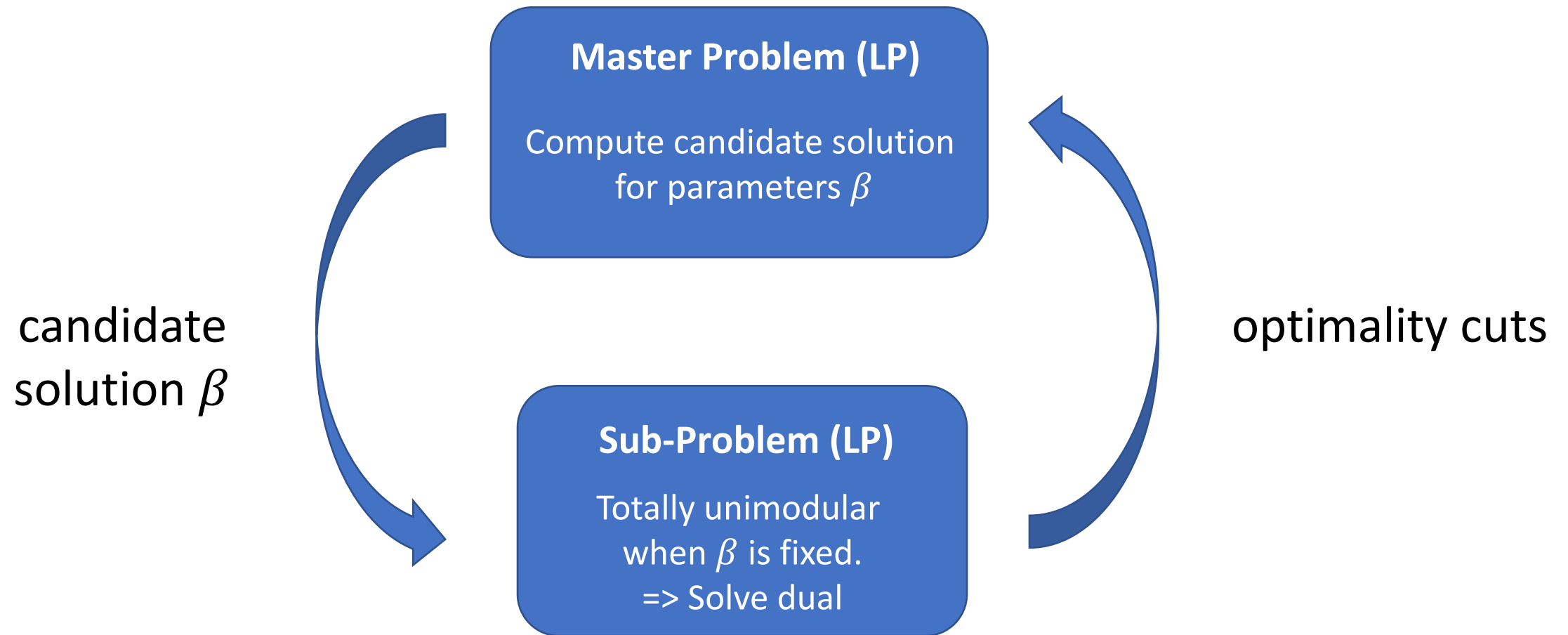
$$\beta, s, z, U, U \in \mathbb{R}$$

Why decomposition?

- Problem: Simulation **increases problem size** by solving **many scenarios**
 ➔ **only small instances** can be solved in reasonable time [1]
- To solve large MILPs efficiently we consider **decomposition methods**



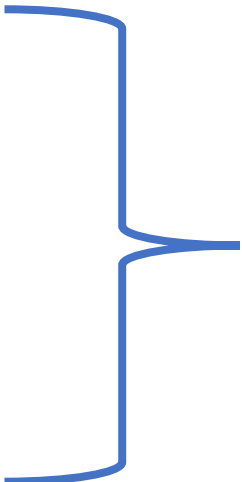
The Benders decomposition



The Benders decomposition

- For a **fixed** β_k the rest of the MILP becomes a **Knapsack-problem**
=> totally unimodular:

- Utilities become fixed $U_{inr} = \sum_k \beta_k^{\text{fixed}} x_{ink} + \epsilon_{inr}$

- Now:
$$\begin{aligned} U_{nr} &= \sum_i U_{inr} \omega_{inr} \\ U_{nr} &\geq U_{inr} \\ \sum_i \omega_{inr} &= 1 \\ \omega_{inr} &\in [0, 1] \end{aligned}$$
 
$$\omega_{i^*nr} = 1$$

for the alternative i^*
with highest utility

The Benders decomposition

- **Typically:**

- The variable to be fixed is **integer**, so that the subproblems are linear
- Thus **MP is an integer program (bottleneck!)**

- **But in our case:**

- The variable to be fixed is **continuous**, but thanks to TU-ness the subproblems are (*technically*) still linear!
- Thus **SP is a linear program**

From solving an MILP to iteratively solving LP's!

The Benders decomposition

- **Difficulty:**

Simply adding the constraint $\beta_k = \beta_k^{\text{fixed}}$ **does not work in our case** because of the **non-linearity** of the problem

The Benders decomposition

- **Constraints:**

Goal: linear in β_k

$$\sum_i \omega_{inr} = 1$$

$$U_{inr} = \sum_k \beta_k x_{ink} + \epsilon_{inr}$$

$$U_{nr} \geq U_{inr}$$

$$U_{nr} = \sum_i U_{inr} \omega_{inr}$$

Non-linear!

$$S_{in} = \sum_r \omega_{inr}$$

$$z_{in} \leq L_r - K_r S_{in}$$

$$\omega_{inr} \in [0, 1]$$

$$\beta, s, z, U, U \in \mathbb{R}$$

$$\forall n, r$$

$$\forall i, n, r$$

$$\forall i, n, r$$

$$\forall n, r$$

$$\forall i, n$$

$$\forall i, n$$

The Benders decomposition

- We design a **quasi**-linearization:

$$\eta_{inrk} \stackrel{!}{=} \beta_k^{\text{fixed}} \omega_{inr} \quad \longrightarrow \quad \begin{aligned} \chi_{inr} + \omega_{inr} &= 1 \\ \eta_{inrk} + \beta_k^{\text{fixed}} \chi_{inr} &= \beta_k^{\text{fixed}} \\ \sum_i \eta_{inrk} &= \beta_k \end{aligned}$$

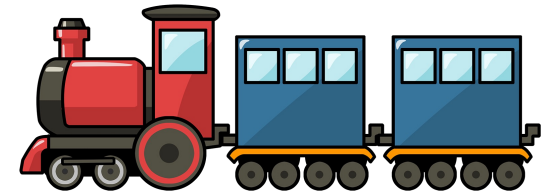
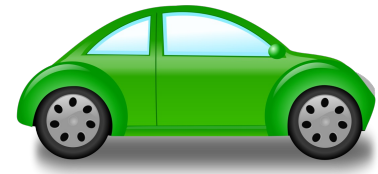
Application to a mode choice problem

- Dataset: **RP** data on **mode choice**, Netherlands, 1987
- Simple **binary logit model**:

choice between two modes – **car** and **rail**

$$U_{\text{car},n} = \beta_{\text{time}} * \text{traveltime}_{\text{car}}$$

$$U_{\text{rail},n} = \beta_{\text{time}} * \text{traveltime}_{\text{rail}}$$



- Compare **decomposition** vs. **undecomposed MILP**

N	R	sLL-M	sLL-D	Gap [%]	T-M	T-D
20	50	-12.607	-12.658	-0.40	64.942	10.061
20	100	-12.212	-12.258	-0.38	403.694	9.902
20	200	-12.283	-12.648	-2.97	1117.064	16.939
50	50	-30.848	-31.030	-0.59	286.679	29.780
50	100	-30.461	-31.040	-1.90	1558.604	65.006
50	200	-30.566	-30.692	-0.41	5375.655	98.206
100	50	-65.204	-65.801	-0.92	2820.229	28.781
100	100	-65.784	-67.419	-2.49	4346.067	274.163
100	200	-65.699	-66.018	-0.49	10800+	295.741
200	50	-123.551	-124.027	-0.39	1476.185	120.579
200	100	-124.000	-124.243	-0.20	10800+	327.253
200	200	-124.707	-124.106	0.48	10800+	1262.755

Application to a mode choice problem

- First **conjecture**: gaps are caused by **log-linearization** in MSLE
- **Remedy**: apply decomposition to *continuous pricing problem (CPP)*
 - ➡ Almost **equivalent** problem structure, **no log-linearization**

Application to a continuous pricing problem

- Continuous pricing problem:

$$\max_{p, \omega, U, H} \sum_n \sum_r \sum_i \frac{1}{R} \theta_{in} p_i \omega_{inr}$$

s.t.

$$\sum_i \omega_{inr} = 1 \quad \forall n, r$$

$$H_{nr} = \sum_i U_{inr} \omega_{inr} \quad \forall n, r$$

$$H_{nr} \geq U_{inr} \quad \forall i, n, r$$

$$U_{inr} = \sum_{k \neq l} \beta_k x_{ink} + \beta_l p_i + \varepsilon_{inr} \quad \forall i, n, r$$

$$\omega \in \{0, 1\}$$

$$p, U, H \in \mathbb{R}$$

Application to a continuous pricing problem

N	R	obj-MILP	obj-D	Gap [%]	P-MILP	P-D	Gap [%]	T-MILP	T-D
20	50	216.407	209.196	3.33	28.475	30.764	-8.04	7	11
20	100	202.642	201.712	0.46	28.302	26.576	6.1	37	21
20	200	200.901	200.185	0.36	30.03	28.721	4.36	205	49
50	50	440.686	437.243	0.78	28.579	29.989	-4.94	55	27
50	100	431.088	426.669	1.03	28.99	27.778	4.18	241	62
50	200	429.605	429.108	0.12	28.574	28.655	-0.28	1022	163
100	50	990.026	988.732	0.13	29.118	28.944	0.6	252	31
100	100	977.606	976.149	0.15	30.099	29.925	0.58	1224	69
100	200	978.589	976.932	0.17	30.106	30.185	-0.26	3039	304
200	50	1906.696	1904.189	0.13	28.977	28.678	1.03	1144	65
200	100	1882.793	1877.641	0.27	29.277	30.052	-2.65	4104	359
200	200	1873.964	1871.614	0.13	29.276	29.343	-0.23	10811	690

Large number of draws (MSLE)

N	R	sLL-M	sLL-D	Gap [%]	T-M	T-D
50	20	-29.417	-29.908	1.67	22	6
50	50	-29.294	-31.173	6.41	279	26
50	100	-28.885	-29.42	1.85	1375	42
50	150	-29.973	-30.092	0.4	2852	70
50	200	-30.091	-30.101	0.03	10800	131
50	250	-30.741	-30.775	0.11	10800	156
50	300	-30.837	-30.843	0.02	10800	133
50	400	-30.632	-30.638	0.02	10800	130
50	600	-30.479	-30.51	0.1	10800	289
50	800		-32.035		10800	319
50	1000		-30.523		10800	349

Ideas for future work

- Improving Benders:
 - Piece-wise linearization
 - Convex-quadratic formulation
- Column generation methods
- Combined column generation + Benders approach

Thanks!

