

# Simulation based Population Synthesis

Michel Berlaire, Bilal Farooq,  
Ricardo Hurtubia, and Gunnar Flötteröd

Transport and Mobility Laboratory  
School of Architecture, Civil and Environmental Engineering  
Ecole Polytechnique Fédérale de Lausanne

February 22, 2013



# Outline

- 1 Motivation
- 2 Introduction to existing literature
- 3 New methodology
- 4 Comparative experiments
- 5 Back to original problem
- 6 Concluding remarks



# Synthetic population



- Agent-based simulation
- Disaggregate models
- Need to access characteristics of agents

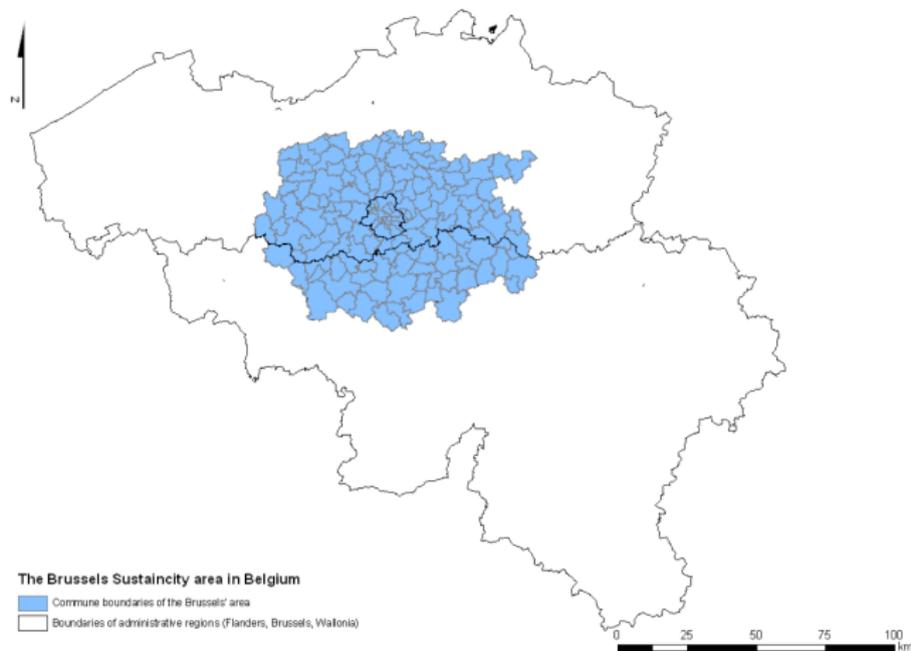
# SustainCity project

## European Union project

- More than 10 major European universities involved
- Aims:
  - Integrated Land Use and Transportation modelling framework
  - Demographics, environment, and multi-scale issues
- Tools: Urbansim, MATSim, Metropolis
- Case studies
  - Paris
  - Zurich
  - Brussels



# SustainCity: Brussels case study



# Brussels case study

## Data sources (extremely limited)

- Incomplete conditionals of households and persons (Census 2001)
- Travel survey of households and individuals (MOBEL 1999)
  - **3063 observations (0.2%)**

## Synthetic household characteristics

- Size, children, workers, cars, income, university education, dwelling type, sector
- *Conventional synthesis procedures were not usable*



# Agent synthesis in transportation

## From *Four-Stage* to *Activity based Integrated* modelling

- Urban Microsimulation

## Forecasting behavior using individual level models

- Lack of individual level data for population
- Synthesis of individual agents and their characteristics

## Initial work

- TRansportation ANalysis SIMulation System (TRANSIMS) [LaRon et al., 1996]
- Focused on synthesis of a small sub-set of characteristics

# Current needs and trends (MATSim, ILUTE, etc.)

## More detailed characteristics

- e.g. social network, dwelling location, employment type

## Associations between different types of agents

- Household, family, person

## Using variety of different data sources at different spatial aggregation

- From samples to aggregate statistics
- Sources: Census, travel survey, household spending survey, labor force survey, statistics from revenue agency, real estate cadaster, ...
- Space: traffic analysis zone, census sector, commune, municipality, ...

# Existing approach

- **Fitting based approach**
  - Iterative proportional fitting
  - Combinatorial optimization
- Adjusting sample weights to fit the aggregate statistics



# Iterative Proportional Fitting (IPF) [Beckman et al., 1996]

## Contingency Table (CT) from sample

- Categorization of variables of interest
- Totals for each cell of the resulting multi-way table

## Fitting

- Sample used to initialize the contingency table
- Use marginal as dimensional totals
- Adjust the cell probabilities to fit dimension totals
- Iterate while the error is large
- Odd-ratio is maintained

## Generation of agents based on fitted weights

- Monte Carlo simulation for fractions

# Iterative Proportional Fitting

Most widely used in transportation literature

## Historic improvements in IPF

Zero-cell issue, dimensionality, fractions in CT, zone-by-zone vs entire area, associations, sample-less, ...

## Literature

Beckman et al. (1996), Frick and Axhausen (2004), Arentze et al. (2007), Guo and Bhat (2007), Pritchard and Miller (2009), Ye (2009) IPU, Auld et al. (2010), Barthelemy and Toint (2012), Müller and Axhausen (2012) HIPF, ...



# Combinatorial Optimization (CO) [Williamson, 1998]

- Zone-by-zone
- 0-1 weights for each row in the sample
- Optimizing the weights to fit zonal marginals
- Use of hill-climbing, simulated annealing, and genetic algorithm to estimate the best set of obs. weights for each zone



# Key issues



## Optimization resulting in one synthetic population

- Data are incomplete and purposely tampered with sophisticated anonymizing techniques
- There can be any number of solutions

## Cloning of data

- Amplification of errors
- Lack of heterogeneity

## Focus on fitting marginals

- No emphasis on the correlation structure



# Key issues



**SMALL SAMPLE SIZE**

- Over reliance on the accuracy of the microdata, without serious consideration to the sampling process and assumptions
- Large enough sample size
- Inefficient use of the available data
- Discrete agent characteristics only
- Scalability issues

# Problem statement

## True population

- Individual agents defined as a set of characteristics  
 $X = (X^1, X^2, \dots, X^n)$
- Discrete (e.g. marital status) or continuous (e.g. income)
- Unique joint distribution represented by  $\pi_X(x)$

## Complex distribution

- No direct access to  $\pi_X(x)$
- Hard to draw from

## Partial views of $\pi_X(x)$

- Marginals, conditional-marginals, and samples

# Problem statement

Develop a synthesis procedure that...

- draws a synthetic population
- uses the partial views as if we were drawing from  $\pi_X(x)$
- generates an empirical distribution  $\pi_{\hat{X}}(\hat{x})$  of  $\hat{X}$  close to  $\pi_X(x)$



# Simulation based approach

## Gibbs sampling

- Markov Chain Monte Carlo [Geman & Geman, 1984]
- draws from complex distributions:  $\pi_X(x)$
- exploits conditionals

$$\pi(X^i | X^j = x^j, \text{ for } j = 1 \dots n \ \& \ i \neq j) = \pi(X^i | X^{-i}), \text{ for } i = 1, \dots, n$$



# Key challenges

- Preparation of the conditional distributions for characteristics from available data sources
- Full-conditionals rarely available



# Completing conditionals by assumptions

Example: (Age | Sex, Income)

- From data only (Age | Income) available
- Assume that for all values of Sex:

$$(Age | Sex, Income) = (Age | Income)$$

- No matter the Sex of a person, Age is only dependent on Income



# Completing conditionals by assumptions

Required

$$\pi(\mathcal{X}^1 | \mathcal{X}^{-1}) = \pi(\mathcal{X}^1 | \mathcal{X}^{(2\dots k)}, \mathcal{X}^{((k+1)\dots n)})$$

Available

$$\pi(\mathcal{X}^1 | \mathcal{X}^{(2\dots k)})$$

Assumption

$$\pi(\mathcal{X}^1 | \mathcal{X}^{-1}) = \pi(\mathcal{X}^1 | \mathcal{X}^{(2\dots k)}), \forall \mathcal{X}^{((k+1)\dots n)}$$

Worst case

$$\pi(\mathcal{X}^1 | \mathcal{X}^{-1}) = \pi(\mathcal{X}^1)$$

# Completing conditionals by domain knowledge

## Example: (Income | Sex, Age)

- From data only (Income | Sex) available
- Known: Infants do not have income, students have low income
  - (Income | Sex, Age) =  $\alpha$  (Income | Sex) for Age = 1...12
  - (Income | Sex, Age) =  $\beta$  (Income | Sex) for Age = 13...18
  - (Income | Sex, Age) =  $\gamma$  (Income | Sex) for Age > 18
  - $\alpha + \beta + \gamma = 1$  and  $\alpha < \beta < \gamma$

# Completing conditionals by domain knowledge

## Expert's assumptions

$$\pi(X^1 | X^{(2\dots k)}, X^{((k+1)\dots n)} = a) = \pi^a(X^1 | X^{(2\dots k)}),$$

$$\pi(X^1 | X^{(2\dots k)}, X^{((k+1)\dots n)} = b) = \pi^b(X^1 | X^{(2\dots k)}),$$

...



# Completing conditionals by parametric models

Example: (Dwelling | Income , Sex, Age)

- In sample
  - Characteristics of person  $n$  (Dwelling , Age, Sex) $_n$
  - Attributes of zone ( $z$ )
- Dwelling choice model can be estimated
  - Choice set:  $dwel\_typ = (attached, semidetached, detached, apartment)$
  - Utility:
 
$$V_{(n,z)}^i = ASC^i + \beta_{age_n}^i \times Age + \beta_{av\_inc_z}^i \times av\_inc_z + interactions + \dots$$
  - Model: logit.

# Completing conditionals by parametric models

Logit model

$$\pi(X_l^1 | X_m^{-1}) = \frac{e^{(V_{X_l^1} | X_m^{-1})}}{\sum_{p=1}^L (e^{(V_{X_p^1} | X_m^{-1})})}$$



# Population from Swiss Census

- Access to Swiss Census for 2000
  - Person and household characteristics (Except for Income)
- Selected area: postal code in Lausanne
  - CH-1004
  - 28,533 persons
- Four person characteristics (384 combinations)
  - Age (<15, 15-24, 25-34, 35-44, 45-54, 55-64, 65-74, >74)
  - Sex (Female, Male)
  - Household size (1, 2, 3, 4, 5, 6 or more)
  - Education level (none, primary, secondary, university/college)



# Comparison between IPF and Simulation

- Criteria: how well the joint distribution is reproduced?



# Data preparation

- Prepared same type of datasets as commonly available
  - Individual level microsample
    - Drawing from Census: Uniformly, without replacement
    - No sampling-zero
  - Zonal level conditionals (with various level of completion)
    - By counting from Census



# List of available sample sizes

| No. | Sample Size |
|-----|-------------|
| 1   | 20%         |
| 2   | 10%         |
| 3   | 5%          |
| 4   | 3%          |
| 5   | 1%          |

# List of available sample sizes

| No. | Sample Size |
|-----|-------------|
| 1   | 20%         |
| 2   | 10%         |
| 3   | 5%          |
| 4   | 3%          |
| 5   | 1%          |

- In practice the sample size is 5% or less
- Larger sizes used to investigate representativeness



## List of available conditionals

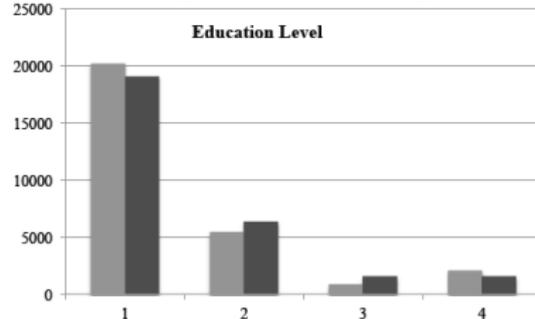
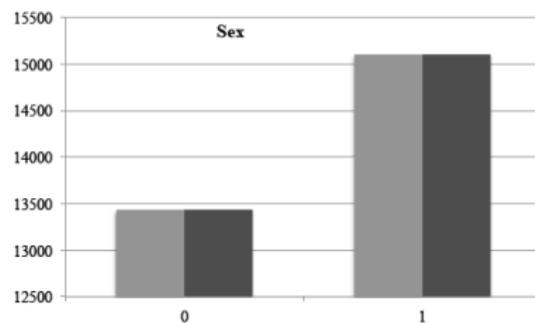
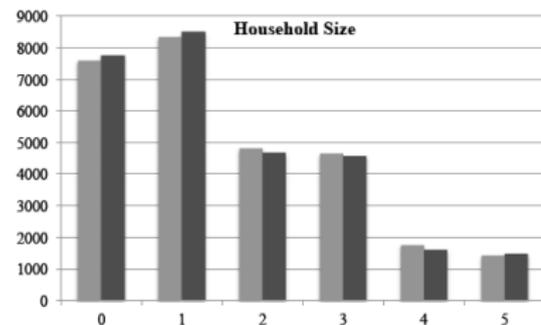
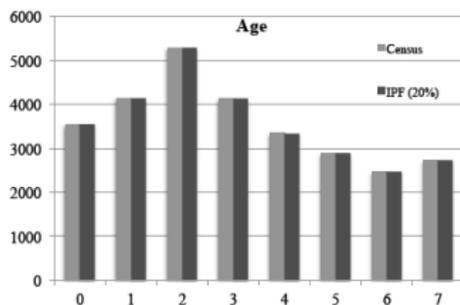
| No. | ID               | Conditionals   |
|-----|------------------|--|
| 1   | <i>FullCond</i>  | $\pi(\text{age} \text{sex}, \text{hhld\_size}, \text{edu\_level})$<br>$\pi(\text{sex} \text{age}, \text{hhld\_size}, \text{edu\_level})$<br>$\pi(\text{hhld\_size} \text{age}, \text{sex}, \text{edu\_level})$<br>$\pi(\text{edu\_level} \text{age}, \text{sex}, \text{hhld\_size})$ |
| 2   | <i>Partial_1</i> | $\pi(\text{age} \text{sex}, \text{hhld\_size}, \text{edu\_level})$<br>$\pi(\text{sex} \text{age}, \text{hhld\_size}, \text{edu\_level})$<br>$\pi(\text{hhld\_size} \text{age}, \text{sex}, \text{edu\_level})$<br>$\pi(\text{edu\_level} \text{age}, \text{sex}, \text{hhld\_size})$ |
| 3   | <i>Partial_2</i> | $\pi(\text{age} \text{sex}, \text{hhld\_size}, \text{edu\_level})$<br>$\pi(\text{sex} \text{age}, \text{hhld\_size}, \text{edu\_level})$<br>$\pi(\text{hhld\_size} \text{age}, \text{sex}, \text{edu\_level})$<br>$\pi(\text{edu\_level} \text{age}, \text{sex}, \text{hhld\_size})$ |
| 4   | <i>Partial_3</i> | $\pi(\text{age} \text{sex}, \text{hhld\_size}, \text{edu\_level})$<br>$\pi(\text{sex} \text{age}, \text{hhld\_size}, \text{edu\_level})$<br>$\pi(\text{hhld\_size} \text{age}, \text{sex}, \text{edu\_level})$<br>$\pi(\text{edu\_level} \text{age}, \text{sex}, \text{hhld\_size})$ |

# Data preparation

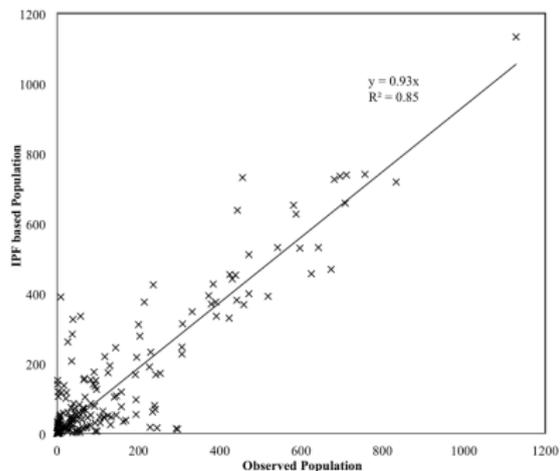
- Based on sample-conditional combinations
  - 20 possibilities
- IPF can use marginals only
  - Number of experiments collapses to 5
- Simulation based synthesis
  - Used conditionals only (used lesser information)
  - Number of experiments collapses to 4



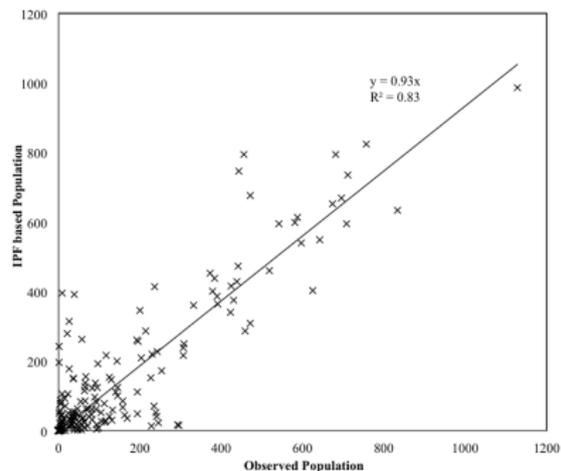
# Results: IPF and Census marginals



# Results: Fit of IPF with Census joint distribution



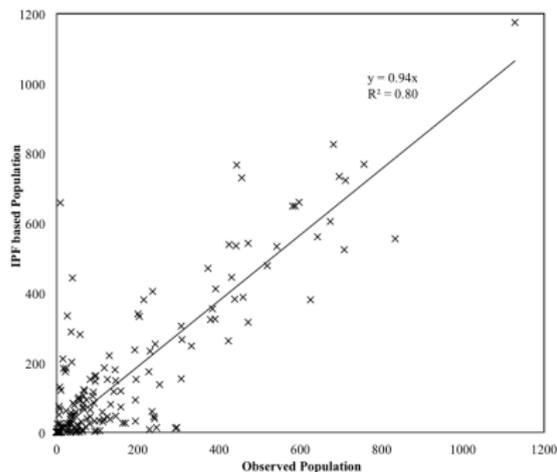
IPF with 20% sample



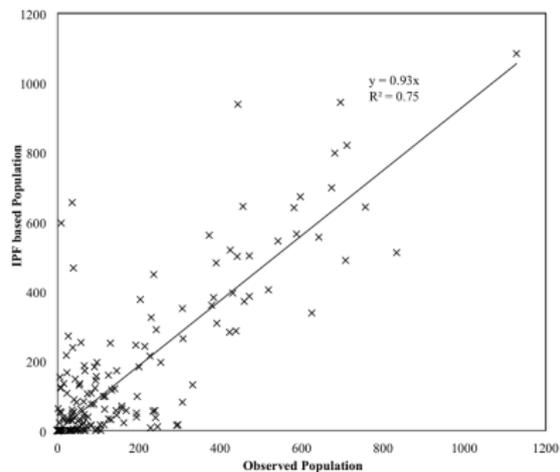
IPF with 10% sample



# Results: Fit of IPF with Census joint distribution



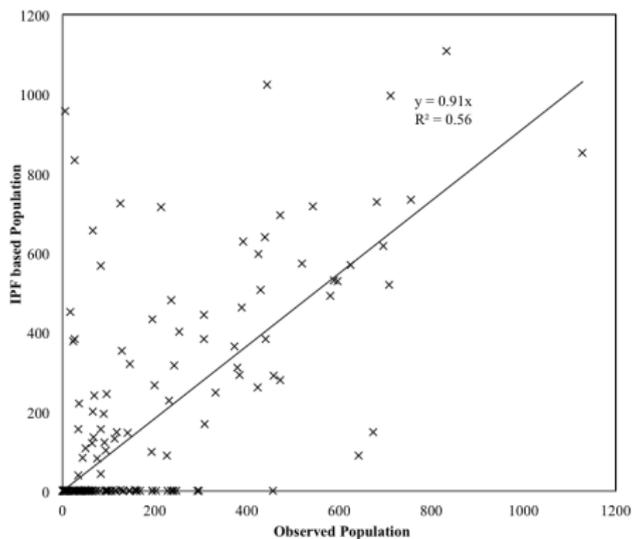
IPF with 5% sample



IPF with 3% sample



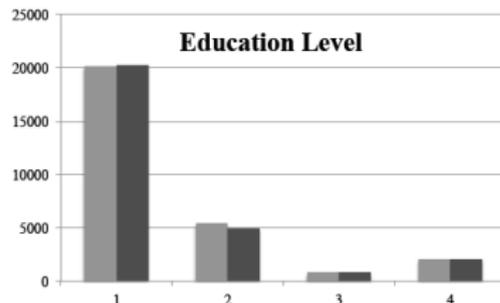
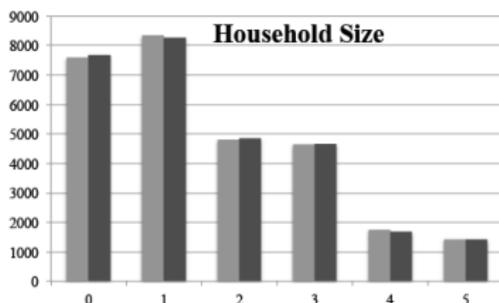
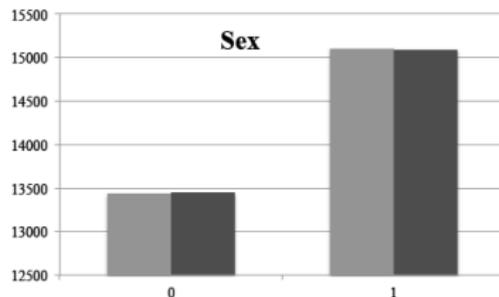
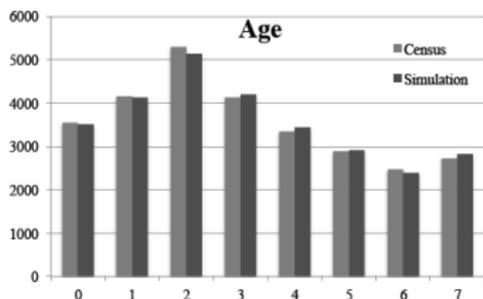
# Results: Fit of IPF with Census joint distribution



IPF with 1% sample

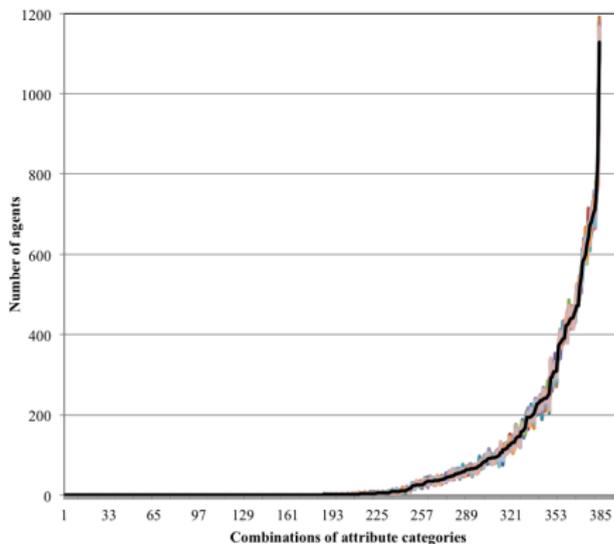


# Results: Simulation and Census marginals



Using full-conditionals (*FullCond*)

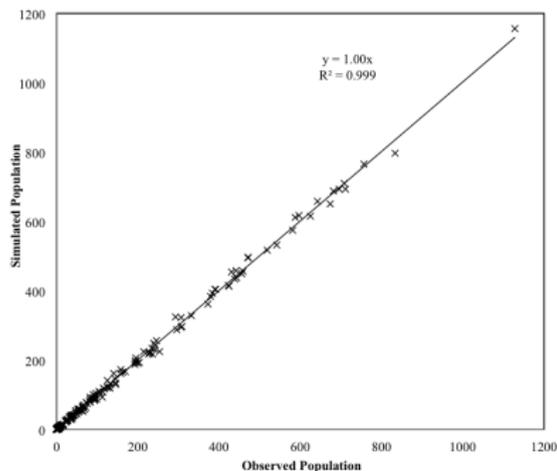
# Results: Simulation and Census joint dist.



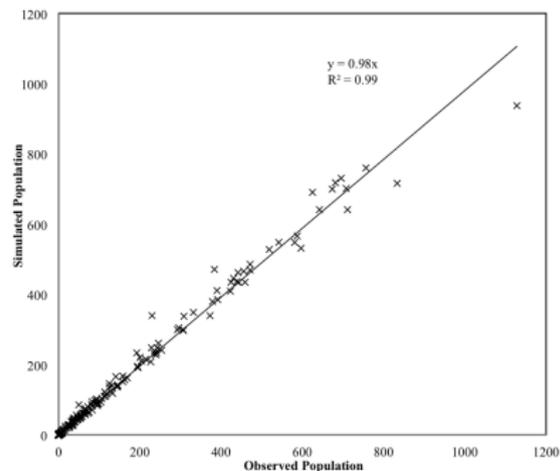
20 runs based on *FullCond* with real population superimposed



# Results: Fit of Simulation with Census joint dist.



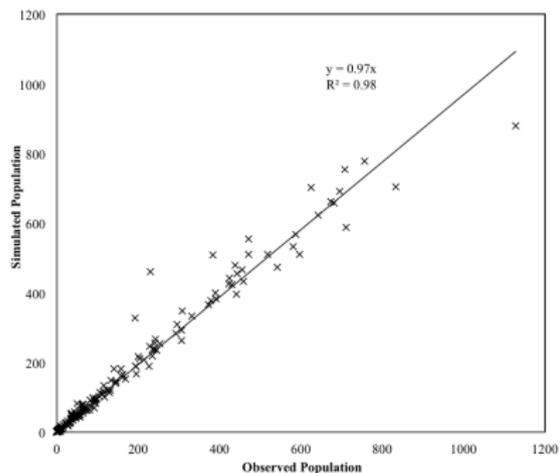
FullCond



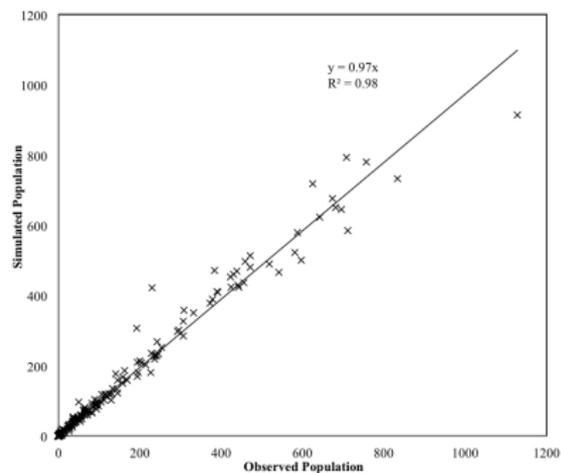
Partial\_1 (Sex missing in 1 conditional)



# Results: Fit of Simulation with Census joint dist.



*Partial.2 (Sex missing in 2 conditionals)*



*Partial.3 (Sex missing in all conditional)*

# Comparison: Standard Root Mean Square Error

$$SRSME = \frac{[\sum_{i=1}^m \dots \sum_{j=1}^n (R_{i\dots j} - T_{i\dots j})^2 / N]^{1/2}}{\sum_{i=1}^m \dots \sum_{j=1}^n (T_{i\dots j}) / N}$$



# Comparison: Standard Root Mean Square Error

$$SRSME = \frac{[\sum_{i=1}^m \dots \sum_{j=1}^n (R_{i\dots j} - T_{i\dots j})^2 / N]^{1/2}}{\sum_{i=1}^m \dots \sum_{j=1}^n (T_{i\dots j}) / N}$$

| Input     | IPF   | Simulation |
|-----------|-------|------------|
| 20%Sample | 0.853 | -          |
| 10%Sample | 0.928 | -          |
| 5%Sample  | 1.020 | -          |
| 3%Sample  | 1.160 | -          |
| 1%Sample  | 1.730 | -          |
| FullCond  | -     | 0.130      |
| Partial_1 | -     | 0.240      |
| Partial_2 | -     | 0.340      |
| Partial_3 | -     | 0.350      |

# Comparison: Standard Root Mean Square Error

$$SRSME = \frac{[\sum_{i=1}^m \dots \sum_{j=1}^n (R_{i\dots j} - T_{i\dots j})^2 / N]^{1/2}}{\sum_{i=1}^m \dots \sum_{j=1}^n (T_{i\dots j}) / N}$$

| Input     | IPF          | Simulation   |
|-----------|--------------|--------------|
| 20%Sample | <b>0.853</b> | -            |
| 10%Sample | 0.928        | -            |
| 5%Sample  | 1.020        | -            |
| 3%Sample  | 1.160        | -            |
| 1%Sample  | 1.730        | -            |
| FullCond  | -            | 0.130        |
| Partial_1 | -            | 0.240        |
| Partial_2 | -            | 0.340        |
| Partial_3 | -            | <b>0.350</b> |

# Comparison: Standard Root Mean Square Error

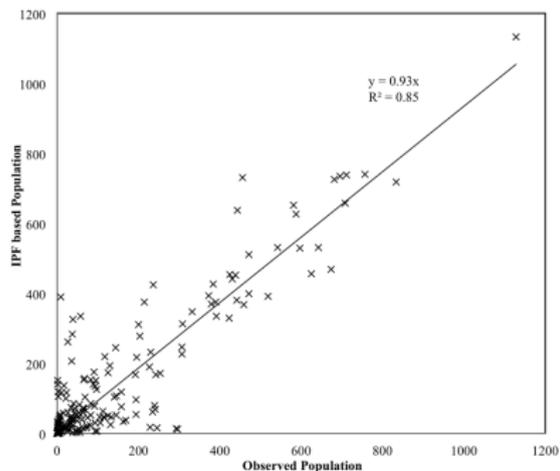
$$SRSME = \frac{[\sum_{i=1}^m \dots \sum_{j=1}^n (R_{i\dots j} - T_{i\dots j})^2 / N]^{1/2}}{\sum_{i=1}^m \dots \sum_{j=1}^n (T_{i\dots j}) / N}$$

| Input     | IPF          | Simulation   |
|-----------|--------------|--------------|
| 20%Sample | <b>0.853</b> | -            |
| 10%Sample | 0.928        | -            |
| 5%Sample  | 1.020        | -            |
| 3%Sample  | 1.160        | -            |
| 1%Sample  | 1.730        | -            |
| FullCond  | -            | 0.130        |
| Partial_1 | -            | 0.240        |
| Partial_2 | -            | 0.340        |
| Partial_3 | -            | <b>0.350</b> |

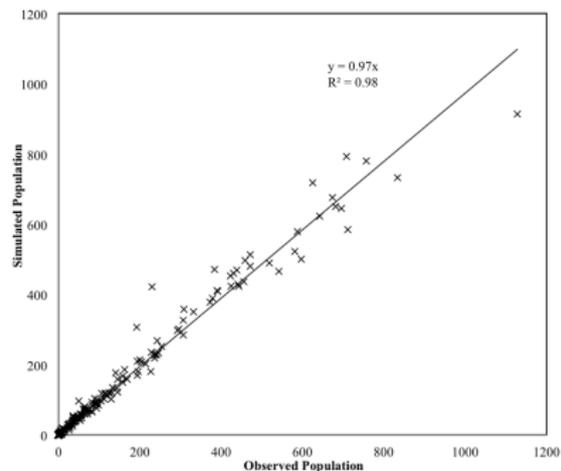
- For Marginals only, both methods give the same fit



# Best case IPF and worst case Simulation



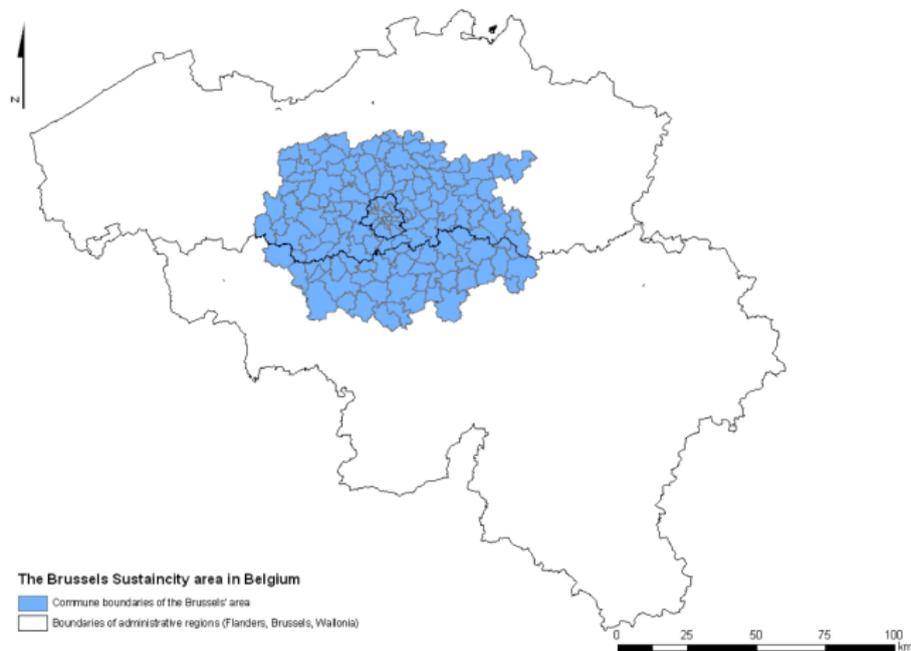
*IPF with 20% sample*



*Partial\_3 (Sex missing from all the conditionals)*



# Back to Brussels case study



# Brussels case study

- Data sources (extremely limited)
  - Incomplete conditionals of households and persons (Census 2001)
  - Travel survey of households and individuals (MOBEL 1999)
    - **3063 observations (0.2%)**
- Synthetic household attributes
  - Size, children, workers, cars, income, university education, dwelling type, sector



# Brussels case study

- Data sources (extremely limited)
    - Incomplete conditionals of households and persons (Census 2001)
    - Travel survey of households and individuals (MOBEL 1999)
      - **3063 observations (0.2%)**
  - Synthetic household attributes
    - Size, children, workers, cars, income, university education, dwelling type, sector
- 
- Data Preparation
    - Aggregation
      - Spatial
      - Categorical
    - Model based conditionals (Logit)
      - Income, univ edu, cars, and dwelling type

# Income level model (5 levels)

$$V_{(hh,z)}^1 = 0$$

$$V_{(hh,z)}^2 = ASC^2 + \beta_{zonal\_inc_z}^2 \times zonal\_inc_z + \beta_{cars_{hh}}^2 \times cars_{hh} + \beta_{workers_{hh}}^2 \times workers_{hh}$$

$$V_{(hh,z)}^3 = ASC^3 + \beta_{educ_{hh}}^3 \times educ_{hh} + \beta_{zonal\_inc_z}^3 \times zonal\_inc_z + \beta_{cars_{hh}}^3 \times cars_{hh} \\ + \beta_{house_{hh}}^3 \times house_{hh} + \beta_{workers_{hh}}^3 \times workers_{hh}$$

$$V_{(hh,z)}^4 = ASC^4 + \beta_{educ_{hh}}^4 \times educ_{hh} + \beta_{zonal\_inc_z}^4 \times zonal\_inc_z + \beta_{cars_{hh}}^4 \times cars_{hh} \\ + \beta_{house_{hh}}^4 \times house_{hh} + \beta_{workers_{hh}}^4 \times workers_{hh}$$

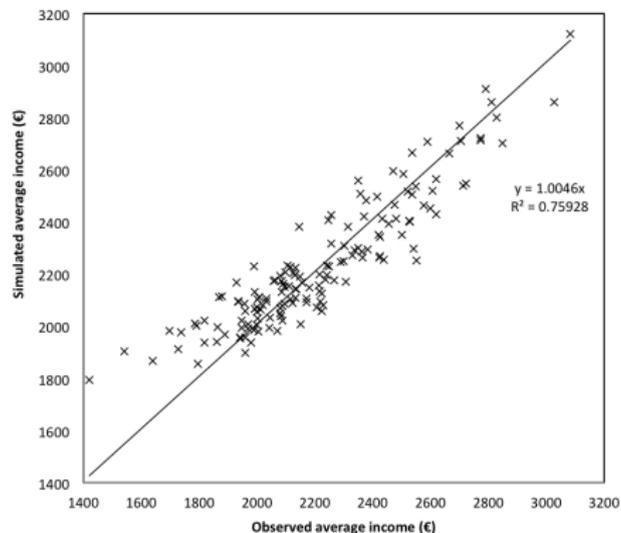
$$V_{(hh,z)}^5 = ASC^5 + \beta_{educ_{hh}}^5 \times educ_{hh} + \beta_{zonal\_inc_z}^5 \times zonal\_inc_z + \beta_{cars_{hh}}^5 \times cars_{hh} \\ + \beta_{house_{hh}}^5 \times house_{hh} + \beta_{workers_{hh}}^5 \times workers_{hh}$$



## Income level model

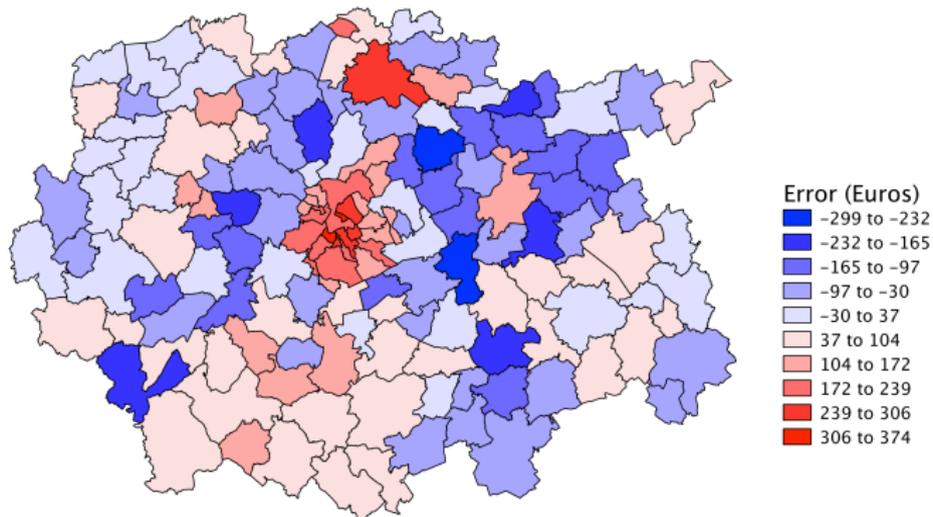
| Parameter              | Variable  | Value  | Std err | t-test |
|------------------------|---|--------|---------|--------|
| $ASC^2$                | constant for income level 2                             | -0.86  | 0.789   | -1.09  |
| $ASC^3$                | constant for income level 3                             | -4.64  | 0.901   | -5.14  |
| $ASC^4$                | constant for income level 4                             | -8.31  | 1.12    | -7.39  |
| $ASC^5$                | constant for income level 5                             | -10.6  | 1.55    | -6.82  |
| $\beta_{educ}^3$       | dummy for presence of people with higher educ in the hh | 0.831  | 0.177   | 4.69   |
| $\beta_{educ}^4$       | dummy for presence of people with higher educ in the hh | 1.72   | 0.314   | 5.49   |
| $\beta_{educ}^5$       | dummy for presence of people with higher educ in the hh | 1.92   | 0.656   | 2.93   |
| $\beta_{zonal\_inc}^2$ | average zonal income                                    | 0.0008 | 0.0004  | 1.84   |
| $\beta_{zonal\_inc}^3$ | average zonal income                                    | 0.0012 | 0.0005  | 2.55   |
| $\beta_{zonal\_inc}^4$ | average zonal income                                    | 0.0016 | 0.0005  | 3.09   |
| $\beta_{zonal\_inc}^5$ | average zonal income                                    | 0.0016 | 0.0006  | 2.47   |
| $\beta_{cars}^2$       | number of cars in the household                         | 1.16   | 0.265   | 4.39   |
| $\beta_{cars}^3$       | number of cars in the household                         | 1.92   | 0.299   | 6.41   |
| $\beta_{cars}^4$       | number of cars in the household                         | 2.33   | 0.341   | 6.83   |
| $\beta_{cars}^5$       | number of cars in the household                         | 3.2    | 0.466   | 6.87   |
| $\beta_{house}^3$      | dummy for dwelling being a house                        | 0.45   | 0.193   | 2.34   |
| $\beta_{house}^4$      | dummy for dwelling being a house                        | 0.485  | 0.294   | 1.65   |
| $\beta_{house}^5$      | dummy for dwelling being a house                        | 0.485  | 0.294   | 1.65   |
| $\beta_{workers}^2$    | number of workers in the household                      | 1.14   | 0.277   | 4.11   |
| $\beta_{workers}^3$    | number of workers in the household                      | 2.22   | 0.295   | 7.53   |
| $\beta_{workers}^4$    | number of workers in the household                      | 2.46   | 0.345   | 7.13   |
| $\beta_{workers}^5$    | number of workers in the household                      | 1.74   | 0.428   | 4.07   |

# Results: Brussels case study



Fit between simulation based and observed average commune-level income

# Results: Brussels case study



Spatial distribution of error in average income

- More zonal level demographic statistics are required to further decrease the error

# Concluding remarks

- From single solution optimization problem to sampling from joint distribution
  - Output of Land Use and Transport models

$$O = \int_{p_{syn}} \text{microsim}(p_{syn}) f(p_{syn}) dp_{syn}.$$

- Focus on reproducing not just marginals, but the whole joint distribution
- Heterogeneous not cloned population
- Population synthesis as part of microsimulation
  - Sensitivity analysis in a coherent way



# Concluding remarks

- Separation of data preparation from agent generation
  - Data, models, assumptions
- Mix of sampling process can be utilized based on the situation
- Works both for continuous and discrete or mixture of conditionals
- Computationally efficient and scalable
  - Clean and simple
- Generic data fusion technique
- Technical report URL: [http://transp-or.epfl.ch/documents/technicalReports/FaroBierHurtFloe\\_PopSyn2013.pdf](http://transp-or.epfl.ch/documents/technicalReports/FaroBierHurtFloe_PopSyn2013.pdf)



# Acknowledgments

- This research is funded by
  - European Commission's Seventh Framework Programme
  - Swiss National Science Foundation
  - Danish Council of Strategic Research
- Many thanks to
  - *Tomáš Robenek* for help in developing Java version
  - *Sohrab Sahaleh* for data processing
  - *Lovisa Arnesson* for data processing and running experiments

