# Synthetic Population Projections and Unforeseen Events: Hybrid Simulator for Capturing Dynamics

Marija Kukic     Michel Bierlaire

05 February, 2024

TRANSP-OR

EPFL

# Outline

Hybrid Simulator for Capturing Dynamics

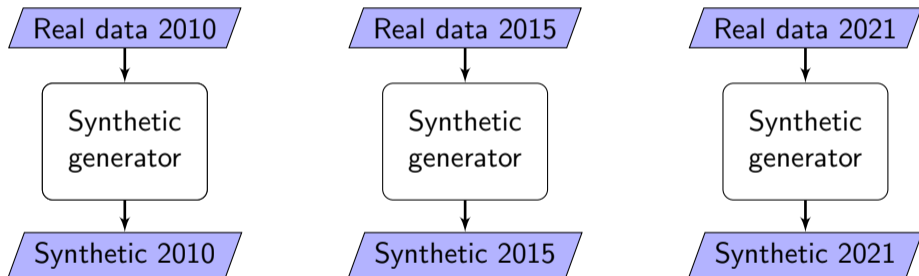# Synthetic Population in Transportation: Why?

### Real Data
- High cost of data collection.
- Lack of representativity.
- Data privacy constraints.
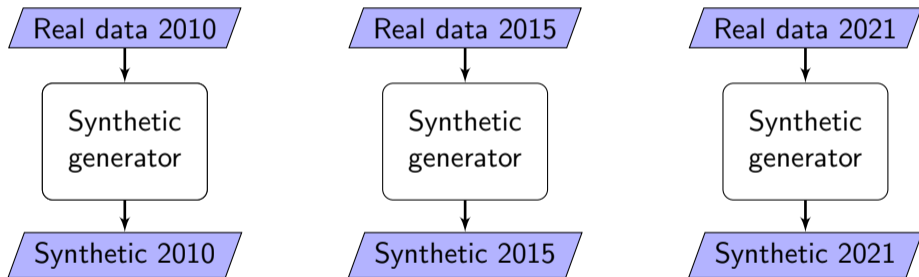
### Synthetic Data
- Open source.
- Bias correction.
- Privacy preservation.

**Synthetic Population** = tabular data on individuals and households
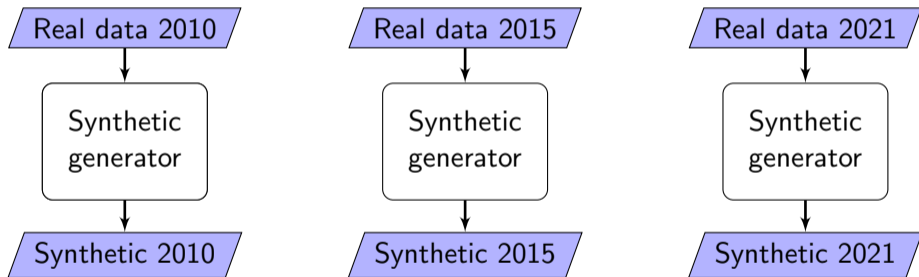
# Snapshot of Synthetic Population: Problems

# Snapshot of Synthetic Population: Problems



**Outdated sample**

# Snapshot of Synthetic Population: Problems
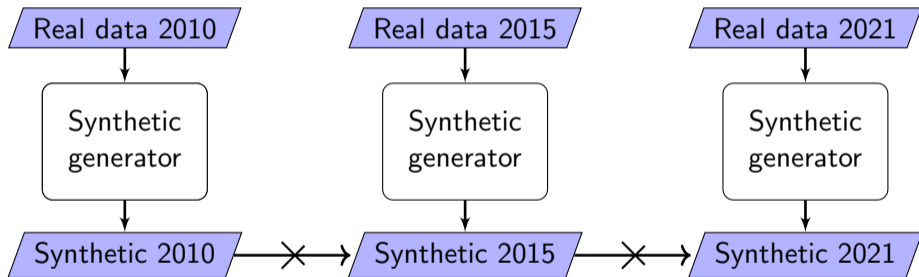
**Complicated and Repetitive**



**Outdated sample**

# Snapshot of Synthetic Population: Problems

**Complicated and Repetitive**



**Outdated sample**                    **Costly**

# Synthetic Population Projections



**Step 1: Generate**

**Step 2: Project**

# Outline

# Literature review - Generation and Projection

|  | Dynamic projection | Static projection | Resampling |
|---|---|---|---|
| **Synthetic reconstruction** | Fatmi et al.[1] **2017** | Lomax et al.[2] **2022** | Prédhumeau et al.[3] **2023** |
| **Combinatorial optimisation** | Namazi-Rad et el.[4] **2014** | **X** | **X** |
| **Statistical learning** | **Hybrid Simulator for Capturing Dynamics Model-driven** | **X** | **Hybrid Simulator for Capturing Dynamics Data-driven** |

# Literature Gaps

**Dynamic projection**

- Evolves population.
- Heterogeneous sample.

**Re-sampling**

- Copying data instead of evolving.
- Lack of heterogenity over time.

## Literature Gaps

### Dynamic projection

- Evolves population.
- Heterogeneous sample.

<br>

- Propagation of the generation bias.
- Increase of the error over time.
- **Not robust to the unusual events.**
- **Dependent on input rates.**

### Re-sampling

- Copying of data instead of evolving.
- Lack of heterogeneity over time.

<br>

- Can achieve a perfect fit.

# Literature Gaps

**Dynamic projection requires demographic rates to simulate events!**

Demographers provide reports on **demographic rates** every five years.

**Assumption:** Population trends remain **stable** over time.

**Problems**

Although **rates** are frequently **updated**, **synthetic datasets** made using them **are not**.

**Unforeseen events** can result in projections that **do not represent** the real population.

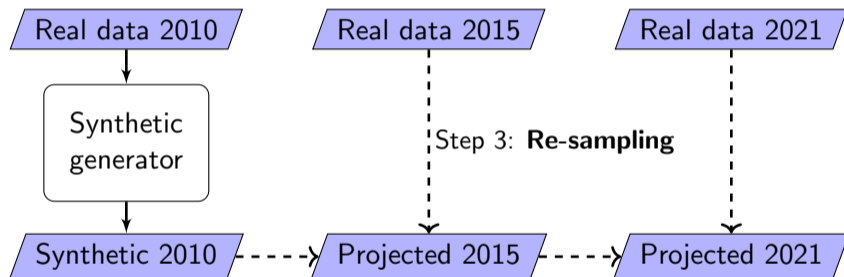Affects the **outcomes of transportation models** employing these samples.

Problematic for **long-term** forecast.

# Outline

# Contribution - Previous work

- Combine **dynamic projection** and **resampling** at the level of **individuals** [5].



Step 1: **Generation**   Step 2: **Dynamic Projection**   Step 4: **Validation**

**Model-based and Data-driven approach**

# Contribution - Previous work

- Simulate effect of **death, birth, migration** to synthetic individuals described by **age, gender, employment**.

**What we showed?**
- Maintenance of synthetic samples without regenerating.
- Access to up-to-date data and making use of the past.
- Trade-off between accuracy and efficiency.

# Contribution - Current work

- Expand the method from the level of **individuals** to the **household** level.
- Evaluate **robustness** of hybrid simulator to **unforeseen events** (i.e., COVID-19) compared to state-of-the-art methods.

# Outline

# Hybrid Simulator for Capturing Dynamics

## Step 1: Generation

Markov Chain Monte Carlo Simulation. [6]

Synthetic households of size $N$, $X = (X_{\text{type}}, X_{\text{nb\_cars}}, [individuals_i]_{i \in [1...N]})$.

Synthetic individual described by $X_{\text{age}}, X_{\text{gender}}, X_{\text{empl}}, X_{\text{marital}}, X_{\text{dl}}$.

Bootstrap and convergence monitoring.

## Step 2: Dynamic projection

When disaggregated data are not available.

Simulate events: **birth, death, migration, marriage, divorce, leaving the house**.

Use the rates provided by the Swiss Federal Office (BFS) [7].

# Hybrid Simulator for Capturing Dynamics

**Step 3: Re-sampling**

When disaggregated data are available.

Compare projected household-type marginals with real data.

**Add or delete** households to achieve desired fit.

**Step 4: Validation**

Compare marginal and sub-distributions with real data.

Statistics (e.g., SRMSE) and Visualization.

# Evaluate projections to unforeseen events

**Test two scenarios:**

**Pre-pandemic**: Using rates from the report from **2010 without knowing** about the pandemic.

**Post-pandemic:** Using rates from the report from **2021 knowing** about the pandemic.

**Goal:**

Compare **dynamic projection** and **hybrid simulator** for these two scenarios by projecting samples from **2010** to **2021**.

# Outline

# Generation of synthetic sample 2010 - Household level

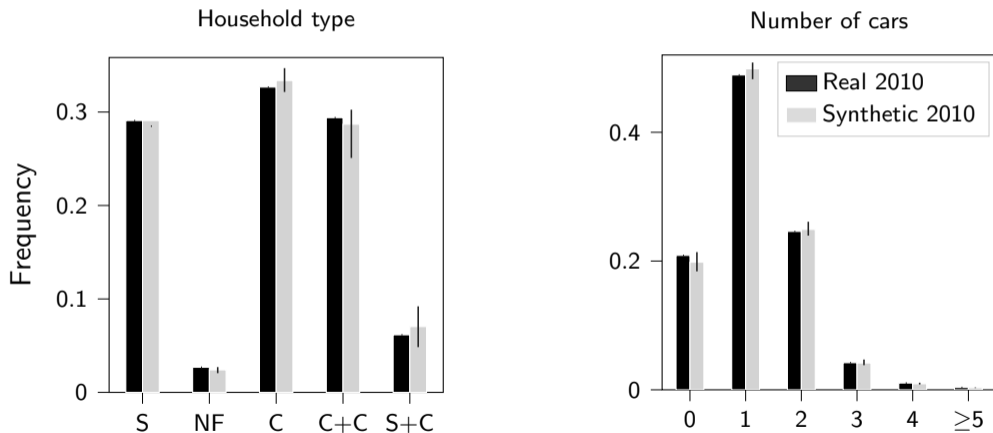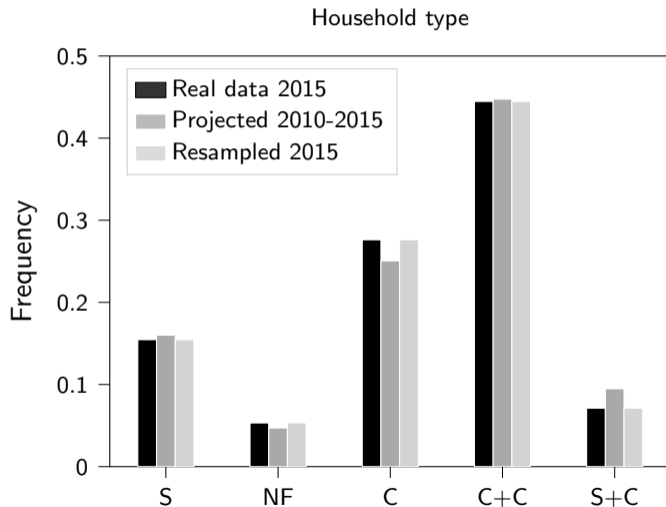Reference data: weighted **MTMC 2010, 2015, 2021** [BFS]



Figure: The comparison of household marginals between synthetic and real sample from 2010

# Dynamic Projection (2010-2015) and Re-sampling (2015)



Household type

# Comparison of Dynamic Projection and Hybrid Approach - 2021

| Variable | Pre-pandemic scenario | | Post-pandemic scenario | |
|---|---|---|---|---|
| | Dynamic projection | Hybrid simulator | Dynamic projection | Hybrid simulator |
| Household size | 0.22 | 0.15 | 0.19 | 0.12 |
| Household type | 0.24 | 0.1 | 0.15 | 0.08 |
| Number of cars | 0.32 | 0.18 | 0.24 | 0.12 |
| Age | 0.24 | 0.07 | 0.04 | 0.02 |
| Gender | 0.01 | 0.01 | 0.01 | 0.01 |
| Driving licence | 0.1 | 0.1 | 0.1 | 0.1 |
| Marital status | 0.07 | 0.06 | 0.07 | 0.06 |
| Employment | 0.26 | 0.25 | 0.16 | 0.15 |
| Average SRMSE | 0.18 | 0.11 | 0.12 | 0.08 |

# Comparison of Dynamic Projection and Hybrid Approach - 2021

The hybrid simulator achieved a better score (i.e., lower) for each attribute in both scenarios.

| Variable | Pre-pandemic scenario | | Post-pandemic scenario | |
|---|---|---|---|---|
| | Dynamic projection | Hybrid simulator | Dynamic projection | Hybrid simulator |
| Household size | 0.22 | **0.15** | 0.19 | **0.12** |
| Household type | 0.24 | **0.1** | 0.15 | **0.08** |
| Number of cars | 0.32 | **0.18** | 0.24 | **0.12** |
| Age | 0.24 | **0.07** | 0.04 | **0.02** |
| Gender | 0.01 | 0.01 | 0.01 | 0.01 |
| Driving licence | 0.1 | 0.1 | 0.1 | 0.1 |
| Marital status | 0.07 | **0.06** | 0.07 | **0.06** |
| Employment | 0.26 | **0.25** | 0.16 | **0.15** |
| Average SRMSE | 0.18 | **0.11** | 0.12 | **0.08** |

# Comparison of Dynamic Projection and Hybrid Approach - 2021

Some attributes are not affected by unforeseen events.

| Variable | Pre-pandemic scenario | | Post-pandemic scenario | |
|---|---|---|---|---|
| | Dynamic projection | Hybrid simulator | Dynamic projection | Hybrid simulator |
| Household size | 0.22 | 0.15 | 0.19 | 0.12 |
| Household type | 0.24 | 0.1 | 0.15 | 0.08 |
| Number of cars | 0.32 | 0.18 | 0.24 | 0.12 |
| Age | 0.24 | 0.07 | 0.04 | 0.02 |
| **Gender** | **0.01** | **0.01** | **0.01** | **0.01** |
| **Driving licence** | **0.1** | **0.1** | **0.1** | **0.1** |
| Marital status | 0.07 | 0.06 | 0.07 | 0.06 |
| Employment | 0.26 | 0.25 | 0.16 | 0.15 |
| Average SRMSE | 0.18 | 0.11 | 0.12 | 0.08 |

# Comparison of Dynamic Projection and Hybrid Approach - 2021

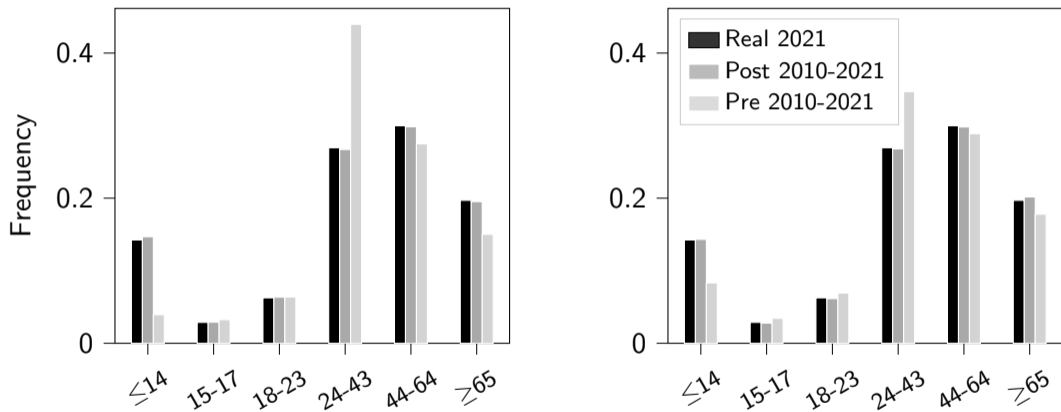Using updated rates leads to better results for both methods.

| Variable | Pre-pandemic scenario | | Post-pandemic scenario | |
|---|---|---|---|---|
| | Dynamic projection | Hybrid simulator | Dynamic projection | Hybrid simulator |
| Household size | 0.22 | 0.15 | 0.19 | 0.12 |
| Household type | 0.24 | 0.1 | 0.15 | 0.08 |
| Number of cars | 0.32 | 0.18 | 0.24 | 0.12 |
| Age | 0.24 | 0.07 | 0.04 | 0.02 |
| Gender | 0.01 | 0.01 | 0.01 | 0.01 |
| Driving licence | 0.1 | 0.1 | 0.1 | 0.1 |
| Marital status | 0.07 | 0.06 | 0.07 | 0.06 |
| Employment | 0.26 | 0.25 | 0.16 | 0.15 |
| **Average SRMSE** | 0.18 | 0.11 | **0.12** | **0.08** |

# Comparison of Dynamic Projection and Hybrid Approach - 2021

The difference between pre and post-pandemic scenarios is smaller for the hybrid simulator.

| Variable | Pre-pandemic scenario | | Post-pandemic scenario | |
|---|---|---|---|---|
| | Dynamic projection | Hybrid simulator | Dynamic projection | Hybrid simulator |
| Household size | 0.22 | 0.15 | 0.19 | 0.12 |
| Household type | 0.24 | 0.1 | 0.15 | 0.08 |
| Number of cars | 0.32 | 0.18 | 0.24 | 0.12 |
| Age | 0.24 | 0.07 | 0.04 | 0.02 |
| Gender | 0.01 | 0.01 | 0.01 | 0.01 |
| Driving licence | 0.1 | 0.1 | 0.1 | 0.1 |
| Marital status | 0.07 | 0.06 | 0.07 | 0.06 |
| Employment | 0.26 | 0.25 | 0.16 | 0.15 |
| **Average SRMSE** | 0.18 | **0.11** | 0.12 | **0.08** |

# Comparison of Dynamic Projection and Hybrid approach - 2021



Figure: Marginal distribution of the age using pre and post-pandemic rates compared to the real data - (left) dynamic projection; (right) hybrid simulator

# Outline

# Conclusion and Future Work

**We show:**

- Resampling step helps reduce the accumulated projection error of dynamic projection.
- The hybrid simulator is more robust to unforeseen events than the dynamic projection.
- The significance of validating and updating synthetic projected samples.

**Future work**

- How to model synthetic individuals over time using Gibbs Sampler?
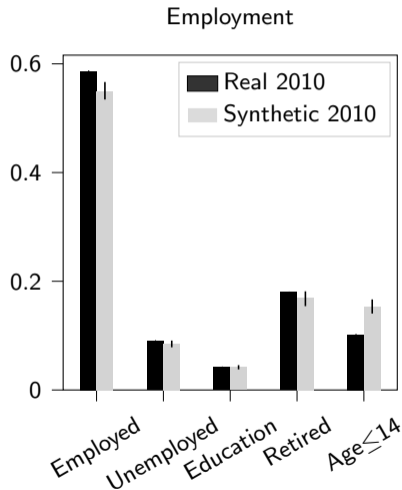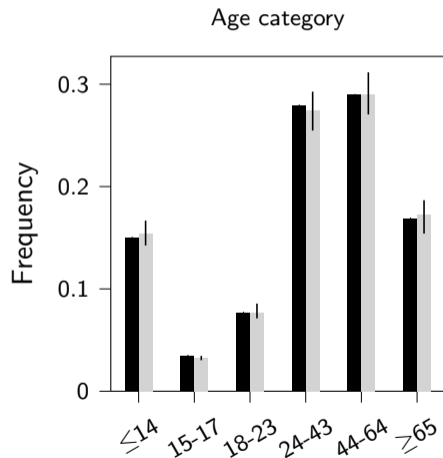
# Thank you! Questions?



Contact: **marija.kukic@epfl.ch**

**References**

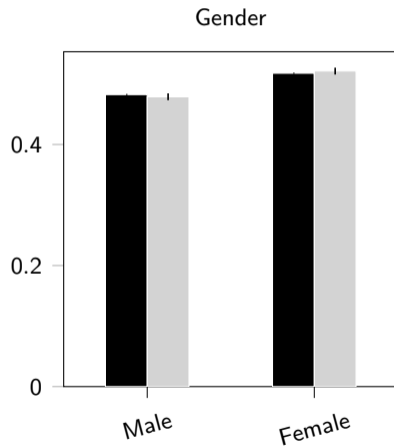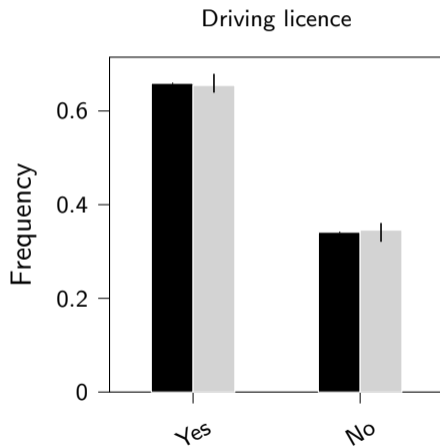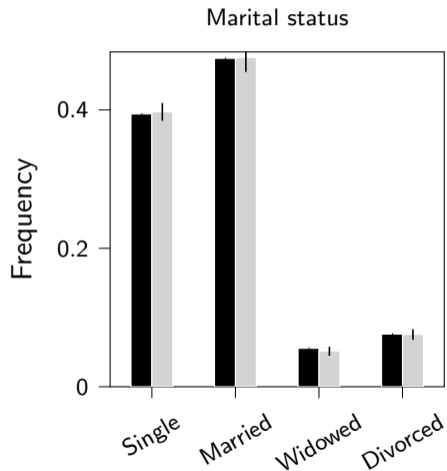# Backup slides

# Backup slides

# Backup slides



Marital status

# Backup slides - Dynamic projection - Births

---

**Algorithm** Births simulation

---

**Require:** $P$ - synthetic population
1: **for** $(a, m, o)$ in [age classes, marital statuses, birth orders] **do**
2:     Extract mother candidates $M$ in $P$ with attributes $(a, m, o)$
3:     Get the number $B$ of births with attributes $(a, m, o)$ {From BFS data}
4:     Draw $B$ mothers from $M$
5:     Add newborn in mothers' households
6: **end for**

---

# Backup slides - Dynamic projection - Migrations

**Algorithm** Migration simulation

1: $P$ - synthetic population
2: **for** $(a, g)$ in [ages, genders] **do**
3:    Get the net migration $N$ for attributes $(a, g)$ {From BFS data}
4:    **if** $N \geq 0$ **then**
5:       Draw $N$ individuals with attributes $(a, g)$ from $P$ {With replacement}
6:       Duplicate the $N$ individuals
7:       Build households from new individuals
8:    **else**
9:       Remove $N$ individuals with attributes $(a, g)$ from $P$
10:      Adapt modified households
11:   **end if**
12: **end for**

# Backup slides - Dynamic projection - Marriages

---

**Algorithm** Marriages simulation

---

1: $P$ - synthetic population
2: **for** $(h, w)$ in [husband ages, wife ages] **do**
3:    Get marriage count $N$ for attributes $(h, w)$ {From BFS data}
4:    Extract husband candidates $H$ from $P$
5:    Extract wife candidates $W$ from $P$
6:    Draw $N$ couples from product set $H \times W$
7:    Create new households for each couple
8:    Change couple marital status to "Married"
9:    Adapt modified households
10: **end for**

---

**Algorithm** Leaving the house simulation

1: $P$ - synthetic population
2: $r$ - official percentage of children in parental house
3: Extract individuals $C$ from $P$ with age in [15-28]
4: Extract individuals $C_{\text{parent}}$ from $C$ living in parental house
5: Compute the current percentage $r_{\text{cur}} = \frac{|C_{\text{parent}}|}{|C|}$
6: **if** $r_{\text{cur}} > r$ **then**
7:     $N \leftarrow \lfloor (r_{\text{cur}} - r) \cdot |P| \rfloor$
8:     Assign weights by age to $C_{\text{parent}}$
9:     Sample $N$ candidates from $C_{\text{parent}}$ with weights
10:     **for** each $c$ in candidates **do**
11:         **if** $c$ has children **then**
12:             Create a new house with type "Single-parent"
13:         **else**
14:             Create a new single household
15:         **end if**
16:     **end for**
17:     Adapt impacted household
18: **end if**

# Backup slides - Resampling

**Algorithm** Resampling procedure

1: **Input:**
2: *counts_real* - an array of frequency counts per household type in the reference sample
3: *counts_projected* - an array of frequency counts per household type in the projected sample
4: *list_of_types* - an array of existing household types
5: *num* - total number of household types
6: **Function** Resample(*counts_real, counts_projected, list_of_types, num, projected_sample*)
7: result_sample ← projected_sample
8: **for** $i \leftarrow 1$ **to** *num* **do**
9:     nb_of_observation ← abs(*counts_real*[i] - *counts_projected*[i])
10:     **if** *counts_real*[i] - *counts_projected*[i] $< 0$ **then**
11:         Delete *list_of_types*[i], nb_of_observation from result_sample
12:     **else**
13:         Add *list_of_types*[i], nb_of_observation to result_sample
14:     **end if**
15: **end for**
16: **end Function**