

Mobile Phone Data for Official Statistics: Elements for a Production Framework

BigSurv 2018

David Salgado (coord.) (INE, Spain)

Marc Debusschere (Statistics Belgium, Belgium)

Ossi Nurmi, Pasi Piela (Tilastokeskus, Finland)

Elise Coudin, Benjamin Sakarovitch (INSEE, France)

Sandra Hadam, Markus Zwick (DESTATIS, Germany)

Roberta Radini (ISTAT, Italy)

Martijn Tennekes (CBS, Netherlands)

Ciprian Alexandru, Bogdan Oancea (INSSE, Romania)

Elisa Esteban, Soledad Saldaña, Luis Sanguiao (INE, Spain)

Susan Williams (ONS, UK)

Barcelona, 24-26 October, 2018



Overview

1. Deliverable 5.1:

Current **status of access** to mobile phone data in the ESS (July, 2016)

2. Deliverable 5.2:

Guidelines for the access to mobile phone data within the ESS (May, 2017)

3. Deliverable 5.3:

Proposed elements for a **methodological framework** for the production of official statistics with mobile phone data (Feb, 2018)

4. Deliverable 5.4:

Some **IT elements** for the use of mobile phone data in the production of official statistics (March, 2018)

5. Deliverable 5.5:

Some **quality aspects and future prospects** for the production of official statistics with mobile phone data (May, 2018)

+ Meeting **Minutes** + **Github** Page

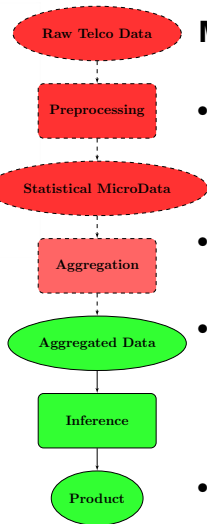
Big Data definition revisited

Big Data *for Official Statistics*

- refer to **third people** and not to data holders;
- are **central in their economic activity**;
- **lack statistical metadata** (since they are generated for very different purposes).

UNECE (wider) Definition for **Admin Data**:

“Data collected by sources **external to statistical offices.**”



MobPhone Data (WP5) Admin Data (Reid et al., 2017)

- Phase 1: Raw Telco Data Generation

- Phase 1: Admin Data Generation

- Phase 2: Statistical MicroData Generation

- Phase 2: Statistical MicroData Generation

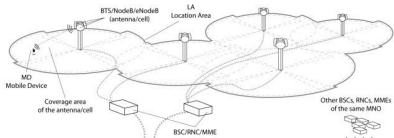
- Phase 3: Aggregated Data Generation

- Phase 3: Inference to Target Population

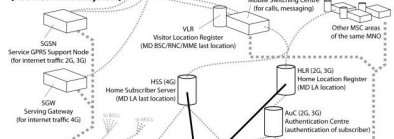
- More phases: Inference

Data generation: Phase 1

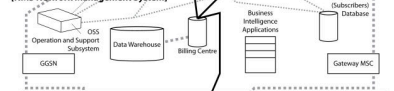
[BSS Base Station Subsystem]



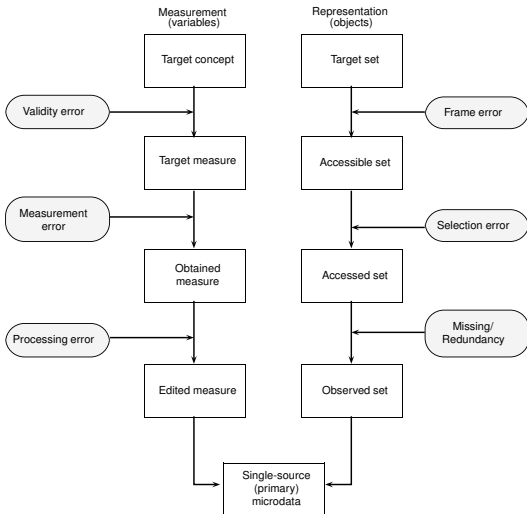
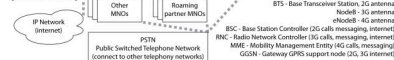
[NSS Network Subsystem]



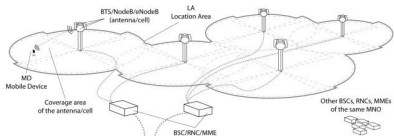
[NMS Network Management System]



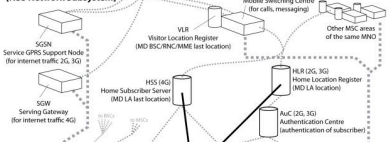
[Outside World]



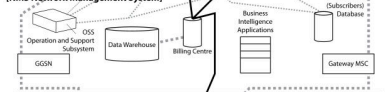
[BSS Base Station Subsystem]



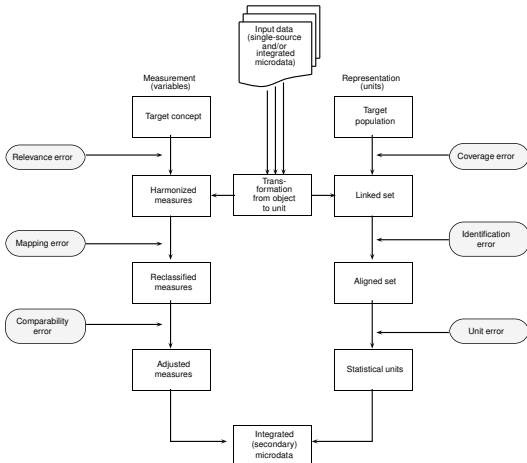
[NSS Network Subsystem]



[NMS Network Management System]



[Outside World]



From
pseudonymised ID, time attributes, space attributes (...)
to

COR
subscriber_id : BIGINT
from_month : DATE
to_month : DATE
iso_a2 : VARCHAR(2)

anchors
subscriber_id : BIGINT
lau3_code : INT
month : DATE
is_ptr : BOOL
is_wtap : BOOL
is_main_wtap : BOOL
is_shap : BOOL
is_trap : BOOL

usual environment
subscriber_id : BIGINT
month : DATE
lau3_code : INT
iso_a2 : VARCHAR(2)

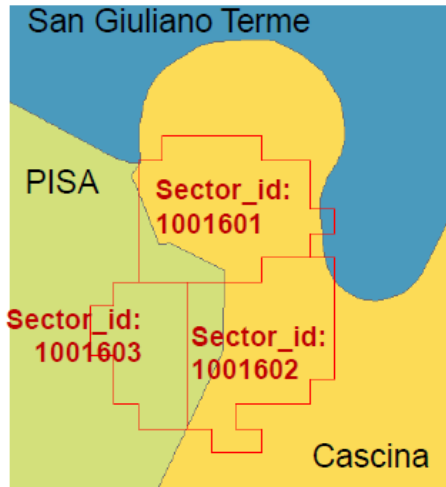
trips
subscriber_id : BIGINT
trip_id : BIGINT
from_time : TIMESTAMP
to_time : TIMESTAMP

stay sections
subscriber_id : BIGINT
stay_id : BIGINT
from_time : TIMESTAMP
to_time : TIMESTAMP
lau3_code : INT [NULL if outbound]
iso_a2 : VARCHAR(2) is_ptr : BOOL
is_wtap : BOOL
is_main_wtap : BOOL
is_shap : BOOL
is_trap : BOOL
trip_id : BIGINT

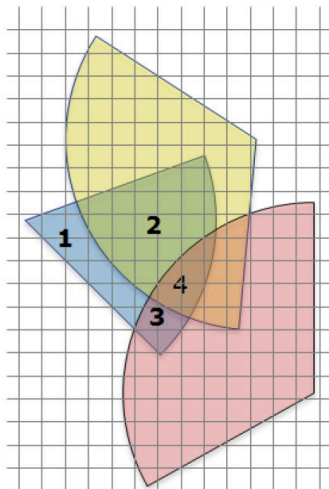
movement sections
subscriber_id : BIGINT
movement_id : BIGINT
from_lau3_code : INT [NULL if outbound]
to_lau3_code : INT [NULL if outbound]
from_time : TIMESTAMP
to_time : TIMESTAMP
from_iso_a2 : VARCHAR(2)
to_iso_a2 : VARCHAR(2)
from_stay_id : INT
to_stay_id : INT
travel_time : INT [in seconds]
intermediate_lau3_codes : INT ARRAY [NULL if outbound]
transport_mode : INT

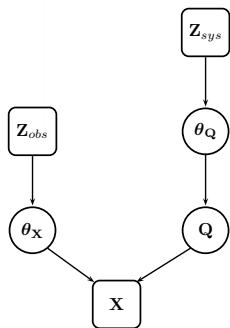
administrative units
lau_level : INT [0=country, 1=county, 2=municipality, 3=village,...]
lau_code : INT [national code of the LAU unit]
name : VARCHAR(255)
parent_lau_code : INT
geom : MULTIPOLYGON GEOMETRY

- **Best Service Area** Approach

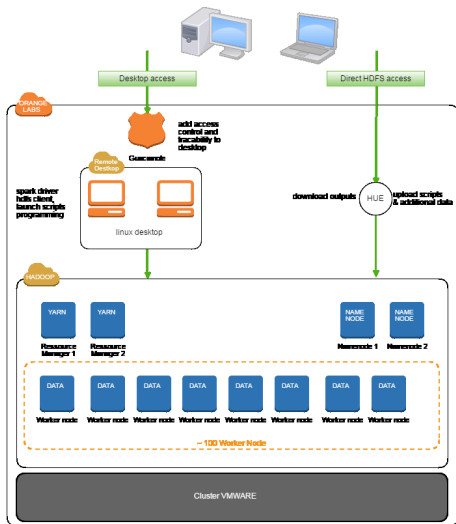


- **Bayesian Approach**





- **Probabilistic sampling** not possible: the **curse of representativity**
- Adaptation of **species abundance problem** in ecological sampling:
 - **Hierarchical** model
 - Two working assumptions:
 - At t_0 individuals are assumed to be physically in the **territorial cell of auxiliary admin/survey data**.
 - Mobility patterns of individuals **do not depend on the concrete MNO** they are subscribed to.
 - **Bayesian** approach:
 - **Computational** scalability.
 - **Integration** with other data sources.



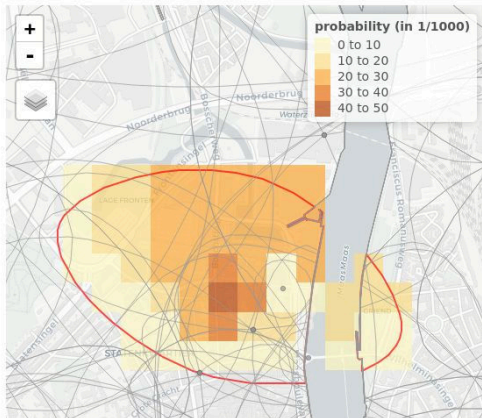
Cell Inspection Tool

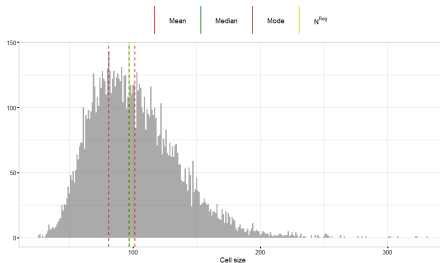
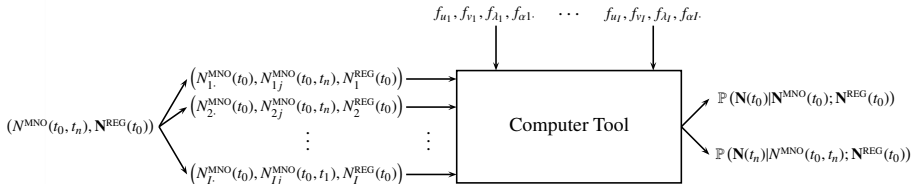
Variable

- Probability
- Signal strength (dB)

Selected cells

- Cell 1
- Cell 2
- Cell 3
- Cell 4
- Cell 5
- Cell 6
- Cell 7
- Cell 8
- Cell 9
- Cell 10
- Cell 11
- Cell 12
- Cell 13
- Cell 14





Quality issues

- **CoP** affected by two generic facts:
 - MNOs **active part of the production process**.
 - Change of **inferential paradigm**.
 - **Higher** degree of **breakdown**.

- Example: **accuracy** dimension

From **confidence** intervals to **credible** intervals

Model **checking** and model **assessment**

Main conclusions

- **Access blocked:** further work on **perceived risks** and **collaboration**.
- **Total Survey Error** paradigm still **valid**.
- **Geolocation**.
- **No probability sampling:** hierarchical **models?**
- **Computational complexity** for storing, accessing, and computing **in situ**.
- **Quality Assurance Framework** needs revision: **active role** of data holders and change of **inferential paradigm**.

