

Big Data Analytics

FOSCA GIANNOTTI AND LUCA PAPPALARDO

[HTTP://DIDAWIKI.DI.UNIPI.IT/DOKU.PHP/BIGDATAANALYTICS/BDA/](http://didawiki.di.unipi.it/doku.php/bigdataanalytics/bda/)

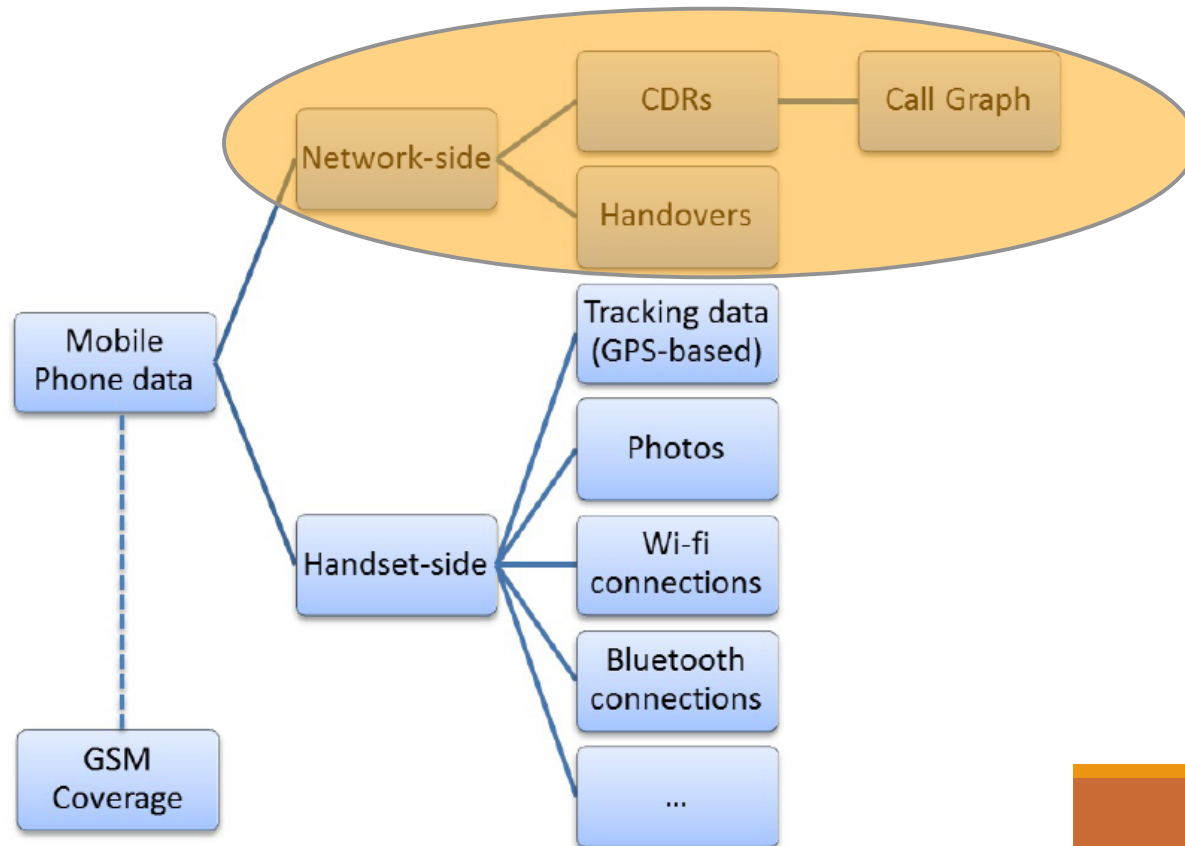
DIPARTIMENTO DI INFORMATICA - Università di Pisa
anno accademico 2018/2019

Mobility Data Mining

CITY DYNAMICS WITH GSM DATA

What are GSM data

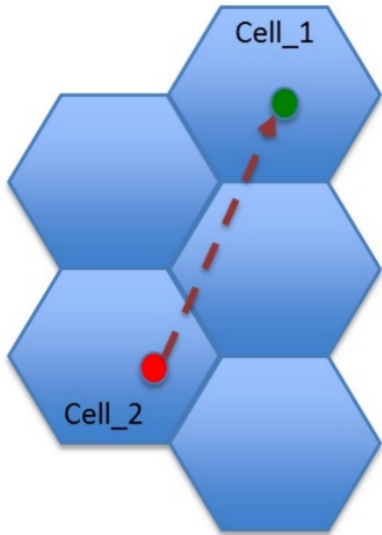
- Most popular resource for mobile phone data
- In principle, several kinds of data



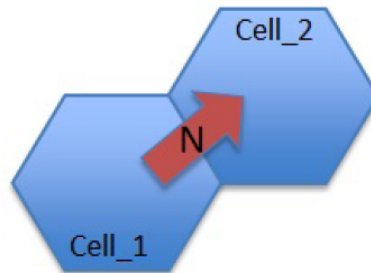
GSM data types

CDR

Who calls, **where** and **when**

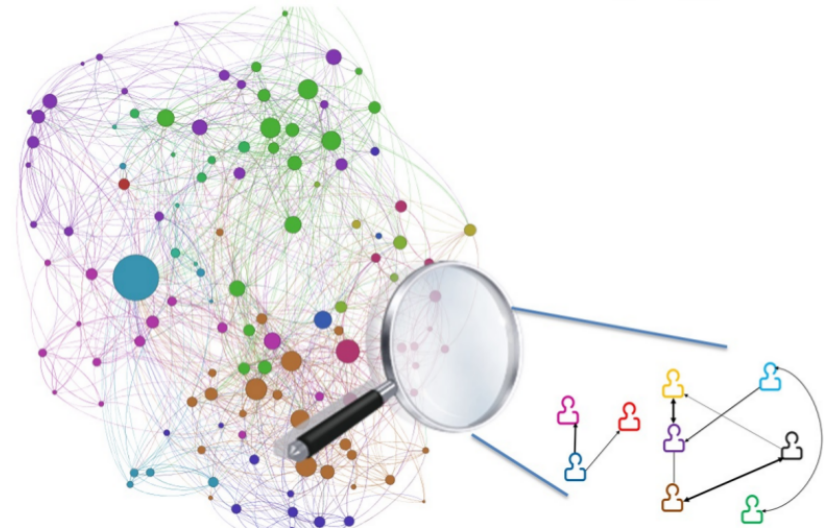


Hand over
Inter-cell flow **counts**



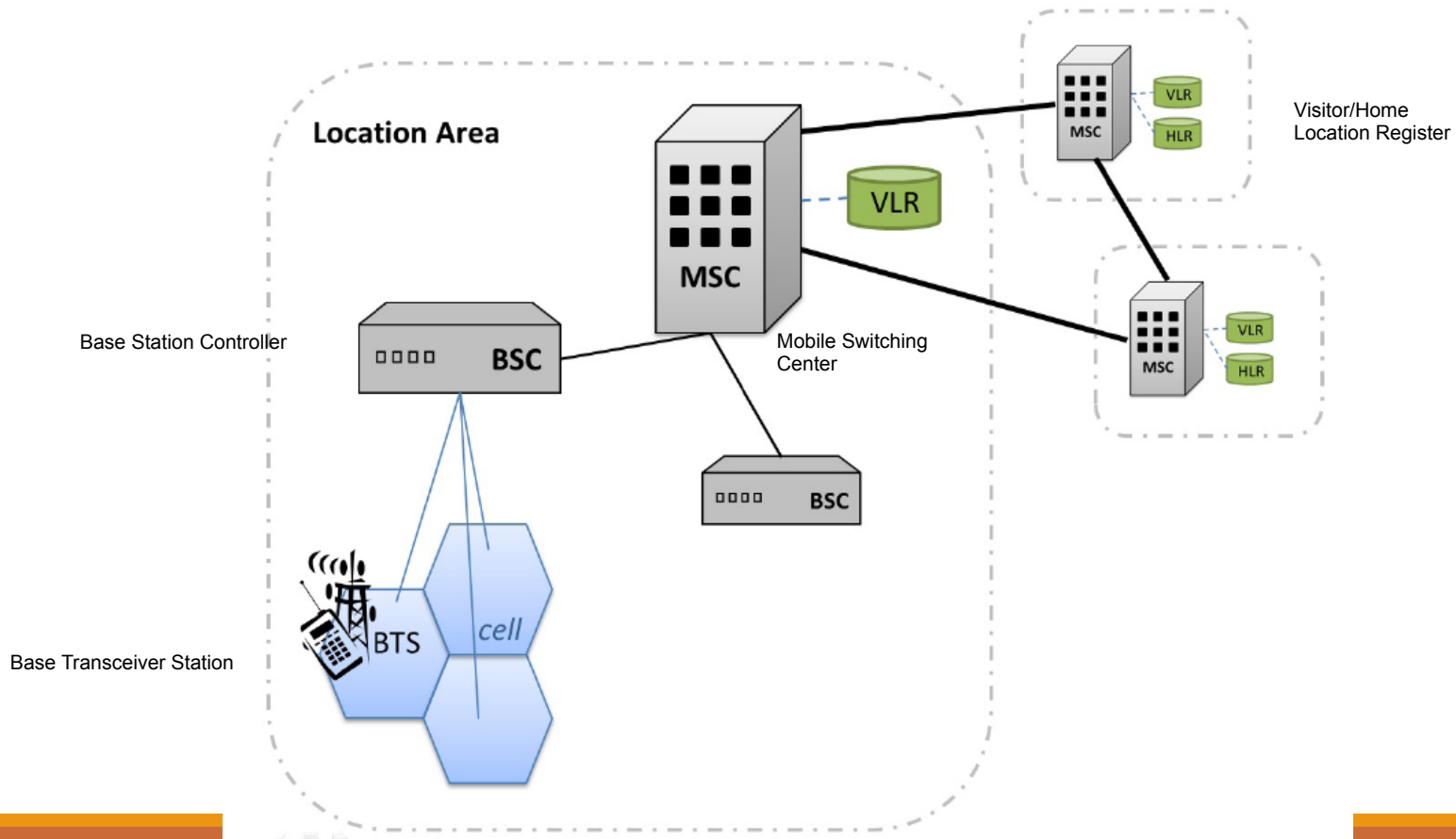
Call Graph

Who calls **whom** and **when**



GSM infrastructure

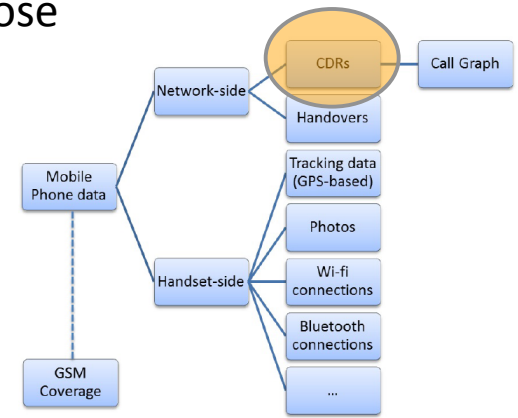
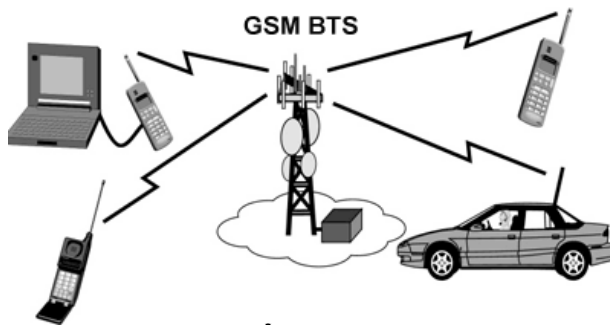
- Aimed at providing voice/data telecom.



GSM data - Description

Call Data Record (CDR)

Data gathered from mobile phone operator for billing purpose

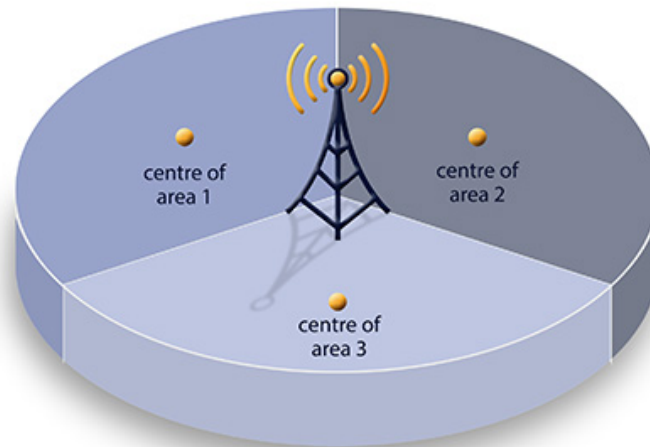


User id	Time start	Cell start	Cell end	Duration
10294595	"2014-02-20 14:24:58"	"PI010U2"	"PI010U1"	48
10294595	"2014-02-20 18:50:22"	"PI002G1"	"PI010U2"	78
10294595	"2014-02-21 09:19:51"	"PI080G1"	"PI016G1"	357

GSM data - Description

- Distinction between antenna and tower
 - Usually one “tower” carries 3 directional antennas
- Which one is in the data depends...

cell tower with 3 cells, each with 120° angle



Pros and cons of using GSM data

Pros

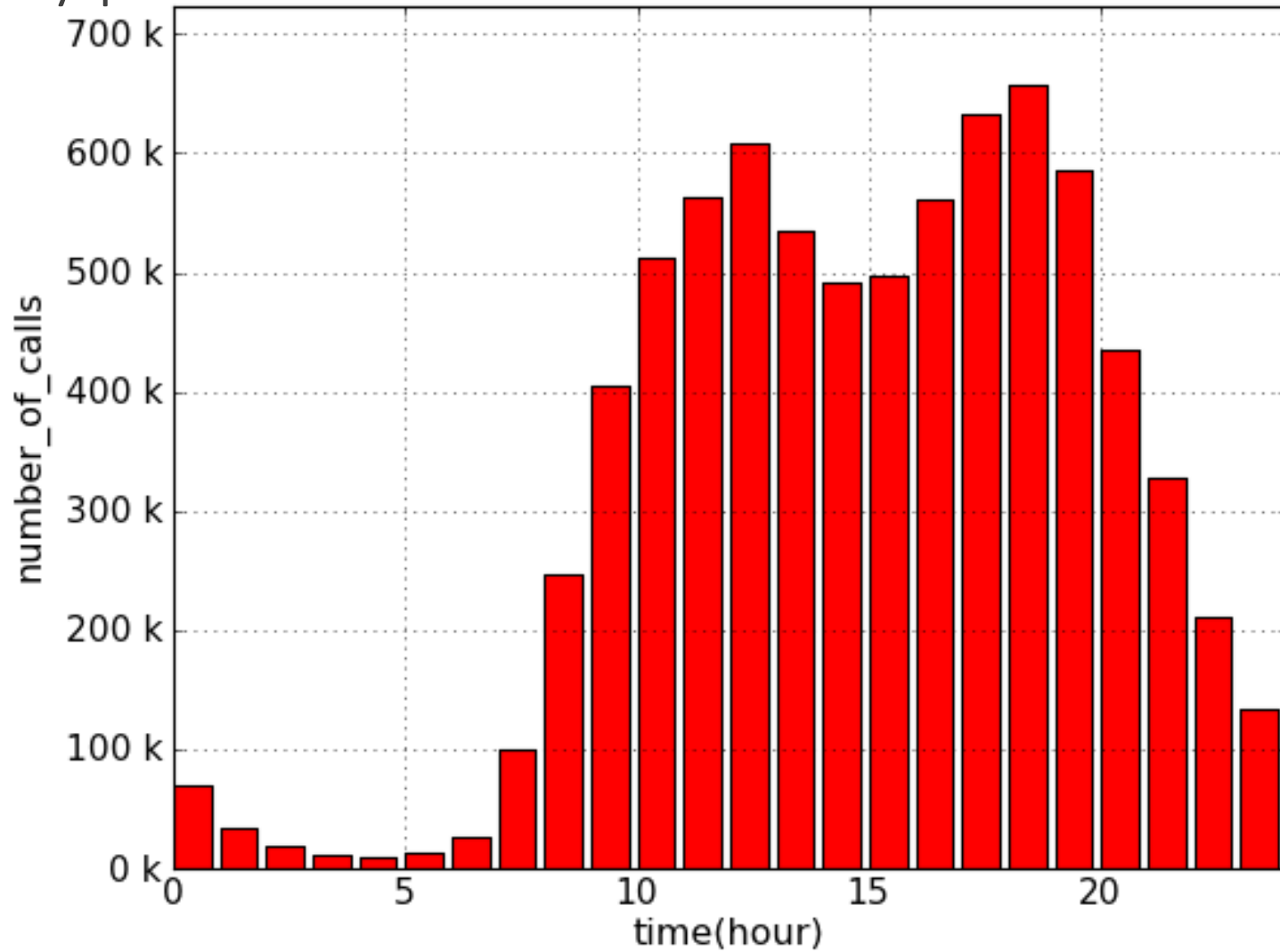
- Passive sensing: does not require an active contribution of the users
- Contains huge amount of information of how, when, with whom we communicate
- Same data format in all the world

Cons

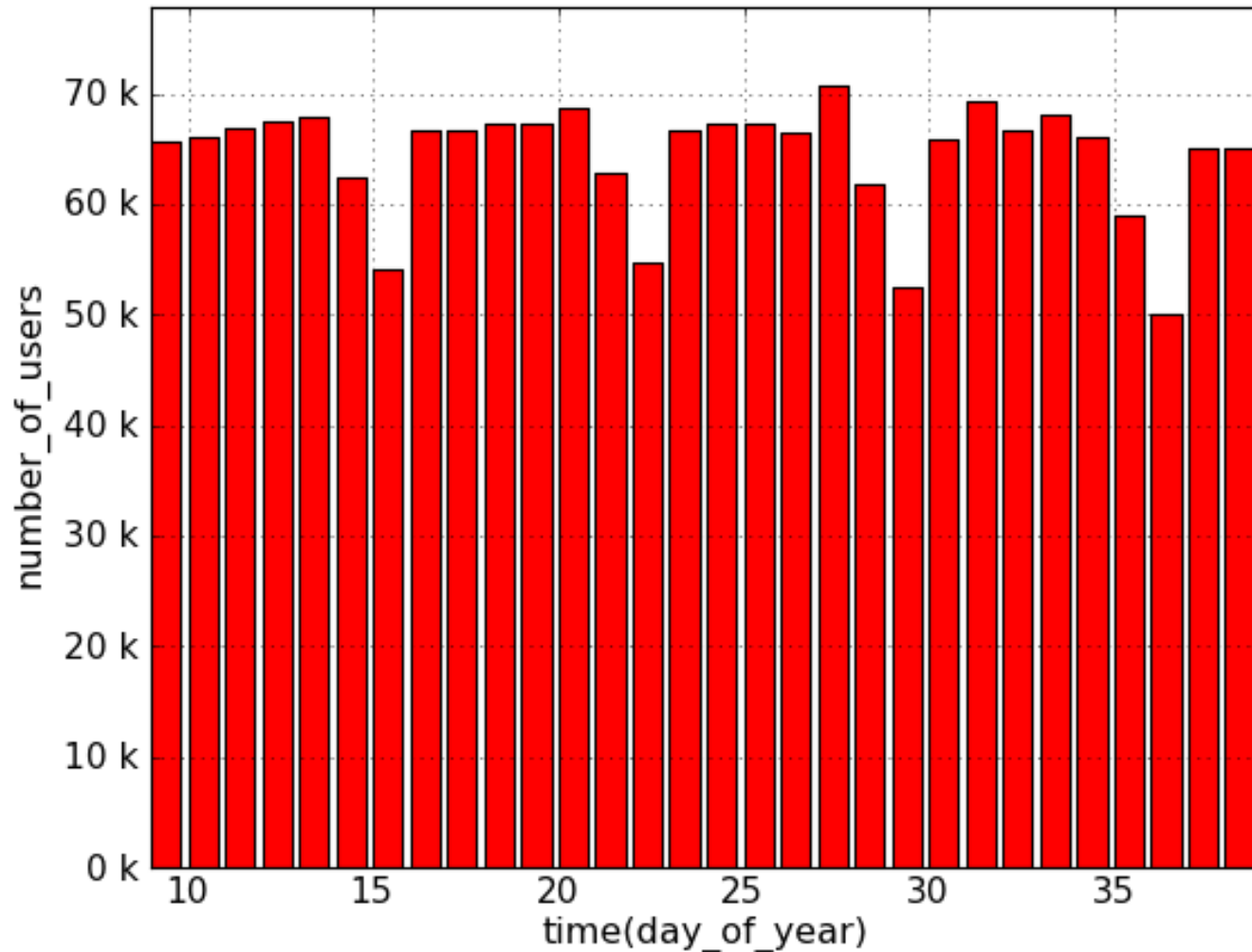
- Poor demographic and economic data
- Privacy concern: different legislations for different countries
- Low sampling: few events of calls for a considerable amount of users

Simple CDR-based statistics

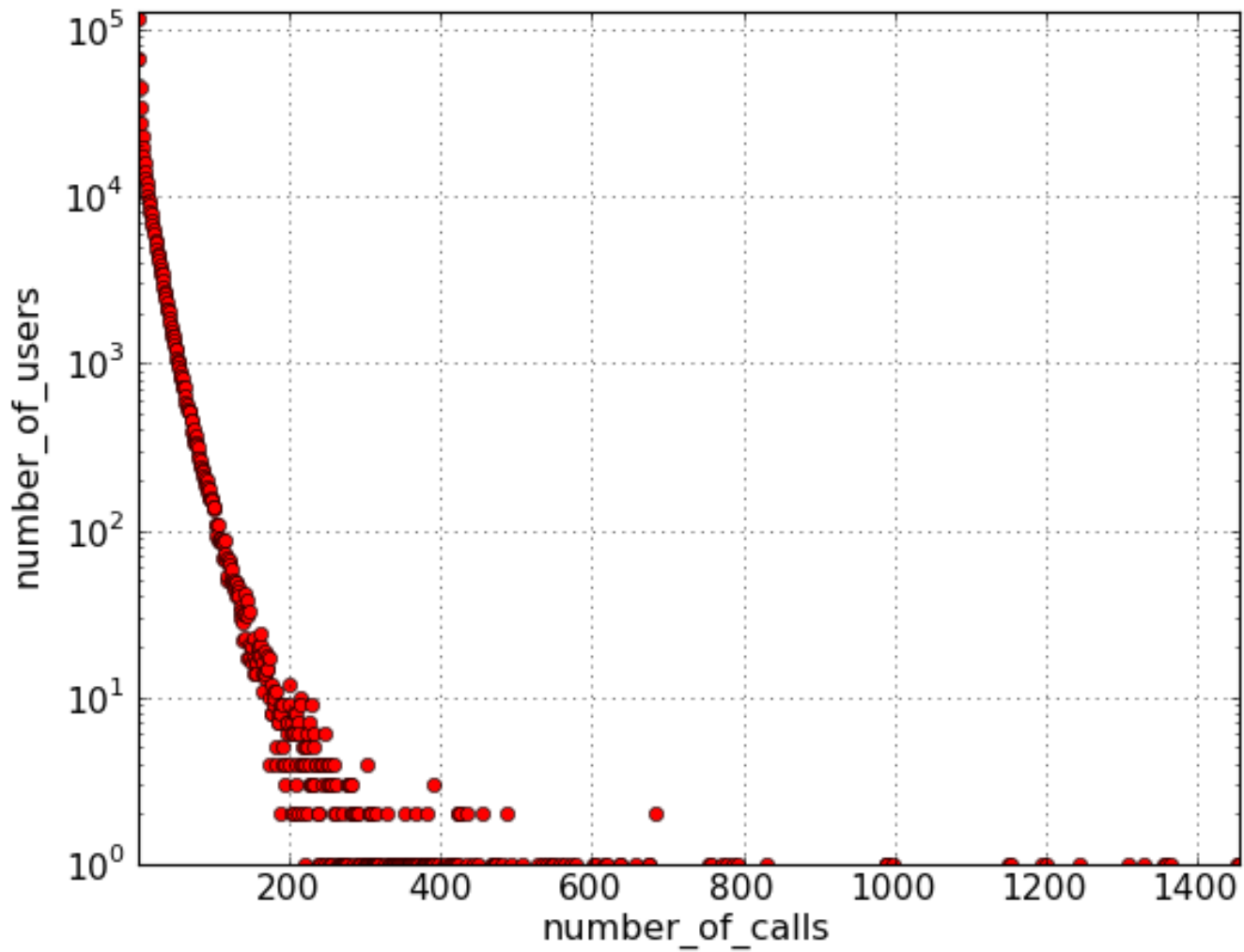
Daily pattern behavior



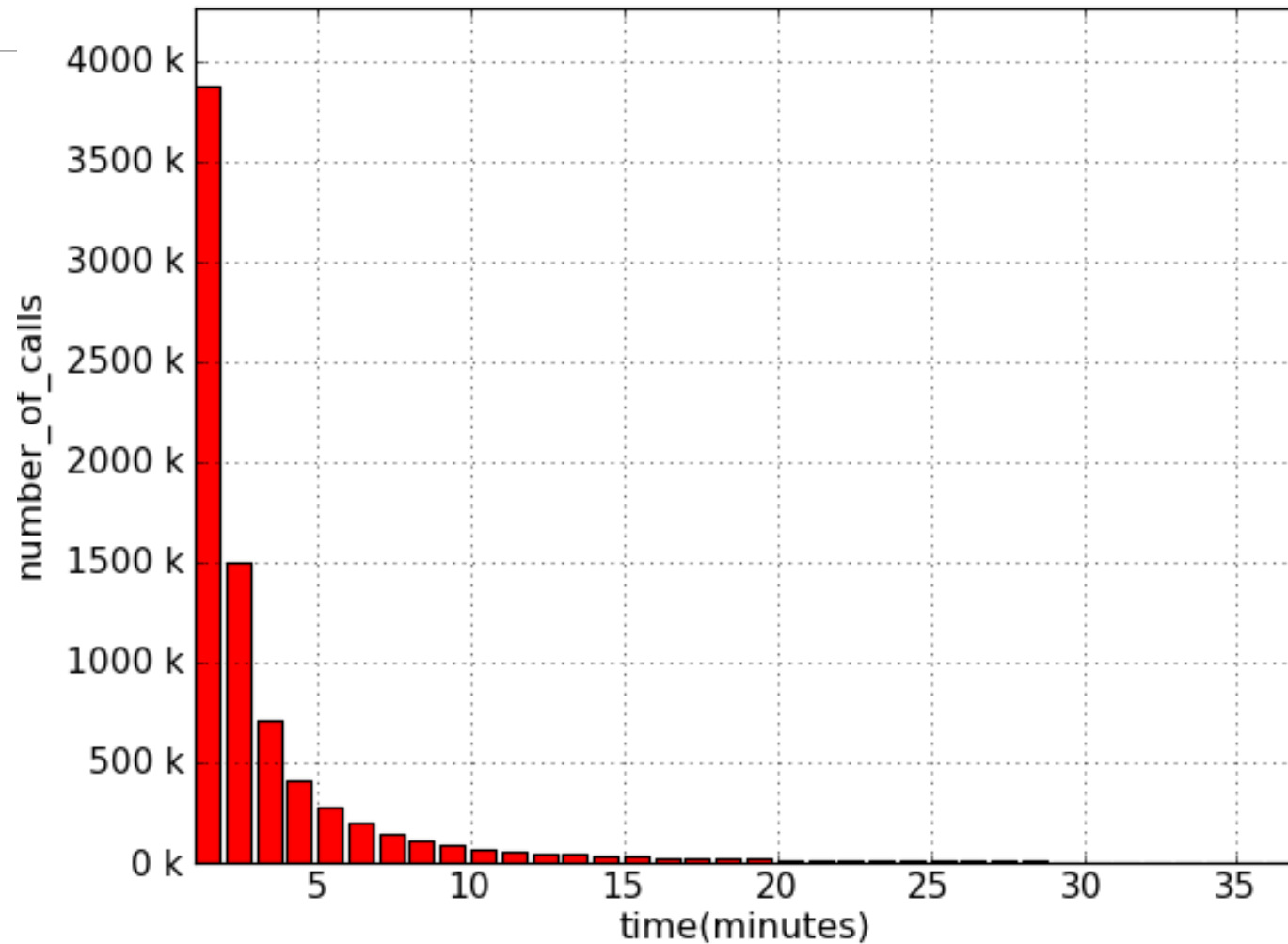
Weekly pattern behavior



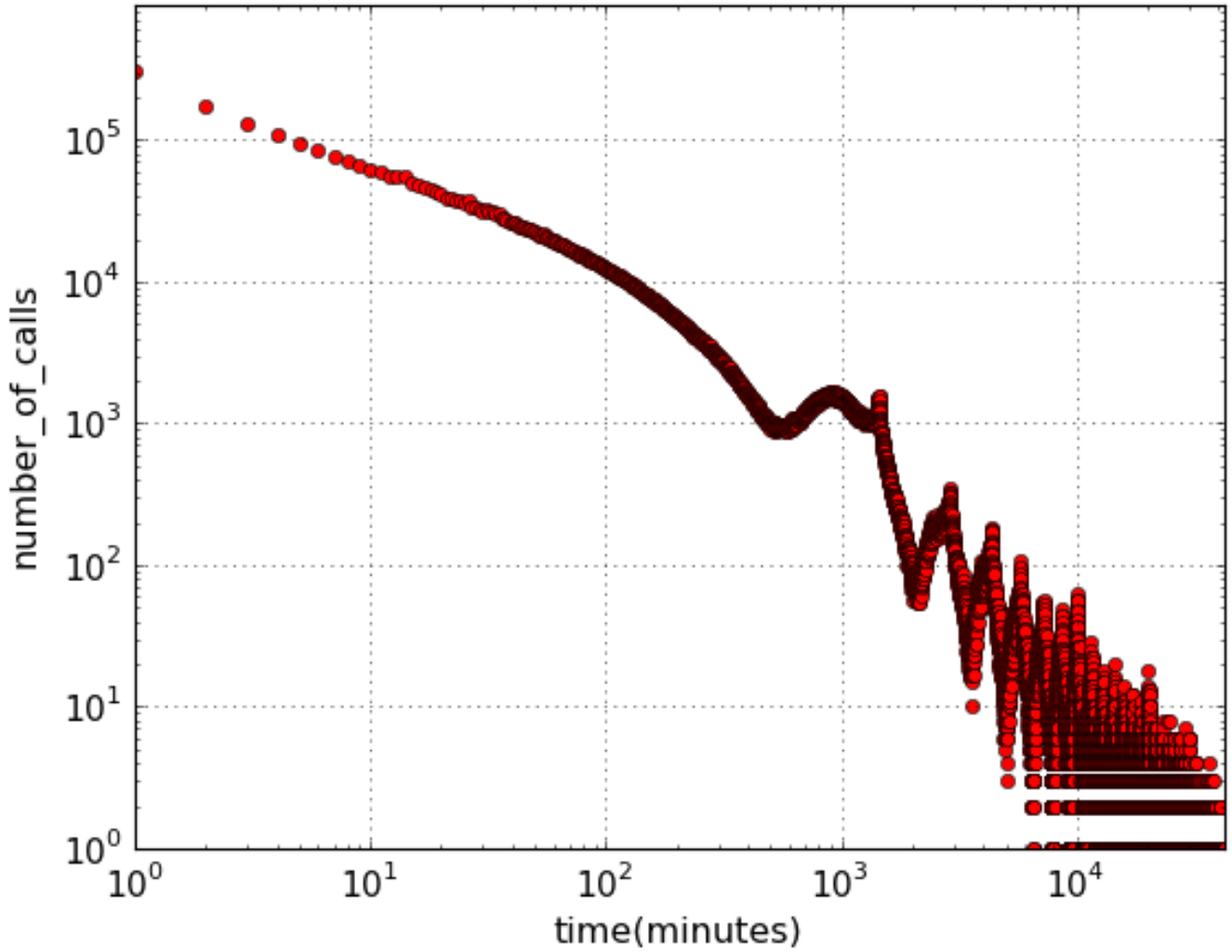
How many times we call?



How long we talk on the phone?



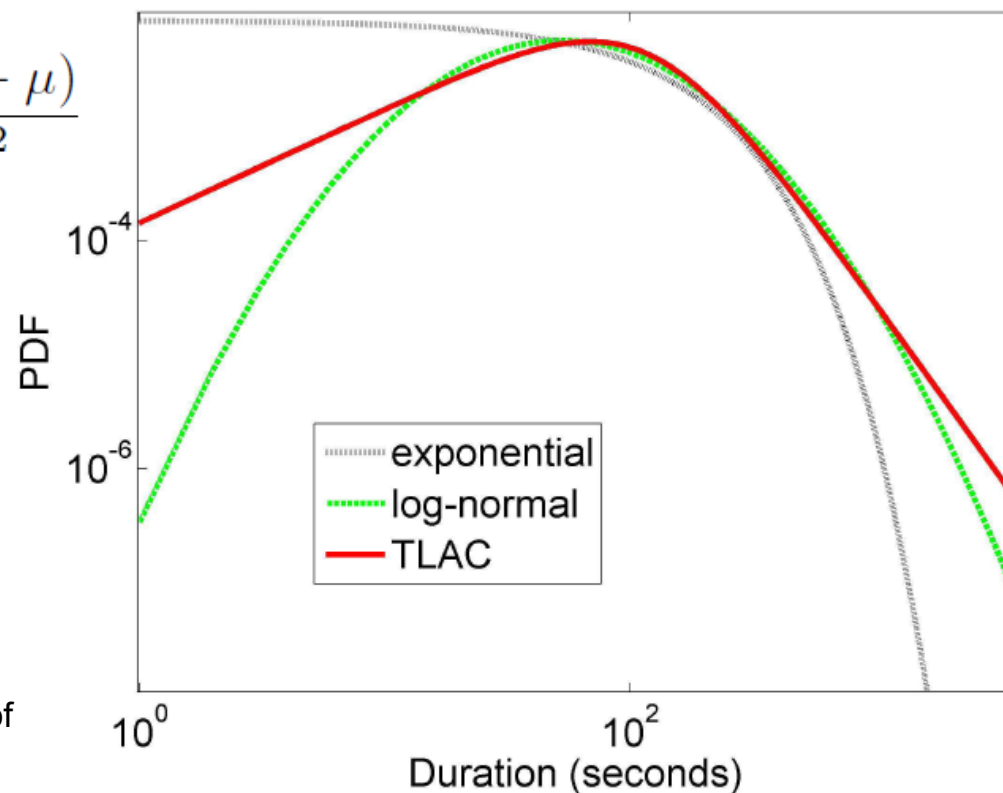
How many minutes goes by a call to the next?



Theoretical model of call durations

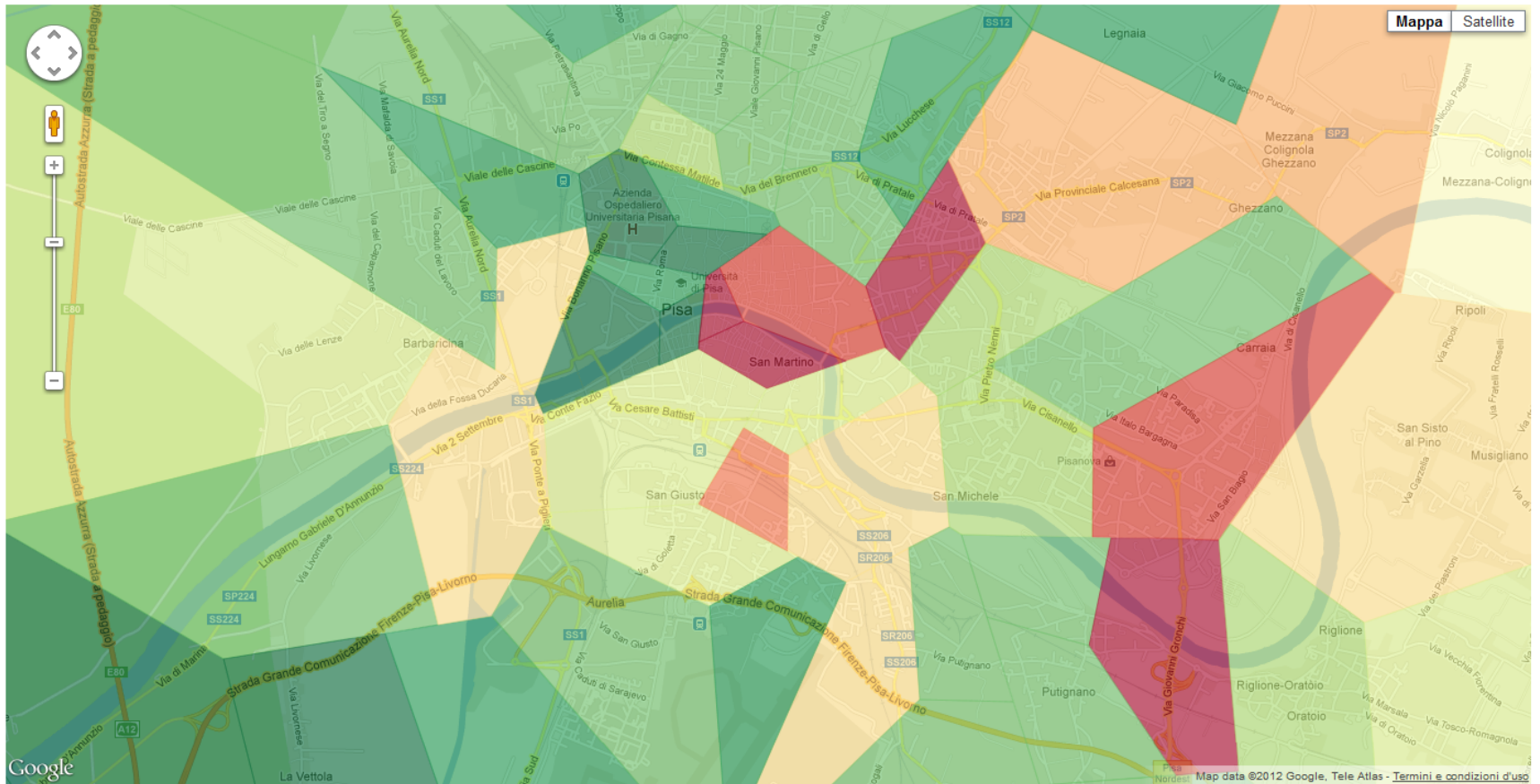
- Truncated Lazy Contractor (TLAC)

$$PDF_{TLAC}(x) = \frac{\exp(z(1 + \sigma) - \mu)}{(\sigma(1 + e^z))^2}$$



Join the **spatial** part of the mobile phone data

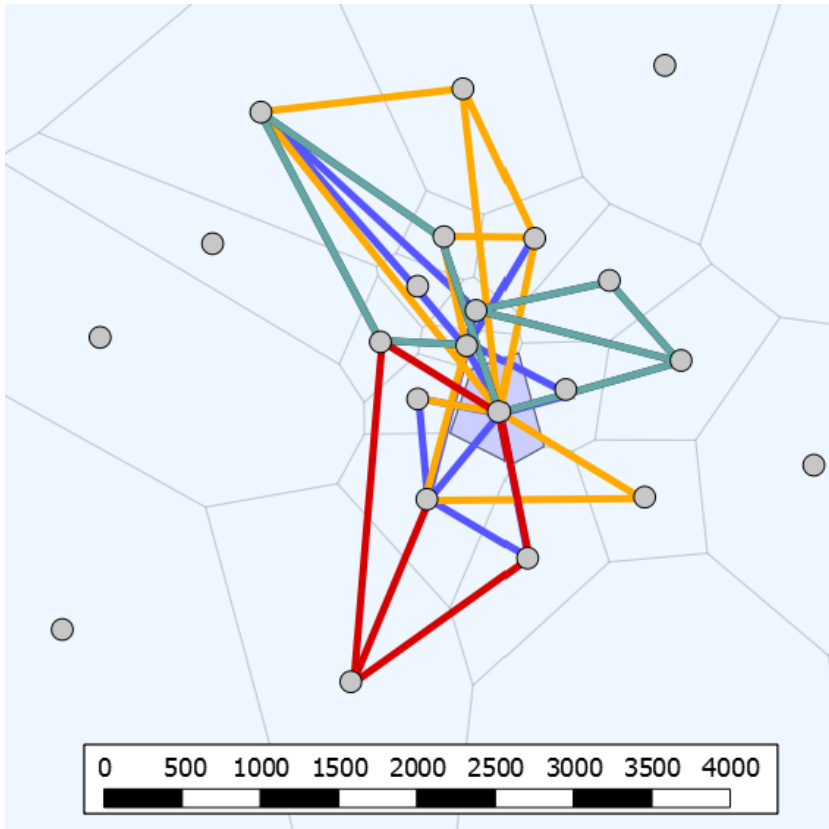
Spatial distribution of calls



ple within the working area of Pisa

Observing the **mobility** of individuals

Mobility Behaviours



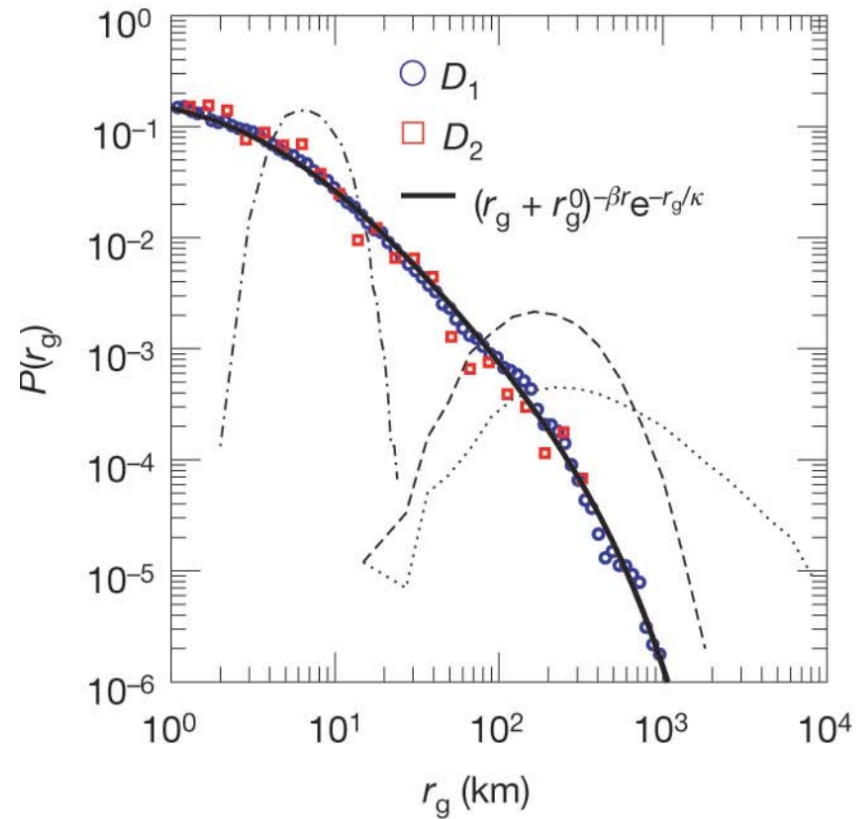
- The phone towers are shown as grey dots
- The trajectory describes the user's movements during 4 days (each day in a different color).

users move within a territory

Characteristic distance traveled by an individual

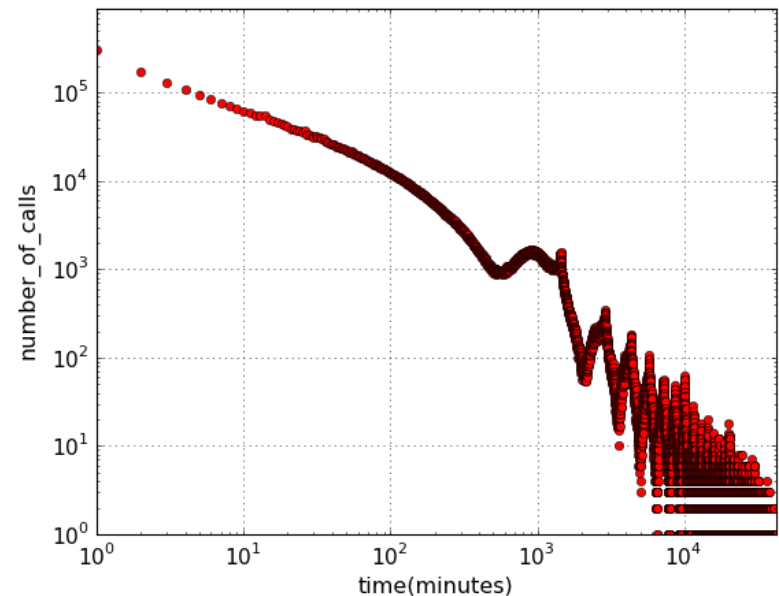
of gyration
has heavy tails

$$r_g = \sqrt{\frac{1}{N} \sum_{i \in L} n_i (\vec{r}_i - \vec{r}_{cm})^2},$$

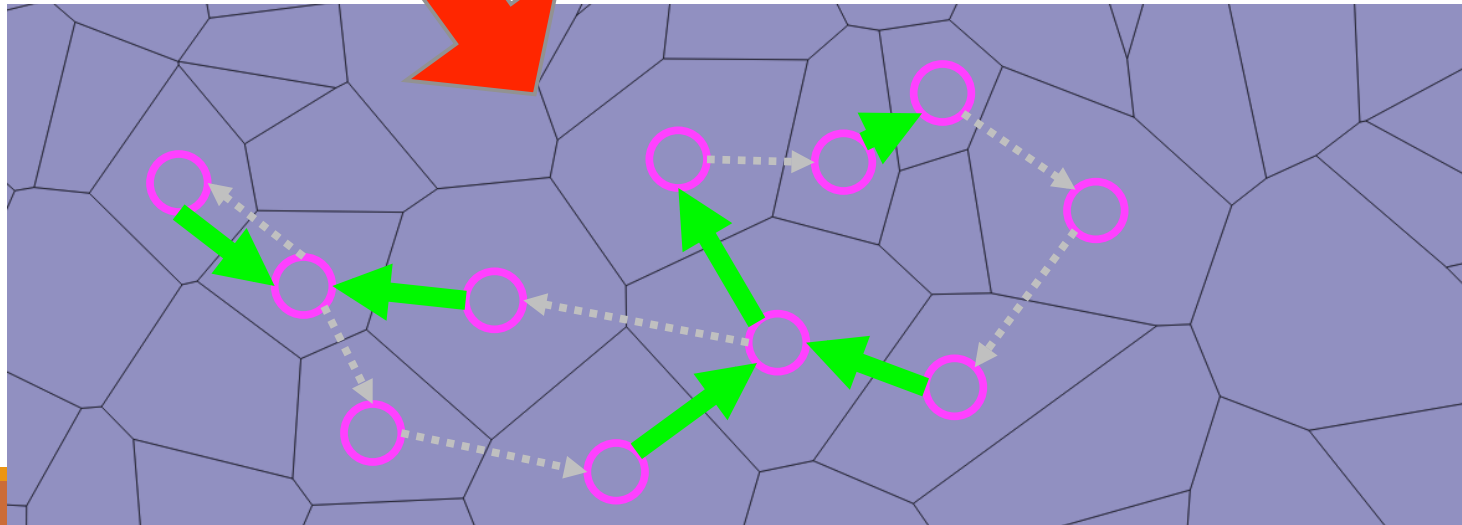
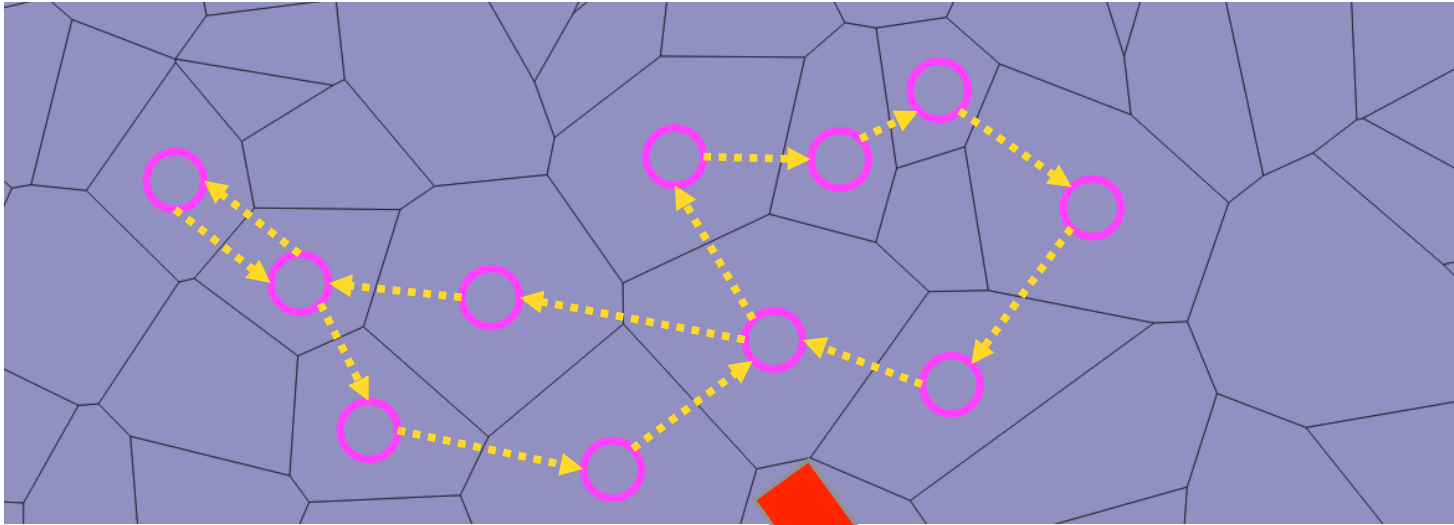


Estimating movements

- Reconstruct individual mobility through consecutive locations (individual flows)
- If $|\text{time}(\text{Call_1}) - \text{time}(\text{Call_2})| < \Delta T$
then consider movement $\text{Call_1} \rightarrow \text{Call_2}$
- Issue: how to choose threshold?
 - Large $\Delta T \Rightarrow$ spurious data
 - Small $\Delta T \Rightarrow$ miss data

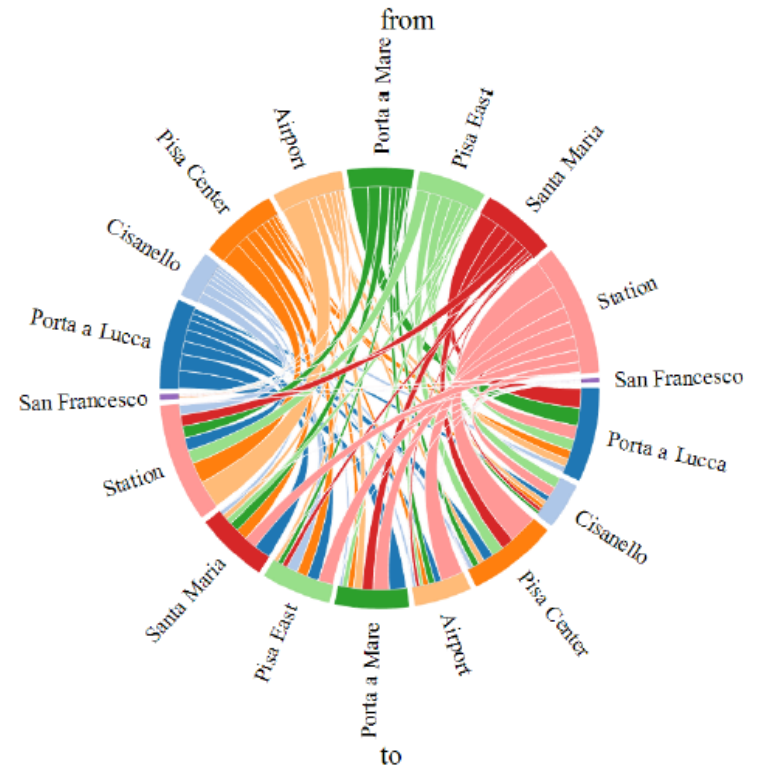
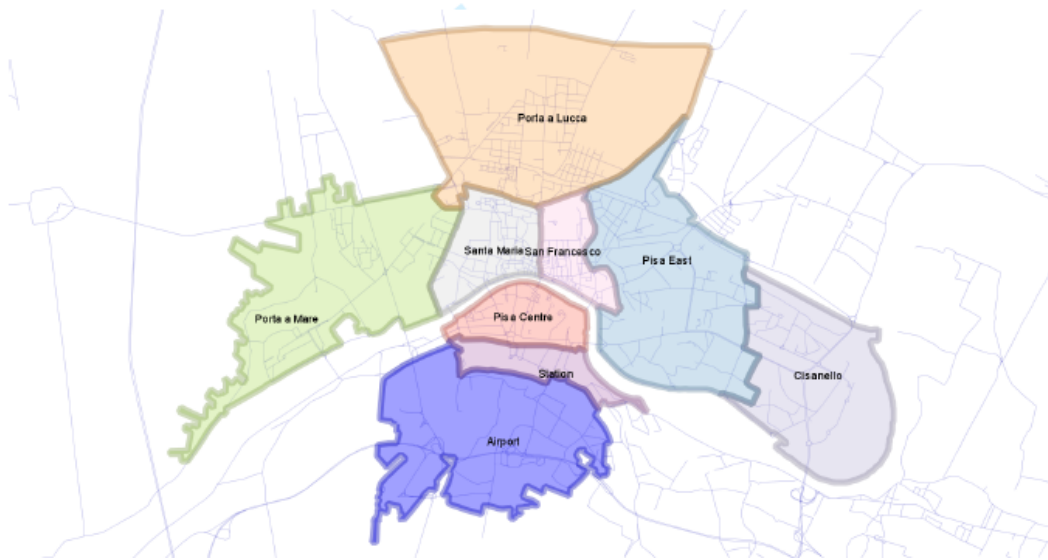


Estimating movements



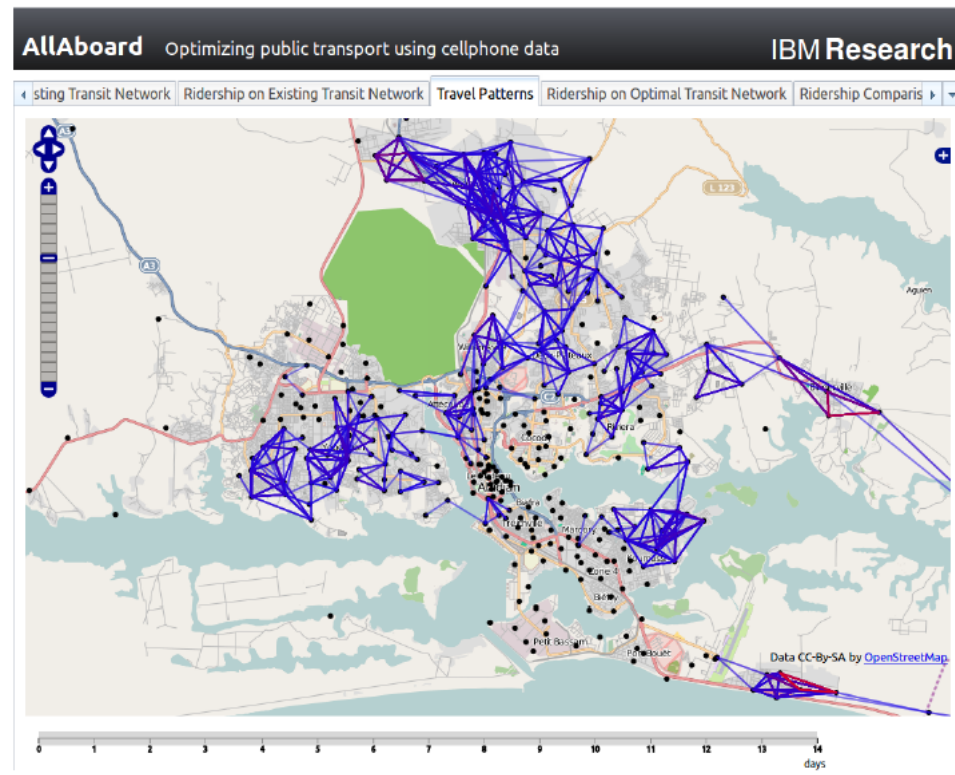
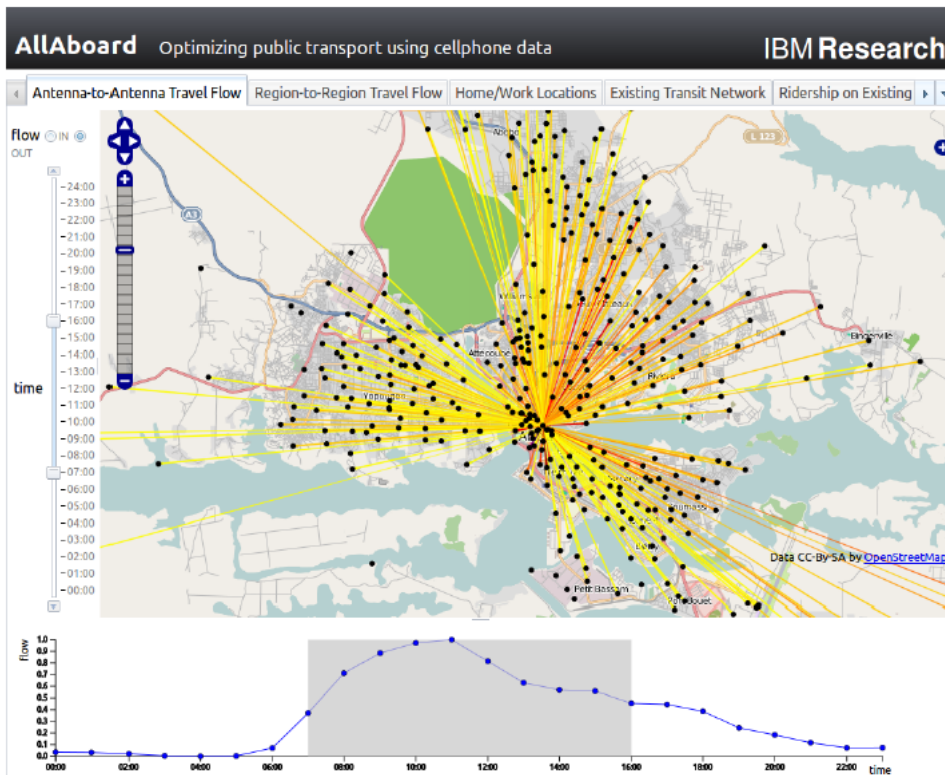
Estimating movements

- Example on Pisa city



Estimating movements

- Example on Abidjan (Ivory Coast)



Michele Berlingiero, Francesco Calabrese, Giusy Di Lorenzo, Rahul Nair, Fabio Pinelli, Marco Luca Sbodio.

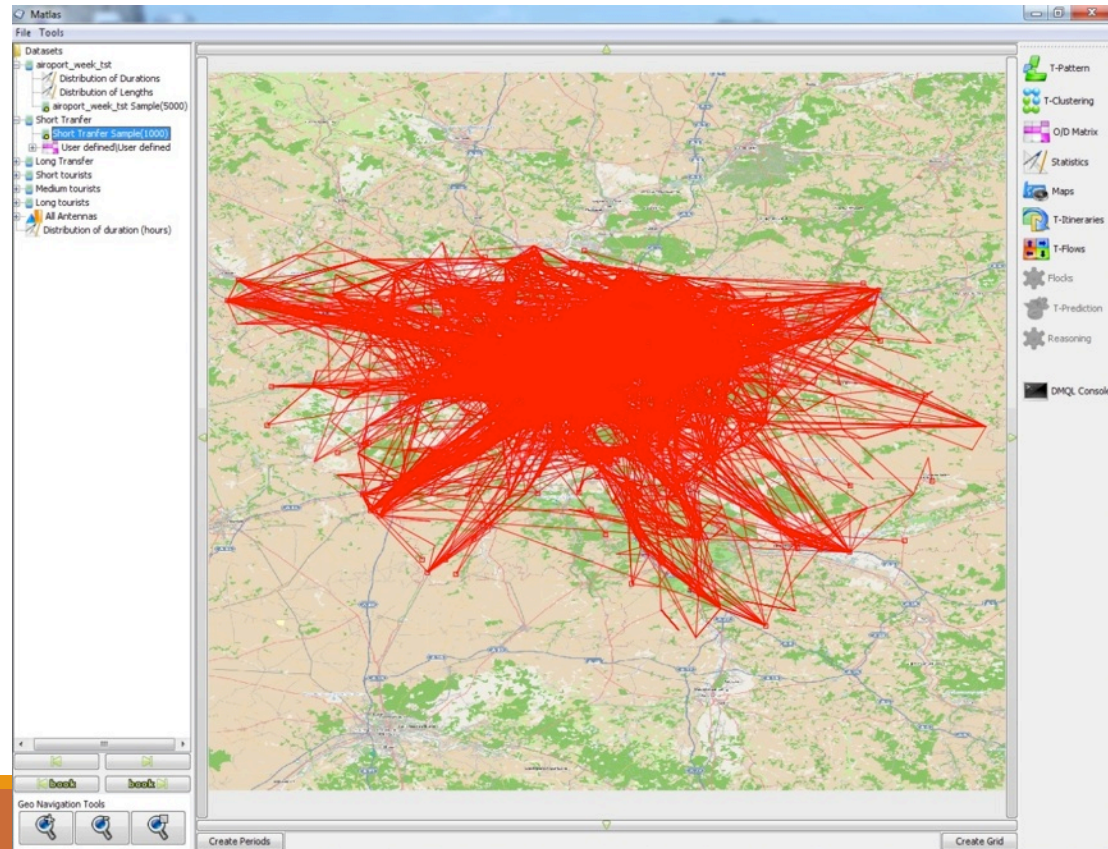
AllAboard: a system for exploring urban mobility and optimizing public transport using cellphone data.

http://researcher.watson.ibm.com/researcher/view_group_subpage.php?id=1716

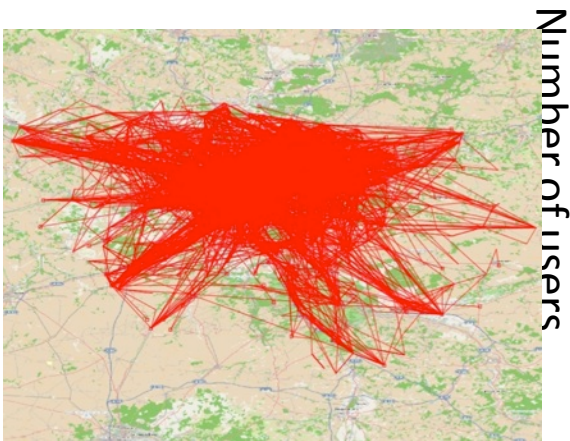
Sample application: Analyzing tourist data

- Case study of foreign (roaming) visitors of Paris area
- Users arriving and leaving at CDG airport

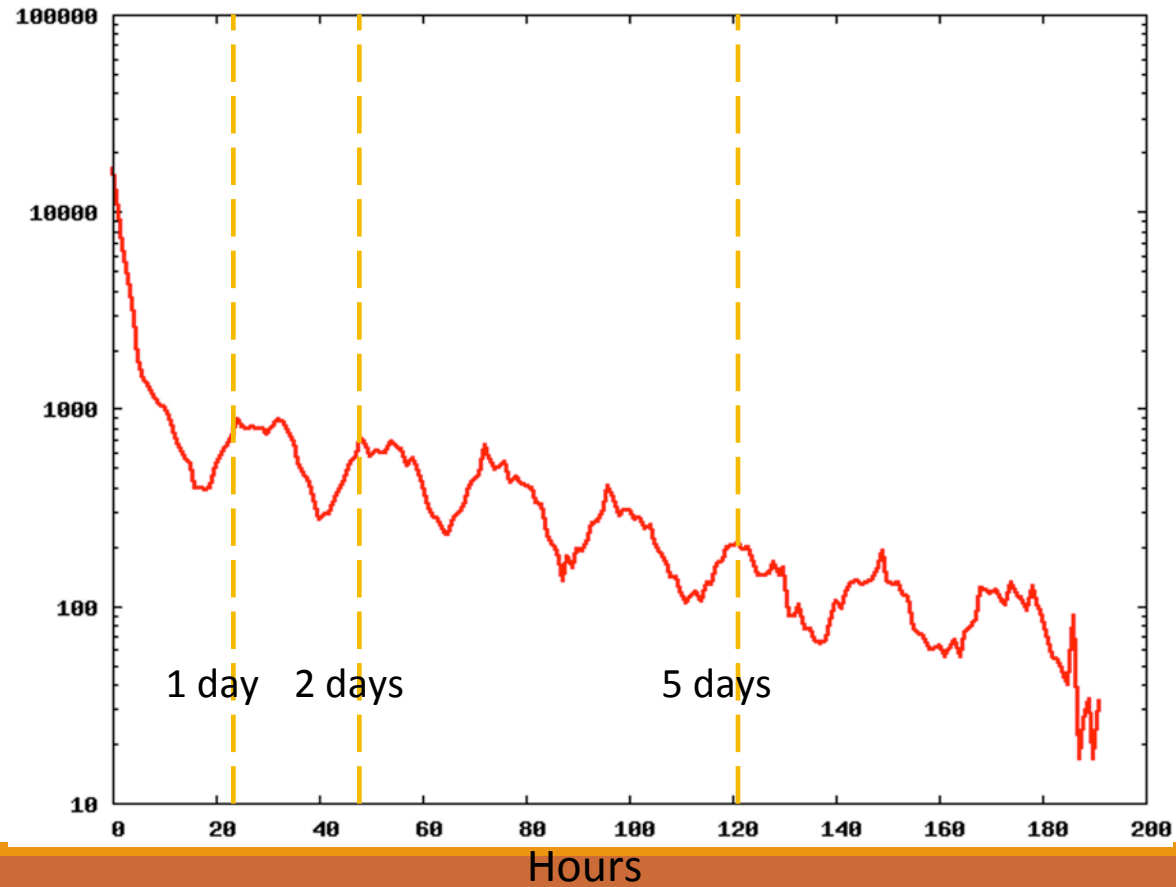
106 000 Users



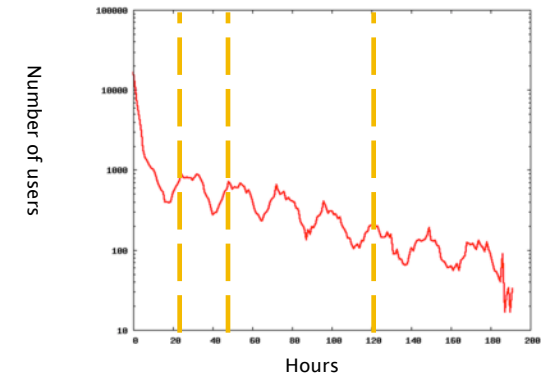
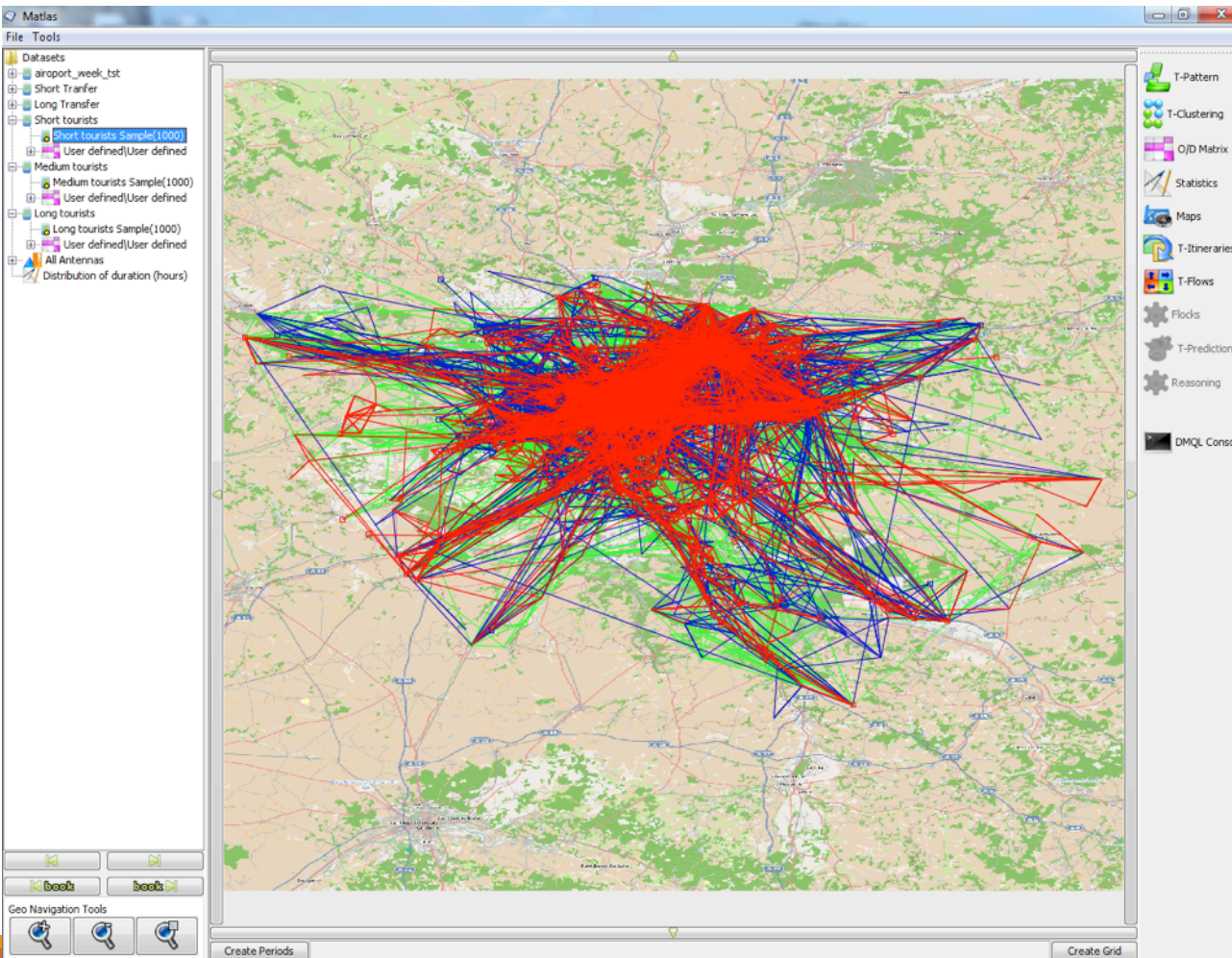
Distribution of visiting time



Number of users

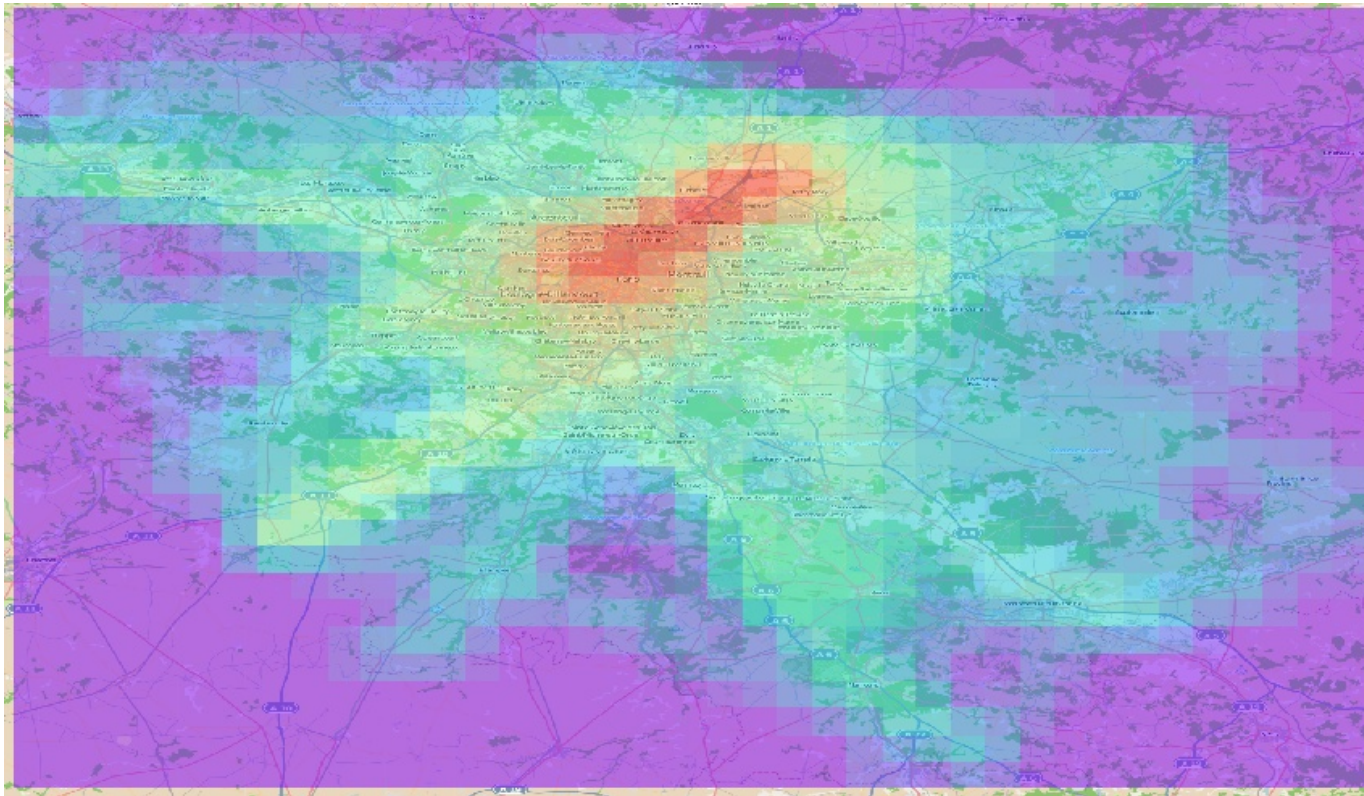


Categorization of tourists



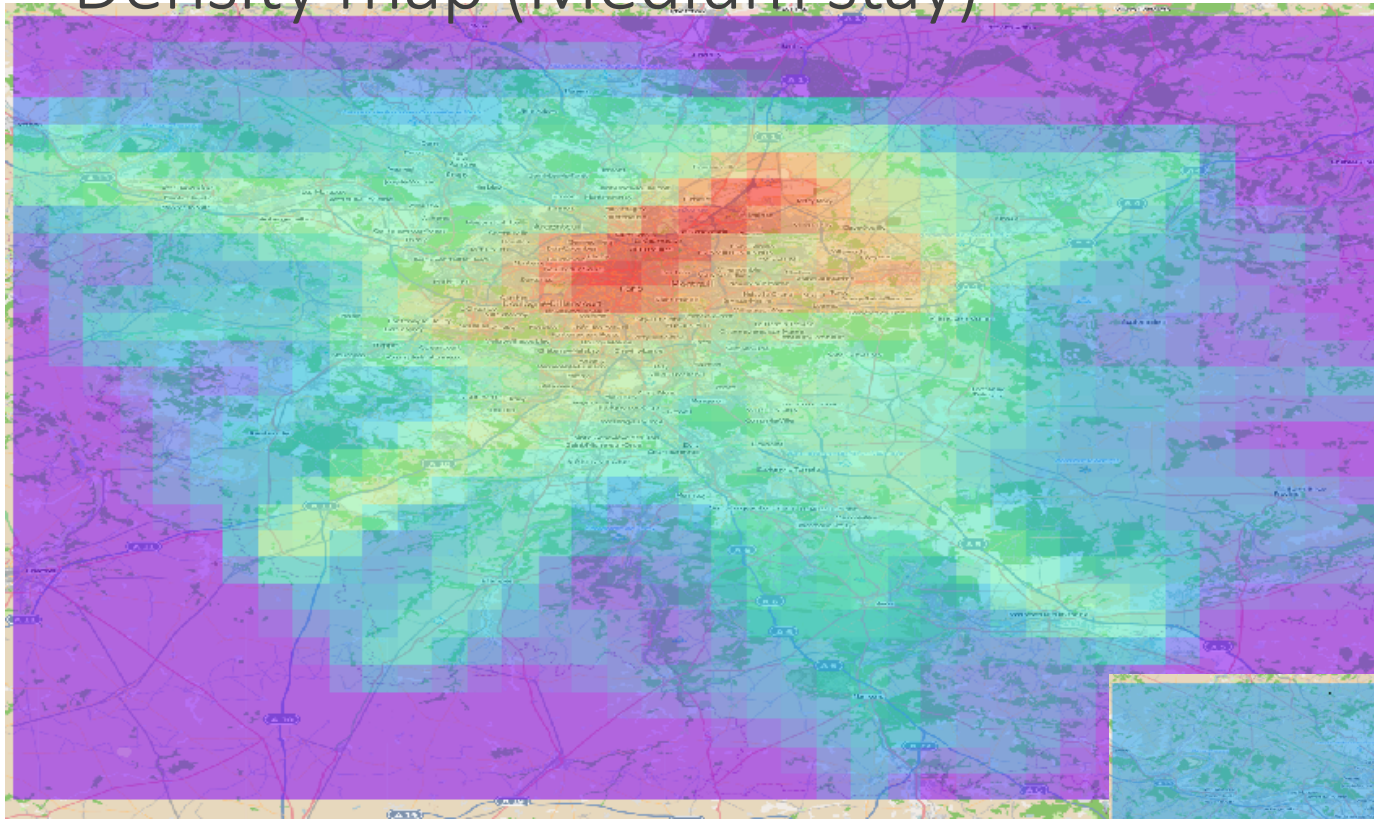
- Short period stay Tourist (1 day → 2 days)
- Medium period stay Tourist (2 day → 5 days)
- Long period stay Tourist (5 day → 7 days)

Density map (Short stay)



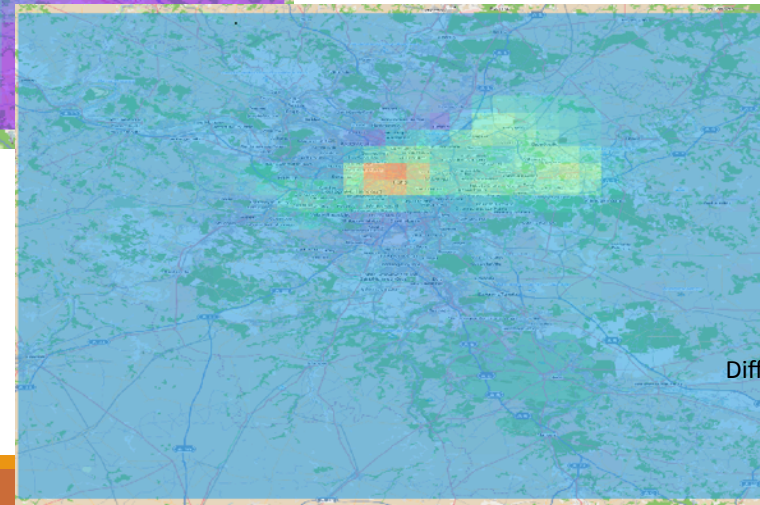
Short stay tourists visi

Density map (Medium stay)



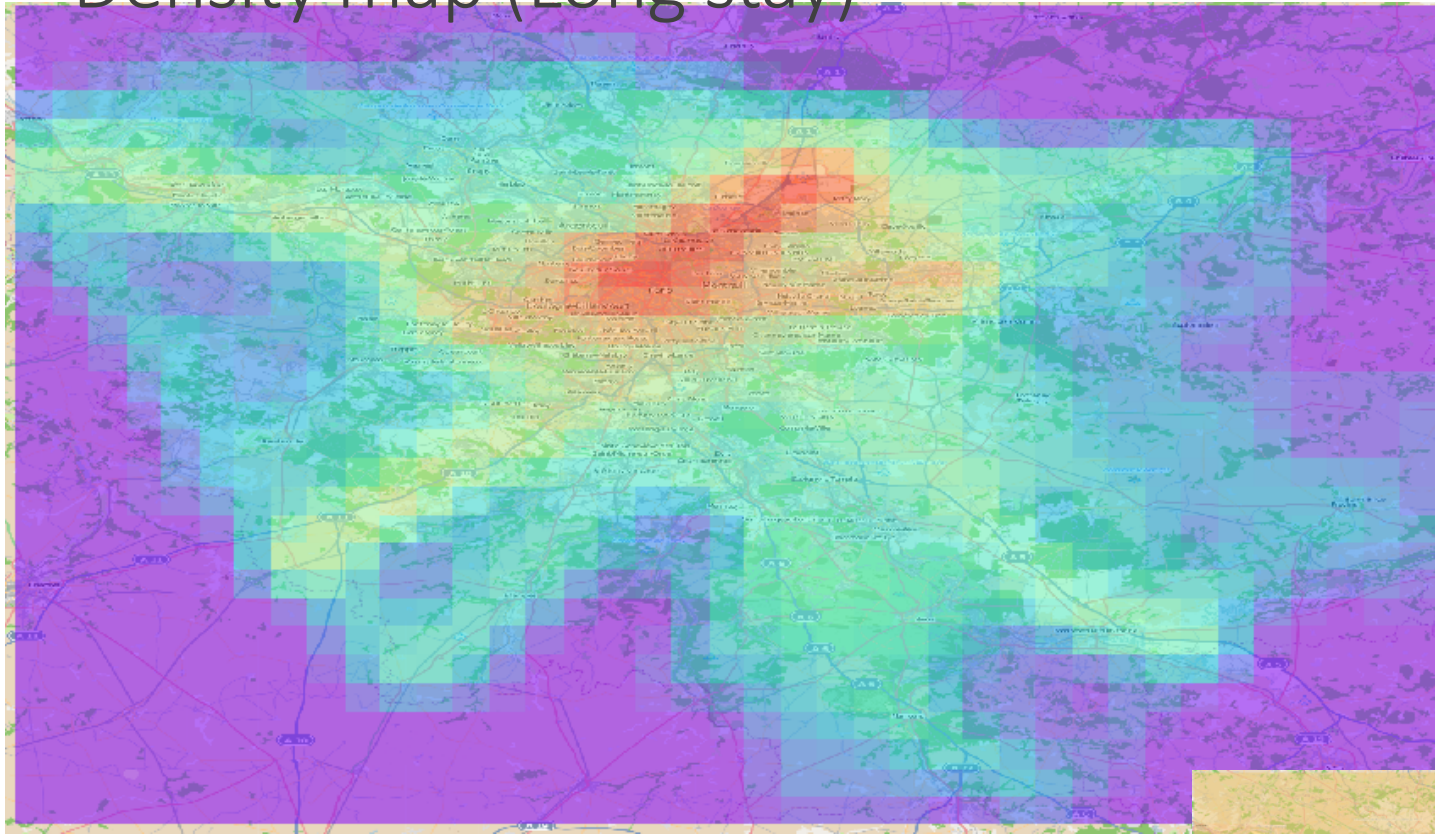
Medium stay tourists visit

Green = Disneyland Paris
Red = Versailles



Differ

Density map (Long stay)



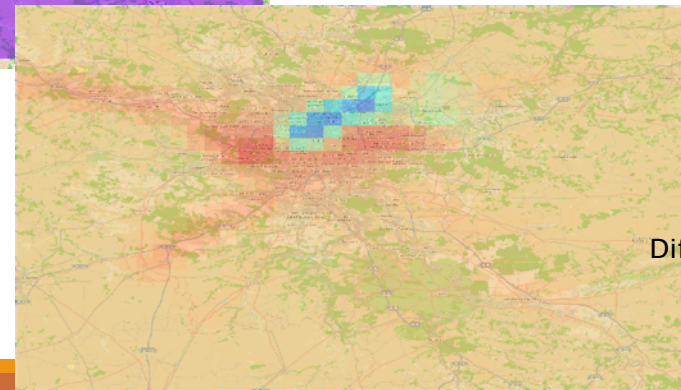
Long stay tourists visit t

Green = Disneyland Paris

Red = Versailles

Blue = Highway/Train to Mante la Jolie

Black = Highway to South-West



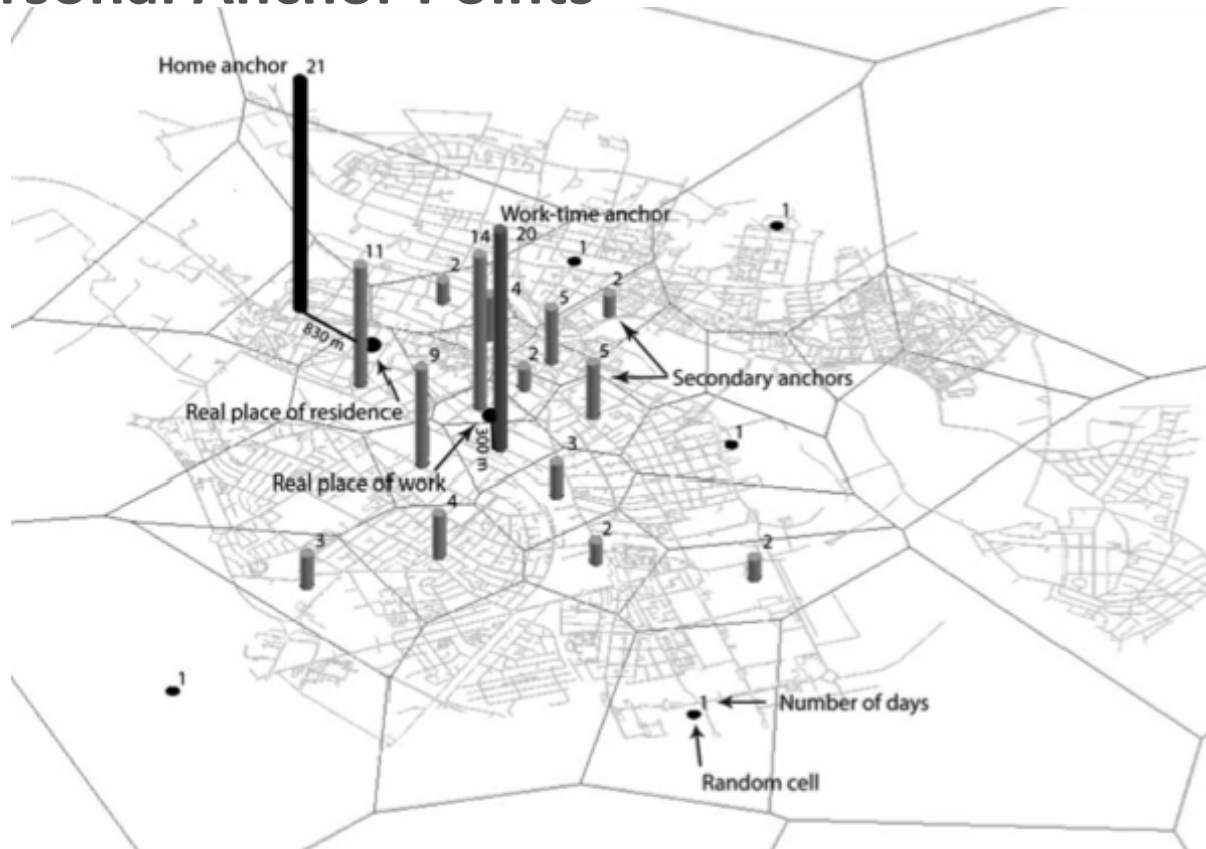
Diff

Identifying important locations

- Home (residence) and Work play an important role in understanding urban mobility
- **“Personal Anchor Points”**: high-frequency visited places of a user
 - Select top 2 cells with max number of days with calls
 - Determine home and work through time constraints:
 - average start time of calls and its deviation

Identifying important locations

- **“Personal Anchor Points”**

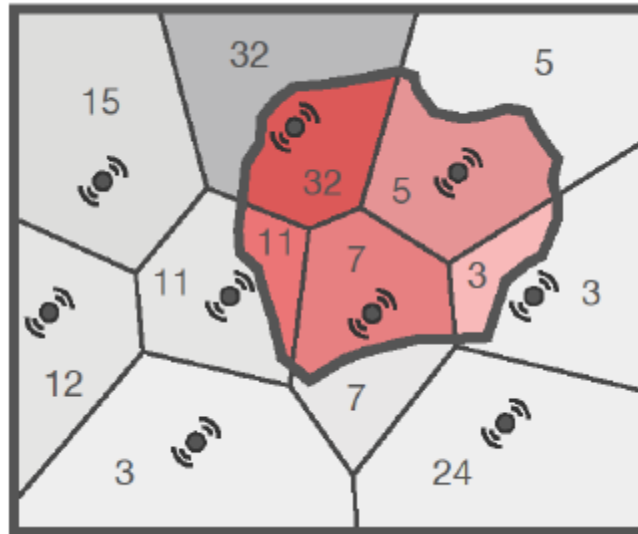


Identifying important locations

- Estimating users' **residence through night activity**
 - Home = region with highest frequency of calls during nighttime
- First issue: cells might not correspond perfectly to the regions to measure
- Second issue: cells might not have uniform density of population

Identifying important locations

- First issue: cells might not correspond perfectly to the regions to measure



$$\sigma_{c_i} = \frac{1}{A_{c_i}} \sum_{v_j} \sigma_{v_j} A_{(c_i \cap v_j)}$$

- Approach: each cell contributes proportionally to its overlap with the region

Identifying important locations

- Second issue: cells might not have uniform density of population



$$\rho_i^{RS} = \frac{W_i}{\sum_j W_j} P_i$$

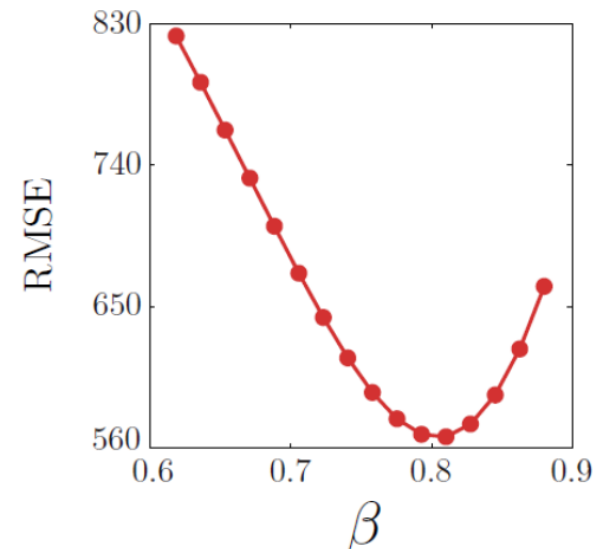
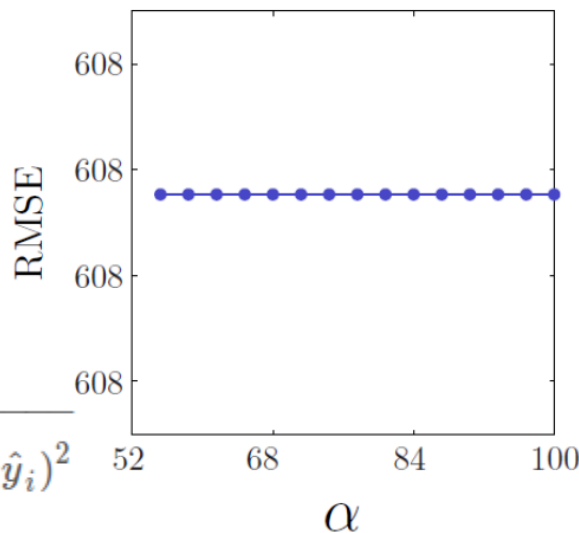
- Approach: integrate external indicators of relative density – e.g. from environment and infrastructures – to distribute cells' contrib.

Identifying important locations

- Linear or superlinear relation?

- ρ_c = population density
- σ_c = mobile phone residents
- P = national population (real vs. estimated)

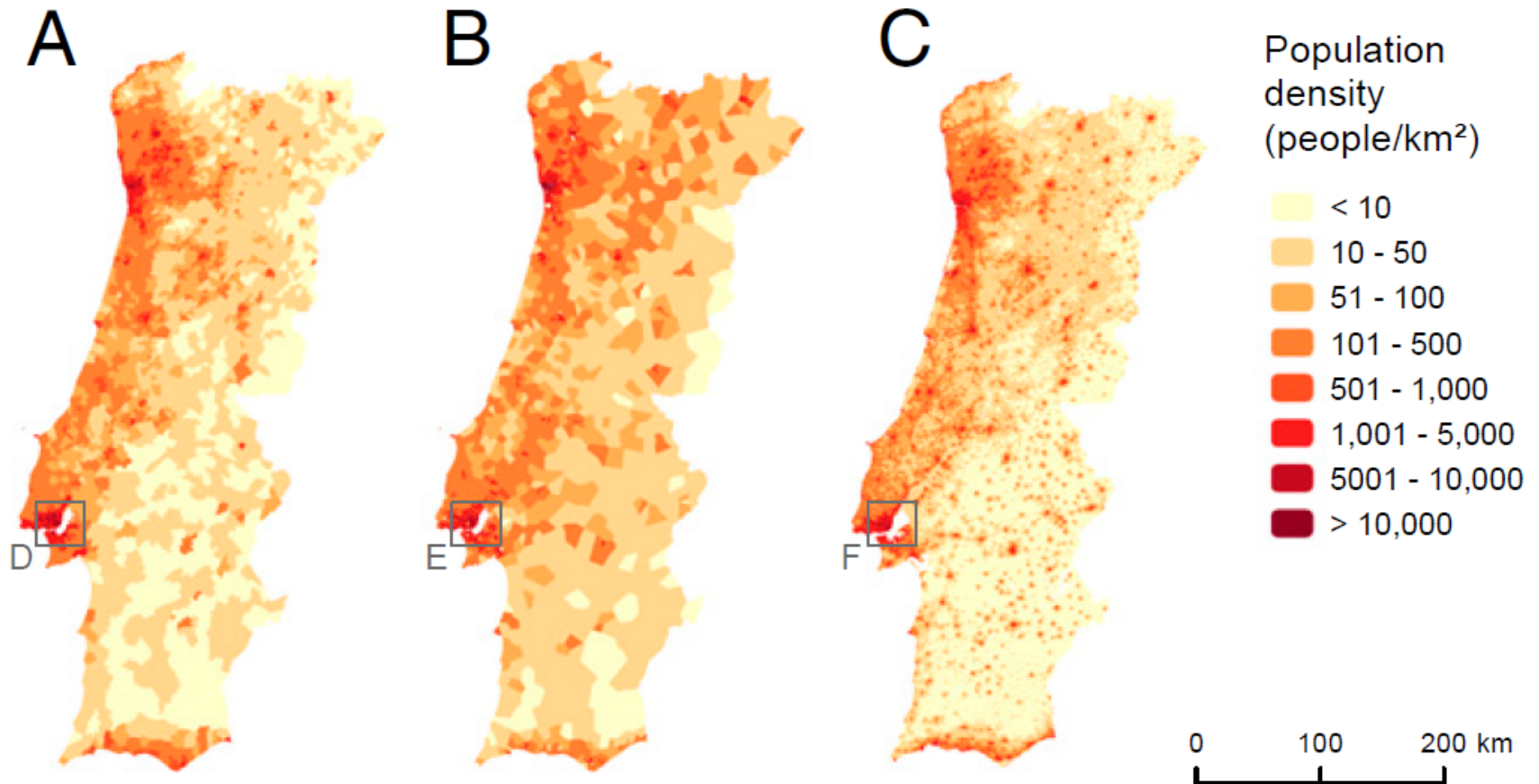
$$\rho_c = \frac{P}{\hat{P}} \alpha \sigma_c^\beta$$



$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Identifying important locations

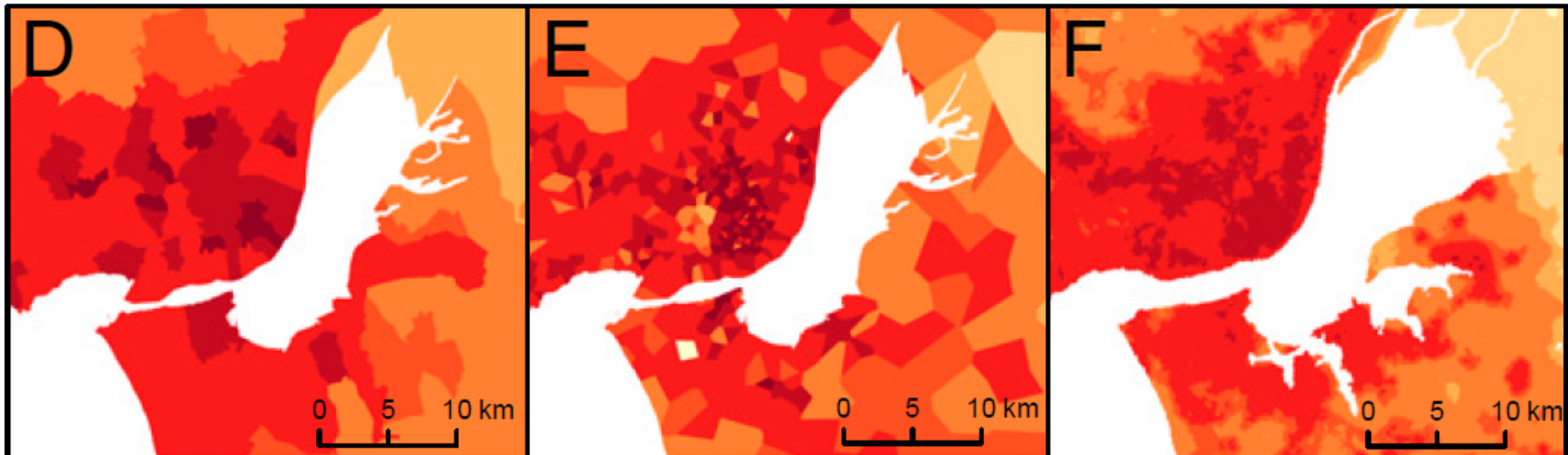
- Sample results on Portugal



A = Census B = GSM data C = Environment/Infrastructures-based

Identifying important locations

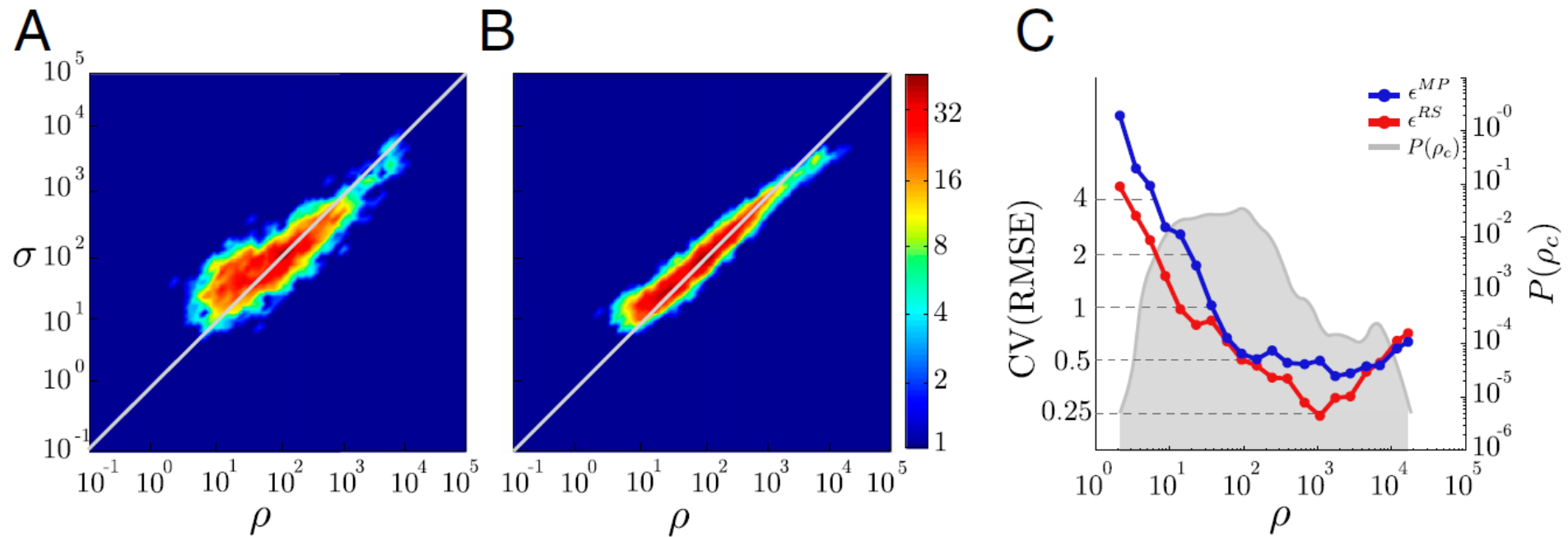
- Sample results on Portugal (close-up)



D = Census E = GSM data F = Environment/Infrastructures-based

Identifying important locations

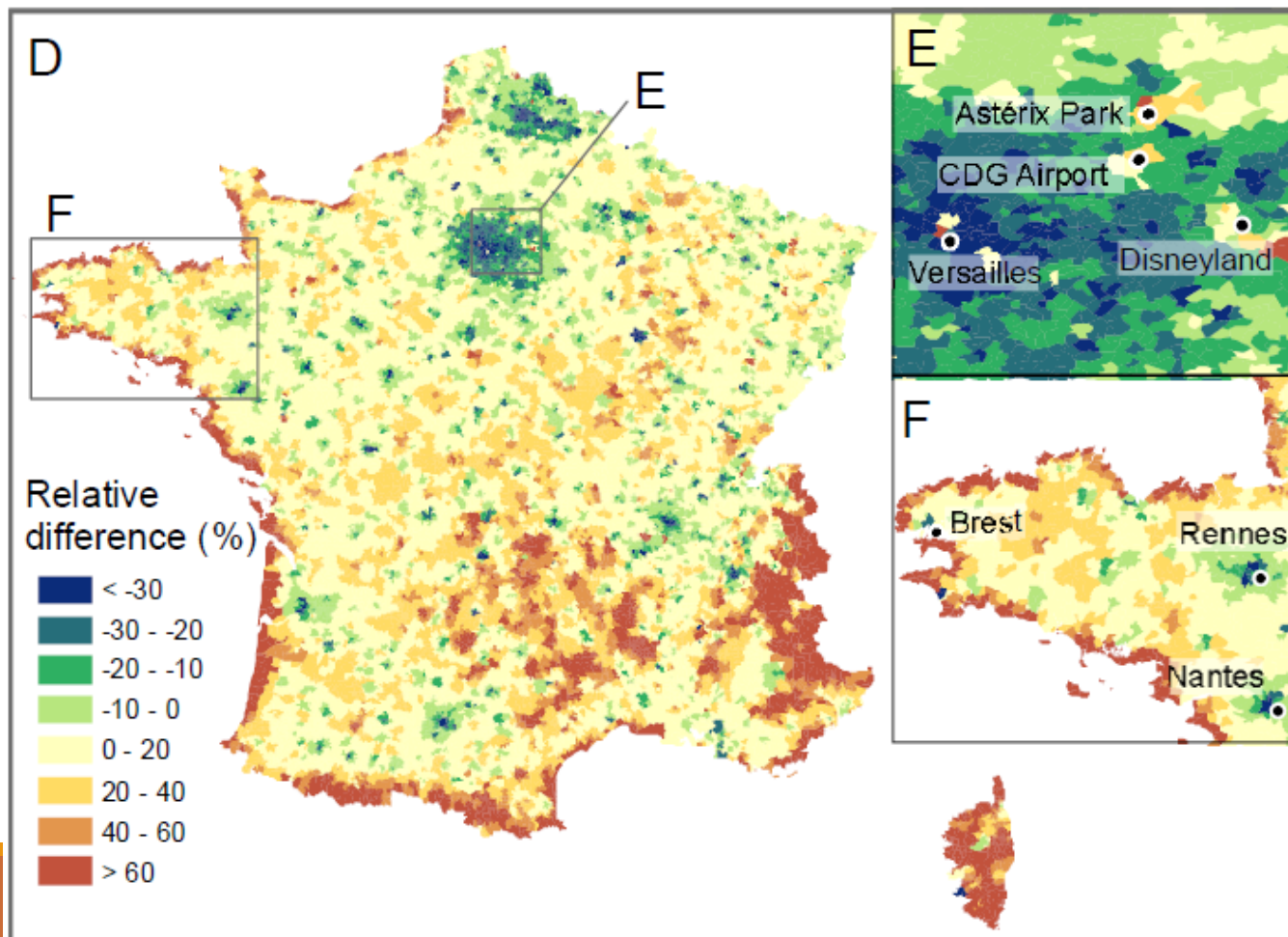
- Sample results



A = GSM data B = Environment/Infrastructures-based

Identifying important locations

- Sample usage: evaluate seasonal changes
 - Summer variations vs. Winter period



Classifying into **city users** categories



Basic methodology: Sociometer

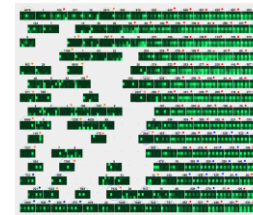
- GSM calls used as proxy of users' presence in a specific area
- 3 categories used: Residents, Commuters, Visitors

GSM Calls

Mo	Tu	We	Th	Fr	Sa	Su
5	4		3	2	1	5
	4	4		1	1	1



Computation



Profile Map



Commuters

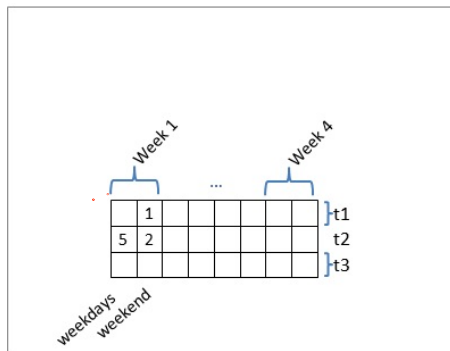


Visitors/Tourists

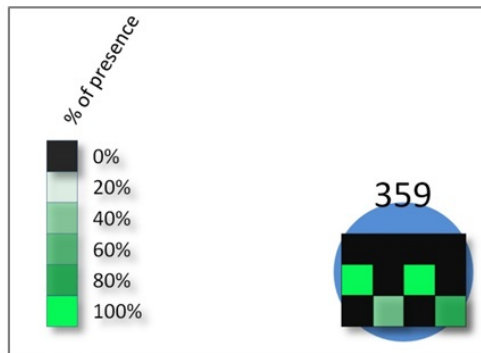


Residents

Temporal Profile



(a)

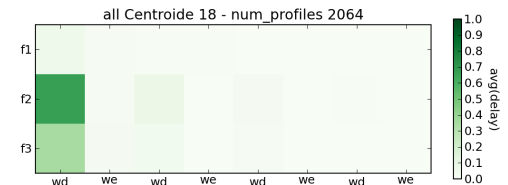
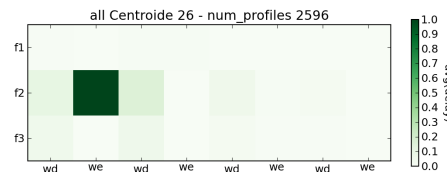
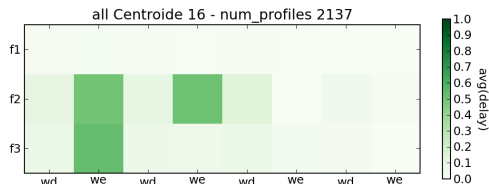
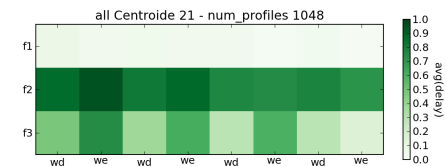
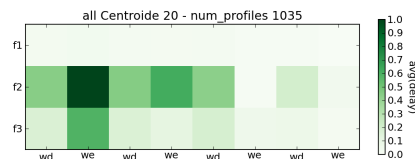


(b)

Sociometer 2.0

Step 2: find representative profiles across all dataset

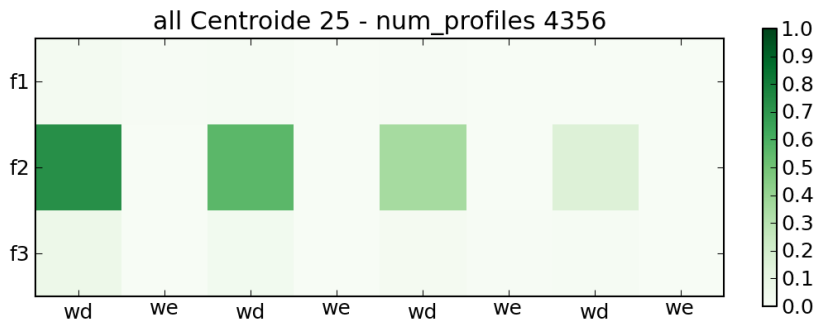
- Based on clustering
 - simple k-means: start with K random representatives, and iteratively refine them
 - in our experiments, k=100
- Output: set of reference (unlabelled) profiles



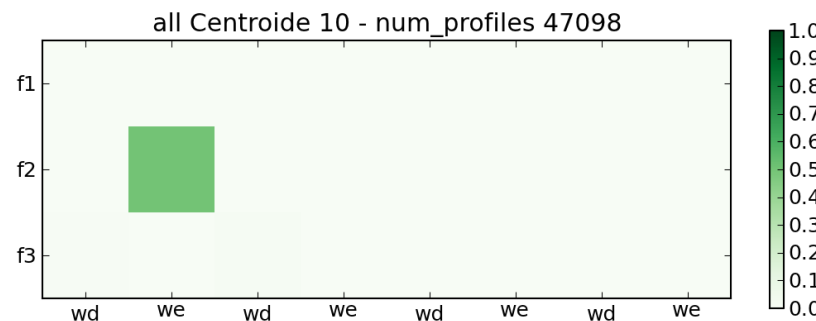
Sociometer 2.0

Step 3: associate representative profiles to categories

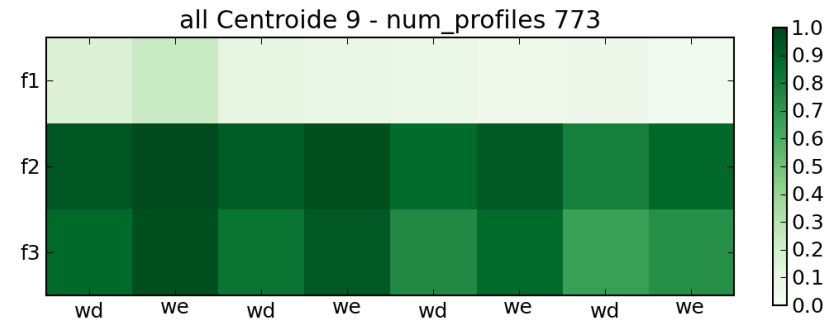
- Manual labelling
 - Use fuzzy rules, difficult to formalize
 - Crisp classification, no weights (reliability of labels)



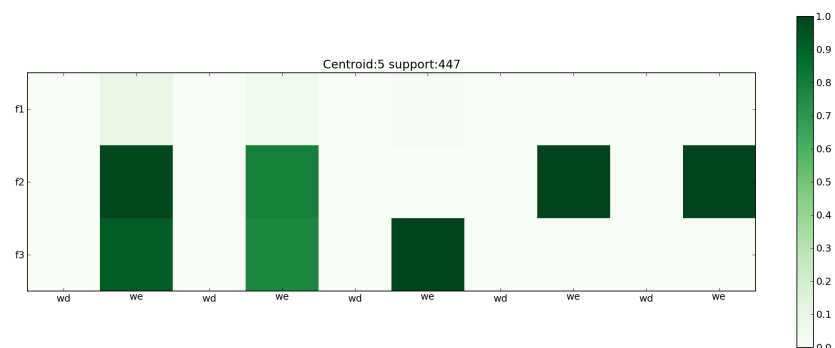
Commuter



Occasional



“Static” resident

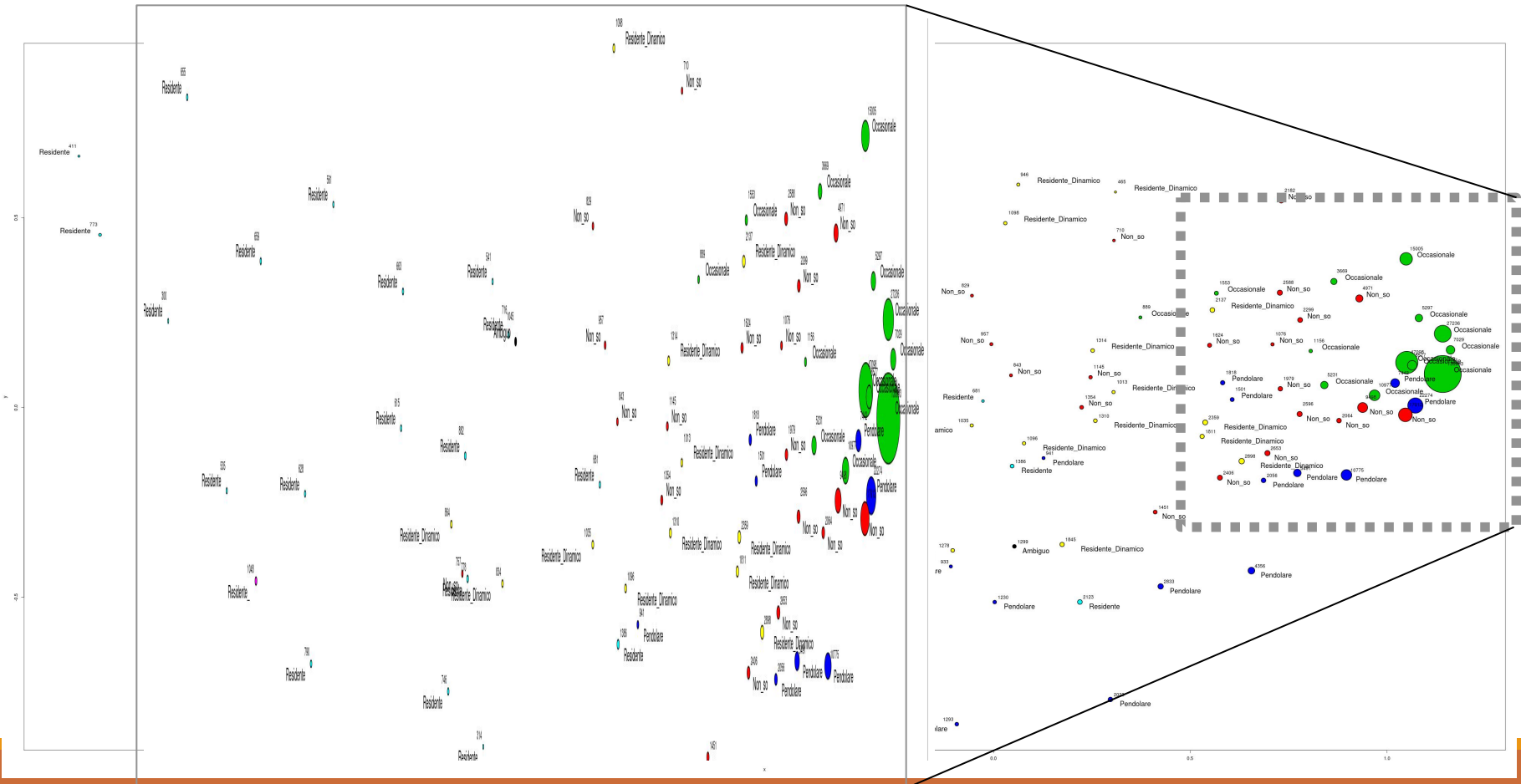


“Dynamic” resident

Sociometer 2.0

Step 3bis: consistency check / labels distribution / fix bugs

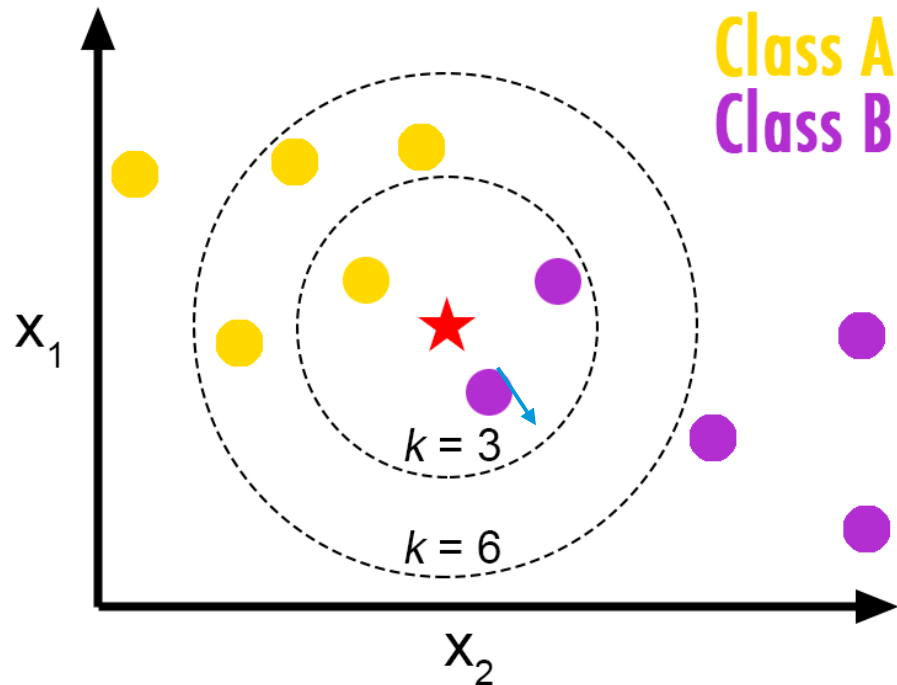
- Profiles (individual and representative) are 24-dimensional
- MDS (24 → 2) to visualize them



Sociometer 2.0

Step 4: label propagation

- Simple k-NN classification, $k=1$
 - Associates each individual profile to the closest representative profile
- So far, no voting schema ($k>1$) was used



Sociometer 2.0

Step 5: aggregate into presence stats and O/D flows

- Presence aggregates
 - Residents = Static + Dynamic residents
- Kind of flows represented:
 - Dynamic residence → sites of commuting
 - Dynamic residence → sites of occasional visits

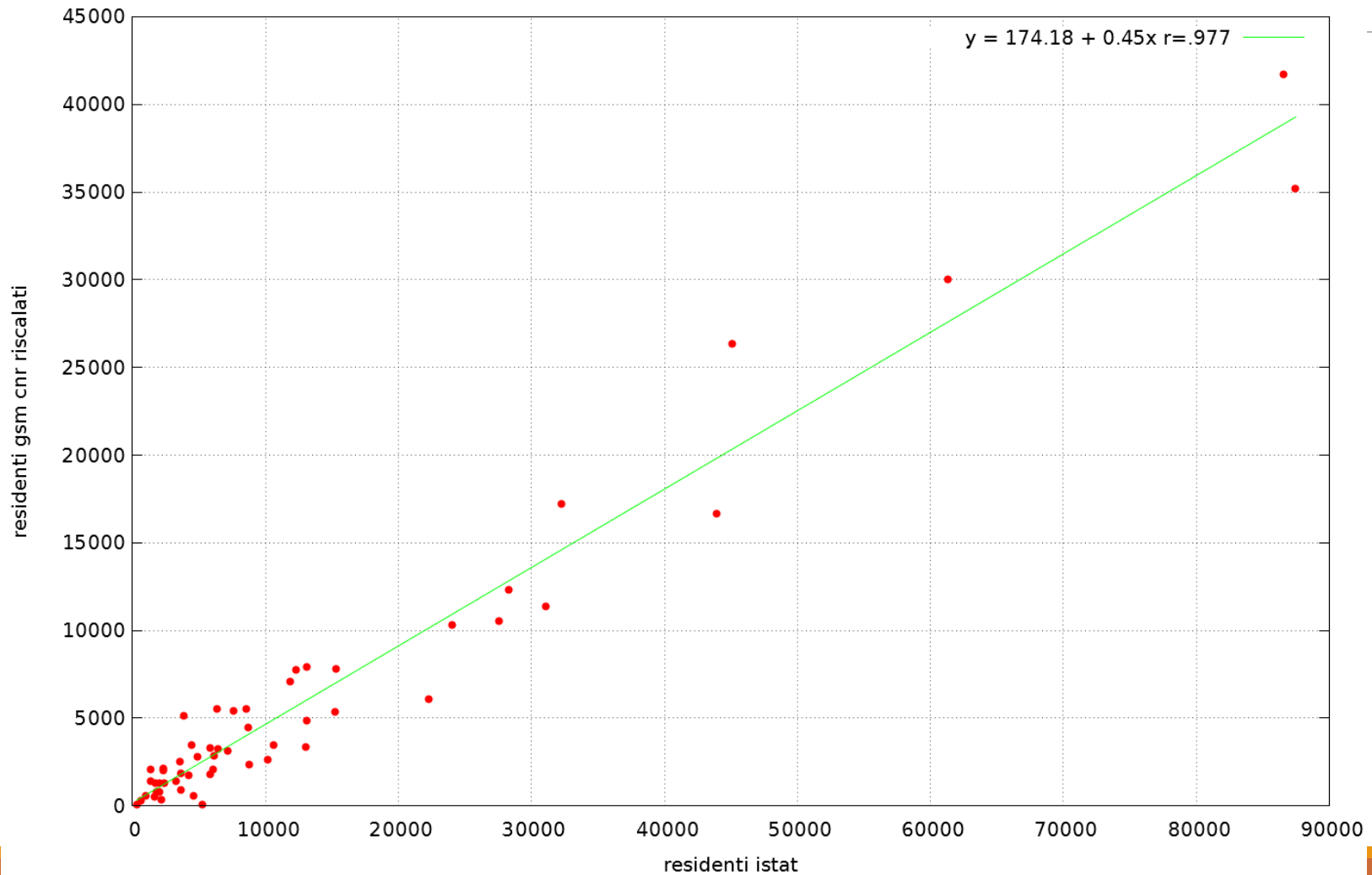
ISTAT Persons & Places project



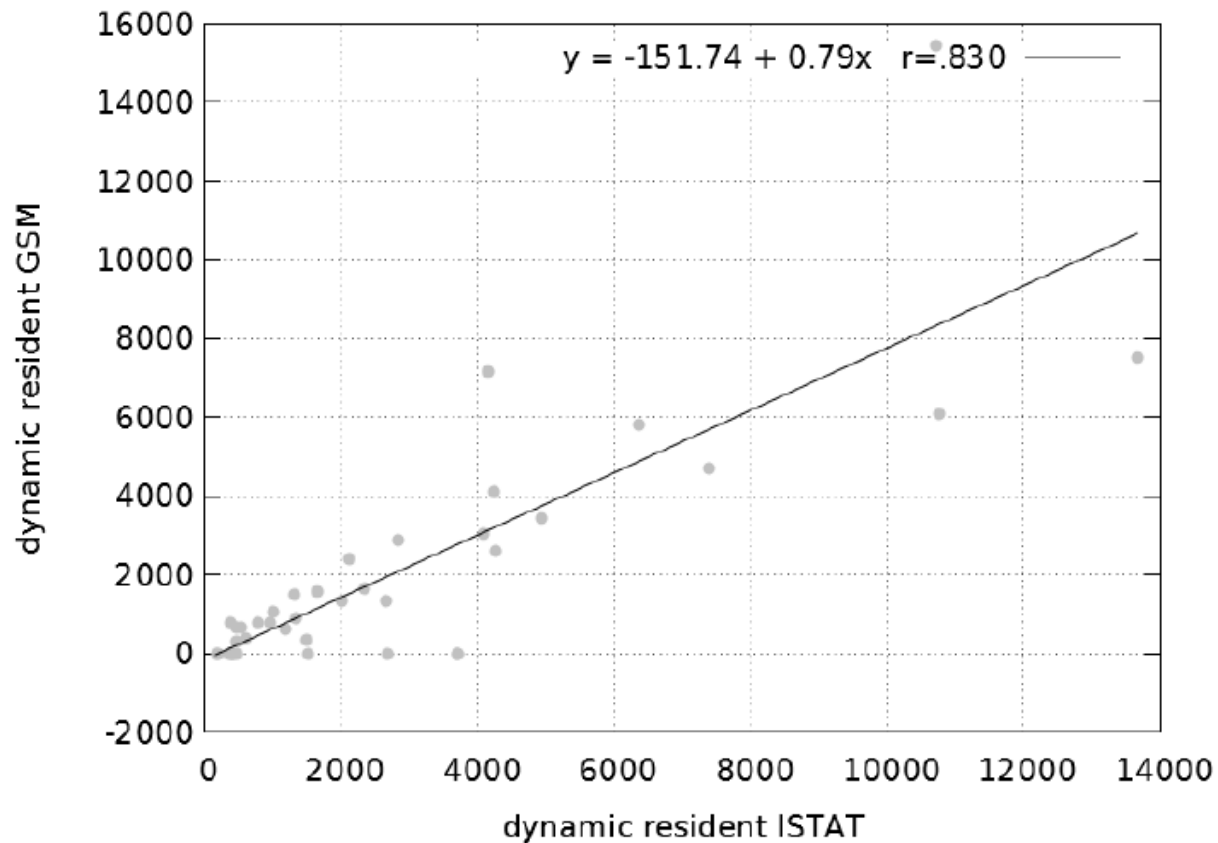
- Ultimate goal: Use Big GSM data to
 - Estimate user categories on a given territory
 - Infer O/D matrix across municipalities
- Goal of this project:
 - Apply/adapt GSM-based user categorization (Sociometer) on municipalities of a large territory
 - Infer partial O/D matrix
 - Direct/Indirect comparison against official data
- GSM 4-weeks Dataset on Pisa and Lucca provinces

Static residents GSM

Correlazione residenti GSM riscalati residenti ISTAT

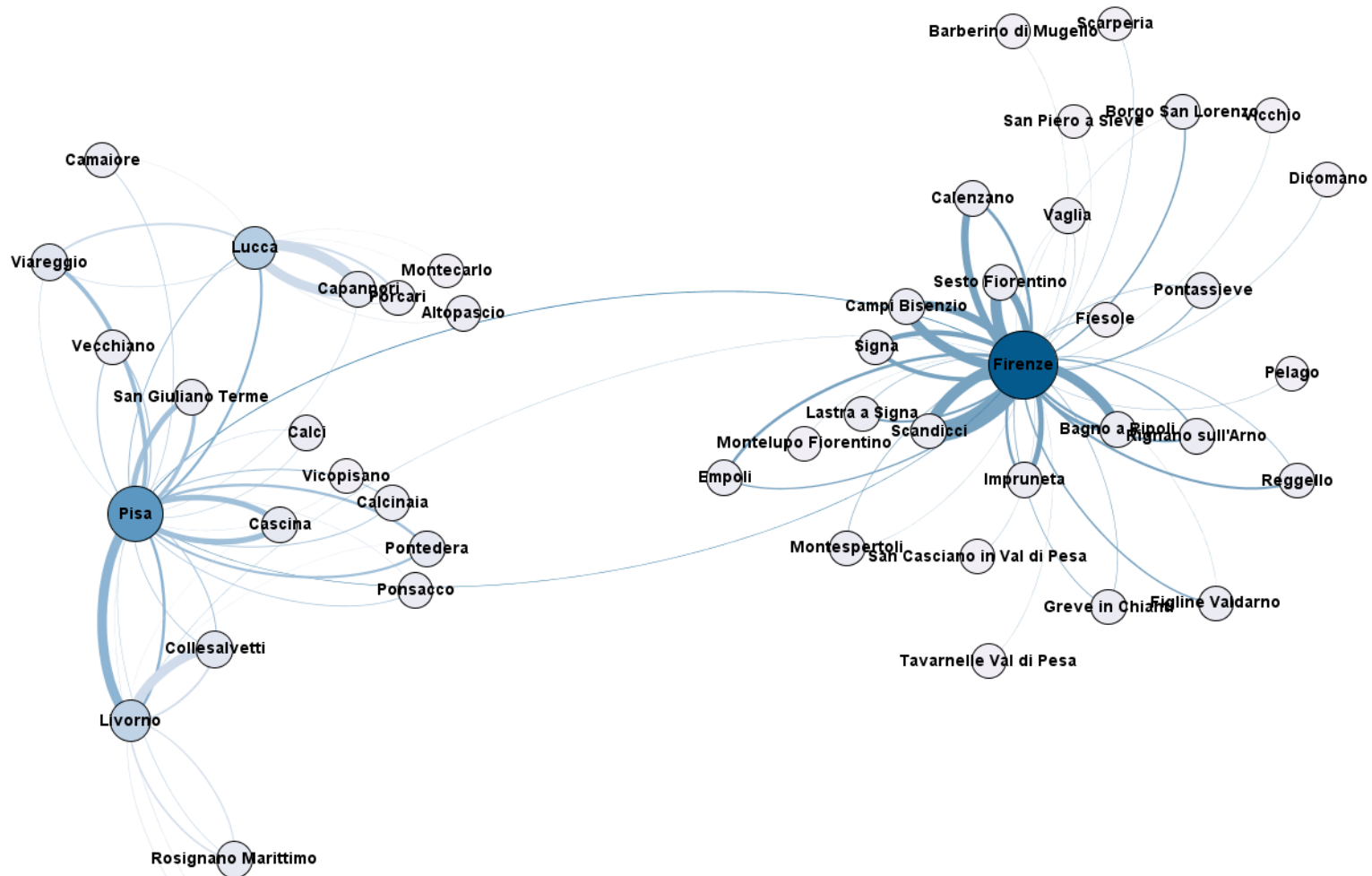


Dynamic residents (outgoing)



Sample results / 1

Home-Work

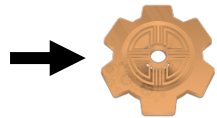
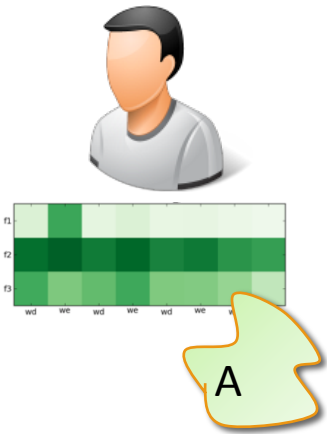


The background features a dense, overlapping pattern of stylized human figures and circles in various colors including brown, purple, green, yellow, blue, pink, and grey. The figures are composed of simple geometric shapes like rectangles and circles, creating a sense of a diverse crowd. The text 'REAL TIME DEMOGRAPHY' is centered over this pattern in a large, black, sans-serif font. A thin horizontal line is positioned directly below the text.

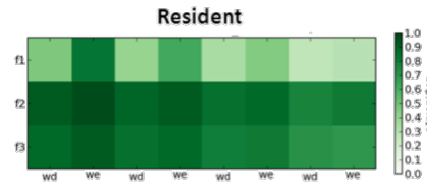
REAL TIME DEMOGRAPHY

Sociometer with Mobile Phone Data.

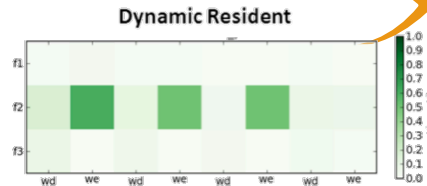
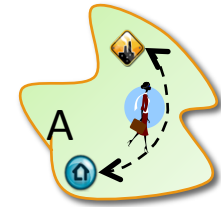
Users' Call Profiles



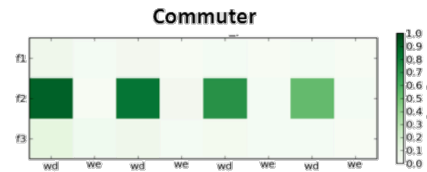
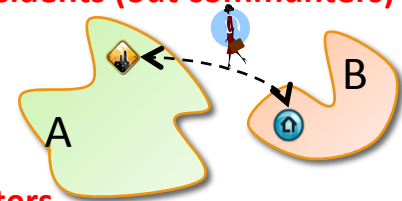
Classification Algorithm



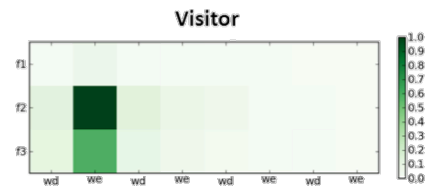
Residents



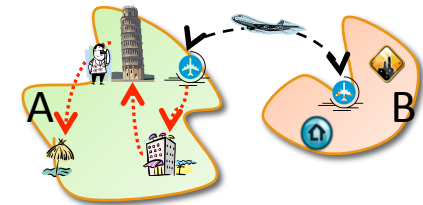
Dynamic Residents (out commuters)



(in) Commuters

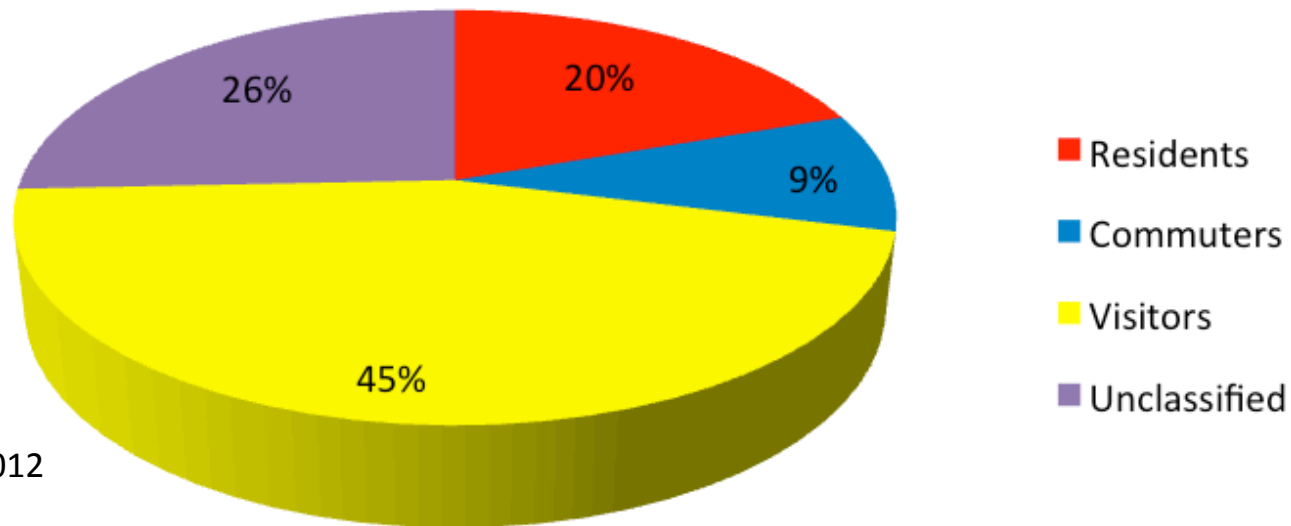


Visitors



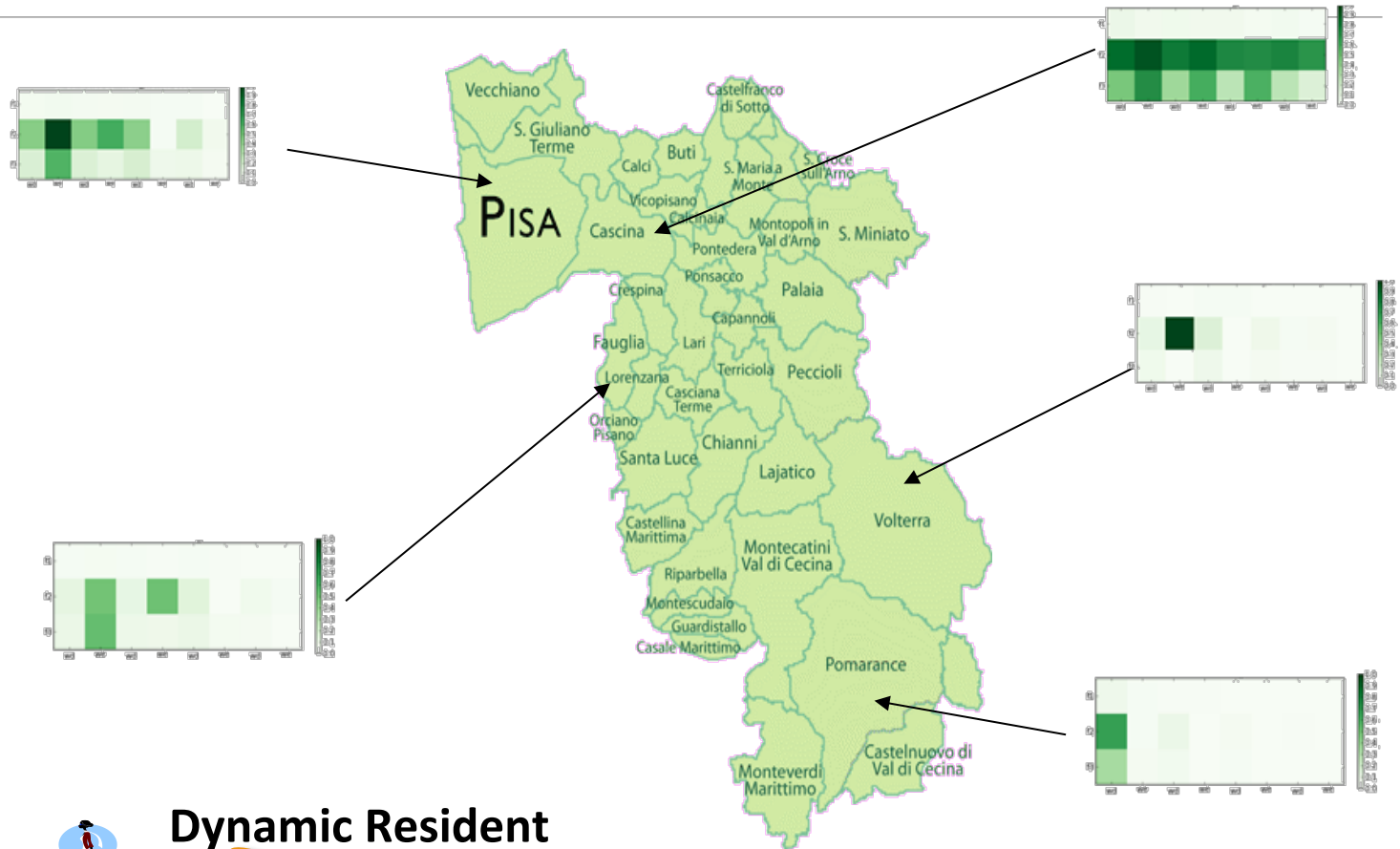
Sociometer: the city user meter

Classification outcome



Pisa, January 2012

The many profiles of an individual

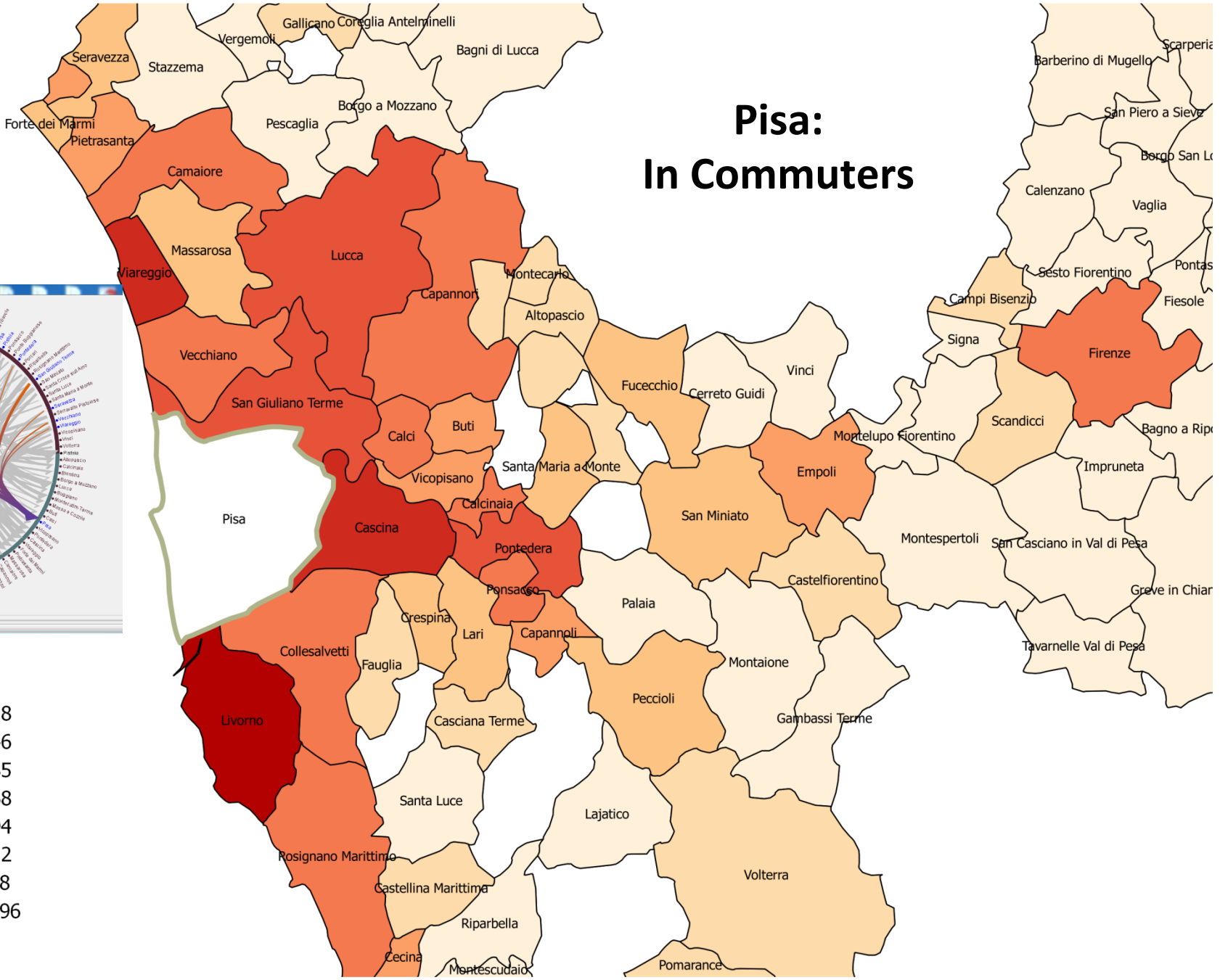
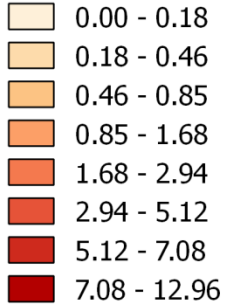
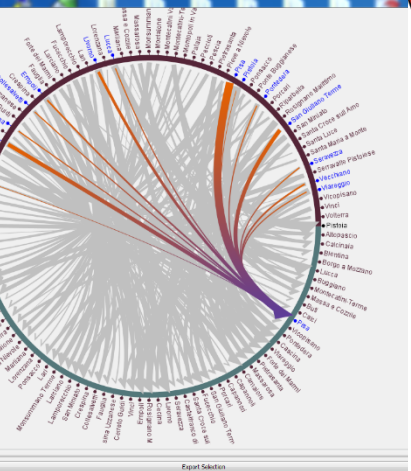


Commuter

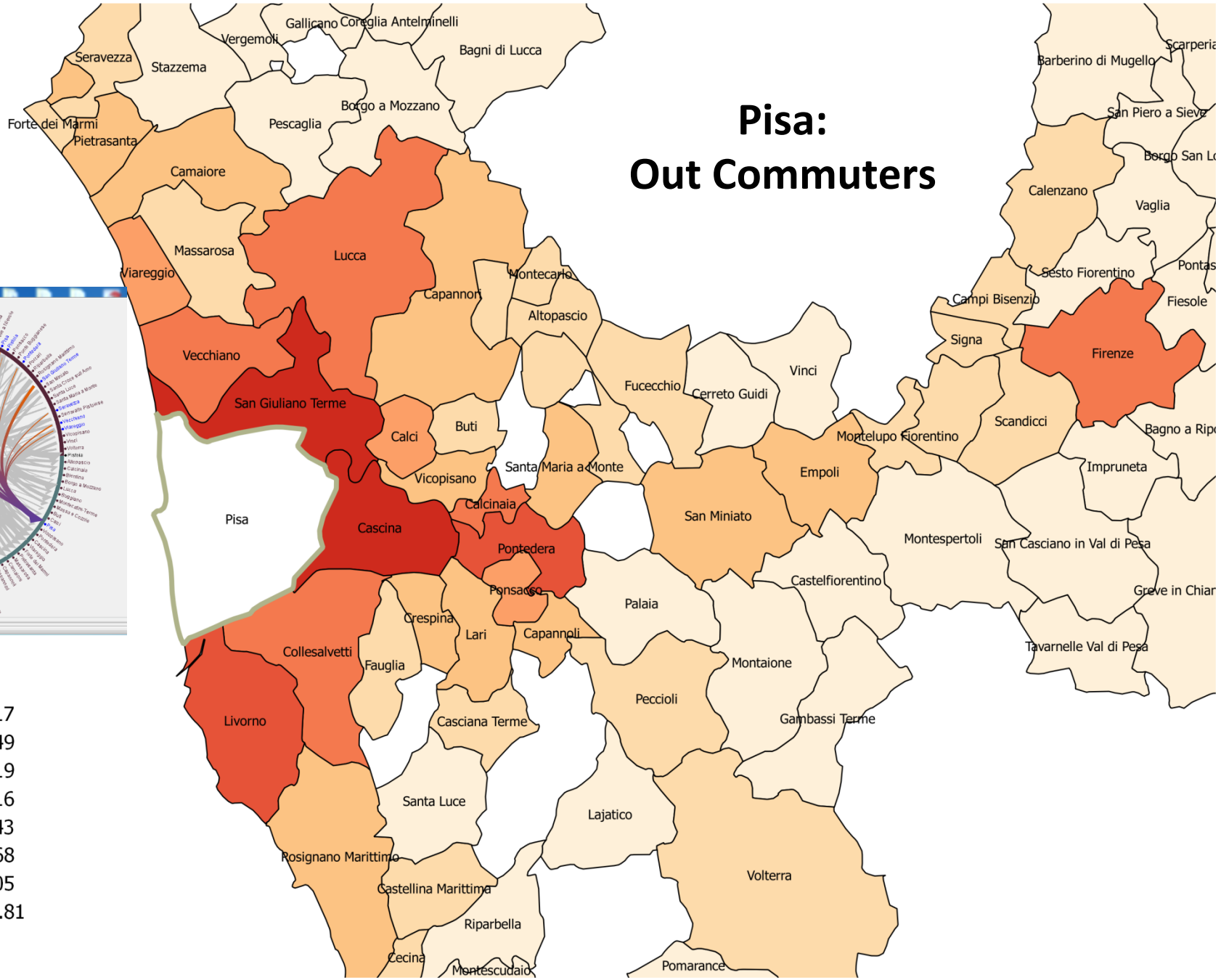
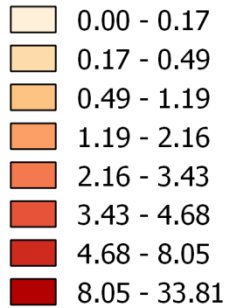
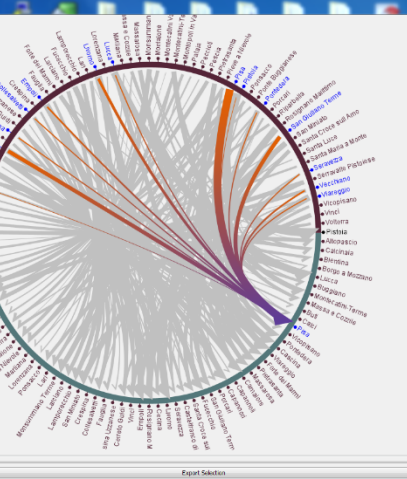
Dynamic Resident



Pisa: In Commuters



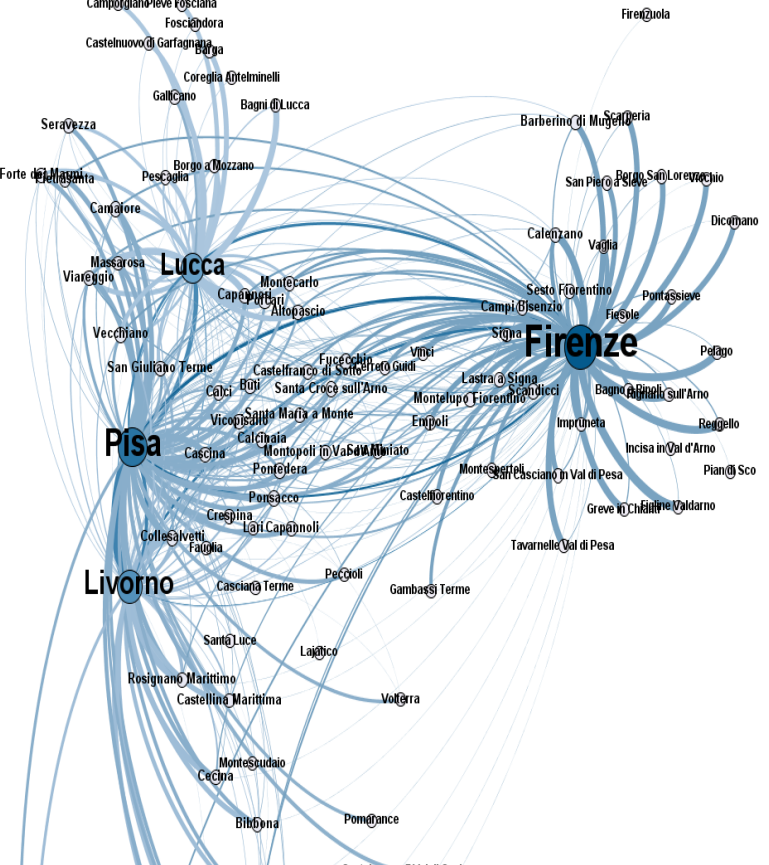
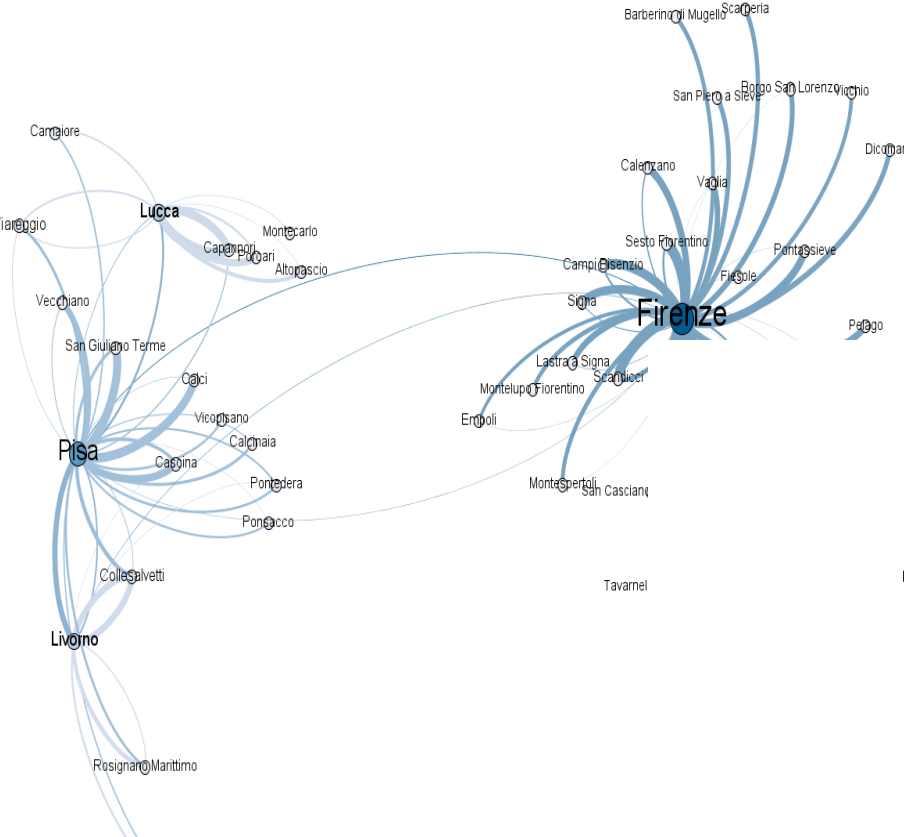
Pisa: Out Commuters



INTER-CITY FLOWS, WITH SEMANTICS

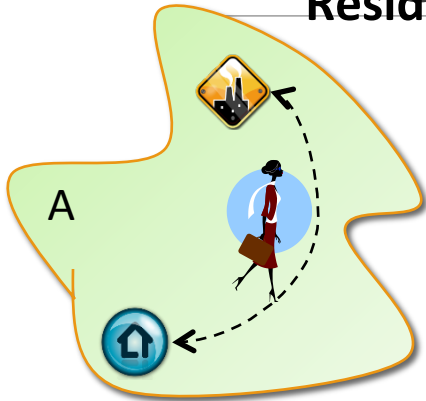
COMMUTER NETWORK

VISITOR NETWORK

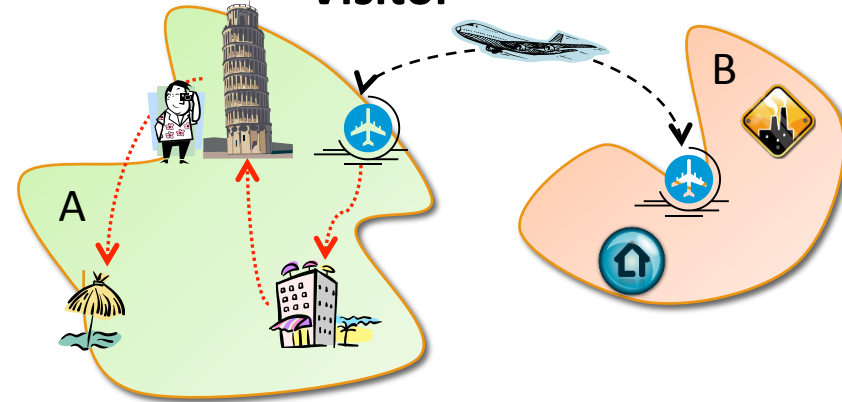


Sociometer: Estimating User Category from mobile phone

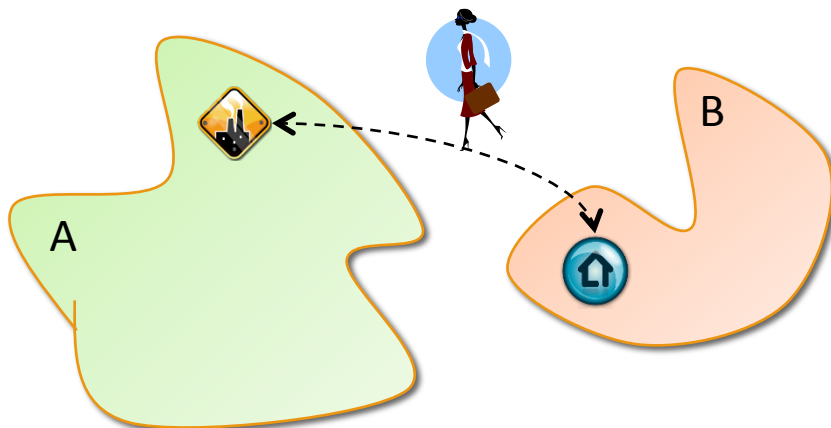
Resident



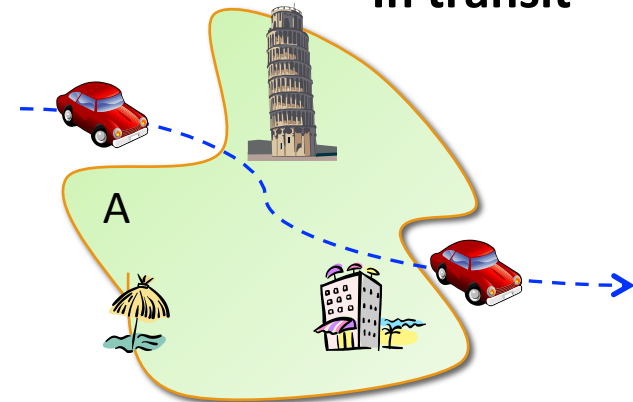
Visitor



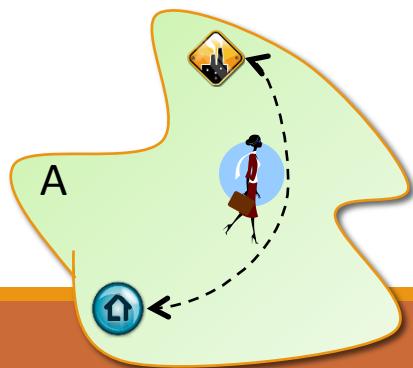
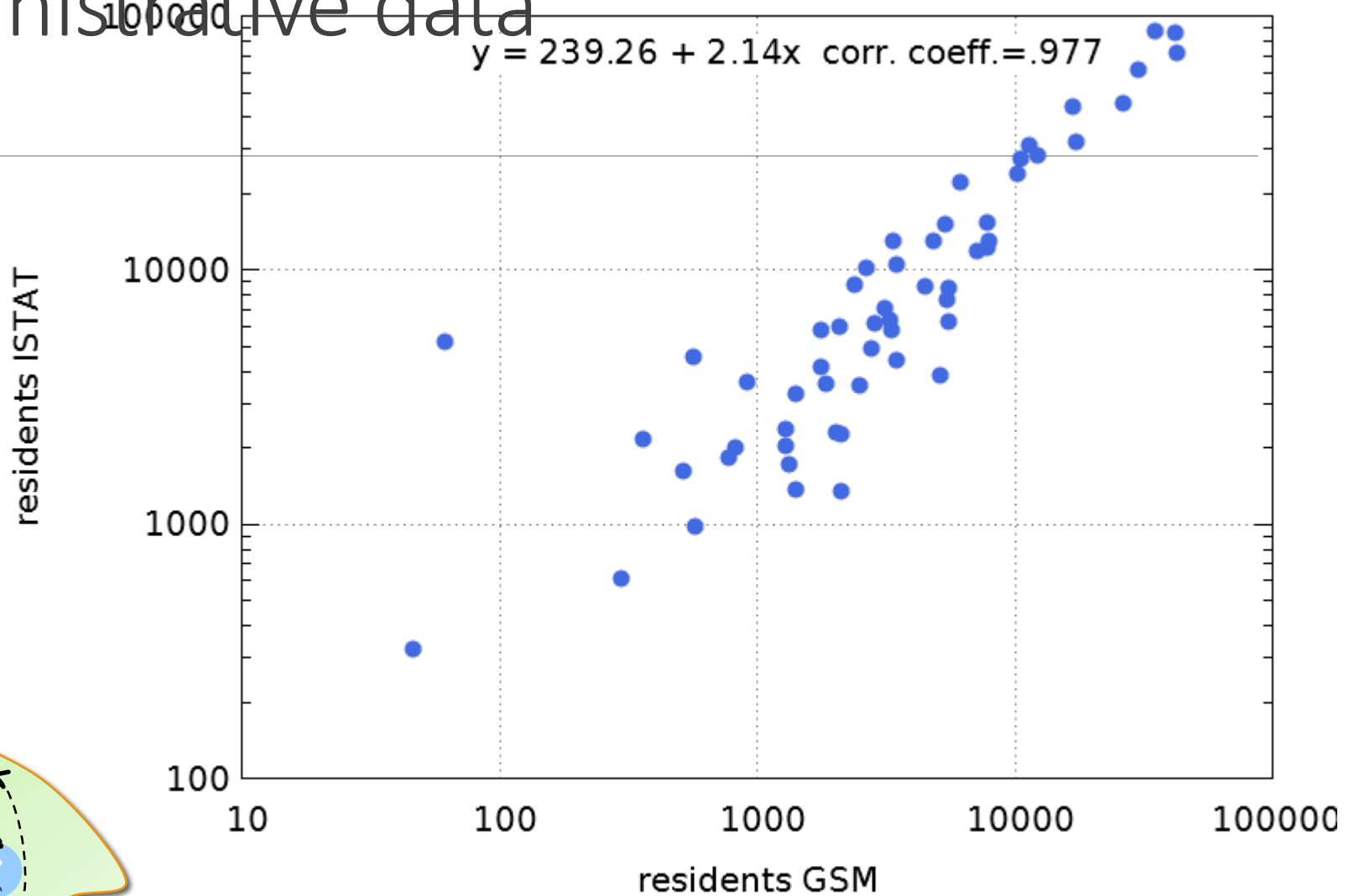
Commuter



In transit

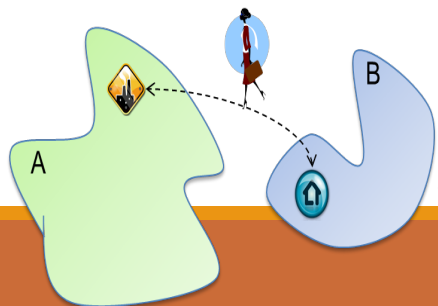
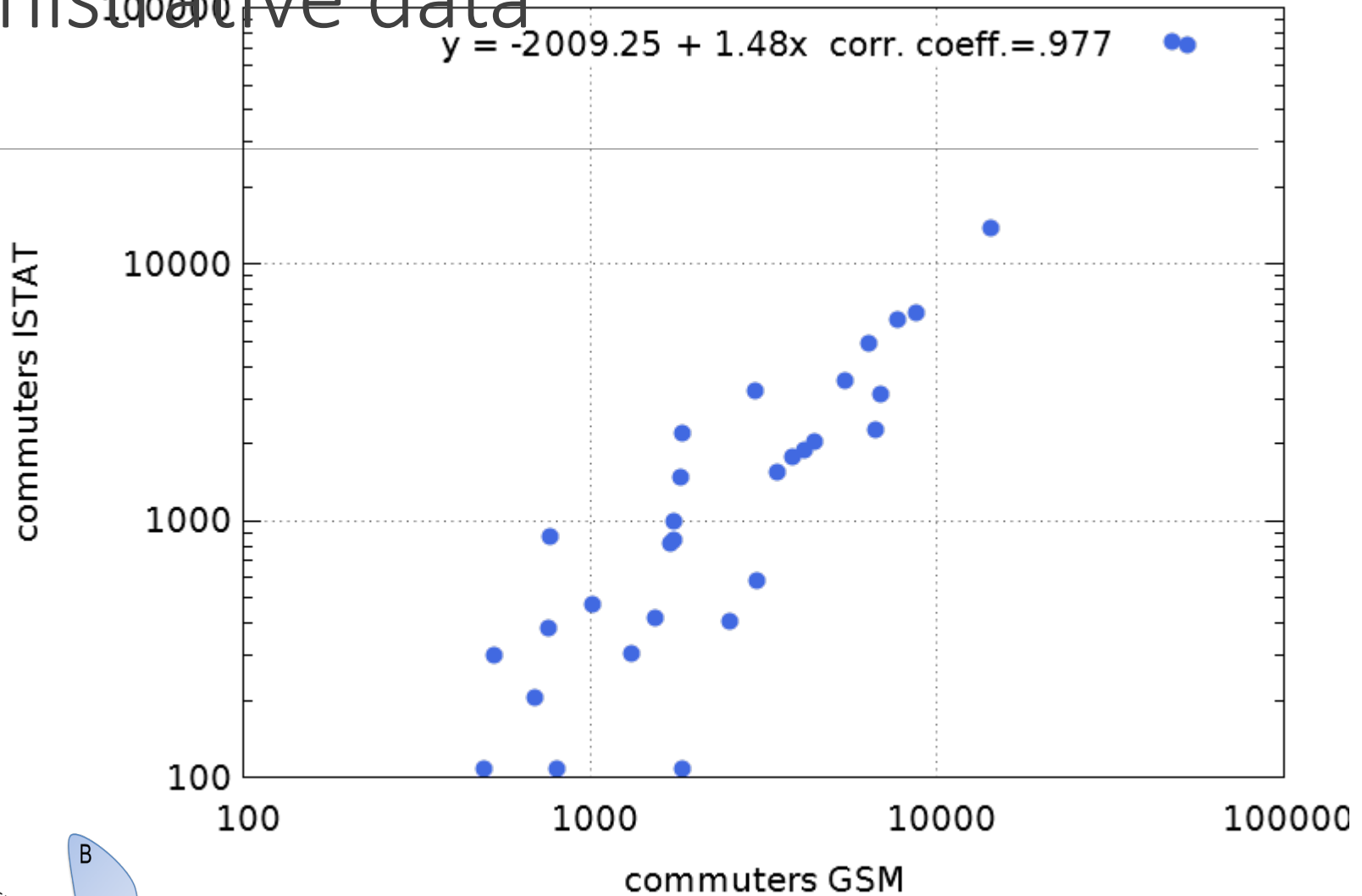


Residents – validation with administrative data



Join work with ISTAT: Barbara Furletti, Lorenzo Gabrielli, Giuseppe Garofalo, Fosca Giannotti, Letizia Milli, Mirco Nanni, Dino Pedreschi, Roberta Vivio. Use of mobile phone data to estimate mobility flows. Measuring urban population and intercity mobility using big data in an integrated approach. Italian Symposium on Statistics (2014).

Commuters - Validation with administrative data



Join work with ISTAT: Barbara Furletti, Lorenzo Gabrielli, Giuseppe Garofalo, Fosca Giannotti, Letizia Milli, Mirco Nanni, Dino Pedreschi, Roberta Vivio. Use of mobile phone data to estimate mobility flows. Measuring urban population and intercity mobility using big data in an integrated approach. Italian Symposium on Statistics (2014).

Measuring exceptional events

- Presences during Jubilee in Rome (December 2015)
- Continuous monitoring

Call Data Records

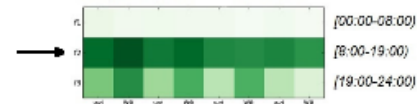
(UserID, Cell, Timestamp)

1256, Cell1, 12/03/2015 09:10

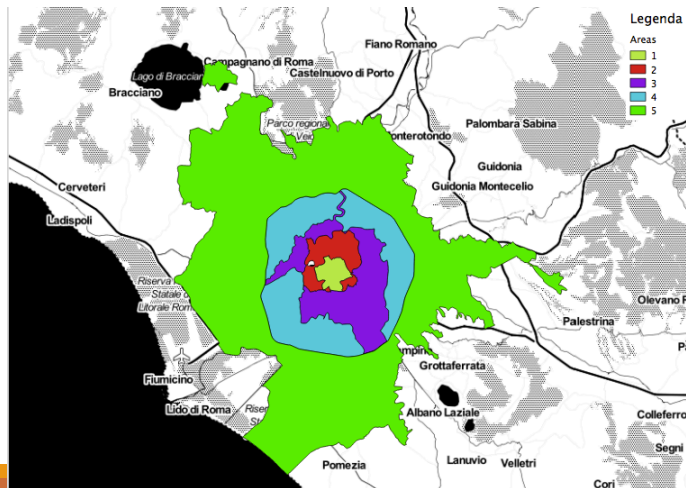
1256, Cell2, 12/03/2015 18:45

...

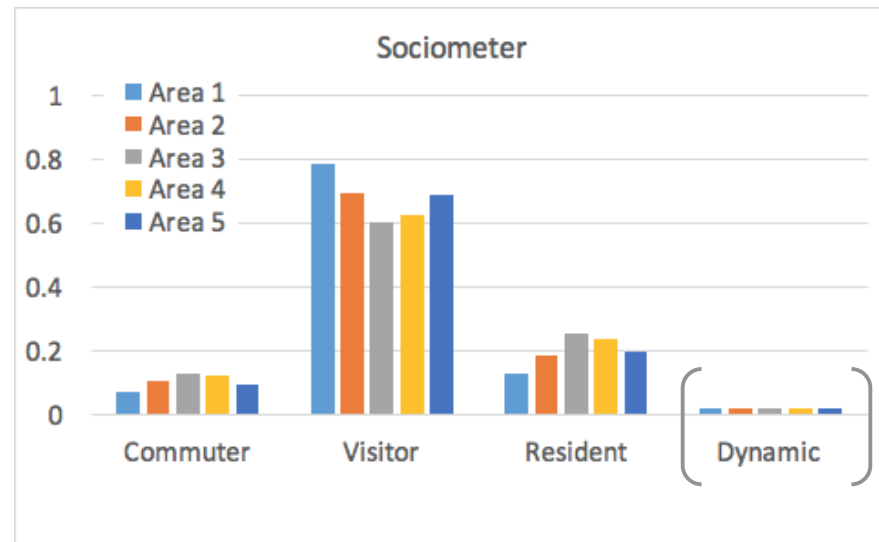
1256, Cell1, 18/03/2015 23:03



+



=



San Pietro Square

