

How to use #DataScience  
and #BigData to nowcast  
well-being?

How to use #DataScience  
and #BigData to nowcast  
**well-being?**

Wealth

Health

Security

Education

Life balance

Environment

Happiness



well-being

Wealth  
Health  
Security  
Education  
Life balance  
Environment  
Happiness



Retail market  
Phone data  
Eating habits  
Social media  
Scientific prod.  
Fitness  
Sociality  
Crimes  
Pollution

Wealth

Health

Security

Education

Life balance

Environment

Happiness

Retail market

Phone data

Eating habits

Social media

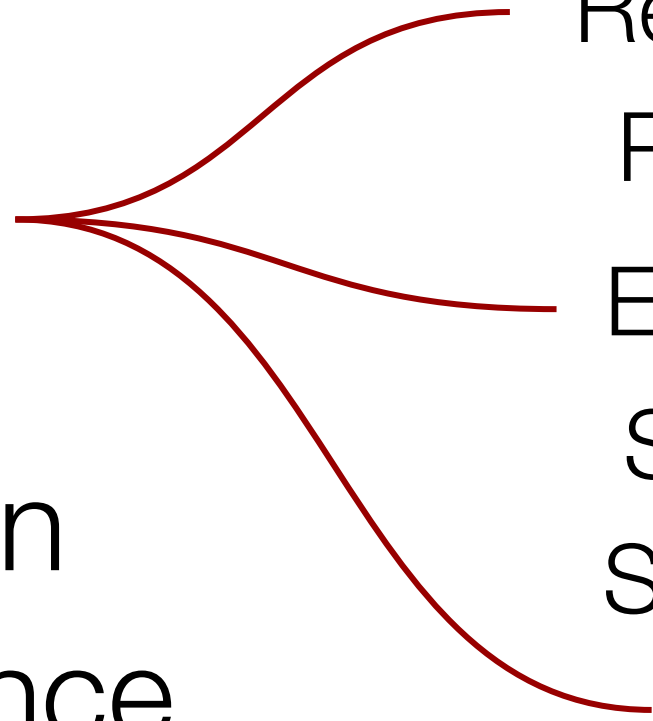
Scientific prod.

Fitness

Sociality

Crimes

Pollution



Wealth

Health

Security

Education

Life balance

Environment

Happiness

Retail market

Phone data

Eating habits

Social media

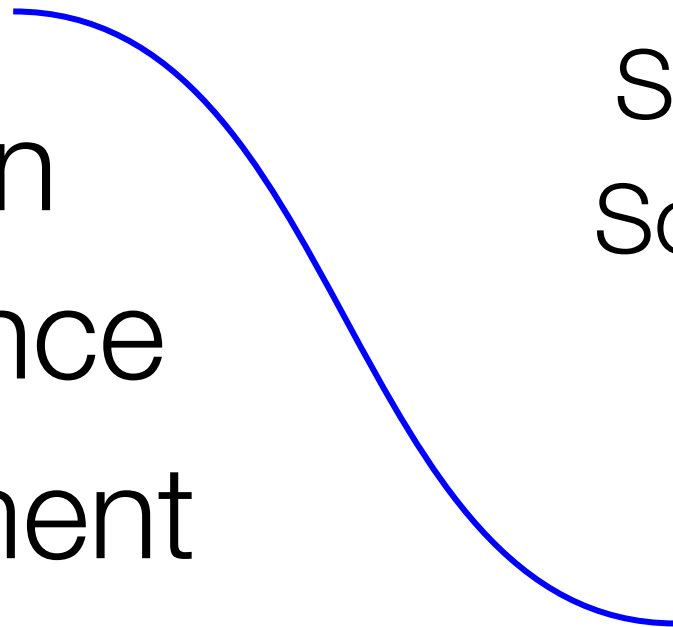
Scientific prod.

Fitness

Sociality

Crimes

Pollution



Wealth

Health

Security

Education

Life balance

Environment

Happiness

Retail market

Phone data

Eating habits

Social media

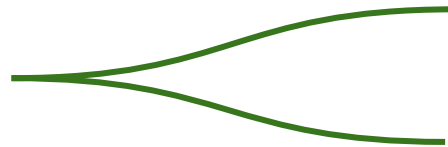
Scientific prod.

Fitness

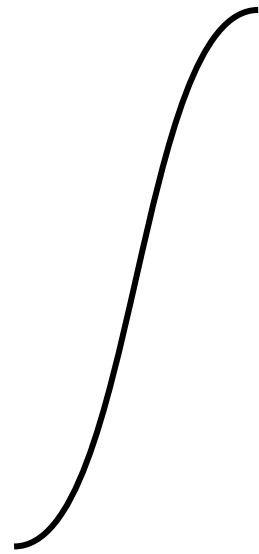
Sociality

Crimes

Pollution



Wealth  
Health  
Security  
Education  
Life balance  
Environment  
Happiness



Retail market  
Phone data  
Eating habits  
Social media  
Scientific prod.  
Fitness  
Sociality  
Crimes  
Pollution



Wealth  
Health  
Security  
Education  
Life balance  
Environment  
Happiness



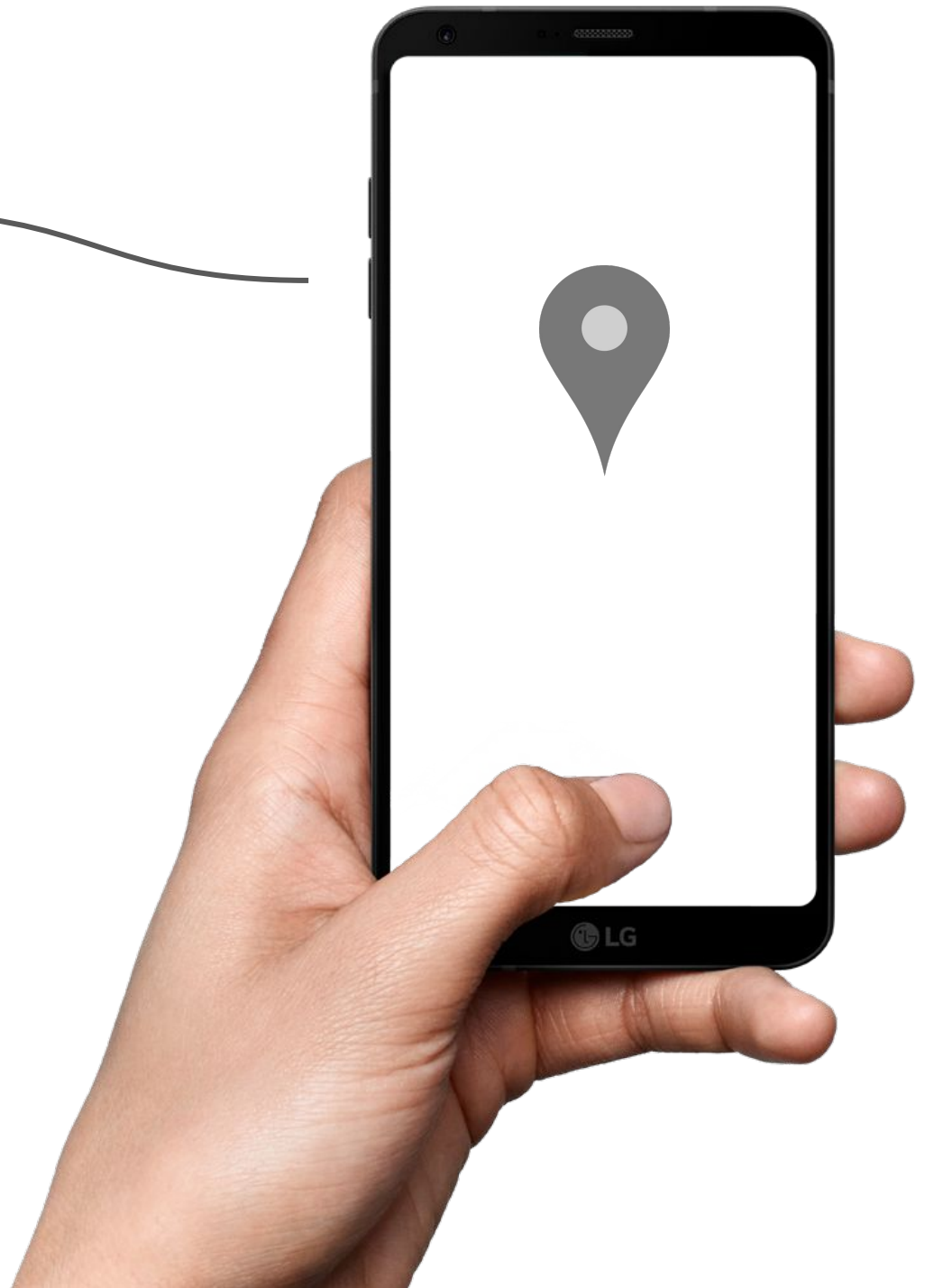
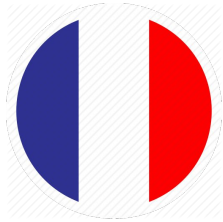
Retail market  
Phone data  
Eating habits  
Social media  
Scientific prod.  
Fitness  
Sociality  
Crimes  
Pollution

Wealth  
Health  
Security  
Education  
Life balance  
Environment  
Happiness



Retail market  
Phone data  
Eating habits  
Social media  
Scientific prod.  
Fitness  
Sociality  
Crimes  
Pollution

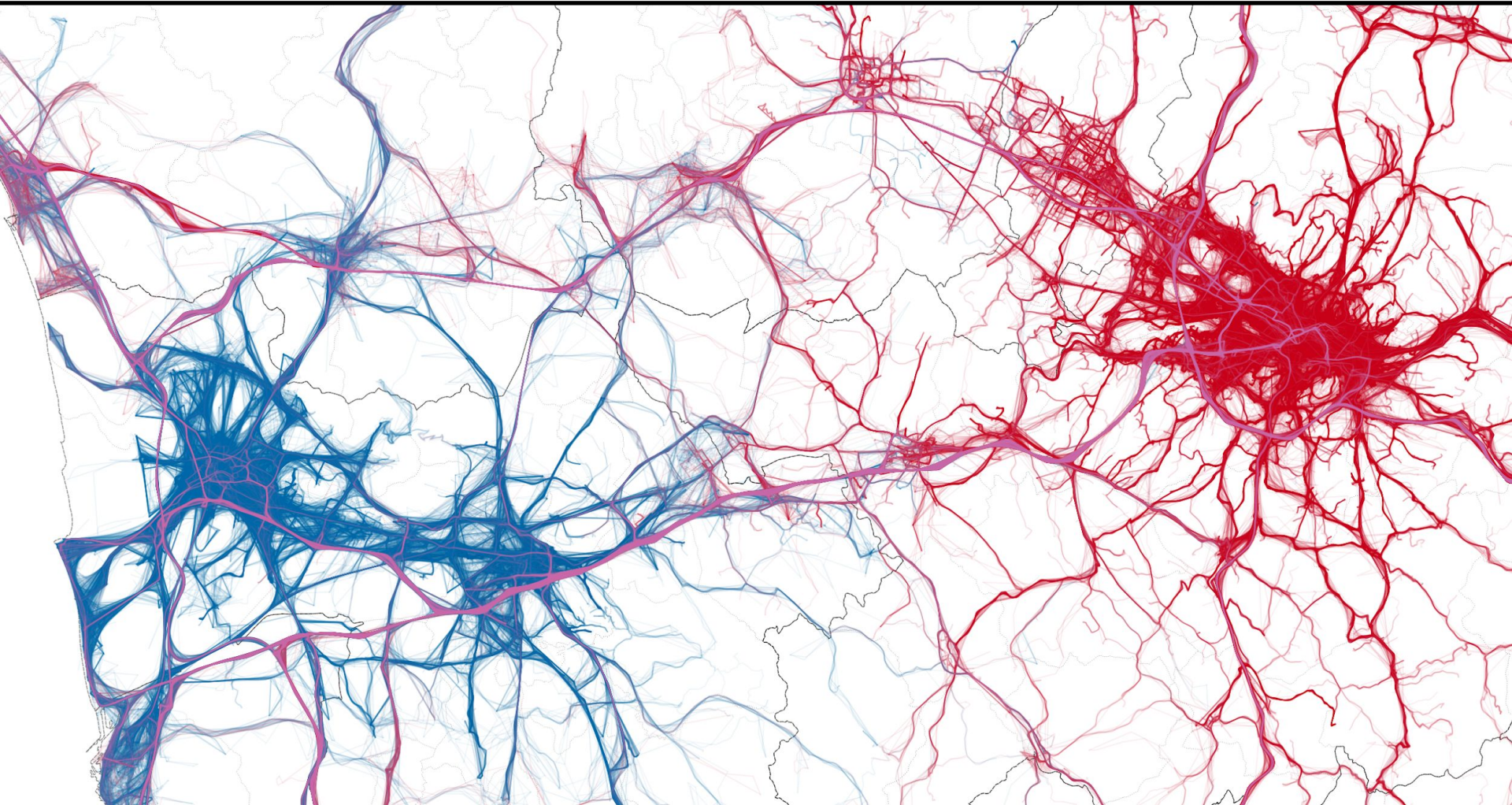
# Wealth



# Surveys have limits

- 1.** They do not scale to small territories and they are not dynamic
- 2.** They do not adapt to Pareto distributions:
  - We cannot use samples
  - 50 people are as wealthy as 3.5 billion people
- 3.** Young people do not respond to surveys

# Big Data

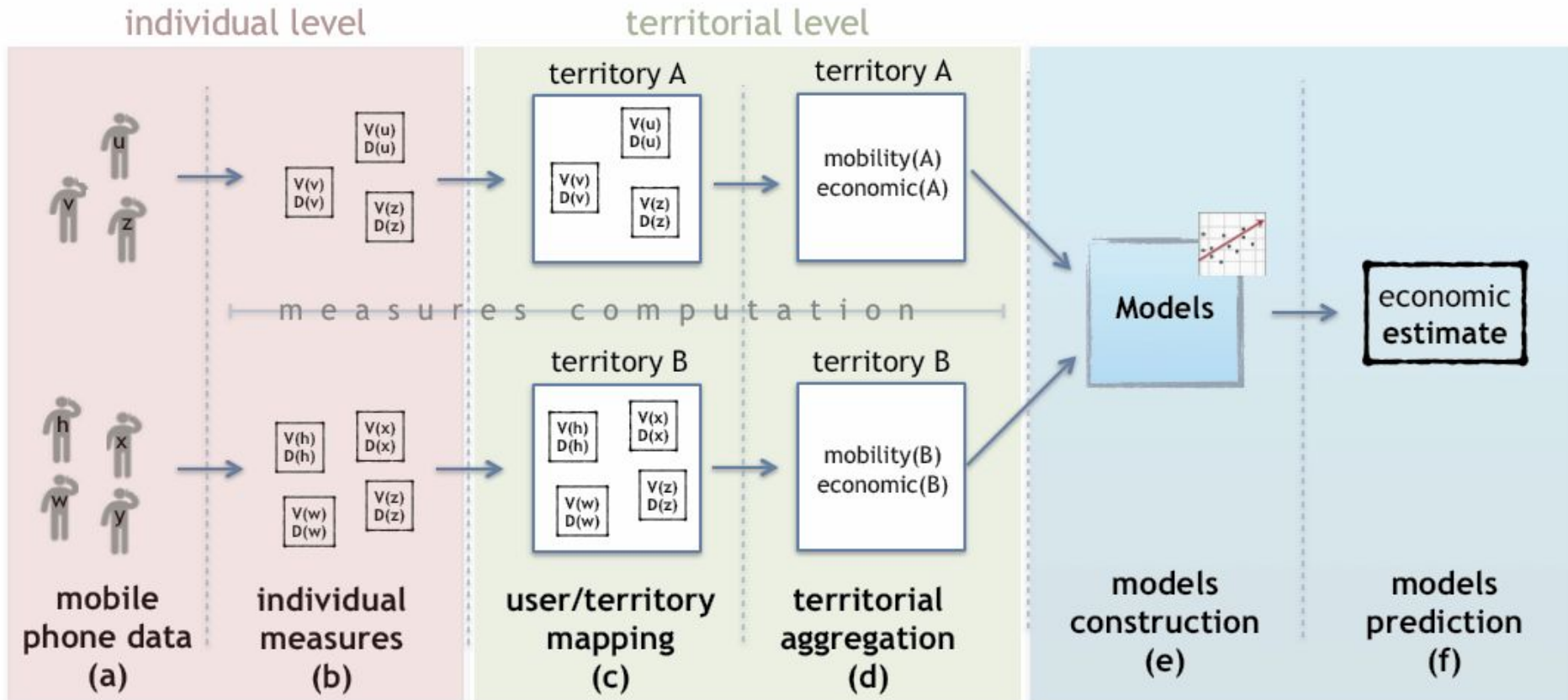


# Advantages of Big Data

- They capture the complexity of social systems in their entirety
- They allow for the observation of complex phenomena, like diversity, resilience and equality

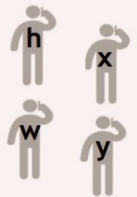


# An analytical framework

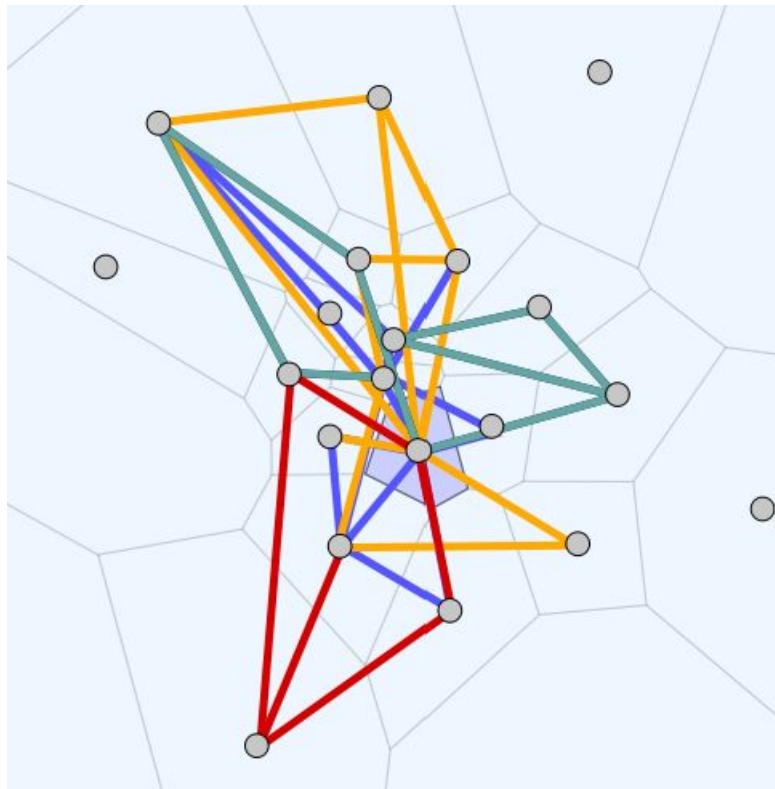


# Data format

timestamp	coords	caller	callee
04/01 23:45:00	132.56, 23.64	145323	452300
04/02 06:02:00	143.28, 54.22	145323	5602
...	...	...	...



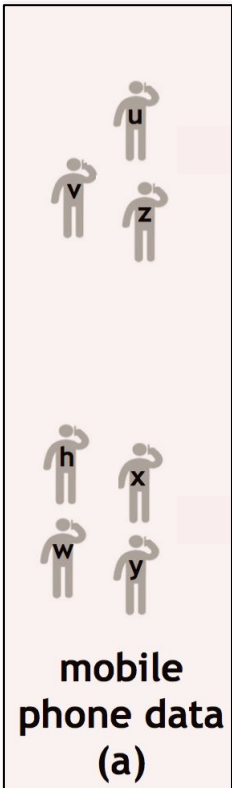
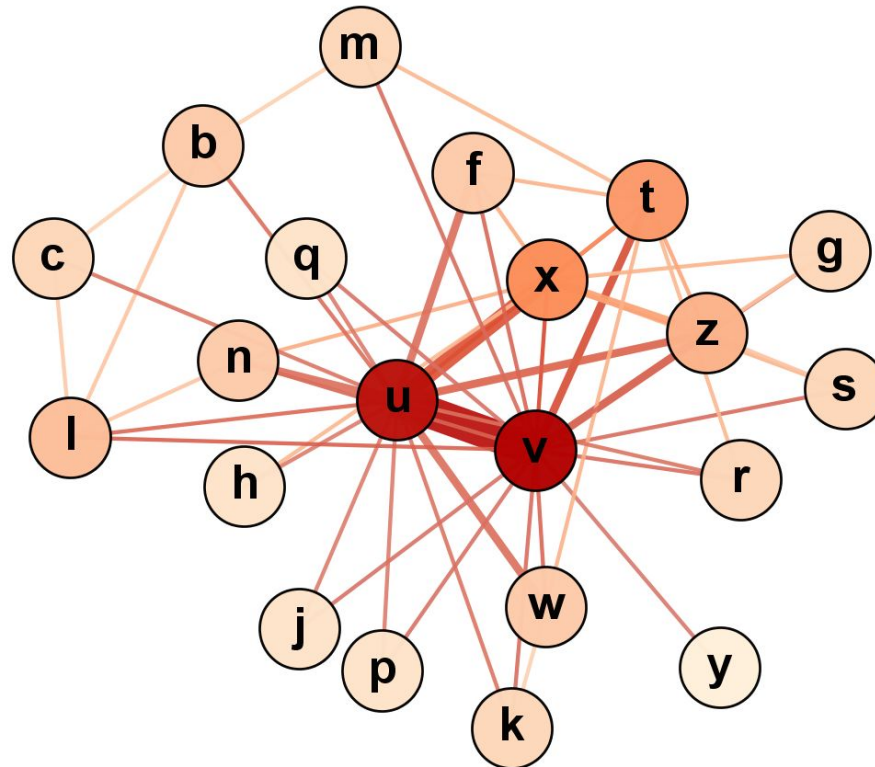
mobile  
phone data  
(a)





# Data format

timestamp	coords	caller	callee
04/01 23:45:00	132.56, 23.64	145323	452300
04/02 06:02:00	143.28, 54.22	145323	5602
...	...	...	...



# Orange dataset



20M users



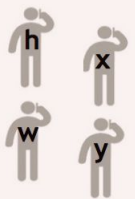
5.7G calls



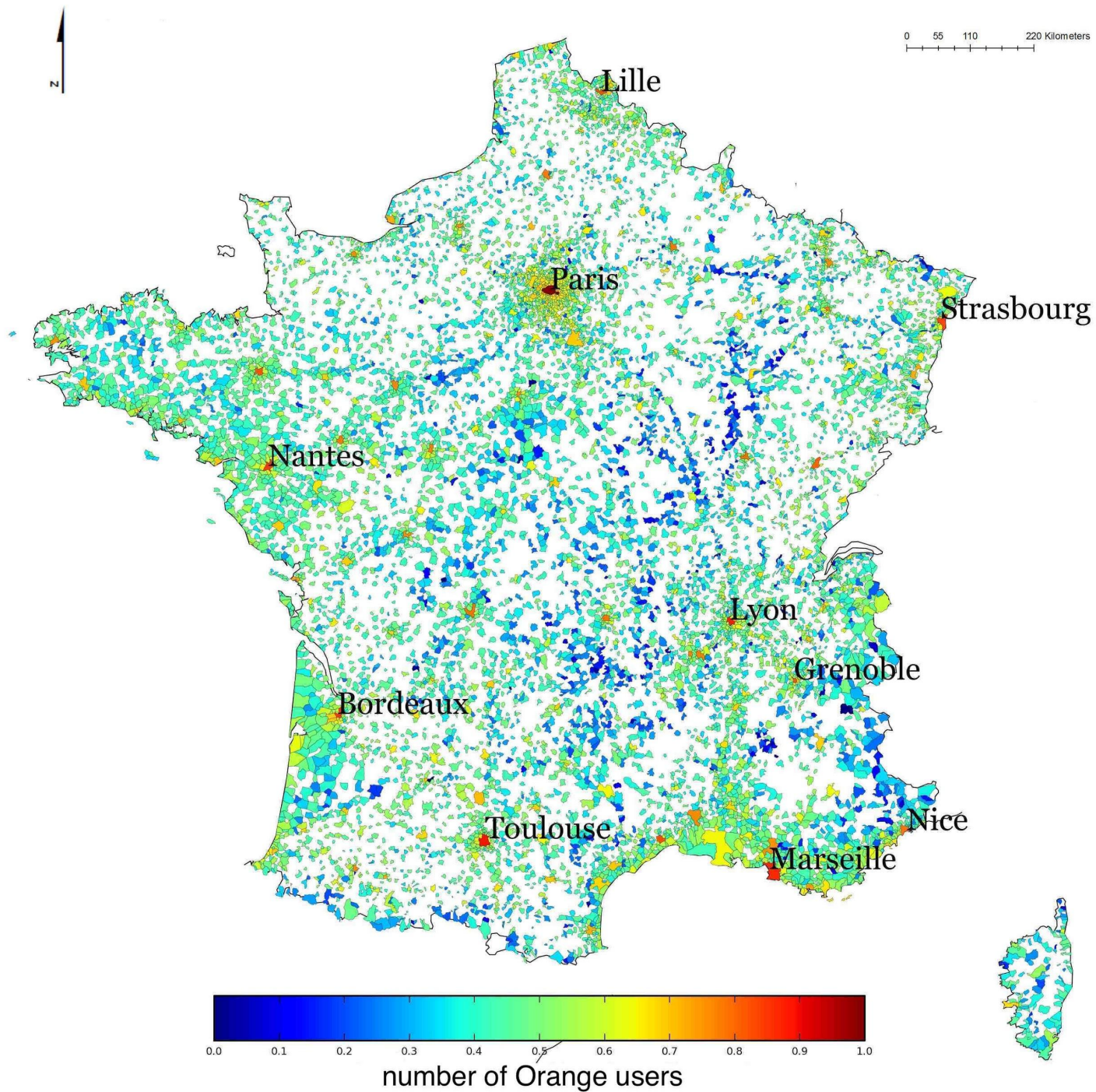
87K towers



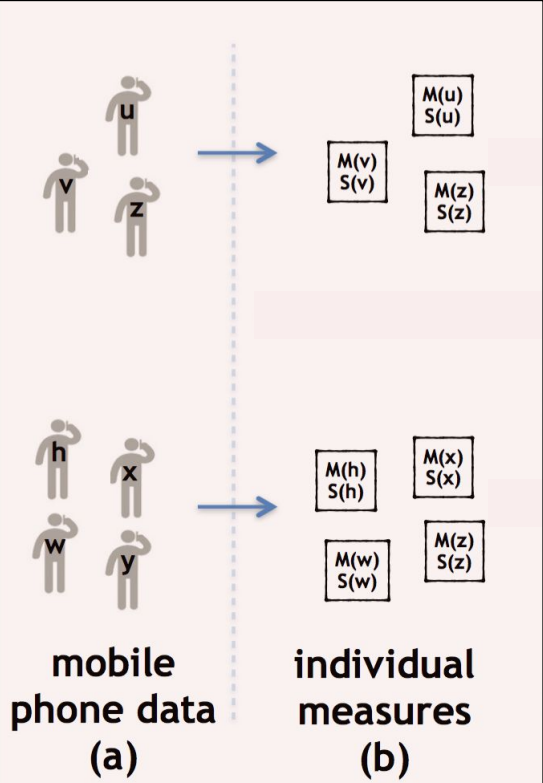
45days



mobile  
phone data  
(a)



individual level



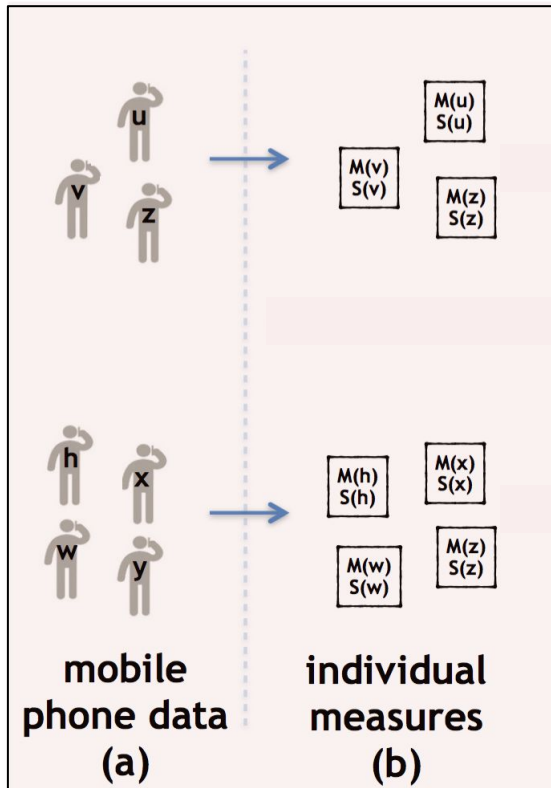








individual level

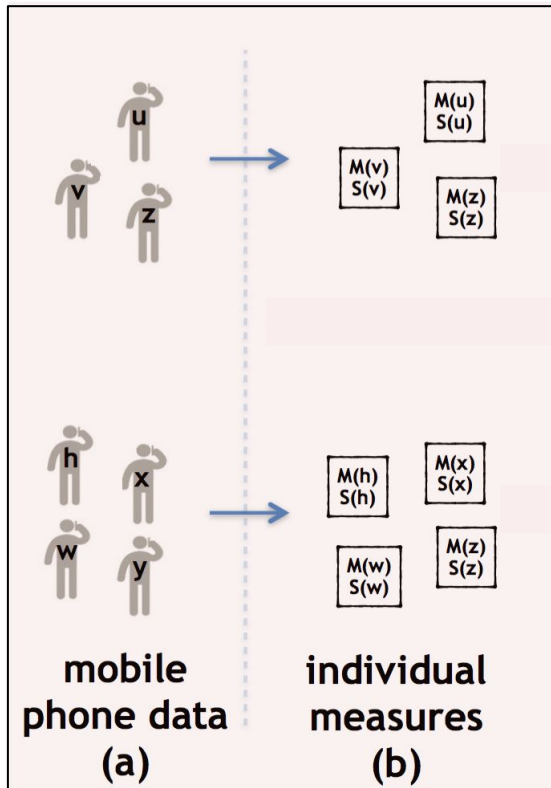


- Mobility volume:  
characteristic distance traveled by  
an individual

$$r_g = \sqrt{\frac{1}{N} \sum_{i=1}^N w_i (r_i - r_{cm})^2}$$

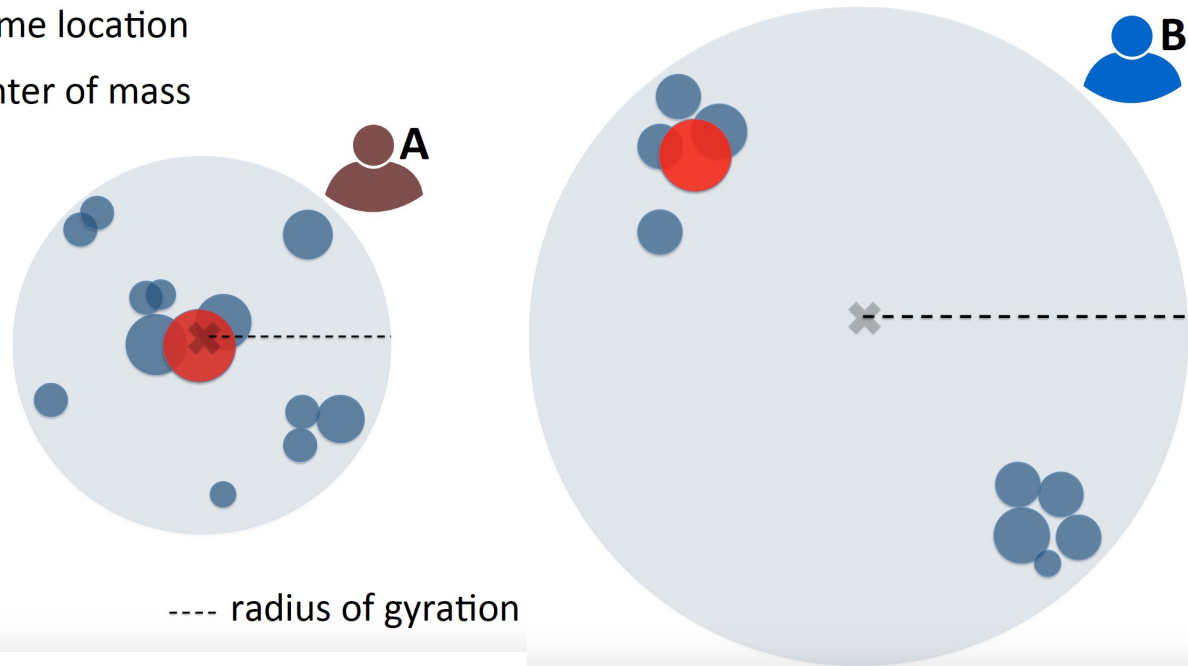
$$r_{cm} = \frac{1}{N} \sum_{i=1}^n w_i r_i \quad N = \sum_{i=1}^n w_i$$

individual level



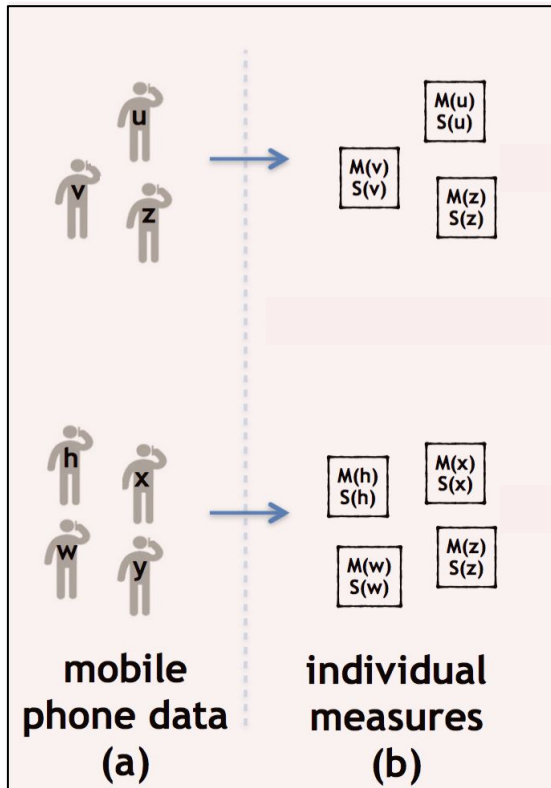
– Mobility volume:  
characteristic distance traveled by  
an individual

- home location
- ✕ center of mass





individual level



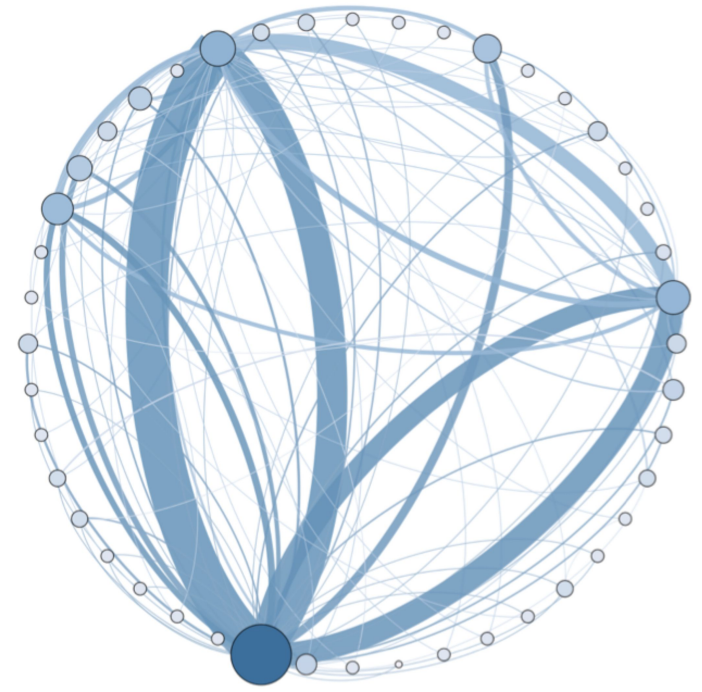
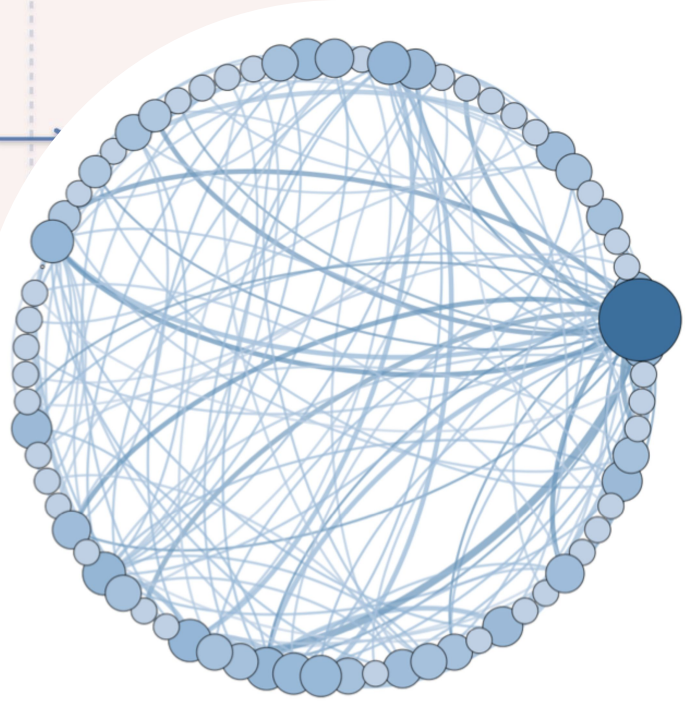
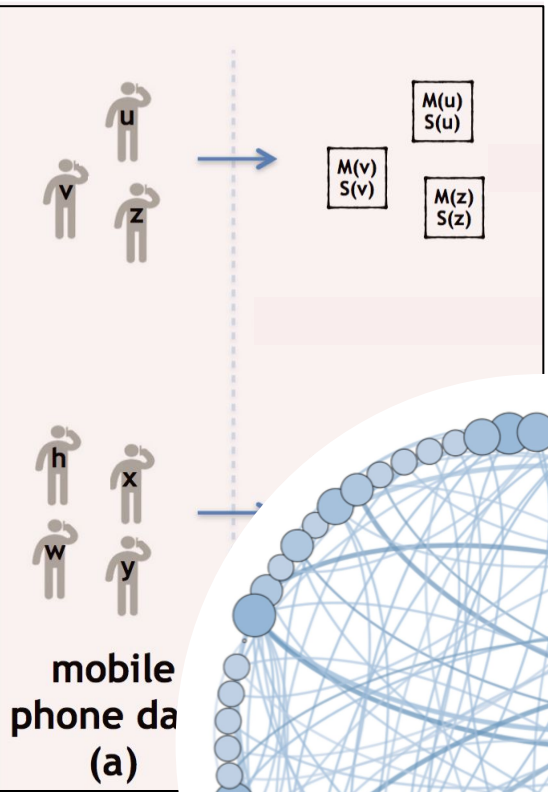
- Mobility diversity  
predictability of an individual's movements

$$S^{unc} = - \sum_{i=1}^n p_i \log_2 p_i$$

A red arrow points from the term  $p_i$  in the equation to the expression  $w_i/N$  written above it.

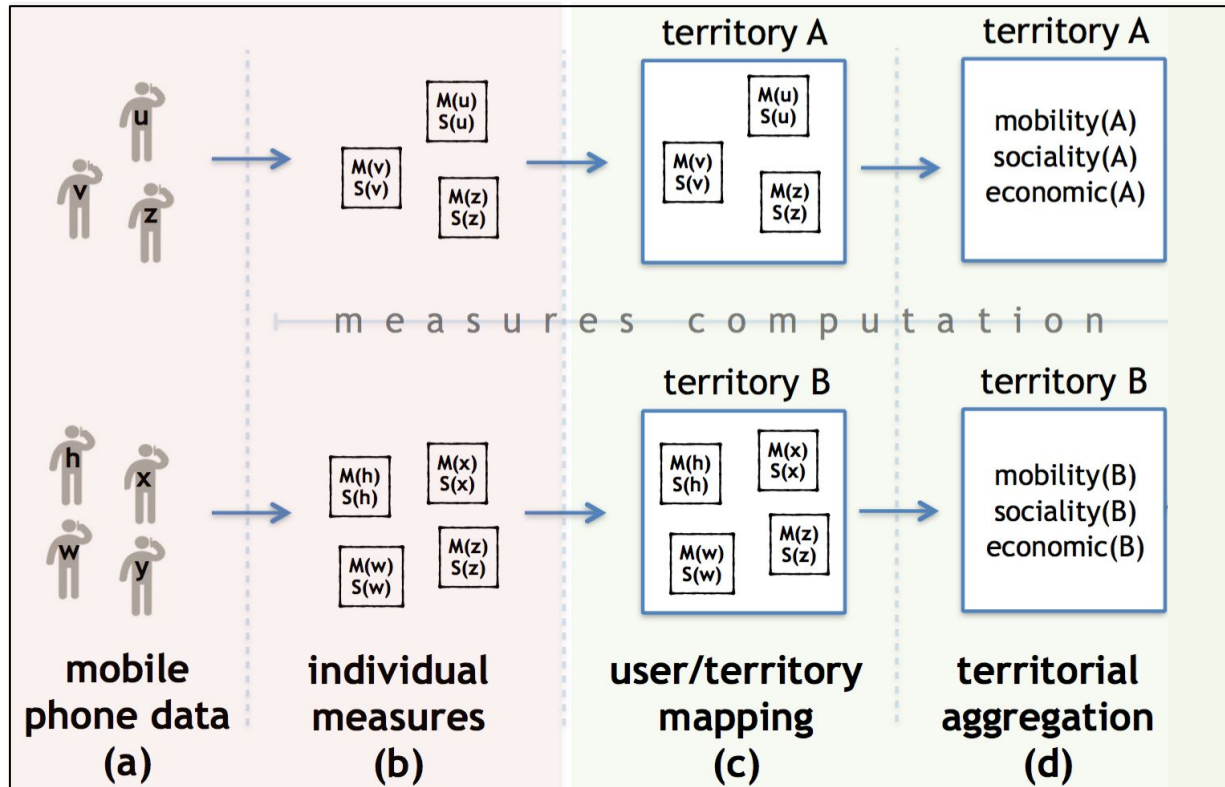
individual level

– Mobility diversity  
predictability of an individual's  
movements



individual level

territorial level



- Home location is the most frequent during nighttime (8 pm - 3 am)

- Every individual's home is assigned to its municipality

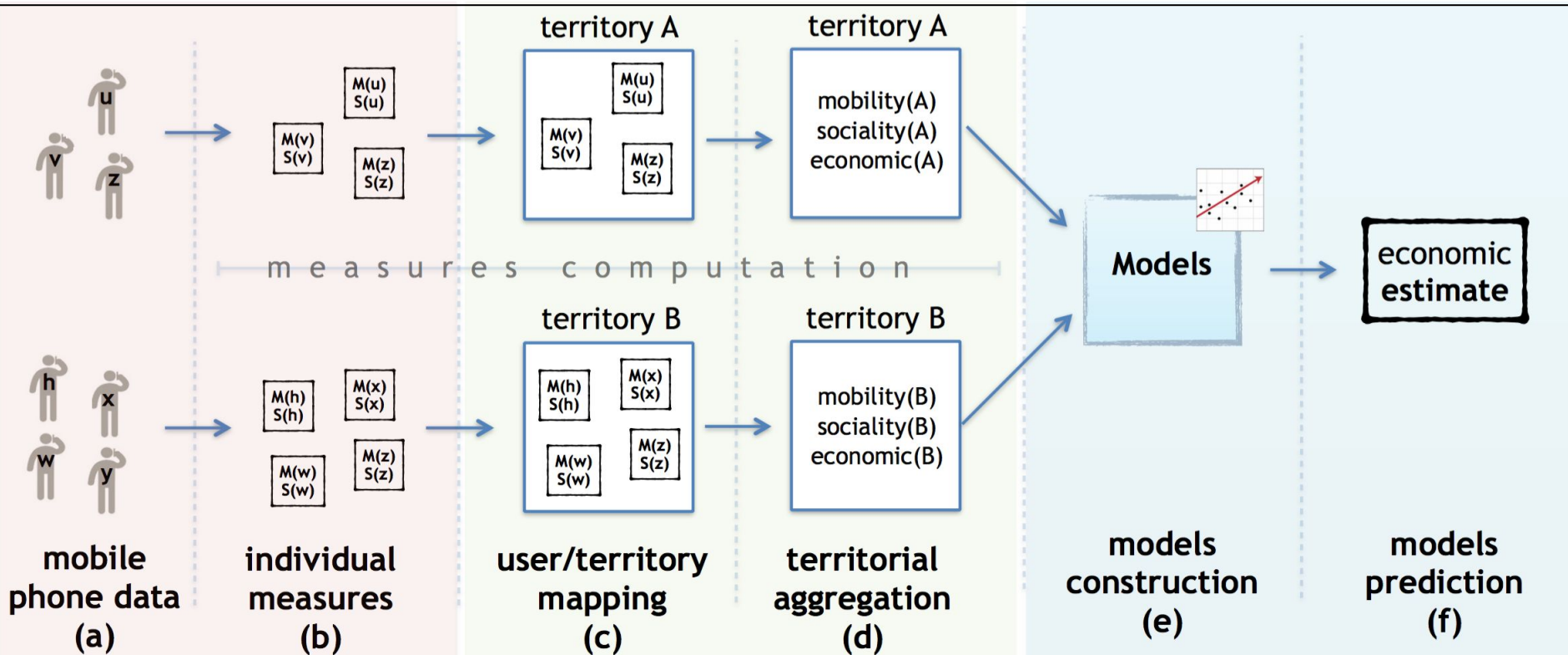
- Measures are aggregated at municipality level

- Economic measures are collected

- Deprivation index
- Per capita income

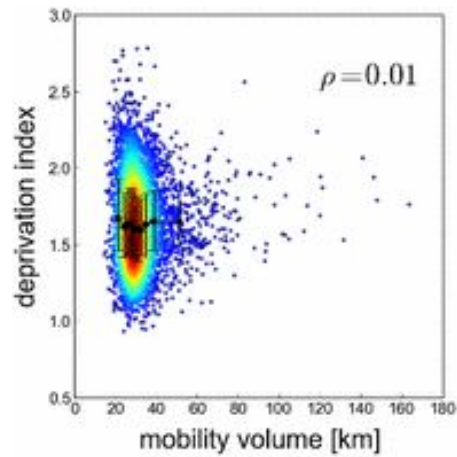
individual level

territorial level

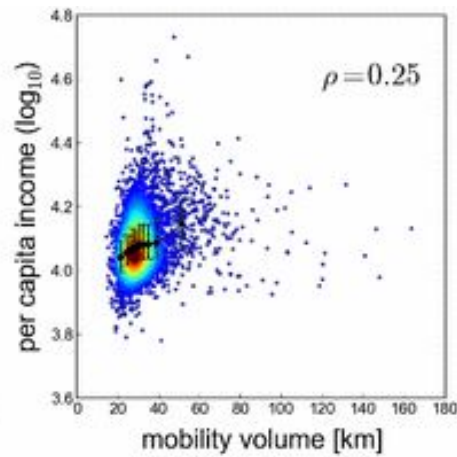




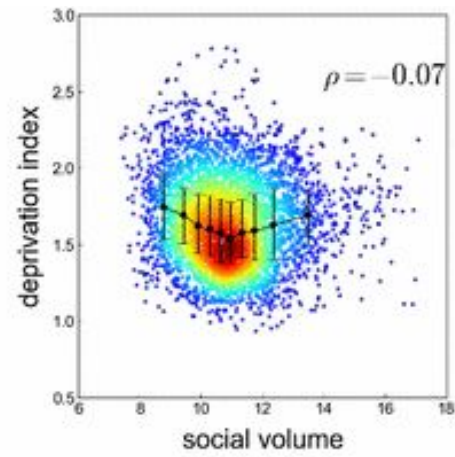
# Mobility diversity and well-being



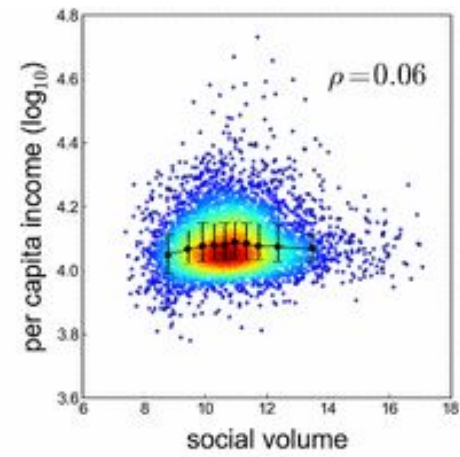
(a)



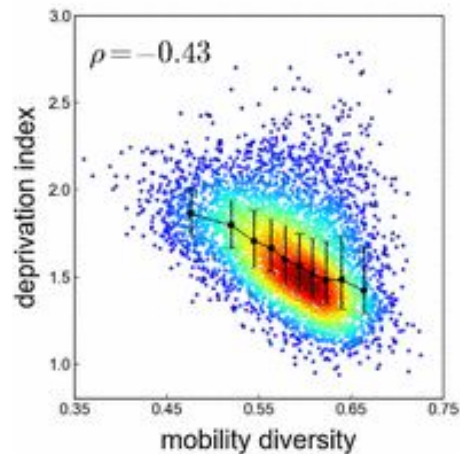
(b)



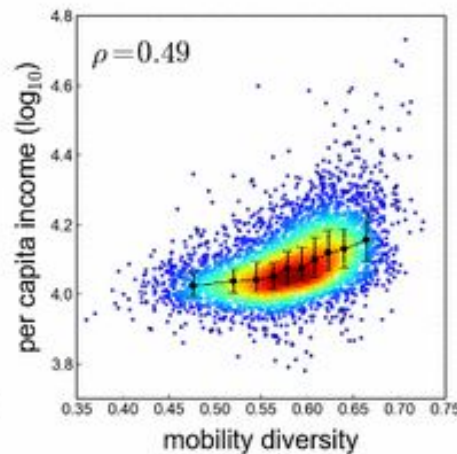
(c)



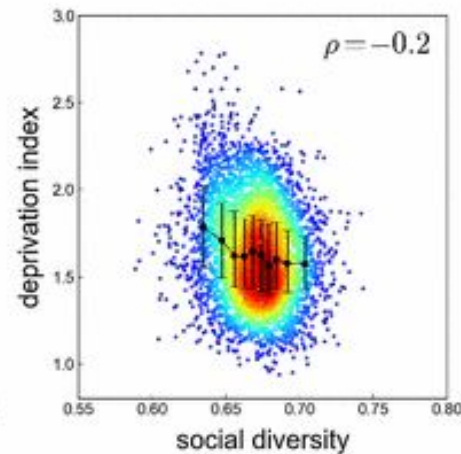
(d)



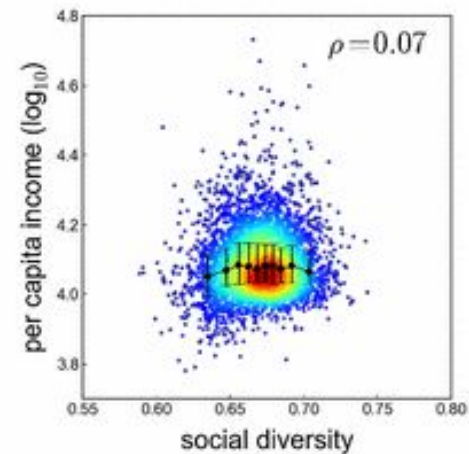
(e)



(f)

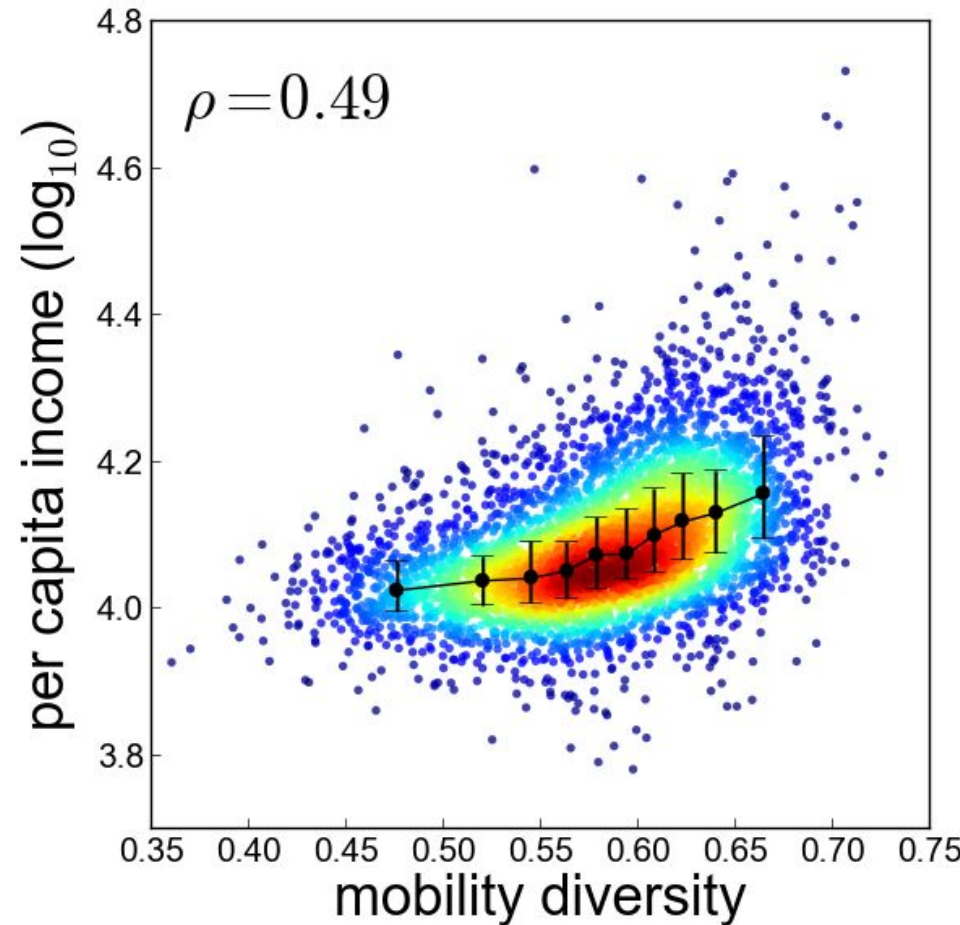
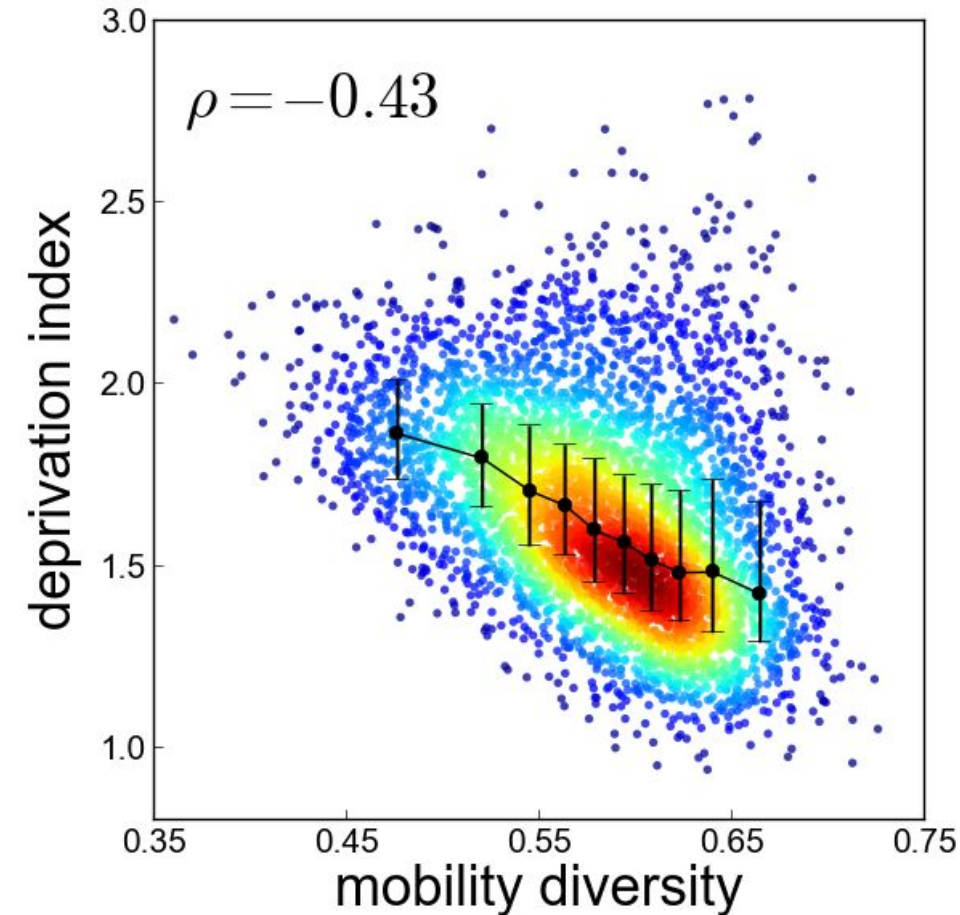


(g)

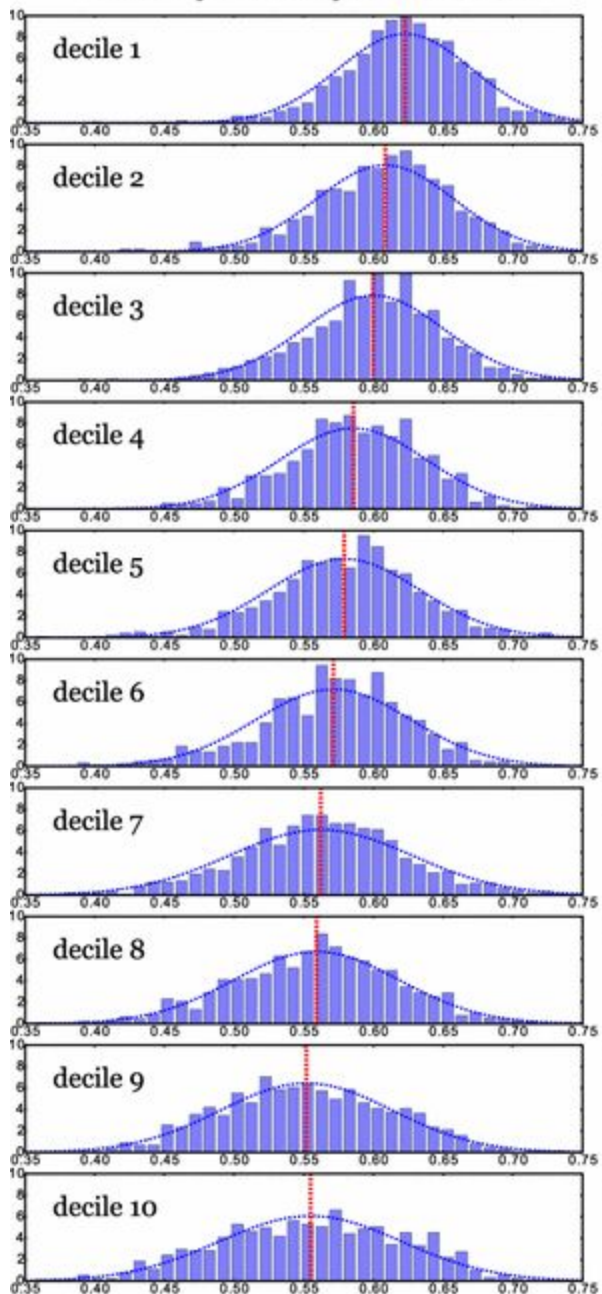


(h)

# Mobility diversity and well-being

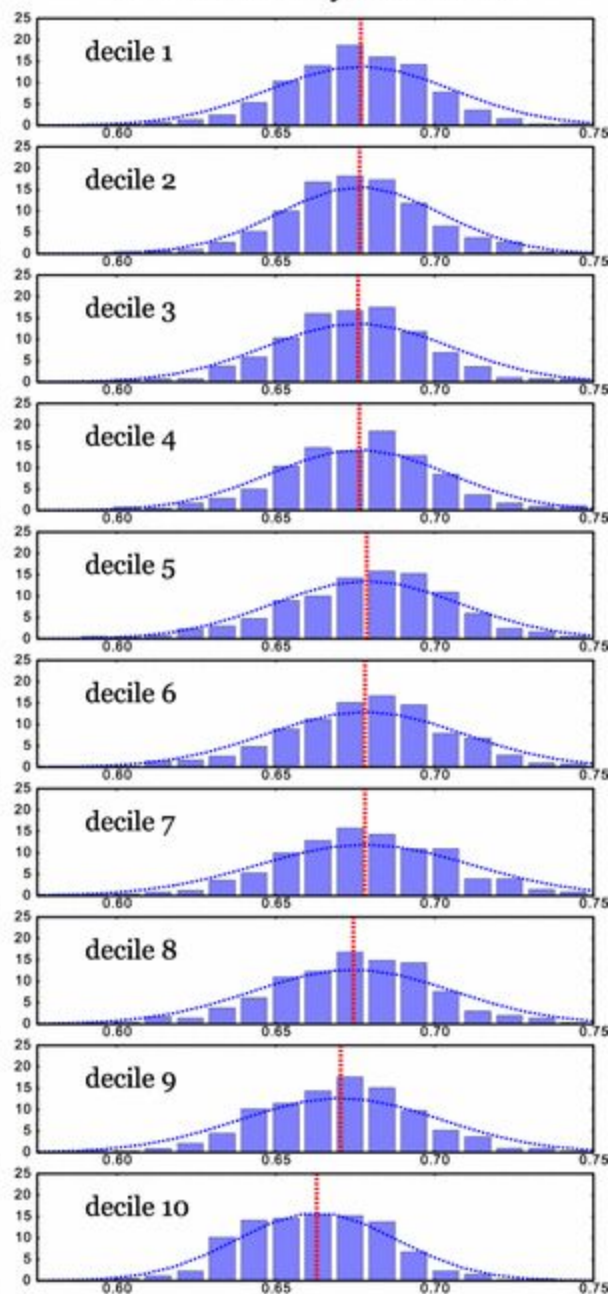


mobility diversity distribution



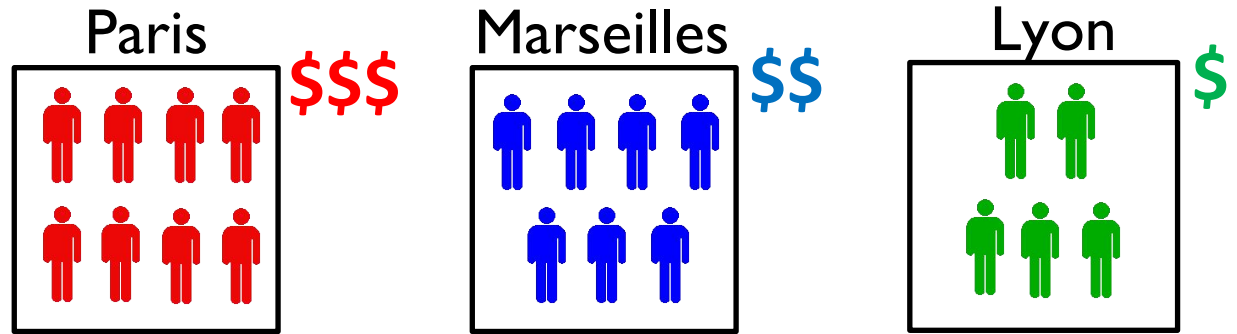
(a)

social diversity distribution



(b)

Data



NM1

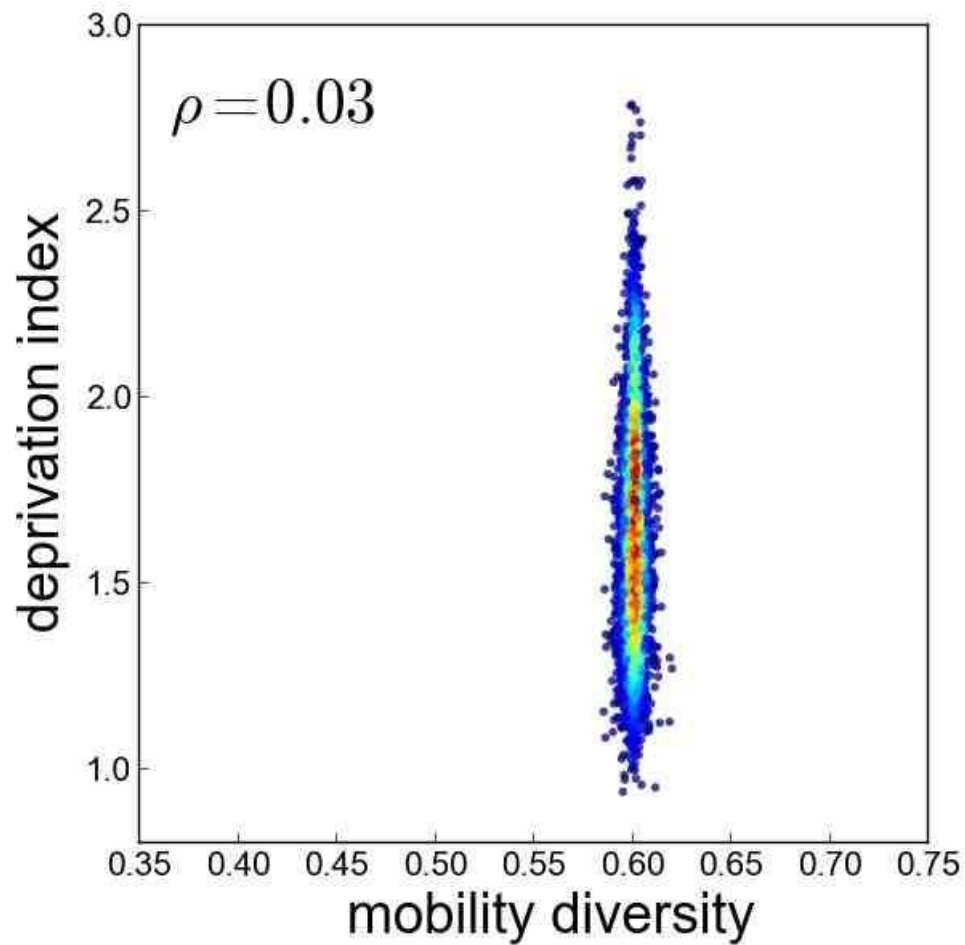


NM2

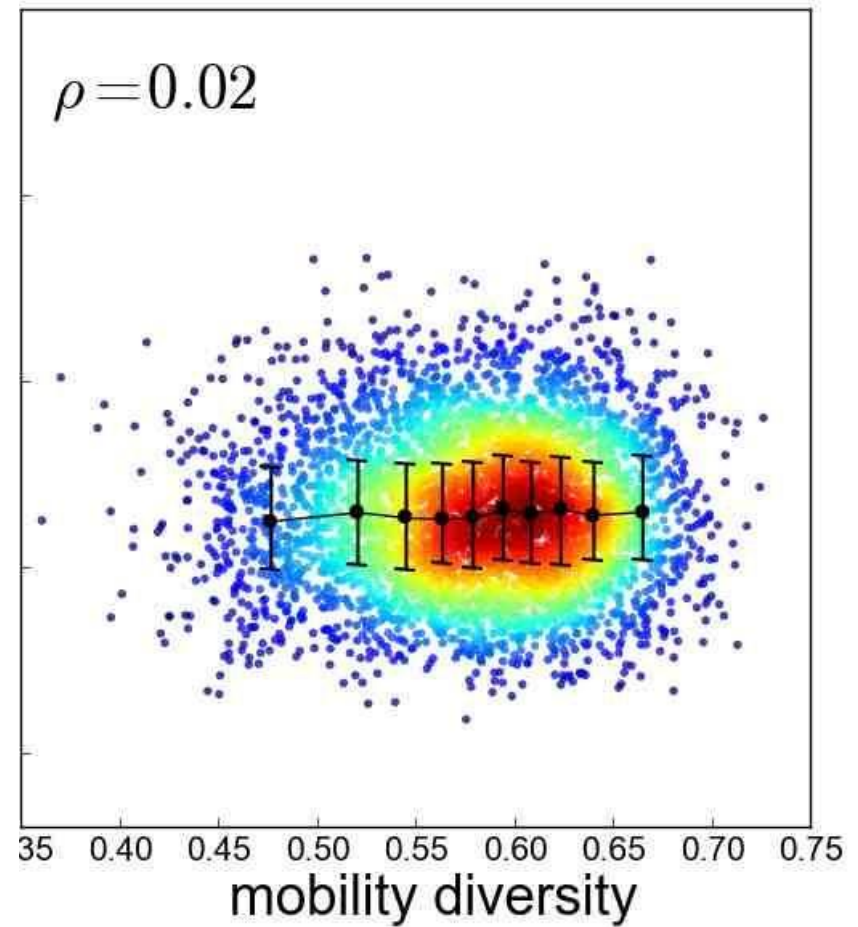




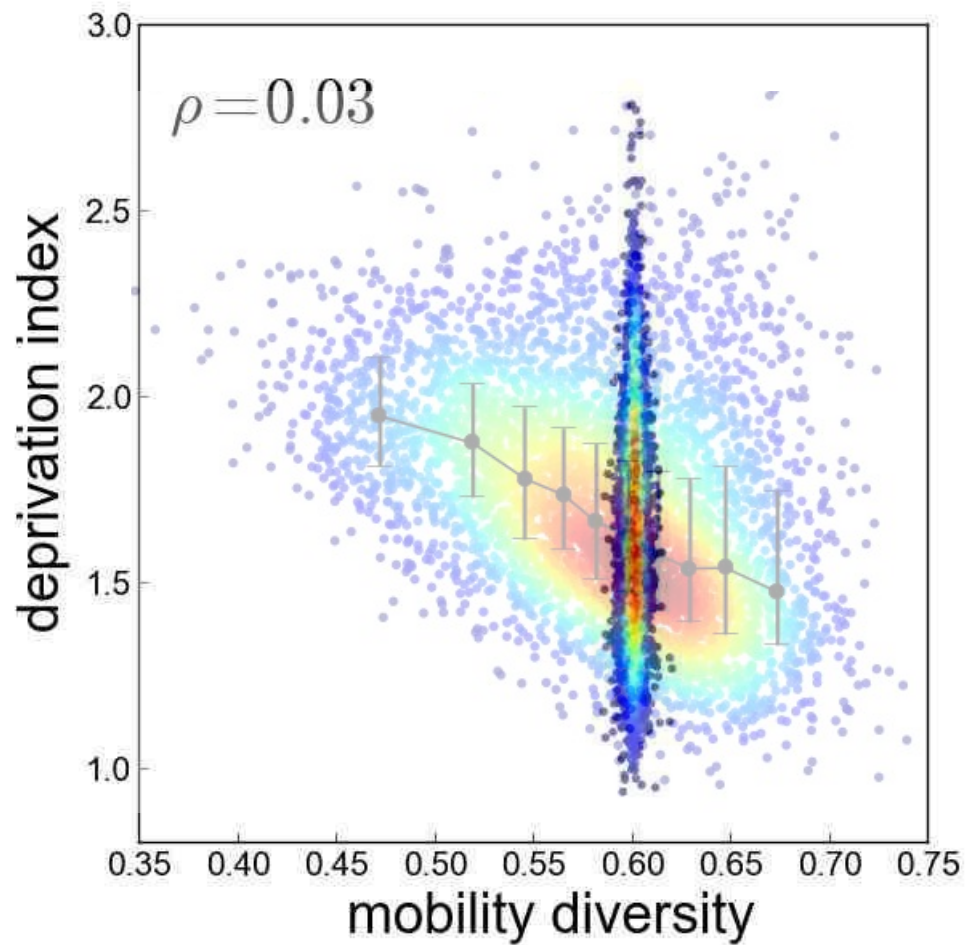
# NM1



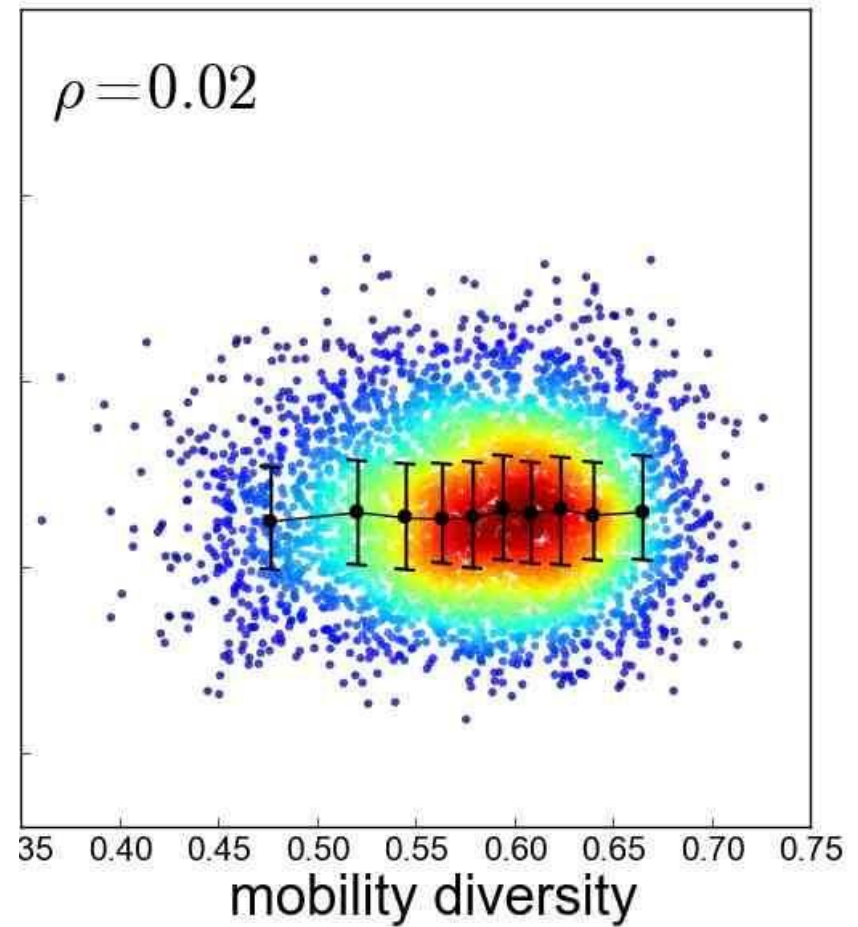
# NM2

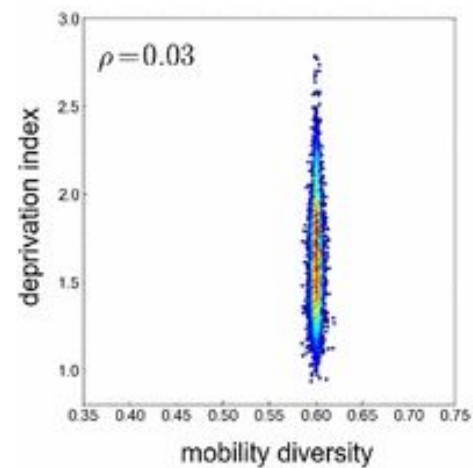


# NM1

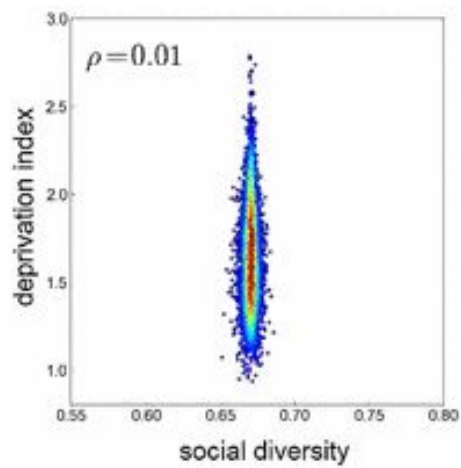


# NM2

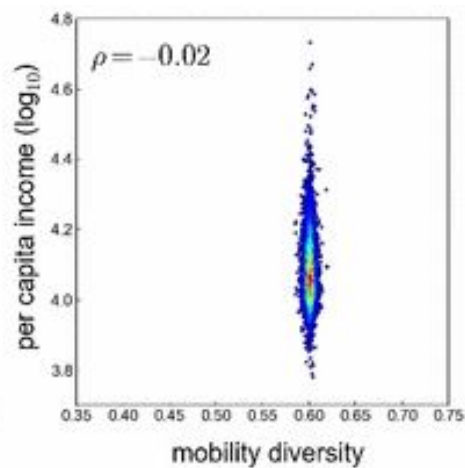




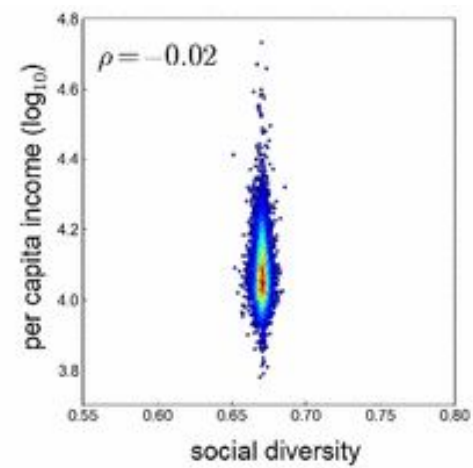
(a)



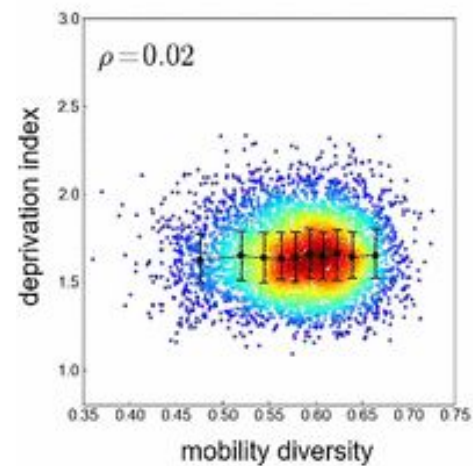
(b)



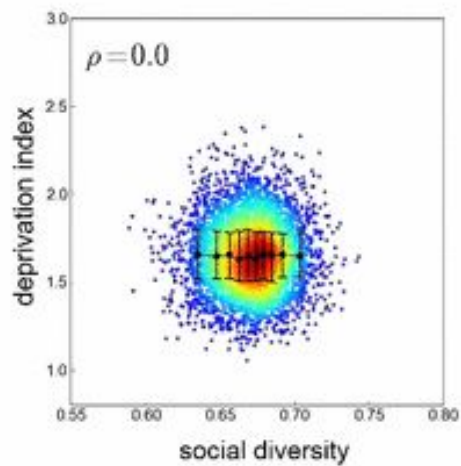
(c)



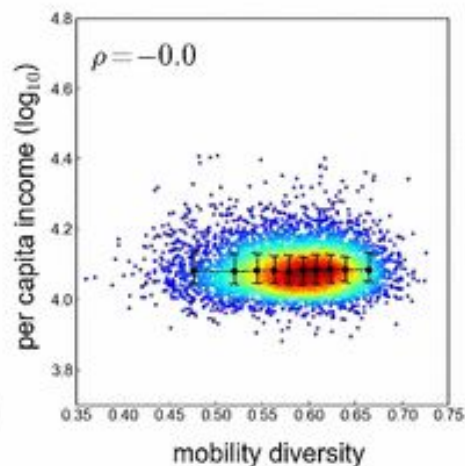
(d)



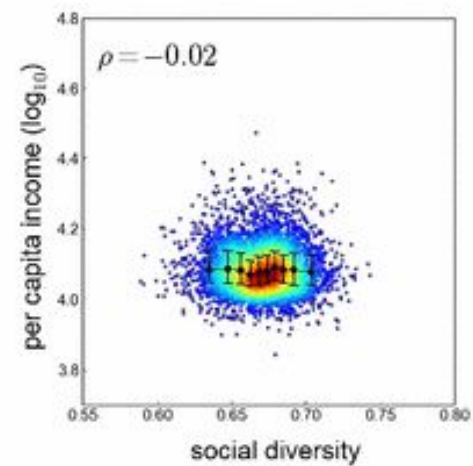
(e)



(f)



(g)



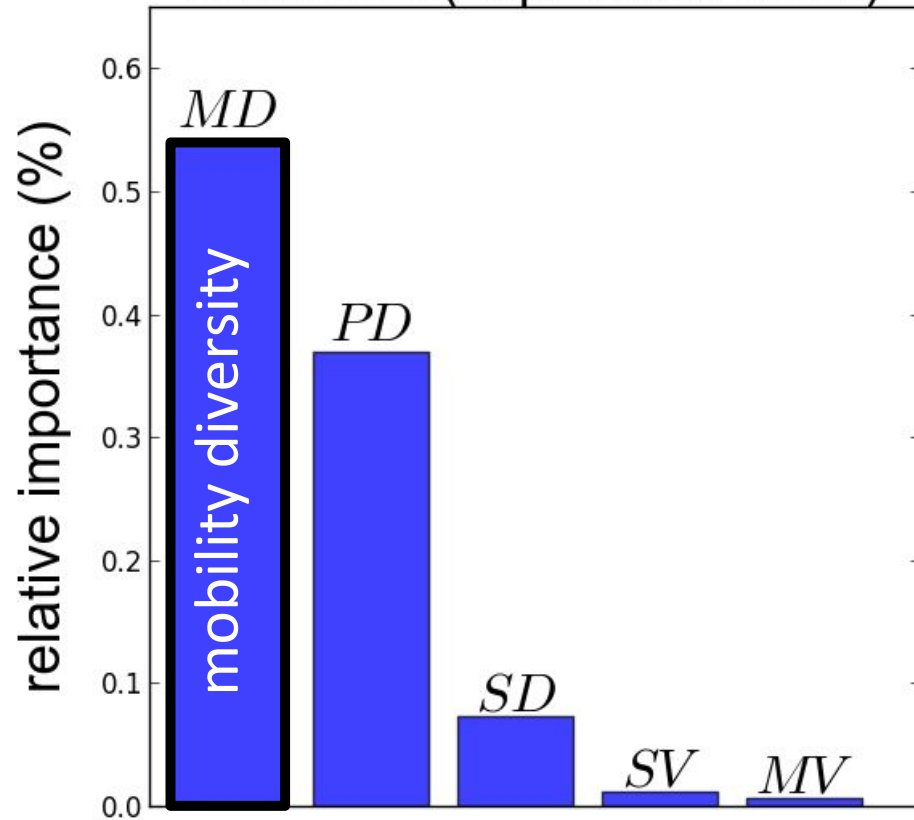
(h)

# Predicting well-being

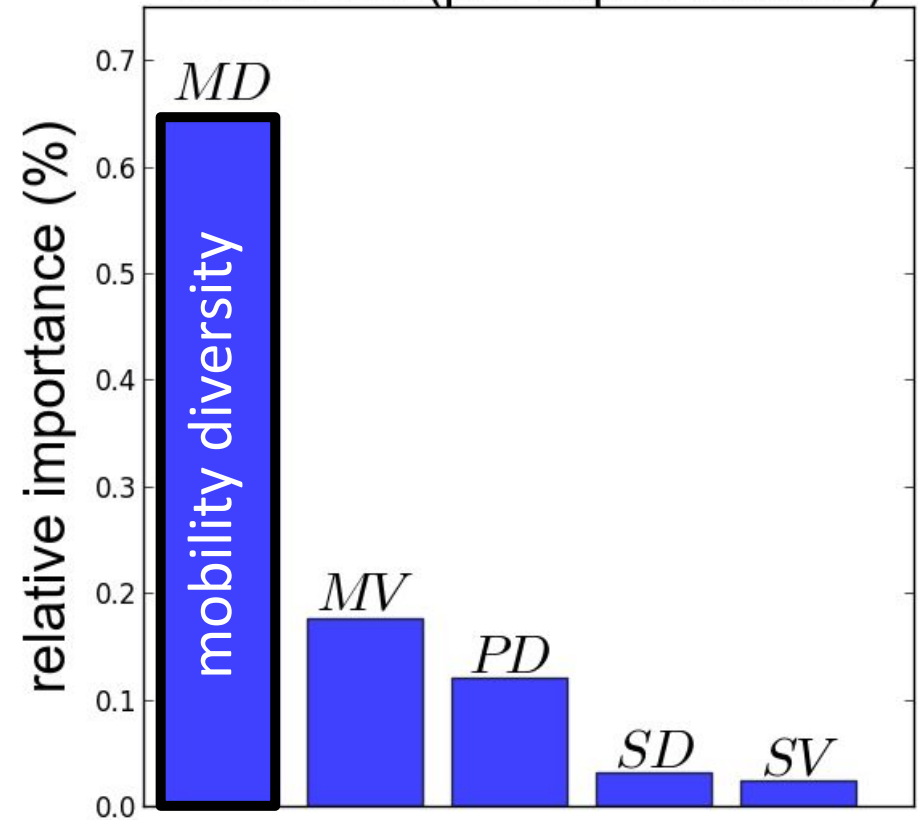
- Multivariate Regression  
predicting the exact value  
 $R^2 = 0.42$  (deprivation)  
 $R^2 = 0.25$  (income)
- Classification:  
predicting class of well-being (low, medium, high)  
acc = 0.61 (deprivation)  
acc = 0.54 (income)

# Diversity matters

Model M1 (deprivation index)

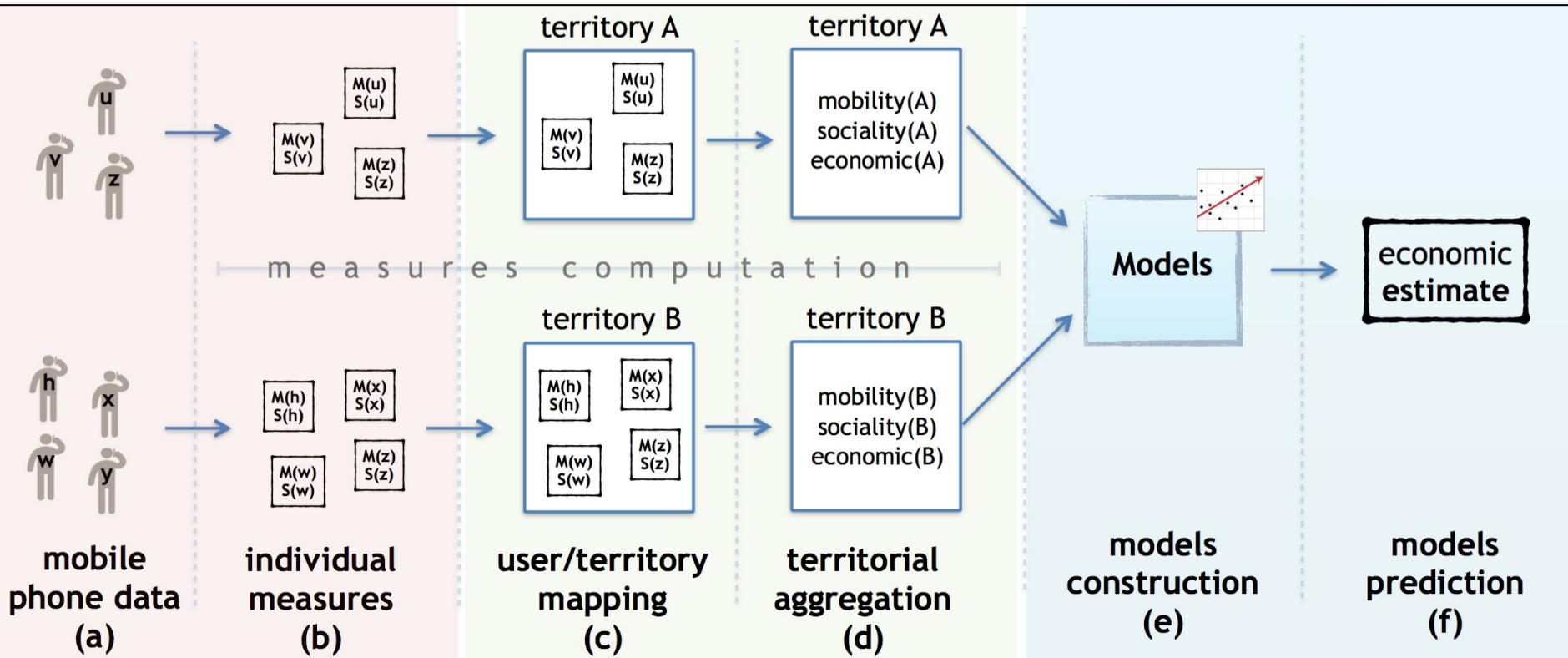


Model M2 (per capita income)



individual level

territorial level





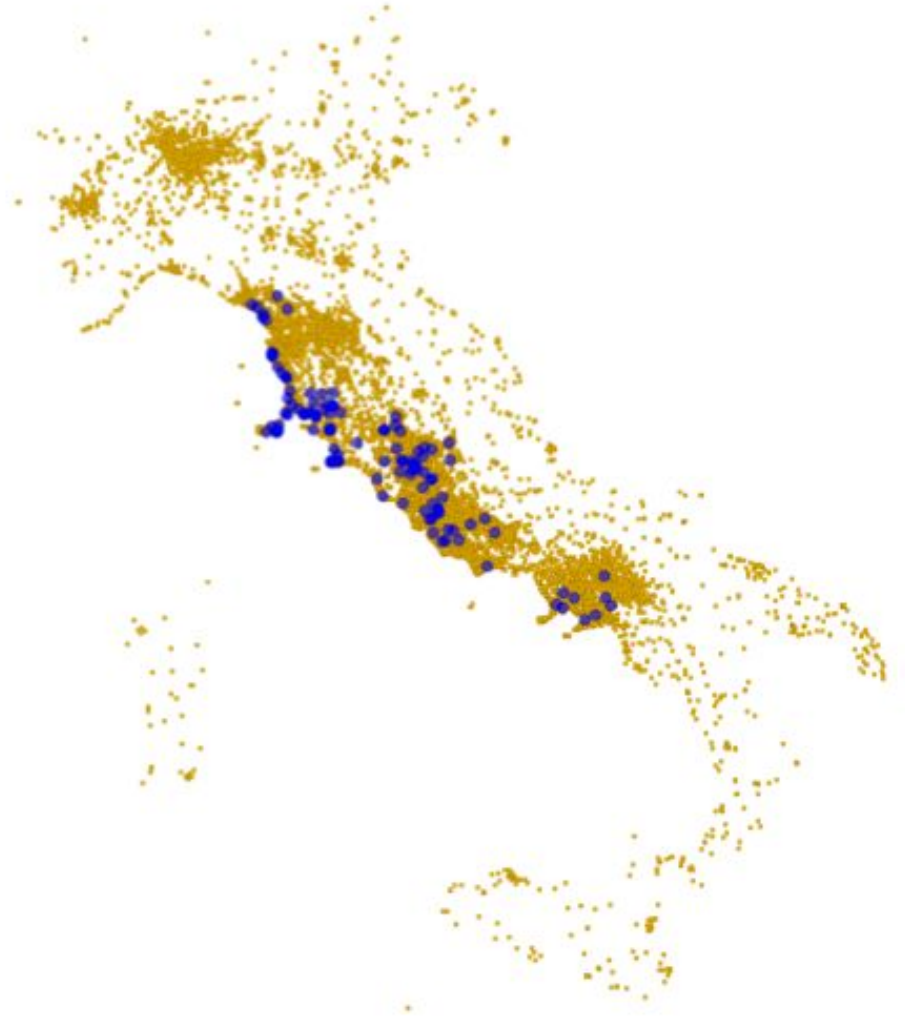
Wealth

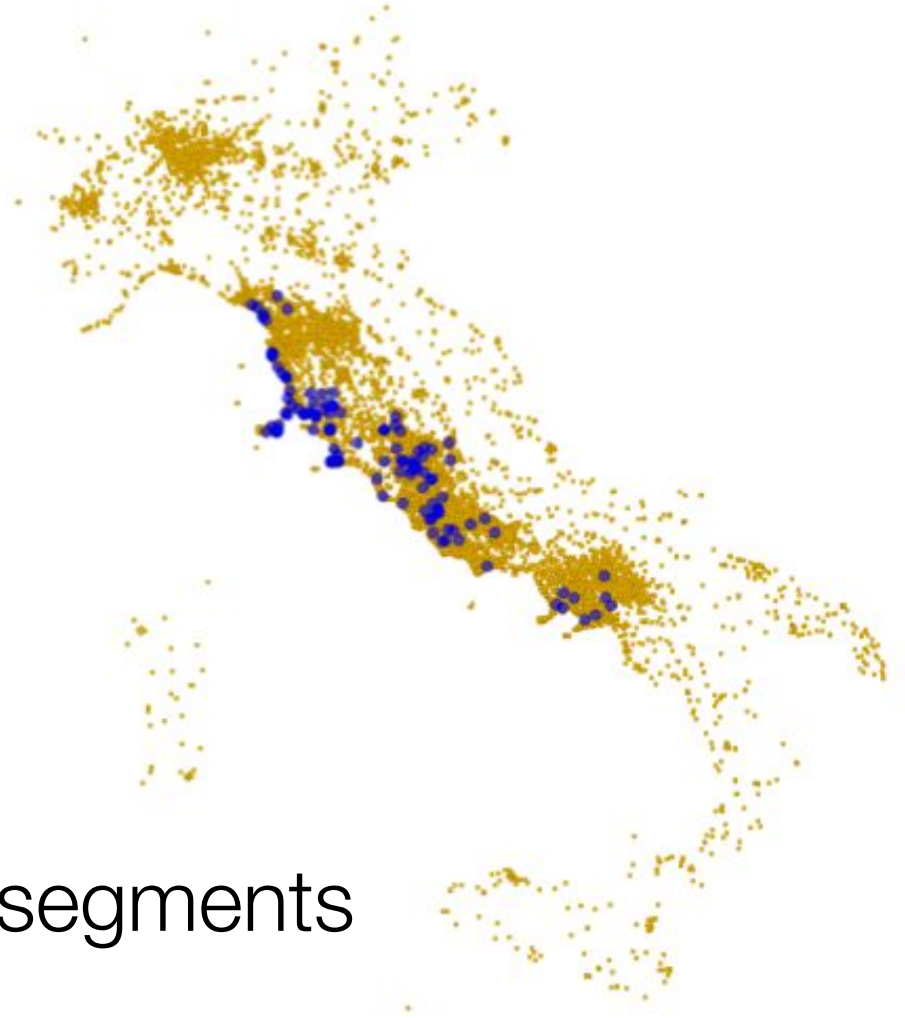
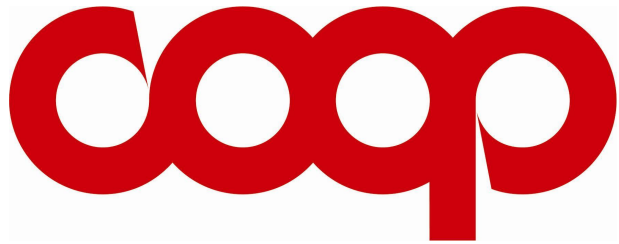


Retail market



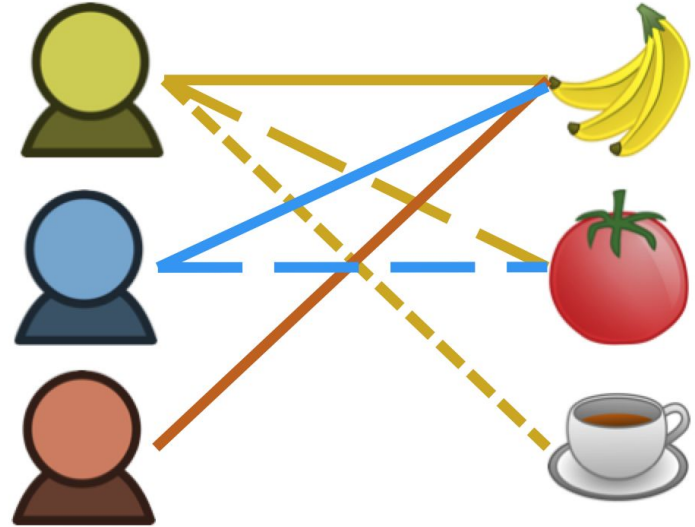
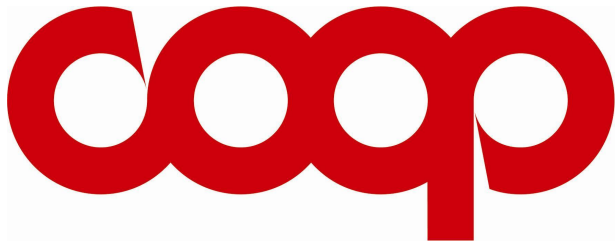
Monitoring  
GDP with  
retail data





- 120 stores
- 1 million customers
- 300K items → 4500 segments
- from 2007 to now



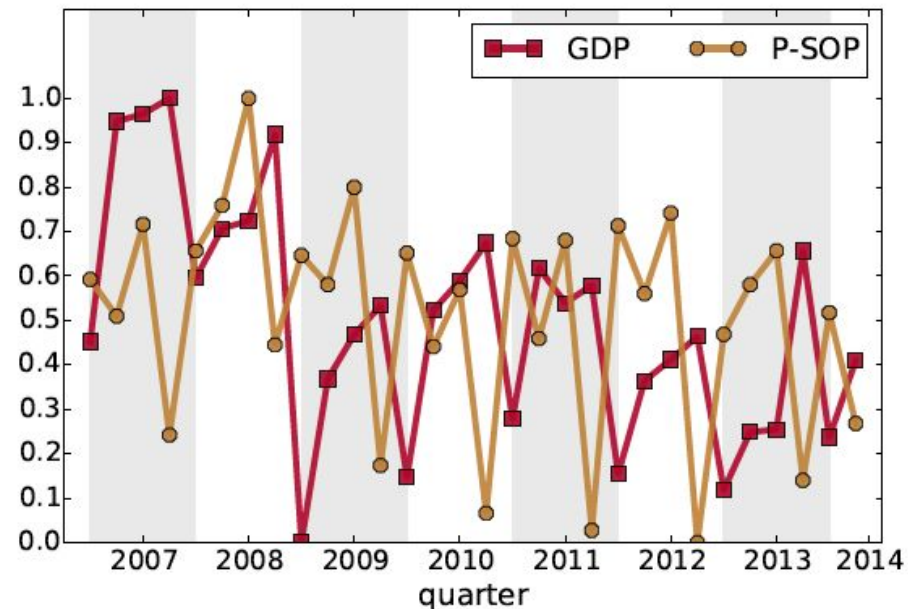
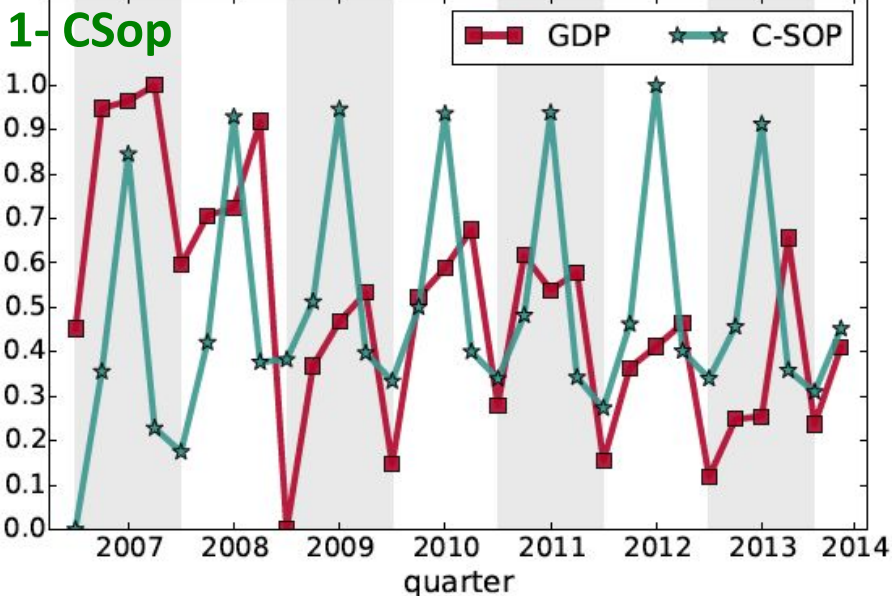


We compute a measure of sophistication  
of customers and products  
(inspired by Google PageRank)

# GDP and sophistication

- Rileviamo una relazione tra la sofisticatezza di clienti e prodotti e il PIL

-2 shift



# Health

## Eating habits



- 80K students
- 6 years

Meals	10,034,413
Students	82,871
<i>with grant</i>	19,141
<i>free meals</i>	4,730,658
Dishes	950
Food categories	41
Period	2,551 days
<i>from</i>	01/01/2010
<i>to</i>	12/26/2016

dish	attribute	description	long description
Bolognese pasta	a <sub>11</sub>	meat flours	flours (pasta, couscous, dumplings) with meat/cheese/eggs
Pasta with pesto	a <sub>13</sub>	veg flours	flours (pasta, couscous, dumplings) with vegetables
Pasta with zucchini	a <sub>13</sub>	veg flours	flours (pasta, couscous, dumplings) with vegetables
...	...	...	...
Saffron and potato soup	a <sub>34</sub>	legumes soup	potato and legumes soups
Hamburger with mushrooms	a <sub>51</sub>	red meat	red meat / salami
...	...	...	...
Green salad	a <sub>81</sub>	raw veg	raw vegetables
...	...	...	...
...	...	...	...
Fruit	a <sub>415</sub>	fruit	fruit
Cheesecake	a <sub>416</sub>	dessert	dessert

Student_id	timestamp	Dishes
A4578A	18/04/2015 12:42:00	pasta with tomato sauce, chicken breast, fruit
G23T20	18/04/2015 12:43:00	mushroom risotto, salad, fruit
GE54Y7	18/04/2015 12:44:01	pasta with tomato, fruit



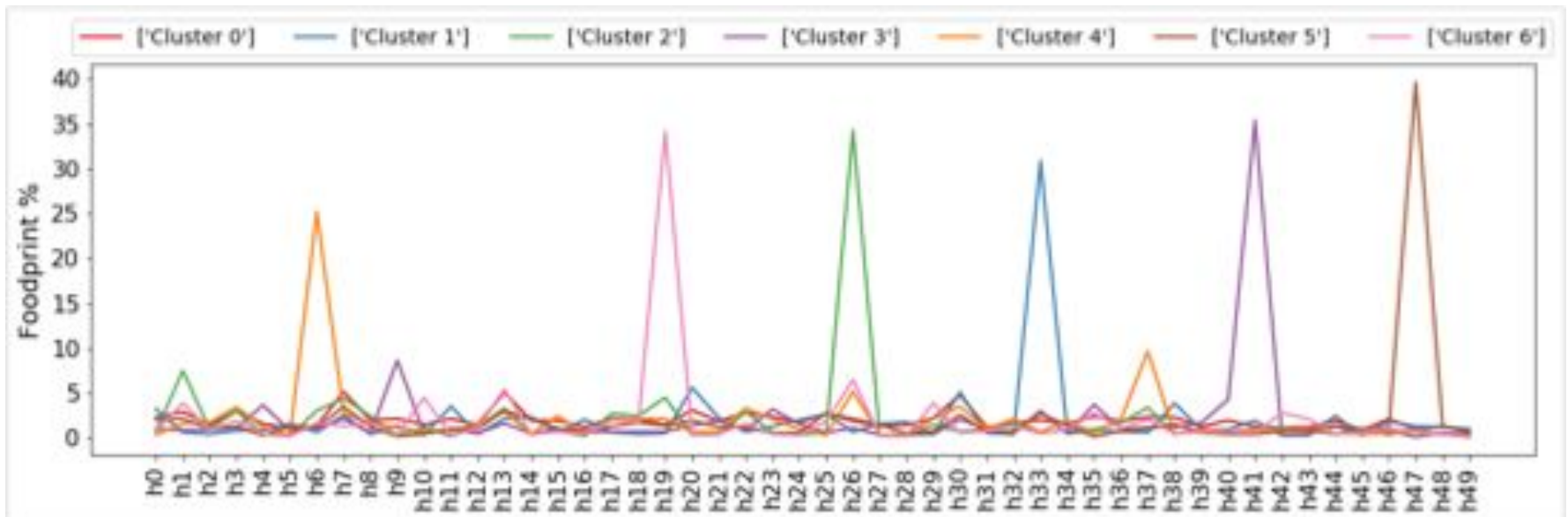
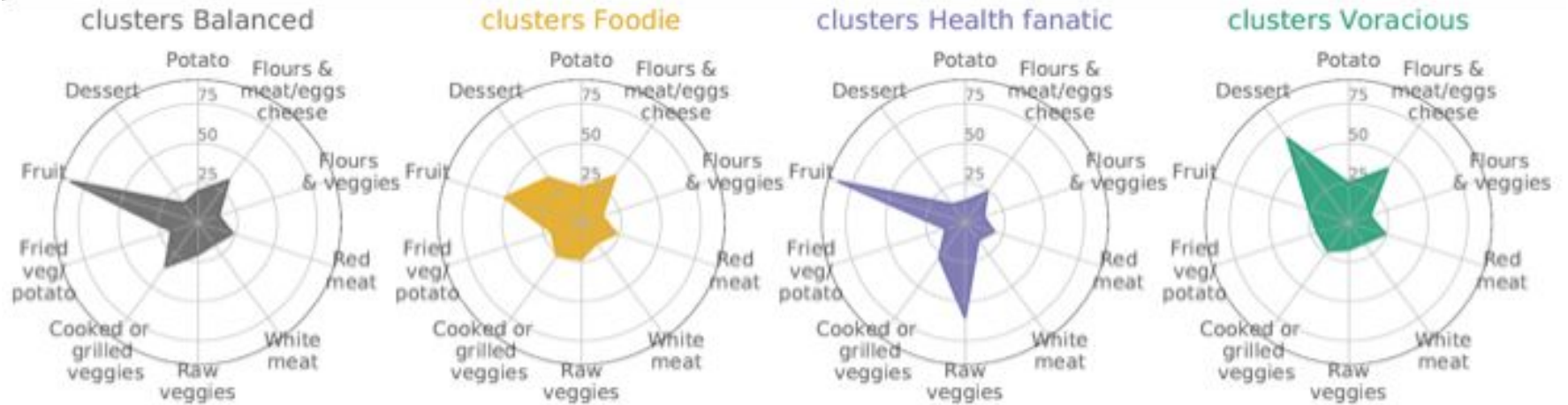
# Health

- Obesity is related to GERD
- Students with GERD: consumes no legumes
- Regular students with GERD: consumes -10% pasta and rice





# Health



# In summary ...

- Analytical frameworks can be defined to extract complex indices from Big Data
- Associations between these indices and well-being indicators can be investigated
- Predictive models can be constructed to nowcast well-being indicators
- This overcomes the limitations of surveys and censuses