Experience with Statistical Analysis of Covid data

INFN

Gaetano Salina INFN Sezione di Roma Tor Vergata



INFN School of Statistics 2022, Paestum 17 May 2022

Summary:

• Data Definition, Data Acquisition & Computational Models Is it a well-defined problem ? Maybe not !

• Epidemic dynamics on Graphs

Epidemic Models & Complex Systems: the role of topology

• Experience with Covid Data Analysis

Creating a Knowledge Base for Covid Data: a Nightmare

I apologize for my broken English and for the large number of slides

The new gospel:

With enough data, the numbers speak for themselves.

Chris Anderson

An English-American author and entrepreneur. He was the WIRED editor-in-chief until 2012.





Illustration: Marian Bantjes

"All models are wrong, but some are useful."

This article has been reproduced in a new format and may be missing content or contain faulty links. Contact wiredlabs@wired.com to report an issue.

So proclaimed statistician George Box 30 years ago, and he was right. But what choice did we have? Only models, from cosmological equations to theories of human behavior, seemed to be able to consistently, if imperfectly, explain the world around us. Until now. Today companies like Google, which have grown up in an era of massively abundant data, don't have to settle for wrong models. Indeed, they don't have to settle for models at all.

https://www.wired.com/2008/06/pb-theory/



The temptation: Correlation is enough !

- Peter Norvig, Google's research director, offered an update to George Box's maxim: "All models are wrong, and increasingly you can succeed without them."
- Forget taxonomy, ontology, and psychology. Who knows why people do what they do? The point is they do it, and we can track and measure it with unprecedented fidelity. With enough data, the numbers speak for themselves.
- We can stop looking for models. We can analyze the data without hypotheses about what it might show. We can throw the numbers into the biggest computing clusters the world has ever seen and let statistical algorithms find patterns where science cannot.

Chris Anderson

An English-American author and entrepreneur. He was the WIRED editor-in-chief until 2012.





Illustration: Marian Bantjes

"All models are wrong, but some are useful."

This article has been reproduced in a new format and may be missing content or contain faulty links. Contact wiredlabs@wired.com to report an issue.

So proclaimed statistician George Box 30 years ago, and he was right. But what choice did we have? Only models, from cosmological equations to theories of human behavior, seemed to be able to consistently, if imperfectly, explain the world around us. Until now. Today companies like Google, which have grown up in an era of massively abundant data, don't have to settle for wrong models. Indeed, they don't have to settle for models at all.

https://www.wired.com/2008/06/pb-theory/

Gaetano Salina



The Knowledge Discovery Process (KDP)

The Knowledge Discovery Process has the aim of providing decisions based on data. The ETL (Extract, Transform and Load, i.e. steps 1-3) is important to have an integrated Knowledge Base, with *unique, well-defined and quality data*, without redundancy.

- Selection: Capturing relevant prior knowledge, identifying the data-mining goal and developing and understanding of the application domain. Based on that, proper data samples as well as relevant variables can be selected.
- Pre-processing: The selected data are processed. Handling of missing values, the noise and errors identification (and correction), the elimination of duplicates, as well as the matching, fusion, and conflict resolution for data taken from different sources are done.
- Transformation: The clean dataset is transformed into a form suitable for data mining algorithms analysis. To improve the analysis performance dimensionality reduction methods can also be applied.



- Data mining: Goals of data mining are prediction and description. In prediction some variables and fields in the dataset are used to predict unknown values of other variables of interest, and description helps in finding human-understandable patterns describing the data. Or: it involves the use of statistical and machine learning algorithms to extract structures, correlations, patterns and rules within data.
- Evaluation and interpretation: The patterns and models derived are analysed with respect to their validity. The user assesses the usefulness of the found knowledge. A data visualization of the extracted patterns and models is involved.



Knowledge



Gaetano Salina

INFN School of Statistics 2022, Paestum 17 May 2022

Experience with Statistical Analysis of Covid Data



We can throw the numbers into the biggest computing clusters the world has ever seen and let statistical algorithms find patterns where science cannot.

Chris Anderson



When is a KDP successful?



E-Promotions:

Based on your current location, on your purchase history, on what you like \rightarrow the store located next to you sends promotions right now to you.



Healthcare monitoring:

Sensors monitoring your activities and body \rightarrow any abnormal measure requires immediate reaction.

Gaetano Salina

Data Definition, Data Acquisition & Computational Models Extracting Knowledge means dev

Extracting Knowledge means developing Computational Models

Computational models are mathematical models that are simulated using computation to study complex systems. ... The parameters of the mathematical model are adjusted using computer simulation to study different possible outcomes.

www.nature.com/subjects/computational-models

This may be a restrictive definition



Extracting Knowledge means developing Computational Models

Knowledge \rightarrow Computational Model means:

Knowledge \rightarrow Theory \rightarrow Algorithms.

We know the relevant observables for describing a system and their interaction laws. Computational models allow us to solve system with a large number of degree of freedom and complex dynamic or with non linear and/or non local interaction. Computational models are used when analytical tools fail or are unsuitable.

or

Knowledge \rightarrow Heuristic \rightarrow Algorithms.

No theory is available and/or we are not sure that the observables are relevant and/or the system is described by non-numerical observables. We use trial errors heuristic approach to define the observables interactions and define some numerical algorithms. Having a large number of degree of freedom and semantic different data we are playing in the field of Big Data. Computational models help us do define the correct observables and the system *dynamic*.



Extracting Knowledge means developing Computational Models

Data Definition, Data Acquisition & Computational Models

- Theoretical Models
- Measurement Theory
- Development of Detectors and Measuring Instruments





its orbit.

Data Definition, Data Acquisition & Computational Models Experience with Statistical Analysis of Covid Data

How Physics tamed and tames (Big) Data

Knowledge \rightarrow Heuristic \rightarrow Algorithms.

No theory is available and/or we are not sure that the observables are relevant and/or the system is described by non-numerical observables. We use trial errors heuristic approach to define the observables interactions and define some numerical algorithms. Having a large number of degree of freedom and semantic different data we are playing in the field of Big Data. Computational models help us do define the correct observables and the system dynamic.

Knowledge \rightarrow Theory \rightarrow Algorithms.

We know the relevant observables for describing a system and their interaction laws. Computational models allow us to solve systems with a large number of degree of freedom and complex dynamics or with non linear, non local interactions. Computational models are used when analytical tools fail or are unsuitable.



- A line segment joining a planet and

the Sun sweeps out equal areas during equal intervals of time.

- The square of a planet's orbital

period is proportional to the cube of

the length of the semi-major axis of

Brahe was a towering figure. He ran a huge research program with a castle like observatory, a NASA-like budget, and the

finest instruments and best assistants money could buy.



Keplero

Brahe

Theoretical

Measurement

Development

and Measuring

of Detectors

Instruments

Models

Theory

A Data Ontology & Data Semantic is necessary

Data Definition, Data Acquisition & Computational Models

Semantics are only one part of the solution and often not the end-product so the focus of the design should be on creating effective methods, tools and APIs to handle and process the semantics. The concept of semantic data have changed the approach to data managing, due to the heterogeneous and unstructured nature of data sets. The ontological & semantic approach has been thought to manage data sets of this kind, i.e. a way to Knowledge Discovery Process.



It means developing Computational Models

Experience with Statistical Analysis of Covid Data



Gaetano Salina

INFN School of Statistics 2022, Paestum 17 May 2022

Experience with Statistical Analysis of Covid Data



Gaetano Salina

Lecce 21 December 202

INFN School of Statistics 2022, Paestum 17 May 2022

Experience with Statistical Analysis of Covid Data

The Knowledge Discovery Process (KDP)

Big Data between myth and reality

FORECASTING IN LIGHT OF BIG DATA

An extreme inductivist approach to forecasting using Big Data

According to a vaguely defined yet rather commonly held view big data may lead to dispense with theory, modeling or even hypothesizing.

All of this would be encompassed, across domains, by smart enough machine learning algorithms operating on large enough data sets.

We are interested in forecasts such that future states of a systems are predicted solely on the basis of known past states.

Two hypotheses are needed to give an affirmative answer:

- Similar premises lead to similar conclusions (Analogy);
- Systems which exhibit a certain behavior, will continue doing so (Determinism).

In more formal terms, given the series $\{x_1, ..., x_M\}$, where x_i is the vector describing the state at time $j\Delta t$, we look in the past for an analogous state, that is a vector x_k with k < M "near enough" (i.e. $|x_k - x_M| < \varepsilon$, being ε the desired degree of accuracy).

Once we find such a vector, we "predict" the future at times M + n > M by simply assuming for $x_{M\pm n}$ the state $x_{k\pm n}$. It all seems quite easy, but it is not at all obvious that an analog can be found.



FORECASTING IN LIGHT OF BIG DATA arXiv:1705.11186v1 [physics.soc-ph] 31 May 2017

HYKEL HOSNI AND ANGELO VULPIANI

ABSTRACT. Predicting the future state of a system has always been a natural motivation for science and practical applications. Such a topic, beyond its obvious technical and societal relevance, is also interesting from a conceptual point of view. This owes to the fact that forecasting lends itself to two equally radical, yet opposite methodologies. A reductionist one, based on the first principles, and the naïve-inductivist one, based only on data. This latter view has recently gained some attention in response to the availability of unprecedented amounts of data and increasingly sophisticated algorithmic analytic techniques. The purpose of this note is to assess critically the role of *big data* in reshaping the key aspects of forecasting and in particular the claim that *bigger* data leads to *better* predictions. Drawing on the representative example of weather forecasts we argue that this is not generally the case. We conclude by suggesting that a clever and context-dependent compromise between modelling and quantitative analysis stands out as the best forecasting strategy, as anticipated nearly a century ago by Richardson and von Neumann.

Nothing is more practical than a good theory (L. Boltzmann)

In particular we ask whether using our knowledge of the past states of a system – and without the use of models for the evolution equation – meaningful predictions about the future are possible.

The key to understanding the "proper level" of abstraction lies with identifying the "relevant variables" and the effective equations which rule their time evolution.

INFN School of Statistics 2022, Paestum 17 May 2022

Gaetano Salina

Gaetano Salina

The Knowledge Discovery Process (KDP) Big Data between myth and reality INFN FORECASTING IN LIGHT OF BIG DATA An extreme inductivist approach to forecasting using Big Data **Time Evolution** Poincare' recurrence theorem After a suitable time, a deterministic system with a bounded phase space returns to a $|\mathbf{x}_k - \mathbf{x}_M| < \varepsilon$ state near to its initial condition. Thus an analog surely exists. How long do we have to go back to find it? Kac who proved a Lemma to the effect that $D \approx 10$ the average return time in a region A is $\varepsilon \approx 0.01$ proportional to the inverse of the probability $\langle T_{R} \rangle = O(10^{20})$ P(A) that the system is in A.



Consider a system of dimension D. The probability P(A) of being in A is

$$P(A) \propto \varepsilon^{D} \mapsto \langle T_{R} \rangle = O(\varepsilon^{-D})$$

Gaetano Salina

The return time is so large that in practice a recurrence is never observed. So the required analog, whose existence is guaranteed in theory, sometimes cannot be expected to be found in practice, even if

complete and precise information about the system.

Lecce 21 December 2021

FORECASTING IN LIGHT OF BIG DATA arXiv:1705.11186v1 [physics.soc-ph] 31 May 2017 HYKEL HOSNI AND ANGELO VULPIANI ABSTRACT. Predicting the future state of a system has always been a natural motivation for science and practical applications. Such a topic, beyond its obvious technical and societal relevance, is also interesting from a conceptual point of view. This owes to the fact that forecasting lends itself to two equally radical, yet opposite methodologies. A reductionist one, based on the first principles, and the naïve-inductivist one, based only on data. This latter view has recently gained some attention in response to the availability of unprecedented amounts of data and increasingly sophisticated algorithmic analytic techniques. The purpose of this note is to assess critically the role of *big data* in reshaping the key aspects of forecasting and in particular the claim that *bigger* data leads to *better* predictions. Drawing on the representative example of weather forecasts we argue that this is not generally the case. We conclude by suggesting that a clever and context-dependent compromise between modelling and quantitative analysis stands out as the best forecasting strategy, as anticipated nearly a century ago by Richardson and von Neumann

Nothing is more practical than a good theory (L. Boltzmann)

In particular we ask whether using our knowledge of the past states of a system – and without the use of models for the evolution equation – meaningful predictions about the future are possible.

The key to understanding the "proper level" of abstraction lies with identifying the "relevant variables" and the effective equations which rule their time evolution.

Gaetano Salina



The Knowledge Discovery Process (KDP)

In socio-economical systems the gap between data and our scientific ability to actually understanding them is typically enormous. Surely the availability of huge amounts of data, sophisticated methods for its retrieval and unprecedented computational power available for its analysis will undoubtedly help moving science and technology forward.

But in spite of a persistent emphasis on a fourth paradigm (beyond experiment, theory and computation) based only on data, there is as yet no evidence data alone can bring about scientifically meaningful advance.

We therefore conclude that the big data revolution is by all means a welcome one for the new opportunities it opens.

However the role of modeling cannot be discounted.

FORECASTING IN LIGHT OF BIG DATA arXiv:1705.11186v1 [physics.soc-ph] 31 May 2017

HYKEL HOSNI AND ANGELO VULPIANI

ABSTRACT. Predicting the future state of a system has always been a natural motivation for science and practical applications. Such a topic, beyond its obvious technical and societal relevance, is also interesting from a conceptual point of view. This owes to the fact that forecasting lends itself to two equally radical, yet opposite methodologies. A reductionist one, based on the first principles, and the naïve-inductivist one, based only on data. This latter view has recently gained some attention in response to the availability of unprecedented amounts of data and increasingly sophisticated algorithmic analytic techniques. The purpose of this note is to assess critically the role of *big data* in reshaping the key aspects of forecasting and in particular the claim that *bigger* data leads to *better* predictions. Drawing on the representative example of weather forecasts we argue that this is not generally the case. We conclude by suggesting that a clever and context-dependent compromise between modelling and quantitative analysis stands out as the best forecasting strategy, as anticipated nearly a century ago by Richardson and von Neumann.

Nothing is more practical than a good theory (L. Boltzmann)

A clever and context-dependent trade-off between modeling and quantitative analysis stands out as the best strategy for meaningful prediction.



Experience with Statistical Analysis of Covid Data

Big data undoubtedly constitute a great opportunity for scientific and technological advance, with a potential for considerable socio-economic impact. To make the most of it, however, the ensuing developments at the interface of statistics, machine learning and artificial intelligence, must be coupled with adequate methodological foundations.

H. Hosni & A. Vulpiani



Gaetano Salina

INFN School of Statistics 2022, Paestum 17 May 2022



The Knowledge Discovery Process (KDP)

A ontological & semantic approach to heterogeneous and unstructured datasets as way to Data Definition and Acquisition.

The develop of Computational Models, i.e.

- Heuristic Models: the goal of the is matched to a particular method, such as classification, regression, or clustering the transformed data. A decision about which models and parameters might be appropriate must be done and matching a particular data mining method with the overall criteria of the KDP in Datasets process is necessary.
- Theoretical Models: the transformed data are used to derive the "interaction" parameters of a theoretical dynamic model. The KDP coincides with the understanding of the system's physical dynamics.

as Data Mining

It can be, roughly speaking, the way of *physics* to try to tame complex systems

Gaetano Salina

. . .



Experience with Statistical Analysis of Covid Data

A Contribution to the Mathematical Theory of Epidemics. By W. O. KERMACK and A. G. MCKENDRICK. (Communicated by Sir Gilbert Walker, F.R.S.—Received May 13, 1927.)

(From the Laboratory of the Royal College of Physicians, Edinburgh.)



Mathematical Theory of Epidemics.

715

We are, in fact, assuming that plague in man is a reflection of plague in rats, and that with respect to the rat (1) the uninfected population was uniformly susceptible; (2) that all susceptible rats in the island had an equal chance of being infected; (3) that the infectivity, recovery, and death rates were of constant value throughout the course of sickness of each rat; (4) that all cases ended fatally or became immune; and (5) that the flea population was so large that the condition approximated to one of contact infection.



The accompanying chart is based upon figures of deaths from plague in the island of Bombay over the period December 17, 1905, to July 21, 1906. The ordinate represents the number of deaths per week, and the abscissa denotes the time in weeks. As at least 80 to 90 per cent. of the cases reported terminate fatally, the ordinate may be taken as approximately representing dz/dt as a function of t. The calculated curve is drawn from



- An epidemic can be rephrased as a stochastic reaction-diffusion process.
- Individuals belonging to the different compartments can be represented as different kinds of "particles" or "species", that evolve according to a given set of mutual interaction rules.
- o In the continuous time limit each reaction (transition) is defined by an appropriate reaction rate.

https://arxiv.org/pdf/1408.2701.pdf

Experience with Statistical Analysis of Covid Data

Classic Approach: Compartmental Models

We infer the model parameters based on the official data (...) about the evolution of the epidemic in Italy from 20/2/2020 through 5/4/2020.





<u>Nat Med.</u> 2020 Apr 22 : 1–6. doi: <u>10.1038/s41591-020-0883-7</u> [Epub ahead of print]	PMCID: PMC7175834 PMID: <u>32322102</u>
Modelling the COVID-19 epidemic and implementation of populin Italy	lation-wide interventions
Giulia Giordano, ⁸¹ Franco Blanchini, ² Raffaele Bruno, ^{3,4} Patrizio Colaneri, ^{5,6} Alessand and <u>Marta Colaneri</u> ³	dro Di Filippo, ³ Angela Di Matteo, ³
Author information Article notes Copyright and License information Disclaimer	
This article has been <u>cited by</u> other articles in PMC.	
Associated Data	
Supplementary Materials	
Data Availability Statement	
	Cotor

In Italy, 128,948 confirmed cases and 15,887 deaths of people who tested positive for SARS-CoV-2 were registered as of 5 April 2020. Ending the global SARS-CoV-2 pandemic requires implementation of multiple population-wide strategies, including social distancing, testing and contact tracing. We propose a new model that predicts the course of the epidemic to help plan an effective control strategy. The model considers eight stages of infection: susceptible (S), infected (I), diagnosed (D), ailing (A), recognized (R), threatened (T), healed (H) and extinct (E), collectively termed SIDARTHE. Our SIDARTHE model discriminates between infected individuals depending on whether they have been diagnosed and on the severity of their symptoms. The distinction between diagnosed and non-diagnosed individuals is important because the former are typically isolated and hence less likely to spread the infection. This delineation also helps to explain misperceptions of the case fatality rate and of the epidemic spread. We compare simulation results with real data on the COVID-19 epidemic in Italy, and we model possible scenarios of implementation of countermeasures. Our results demonstrate that restrictive social-distancing measures will need to be combined with widespread testing and contact tracing to end the ongoing COVID-19 pandemic.

Subject terms: Epidemiology, Infectious diseases, Computational models, Differential equations

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7175834/

INFN School of Statistics 2022, Paestum 17 May 2022

Gaetano Salina

A New Approach: Complex Networks Models

The Complex Networks Zoo

INFN



FIG. 3 A zoo of complex networks. In this qualitative space, three relevant characteristics are included: randomness, heterogeneity and modularity. The first introduces the amount of randomness involved in the process of network's building. The second measures how diverse is the link distribution and the third would measure how modular is the architecture. The position of different examples are only a visual guide. The domain of highly heterogeneous, random hierarchical networks appears much more occupied than others. Scale-free like networks belong to this domain.





Distribution of connections and the mean degree of a randomly chosen vertex (left) differ sharply from those of end vertices of a randomly chosen edge (right)



(q'-1)/q' ratio is due to the fact that an infected vertex in this model cannot infect back its infector, and so one of the q edges is effectively blocked.



The network is completely described by the joint degree-degree distribution P(q',q) and by N. It is convenient to use a conditional probability P(q'|q) that if an end vertex of an edge has degree q, then the second end has degree q'.

The SIS, SIR, SI, and SIRS models

Four basic models of epidemics are widely used: the SIS, SIR, SI, and SIRS models;

 $S \rightarrow susceptible$

 $I \rightarrow infective$

$R \rightarrow \mbox{ recovered or removed}$

In the network context, vertices are individuals which are in one of these three (S, I, R) or two (S, I) states, and infections spread from vertex to vertex through edges. Note that an ill vertex can infect only its nearest neighbors: S \Rightarrow I.

The SIS model describes infections without immunity, where recovered individuals are susceptible. In the SIR model, recovered individuals are immune forever, and do not infect. In the SI model, recovery is absent. In the SIRS model, the immunity is temporary. The SIS, SIR, and SI models are particular cases of the more general SIRS model. We touch upon only the first three models. (Pastor-Satoras and Vespienal 2002, 200, Newna 2002a, Kenah and Robins 2007)

Complex and neural networks

Lesson 19, AA 2016-201

Complex and neural networks

The SIS, SIR, SI, and SIRS models

Consider the evolution of the probabilities:

 $i_q(t); s_q(t); r_q(t)$

 $i_q(t) = \frac{n_i(q)}{NP(q)}$

 $i_a(t) + s_a(t) + r_a(t) = 1$

that a vertex of degree q is in the I, S, and R states.

Let $\underline{\lambda}$ be the infection rate, a susceptible vertex becomes infected with the probability $\underline{\lambda}$ per unit time if at least one of the nearest neighbors is infected.

 $\underline{\lambda}$ is the only parameter in the SIS and SIR models (other parameters can be easily set to 1 by rescaling)

Lesson 19. AA 2016-2017

Gaetano Salina

INFN School of Statistics 2022, Paestum 17 May 2022

Experience with Statistical Analysis of Covid Data

A New Approach: Complex Networks Models

Epidemic thresholds:

INFN

The epidemic threshold λ_c is a basic notion in epidemiology, the fractions of infected and recovered or removed vertices in the final state (stationary) are defined as:



The Geo-Random has a finite epidemic threshold λ_c . If the spreading rate λ exceeds λ_c , the pathogen becomes endemic. For a Scale-Free network we have $\lambda_c=0$, hence even viruses with a very

small spreading rate λ can persist in the population.

The SIS models on Uncorrelated Network





Gaetano Salina

INFN School of Statistics 2022, Paestum 17 May 2022

Experience with Statistical Analysis of Covid Data



Gaetano Salina

0

0

Ο

Ο



Experience with Statistical Analysis of Covid Data



Gaetano Salina

INFN School of Statistics 2022, Paestum 17 May 2022

Experience with Statistical Analysis of Covid Data

SIR models on 2D Networks

Go beyond the time evolution: Metric Spaces and Network Topology



Gaetano Salina

INFN

INFN School of Statistics 2022, Paestum 17 May 2022



Small World Links; Size = 100; λ = 0.5; +Links: 30 %

Gaetano Salina

Epidemic dynamics on Graphs

INFN School of Statistics 2022, Paestum 17 May 2022

Experience with Statistical Analysis of Covid Data



Experience with Statistical Analysis of Covid Data

SIR models on 2D Networks



Different spread in Time

Gaetano Salina



INFN School of Statistics 2022, Paestum 17 May 2022

Gaetano Salina

SIR models on 2D Networks





The probability to have k-loops goes to zero in thermodynamic limit, i.e. $p(k-loops) \rightarrow 0$. Statistically irrelevant giant loops are present.



Adding SW links, the open m-paths grows, while the number of loops remains constant (V 4-loops). Due to SW links, statistically irrelevant giant loops are present.

A graph's diameter D is the largest number of vertices which must be traversed in order to travel from one vertex to another when paths which backtrack, detour, or loop are excluded from consideration.

In the limit of infinite links added, SW Networks behaves as ER Networks.

SW Networks quickly forget the 2D Square Lattice substrate.

SIR models on 2D Networks

Ising Model on 2D Networks



In different physical models we observe the same effects due to the topological and metric properties of the network underlying the dynamics of the models.

Experience with Statistical Analysis of Covid Data

SIR models on 2D Networks

- Open and Closed Paths
- Clusters Perimeter versus Cluster Volume



Go beyond the time evolution:

Metric Spaces and Network Topology



2D SW Network

The probability to have k-loops goes to zero in thermodynamic limit, i.e. $p(k-loops) \rightarrow 0$. Statistically irrelevant giant loops are present. ER Networks can be seen as Trees.



Adding SW links, the open m-paths grows, while the number of loops remains constant (we have V 4-loops). Due to SW links, statistically irrelevant giant loops are present. SW Networks can be seen as Trees.

2D Square Lattice In thermodynamic limit the ratio k-loops over k-paths remains constant, so that the effects of closed loops cannot be neglected.



Gaetano Salina

INFN School of Statistics 2022, Paestum 17 May 2022

Experience with Statistical Analysis of Covid Data

Size = 100:

λ = 1.6; +Links: 30 %

Go beyond the time evolution: Metric Spaces and Network Topology

SIR models on 2D Networks

2D NN Network

Clusters Perimeter versus Cluster Volume

The impact of heterogeneous connectivity patterns, reflected by an underlying network topology, on epidemic behavior is a key point and must be understood.



Gaetano Salina

INFN

INFN School of Statistics 2022, Paestum 17 May 2022

Go beyond the time evolution: Metric Spaces and Network Topology

Time-dependent hierarchy structure, with both temporal and spatial scales to be determined.



Modeling human mobility is madatory for the construction of a realistic theory of the space-time spread of an epidemic.

The Hong Kong flu was a flu pandemic whose outbreak in 1968 and 1969 killed between one and four million people globally.

The first recorded instance of the outbreak appeared on 13 July 1968 in British Hong Kong.

The first recorded case of the outbreak appeared in Italy was in January 1969, **six months** later.



The complexity and heterogeneity of the present time human mobility network favor considerably the global spreading of infectious diseases. Only unfeasible mobility restrictions reducing the global travel fluxes by 90% or more would be effective.

Experience with Statistical Analysis of Covid Data

Epidemic on complex networks and definition of a Knowledge Base for the analysis of the Covid19 in Italy

> G. Salina INFN Sezione di Roma Tor Vergata

G. Cimini, A. Desiderio and L. Rossi Mori

Physics Department, Tor Vergata University, Rome

A. Gabrielli

Engineering Department, Roma Tre University, Rome

F. Sylos Labini Enrico Fermi Research Center, CREF

L. Vigliano Mathematics Department, Tor Vergata University, Rome



Knowledge Base



INFN School of Statistics 2022, Paestum 17 May 2022

Gaetano Salina

Experience with Statistical Analysis of Covid Data

Experience with Covid Data Analysis

INFN '

Networks Dynamics

Multiplex mobility network and metapopulation epidemic simulations of Italy based on Open Data

The patterns of human mobility play a key role in the spreading of infectious diseases and thus represent a key ingredient of epidemic modeling and forecasting.

Unfortunately, as the Covid-19 pandemic has dramatically highlighted, for the vast majority of countries there is no availability of granular mobility data (not only).

We build a multiplex mobility network **based** solely on **open data**, and implement a SIR metapopulation model that allows scenario analysis through data-driven stochastic simulations.

The mobility flows that we estimate are in agreement with **real-time proprietary data** from smartphones (as Facebook Data for Good).

Our modeling approach can thus be useful in contexts where highresolution mobility data is not available.

	on Ope	n Data								
	Antonio Physics Department, University of Rome Centro Ricerche Enrico F	Desiderio "Tor Vergata", 00133 Rome (Italy) and ermi, 00184 Rome (Italy)								
	Giulio Physics Department, University of Ror Centro Ricerche Enrico Ferr Istituto Nazionale di Fisica Nucleare, Sezion	Cimini ne "Tor Vergata", 00133 Rome (Italy) mi, 00184 Rome (Italy) and e Roma "Tor Vergata", 00133 Rome (Italy)								
7707	Gaetano Salina Istituto Nazionale di Fisica Nucleare, Scione Roma "Tor Veryata", 00133 Rome (Italy) (Dateci. May 10, 0202)									
OC-PIIJ / INTAY	The patterns of human mobility play a key role in the spreading of infections diseases and thus represent a key ingredient of epidenic modeling and forcexating. Unfortunately, as the Covid-19 pandemic has dramatically highlighted, for the vast majority of countries there is no availability of granular mobility data. This hindres the possibility of developing computational frameworks to monitor the evolution of the disease and to adopt timely and adequate prevention policies. Here we show how this problem can be addressed in the case study of Italy. We build a multiplex mobility network based solely on open data, and implement a SIR metapopulation model that allows scenario analysis through data-driven stochastic simulations. The mobility flows that we estimate are in agreement with reak-time proprietary data from smartphones. Our modeling approach can thus be useful in context where high-resolution mobility data is not veniable.									
SINTERNIN I VECUCULUATION SILVIN	<section-header><section-header><text><text><text></text></text></text></section-header></section-header>	gion can spread at different scales [45:47]. Mobility data are crucial in this respect [48]. Unfortunately, in many cases the health authorities had to deal with the scarcity and incompleteness of this type of data, which especially in the early stages of the gandenine prevented them from understanding the spreading patterns of the disease and between the various containment measures imple- tance. The second second second second second second the various containment measures imple- tance. The second second second second second second the various containment measures imple- ments. The second second second second second second the second second second second second second second the second second second second second second second time second second second second second second second time second second second second second second second the second second second second second second second the second second second second second second second test is a second second second second second second second test is a second second second second second second test of study second second second second second test of study mobility at the geographic level of regions or protines. In particular, the authors of [58] have to low is second second second second second second second second second second second second second second second second second second test or study mobility at the geographic level of regions or protines. In particular, the authors of [58] have to low is second second second second second second second second s								

https://arxiv.org/pdf/2205.03639.pdf

Networks Dynamics

The mobility network is defined by representing each of the M = 7897 municipalities as a node. We use ISTAT data to collect information on identification number, reference province, population, surface, latitude and longitude of each municipality. https://www.istat.it/it/archivio/6789

Modelling the Municipalities



Modelling the mobility of individuals At each time step of the dynamics we:

- a) A sample a fraction $p_M = 0.1$ of the population of each municipality is choosen
- b) Move it to other municipalities. This moving population is distributed over the four layers with the p_{Laver}
- c) Execute a Metapopulation Dynamic step
- d) At the end of the current time step, the moving population returns to the municipality of departure.



Gaetano Salina

The connectivity of the Intra-Province and Inter-Province layers is set by the territorial contiguity matrix collected by ISTAT. https://www.istat.it/it/ archivio/157423



i,j = 1,2,3, ...,7897

- 1) Intra-Province Layer All $M_i \iff M_j$ if $c_{i,j} = 0$ and $Pr_i = Pr_j$
- 2) Inter-Province Layer

All $M_i \iff M_j$ if $c_{i,j} = 0$ and $Pr_i \neq Pr_j$

Experience with Statistical Analysis of Covid Data

Modelling the short range networks: Intra-province and Inter-province layers

Connections in these layers are weighted using gravity-like assumptions:

$$w^{\rm L}_{ij} = \frac{\rho_i \rho_j}{r_{ij}} c^{\rm L}_{ij}$$

where ρ_i is the population density of node *i*, $r_{i,j}$ is the distance between the two municipalities and $L \in \{$ Intra, Inter $\}$ denotes the layer.

In order to transform such weights into Markovian connection probabilities we use the following normalization

$$\tilde{w}_{ij}^{\mathrm{L}} = \frac{w_{ij}^{\mathrm{L}}}{\sum_{l \in n_i^{\mathrm{L}}} w_{il}^{\mathrm{L}}} c_{ij}^{\mathrm{L}}$$

where n_i^L are the territorial neighbours of node *i* in the corresponding layer.



Gaetano Salina

INFN School of Statistics 2022, Paestum 17 May 2022

Networks Dynamics

The Trains layer

INFN

We use the RESTful API of Viaggiatreno to retrieve information. http://www.viaggiatreno.it/

- $\circ~$ Each station is identified by a id, latitude and longitude
- Each train departing from station i we have the id j of the next stop station and the type of connection *e* (Frecciarossa, Intercity, Metropolitana, EuroNight, Regionale, EuroCity)
- We associated each station with the municipalities that are located no more than 5 km away (in this way, on average a single station corresponds to three municipalities).



 au_{ij}^e is boolean indicator for the presence of a connection g^e is the fraction of total travellers for train type *e*. Dati ISTAT

https://www.istat.it/it/archivio/13995

Experience with Statistical Analysis of Covid Data

Modelling the long range networks: Train layer



INFN School of Statistics 2022, Paestum 17 May 2022

Gaetano Salina

Networks Dynamics Mobility Network

The Flights layer

INFN

- To build the Flight layer, we collect IATA and ICAO codes, latitudes and longitudes of airports from Open- Flights https://openflights.org/data.html
- From the annual ENAC 2019 report we collected the air traffic data: origin and destination airports, with IATA and ICAO Ids, and total passengers moving between airports. https://www.enac.gov.it/pubblicazioni/dati-di-traffico-2019
- We associated each airport with all municipalities located no more than 15 km away (on average a single airport is mapped on ten municipalities).

$$w_{ij}^{ ext{Flight}} = g_{ij}$$
 $ilde{w}_{ij}^{ ext{L}} = rac{w_{ij}^{ ext{L}}}{\sum_{l \in n_i^{ ext{L}}} w_{il}^{ ext{L}}}$ L = Flights

 g_{ij} is the fraction of total travellers *i-j*. Dati ENAC

Experience with Statistical Analysis of Covid Data

Modelling the long range networks: Flights layer



INFN School of Statistics 2022, Paestum 17 May 2022

Gaetano Salina



Gaetano Salina

INFN School of Statistics 2022, Paestum 17 May 2022

Experience with Statistical Analysis of Covid Data

Building a Knowledge Base



Federico Zuccari, Inferno Canto III.

CovKB: A Knowledge Base for the analysis of the Covid19 epidemic in Italy

G. Salina INFN Sezione di Roma Tor Vergata

G. Cimini and A. Desiderio Physics Department, Tor Vergata University, Rome

> F. Sylos Labini Enrico Fermi Research Center, CREF

L. Vigliano Mathematics Department, Tor Vergata University, Rome

General Framework

The outbreak of an epidemic crisis has similarities with the occurrence of an earthquake, where we do not have access to the state of the system, and thus we cannot make a precise forecast, but where we can monitor what happens and adopt adequate prevention policies. Indeed, that a global epidemic of the Covid-19 type could be triggered with a non-negligible probability was known since some decades. This is due to a series of factors, from the increase in promiscuity between humans and wild species to the development of a global interconnected world. However, the prediction of when and where an epidemic of this kind would have erupted is not possible because there are no data available to predict such an event.

Instead, it is possible to intervene promptly when the epidemic has started to spread because it is only then that it is possible to have data on its development. Computational epidemiology uses diffusion models on complex networks to make forecasts if big data concerning the number of patients, the mobility of people, etc are available. In this sense, the forecasts have a similar character to those of meteorology in which the system is observed at a certain time and the equations of fluid dynamics and thermodynamics are appropriately simplified to calculate the weather tomorrow. The combination of computing power and big data plays a fundamental role in both cases.

The problem in calculating the spread of the virus concerns the sensitive dependence on microscopic conditions and the extreme heterogeneity in its propagation. This intrinsic chaoticity in the diffusion process complicates the outlook for the epidemic a lot, especially when the diffusion is far from uniform, even when big data, computing power and mathematical models are available. The role of forecasts in the case of the spread of an epidemic is to help the political decision maker to make the most appropriate choices. The pandemic crisis cannot be avoided, but it can be avoided that it has disastrous effects with appropriate interventions, a priori and posteriori, and it is possible to change its evolution with appropriate policies. The lessons for policy makers are that of developing prevention and support systems that helps at the right time to face an epidemic. The aim of this project is that of providing a structured dataset to give a rational basis to political decision makers

Definition of a Knowledge Base (KB) for the analysis of the Covid19 epidemic in Italy (CovKB) The defection of the CovKB will permit the integration of empirical mobility networks (short and long range) with the other subpopulation attributes (population age structure and density, health and economic infrastructures, etc.). Such a networked data structure, modulated by the timeline of

Experience with Statistical Analysis of Covid Data

Building a Knowledge Base

The Knowledge Discovery Process (KDP)

- Selection: Capturing relevant prior knowledge, identifying the data-mining goal and developing and understanding of the application domain. Based on that, *proper data samples as well as relevant variables* can be selected.
- Pre-processing: The selected data are processed. Handling of missing values, the noise and errors identification (and correction), the elimination of duplicates, as well as the matching, fusion, and conflict resolution for data taken from different sources are done.
- Transformation: The clean dataset is transformed into a form suitable for data mining algorithms analysis. To improve the analysis performance dimensionality reduction methods can also be applied.
- o Data mining:
 - Heuristic Models: the goal of the is matched to a particular method, such as classification, regression, or clustering the transformed data. A decision about which models and parameters might be appropriate must be done and matching a particular data mining method with the overall criteria of the KDP in Datasets process is necessary.
 - Theoretical Models: the transformed data are used to derive the "interaction" parameters of a theoretical dynamic model. The KDP coincides with the understanding of the system's physical dynamics.
- Evaluation and interpretation: The patterns and models derived are analysed with respect to their validity. The user assesses the usefulness of the found knowledge. A data visualization of the extracted patterns and models is involved.

The first three steps define the ETL process: these will allow one to build a solid structure for the KB, defining its fundamental variables, computational models and logical and physical architecture, i.e. the whole set of information (dataset and metadata). This is particularly relevant, given that the analysis of the Covid19 epidemic has highlighted the lack of a solid KB useful for studying the phenomenon itself: few data available, affected by systematics that are sometimes not well defined and of poor quality.

The goal is define and implement a KB that:

- \circ $\,$ identifies the different domains and the different data sources
- is integrated and standardized between the data of the various sources, with high quality, unambiguous and homogeneous data
- is useful for studying the Covid19 in Italy and for discovering the factors that determined its evolution
- supports the simulation of different scenarios, regarding different containment strategies
- supports policymakers' decisions to be taken based on a Data-Driven approach
- is useful for defining a response strategy for any future epidemics by means of a predictive model

To reread them today, the last two points are pure Utopia

Experience with Statistical Analysis of Covid Data

Building a Knowledge Base

The KB must contain Data and Metadata of different types:

- Concerning the parameters of the epidemic of the individuals (provinces as elementary entities)
- Concerning to the description of the spatial distribution of the elementary sets of then individuals (municipalities as elementary entities) and of all the characteristics of these elementary sets (geolocation, inhabitants, surface, urbanization level, demographic distribution, mortality data, ethnic distribution, geophysical variables, etc.)
- Concerning the structure of the National Health System (territorial location of hospitals and treatment centres and their characteristics as beds, hospitalizations, etc.). The following time series should be at a daily detail: emergency interventions and their description, clinical trials aimed at monitoring certain diseases.
- Concerning response operations aimed at containing the epidemic (government's decrees, identification of health structures aimed at curing Covid19, identification of municipalities subject to further restrictions, red areas, etc.)
- Concerning to mobility (origin destination matrices of individual movements at municipal level, road and motorway distance matrix, road and motorway flows, railway, air and naval network and respective flows)
- Concerning to the productive structure (number of industries and number of employees at the municipal level by Ateco class, matrix of goods/people flows generated by the production processes).

Epidemic Data:

- o github.com/pcm-dpc/COVID-19, PL
- o https://github.com/collections/open-data, ML Some
- o ISS DataSet, PL Not Open

Demographic and Environmental Data:

- o www.istat.it, ML
- Open Data Sites of ARPA Regional Agencies, ML
- Open Data Sites of Italian Regional Administrations, ML

National Health System Data:

- o www.dati.salute.gov.it/dati/homeDataset.jsp, ML
- o www.iss.it/basi-di-dati, ML
- Open Data Sites of Italian Regional Administrations, ML

Mobility Data:

- o www.istat.it, ML
- o dati.mit.gov.it/catalog/dataset, ML
- $\circ \quad www.stradeanas.it/it/le-strade/osservatorio-del-traffico$
- o www.enac.gov.it/open-data
- o dataforgood.fb.com/docs/covid19/, PL & ML
- https://news.google.com/covid19/map?hl=it&gl=IT&ceid=IT%3Ai, PL

Open Data Regional Administrations Economic Data:

- o www.istat.it, ML
- o www.camcom.gov.it/P43K973O0/open-data.htm, ML

Generic Data:

- www.dati.gov.it
- www.datiopen.it/it

Main Open Data Resouces

Administrative Elementary Entities: ML = Municipal Level. PL = Provincial Level .

Building a Knowledge Base

The KB must contain Data and Metadata of different types:

- Concerning the parameters of the epidemic of the individuals (provinces as elementary entities)
- Concerning to the description of the spatial distribution of the elementary sets of then individuals (municipalities as elementary entities) and of all the characteristics of these elementary sets (geolocation, inhabitants, surface, urbanization level, demographic distribution, ethnic distribution, geophysical variables, etc.)
- Concerning the structure of the National Health System (territorial location of hospitals and treatment centres and their characteristics as beds, hospitalizations, etc.). The following time series should be at a daily detail: emergency interventions and their description, clinical trials aimed at monitoring certain diseases.
- Concerning response operations aimed at containing the epidemic (government's decrees, identification of health structures aimed at curing Covid19, identification of municipalities subject to further restrictions, red areas, etc.)
- Concerning to mobility (origin destination matrices of individual movements at municipal level, road and motorway distance matrix, road and motorway flows, railway, air and naval network and respective flows)
- Concerning to the productive structure (number of industries and number of employees at the municipal level by Ateco class, matrix of goods/people flows generated by the production processes).

Experience with Statistical Analysis of Covid Data

8	Structure 🥂 SQL 🔑 Sea	rch	A Query	Bexport A	Import %	Operations	B 😭 Privil	eges 🚲	Routines 🕑	Events	🞎 Triggers 🖓 D	esigner	
	Table 🔺				Action				Rows (7)	Туре	Collation	Size	Overhead
	Areoporti	索	Browse	Structure	Search	3-é Insert	Empty	🗙 Drop	~68,840	InnoDB	utf8_general_ci	8.5 MiB	
	CovidMortiGiorno	숥	Browse	Structure	E Search	s-i Insert	Empty	X Drop	10,368	InnoDB	latin1_swedish_ci	480 KiB	
	CovidMortiGiorno-090420	余	Browse	Structure	Search	-i Insert	Empty	X Drop	8,942	InnoDB	latin1_swedish_ci	416 KiB	
	CovidProvMortiGiorno	숥	T Browse	Structure	🔛 Search	3-é Insert	Empty	🗙 Drop	~91,787	InnoDB	latin1_swedish_ci	13.5 MiB	
	CovidRegMortiGiorno	余	Browse	Structure	Search	3-i Insert	Empty	🗙 Drop	16,968	InnoDB	latin1_swedish_ci	3.5 MiB	
	DataClusterMean	숥	Browse	Structure	Search	3-i Insert	Empty	X Drop	~3,894,502	InnoDB	utf8_general_ci	219.8 MiB	
	DataSimulations	余	Browse	Structure	E Search	3-i Insert	Empty	× Drop	~89,279	InnoDB	utf8_general_ci	19.5 MiB	
	DPMapsColocation		Browse	Structure	Search	s-i Insert	Empty	X Drop	~979,605	InnoDB	utf8_general_ci	309.9 MiB	
	DPMapsGBMovAdmin	會	Browse	Structure	E Search	3-i Insert	Empty	× Drop	~25,171,578	InnoDB	utf8_general_ci	11 GiB	
	DPMapsGEMovAdmin	\$	T Browse	Structure	🔝 Search	s-i Insert	The Empty	X Drop	0	InnoDB	utf8_general_ci	16 KiB	
	DPMapsMovAdmin	會	Browse	Structure	E Search	3-i Insert	Empty	× Drop	~1,582,854	InnoDB	utf8_general_ci	711 MiB	
	DPMapsMovRange	\$	Browse	Structure	E Search	3-i Insert	Empty	× Drop	16,400	InnoDB	utf8_general_ci	1.5 MiB	
	DPMapsPopAdmin	全	Browse	PS Structure	Search	i Insert	Empty	× Drop	~142,507	InnoDB	utf8_general_ci	72.6 MiB	
	DPMapsPopTiles	\$	Browse	Structure	Search	se Insert	Empty	× Drop	~44,394,697	InnoDB	utf8_general_ci	18.8 GiB	
	FlightMap	會	Browse	Structure	📜 Search	s-i Insert	The Empty	X Drop	0	InnoDB	utf8_general_ci	16 KiB	
	GoMaps	\$	Browse	Structure	Search	se Insert	Empty	× Drop	~28,566	InnoDB	utf8_general_ci	3.5 MiB	
	IndustriaComuni	*	Browse	Structure	Search	3-i Insert	Empty	× Drop	4,573	InnoDB	latin1_swedish_ci	256 KiB	
	italy_cities		Browse	Structure	Search	si Insert	Empty		7,911	InnoDB	latin1_swedish_ci	1.5 MiB	
	italy_geo	*	Browse	Structure	Search	3-i Insert	Empty		7,911	InnoDB	latin1_swedish_ci	1.5 MiB	
	italy_geo_old	\$	Browse	Structure	E Search	s-i Insert	Empty	X Drop	7,954	InnoDB	latin1_swedish_ci	1.5 MiB	
	italy_metropolitane	\$	Browse	Structure	Search	3-i Insert	Empty		14	InnoDB	latin1_swedish_ci	16 KiB	
	italy_provincies	\$	Browse	Structure	E Search	s-i Insert	Empty	X Drop	107	InnoDB	latin1_swedish_ci	16 KiB	
	italy_regions	*	Browse	Structure	Search	3-i Insert	Empty		20	InnoDB	latin1_swedish_ci	16 KiB	
	LombardiaRsa		Browse	Structure	Search	3-i Insert	Empty	× Drop	706	InnoDB	latin1_swedish_ci	96 KiB	
	MatriceComuni	*	Browse	Structure	Search	s-i Insert	Empty	× Drop	~46,328	InnoDB	latin1_swedish_ci	3.5 MiB	
	MortiGiornoEta	\$	Browse	Structure	Search	3-é Insert	Empty	× Drop	~3,675,652	InnoDB	latin1_swedish_ci	452.9 MiB	
	MortiGiornoRegioni	*	Browse	Structure	Search	3-i Insert	Empty	× Drop	7,644	InnoDB	utf8_general_ci	1.5 MiB	
	SourceMatriciOD	\$	Browse	Structure	Search	3-é Insert	Empty		~5,969,320	InnoDB	utf8_general_ci	669 M1B	
	TurismoComuni	*	Browse	Structure	Search	3-i Insert	Empty	× Drop	7,906	InnoDB	utf8_general_ci	448 KiB	
	ZoneCovid	4	Browse	Structure	Search	i Insert	Empty	× Drop	158	InnoDB	utf8_general_ci	16 K18	
	30 tables			-	Sum				~86.233.089	InnoDB	utt8 general ci	32.2 GiB	

~ 90.000.000 Records



Some Analysis

Mortality March 2020 \leftrightarrow November 2020



Mortality Istat (CovKB): Source https://www.istat.it/it/archivio/240401



RSA are Nursing home for the elderly

Experience with Statistical Analysis of Covid Data

Gaetano Salina

Experience with Covid Data Analysis

Some Analysis

INFN

Effect of the Vaccination Campaign: January 2020 ←→ November 2021

- In blue the incidence (number of infections per 100,000 individuals) processed by the data of the Civil Protection from March 2020 for all age groups, Full Population.
- In red (right axis) the ratio between the Susceptible Population and the Full Population. This ratio is less than one since the vaccination campaign.
- The Susceptible Population is evaluated starting from the data on the temporal evolution of vaccines. It is the sum of the number of people not vaccinated plus the fraction of the cohorts (vaccinated one dose, vaccinated two doses, vaccinated one dose + previous infection) multiplied by $1 - \alpha$, α is the average efficacy factor of the vaccines (from ISS). In orange the incidence calculated on the Susceptible Population.



Experience with Covid Data Analysis Experience with Statistical Analysis of Covid Data INFN Mobility Data FACEBOOK Data for Good Colocation Matrix: $M_{A \rightarrow B} A, B$ Italian Provinces Mobility Data Data for Good at Meta Datasets ending soor After May 24, 2022, many datasets that rely on location data will no longer ☆ Home Find datasets by rus Disease Prevention Map Feb 24 2020 Id Q Explore datasets Italy Coronavirus Dataset type Dataset date ∫∃ File Population Density Maps 1 mar 2020 - 12 mag 2022 Facebook Population (Administrative Region) Documentazione Azzera i filtri Filtri Vector: P_A A Italian Province Dataset type · Facoltativo Our Population Density Maps help nonprofit and multilateral agencies plan vaccination campaigns, respond to natural disasters, and evaluate rural Tiles 1.6 x 1.6 Km² Q electrification plans. These maps help researchers assess the ways in which Italy Coronavirus Disease Prevention Map Feb 24 2020 Id **NW Italy** Paese · Facoltativo climate change and urbanization impact where people live. Dataset date Q Dataset type 1 mar 2020 - 12 mag 2022 Facebook Population (Tile Level Vector: P_A A Tiles Movement Maps Aggregate information from people using Facebook on their mobile Italy Coronavirus Disease Prevention Map Feb 24 2020 Id hones with Location History enabled, showing movement between two Dataset date Dataset type points. The methodology can be found here 1 mar 2020 - 10 mag 2022 Movement Between Tiles OD Matrix: $M_{A \rightarrow B} A.B$ Tiles **Co-location Maps** These improved co-location insights help researchers better understand Italy Coronavirus Disease Prevention Map Feb 24 2020 Id general movement patterns like commute patterns, and the probability Dataset date Dataset type that people in areas with disease outbreaks will come in contact with 3 mar 2020 - 3 mag 2022 Colocation new populations. These maps are important to understand the actual chance that a disease will be spread by human-to-human contact. OD Matrix: $M_{A \rightarrow B} A, B$ Italian Provinces Movement Maps Italy Coronavirus Disease Prevention Map Feb 24 2020 Id Aggregate information from people using Facebook on their mobile Dataset date Dataset type phones with Location History enabled, showing movement between two 1 mar 2020 - 10 mag 2022 Movement Between Administrative Regions points. The methodology can be found here OD Matrix: $M_{A \rightarrow B} A, B$ Italian Provinces

INFN Mobility Data

Mobility Data FACEBOOK Data for Good





Population Density Maps

Our Population Density Maps help nonprofit and multilateral agencies plan vaccination campaigns, respond to natural disasters, and evaluate rural electrification plans. These maps help researchers assess the ways in which climate change and urbanization impact where people live.

From My post on FB.

Tonight I was "lucky" to hear an interview with Flavio Briatore. He rattled off numbers with extreme ease about the economic situation. The message is always the same: "The economic crisis will kill more than the coronavirus, and that closing restaurants and nightclubs as well as being crazy is useless in order to fight the epidemic."

I would like to contrast the following graph with Briatore's numbers. In blue the infections in the province of Sassari. In orange the FB index proportional to the population present in the province of Sassari normalized to the average of the index value in the first 15 days of March. In red the same quantity in the Porto Cervo area where II Billionare is located (it is a small square of 1.6 Km x 1.6 Km).

I don't think any comments are needed.

Experience with Statistical Analysis of Covid Data



blothilly Data 2020-00-03 Holdhilly Data 2020-06-03 Holdhilly Data 2020-0

Co-location Maps

These improved co-location insights help researchers better understand general movement patterns like commute patterns, and the probability that people in areas with disease outbreaks will come in contact with new populations. These maps are important to understand the actual chance that a disease will be spread by human-to-human contact.

The denominator is the maximum number of colocations observed in the week. The link values are hence always symmetric as well



INFN School of Statistics 2022, Paestum 17 May 2022

Gaetano Salina

Some Analysis

Experience with Statistical Analysis of Covid Data





Human mobility and the spread of the pandemic

Some of my Considerations

- *"Few data are available, affected by systematics that are sometimes not well defined and of poor quality". It was true two years ago, remains true today.*
- Paradoxically, in the first months of 2020 there was a greater availability of data. Movements for OpenData released epidemic data at the municipal level obtained with scraping techniques on institutional servers. When the fact was publicly known all the servers were removed from Internet.
- The lack of granular data for the outbreak is a purely political choice. There was no possibility of agreeing anything in the ISS-Lincei and ISS-INFN agreements. This makes the ISS dataset, in fact, no more useful than that of civil protection.
- Facebook offered better mobility data than those provided by the big tech companies, but their metrics do not allow easy integration with the models.
- You should start with the data of the telephone operators, but they are expensive and the work to arrive at the mobility data is enormous.
- Building a knowledge base for Covid was a beautiful illusion, but the task is bigger than the capabilities of our small group.
- However, the built database is very useful as a support to the study of realistic diffusion models on complex networks.

Gaetano Salina

Conclusion

INFN



I go back to simulating Physical Systems on complex topology networks.

At least, in case of missing or poor quality data, I know who to contact.

Grazie !