

STEP 1: Literature Review and Formulation of the Definition

- Over the past several years we have developed a case-based, mixed-methods, density approach to modeling the temporal and spatial complexities of big data.

STEP 2: Methods

- The platform for this approach is called the SACS Toolkit. In terms of simplifying assumptions, the Toolkit employs three novel solutions:

- (1) it conceptualizes the complex causal organization of a system as a set of microscopic cases (k-dimensional vectors spaces);
- (2) it clusters/groups cases to identify major and minor profiles and (discrete or continuous) trajectories
- (3) it translates their high-dynamic microscopic trajectories into the linear movement of macroscopic, low-dynamic densities.

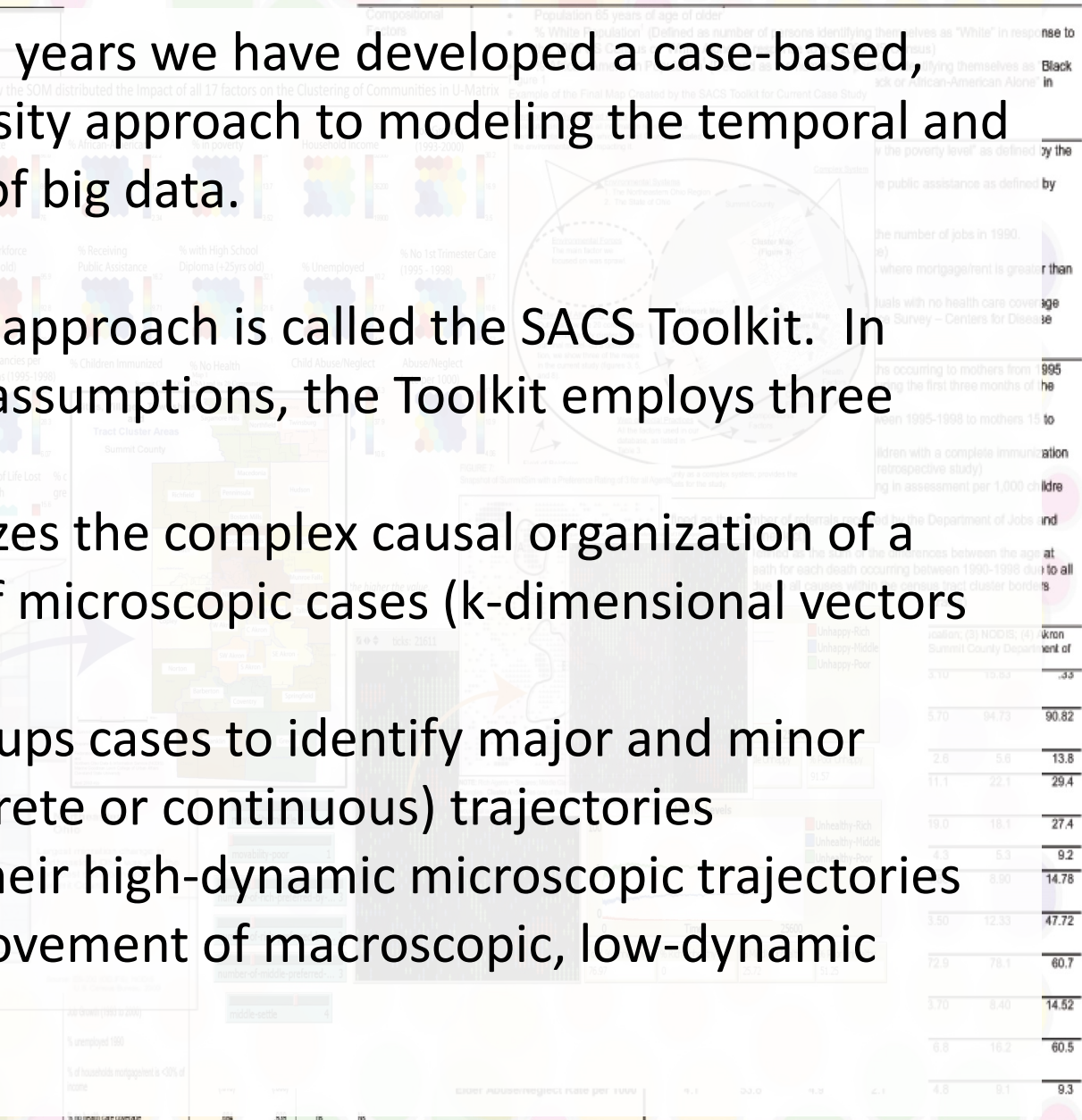
STEP 3: Run Test

STEP 4: Determine Results

QUESTION SET
determined?
7. Did the test s
8. Did the test s
9. In terms of the
a. Did the
b. Did the
c. Does the



NOTE: Distances between clusters are based on Euclidian distances arrived at through k-means analysis. Distances within clusters for each community are based on within-cluster measures. All measures are non-standardized.



COLUMN 1 provides zero-order, pairwise correlations for all compositional and contextual factors listed in Table 3 with two health outcomes, years of life lost per death and Teen Birth Rate. In this column, (*) is the correlation coefficient, and (**) is a two-tailed, significance level.

COLUMN 2 provides the results of our hierarchical analysis of the "independent" relationships all compositional and contextual factors listed in Table 3 with two health outcomes, years of life lost per death and Teen Birth Rate. In this column (a) is a non significant partial correlation coefficient, "" is a significant partial correlation coefficient for a two-tailed significance level.

Years Lost per Death 1998	13.83	16.40	13.98	10.50	10.60	14.40	15.18
1. (*)	13.83	16.40	13.98	10.50	10.60	14.40	15.18
2. Community Membership for each of the 7 Clusters is as follows:							
Cluster 1: Slow/Silverlake, Northfield/Lacedonia/Sagamore, and Richfield/Penninsula.							
Cluster 2: Central Akron.							
Cluster 3: Twinsburg, Northwest Akron, Munroe Falls/Tallmadge, Norton and Franklin.							
Cluster 4: Hudson.							
Cluster 5: Copley/Bath/Fairlawn.							
Cluster 6: Springfield, Coventry/Green and Cuyahoga Falls.							
Cluster 7: North, West, Southwest, South and Southeast Akron and Barberton City.							

STEP 1: Literature Review and Formulation of the Definition

- The strengths of this approach are several. It allows researchers to:
 - Model complex systems as sets of cases.
 - Explore these systems at multiple levels.
 - Examine the interactions between system and environment.
 - Explore the relationships amongst the cases (networks).
 - Address and combine both structure (organizational pattern) and agency.
 - Study complex causal structure.
 - Use small to big data.
 - Model these systems as static or longitudinal.
 - In terms of longitudinal, we can model as discrete or continuous
 - In terms of continuous modeling, we can:
 - map the complex, nonlinear evolution of ensembles (or densities) of cases;
 - classify major and minor clusters and time-trends;
 - visually identify dynamical states, such as saddles and attractor points;
 - plot the speed of cases along different states;
 - detect the non-equilibrium clustering of case trajectories during key transient times;
 - construct multiple models to fit novel data;
 - predict future time-trends and dynamical states; and, finally, in terms of impact,
 - generate results that are visually and conceptually intuitive to private/public sector users and policy makers.

NOTE: Distances between clusters are based on Euclidian distances arrived at through k-means analysis. Distances within clusters for each community are based on within-cluster measures. All measures are non-standardized.

COLUMN 2 provides the results of our hierarchical analysis of the "independent" relationships all compositional and contextual factors listed in Table 3 with two health outcomes: years of life lost per death and Teen Birth Rate. In this column (a) is a non significant partial correlation coefficient. "" is a significant partial correlation coefficient for a two-tailed significance level.

Clusters is as follows: Cluster 1: Slow/Silverlake, Northfield/Lacedonia/Sagamore, and Richfield/Peninsula; Cluster 2: Central Akron; Cluster 3: Twinsburg, Northwest Akron, Munroe Falls/Tallmadge, Norton and Franklin; Cluster 4: Hudson; Cluster 5: Copley/Bain/Fairlawn; Cluster 6: Springfield, Coventry/Green and Cuyahoga Falls; Cluster 7: North, West, Southwest, South and Southeast Akron and Barberton City.

STEP 1: Literature Review and Formulation of the Definition

FIGURE 2
Final SOM Solution for 20 Communities in Summit County

Cases Are Complex Systems

- Researchers in the social sciences currently employ a variety of mathematical/ computational models for studying complex systems.

STEP 2: Methods

- Despite the diversity of these models, the majority can be grouped into one of four types:

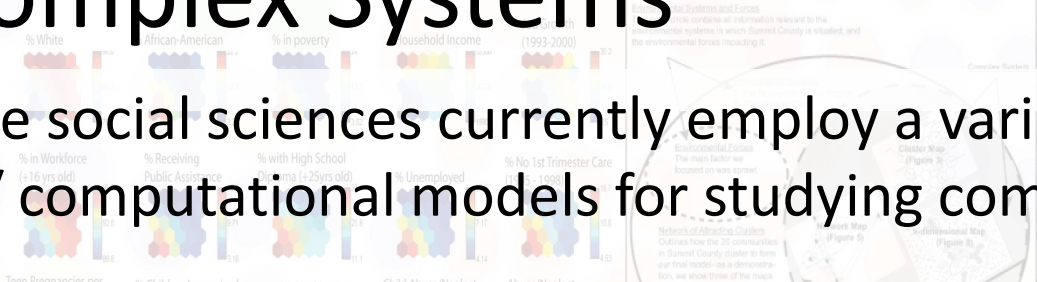
STEP 3: Run Test

STEP 4: Determine Results

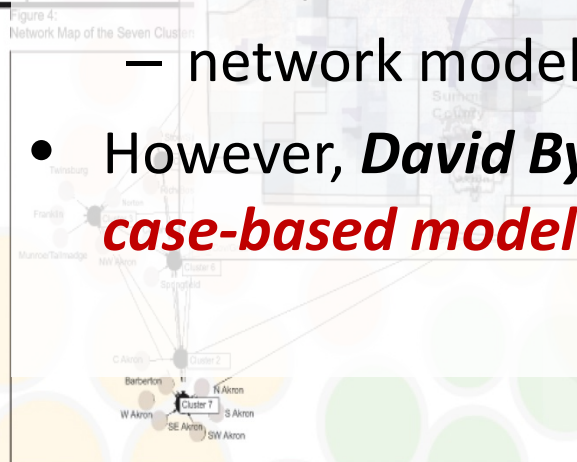
- equation-based modeling,
- stochastic (statistical) modeling,
- computational modeling
- network modeling.

- However, **David Byrne** and colleagues have added a fifth type: **case-based modeling**

Compositional Factors	Variables
	• Population 65 years of age or older
	• % White Population (Defined as number of persons identifying themselves as "White" in response to the 1990 US Census or "White Alone" in response to the 2000 US Census)
	• % African-American Population (Defined as the number of persons identifying themselves as "Black" or "African-American Alone" in response to the 1990 US Census or "Black or African-American Alone" in response to the 2000 US Census)



Cluster	Unhappy-Rich	Unhappy-Middle	Unhappy-Poor
1	3.10	19.93	3.39
2	5.70	94.73	90.82
3	2.6	5.6	13.8
4	11.1	22.1	29.4
5	19.0	18.1	27.4
6	4.3	5.3	9.2
7	4.80	8.90	14.78
8	12.33	47.72	60.7
9	3.70	8.40	14.52
10	6.8	16.2	60.5
11	4.8	9.1	9.3



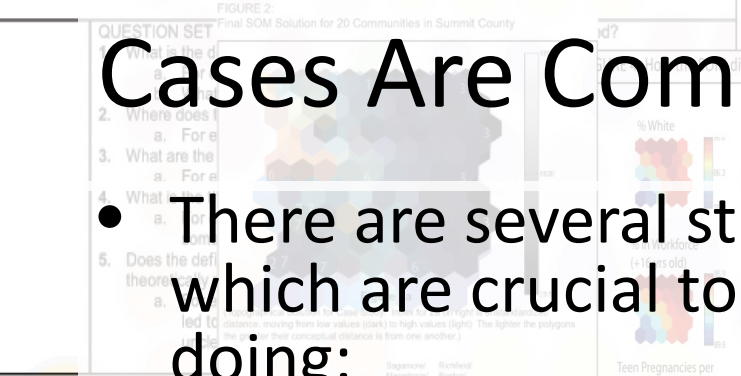
COLUMN 1 provides zero-order, pairwise correlations for all compositional and contextual factors listed in Table 3 with the health outcomes, years of life lost per death and Teen Birth Rate. In this column, (**) is the correlation coefficient, and (**) is a two-tailed, significance level.

COLUMN 2 provides the results of our hierarchical analysis of the "independent" relationships all compositional and contextual factors listed in Table 3 with the health outcomes, years of life lost per death and Teen Birth Rate. In this column (**) is a non significant partial correlation coefficient, (**) is a significant partial correlation coefficient for a two-tailed significance level.

Years Lost per Death 1998	13.83	16.40	13.98	10.50	10.60	14.40	15.18
---------------------------	-------	-------	-------	-------	-------	-------	-------

1. (*) The values listed in the columns for all 7 clusters represent the average value/measurement that the communities in that cluster scored for each variable listed in Column 1. In cluster analysis, these averages are called the cluster's centroids. 2. Community Membership for each of the 7 Clusters is as follows: Cluster 1: Silesville, Northfield/Lacedonia/Sagamore, and Richfield/Penninsula; Cluster 2: Central Akron; Cluster 3: Twinsburg, Northwest Akron, Munroe Falls/Tallmadge, Norton and Franklin; Cluster 4: Hudson; Cluster 5: Copley/Bath/Fairlawn; Cluster 6: Springfield, Cuyahoga/Green and Cuyahoga Falls; Cluster 7: North, West, Southwest, South and Southeast Akron and Barberton City.

STEP 1: Literature Review and Formulation of the Definition



Cases Are Complex Systems

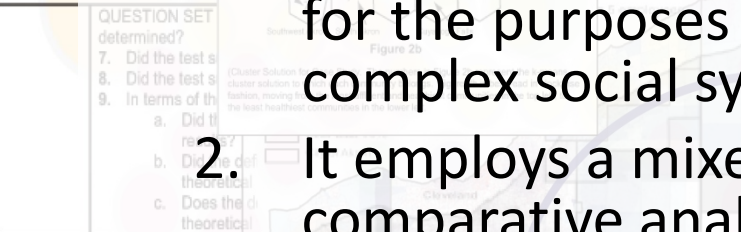
• There are several strengths to this approach, **three** of which are crucial to the work my colleagues and I are doing:

STEP 2: Methods

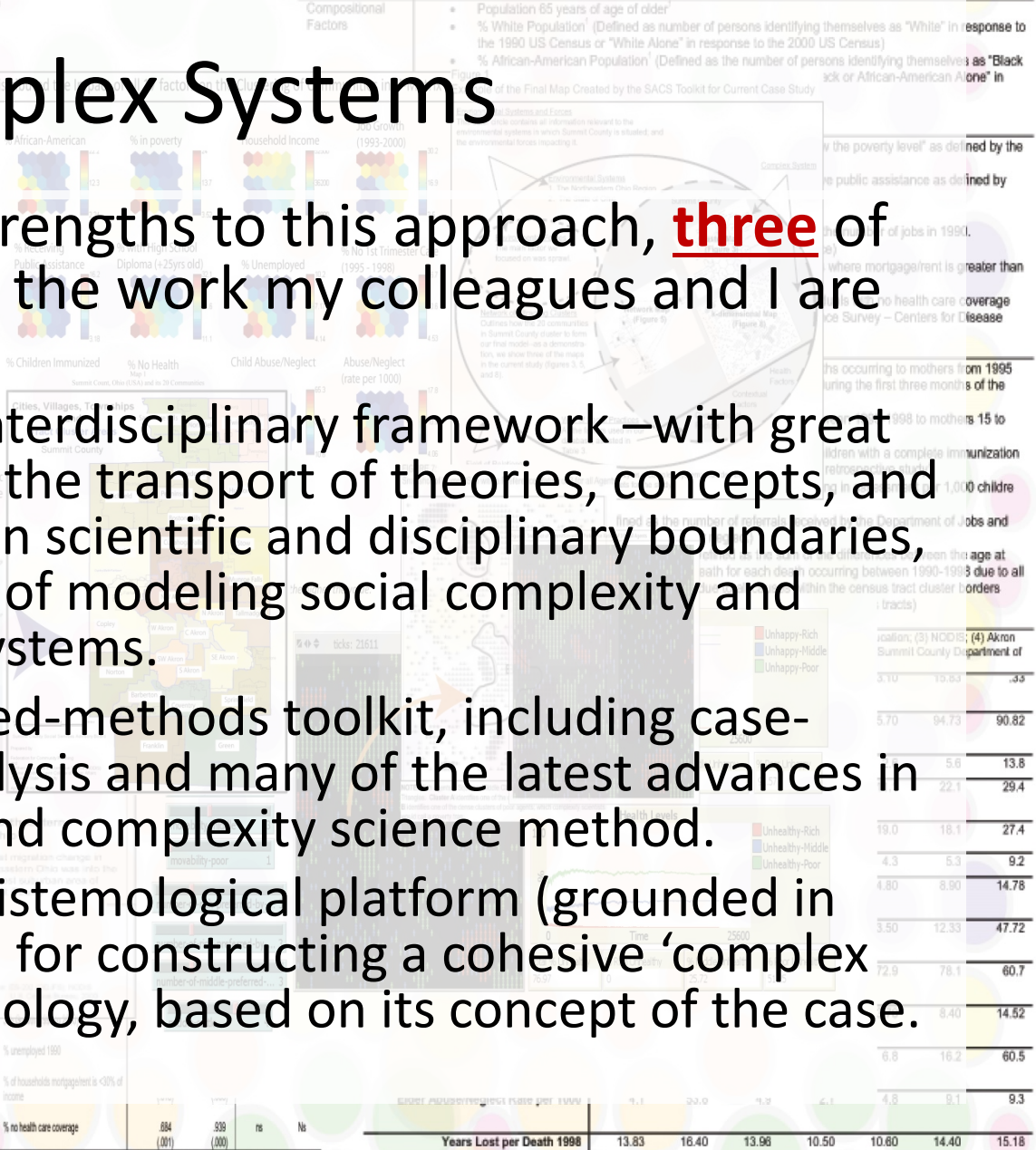


STEP 3: Run Test

STEP 4: Determine Results



1. It embraces an interdisciplinary framework –with great thought given to the transport of theories, concepts, and methods between scientific and disciplinary boundaries, for the purposes of modeling social complexity and complex social systems.
2. It employs a mixed-methods toolkit, including case-comparative analysis and many of the latest advances in computational and complexity science method.
3. It provides an epistemological platform (grounded in complex realism) for constructing a cohesive ‘complex systems’ methodology, based on its concept of the case.



NOTE: Distances between clusters are based on Euclidian distances arrived at through k-means analysis. Distances within clusters for each community are based on within-cluster measures. All measures are non-standardized.

COLUMN 1 provides zero-order, pairwise correlations for all compositional and contextual factors listed in Table 3 with two health outcomes: years of life lost per death and Teen Birth Rate. In this column, (*) is the correlation coefficient, and (**) is a two-tailed, significance level.
 COLUMN 2 provides the results of our hierarchical analysis of the "independent" relationships of all compositional and contextual factors listed in Table 3 with two health outcomes: years of life lost per death and Teen Birth Rate. In this column (ns) is a non-significant partial correlation coefficient, "*" is a significant partial correlation coefficient for a two-tailed significance level.

1. (*) The values listed in the columns for all 7 clusters represent the average value/measurement that the communities in that cluster scored for each variable listed in Column 1. In cluster analysis, these averages are called the cluster's centroids. 2. Community Membership for each of the 7 Clusters is as follows: Cluster 1: Snow/Silverlake, Northfield/Macedonia/Sagamore, and Richfield/Penninsula; Cluster 2: Central Akron; Cluster 3: Twinsburg, Northwest Akron, Munroe Falls/Talmadge, Norton and Franklin; Cluster 4: Hudson; Cluster 5: Copley/Baldwin; Cluster 6: Springfield, Coventry/Green and Cuyahoga Falls; Cluster 7: North, West, Southwest, South and Southeast Akron and Barberton City.

STEP 1: Literature Review and Formulation of the Definition

FIGURE 2
Final SOM Solution for 20 Communities in Summit County

Cases Are Complex Systems

To begin, we have introduced two new terms:

— **case-based complexity science** is the attempt to actively integrate case-based method with the latest developments in the complexity and social sciences for the purpose of modeling complex social systems as sets of cases.

- It also revolves around a particular set of epistemological assumptions:
- Complexity theory is not so much a substantive theory, as much as it is an epistemologically explicit attempt to model social life in complex systems terms.
- It also revolves around complex realism

— In turn, **case-based modeling** is the mixed-methods set of techniques scholars use to engage in case-based complexity science, particularly the latest developments in the computational and complexity sciences.

- The key to this approach is that the methods serve the purpose of case-comparative analysis, from small to big data!

STEP 2: Methods

STEP 3: Fun Test

STEP 4: Determine Results

Figure 4
Network Map of the Seven Clusters



NOTE: Distances between clusters are based on Euclidian distances arrived at through k-means analysis. Distances within clusters for each community are based on within-cluster measures. All measures are non-standardized.

COLUMN 2 provides the results of our hierarchical analysis of the "independent" relationships all compositional and contextual factors listed in Table 3 with two health outcomes: years of life lost per death and Teen Birth Rate. In this column (a) is a non significant partial correlation coefficient, "" is a significant partial correlation coefficient for a two-tailed significance level.

1. (*) The values listed in the columns for all 7 clusters represent the average value measurement that the communities in that cluster scored for each variable listed in Column 1. In cluster analysis, these averages are called the cluster's centroids. 2. Community Membership for each of the 7 Clusters is as follows: Cluster 1: Snow/Silverlake, Northfield/Lacedonia/Sagamore, and Richfield/Penninsula; Cluster 2: Central Akron; Cluster 3: Trumborg, Northwest Akron, Munroe Falls/Tallmadge, Norton and Franklin; Cluster 4: Hudson; Cluster 5: Copley/Bain/Fairlawn; Cluster 6: Springfield, Coventry/Green and Cuyahoga Falls; Cluster 7: North, West, Southwest, South and Southeast Akron and Barberton City.

STEP 1: Literature Review and Formulation of the Definition

FIGURE 2
Final SOM Solution for 20 Communities in Summit County

- QUESTION SET
1. What is the definition of a complex system?
 2. Where do complex systems exist?
 3. What are the characteristics of complex systems?
 4. What is the theoretical basis for complex systems?
 5. Does the definition of a complex system meet the criteria?

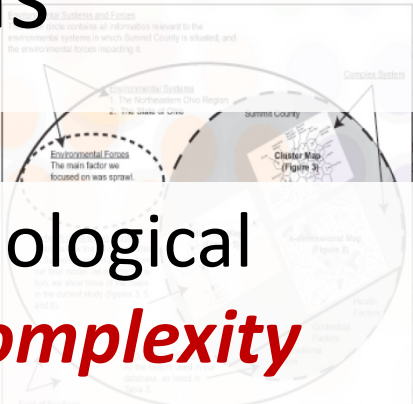
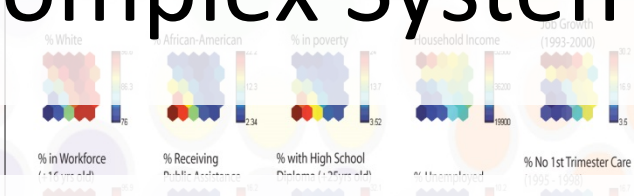


Cases Are Complex Systems

- We also introduce a new methodological framework: the **Sociology and Complexity Science (SACS) Toolkit**.
- **The SACS Toolkit** is a the case-based, mixed-methods, computationally-grounded platform for modeling socio-biological complexity and, more specifically, complex socio-biological systems.

Compositional Factors

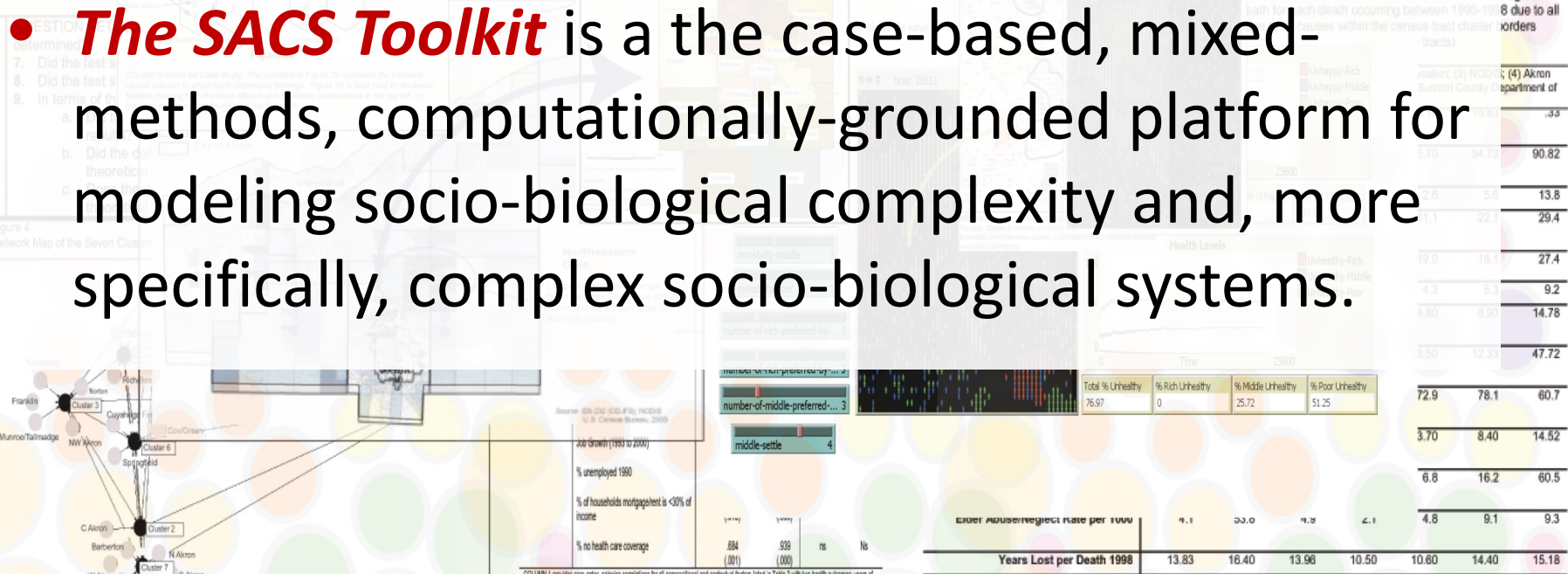
- Population 65 years of age or older
- % White Population (Defined as number of persons identifying themselves as "White" in response to the 1990 US Census or "White Alone" in response to the 2000 US Census)
- % African-American Population (Defined as the number of persons identifying themselves as "Black" or "African-American Alone" in response to the 1990 US Census or "Black or African-American Alone" in response to the 2000 US Census)



STEP 2: Methods

STEP 3: Run Test

STEP 4: Determine Results



NOTE: Distances between clusters are based on Euclidian distances arrived at through k-means analysis. Distances within clusters for each community are based on within-cluster measures. All measures are non-standardized.

COLUMN 1 provides zero-order, pairwise correlations for all compositional and contextual factors listed in Table 3 with two health outcomes, years of life lost per death and Teen Birth Rate. In this column, (r) is the correlation coefficient, and (**) is a two-tailed, significance level.
COLUMN 2 provides the results of our hierarchical analysis of the "independent" relationships all compositional and contextual factors listed in Table 3 with two health outcomes, years of life lost per death and Teen Birth Rate. In this column (s) is a non significant partial correlation coefficient, (**) is a significant partial correlation coefficient for a two-tailed significance level.

1. (**) The values listed in the columns for all 7 clusters represent the average value/measurement that the communities in that cluster scored for each variable listed in Column 1. In cluster analysis, these averages are called the cluster's centroids. 2. Community Membership for each of the 7 Clusters is as follows: Cluster 1: Silesville/Silverlake, Northfield/Lacedonia/Sagamora, and Richfield/Penninsula; Cluster 2: Central Akron; Cluster 3: Taineburg, Northwest Akron, Munroe Falls/Tallmadge, Norton and Franklin; Cluster 4: Hudson; Cluster 5: Copley/Bath/Fairlawn; Cluster 6: Springfield, Coventry/Green and Cuyahoga Falls; Cluster 7: North, West, Southwest, South and Southeast Akron and Barberton City.

SACS Toolkit

Comput Math Organ Theory
DOI 10.1007/s10588-012-9114-1

Case-based modeling and the SACS Toolkit: a mathematical outline

Brian Castellani · Rajeev Rajaram

© Springer Science+Business Media, LLC 2012

Abstract Researchers in the social sciences currently employ a variety of mathematical/computational models for studying complex systems. Despite the diversity of these models, the majority can be grouped into one of three types: agent (rule-based) modeling, dynamical (equation-based) modeling and statistical (aggregate-based) modeling. The purpose of the current paper is to offer a fourth type: case-based modeling. To do so, we review the SACS Toolkit: a new method for quantitatively modeling complex social systems, based on a case-based, computational approach to data analysis. The SACS Toolkit is comprised of three main components: a theoretical blueprint of the major components of a complex system (*social complexity theory*); a set of case-based instructions for modeling complex systems from the ground up (*assemblage*); and a recommended list of case-friendly computational modeling techniques (*case-based toolset*). Developed as a variation on Byrne (in Sage Handbook of Case-Based Methods, pp. 260–268, 2009), the SACS Toolkit models a complex system as a set of k -dimensional vectors (cases), which it compares and contrasts, and then condenses and clusters to create a low-dimensional model (map) of a complex system's structure and dynamics over time/space. The assembled nature of the SACS Toolkit is its primary strength. While grounded in a defined mathematical framework, the SACS Toolkit is methodologically open-ended and therefore adaptable and amenable, allowing researchers to employ and bring together a wide variety of modeling techniques. Researchers can even develop and modify the SACS Toolkit for their own purposes. The other strength of the SACS Toolkit, which makes it a very effective technique for modeling large databases, is its ability to compress data matrices while preserving the most important aspects of a complex system's structure and

B. Castellani (✉)
Dept. of Sociology, Kent State University, Ashtabula, OH 44004, USA
e-mail: bcastel3@kent.edu

R. Rajaram
Dept. of Mathematical Sciences, Kent State University, Ashtabula, OH 44004, USA

SACS Toolkit

1. First, it is comprised of a theoretical blueprint for studying complex systems called it social complexity theory. Social complexity theory is not a substantive theory; instead, it is a theoretical framework comprised of a series of key concepts necessary for modeling complex systems. These concepts include field of relations, network of attracting clusters, environmental forces, negotiated ordering, social practices, and so forth. Together, these concepts provide the vocabulary necessary for modeling a complex system.

2. Second, it is comprised of a set of case-based instructions for modeling complex systems from the ground up called it assemblage. Regardless of the methods or techniques used, assemblage guides researchers through a seven-step process of model building which we review below starting with how to frame ones topic in complex systems terms, moving on to building the initial model, then on to assembling the working model and its various maps to finally ending with the completed model.

3. Third, it is comprised of a recommend list of case-friendly modeling techniques called the *case-based toolset*. The case-based toolset capitalizes on the strengths of a wide list of techniques, using them in service of modeling complex systems as a set of cases. Our own repertoire of techniques include k-means cluster analysis, the self-organizing map neural net, Ragins QCA, network analysis, agent-based modeling, hierarchical regression, factor analysis, grounded theory method, and historical analysis.

SACS Toolkit

We begin our review of the SACS Toolkit with five opening points:

- (1) For the SACS Toolkit, case-based modeling is the study of a complex system S as a set of cases c_i such that:

$$S = \{c_i : c_i \text{ is a case relevant to the system under study}\}. \quad (1)$$

- (2) At minimum, S is comprised of one case c_i .

(3) While there is no theoretical maximum number of cases that can be included in S , practically speaking the upper limit will be bounded, based on the particular set of cases identified for study—which is always an empirical issue.

- (4) We denote the number of cases being studied by n .

- (5) Each case c_i in S is a k dimensional row vector $c_i = [x_{i1}, \dots, x_{ik}]$, where each x_{ij} represents a measurement on one of the variables being used to model a complex system.

SACS Toolkit

TABLE 3
Variables Analyzed for the 20 Communities in the Summit County Database

Compositional Factors	<ul style="list-style-type: none"> Population 65 years of age or older¹ % White Population¹ (Defined as number of persons identifying themselves as "White" in response to the 1990 US Census or "White Alone" in response to the 2000 US Census) % African-American Population¹ (Defined as the number of persons identifying themselves as "Black or African-American" in response to the 1990 US Census or "Black or African-American Alone" in response to the 2000 US Census) Median Household Income¹
Contextual Factors	<ul style="list-style-type: none"> Overall Poverty¹ (Defined as the number of persons living "below the poverty level" as defined by the U.S. Census) Public Assistance¹ (Defined as the number of households receive public assistance as defined by

Place and Health as Complex Systems: A Case Study and Empirical Test

Brian Castellani(1) · Rajeev Rajaram(2) · J Galen Buckwalter(3) · Michael Ball(4) · Frederic Hafferty(5)

the number of jobs in 1990.
 (ce)
 is where mortgage/rent is greater than
 iduals with no health care coverage
 nce Survey – Centers for Disease

irths occurring to mothers from 1995
 during the first three months of the
 :tween 1995-1998 to mothers 15 to

- Childhood Immunization Rate⁶ (Defined as the percentage of children with a complete immunization series 4:3:1 by their second birthday based on the kindergarten retrospective study)
- Child Abuse/Neglect⁶ (Defined as the number of referrals resulting in assessment per 1,000 children under 18 years of age)
- Elder Abuse/Neglect⁷ (Defined as the number of referrals received by the Department of Jobs and Family Services for abuse, exploitation, or neglect)
- Years of Potential Life Lost per Death⁵ (Defined as the sum of the differences between the age at death and the life expectancy at age of death for each death occurring between 1990-1998 due to all causes divided by the number of deaths due to all causes within the census tract cluster borders where those borders are defined by United States Census Bureau census tracts)

Data Sources: (1) United States Census Bureau 1990 and 2000 Decennial Censuses; (2) Ohio Department of Education; (3) NODIS; (4) Akron City Health Department, Office of Epidemiology; (5) Ohio Department of Health; (6) Children's Services Board; (7) Summit County Department of Jobs and Family Service.

SACS Toolkit

Because S consists of n cases $\{c_i\}_{i=1}^n$, and each case c_i has a vector configuration of k dimensions, it is natural to represent S , at least initially and at its most basic, in the form of a data matrix D as follows:

$$D = \begin{bmatrix} c_1 \\ \vdots \\ c_n \end{bmatrix} = \begin{bmatrix} x_{11} & \dots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nk} \end{bmatrix}. \quad (6)$$

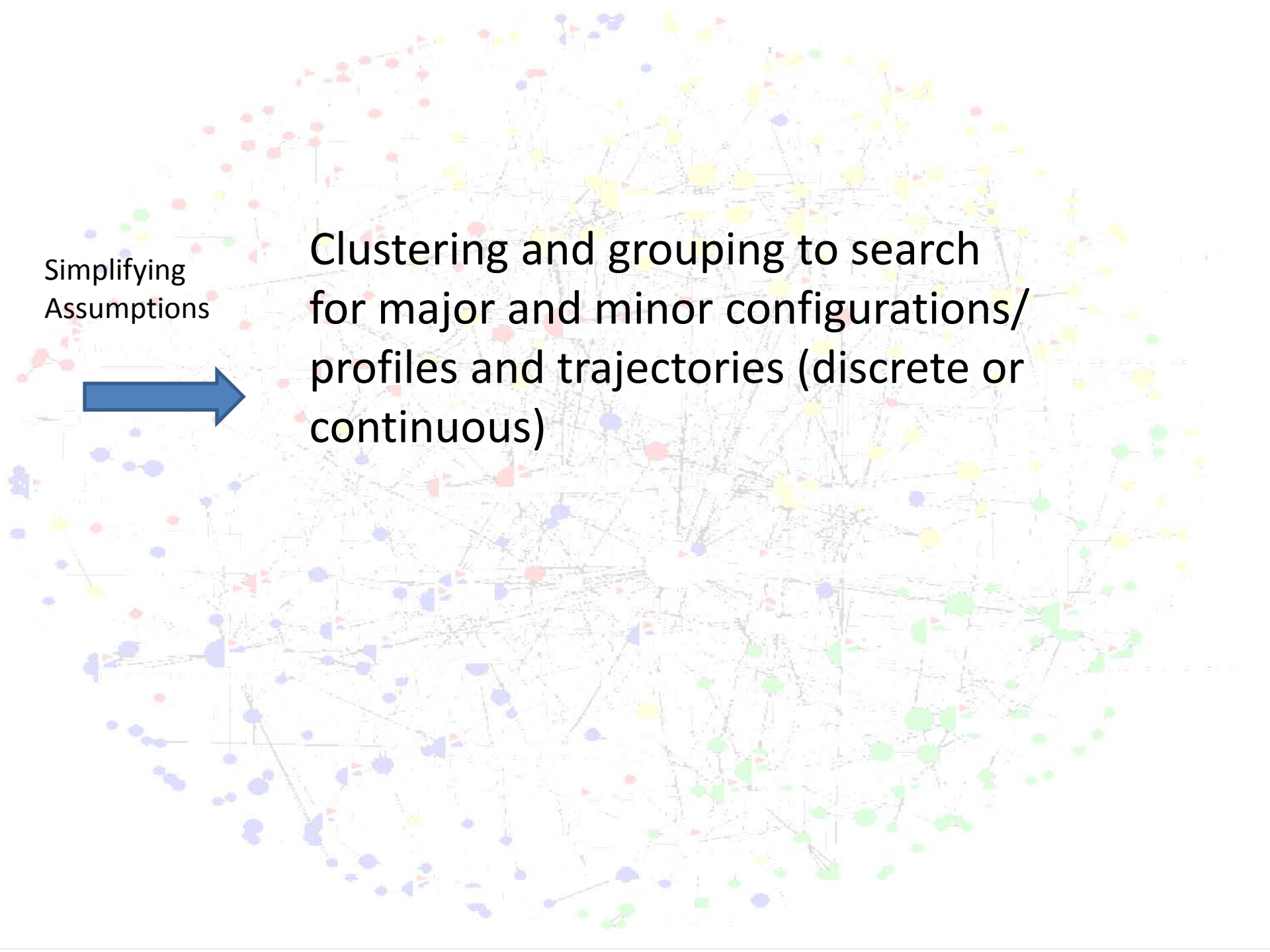
In the notation above, the n rows in D represent the set of cases $\{c_i\}$ in S , and the k columns represent the measurements on some finite partition $\bigcup_{i=1}^p O_i$ of W_S and E_S as defined in Eq. (5) that couple to form the vector configuration for each c_i .



Simplifying
Assumptions



Clustering and grouping to search
for major and minor configurations/
profiles and trajectories (discrete or
continuous)



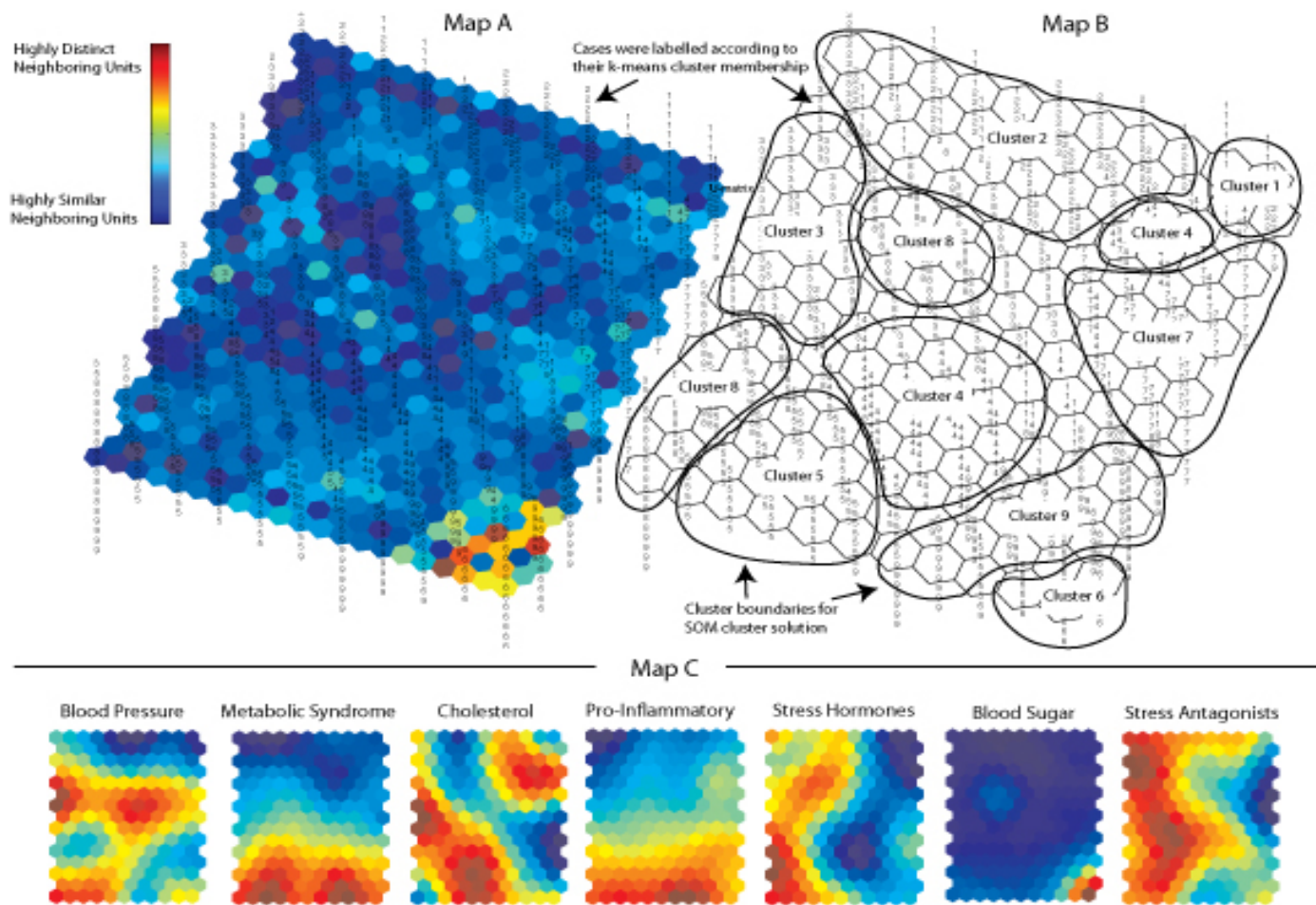


Figure 1.

Map A and *Map B* are graphic representations of the cluster solution arrived at by the Self-Organizing Map (SOM) Neural Net, referred to as the *U-Matrix*. In terms of the information they provide, *Map A* is a three-dimensional (topographical) *u-matrix*: for it, the SOM adds hexagons to the original 15X11 map to allow for visual inspection of the degree of similarity amongst neighboring map units; the dark blue areas indicate neighborhoods of cases that are highly similar; in turn, bright yellow and red areas, as in the lower right corner of the map, indicate highly distinct cluster boundaries. *Map B* is a two-dimensional version of *Map A* that allows for visual inspection of how the SOM clustered the individual cases. Cases on this version of the *u-matrix* (as well as *Map A*) were labelled according to their k-means cluster membership (The 9 cluster solution shown in Table 2) to see if the SOM would arrive at a similar solution. *Map C* is a graphic representation of the relative influence that the seven factors (shown in Table 1) had on the SOM cluster solution. The SOM generates a mini-map for the seven factors, each of which can be overlaid across maps A and B. Each of these mini-maps can then be inspected visually to examine what its rates are across the different neighborhoods (clusters of cases). Dark blue areas indicate the lowest rates for a factor; and the bright red areas indicate the highest rates for a factor. For example, looking at the mini-map for Factor 6 (*Blood Sugar*), its rates are extremely low across most of the map, except for the lower right corner, which is where (looking at *Map A* and *Map B*) the SOM placed Cluster 6.

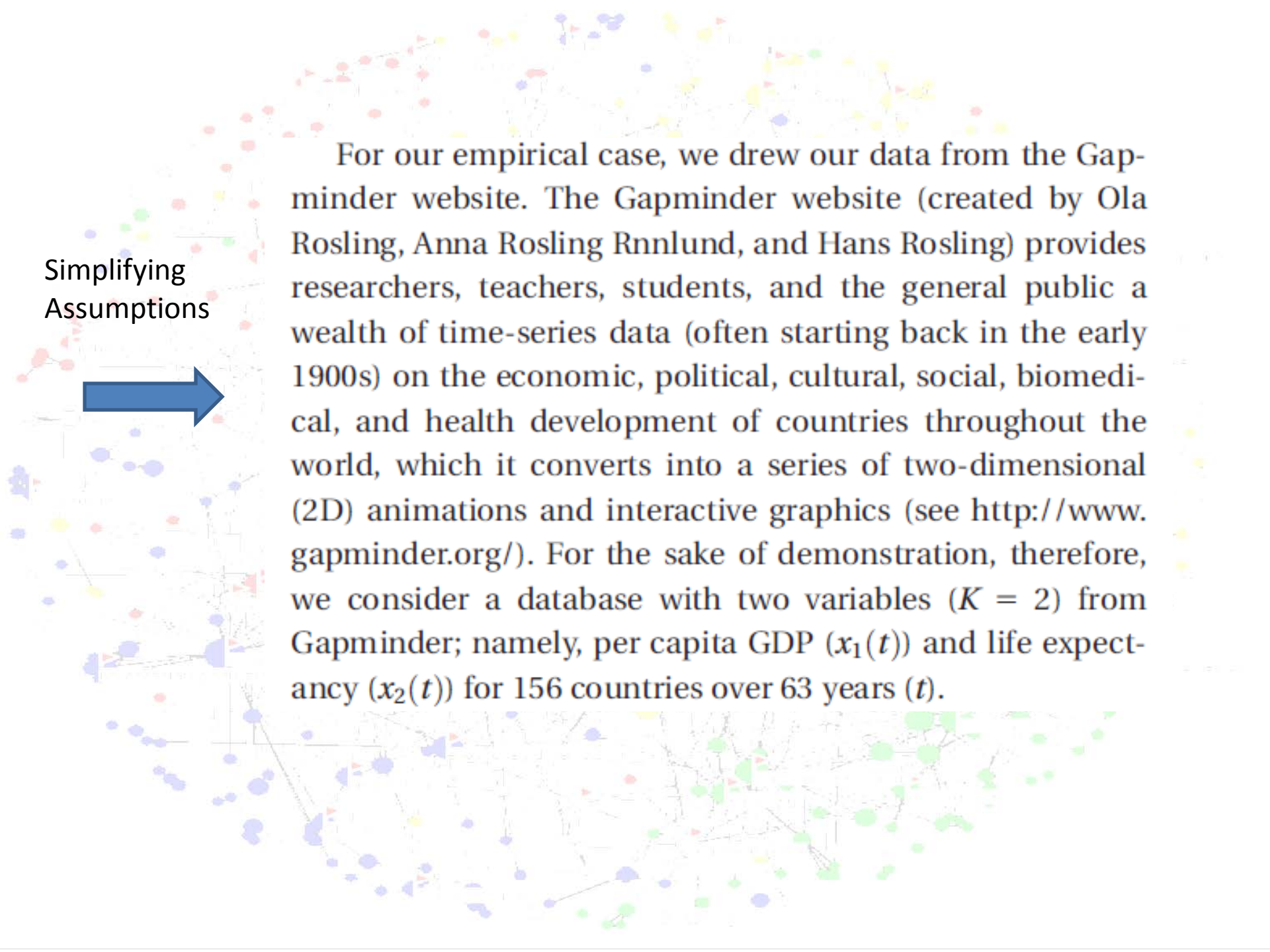


Simplifying
Assumptions

Our approach (which combines what is known in physics and applied mathematics as the inverse and direct problem) is novel in four important ways: first, we take a unique, data-driven view of the cases in a cohort, which we define as K dimensional vectors, where the velocity vector for each case is computed according to its particular measurements on some set of empirically defined social, psychological, or biological variables.

Second, we translate the data-driven, nonlinear trajectories of these microscopic cohort constituents (cases) into the linear movement of macroscopic trajectories, which take the form of densities.

Here, we are drawing on Haken's synergetics and the idea that self-organizing macroscopic trajectories are less dynamic, generally speaking, than microscopic trajectories, which are high dynamic, out of which the former emerge.



For our empirical case, we drew our data from the Gapminder website. The Gapminder website (created by Ola Rosling, Anna Rosling Rnnlund, and Hans Rosling) provides researchers, teachers, students, and the general public a wealth of time-series data (often starting back in the early 1900s) on the economic, political, cultural, social, biomedical, and health development of countries throughout the world, which it converts into a series of two-dimensional (2D) animations and interactive graphics (see <http://www.gapminder.org/>). For the sake of demonstration, therefore, we consider a database with two variables ($K = 2$) from Gapminder; namely, per capita GDP ($x_1(t)$) and life expectancy ($x_2(t)$) for 156 countries over 63 years (t).

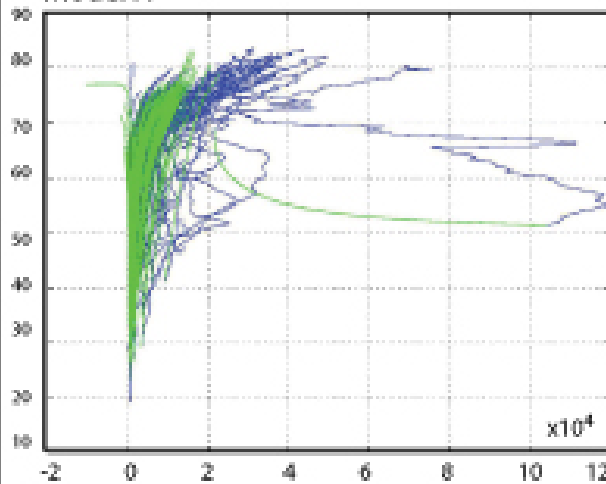
Simplifying
Assumptions



FIGURE 4

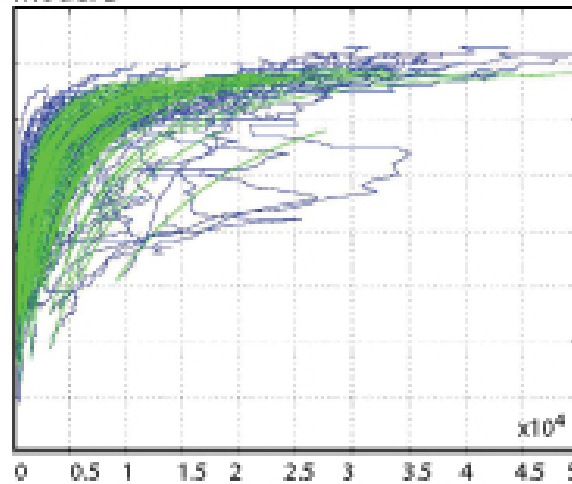
State Space Fit for Best Model

Model A



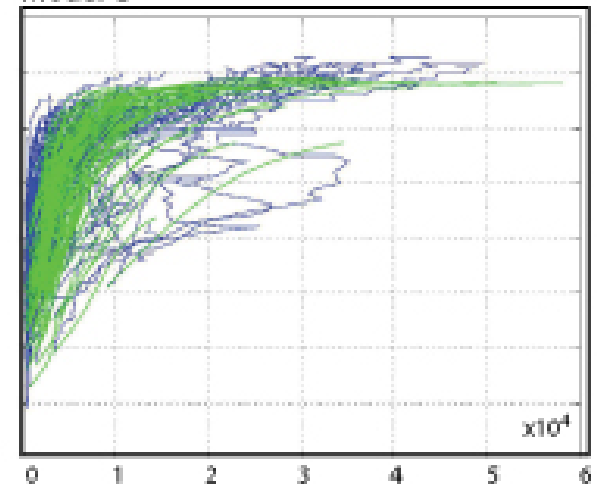
State space fit for the best model.

Model B



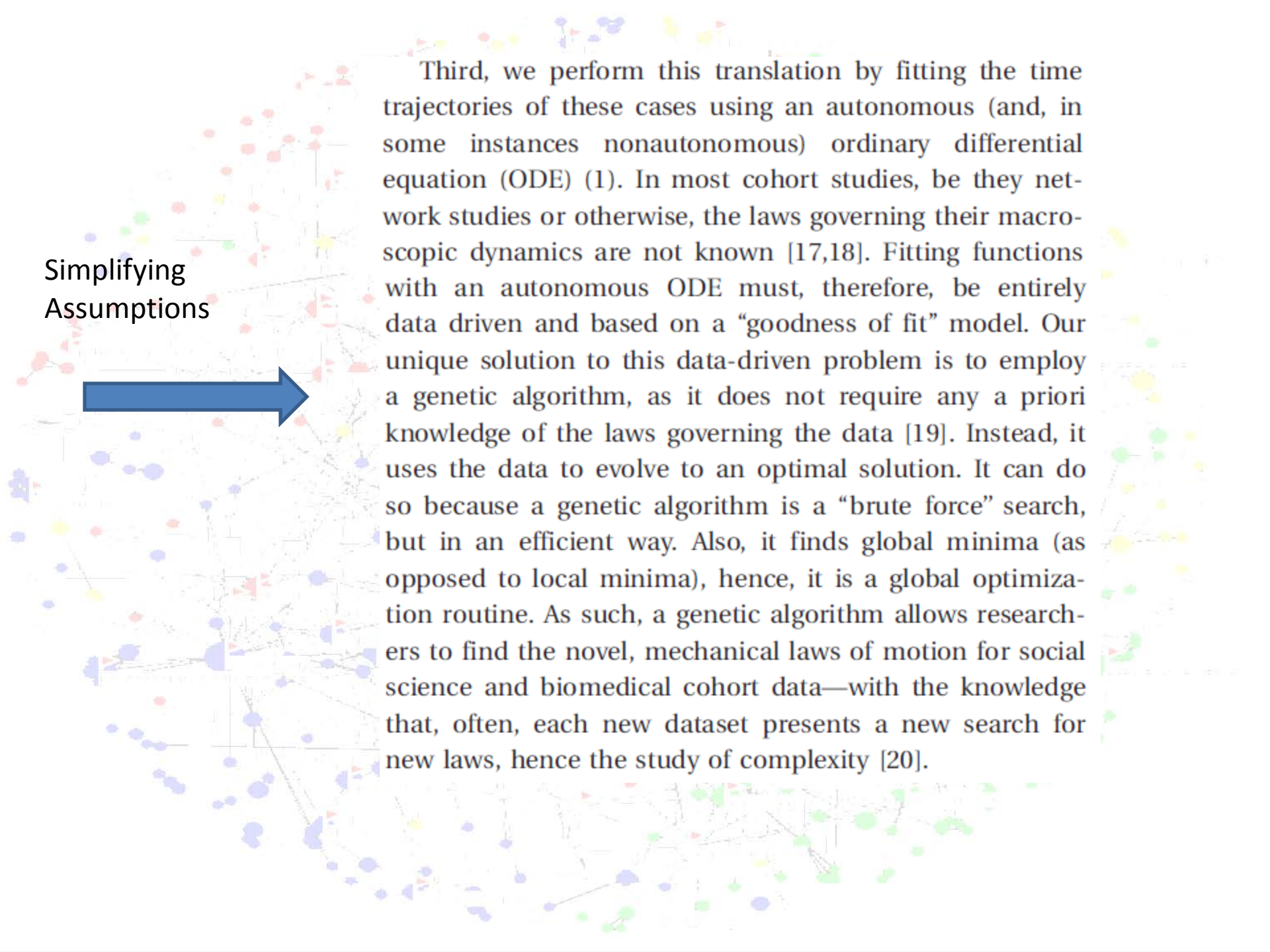
State space fit for the best model without Kuwait or Luxembourg.

Model C



State space fit for the best model without Kuwait or Luxembourg, but with time as an independent variable.

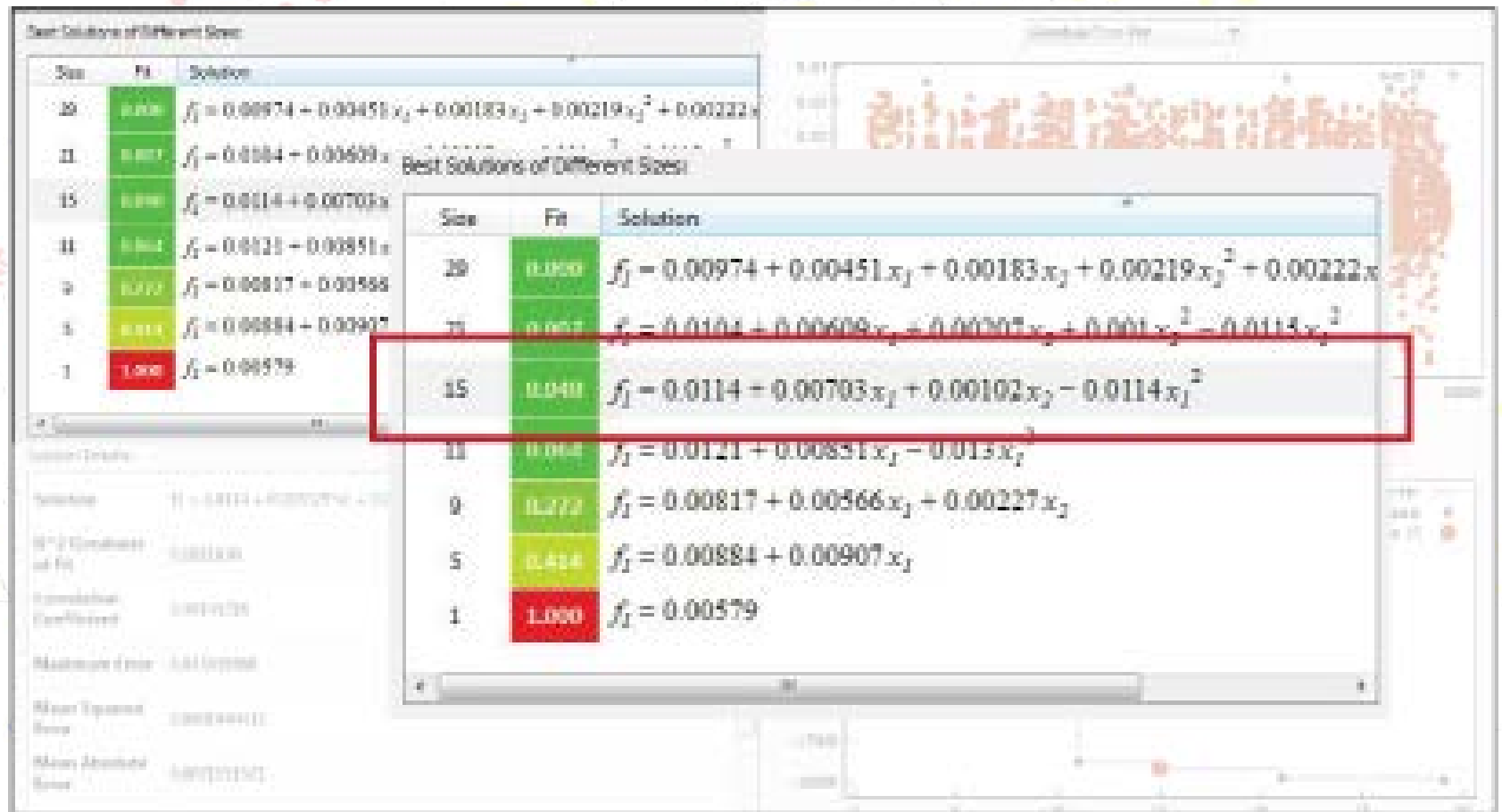
Shown here are several computed Matlab models for the first component of velocity vector f_1 . Models were created using the ordinary differential equation solution from Eureka. In all three models, the X-axis represents **GDP**, and the Y-axis represents **Life Expectancy**. In the models, the blue trajectories are from the data; green trajectories are the fitted model



Simplifying
Assumptions

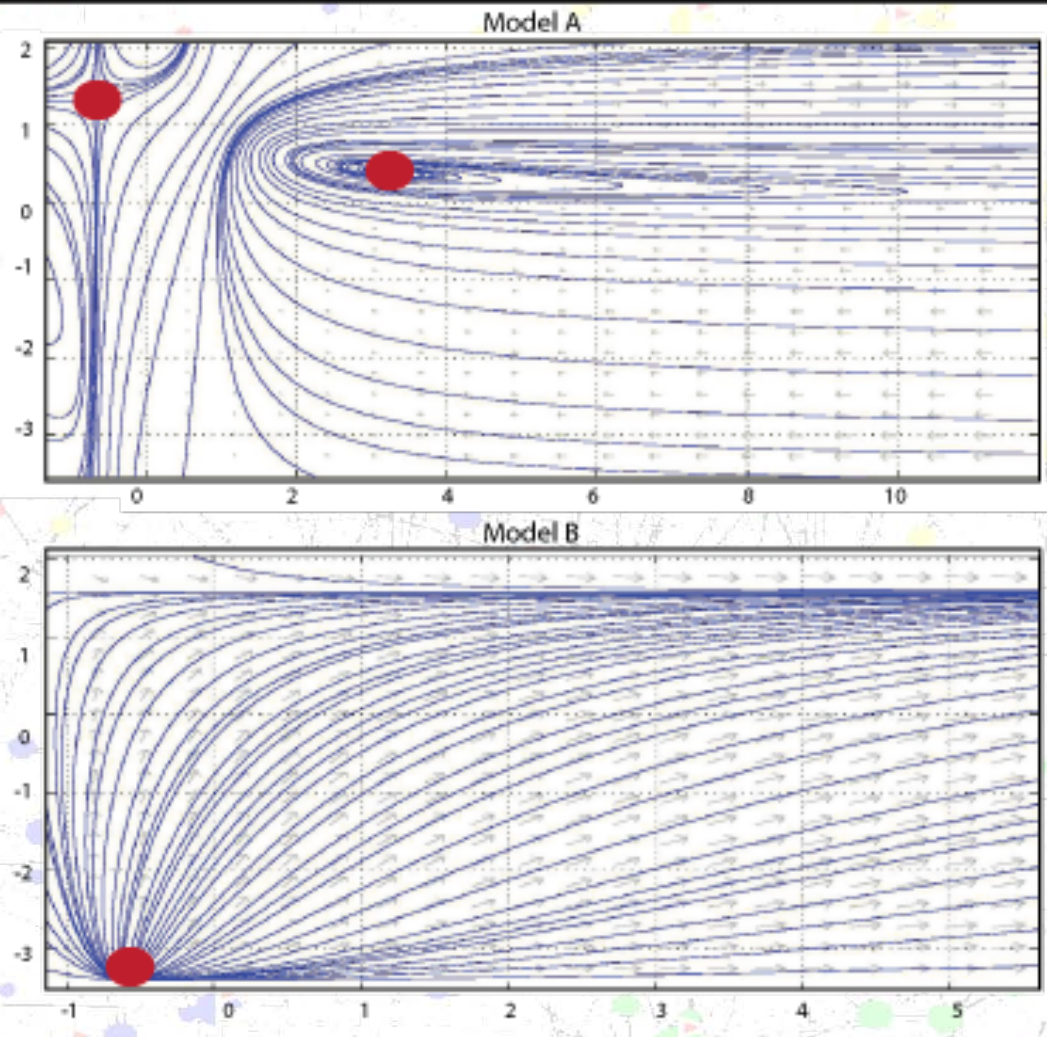
Third, we perform this translation by fitting the time trajectories of these cases using an autonomous (and, in some instances nonautonomous) ordinary differential equation (ODE) (1). In most cohort studies, be they network studies or otherwise, the laws governing their macroscopic dynamics are not known [17,18]. Fitting functions with an autonomous ODE must, therefore, be entirely data driven and based on a “goodness of fit” model. Our unique solution to this data-driven problem is to employ a genetic algorithm, as it does not require any a priori knowledge of the laws governing the data [19]. Instead, it uses the data to evolve to an optimal solution. It can do so because a genetic algorithm is a “brute force” search, but in an efficient way. Also, it finds global minima (as opposed to local minima), hence, it is a global optimization routine. As such, a genetic algorithm allows researchers to find the novel, mechanical laws of motion for social science and biomedical cohort data—with the knowledge that, often, each new dataset presents a new search for new laws, hence the study of complexity [20].

FIGURE 2



*Eureka gives multiple models for the vector field of velocities. Figure 2 shows several computed models for the first component of velocity vector f_1 . The best fit model (#15 in our case, shown above) is usually the one that has a mid-level complexity in terms of number of polynomial symbols and the error values in the mid range amongst all models.

FIGURE 3

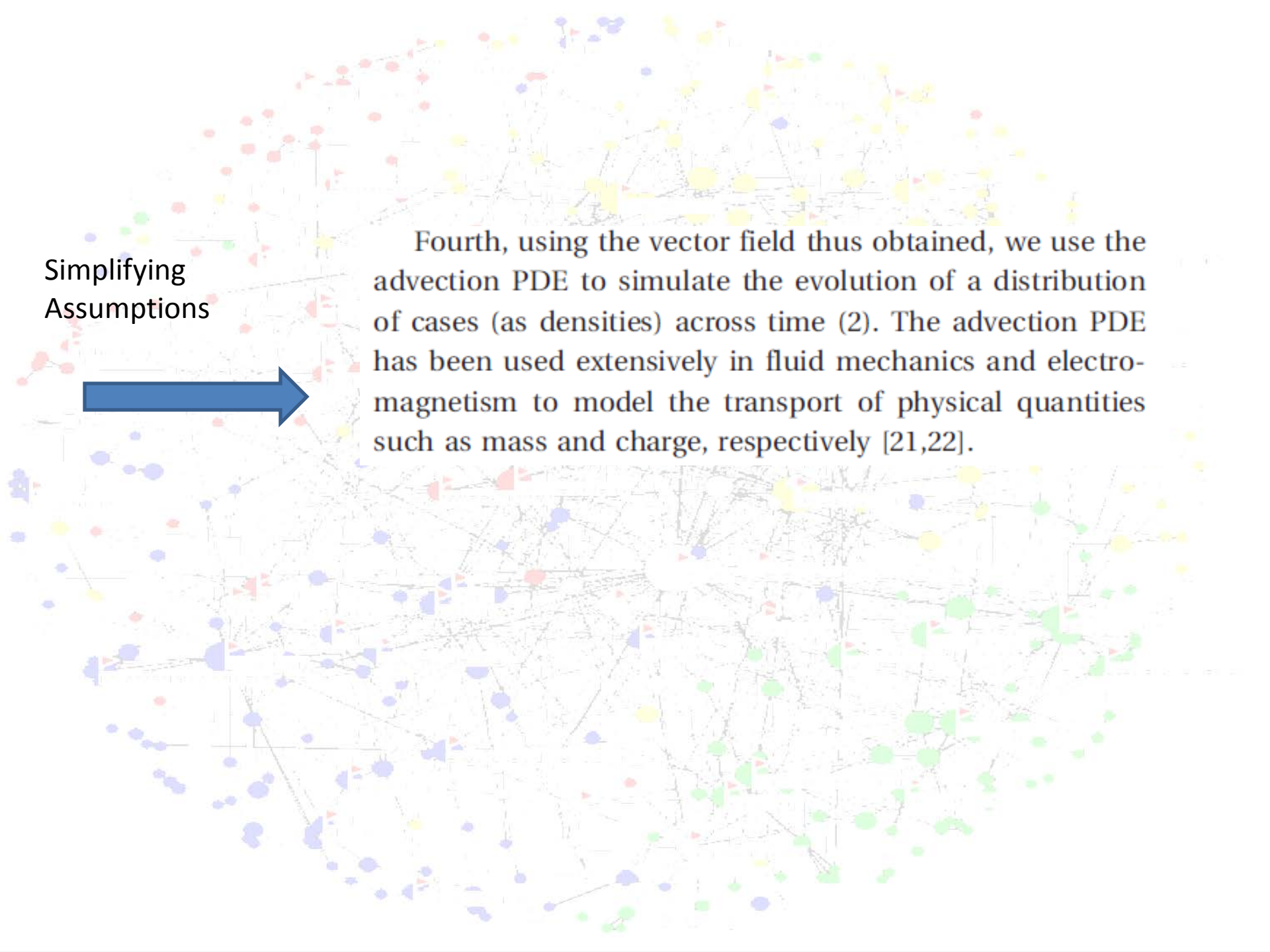


Shown here is the state space trajectories for two of the models we settled on using Eureka. In both models (A and B), the X-axis represents **GDP** (converted to z-scores); and the Y-axis represents **Life Expectancy** (converted to z-scores). In both models, the arrows show the direction of the trajectories; the larger the arrow the higher the vector's velocity. **Model A:** In this model, all countries are included; the red dot located in the top-left section of Figure 3 shows a saddle point; and the red dot located in the top-middle section of Figure 3 shows a spiraling source. **Model B:** In this model, the minority trajectories of Luxembourg and Kuwait were removed; the red dot here is a source.

Simplifying
Assumptions



Fourth, using the vector field thus obtained, we use the advection PDE to simulate the evolution of a distribution of cases (as densities) across time (2). The advection PDE has been used extensively in fluid mechanics and electromagnetism to model the transport of physical quantities such as mass and charge, respectively [21,22].



Advection equation – transport of density of cases

- Transforms the motion of individual cases to the motion of a density of cases.

$$P(t) = \iint_{\Omega} \rho(x, y, t) dx dy = \iint_{\Omega} \rho_0(x, y) dx dy = P_0. \quad (7)$$

- Requires the initial distribution of case profiles, and the velocity vector field of cases (same as the one used in the ODE), and can compute the motion of the initial density assuming that the total number of cases is a constant (called mass conservation property).

$$P_t \rho = \rho(\phi_{-t}(x)) \left| \frac{d\phi_{-t}(x)}{dx} \right|, \quad (8)$$

- Used in modeling of transport phenomena such as fluid dynamics (oil spill), traffic on streets.

$$\rho_t + \nabla \cdot (f \rho) = 0; \quad \rho|_{\Gamma_1} = 0; \quad \rho(x, y, 0) = \rho_0(x, y).$$



AP

Advection equation – transport of density of cases

- Notion of transport is applicable to a variety of topics in sociology such as residential mobility and health trajectories.
- **Residential mobility** – variables are actual geographical ones. Trajectories are in physical coordinate space.
- **Health trajectories** – Variables are biological, sociological markers – state space is more abstract

$$\rho_t + \nabla \cdot (f \rho) = 0; \quad \rho|_{\Gamma_1} = 0; \quad \rho(x, y, 0) = \rho_0(x, y).$$

FIGURE 7

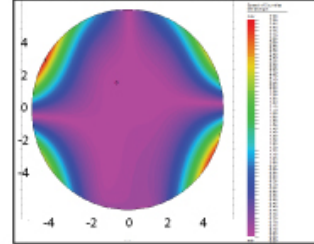
FlexPDE Contour Plots for GDP and Longevity

Model A: Full Model

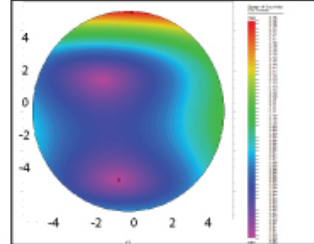
Model B: Without Kuwait or Luxembourg

Model C: Without Kuwait or Luxembourg, but with time as an independent variable.

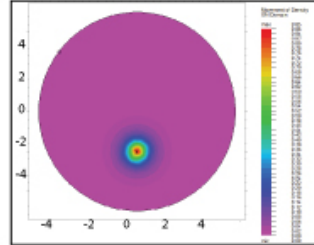
Contour Plot for Speed of Cases



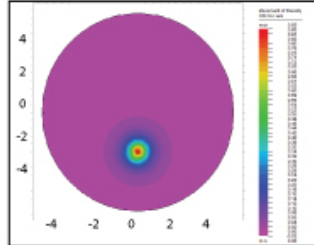
Contour Plot for Speed of Cases



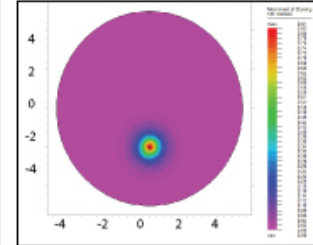
Contour Plot for Distribution of Cases at $t=0$



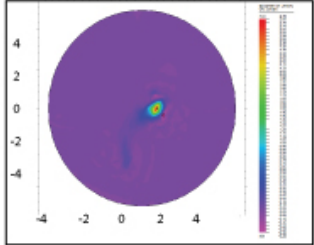
Contour Plot for Distribution of Cases at $t=0$



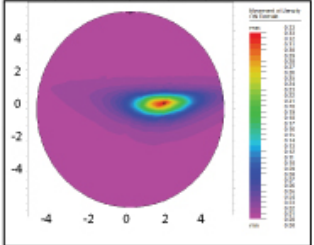
Contour Plot for Distribution of Cases at $t=0$



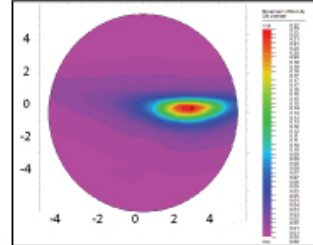
Contour Plot for Distribution of Cases at $t=40$



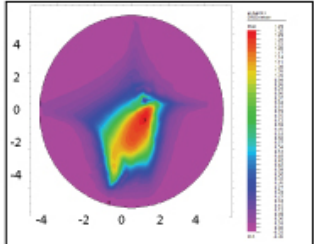
Contour Plot for Distribution of Cases at $t=40$



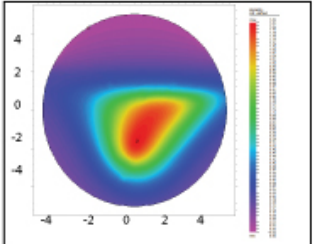
Contour Plot for Distribution of Cases at $t=40$



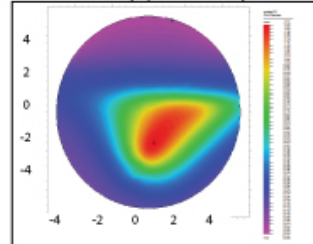
Contour Plot of Lyapunov Density



Contour Plot of Lyapunov Density



Contour Plot of Lyapunov Density



Shown here are several computed models for FlexPDE. Models were created using the advection equation. In all three models, the X-axis represents GDP and the Y-axis represents Life Expectancy; both scores on the axes were converted to z-scores for normalization and comparison. NOTE: In the Lyapunov Density, higher values mean more cases go through that region.

Uniqueness of our approach

$$\dot{x} = f(x); x(0) = x_0$$

$$\rho_t + \nabla \cdot (f\rho) = 0; \rho(x, 0) = \rho_0(x); \rho|_{\Gamma_i} = 0$$

- Continuous time modeling
- Deterministic modeling
- Differential equations (both ODE and PDE)
- Gradation of state space based on velocity of motion
- Non-equilibrium clustering using the Lyapunov density plot

$$\dot{x} = f(x); x(0) = x_0$$

$$f(x, 0) = \rho_0(x); \rho|_{\Gamma_i} = 0$$

Strengths

- Prediction of longitudinal evolution of cases with multiple variables across time
- Studying complexity in dynamical motion of cases in the form of saddles, sources, sinks, or periodic orbits
- Gradation of the state space into regions where cases move faster (or slower) from the velocity contour plot
- Non-equilibrium clustering of trajectories from the Lyapunov density plot (higher values mean more trajectories have squeezed through)

$$\dot{x} = f(x); x(0) = x_0$$
$$\rho(x, t) = \rho_0(x); \rho|_{\Gamma_i} = 0$$

Strengths

- Prediction of majority trends in trajectories for novel choices of initial profiles or densities
- Multiple models to describe the same phenomena allowing for a choice of better ones
- Ease of incorporation of new data into the modeling process to fit the database as it grows



Thanks!