

Data Fusion Reconstruction of Spatially Embedded Complex Networks

Jie Sun,^{1,2,3,4,*} Fernando J. Quevedo,^{1,2,5} and Erik Bollt^{1,2,6}

¹Clarkson Center for Complex Systems Science, Clarkson University, Potsdam, New York, 13699, USA

²Department of Mathematics, Clarkson University, Potsdam, New York, 13699, USA

³Department of Physics, Clarkson University, Potsdam, New York, 13699, USA

⁴Department of Computer Science, Clarkson University, Potsdam, New York, 13699, USA

⁵Department of Mechanical Engineering, Clarkson University, Potsdam, New York, 13699, USA

⁶Department of Electrical & Computer Engineering, Clarkson University, Potsdam, New York, 13699, USA

We introduce a kernel Lasso (kLasso) optimization that simultaneously accounts for spatial regularity and network sparsity to reconstruct spatial complex networks from data. Through a kernel function, the proposed approach exploits spatial embedding distances to penalize overabundance of spatially long-distance connections. Examples of both synthetic and real-world spatial networks show that the proposed method improves significantly upon existing network reconstruction techniques that mainly concerns sparsity but not spatial regularity. Our results highlight the promise of data fusion in the reconstruction of complex networks, by utilizing both microscopic node-level dynamics (e.g., time series data) and macroscopic network-level information (metadata).

PACS numbers: 89.75.Hc, 05.45.Tp, 02.50.Tt

Reconstructing a complex network from observational data is an outstanding problem. Successful network reconstruction can reveal important topological and dynamical features of a complex system and facilitate system design, prediction, and control, as demonstrated in several recent studies across multiple disciplines [1–7]. In many applications, such as material science, infrastructure engineering, neural sensing and processing, and transportation, the underlying complex networks are often spatially embedded (see [8] for an excellent review). The spatial embedding adds yet another dimension to the problem of complex network reconstruction.

A common property of spatially embedded networks is *spatial regularity*, which manifests itself as an inverse dependence of connection probability on spatial distance: generally, the larger the spatial distance is between two nodes, the less likely there exists an edge connecting these nodes [8]. This feature of spatial networks, which can be attributed to physical, financial, or other constraints, has been observed in various types of spatial networks from several different studies, including street patterns [9], mobile communication [10], and social networks [11]. Indeed, the interdependence between network structure and spatial distance is a key ingredient in many widely used models and important studies of spatially embedded networks [12–22].

In this Letter, we show that spatial regularity can be exploited to significantly enhance the accuracy of network reconstruction. In particular, in view of the often limited amount of data available for the inference of large complex networks, the central challenge has always been to better utilize information that potentially arise from distinct sources. To this end, we propose *data fusion reconstruction (DFR)* as a principal framework to infer networks in the presence of both microscopic dynamic data (e.g., time series) and metadata (i.e., spatial embedding information). See Fig. 1. To demonstrate the concept of DFR, we developed kernel Lasso (kLasso) as a generalization of the Lasso, the latter is widely used for sparse regression [23]. Using examples of both synthetic and real-

world spatial networks, we show that due to the integration of sparsity and spatial regularity effects, kLasso reconstructs spatial networks significantly better than Lasso.

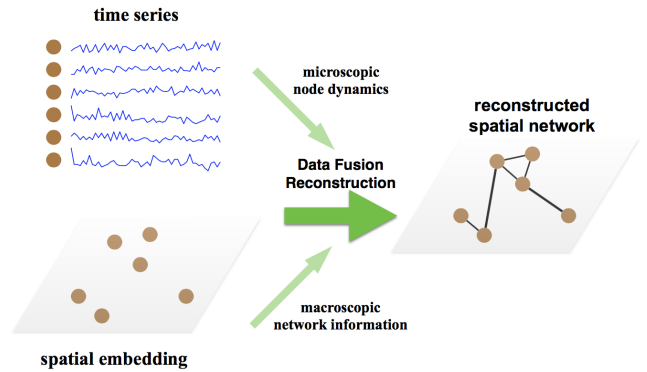


FIG. 1. (color online.) Data fusion reconstruction involves the appropriate “fusion” of information from various sources. For example, such information can arise from both microscopic node-level dynamics (such as the time series of the individual nodes) and macroscopic network-level meta data (e.g., spatial embedding of the nodes).

Mathematically, a spatially embedded network can be represented by a triplet of sets, $G = (V; E; \Phi)$, where $V = \{1, 2, \dots, n\}$ is the set of nodes, $E = \{(i_k, j_k)_{k=1}^m\} \subset V \times V$ is the set of (directed) edges, and $\Phi = [\Phi_1, \Phi_2, \dots, \Phi_n]$ encodes the spatial embedding of the nodes: for example, for a q -dimensional Euclidean embedding, $\Phi_i \in \mathbb{R}^q$.

Problem setup. Given time series data as well as spatial location of the nodes, the problem is to reconstruct the underlying network structure. We represent the time series data as $\{x_j(t)\}_{j=1, \dots, n; t=0, 1, \dots, T}$, where $x_j(t)$ denotes the observed state of node j at the t -th time instance. The sample size T is the number of times each node is observed. In addition, the spatial coordinates of the nodes give rise to an embedded distance $d_{ij} = d(\Phi_i, \Phi_j)$ defined for each pair of nodes (i, j) .

To represent the interactions among the nodes, we employ a standard time series modeling approach [24], by seeking a

(stochastic) linear dependence of the state of each node i at time t on the state of all nodes at time $t - 1$:

$$x_i(t) = \sum_{j=1}^n A_{ij}x_j(t-1) + \xi_i(t), \text{ for } i = 1, 2, \dots, n. \quad (1)$$

The network structure is encoded in the adjacency matrix $A = [A_{ij}]_{n \times n}$, where $A_{ij} \neq 0$ if the state of node i depends on the state of node j , and $A_{ij} = 0$ otherwise. The extra term $\xi_i(t)$ denotes (dynamical) noise. In the case where both $x_i(0)$ and $\xi_i(t)$ are Gaussian, then so is $x_i(t)$, and the vector $\mathbf{x}(t) = [x_1(t), \dots, x_n(t)]^\top$ follows a multivariate Gaussian distribution. Hidden behind the deceptively simple form of Eq. (1) is complexity encoded by the network structure as represented by matrix A , a key factor that enables such standard model to be applicable to broad research topics such as information coding and communication [25], linearization of nonlinear dynamics [26], statistical learning [27], and as a fundamental model form underlying dynamic mode decomposition [28] and Koopman analysis of nonlinear systems [29, 30].

Kernel Lasso. We breakdown the reconstruction problem into the inference of each node's set of neighbors, $\mathcal{N}_i = \{j : A_{ij} \neq 0\}$. Given i , we define $\mathbf{b} = [x_i(1), x_i(2), \dots, x_i(T)]^\top$, $M = [M_{tj}]_{T \times n}$ where $M_{tj} = x_j(t-1)$, and $\mathbf{z} = [A_{i1}, \dots, A_{in}]^\top$. We collect the spatial distance between i and the other nodes into a vector $\mathbf{s} = [s_j]_{n \times 1}$ with $s_j = d_{ij}$. We propose the kernel Lasso (kLasso) optimization problem

$$\min_{\mathbf{z}} (\|\mathbf{M}\mathbf{z} - \mathbf{b}\|_2^2 + \lambda \langle \kappa(\mathbf{s}), |\mathbf{z}| \rangle), \quad (2)$$

where $|\mathbf{z}| = [|z_1|, \dots, |z_n|]^\top$, $\langle \cdot, \cdot \rangle$ denotes inner product, and $\kappa(\mathbf{s}) = [\kappa(s_1), \dots, \kappa(s_n)]^\top$ is obtained by applying a scalar-valued kernel function $\kappa(\cdot)$ to the spatial distances to facilitate preference of spatially short-distance edges over long-distance ones. Finally, the regularization parameter $\lambda \geq 0$ controls the tradeoff between model fit and model regularity. kLasso generalizes the classical Lasso formulation: when the kernel function is a constant, kLasso reduces to Lasso. As we demonstrate later using both synthetic and real-world spatial networks, by explicitly account for spatial embedding information, kLasso generally achieves better reconstruction.

Next we show how to solve kLasso problems. Consider an arbitrary kernel function $\kappa : \mathbb{R} \rightarrow \mathbb{R}^+$. Define matrix $\tilde{M} = [\tilde{M}_{tj}]_{T \times n}$ and vector $\tilde{\mathbf{z}} = [\tilde{z}_1, \dots, \tilde{z}_n]^\top$ as follows:

$$\begin{cases} \tilde{M}_{tj} = M_{tj}/\kappa(s_j), \\ \tilde{z}_j = \kappa(s_j)z_j. \end{cases} \quad (3)$$

Applying these transformations to Eq. (2) converts a kLasso problem into a Lasso problem: $\min_{\tilde{\mathbf{z}}} (\|\tilde{M}\tilde{\mathbf{z}} - \mathbf{b}\|_2^2 + \lambda \|\tilde{\mathbf{z}}\|_1)$, which can be efficiently solved using standard algorithms (such as sequential least squares) found in the literature of computational inverse problems and statistics [27].

Here we focus on a general class of kernel functions of the shifted power-law form

$$\kappa(d) = (d + d_0)^\gamma, \quad (4)$$

where the parameter $d_0 > 0$ ensures that $\kappa(d) > 0$ for all $d \geq 0$ whenever $\gamma \geq 0$. On the other hand, the kernel exponent $\gamma \geq 0$ is used to tune the preference toward spatial regularity: the choice of $\gamma = 0$ recovers the Lasso solution, while the other extreme of $\gamma \rightarrow \infty$ “selects” only the edges that have shortest spatial distance to each node. Intermediate values of γ typically result in a more balanced mix of short-distance edges and long-distance edges appearing in the reconstructed network. Unless otherwise noted, we set the parameters at the default values $\gamma = 1$ and $d_0 = \min_{i \neq j} d_{ij} > 0$.

Synthetic network example: data-enabled inference of random spatial networks. To benchmark the proposed kLasso method, we consider random spatial networks generated by the Waxman model [12]. In particular, for each node pair (i, j) whose spatial distance is d_{ij} , the probability of having an edge between i and j follows

$$P(d_{ij}) = ce^{-\alpha d_{ij}}, \quad (5)$$

where $c > 0$, and $\alpha \geq 0$ (the special case of $\alpha = 0$ produces a classical Erdős-Rényi random network embedded in space) with larger values of α lead to relatively more short-distance edges as compared to long-distance edges. For fixed α , larger values of c generally result in denser networks.

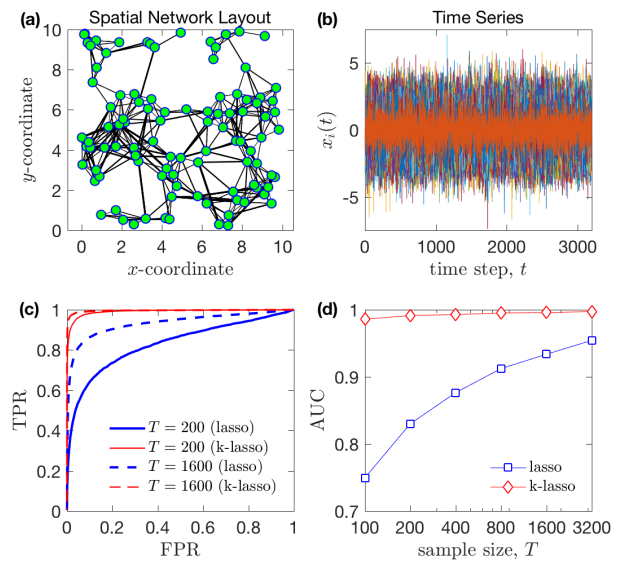


FIG. 2. (color online.) (a) Layout of a random spatial network of $n = 100$ nodes. (b) Typical time series obtained by the stochastic dynamics on the network (see main text for details). (c) Quality of network reconstruction as shown by ROC curves using Lasso versus using kLasso, for two sample sizes. (d) Additional comparison between Lasso and kLasso in network reconstruction based on AUC values across a range of sample sizes. Each data point in (c) and (d) represents an average over 10 independent numerical experiments.

We show an example spatial network in Fig. 2(a). The network contains $n = 100$ nodes that are randomly placed in the 2D spatial domain $(0, \sqrt{n})^2$. The structure of the network is

generated according to Eq. (5) with $\alpha = 2$ and $c = 10$, resulting in a total of $m = 774$ edges. In addition, we generate a self-loop at each node to resemble short-term memory effects. For each edge (i, j) (including the self loops), the weights A_{ij} and A_{ji} are independently drawn from the uniform distribution in $[-1, 1]$. After that, the entire matrix A is scaled $A \rightarrow cA$ with some constant c such that $\rho(A)$, the spectral radius of A , is smaller than 1 to ensure stability of the stochastic process. We select c to yield $\rho(A) = 0.9$. Stochastic time series data is obtained from the network dynamics (1) using iid Gaussian noise $\xi_i(t) \sim \mathcal{N}(0, 1)$. After discarding initial transient, data from T time steps is used for reconstruction. Typical time series of the network is shown in Fig. 2(b).

We compare the results of network reconstruction using kLasso [Eq. (2) with $\gamma = 3$] versus Lasso. To measure the quality of reconstruction, we compute, for each estimate \hat{A} of A , the true positive rate (TPR) and false positive rate (FPR) as

$$\begin{cases} \text{TPR} = |\{(i, j) : \hat{A}_{ij} \neq 0 \ \& \ A_{ij} \neq 0\}| / |\{(i, j) : A_{ij} \neq 0\}|, \\ \text{FPR} = |\{(i, j) : \hat{A}_{ij} \neq 0 \ \& \ A_{ij} = 0\}| / |\{(i, j) : A_{ij} = 0\}|. \end{cases}$$

In Fig. 2(c) we plot the receiver operating characteristic (ROC) curves resulted from kLasso versus Lasso. An ROC curve shows the relationship between TPR and FPR as the regularization parameter λ is varied. Exact, error-free reconstruction corresponds to the upper-left corner of the unit square $[0, 1]^2$ (TPR = 1, FPR = 0), whereas reconstruction by random guesses would yield a diagonal line connecting $(0, 0)$ (empty network) to $(1, 1)$ (complete network). Each ROC curve can be summarized by a scalar defined as the area under the curve (AUC). AUC values are bounded between 0 and 1, with the larger the AUC, generally the closer the ROC curve is to the upper-left error-free corner and the better the reconstruction (AUC value of 1 corresponds to exact reconstruction). As shown in Fig. 2(d) for a wide range of sample sizes, kLasso yields significant improvement over Lasso for network reconstruction. The key reason behind kLasso's success in reconstructing spatial networks lies in its unique capability to incorporate spatial embedding information to "penalize" formation of edges that span over larger spatial distances.

Application: reconstruction from hidden individual dynamics. We now turn to an application of reconstructing a transportation network from observable population-level dynamics data that result from hidden individual trajectories.

The network here is a continent-scale transportation network of Europe, referred to as the E-Road network [31], visualized in Fig. 3(a). A node represents a city of Europe whereas an edge between two nodes represents a highway segment that directly connects the corresponding cities. We compute the embedding distance between each pair of nodes as the shortest distance along the Earth's surface using the corresponding cities' latitude and longitude information.

The dynamical system here describes hidden dynamics on a hidden network, and can be conceptually understood by considering two layers, as illustrated in Fig. 3(b). On the hidden, dynamical layer, there is a total of N individuals, each

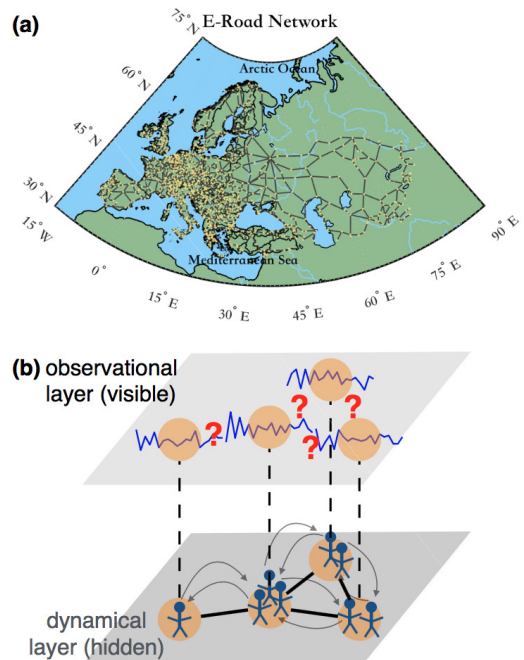


FIG. 3. (color online.) (a) Visualization of the E-Road network, where nodes (marked by light yellow dots) correspond to cities in the Europe and edges (marked by dark gray lines) represent highway segments connecting the cities (water crossings are excluded). This spatial network contains $n = 955$ nodes and $m = 1255$ edges. (b) Two-layer illustration of the hidden dynamics on hidden network, where both the network structure and the dynamics of individuals on the network are hidden (hidden dynamical layer), only the aggregated population dynamics is measured (observational layer).

moving around independently in the spatial network by following a discrete-time random walk [32]. At each time step, an individual at node i moves along one of the edges (i, j) in the network at random to reach node j . On the observational layer, the *aggregated* number of individuals at each node is observed, producing a time series $\{x_i(t)\}$, where $x_i(t)$ is the aggregated number of individual walkers occupying node i at time t . The problem is to reconstruct the hidden spatial network from the observed time series of the population dynamics in the absence of individual trajectories.

Figure 4 shows network reconstruction results using kLasso across a range of kernel exponents γ and sample size T . Excellent reconstruction is generally achieved, with AUC value starts to increase above 0.99 for sample size as low as $T \approx 80$, a number that is surprisingly small compared to the size of the network ($n = 955$ nodes, $m = 1255$ edges). kLasso better reconstructs the network than Lasso (corresponds to $\gamma = 0$) in all parameter combinations, with most significant increase of AUC occurring for $1 \lesssim \gamma \lesssim 3$. For fixed γ , improvement is more significant for smaller T . In addition to further validating the effectiveness of kLasso, the example demonstrates the possibility to reconstruct a hidden spatial network by merely observing aggregated population-level dynamics instead of having to follow detailed individual trajectories.

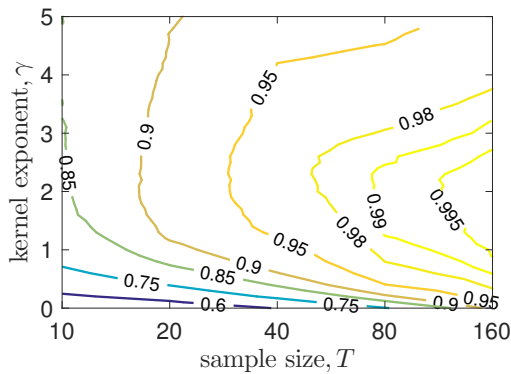


FIG. 4. (color online.) Data fusion reconstruction of the E-Road network using kLasso, shown as a contour plot of AUC values in the (T, γ) plane. Note that for fixed sample size T , Lasso corresponds to choosing $\gamma = 0$, which always yields lower AUC (worse reconstruction) than using kLasso ($\gamma > 0$). Here the optimal choice of γ lies somewhere between 2 and 3.

To summarize, we here developed a kLasso approach for data fusion reconstruction of spatially embedded complex networks. We show that under appropriate linear transformations, kLasso can be converted into a corresponding Lasso problem and thus be efficiently solved using standard Lasso algorithms. We benchmark the effectiveness of kLasso using data from stochastic dynamics on a random spatial network. Furthermore, we consider hidden individual dynamics on E-Road network (a real-world transportation network) where the only observables are the aggregated population dynamics over spatially embedded node locations. kLasso attains excellent reconstruction of the network without the need to fine-tune parameters even for very short time series. These results demonstrate the power of data fusion in the inference in complex systems, in particular the utility of kLasso in the efficient and effective reconstruction of spatially embedded complex networks, when there is both microscopic (e.g., time series data on the nodes) and macroscopic (e.g., metadata of the network) information. Given the flexibility of designing the kernel, it will be interesting to explore other types of metadata for enhanced network reconstruction, such as occupation in social networks. Reconstruction of the E-Road network despite unobservable individual dynamics suggests the possibility of inferring transportation channels from population-level observations without the necessity to trace detailed individual trajectories. This makes kLasso a potentially useful tool for uncovering hidden spatial mobility patterns in practice.

This work was funded in part by the Army Research Office grant W911NF-16-1-0081, the Simons Foundation grant 318812, the Office of Naval Research grant N00014-15-1-2093, and the Clarkson University Space Grant Program an affiliate of the New York NASA Space Grant Consortium.

* sunj@clarkson.edu

- [1] D. Napoletani and T. D. Sauer, Phys. Rev. E **77**, 026103 (2008).
- [2] J. Runge, J. Heitzig, V. Petoukhov, and J. Kurths, Phys. Rev. Lett. **108**, 258701 (2012).
- [3] F. Altarelli, A. Braunstein, L. Dall'Asta, A. Lage-Castellanos, and R. Zecchina, Phys. Rev. Lett. **112**, 118701 (2014).
- [4] N. Antulov-Fantulin, A. Lančić, T. Šmuc, H. Štefančić, and M. Šikić, Phys. Rev. Lett. **114**, 248701 (2015).
- [5] J. Runge, V. Petoukhov, J. F. Donges, J. Hlinka, N. Jajcay, M. Vejmelka, D. Hartman, N. Marwan, M. Paluš, and J. Kurths, Nat. Commun. **6**, 8502 (2015).
- [6] A. S. Ambedgedara, J. Sun, K. Janoyan, and E. Bollt, Chaos **26**, 116312 (2016).
- [7] W. M. Lord, J. Sun, N. T. Ouellette, and E. M. Bollt, IEEE Trans. Mol. Biol. Multi-Scale Commun. **2**, 107 (2017).
- [8] M. Barthélemy, Phys. Rep. **499**, 1 (2011).
- [9] A. P. Masucci, D. Smith, A. Crooks, and M. Batty, Eur. Phys. J. B **71**, 259 (2009).
- [10] R. Lambiotte, V. D. Blondel, C. de Kerchove, E. Huens, C. Prieur, Z. Smoreda, and P. Vandooren, Physica A **387**, 5317 (2008).
- [11] D. Liben-Nowell, J. Nowak, R. Kumar, P. Raghavan, and A. Tomkins, Proc. Natl. Acad. Sci. USA **102**, 11623 (2005).
- [12] B. M. Waxman, IEEE J. Select. Areas. Commun. **6**, 1617 (1988).
- [13] J. Kleinberg, Nature **406**, 845 (2000).
- [14] A. F. Rozenfeld, R. Cohen, D. ben-Avraham, and S. Havlin, Phys. Rev. Lett. **89**, 218701 (2002).
- [15] A. Barrat, M. Barthélemy, and A. Vespignani, J. Stat. Mech. P05003 (2005)
- [16] M. C. González, P. G. Lind, and H. J. Herrmann, Phys. Rev. Lett. **96**, 088702 (2006).
- [17] S. Carmi, S. Carter, J. Sun, and D. ben-Avraham, Phys. Rev. Lett. **102**, 238702 (2009).
- [18] S. Bradde, F. Caccioli, L. Dall'Asta, and G. Bianconi, Phys. Rev. Lett. **104**, 218701 (2010).
- [19] P. Expert, T. S. Evans, V. D. Blondel, and R. Lambiotte, Proc. Natl. Acad. Sci. U.S.A. **108**, 7663 (2011).
- [20] J. A. Henderson and P. A. Robinson, Phys. Rev. Lett. **107**, 018102 (2011).
- [21] G. F. Frasco, J. Sun, H. Rozenfeld, and D. ben-Avraham, Phys. Rev. X **4**, 011008 (2014).
- [22] S. Wickramasinghe, O. Onyerikwu, J. Sun, and D. ben-Avraham, arXiv:1705.07251 (2017).
- [23] R. Tibshirani, J. R. Stat. Soc. B **58**, 267 (1996).
- [24] P. J. Brockwell, *Time Series Analysis: Encyclopedia of Statistics in Behavioral Science*, John Wiley & Sons, Hoboken, NJ, 2005.
- [25] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed., John Wiley & Sons, Hoboken, NJ, 2006.
- [26] A. Lasota and M. C. Mackey, *Chaos, Fractals, and Noise: Stochastic Aspects of Dynamics*, 2nd ed., Springer-Verlag, New York, 1994.
- [27] T. Hastie, R. Tibshirani, and M. Wainwright, *Statistical Learning with Sparsity: The Lasso and Generalizations*, CRC Press, New York, 2015.
- [28] P. J. Schmid, J. Fluid Mech. **656**, 5 (2010).
- [29] I. Mezić, Annu. Rev. Fluid Mech. **45**, 357 (2013).
- [30] M. O. Williams, I. G. Kevrekidis, and C. W. Rowley, J. Nonlinear Sci. **25**, 1307 (2015).
- [31] Data source of the European E-Road Network: <https://eldonk.home.xs4all.nl/eroads/index.htm>
- [32] J. D. Noh and H. Rieger, Phys. Rev. Lett. **92**, 118701 (2004).