



# Randomly distributed embedding making short-term high-dimensional data predictable

Huanfei Ma<sup>a</sup>, Siyang Leng<sup>b,c,d</sup>, Kazuyuki Aihara<sup>b,e,1</sup>, Wei Lin<sup>c,d,f,g,h,1</sup>, and Luonan Chen<sup>ij,k,l,1</sup>

<sup>a</sup>School of Mathematical Sciences, Soochow University, Suzhou 215006, China; <sup>b</sup>Institute of Industrial Science, The University of Tokyo, Tokyo 153-8505, Japan; <sup>c</sup>School of Mathematical Sciences, Fudan University, Shanghai 200433, China; <sup>d</sup>Center for Computational Systems Biology, Institute of Science and Technology for Brain-Inspired Intelligence, Fudan University, Shanghai 200433, China; <sup>e</sup>International Research Center for Neurointelligence, The University of Tokyo Institutes for Advanced Study, The University of Tokyo, Tokyo 113-0033, Japan; <sup>f</sup>Research Institute of Intelligent and Complex Systems, Fudan University, Shanghai 200433, China; <sup>g</sup>Key Laboratory of Mathematics for Nonlinear Sciences (Fudan University), Ministry of Education, Shanghai 200433, China; <sup>h</sup>Key Laboratory of Computational Neuroscience and Brain-Inspired Intelligence (Fudan University), Ministry of Education, Shanghai 200433, China; <sup>i</sup>Key Laboratory of Systems Biology, Center for Excellence in Molecular Cell Science, Shanghai Institute of Biochemistry and Cell Biology, Chinese Academy of Sciences, Shanghai 200031, China; <sup>j</sup>Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences, Kunming 650223, China; <sup>k</sup>School of Life Science and Technology, ShanghaiTech University, Shanghai 200031, China; and <sup>l</sup>Shanghai Research Center for Brain Science and Brain-Inspired Intelligence, Shanghai 201210, China

Edited by Wing Hung Wong, Stanford University, Stanford, CA, and approved September 11, 2018 (received for review February 19, 2018)

Future state prediction for nonlinear dynamical systems is a challenging task, particularly when only a few time series samples for high-dimensional variables are available from real-world systems. In this work, we propose a model-free framework, named randomly distributed embedding (RDE), to achieve accurate future state prediction based on short-term high-dimensional data. Specifically, from the observed data of high-dimensional variables, the RDE framework randomly generates a sufficient number of low-dimensional “nondelay embeddings” and maps each of them to a “delay embedding,” which is constructed from the data of  $a$  to be predicted target variable. Any of these mappings can perform as a low-dimensional weak predictor for future state prediction, and all of such mappings generate a distribution of predicted future states. This distribution actually patches all pieces of association information from various embeddings unbiasedly or biasedly into the whole dynamics of the target variable, which after operated by appropriate estimation strategies, creates a stronger predictor for achieving prediction in a more reliable and robust form. Through applying the RDE framework to data from both representative models and real-world systems, we reveal that a high-dimension feature is no longer an obstacle but a source of information crucial to accurate prediction for short-term data, even under noise deterioration.

memory network (14), and reservoir computing (15–18), have been intensively studied and applied to achieve systems reconstructions and dynamics prediction (19–26). However, based on the neural networks framework (27, 28), the performance of the artificial neural networks crucially and largely relies on the length of the available training data. Thus, these representative methods are effective in accurate prediction only when the training set contains a sufficiently large amount of training data. To handle high-dimensional data, dimension reduction techniques [e.g., various principal component analyses (29, 30), sparse regularization (31–33), and local linearizations] are usually applied for feature extraction. However, the consequence of these applications is likely to overlook interactions (particularly nonlinear interactions) or associations mutually between variables in high-dimensional systems. These interactions in nonlinear dynamics are the crucial information for prediction, remedying the difficulty due to the limited length of observed data, and therefore, the reduction techniques are not always beneficial to accurate prediction of dynamics in complex nonlinear systems (34). Thus, making a good use of the deterministic association or interaction information among the high-dimensional

prediction | nonlinear dynamics | time series | high-dimensional data | short-term data

The big data era has witnessed the accumulation of various types of time series data from microscopic gene expression data through mesoscopic neural activity data to macroscopic ecological or/and atmosphere data (1–5). A challenging task is making accurate forecast or prediction (6, 7) based on such time series datasets, in particular for those datasets with short-term time points but high-dimensional variables. Generally, these two properties are both considered as obstacles for accurate and robust prediction, because short-term datasets always result in fewer statistical patterns for prediction while high-dimensional system variables are likely to bring the curse of dimensionality problem. Specifically, for the model-based methods, such as regression methods (8), or equation-based models (9, 10), taking account of higher-dimensional variables requires a larger number of parameters or weights in the model, making it impractical to estimate these parameters or weights accurately only with short-term data. For the model-free methods, such as the empiricism-based methods where the nearest neighbors in historical data are used to predict the future values (11, 12), short-term data make the depicted attractor sparse in a high-dimensional space, which therefore, yields a problem of the false nearest neighbors. Additionally, machine learning methods, including deep belief network (13), long short-term

## Significance

Making accurate forecast or prediction is a challenging task in the big data era, in particular for those datasets involving high-dimensional variables but short-term time series points, and these datasets are omnipresent in many fields. In this work, a model-free framework, named as “randomly distributed embedding” (RDE), is proposed to accurately predict future dynamics based on such short-term but high-dimensional data. The RDE framework creates the distribution information from the interactions among high-dimensional variables to compensate for the lack of time points in real applications. Instead of roughly predicting a single trial of future values, this framework achieves the accurate prediction by using the distribution information.

Author contributions: H.M., K.A., W.L., and L.C. designed research; H.M., S.L., and L.C. performed research; H.M. and S.L. analyzed data; and H.M., K.A., W.L., and L.C. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

This open access article is distributed under Creative Commons Attribution-NonCommercial-NoDerivatives License 4.0 (CC BY-NC-ND).

<sup>1</sup>To whom correspondence may be addressed. Email: aihara@sat.t.u-tokyo.ac.jp, wlin@fudan.edu.cn, or lchen@sibs.ac.cn.

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.1802987115/-DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1802987115/-DCSupplemental).

Published online October 8, 2018.

variables becomes a pivotal key to designing a useful prediction method (35).

In this work, we propose a model-free framework, named as randomly distributed embedding (RDE), to accurately predict future dynamics based on the observed short-term high-dimensional data. In addition to using the temporal information of each variable, such as the traditional methods usually execute for the long-term data, we exploit the spatial information of the short-term data, such as associations or interactions among the high-dimensional variables. Particularly, the RDE framework can be thought of as an exchange scheme between the spatial information among the observed high-dimensional variables and the time-dependent probability distributions for the temporal dynamics. Thus, it improves the predictability significantly for a target variable. By using the RDE framework to the short-term high-dimensional data produced by both representative models and real-world systems, we show that a high-dimensional feature is no longer an obstacle but a source of information cru-

cial to accurate prediction for short-term data even under noise perturbation.

### RDE Framework

**Delay and Nondelay Embeddings Form Low-Dimensional Attractors.** Usually in a typical high-dimensional nonlinear system, there is a large number of variables interacting with each other; however, the steady dynamics after a transient phase is generally constrained into a low-dimensional subspace due to dissipation. Thus, the state-space technique, based on the embedding theorem, makes it possible to reconstruct a low-dimensional attractor from time series data observed from such a system (36, 37).

As particularly shown in Fig. 1, with the  $n$ -dimensional time series data  $x_i(t), i = 1, 2, \dots, n$ , two kinds of 3D (three-dimensional) attractors can be reconstructed. Specifically, according to the delayed embedding theory (36, 37), one kind is reconstructed in a form of  $\mathcal{M}(x_k(t), x_k(t + \tau), x_k(t + 2\tau))$ ,

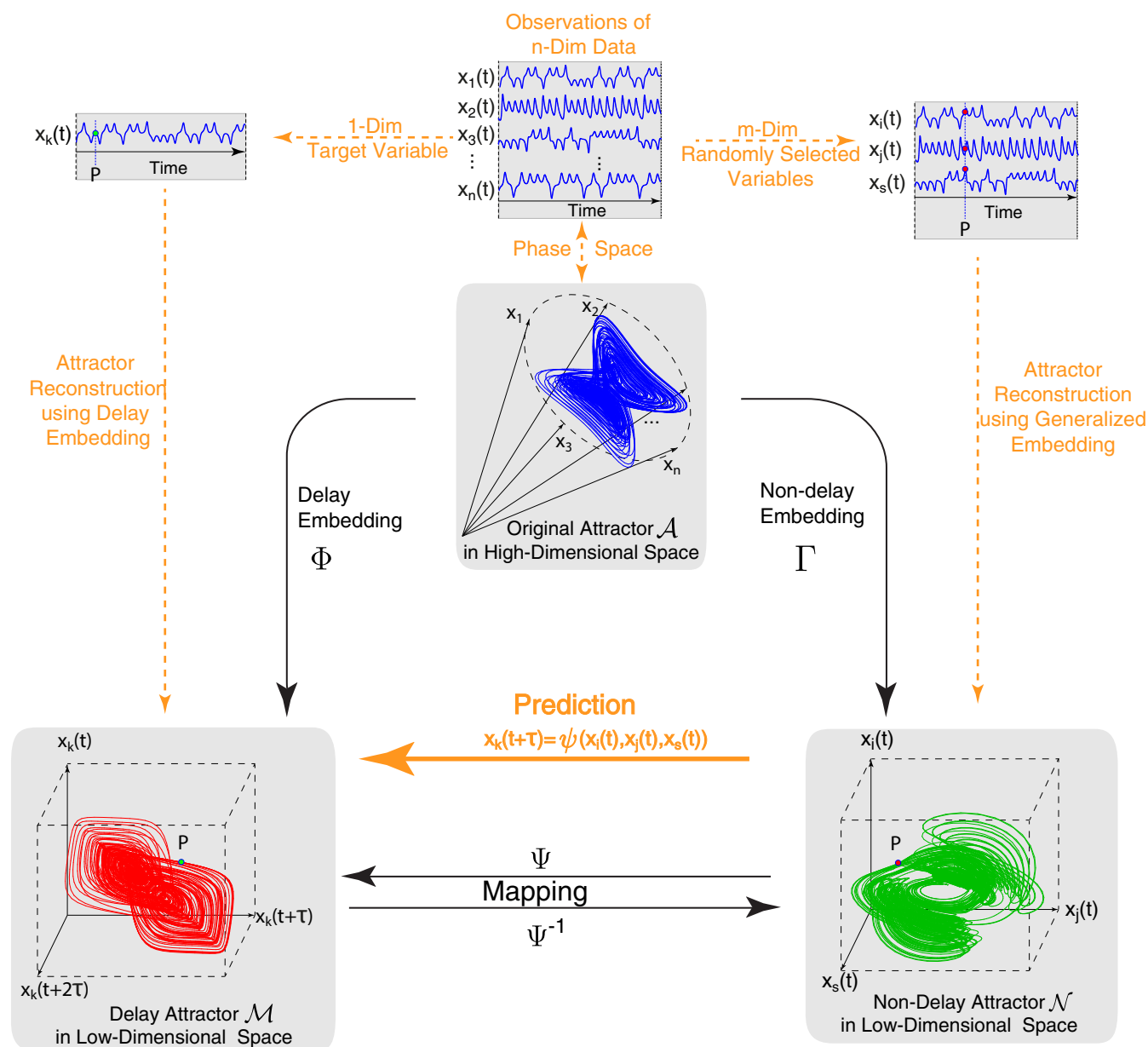


Fig. 1. Sketch of embedding the original attractor in a high-dimensional space into a reconstructed attractor in a low-dimensional space.











can be further alleviated by using parallel computation. For the aggregation scheme, however, we use the in-sample test or the Monte Carlo method with replacement to score candidate random embeddings. In fact, as the number of random embeddings increases, the best in-sample error (or the fitting error) decreases exponentially as shown in *SI Appendix*. Thus, we terminate random embeddings sampling when the in-sample error converges (at the elbow of the exponential decrease), which reduces computational cost and brings good generalization as well.

**Short-Term Data, Robustness, and Comparisons.** Since the RDE framework fully exploits the information embedded in low-dimensional attractors and does not require the coverage of the whole attractor, it is possible to deal with very limited training data. To validate this, we carry out a length test on the coupled Lorenz systems with 15 variables. The test is based on multiple randomly selected sections of measured data. The results are shown in Fig. 7*A* and *B*, where two criteria for one-step predictions are plotted vs. the length of measured data. Compared with other prediction methods for high-dimensional data, the RDE framework particularly works well with very short-term data. Clearly, around 20 time points of the measured data are sufficient for reconstructing system's dynamics. In the literature, both the classic single-variable embedding (SVE) method (11) and the recently proposed multiview embedding (MVE) method (55) can deal with the prediction of high-dimensional data. To make predictions, they both rely on the nearest neighbors in the attractor reconstructed by the historical data, and thus, they may suffer from false nearest neighbors when the length of the time series data is very short. However, the RDE framework does not require that the measured data (training data) cover the whole attractor. It works effectively even when only small segments of the attractor are covered by the measured data as shown in Fig. 5*B*. As clearly shown in Fig. 7*A* and *B*, for the same short-term data (less than 30 points), both methods, MVE and SVE, have poor convergence, while the RDE framework performs well. Indeed, MVE and SVE work well only when the training data become longer (but they are still far from convergence), since longer training data produce better coverage of the nearest neighbors in the attractor.

Noise is inevitable in real applications, and to test the practical robustness of the RDE framework, we also consider the effect of additive white noise in the above 15D coupled Lorenz system with 50 time points as training data. Fig. 7*C* and *D* shows that the RDE framework works well for the signal-to-noise ratio larger than 10, which is as robust as the empirical data-based MVE and SVE methods. Moreover, although both the RDE framework and the RBF (radial basis function) network method proposed in ref. 33 use the inverse embedding technique, the RDE framework fully leverages the information in the distribution of a large amount of random embeddings, while the RBF method uses inverse embedding directly for a high-dimensional system. This difference outstandingly promotes the robustness of the RDE framework against noise deterioration as shown in Fig. 7*C* and *D*.

## Conclusion

In summary, we have established a framework to make predictions from short-term high-dimensional data accurately. The novelty of this RDE framework roots in a full exploitation of the information embedded in a large number of low-dimensional

nondelay attractors as well as in an appropriate use of the exploited distribution of the target variable for prediction. On one hand, the RDE framework creates a distribution, patching all pieces of information from various embeddings into the entire dynamics of the predicted variable. On the other hand, the selection of suitable estimation schemes based on the distribution information thereby significantly increases the prediction reliability and robustness, even for those short-term data with noise deterioration. As validated by datasets produced by both benchmark models and real-world systems, the method is especially effective for the observed short-term high-dimensional time series. This virtue makes the RDE framework potentially useful in mining big datasets from real-world systems.

## Materials and Methods

Given time series data sampled from  $n$  variables of a system with length  $m$  (i.e.,  $\mathbf{x}(t) \in \mathbb{R}^n$ ,  $t = t_1, t_2, \dots, t_m$ , where  $t_i = t_{i-1} + \tau$ ), one can estimate the box-counting dimension  $d$  of the system's dynamics using the false nearest neighbor algorithm (56) and choose embedding dimension  $L > 2d$ . Assume that the target variable to be predicted is represented as  $x_k$ . The RDE algorithm is listed as follows:

- Randomly pick  $s$  tuples from  $(1, 2, \dots, n)$  with replacement, and each tuple contains  $L$  numbers.
- For the  $l$ th tuple  $(l_1, l_2, \dots, l_L)$ , fit a predictor  $\psi_l$  so as to minimize  $\sum_{i=1}^{m-1} \|x_k(t_i + \tau) - \psi_l(x_{l_1}(t_i), x_{l_2}(t_i), \dots, x_{l_L}(t_i))\|$ . Standard fitting algorithms could be adopted. In this paper, Gaussian Process Regression is used.
- Use each predictor  $\psi_l$ , and make one-step prediction  $\tilde{x}_k^l(t^* + \tau) = \psi_l(x_{l_1}(t^*), x_{l_2}(t^*), \dots, x_{l_L}(t^*))$  for a specific future time  $t^* + \tau$ .
- Multiple predicted values form a set  $\{\tilde{x}_k^l(t^* + \tau)\}$ . Exclude the outliers from the set, and use the Kernel Density Estimation method to approximate the probability density function  $p(x)$  of its distribution.
- Calculate the skewness  $\gamma$  of such distribution. In the case  $\gamma < \xi$ , where  $\xi$  is a threshold value, make the final prediction as  $\bar{x}_k(t^* + \tau) = \int xp(x)dx$ . Otherwise, calculate the in-sample prediction error  $\delta_l$  for the fitted  $\psi_l$  using the leave-one-out method. Based on the rank of the in-sample error,  $r$  best tuples are picked out, and the final prediction is given by the aggregated average in the form of  $\bar{x}_k(t^* + \tau) = \sum_{i=1}^r \omega_i \tilde{x}_k^i(t^* + \tau)$ , where the weight  $\omega_i = \frac{\exp(-\delta_i/\delta_1)}{\sum_j \exp(-\delta_j/\delta_1)}$ .

Here, the condition  $\gamma < \xi$  implies that the distribution is nearly symmetric; then, the expectation of the distribution is used as the final prediction. Otherwise, the distribution is asymmetric, indicating that the expectation is not the best choice for the final prediction; then, the aggregation average is used as the final prediction. In this work, we empirically set  $\xi$  as 0.1, and a statistical hypothesis test with shuffling data could be carried out to get a significant level. In this algorithm, the number  $s$  of tuples is determined using a confidence interval or convergence of in-sample errors as given in *SI Appendix*, and the number  $r$  of best tuples is empirically chosen as  $L$ . The RDE algorithm described above is for one-step prediction, but the RDE framework can be extended to multistep prediction. Particularly for the case where  $\psi_l$  is approximated as a linear mapping, the form of  $\psi_l$  can be further explicitly obtained as presented in *SI Appendix*.

**ACKNOWLEDGMENTS.** We thank Qunxi Zhu (Fudan University) for technical support on numerical simulations. We thank the anonymous reviewers for relevant suggestions to improve our work. We also thank the Japan Meteorological Agency, which provided the datasets of wind speeds used in this study (available via the Japan Meteorological Business Support Center). This paper is financially supported by National Key R&D Program of China Grants 2017YFA0505500 and 2018YFC0116600; Strategic Priority Research Program of the Chinese Academy of Sciences Grant XDB13040700; Japan Society for the Promotion of Science, Grants-in-Aid for Scientific Research Grant 15H05707; WPI, Ministry of Education, Culture, Sports, Science and Technology, Japan; National Natural Science Foundation of China Grants 91530320, 11322111, 11771010, and 61773125; and Science and Technology Commission of Shanghai Municipality Grant 18DZ1201000.

1. Lockhart DJ, Winzeler EA (2000) Genomics, gene expression and DNA arrays. *Nature* 405:827–836.
2. De Jong H (2002) Modeling and simulation of genetic regulatory systems: A literature review. *J Comput Biol* 9:67–103.

3. Stein RR, et al. (2013) Ecological modeling from time-series inference: Insight into dynamics and stability of intestinal microbiota. *PLoS Comput Biol* 9:e1003388.
4. Rienecker MM, et al. (2011) Merra: NASA's modern-era retrospective analysis for research and applications. *J Clim* 24:3624–3648.



5. Fan J, Han F, Liu H (2014) Challenges of big data analysis. *Natl Sci Rev* 1:293–314.
6. Clauset A, Larremore DB, Sinatra R (2017) Data-driven predictions in the science of science. *Science* 355:477–480.
7. Subrahmanian V, Kumar S (2017) Predicting human behavior: The next frontiers. *Science* 355:489–489.
8. Hamilton JD (1994) *Time Series Analysis* (Princeton Univ Press, Princeton), Vol 2.
9. Ma H, Lin W (2013) Realization of parameters identification in only locally lipschitzian dynamical systems with multiple types of time delays. *SIAM J Control Optim* 51:3692–3721.
10. Ma H, Lin W (2009) Nonlinear adaptive synchronization rule for identification of a large amount of parameters in dynamical models. *Phys Lett A* 374:161–168.
11. Farmer JD, Sidorowich JJ (1987) Predicting chaotic time series. *Phys Rev Lett* 59:845–848.
12. Wang WX, Lai YC, Grebogi C (2016) Data based identification and prediction of nonlinear and complex dynamical systems. *Phys Rep* 644:1–76.
13. Hinton GE, Osindero S, Teh YW (2006) A fast learning algorithm for deep belief nets. *Neural Comput* 18:1527–1554.
14. Hochreiter S, Schmidhuber J (1997) Long short-term memory. *Neural Comput* 9:1735–1780.
15. Jaeger H (2001) The “echo state” approach to analysing and training recurrent neural networks (German National Research Center for Information Technology GMD, Bonn), Technical Report 148(34):13.
16. Maass W, Natschläger T, Markram H (2002) Real-time computing without stable states: A new framework for neural computation based on perturbations. *Neural Comput* 14:2531–2560.
17. Jaeger H, Haas H (2004) Harnessing nonlinearity: Predicting chaotic systems and saving energy in wireless communication. *Science* 304:78–80.
18. Pathak J, Hunt B, Girvan M, Lu Z, Ott E (2018) Model-free prediction of large spatiotemporally chaotic systems from data: A reservoir computing approach. *Phys Rev Lett* 120:024102.
19. Ma H, Lin W, Lai YC (2013) Detecting unstable periodic orbits in high-dimensional chaotic systems from time series: Reconstruction meeting with adaptation. *Phys Rev E* 87:050901.
20. Kuremoto T, Kimura S, Kobayashi K, Obayashi M (2014) Time series forecasting using a deep belief network with restricted Boltzmann machines. *Neurocomputing* 137:47–56.
21. Xingjian S, et al. (2015) Convolutional lstm network: A machine learning approach for precipitation nowcasting. *Advances in Neural Information Processing Systems*, eds Cortes C, Lawrence ND, Lee DD, Sugiyama M, Garnett R (Curran Associates, Inc., Montreal), pp 802–810.
22. Lu Z, et al. (2017) Reservoir observers: Model-free inference of unmeasured variables in chaotic systems. *Chaos Interdiscip J Nonlinear Sci* 27:041102.
23. Pathak J, Lu Z, Hunt BR, Girvan M, Ott E (2017) Using machine learning to replicate chaotic attractors and calculate Lyapunov exponents from data. *Chaos Interdiscip J Nonlinear Sci* 27:121102.
24. Larger L, et al. (2017) High-speed photonic reservoir computing using a time-delay-based architecture: Million words per second classification. *Phys Rev X* 7:011015.
25. Yeo K, Melnyk I (2018) Deep learning algorithm for data-driven simulation of noisy dynamical system. arXiv:1802.08323.
26. Pathak J, et al. (2018) Hybrid forecasting of chaotic processes: Using machine learning in conjunction with a knowledge-based model. *Chaos Interdiscip J Nonlinear Sci* 28:041101.
27. Haykin S (1994) *Neural Networks: A Comprehensive Foundation* (Macmillan, New York).
28. Dambre J, Verstraeten D, Schrauwen B, Massar S (2012) Information processing capacity of dynamical systems. *Sci Rep* 2:514.
29. Pearson K (1901) On lines and planes of closest fit to systems of points in space. *Lond Edinb Dublin Philos Mag J Sci* 2:559–572.
30. Hotelling H (1933) Analysis of a complex of statistical variables into principal components. *J Educ Psychol* 24:417–441.
31. Tibshirani R (1996) Regression shrinkage and selection via the lasso. *J R Stat Soc Ser B* 73:273–282.
32. Candes EJ, Romberg JK, Tao T (2006) Stable signal recovery from incomplete and inaccurate measurements. *Commun Pure Appl Math* 59:1207–1223.
33. Ma H, Zhou T, Aihara K, Chen L (2014) Predicting time series from short-term high-dimensional data. *Int J Bifurcation Chaos* 24:1430033.
34. Van Der Maaten L, Postma E, Van den Herik J (2009) Dimensionality reduction: A comparative. *J Mach Learn Res* 10:66–71.
35. Kantz H, Schreiber T (2004) *Nonlinear Time Series Analysis* (Cambridge Univ Press, Cambridge, UK), Vol 7.
36. Takens F (1981) Detecting strange attractors in turbulence. *Dynamical Systems and Turbulence, Warwick 1980*, eds Rand DA, Young L-S (Springer, Berlin), pp 366–381.
37. Sauer T, Yorke JA, Casdagli M (1991) Embedology. *J Stat Phys* 65:579–616.
38. Packard NH, Crutchfield JP, Farmer JD, Shaw RS (1980) Geometry from a time series. *Phys Rev Lett* 45:712–716.
39. Deyle ER, Sugihara G (2011) Generalized theorems for nonlinear state space reconstruction. *PLoS One* 6:e18295.
40. Rasmussen C, Williams C (2006) *Gaussian Processes for Machine Learning* (MIT Press, Cambridge, MA).
41. Ho TK (1998) The random subspace method for constructing decision forests. *IEEE Trans Pattern Anal Mach Intell* 20:832–844.
42. Bryll R, Gutierrez-Osuna R, Quek F (2003) Attribute bagging: Improving accuracy of classifier ensembles by using random feature subsets. *Pattern Recognit* 36:1291–1302.
43. Colizza V, Pastor-Satorras R, Vespignani A (2007) Reaction-diffusion processes and metapopulation models in heterogeneous networks. *Nat Phys* 3:276–282.
44. Kondo S, Miura T (2010) Reaction-diffusion model as a framework for understanding biological pattern formation. *Science* 329:1616–1620.
45. Dress A, Hordijk W, Lin W, Serocka P (2010) The ideal storage cellular automaton model. *Structure Discovery in Biology: Motifs, Networks & Phylogenies*, Dagstuhl Seminar Proceedings, eds Apostolico A, Dress A, Parida L (Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, Dagstuhl, Germany), pp 1862–4405.
46. Dress AW, Lin W (2011) Dynamics of a discrete-time model of an “ideal-storage” system describing hetero-catalytic processes on metal surfaces. *Int J Bifurcation Chaos* 21:1331–1339.
47. Na YJ, et al. (2009) Comprehensive analysis of microRNA-mRNA co-expression in circadian rhythm. *Exp Mol Med* 41:638–647.
48. Shen-Orr SS, Milo R, Mangan S, Alon U (2002) Network motifs in the transcriptional regulation network of Escherichia coli. *Nat Genet* 31:64–68.
49. Hirata Y, Aihara K (2016) Predicting ramps by integrating different sorts of information. *Eur Phys J Spec Top* 225:513–525.
50. Wong TW, et al. (1999) Air pollution and hospital admissions for respiratory and cardiovascular diseases in Hong Kong. *Occup Environ Med* 56:679–683.
51. Fan J, Zhang W (1999) Statistical estimation in varying coefficient models. *Ann Stat* 27:1491–1518.
52. Xia Y, Härdle W (2006) Semi-parametric estimation of partially linear single-index models. *J Multivar Anal* 97:1162–1184.
53. Kish L (1995) *Survey Sampling* (John Wiley & Sons, New York).
54. Hogg RV, Craig AT (1995) *Introduction to Mathematical Statistics* (Prentice Hall, Upper Saddle River, NJ), 5th Ed.
55. Ye H, Sugihara G (2016) Information leverage in interconnected ecosystems: Overcoming the curse of dimensionality. *Science* 353:922–925.
56. Kennel MB, Brown R, Abarbanel HDI (1992) Determining embedding dimension for phase-space reconstruction using a geometrical construction. *Phys Rev A* 45:3403–3411.