



School and Workshop on
Topological Data Analysis



Data Types

Methods for Reconstruction of Data Sets of Different Types

Hosein Masoomy

Computational Cosmology Group (CCG) [ccg.sbu.ac.ir]

&

Center for Complex Networks and Social Data Science (CCNSD) [ccnsd.ir]

@

Department of Physics, Shahid Beheshti University (SBU) [sbu.ac.ir]

August 24, 2022



Overview

Data Types

Time Series

Field

Point Cloud

Network (Graph)

Methods for Reconstruction of Data Sets of Different Types

Time Delay Embedding (TDE)

Recurrent Plot (RP)

Visibility Graph (VG)

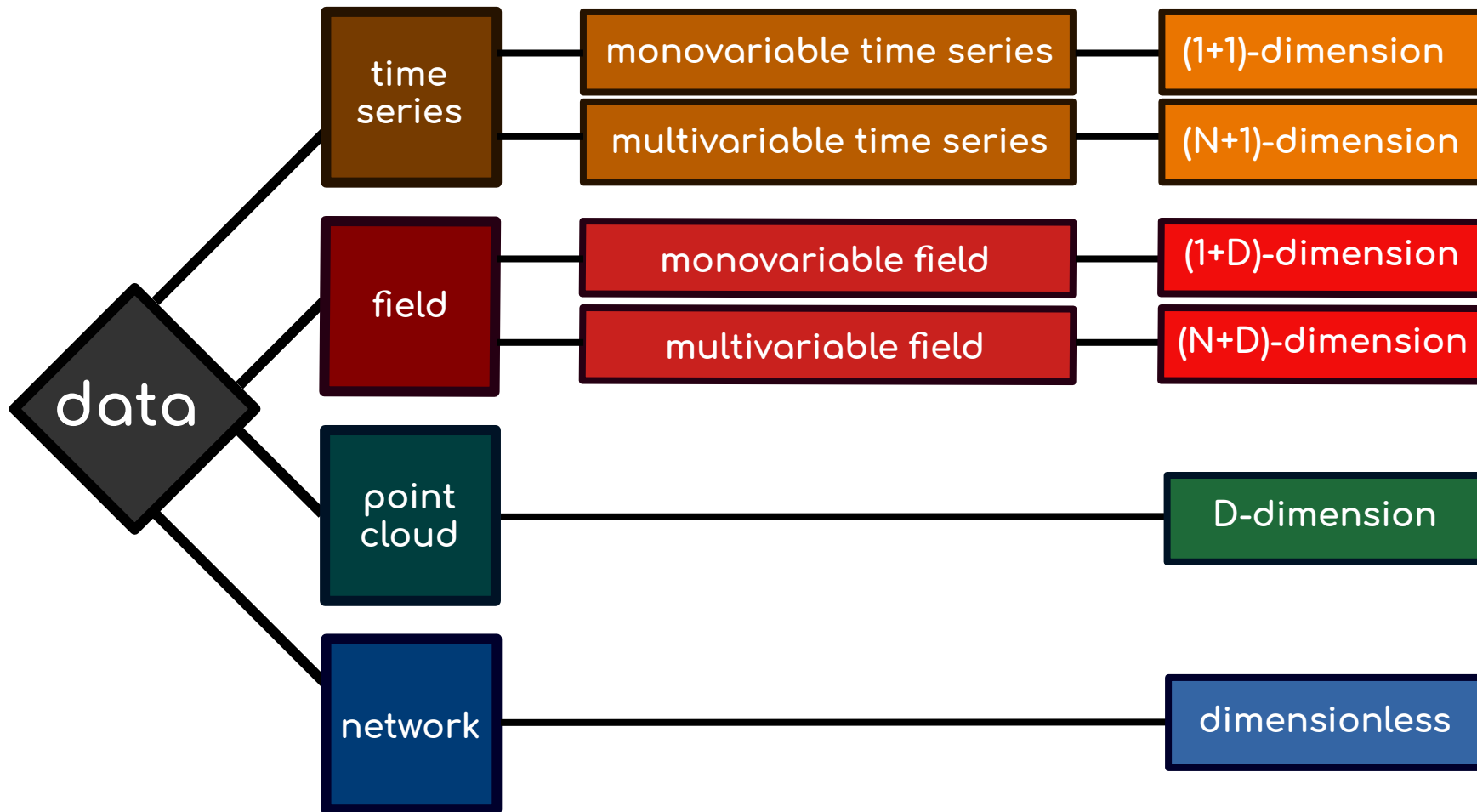
State Space (SS)

Correlation Network (CN)

Recurrent Network (RN)

Excursion Sets (ES)

Data Types



Data Types / Time Series



Data Types / Time Series

time series

monovariable time series

(1+1)-dimension

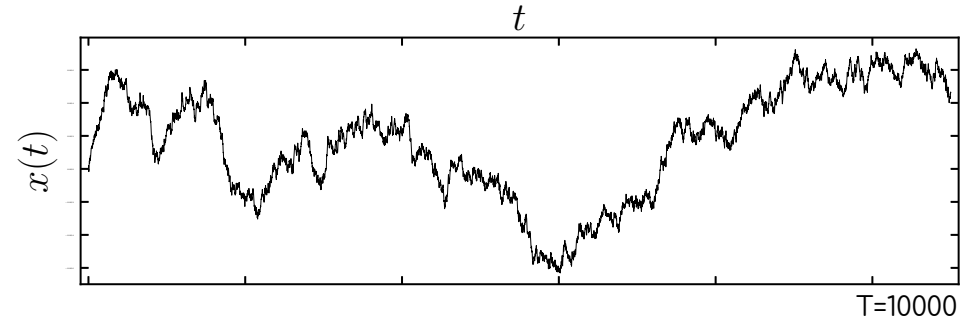
T = length of time series

$\mathcal{T} \subseteq \mathbb{R}^+ \equiv [0, +\infty)$ domain

$X \subseteq \mathbb{R}$ range

$x : \mathcal{T} \rightarrow X$ map

$\vec{x} \equiv \left(x(t_i) \right)_{i=1}^T$



Data Types / Time Series

time series

monovariable time series

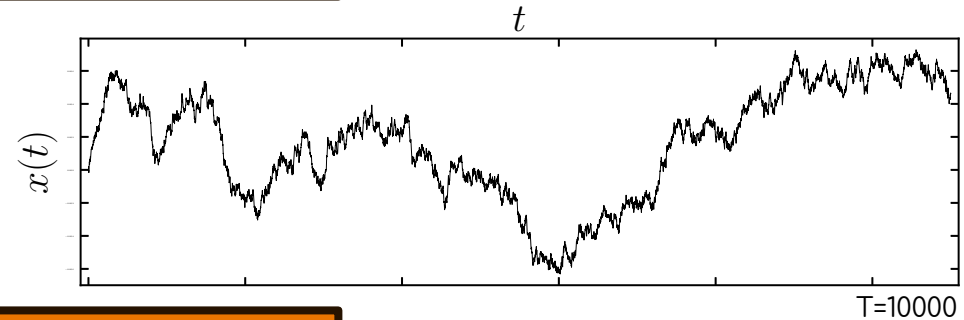
(1+1)-dimension

T = length of time series

$\mathcal{T} \subseteq \mathbb{R}^+ \equiv [0, +\infty)$ domain

$X \subseteq \mathbb{R}$ range $\vec{x} \equiv \left(x(t_i) \right)_{i=1}^T$

$x : \mathcal{T} \rightarrow X$ map



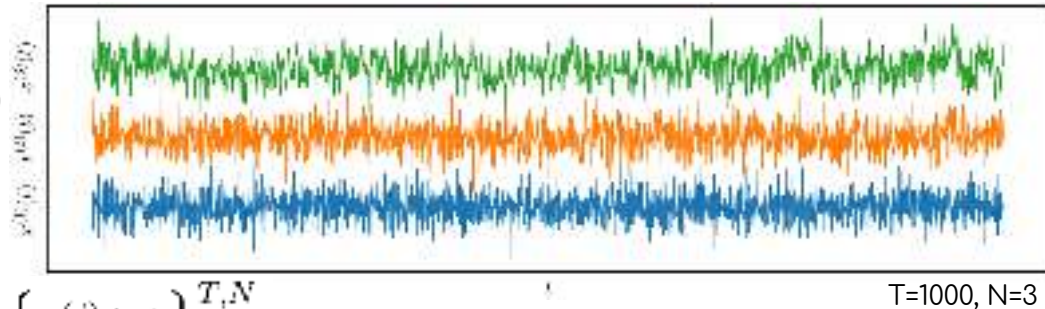
multivariable time series

(N+1)-dimension

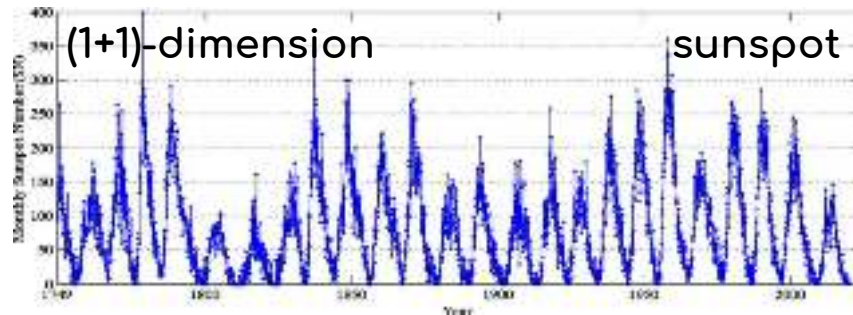
N = number of degrees of freedom
(number of dependent variables)

To express the evolution of
the system, we need N maps:

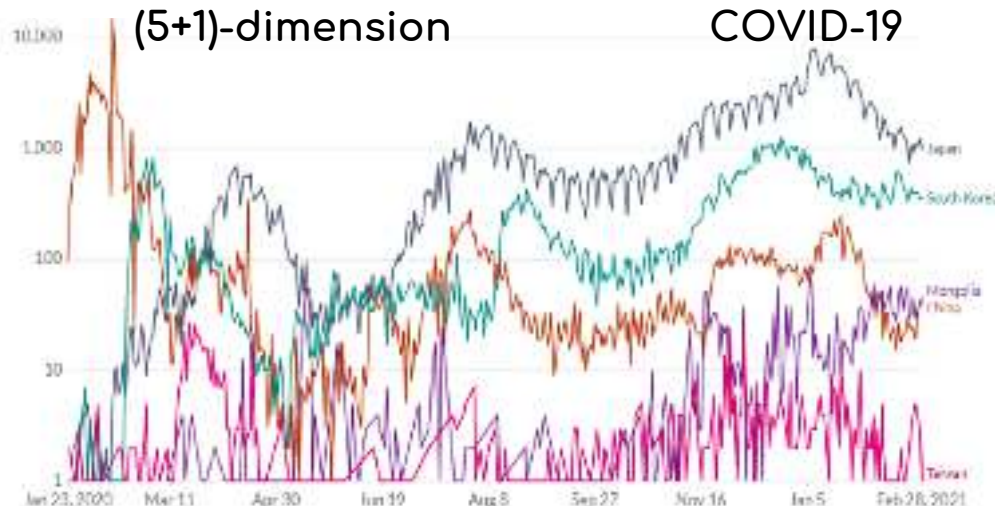
$$\mathcal{X} = \left\{ \vec{x}^{(j)} \mid \vec{x}^{(j)} = \left(x^{(j)}(t_i) \right)_{i=1}^T \right\}_{j=1}^N \equiv \left\{ x^{(j)}(t_i) \right\}_{i,j=1}^{T,N}$$



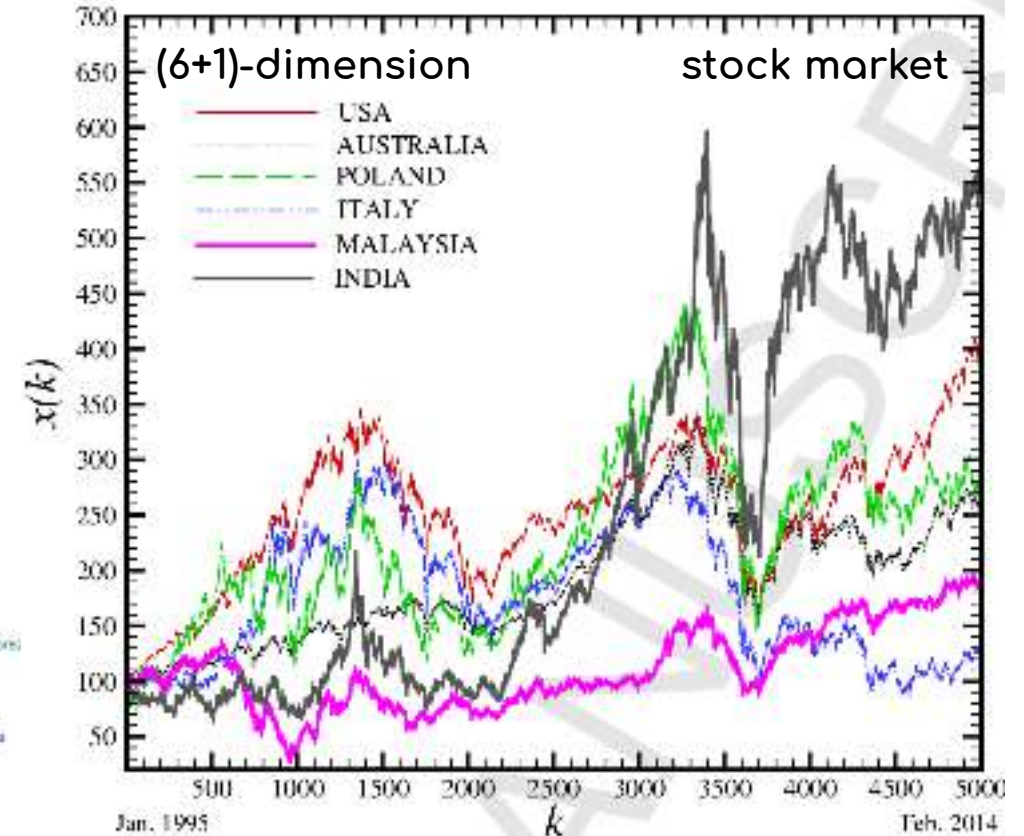
Data Types / Time Series : real data



Panigrahi, Sibarama, Radha Mohan Pattanayak, Prabira Kumar Sethy, and Santi Kumari Behera. "Forecasting of sunspot time series using a hybridization of ARIMA, ETS and SVM methods." *Solar Physics* 296, no. 1 (2021): 1-19.



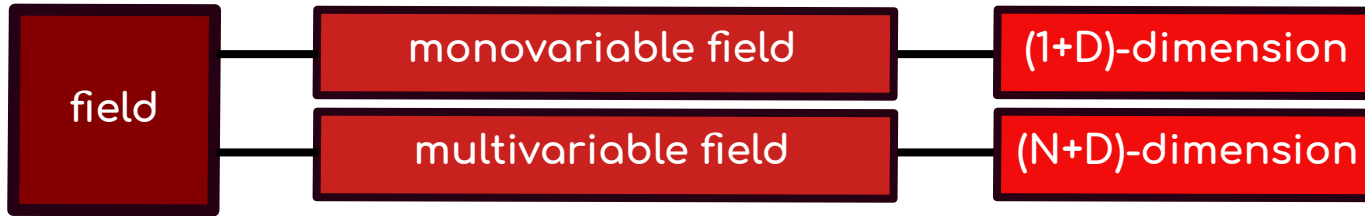
Source: Johns Hopkins University CSSE COVID-19 Data. Last updated 3 March, 2021 (London, UK)



Ferreira, Paulo, Andreia Dionísio, and S. M. S. Movahed. "Assessment of 48 stock markets using adaptive multifractal approach." *Physica A: Statistical Mechanics and its Applications* 486 (2017): 730-750.

CC BY

Data Types / Field



Data Types / Field

field

monovariable field

(1+D)-dimension

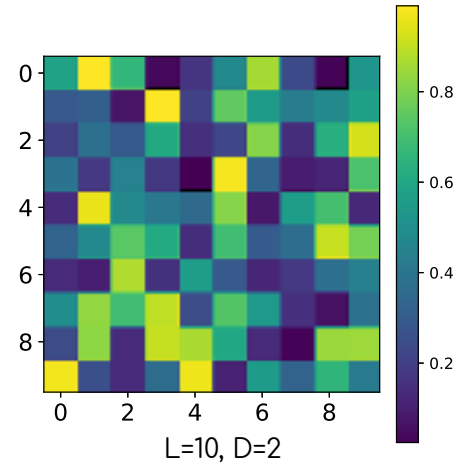
$$\mathcal{F} : \Pi \rightarrow \mathbb{R}$$

L = length of field

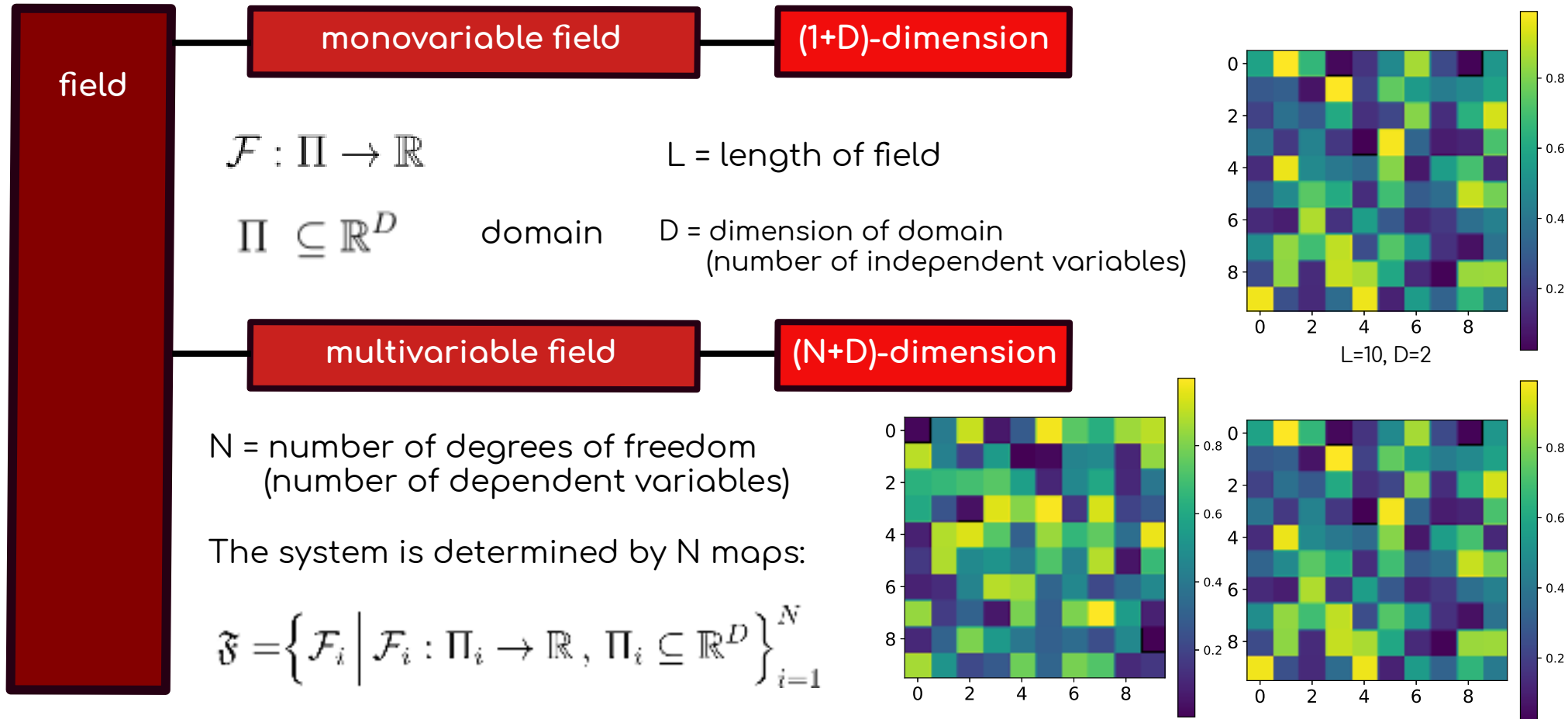
$$\Pi \subseteq \mathbb{R}^D$$

domain

D = dimension of domain
(number of independent variables)

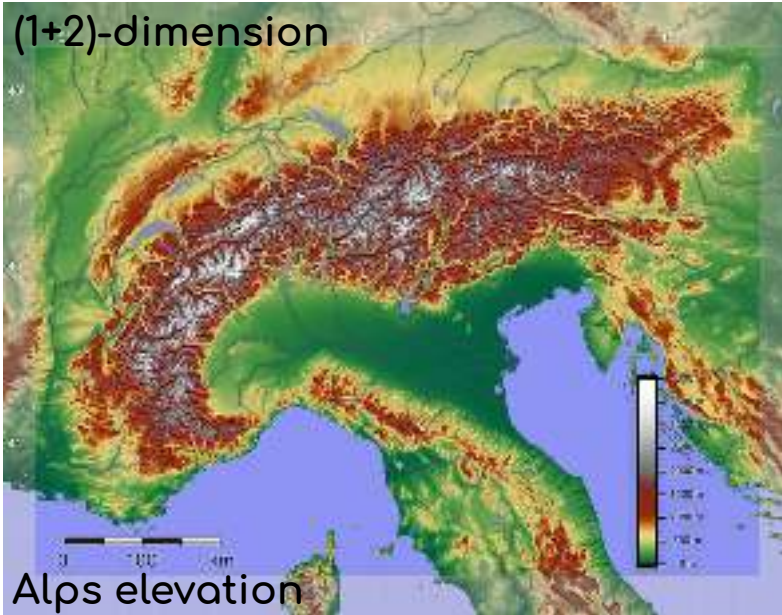


Data Types / Field



Data Types / Field : real data

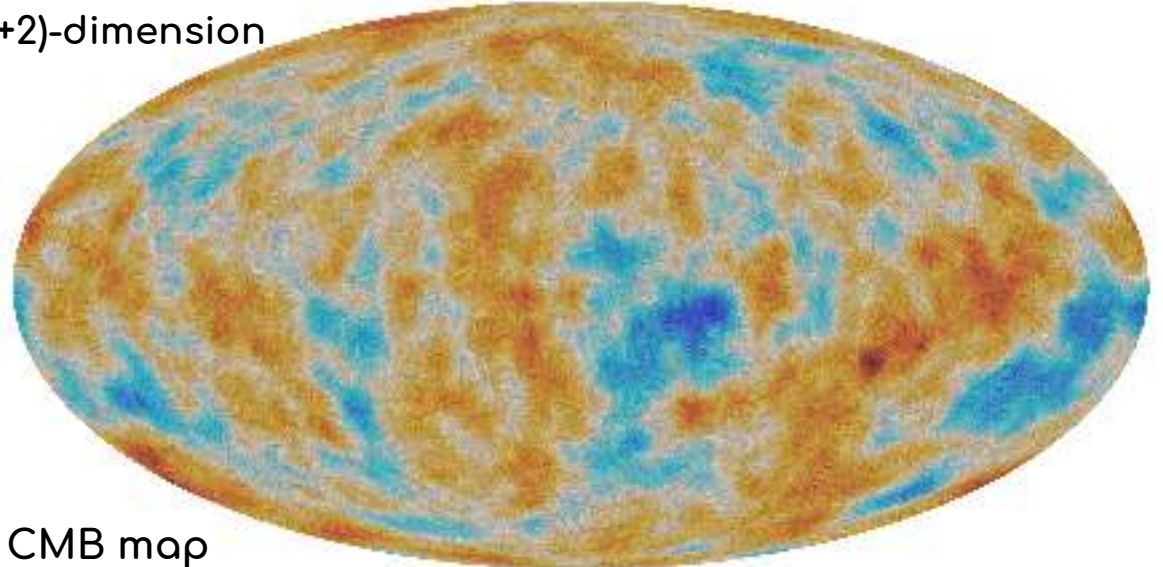
(1+2)-dimension



Alps elevation

<https://en.wikipedia.org/wiki/Alps>

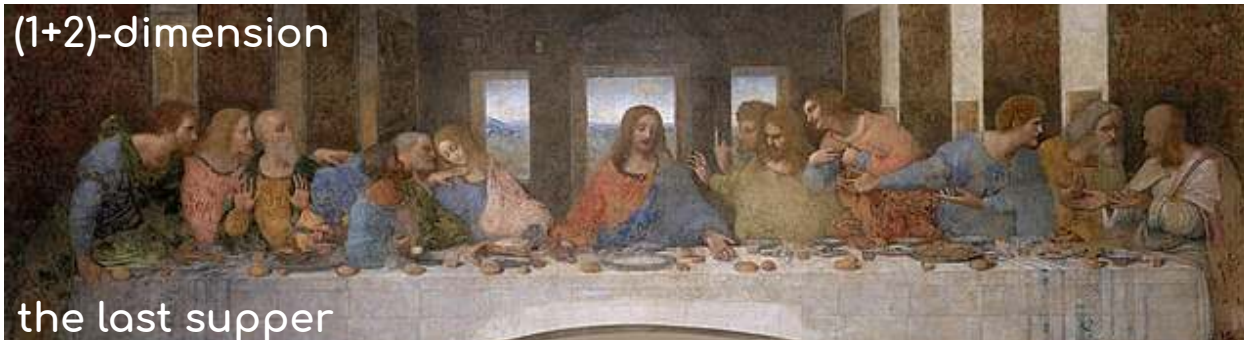
(1+2)-dimension



CMB map

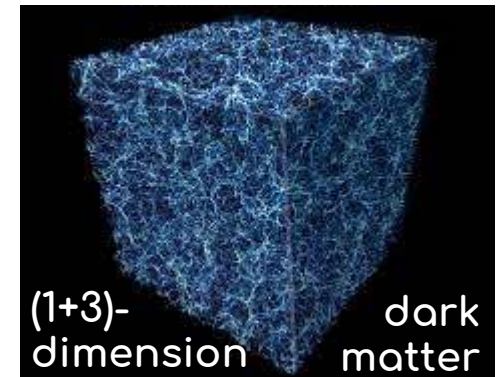
https://www.esa.int/ESA_Multimedia/Images/2015/02/Polarisation_of_the_Cosmic_Microwave_Background

(1+2)-dimension



the last supper

https://en.wikipedia.org/wiki/The_Last_Supper_%28Leonardo%29



(1+3)-dimension dark matter

<https://futurism.com/20-percent-universes-normal-matter-exists-dark-cosmic-voids>

Data Types / Point Cloud

point
cloud

D-dimension

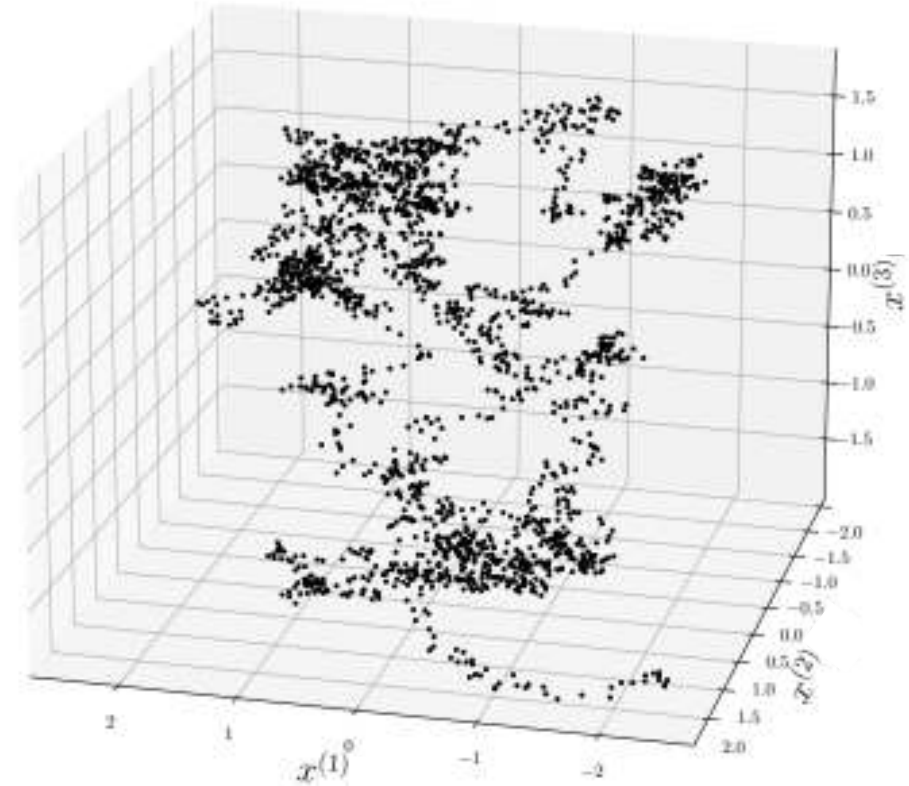
$$\mathbb{X} = \left\{ x_i \mid x_i \equiv (x_i^{(d)})_{d=1}^D, x_i^{(d)} \in \mathbb{R} \right\}_{i=1}^{N \neq \infty}$$

x_i = ith point (ith element of point cloud)

$x_i^{(d)}$ = dth element of ith point

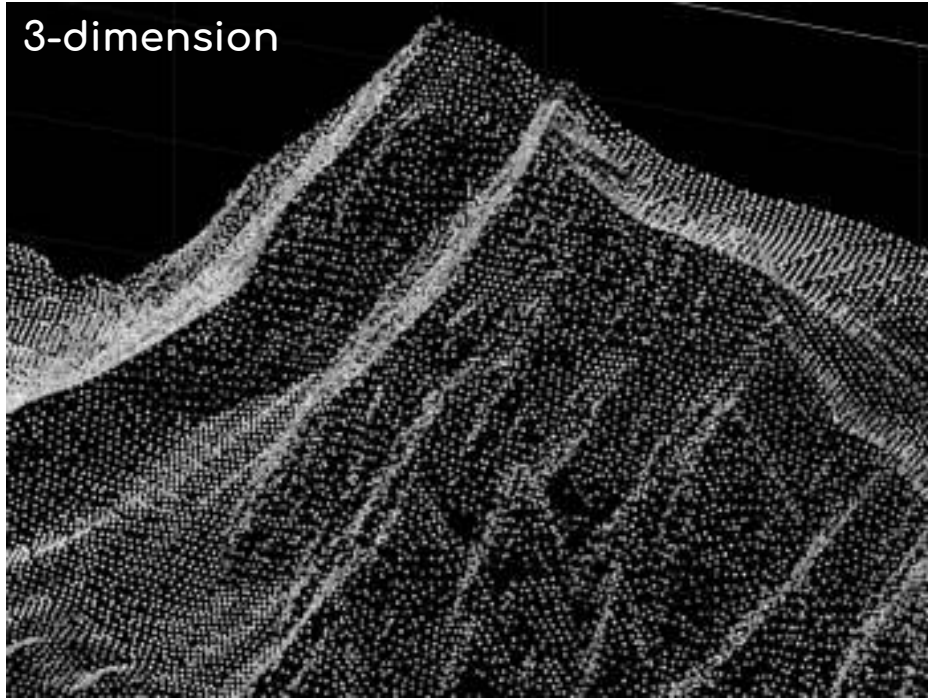
D = dimension of point cloud

N = size of point cloud (number of data points)



N=1000, D=3

Data Types / Point Cloud : real data

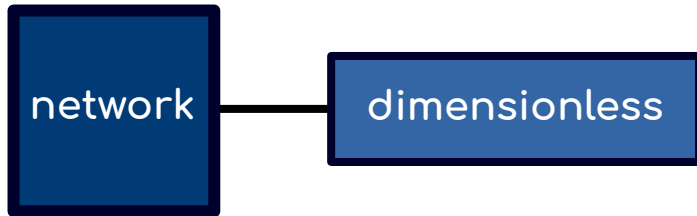


<https://www.geo.tuwien.ac.at/downloads/pg/pctools/publish/pointCloudThinOut/html/pointCloudThinOut.html>



<https://www.wired.com/2014/09/shaun-kardinal-flying-formation/>

Data Types / Network



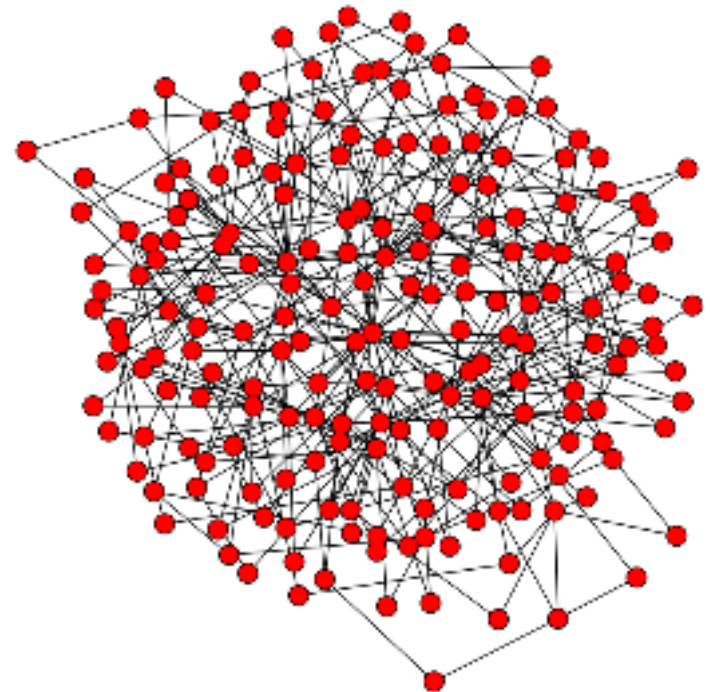
$G = (V, E, w)$ graph (network)

$V = \{v_i\}_{i=1}^N$ vertex (node) set
N = network size (number of nodes)

$E = V \times V$ edge (link) set (Cartesian product)
L = |E| = number of links

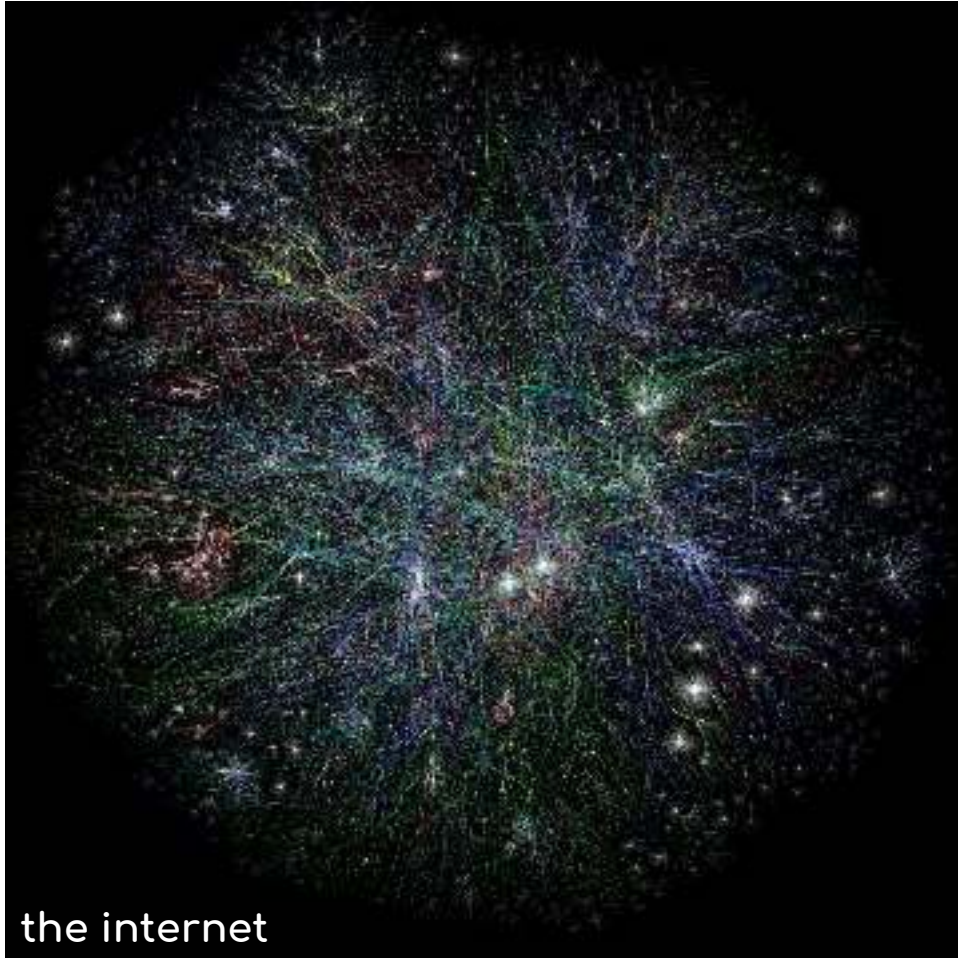
$w : E \rightarrow \mathbb{R}$ weight function (link map)

$w(e_{ij}) \equiv w((v_i, v_j)) = w_{ij}$



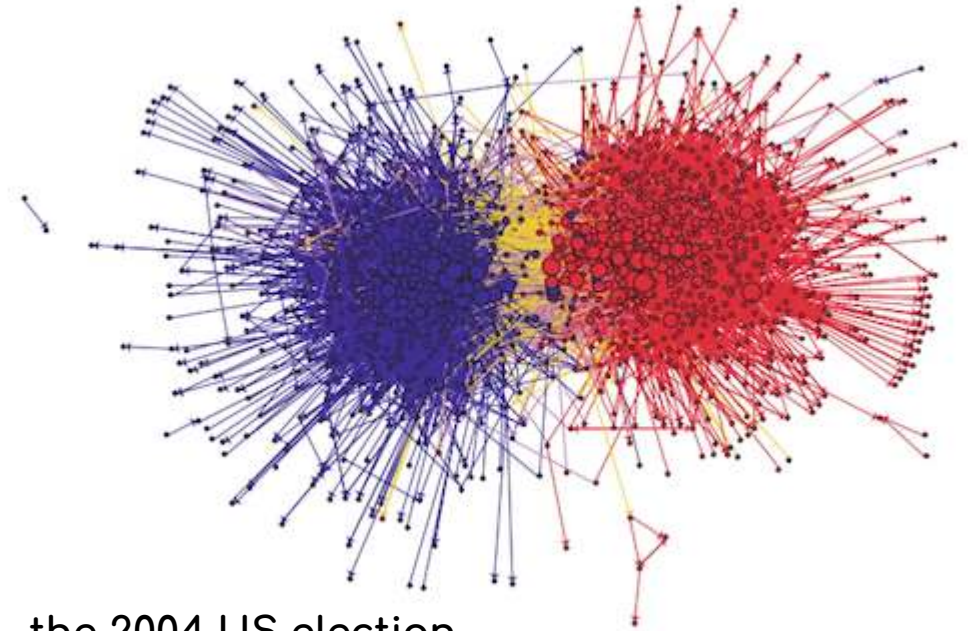
N=100, L=150

Data Types / Network : real data



the internet

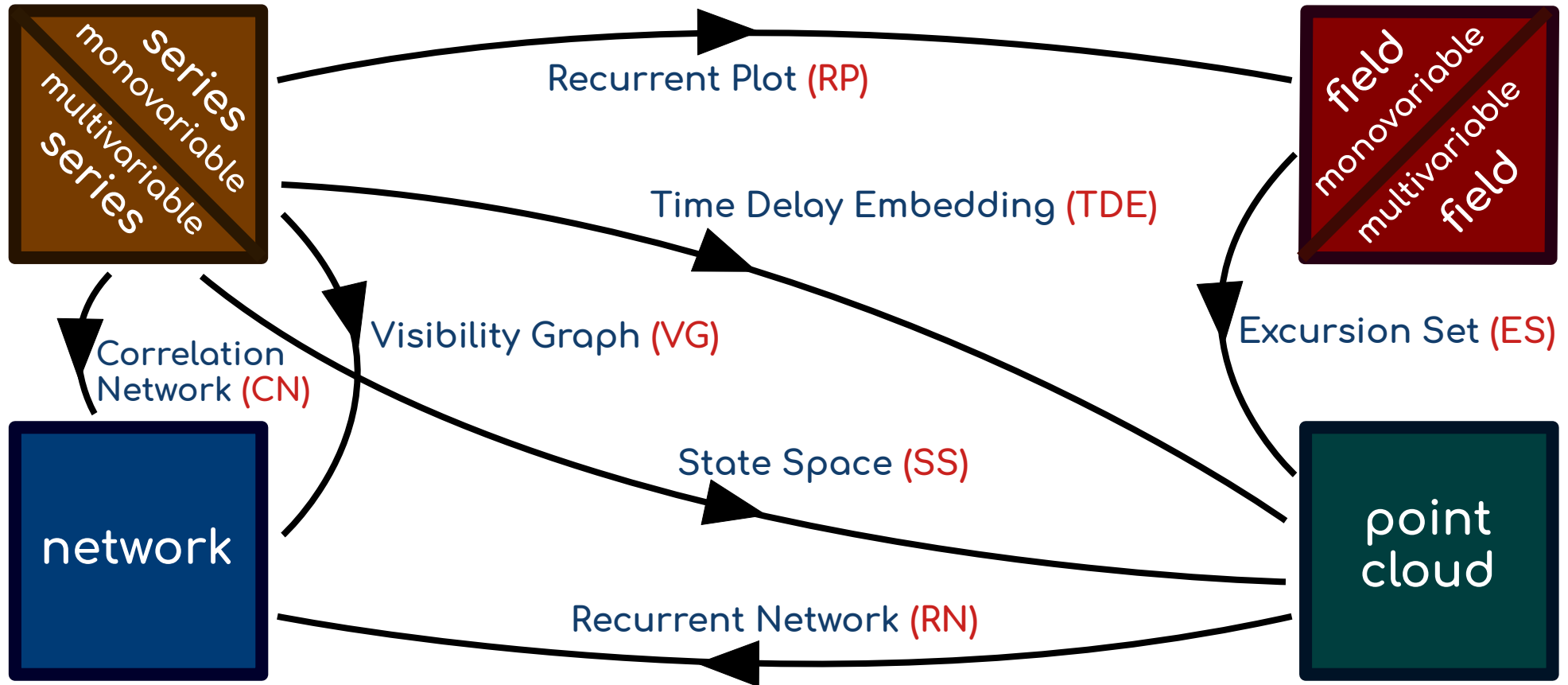
<https://www.kaspersky.com/blog/amazing-internet-maps/10441/>



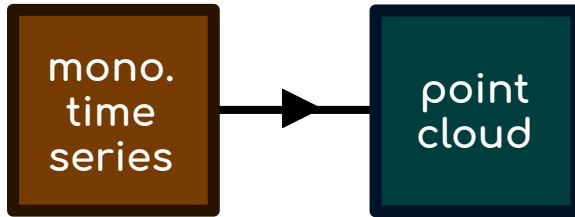
the 2004 US election

Adamic, Lada A., and Natalie Glance. "The political blogosphere and the 2004 US election: divided they blog." In Proceedings of the 3rd international workshop on Link discovery, pp. 36-43. 2005.

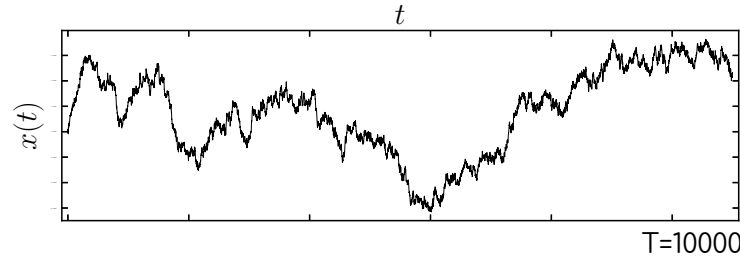
Methods for Reconstruction of Data Sets of Different Types



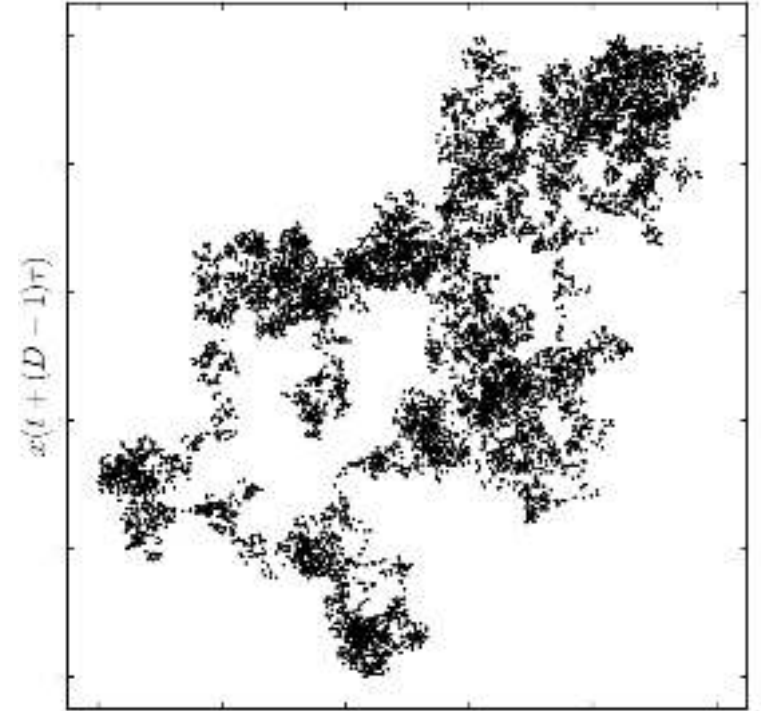
Methods ... / Time Delay Embedding (TDE)



$$\vec{x} \equiv \left(x(t_i) \right)_{i=1}^T$$



D = embedding dimension
(dimension of reconstructed point cloud)
 τ = time-delay



$$\mathbb{X}(\vec{x}) \equiv \left\{ x_i \in \mathbb{K}^D \mid x_i \equiv \left(x(t_i), x(t_i + \tau), x(t_i + 2\tau), \dots, x(t_i + (D-1)\tau) \right) \right\}_{i=1}^{T-(D-1)\tau} \quad x(t) \quad N=9000, D=2, \tau=1000$$

$$|\mathbb{X}(\vec{x})| = |\vec{x}| - \left(\dim(\mathbb{X}(\vec{x})) - 1 \right) \tau \quad |\mathbb{X}(\vec{x})| = N \quad \left[\dim(\mathbb{X}(\vec{x})) = D \quad |\vec{x}| = T \right]$$

Methods ... / Recurrent Plot (RP)



$$\vec{x} \equiv \left(x(t_i) \right)_{i=1}^T$$

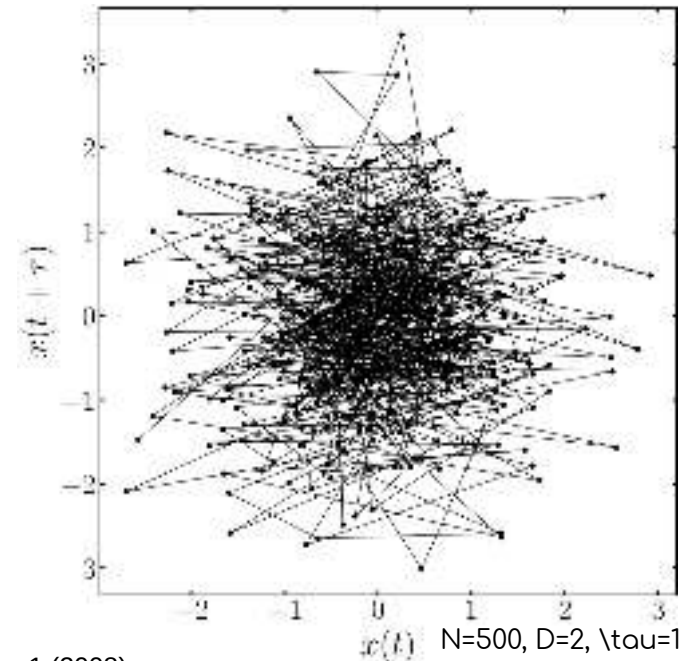
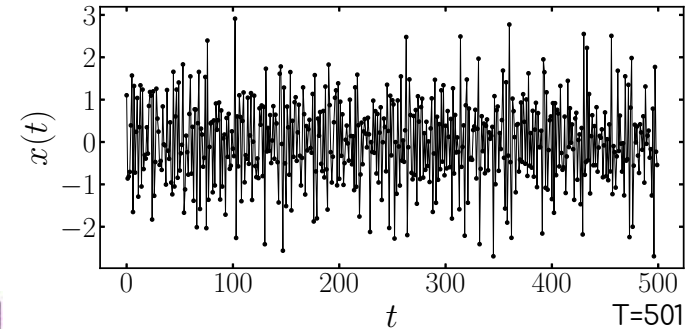
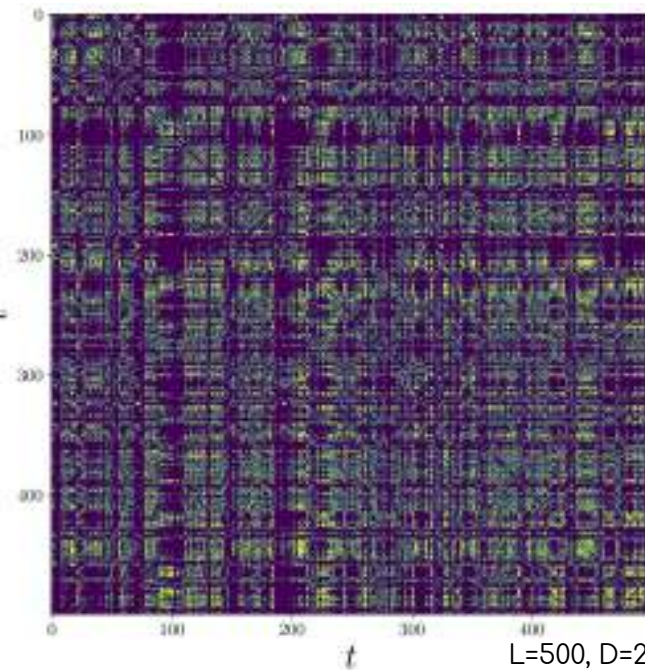
TDE : (D, τ)

$$\mathbf{X}(\vec{x}) \equiv \left\{ x_i \in \mathbb{R}^D \right\}_{i=1}^{T-(D-1)\tau}$$

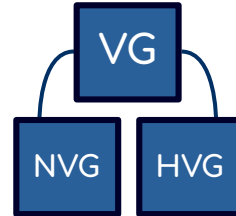
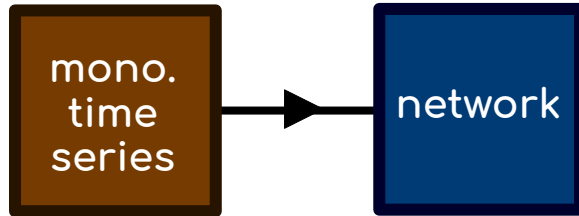
$\left[\begin{array}{l} N = T - (D-1)\tau \\ D \end{array} \right]$ size
dimension

$$\mathcal{F}_{ij}^{(\epsilon)} = \Theta(\epsilon - d(x_i, x_j))$$

$\left[\begin{array}{l} L = N \\ D \end{array} \right]$ length
dimension



Methods ... / Visibility Graph (VG)

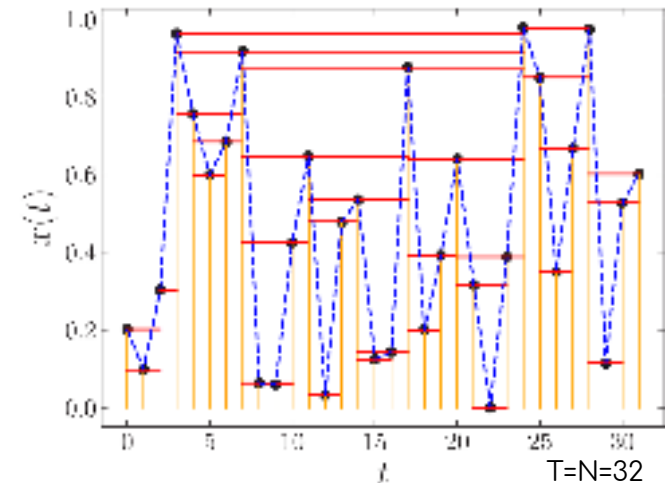
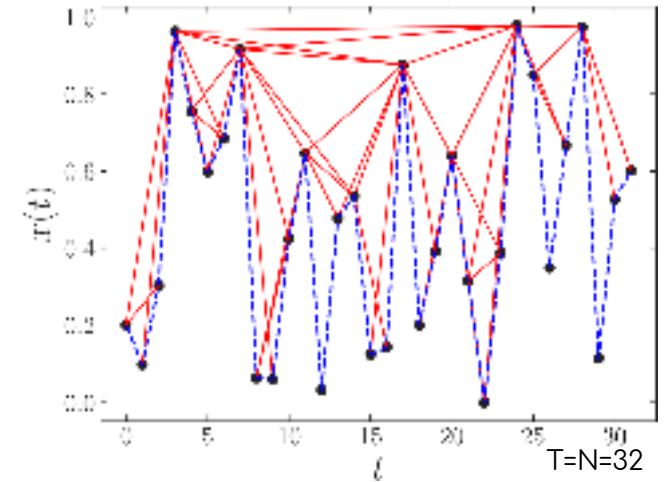


$$\vec{x} \equiv \left(x(t_i) \right)_{i=1}^T \quad \text{time series} \quad G = (V, E, w) \quad \text{network}$$

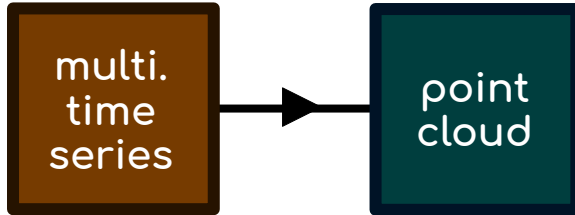
$$f : V \equiv \{v_i\}_{i=1}^N \leftrightarrow \mathcal{T} \equiv (t_i)_{i=1}^T, \quad f(v_i) = t_i \quad \text{bijection: } N = T$$

$$w_{ij}^{(RN)} \equiv \begin{cases} 1, & |f(v_i) - f(v_j)| = 1 \\ \prod_{k=i+1}^{j-1} \Theta(s_{ij} - s_{ik}), & |f(v_i) - f(v_j)| > 1 \end{cases} \quad \text{weight function}$$

$$s_{ij} \equiv \frac{x(f(v_j)) - x(f(v_i))}{f(v_j) - f(v_i)} \quad \text{slope of the visibility line between } i\text{th and } j\text{th data points}$$



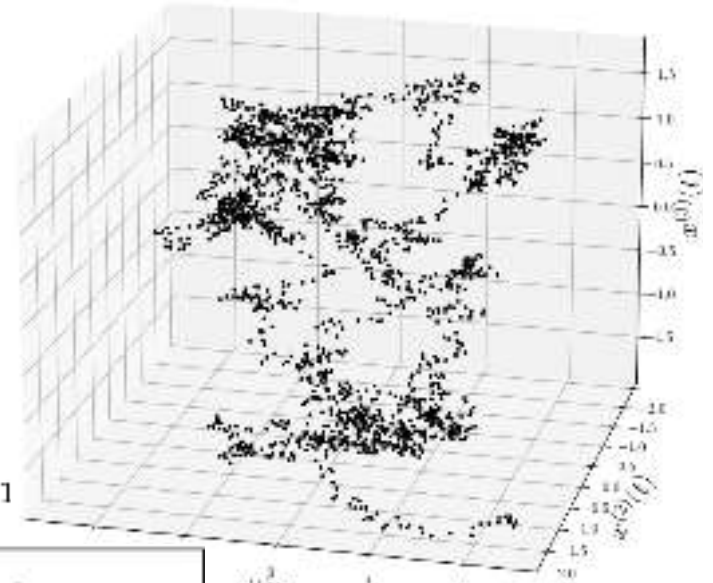
Methods ... / State Space (SS)



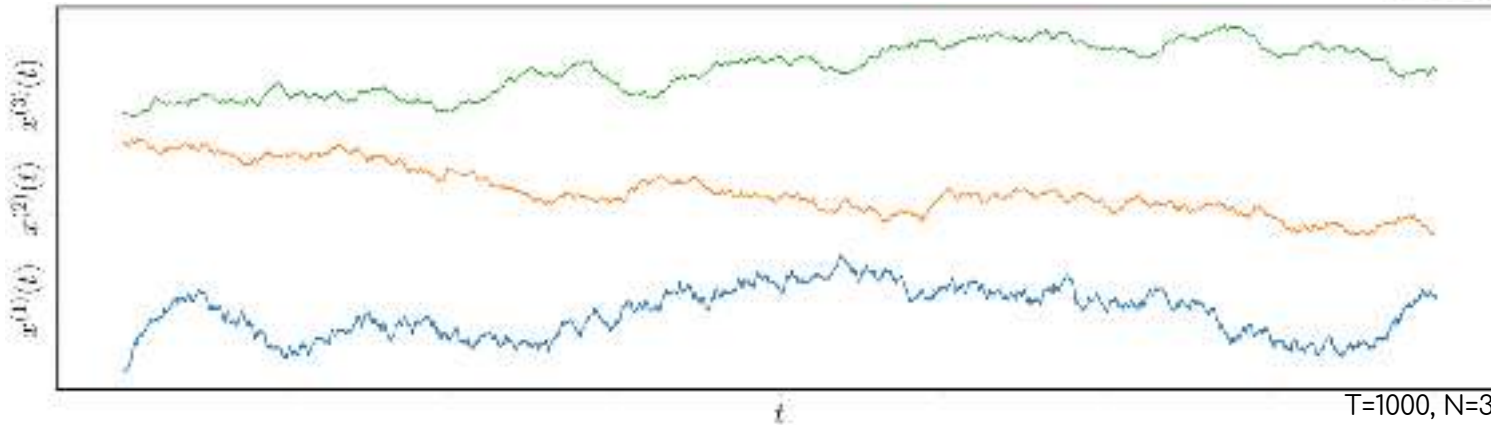
$$\mathcal{X} = \left\{ \vec{x}^{(j)} \mid \vec{x}^{(j)} = \left(x^{(j)}(t_i) \right)_{i=1}^T \right\}_{j=1}^N$$

N = number of dependent variables
T = length of time series

$$\mathbb{X}(\mathcal{X}) = \left\{ \vec{x}(t_i) \equiv \left(x^{(1)}(t_i), x^{(2)}(t_i), \dots, x^{(N)}(t_i) \right) \in \mathbb{R}^N \mid x^{(j)}(t_i) \in \mathbb{R} \right\}_{i=1}^T$$



N=1000, D=3

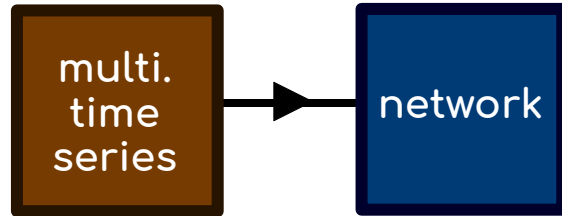


T=1000, N=3

$$D(\mathbb{X}) = N(\mathcal{X})$$

$$N(\mathbb{X}) = T(\mathcal{X})$$

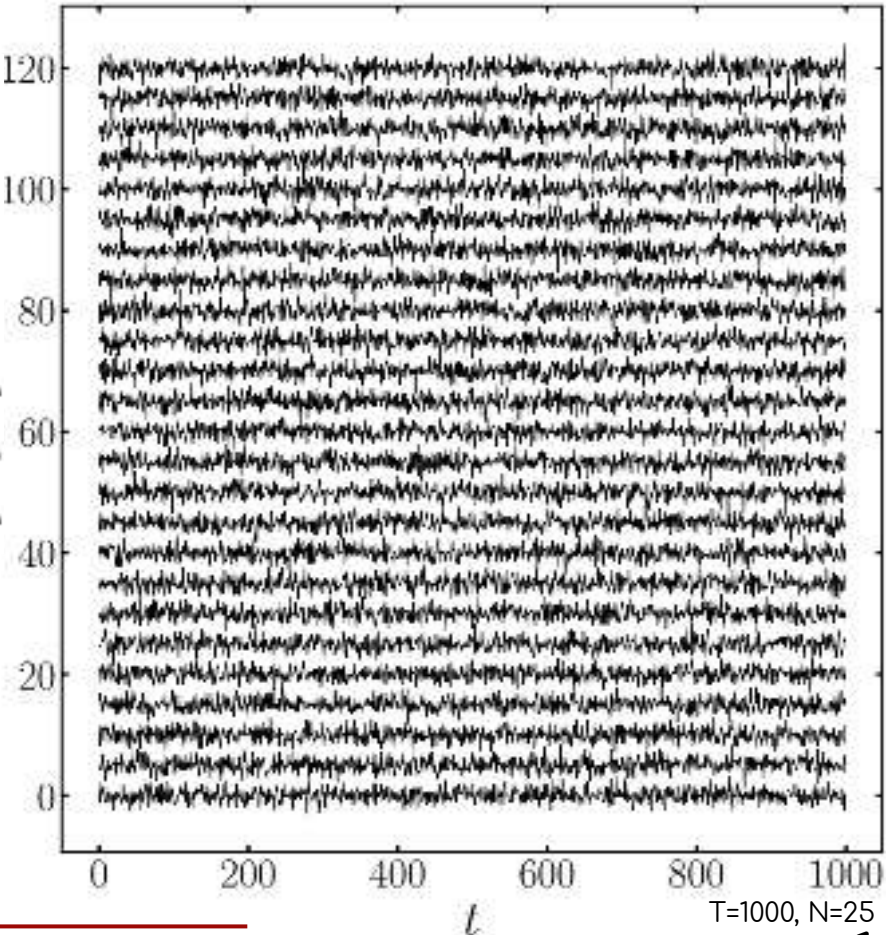
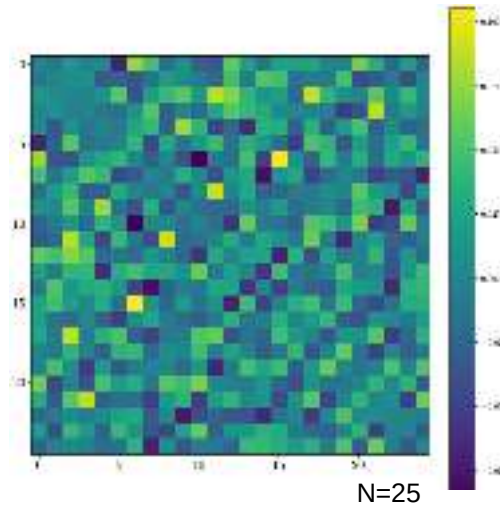
Methods ... / Correlation Network (CN)



$$\mathcal{X} = \left\{ \bar{x}^{(j)} \mid \bar{x}^{(j)} = \left(x^{(j)}(t_i) \right)_{i=1}^T \right\}_{j=1}^N$$

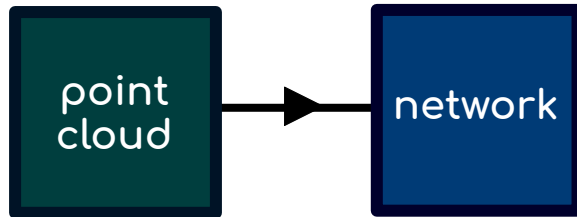
$$G = (V, E, w)$$

$$N : \mathcal{X} \leftrightarrow V \quad N(G) = N(\mathcal{X})$$



$$w_{jk} = \frac{1}{\min\{T_j, T_k\}} \sum_{i=1}^{\min\{T_j, T_k\}} (x^{(j)}(t_i) - \mu_{\bar{x}^{(j)}})(x^{(k)}(t_i) - \mu_{\bar{x}^{(k)}}) \bigg/ \left[\frac{1}{T_j} \sum_{i=1}^{T_j} (x^{(j)}(t_i) - \mu_{\bar{x}^{(j)}})^2 \cdot \frac{1}{T_k} \sum_{i=1}^{T_k} (x^{(k)}(t_i) - \mu_{\bar{x}^{(k)}})^2 \right]^{1/2}$$

Methods ... / Recurrent Network (RN)

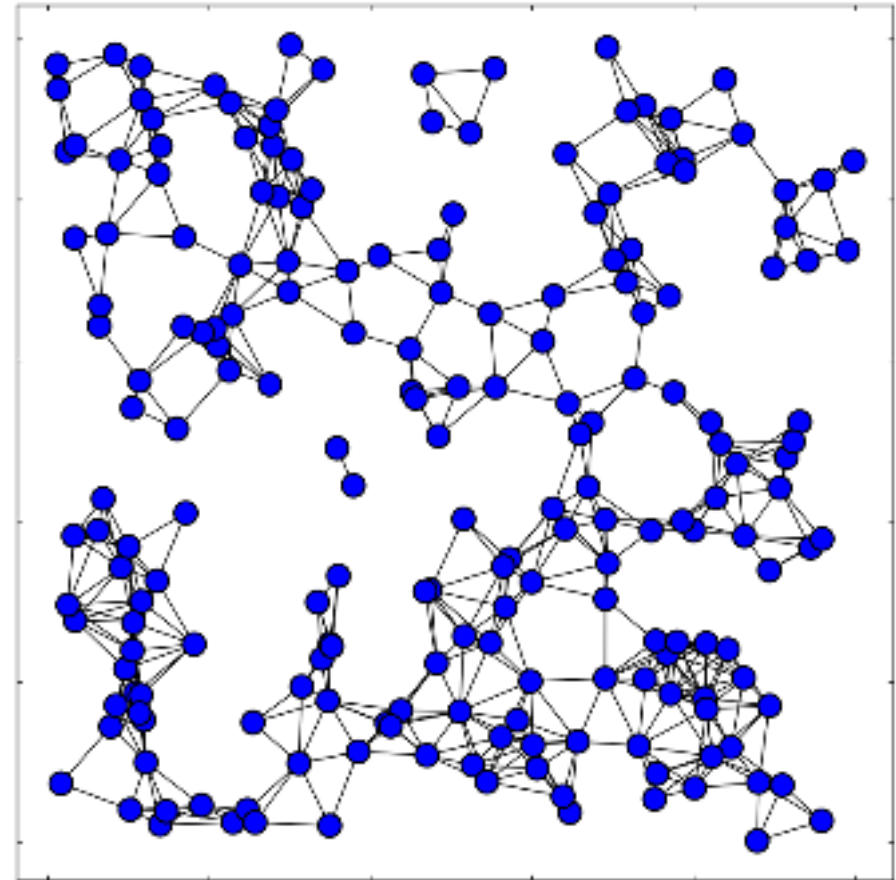


$$\mathbb{X} = \left\{ x_i \mid x_i = (x_i^{(d)})_{d=1}^D, x_i^{(d)} \in \mathbb{R} \right\}_{i=1}^{N \neq \infty}$$

$$G(\mathbb{X}, \epsilon) = (V, E, w^{(\epsilon)})$$

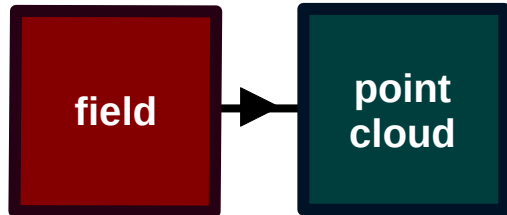
$$f : \mathbb{X} \leftrightarrow V \quad ; \quad f(x_i) = v_i \quad N(G) = N(\mathbb{X})$$

$$w^{(\epsilon)} : E \rightarrow \{0, 1\} \quad ; \quad w_{ij}^{(\epsilon)} = \Theta(\epsilon - d(x_i, x_j))$$



D=2, N=100, \epsilon=0.1, L=150

Methods ... / Excursion Set (ES)



$$\mathcal{F} : \Pi \rightarrow \mathbb{R} \quad (D,L)$$

$$\mathbb{X}(\mathcal{F}) = \left\{ \pi \in \Pi \mid \pi \in \mathcal{E}(\mathcal{F}) \right\}$$

for D=3: [local maxima and minima]

$$\bullet \mathcal{E}_{\max}(\Pi) = \max(\Pi) = \left\{ (i, j, k) \in \Pi \mid \mathcal{F}(i, j, k) > \max\{\mathcal{F}(\mathcal{N}(i, j, k))\} \right\}_{i,j,k=2}^{L-1}$$

$$\bullet \mathcal{E}_{\min}(\Pi) = \min(\Pi) = \left\{ (i, j, k) \in \Pi \mid \mathcal{F}(i, j, k) < \min\{\mathcal{F}(\mathcal{N}(i, j, k))\} \right\}_{i,j,k=2}^{L-1}$$

$$\mathcal{N}(i, j, k) = \left\{ (i-1, j, k), (i+1, j, k), (i, j-1, k), (i, j+1, k), (i, j, k-1), (i, j, k+1), \right.$$

$$(i-1, j-1, k), (i-1, j+1, k), (i-1, j, k-1), (i-1, j, k+1), (i, j-1, k-1),$$

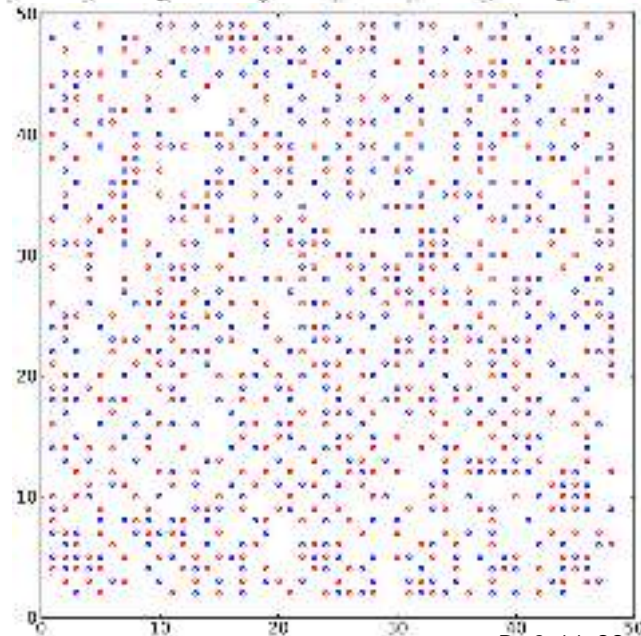
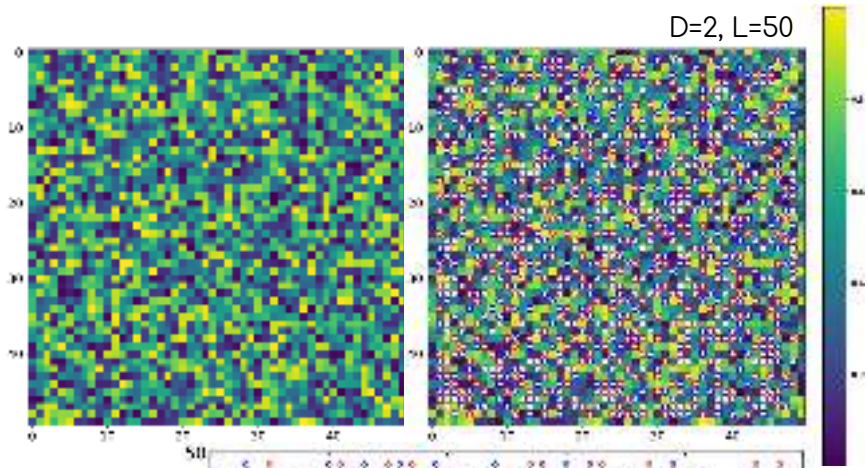
$$(i, j-1, k+1), (i, j+1, k-1), (i, j-1, k+1), (i+1, j-1, k), (i+1, j+1, k),$$

$$(i+1, j, k-1), (i+1, j, k+1), (i-1, j-1, k-1), (i-1, j+1, k-1),$$

$$(i-1, j-1, k+1), (i-1, j+1, k+1), (i+1, j-1, k-1), (i+1, j+1, k-1),$$

$$(i+1, j-1, k+1), (i+1, j+1, k+1) \left. \right\} \subset \Pi(\mathcal{F}) \quad \text{first neighbors}$$

$$D(\mathcal{F}) = D(\mathbb{X})$$

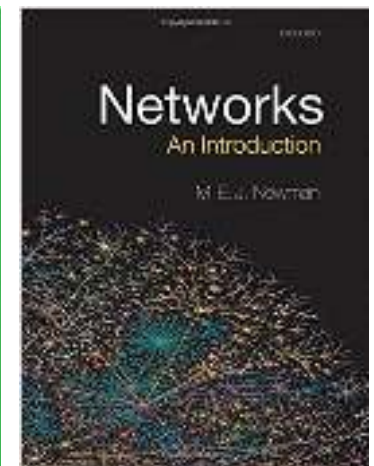
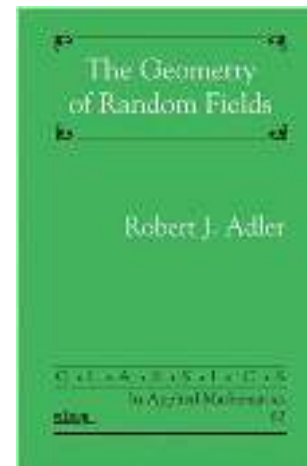
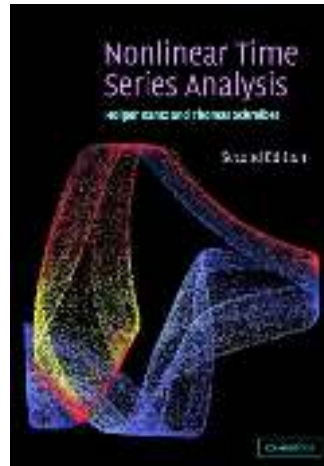


D=2, L=50

D=2, N=?? 18

References

- [1] H. Kantz, "Nonlinear time series analysis", Vol. 7, Cambridge university press, (2004).
- [2] R. J. Adler, "The geometry of random fields", Society for Industrial and Applied Mathematics, (2010).
- [3] M. Newman, "Networks", Oxford university press, (2018).
- [4] F. Takens, "Detecting strange attractors in turbulence", In Dynamical systems and turbulence, Warwick 1980, pp. 366-381. Springer, Berlin, Heidelberg, (1981).
- [5] N. Marwan, "A historical review of recurrence plots", The European Physical Journal Special Topics 164, no. 1 (2008).
- [6] L. Lacasa, et al, "From time series to complex networks: The visibility graph", Proceedings of the National Academy of Sciences 105, no. 13 (2008).
- [7] Y. Zou, et al, "Complex network approaches to nonlinear time series analysis", Physics Reports 787 (2019).
- [8] C. Chen, et al, "Recurrence network modeling and analysis of spatial data", Chaos: An Interdisciplinary Journal of Nonlinear Science 28, no. 8 (2018).
- [9] H. Masoomy, B. Askari, M. N. Najafi, and S. M. S. Movahed, "Persistent homology of fractional Gaussian noise", Physical Review E 104, no. 3, (2021).
- [10] H. Masoomy, and M. N. Najafi, "The Visibility Graphs of Correlated Time Series Violate the Barthelémy's Conjecture for Degree and Betweenness Centralities", arXiv preprint arXiv:2112.07698, (2021).



Thanks :)



hoseingmasoomy@gmail.com