

Learning high-dimensional data

Michel Verleysen
Université catholique de Louvain (Louvain-la-Neuve, Belgium)
Electricity department

January 2002



Acknowledgements

∕ Part of this work has been realized in collaboration with PhD students:

- ∕ Philippe Thissen
- ∕ Jean-Luc Voz
- ∕ Amaury Lendasse
- ∕ John Lee

∕ Some ideas and figures come from

- ∕ D. L. Donoho, High-Dimensional Data Analysis: The Curses and Blessings of Dimensionality. Lecture on August 8, 2000, to the American Mathematical Society "Math Challenges of the 21st Century". Available from <http://www-stat.stanford.edu/~donoho/>.



Data mining

D columns (dimension of space)

| | | | | | | |
|--|------|-----|------|------|-----|------|
| N lines (number of observations) | 1.2 | 7.5 | -1.9 | 2 | ... | 1.9 |
| | -7.6 | 12 | 17.2 | 2.4 | ... | 1.5 |
| | -8.5 | 13 | 14 | 8.5 | ... | -1.9 |
| | 9 | 5.4 | -5.2 | 8.2 | ... | 9.4 |
| | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| | 25.1 | 5.2 | -9.1 | -8.5 | ... | 5.4 |

⚡ Data mining: find information in large databases
(large = D and/or N)

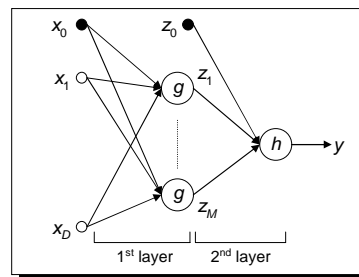


High-Dimensional spaces

- ⚡ Inputs = **High-Dimension (HD)** vectors (D is large)
 - ⚡ many parameters in the model
 - ⚡ local minima
 - ⚡ slow convergence

⚡ The questions

- ⚡ Learning algorithms in HD spaces ?
- ⚡ Local learning or not ?



⚡ The arguments

- ⚡ real data are seldom HD (concept of intrinsic dimension)
- ⚡ local learning is not worse than global learning...



Contents

- ⚡ High-dimensional data
 - ⚡ Surprising results
 - ⚡ Intrinsic dimension
- ⚡ Local learning
 - ⚡ Use of distance measures
- ⚡ Dimension reduction
 - ⚡ non-linear projection
 - ⚡ application to time-series forecasting



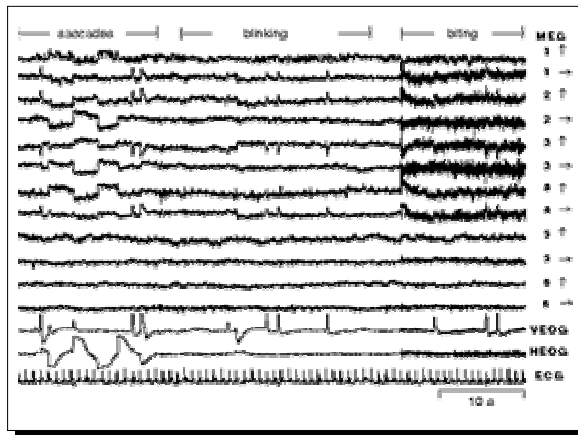
Contents

- ⚡ High-dimensional data
 - ⚡ Surprising results
 - ⚡ Intrinsic dimension
- ⚡ Local learning
 - ⚡ Use of distance measures
- ⚡ Dimension reduction
 - ⚡ non-linear projection
 - ⚡ application to time-series forecasting



Data mining: large databases

⚡ (biomedical) signals

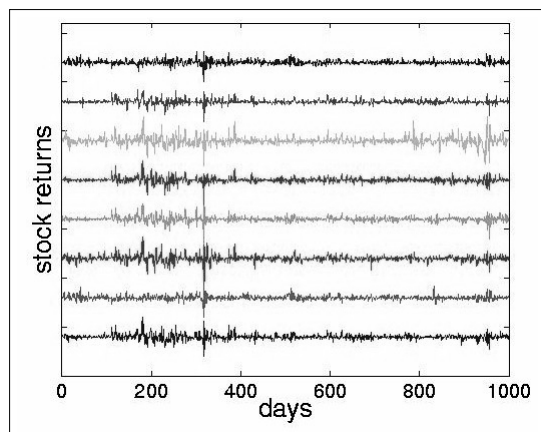


Michel Verleysen

7

Data mining: large databases

⚡ financial data



Michel Verleysen

8

Data mining: large databases

∕ imagery



Data mining: large databases

∕ recordings of consumers' habits (credit cards)

∕ bio- data (human genome, etc.)

∕ satellite images

∕ hyperspectral images

∕ ...



John Wilder Tukey

⚡ The Future of Data Analysis, *Ann. Math. Statis.*, 33, 1-67, 1962.

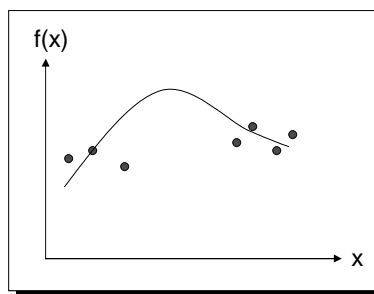
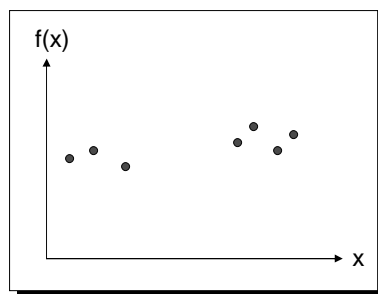
« Analyze data rather than prove theorems... »

⚡ In other words:

- ⚡ data are here
- ⚡ they will be coming more and more in the future
- ⚡ we must analyze them
- ⚡ with very humble means
- ⚡ insistence on mathematics will distract us from fundamental points



Empty space phenomenon

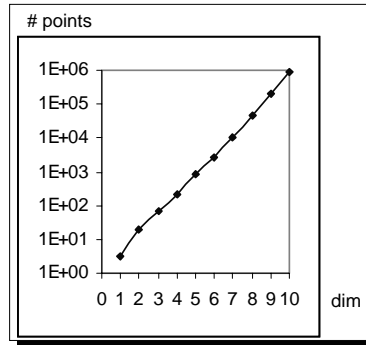


- ⚡ Necessity to *fill* space with learning points
- ⚡ # learning points exponential with dimension



Example: Silvermann's result

- ⚡ How to approximate a Gaussian distribution with Gaussian kernels
- ⚡ Desired accuracy: 90%

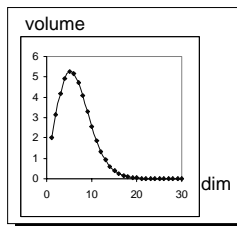


Michel Verleysen

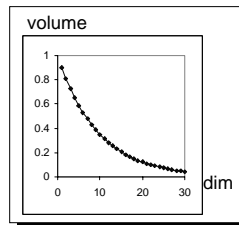
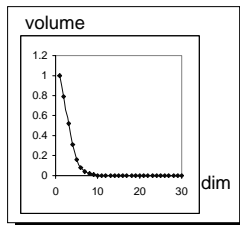
13

Surprising phenomena in HD spaces

- ⚡ Sphere volume



- ⚡ Sphere volume / cube volume
- ⚡ Embedded spheres (radius ratio = 0.9)

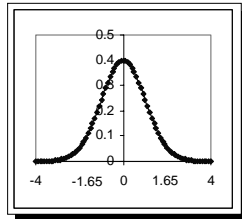


Michel Verleysen

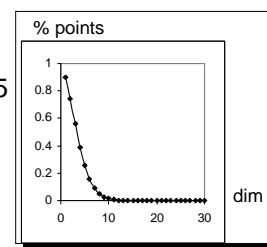
14

Gaussian kernels

⚡ 1-D Gaussian



⚡ % points inside a sphere of radius 1.65



Concentration of measure phenomenon

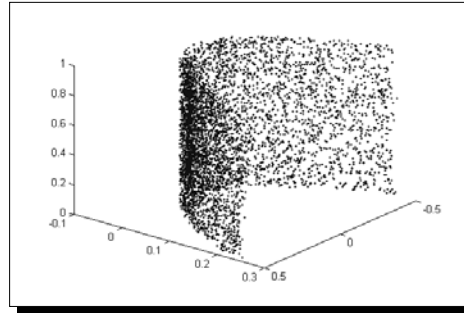
- ⚡ Take all pairwise distances in random data
- ⚡ Compute the average A and the variance V of these distances
- ⚡ If D increases then
 - ⚡ V remains fixed
 - ⚡ A increases
- ⚡ All distances seem to concentrate !!!

- ⚡ Example: Eucliden norm of samples
 - ⚡ average A increases with $(D)^{0.5}$
 - ⚡ variance V remains fixed
 - ⚡ \rightarrow samples seem to be normalized !



Intrinsic dimension

- ⚡ No definition
- ⚡ Example:



From "Analyse de données par réseaux de neurones auto-organisés », P. Demartines, Ph.D. thesis, Institut National Polytechnique de Grenoble (France), 1994.

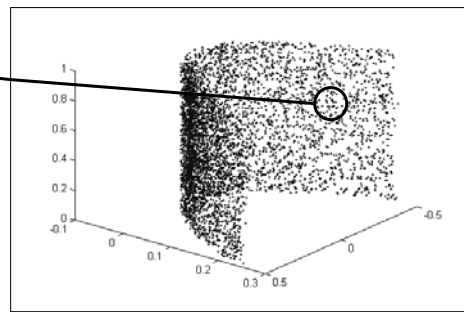


Estimation of intrinsic dimensions 1/4

| | |
|--------------|-------------------------|
| local PCA | a posteriori estimation |
| box counting | Grassberger-Proccacia |

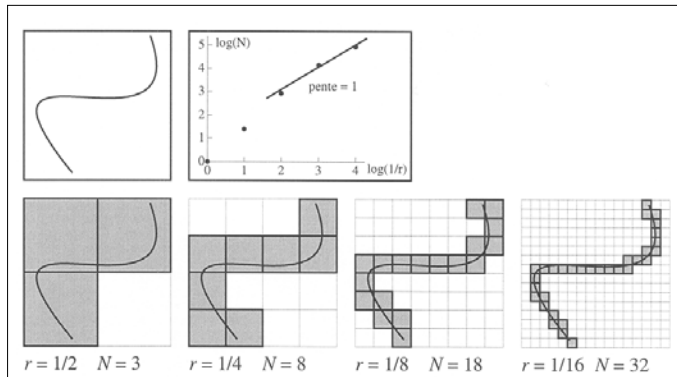
- ⚡ small region → approximately plane

- ⚡ PCA applied on small regions



Estimation of intrinsic dimensions 2/4

| | |
|--------------|-------------------------|
| local PCA | a posteriori estimation |
| box counting | Grassberger-Proccacia |



Michel Verleysen

From P. Demartines,
op. cit.

19

Estimation of intrinsic dimensions 3/4

| | |
|--------------|-------------------------|
| local PCA | a posteriori estimation |
| box counting | Grassberger-Proccacia |

- ⚡ Similar to box counting
- ⚡ Mutual distances between pairs of points
- ⚡ Advantage: $N(N-1)/2$ distances for N points
- ⚡ $\log N - \log r$ graph: number N of pairs of points closer than r

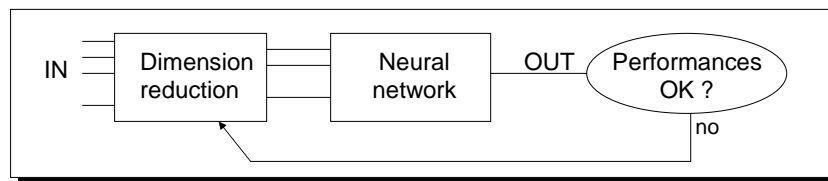


Michel Verleysen

20

Estimation of intrinsic dimensions 4/4

| | |
|--------------|-------------------------|
| local PCA | a posteriori estimation |
| box counting | Grassberger-Proccacia |

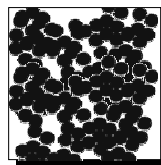
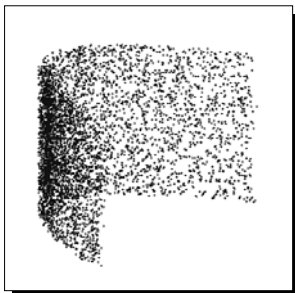


∴ example: forecasting problem

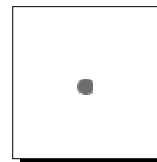


Intrinsic dimension: limitation of the concept

∴ Seen from very close



∴ Seen from very far



Contents

- ⚡ High-dimensional data
 - ⚡ Surprising results
 - ⚡ Intrinsic dimension

- ⚡ Local learning
 - ⚡ Use of distance measures

- ⚡ Dimension reduction
 - ⚡ non-linear projection
 - ⚡ application to time-series forecasting



Local learning

- ⚡ « Local »: by means of local functions
- ⚡ Typical example: Gaussian kernels
- ⚡ Radial Basis Function Networks

$$f(\mathbf{x}) = \sum_{j=1}^P w_j K_j(\mathbf{x}) \quad K_j(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_j\|^2}{h_j^2}\right)$$

- ⚡ Local **and** global: ANN are interpolators, and do not extrapolate to regions without learning data !



Radial-Basis Function Networks (RBFN) for approximation

$$f(\mathbf{x}) = \sum_{j=1}^P w_j K_j(\mathbf{x}) \quad K_j(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_j\|^2}{h_j^2}\right)$$

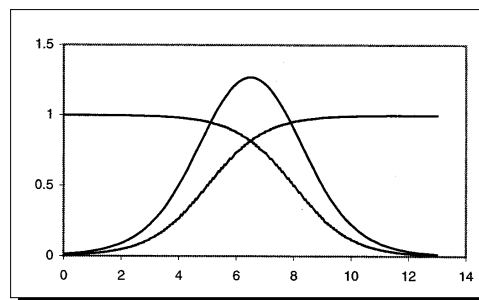
⚡ Advantages (over MLP, ...):

- ⚡ splitted computation of
 - centres \mathbf{c}_j
 - widths h_j
 - multiplying factors w_j
- ⚡ easier learning



Local learning 2/2

⚡ Sum of sigmoids = Gaussian !



Problem with (local ?) learning

- ⚡ Most ANN use *distances* between input and weight vectors:
 - ⚡ RBFN: as argument to radial kernels
 - ⚡ VQ and SOM: to choose the « winner »
 - ⚡ first layer in MLP: $\mathbf{x} \mathbf{w}$
- ⚡ In high-dimensional spaces: all these distances seem identical !
(concentration of measure phenomenon)
- ⚡ → need for
 - ⚡ other neural networks
 - ⚡ same neural networks, with other distance measures



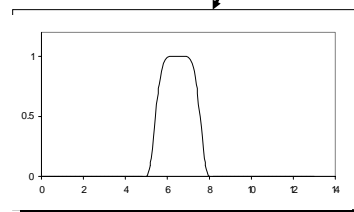
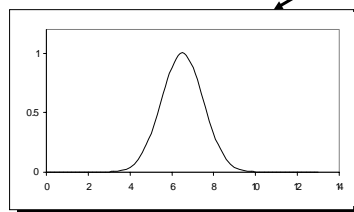
Using new distance measures

- ⚡ Example: RBF

$$f(\mathbf{x}) = \sum_{j=1}^P w_j K_j(\mathbf{x}) \quad K_j(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_j\|^2}{h_j^2}\right)$$

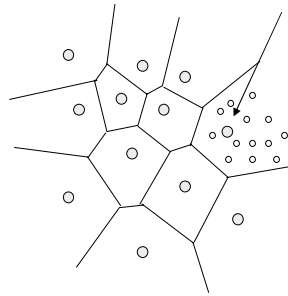
- ⚡ Use of « super-Gaussian » kernels:

$$K_j(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_j\|^r}{h_j^r}\right)$$



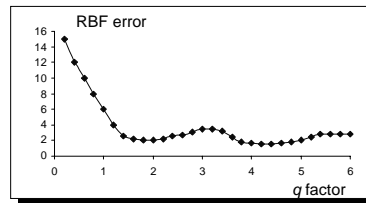
Not assuming evidence...

- ⚡ widths of kernels in RBF are not necessarily equal to the STDV of samples in the Voronoi zone!



$\sigma_j = \text{STDV}(\text{points in Voronoi zone})$

$$K_j(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_j\|^2}{q\sigma_j^2}\right)$$

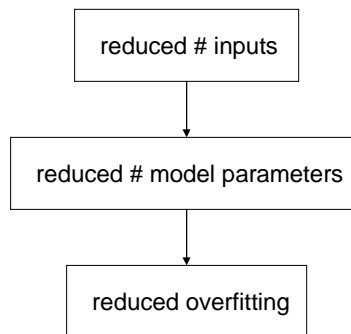


Contents

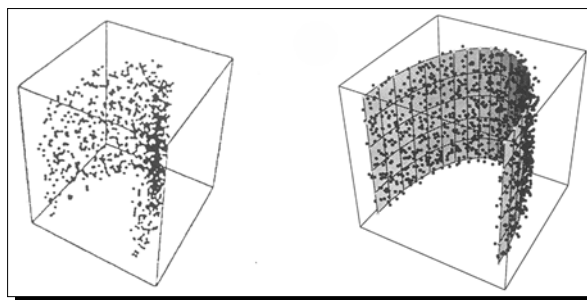
- ⚡ High-dimensional data
 - ⚡ Surprising results
 - ⚡ Intrinsic dimension
- ⚡ Local learning
 - ⚡ Use of distance measures
- ⚡ Dimension reduction
 - ⚡ non-linear projection
 - ⚡ application to time-series forecasting



Dimension reduction



Dimension reduction & intrinsic dimension



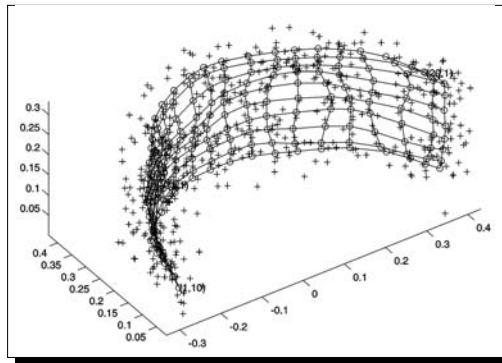
// Questions:

- // intrinsic dimension is unknown
- // non-linear submanifolds



Kohonen maps

⚡ Based on *topology* preservation



Distance-preservation methods

⚡ Many non-linear projection methods:
based on *distance* preservation between input and output pairs

⚡ Examples

- ⚡ Multi-dimensional scaling (MDS)
- ⚡ Sammon's mapping
- ⚡ Curvilinear Component Analysis

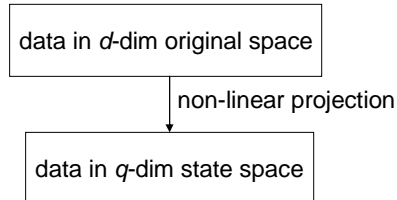
⚡ Principle: place points in the output space, so that pairwise distances are as equal as possible to the corresponding pairwise distances in the input space

⚡ Impossible to respect *all* pairwise distances → insist on small ones (*local* pairs)



Curvilinear Component Analysis

⚡ Principle



⚡ $q < d$

⚡ respect of mutual distance between pairs of points

$$E = \sum_{i=1}^N \sum_{j=1}^N (X_{ij} - Y_{ij})^2 F(Y_{ij}, \lambda)$$



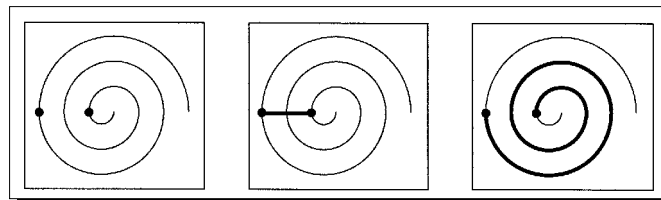
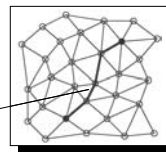
Adjusting parameters in CCA

⚡ Need for curvilinear distance

⚡ VQ → centroids

⚡ linking centroids

⚡ measuring the distances via the links

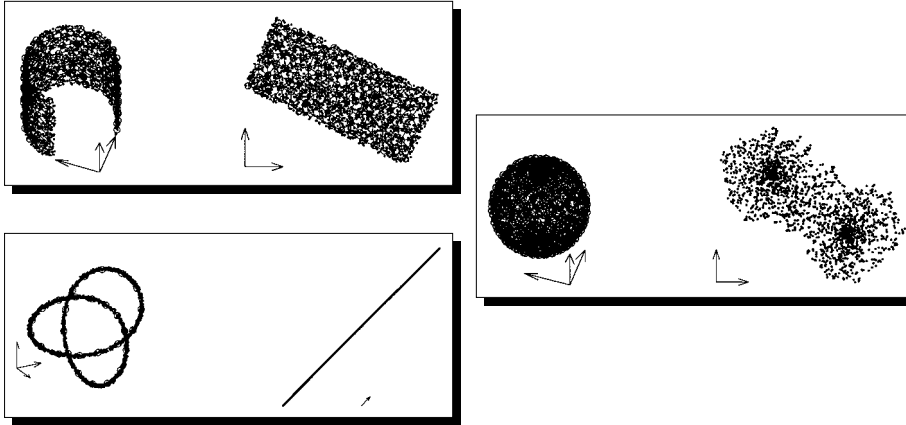


⚡ In practice: generalized distance

$$\Delta_{ij} = (1 - \omega(t))X_{ij} + \omega(t)\delta_{ij}$$



CCA: examples



Michel Verleysen

37

Application of dimension reduction to forecasting tasks

- ⚡ Regressor: - past values $x(t-i)$
- exogenous data $in(j)$

- ⚡ Forecasting

$$x(t+1) = f(x(t), x(t-1), \dots, x(t-k), in(1), in(2), \dots, in(l))$$

- ⚡ Non-linear forecasting:

- ⚡ 1. optimise regressor on linear predictor
- ⚡ 2. use the same regressor with non-linear predictor f
- ⚡ trials and errors (computational load !)



Michel Verleysen

38

Forecasting: selection of input variables

⚡ Starting with many input variables, then reduce their number

⚡ Two options:

1. selection of input variables

➤ interpretability

➤ limited to existing variables

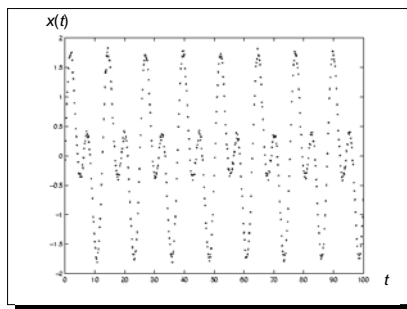
2. projection of input variables

⚡ linear: PCA

⚡ non-linear: CCA, Kohonen, etc.
based on **Takens' theorem**

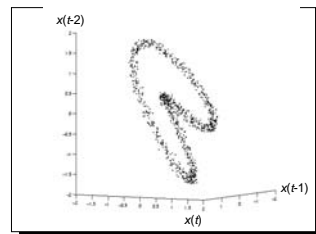
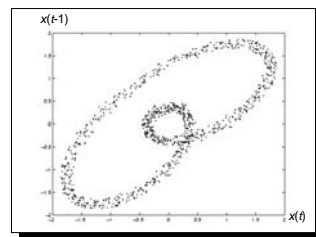


Forecasting: Takens' theorem 1/2



time series

intrinsic dimension (q) = 1



⚡ Takens' theorem:

$q \leq \text{size of regressor} \leq 2q+1$
(AR model)



Forecasting: Taken's theorem 2/2

⚡ Takens' theorem:

$$q \leq \text{size of regressor} \leq 2q+1$$

⚡ In the $2q+1$ space, there exists a q -surface without intersection points

⚡ Projection from $2q+1$ to q possible !

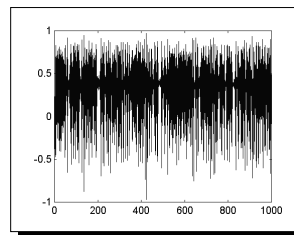


Forecasting: 1st example 1/2

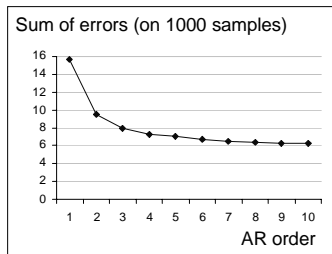
⚡ Artificial series

$$x(t+1) = ax(t)^2 + bx(t-2) + \varepsilon(t)$$

Two past values!

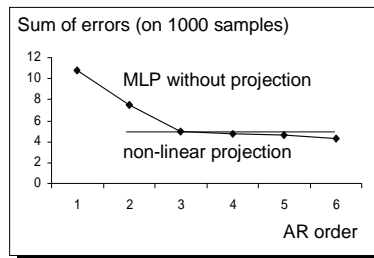


⚡ Linear AR model



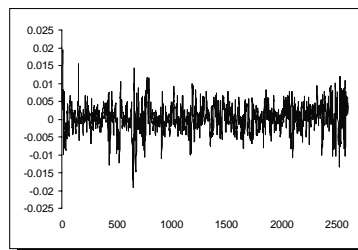
Forecasting: 1st example 2/2

- ⚡ Non-linear AR model
 - ⚡ initial regressor: size=6
 - ⚡ intrinsic dimension: 2
 - ⚡ CCA from dim=6 to dim=2
 - ⚡ MLP on 2-dim data



Forecasting: 2nd example 1/2

- ⚡ Daily returns of BEL20 index



- ⚡ 42 indicators from inputs and exogenous variables:

- ⚡ returns: $x_t, x_{t-10}, x_{t-20}, x_{t-40}, \dots, y_t, y_{t-10}, \dots$
- ⚡ differences of returns: $x_t - x_{t-5}, x_{t-5} - x_{t-10}, \dots, y_t - y_{t-5}$
- ⚡ oscillators: $K(20), K(40), \dots$
- ⚡ moving averages: $MM(10), MM(50), \dots$
- ⚡ exponential moving averages: $MME(10), MME(50), \dots$
- ⚡ etc



Forecasting: 2nd example 2/2

Method:

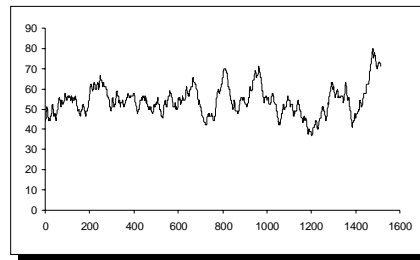
- 42 indicators
- PCA → 25 variables

Grassberger-Proccacia: intrinsic dimension = 9

- CCA → 9 variables
- RBF → forecasting

Result: % of correct approximations of sign (90-days average)

- In average: 57.2% on test set



Conclusion

- High dimensions: > 3 !
- NN in general: also difficulties in high dimensions
- Common problems:
 - Euclidean distance
 - empty space phenomenon
- Towards solutions...
 - Local NN (RBFN, etc.): easier learning
 - Generic methods for non-linear projections: reduce dimension
- Open perspectives:
 - using dimension reduction techniques based on topology (and not on distances)
 - study the possibility of non-Euclidean distances and non-Gaussian kernels

