Towards a low-cost scalable data-driven modeling of coupled systems and non-stationary time series analysis

IPAM Workshop, Los Angeles October 28 - November 1, 2019



<u>Illia Horenko</u>

USI Lugano, Switzerland



Susanne Gerber

JGU Mainz, Germany



Lukáš Pospíšil

VSB-TU Ostrava, Czech Republic

Patrick Gagliardini Will Sawyer Ganna Marchenko Karin Everschor-Sitte Terry O'Kane et al.

Outline

- 1. Non-stationary time series analysis
 - General talk about time series (polynomial regression, AIC,...)
 - Non-stationary modeling
 - K-means: unsupervised clustering algorithm
 - FEM-H1, FEM-BV
- 2. Scalable Probabilistic Approximation (SPA)
 - Exact Law of the Total Probability
 - Clustering = the cure for the curse of dimensionality?
 - Soft (probabilistic) K-means without Jensen inequality
 - Compressing the pipeline:

the combination of optimal discretization with Markov model

Regime-based non-stationary time series analysis with regularization (FEM-H1, FEM-BV)



• time series = data measured in intervals over a period of time

$$x_t \in \mathbb{R}^n, t = 1, \dots, T$$



$$x_t = \mu(t, \Theta) + \varepsilon_t$$

model function noise

• time series = data measured in intervals over a period of time

$$x_t \in \mathbb{R}^n, t = 1, \dots, T$$
noise --->
time --->
black box --->
time series



$$x_t = \mu(t, \Theta) + \varepsilon_t$$

• <u>Example:</u> linear regression

$$\mu(t,\theta_0,\theta_1) = \theta_1 t + \theta_0$$

• time series = data measured in intervals over a period of time

$$x_t \in \mathbb{R}^n, t = 1, \dots, T$$

choose model



• time series = data measured in intervals over a period of time

$$x_t \in \mathbb{R}^n, t = 1, \dots, T$$

- choose model
- find parameters of (optimal) model





- more "sophisticated" parametric models:
 - general polynomial regression
 - autoregressive models
 - Hidden Markov models
 - Neural Networks



$$\bar{\Theta} = \arg\min_{\Theta \in \mathcal{P}} \sum_{t=1}^{T} \rho(x_t, \mu(t, \Theta))$$

- more "sophisticated" parametric models:
 - general polynomial regression
 - autoregressive models
 - Hidden Markov models
 - Neural Networks





- more "sophisticated" parametric models:
 - general polynomial regression
 - autoregressive models
 - Hidden Markov models
 - Neural Networks



$$\bar{\Theta} = \arg\min_{\Theta \in \mathcal{P}} \sum_{t=1}^{T} \rho(x_t, \mu(t, \Theta))$$

- more "sophisticated" parametric models:
 - general polynomial regression
 - autoregressive models
 - Hidden Markov models
 - Neural Networks





- more "sophisticated" parametric models:
 - general polynomial regression
 - autoregressive models
 - Hidden Markov models
 - Neural Networks

• in general:

more parameters = lower modeling error

• in general:

more parameters = lower modeling error

- noise reduction ?
- overfitting ?
- interpretation of results ?

• •

• in general:

more parameters = lower modeling error

more parameters = more unknowns



- larger optimization problem
- harder to solve

• in general:

more parameters = lower modeling error

more parameters = more unknowns



- larger optimization problem
- harder to solve







Among competing hypotheses, the one with the fewest assumptions should be selected.

• Akaike H.: Information theory and an extension of the maximum likelihood principle, (1971)

Example: Moore's law

Moore's law:

"the number of transistors in a dense integrated circuit doubles about every two years."



• H. Khan, D. Hounshell, E. Fuchs: Science and research policy at the end of Moore's law, Nature Electronics, 1: 14–21, (2018)

Example

Filtering Example: signal/noise $\approx 2/1$



MIT Machine Learning Toolbox: Gaussian smoothing and HMM-Gauss in optimal setting



MIT Machine Learning Toolbox: Gaussian smoothing and HMM-Gauss in optimal setting

FEM-H1: demonstration by results



From non-stationary models to clustering

$$\bar{\Theta}(t) = \arg\min_{\Theta} \sum_{t=1}^{T} \|x_t - \mu(t, \Theta(t))\|^2$$
 (non-stationary model)

"constant" regression in regimes



From non-stationary models to clustering

$$\bar{\Theta}(t) = \arg\min_{\Theta} \sum_{t=1}^{T} \|x_t - \mu(t, \Theta(t))\|^2 \qquad \text{regime indicator function:} \\ \gamma_k(t) = \begin{cases} 1 & \text{if } x_t \text{ belongs to } k\text{-th regime indicator function:} \\ \gamma_k(t) = \begin{cases} 1 & \text{if } x_t \text{ belongs to } k\text{-th regime indicator function:} \\ 0 & \text{elsewhere } \end{cases}$$

$$\bar{\nabla}_k(t) = \sum_{t=1}^{T} \sum_{k=1}^{K} \gamma_k(t) \cdot \|x_t - \theta_k\|^2$$

$$\bar{\Theta}_k(t) = \sum_{t=1}^{T} \sum_{k=1}^{K} \gamma_k(t) \cdot \|x_t - \theta_k\|^2$$

$$\bar{\Theta}_k(t) = \sum_{t=1}^{T} \sum_{k=1}^{K} \sum_{t=1}^{T} \sum_{k=1}^{K} \gamma_k(t) \cdot \|x_t - \theta_k\|^2$$

$$\bar{\Theta}_k(t) = \sum_{t=1}^{T} \sum_{k=1}^{T} \sum_{t=1}^{K} \sum_{k=1}^{T} \sum_{t=1}^{T} \sum_{k=1}^{T} \sum_{t=1}^{T} \sum_{k=1}^{T} \sum_{t=1}^{T} \sum_{k=1}^{T} \sum_{t=1}^{T} \sum_{k=1}^{T} \sum_{t=1}^{T} \sum_{k=1}^{T} \sum_{t=1}^{T} \sum$$

K-means algorithm: unsupervised clustering



• Lloyd, Stuart P. : Least squares quantization in PCM. Information Theory, IEEE Transactions on 28.2 (1982): 129-137.

K-means algorithm: unsupervised clustering



- we don't have any apriori classification for supervised learning
- each cluster consists of *similar* points (points in cluster are close to each other)
- each cluster can be characterised by *mean value* ("centroid") of points inside it

K-means algorithm



set feasible initial approximation Γ_0

while
$$||L(\Gamma_{it}, \theta_{it}) - L(\Gamma_{it-1}, \theta_{it-1})|| > \varepsilon$$

solve $\theta_{it} = \arg\min_{\theta} L(\theta, \Gamma_{it})$ (with fixed Γ_{it})
solve $\Gamma_{it} = \arg\min_{\Gamma \in \Omega_{\Gamma}} L(\theta_{it}, \Gamma)$ (with fixed θ_{it})
 $it = it + 1$
endwhile

K-means algorithm

$$\Gamma_{it} = \arg\min_{\Gamma \in \Omega_{\Gamma}} \sum_{k=1}^{K} \sum_{t=1}^{T} \gamma_k(x_t) \|x_t - \mu_k\|^2$$

subject to $\gamma_k(x_t) \in \{0, 1\}, \forall t : \text{exactly one of } \gamma_k(x_t) \text{ is equal to } 1$

$$\hat{\gamma}_{\hat{k}}(x_t) := \begin{cases} 1 & \text{if } \hat{k} = \arg_k \min \|x_t - \mu_k\|^2 \\ 0 & \text{otherwise} \end{cases}$$



set feasible initial approximation Γ_0

while
$$||L(\Gamma_{it}, \theta_{it}) - L(\Gamma_{it-1}, \theta_{it-1})|| > \varepsilon$$

solve $\theta_{it} = \arg\min_{\theta} L(\theta, \Gamma_{it})$ (with fixed Γ_{it})
solve $\Gamma_{it} = \arg\min_{\Gamma \in \Omega_{\Gamma}} L(\theta_{it}, \Gamma)$ (with fixed θ_{it})
 $it = it + 1$
endwhile

K-means algorithm











- K-means ignores (doesn't take into account) the time

 we want to implement "rationality" of the cluster changes during time (i.e. the time series is not jumping between clusters in crazy way)



caused by the change of cluster (model)

typically caused by the noise

<u>Assumption</u>: smooth $\gamma_k(t) \Rightarrow$ regularisation

• enforce the persistency of underlying hidden switching process

$$\begin{aligned} \|\gamma_k\|_{\mathcal{H}_1} &:= & \sqrt{\sum_{t=1}^{T-1} (\gamma_k(t+1) - \gamma_k(t))^2} &\leq C_k \\ \|\gamma_k\|_{BV} &:= & \sum_{t=1}^{T-1} |\gamma_k(t+1) - \gamma_k(t)| &\leq C_k \end{aligned}$$
 FEM-BV

Examples:



 $[\Theta, \Gamma(t)] = \arg\min_{\Theta, \Gamma} \sum_{t=1}^{T} \sum_{k=1}^{K} \gamma_k(t) \cdot \|x_t - \theta_k\|^2 + \varepsilon^2 \sum_{k=1}^{K} \sum_{t=1}^{T-1} (\gamma_k(t+1) - \gamma_k(t))^2$

From non-stationaty models to clustering

$$[\bar{\Theta}, \bar{\Gamma}] = \arg\min_{\Theta, \Gamma} \sum_{t=1}^{T} \sum_{k=1}^{K} \gamma_k(t) \cdot \|x_t - \mu(t, \theta_k)\|^2 + \varepsilon^2 \sum_{k=1}^{K} \sum_{t=1}^{T-1} (\gamma_k(t+1) - \gamma_k(t))^2$$

s.t. $\gamma_k(t) \in \{0, 1\}, \ \forall t : \sum_{k=1}^{K} \gamma_k(t) = 1$

<u>Assumption</u>: continuous real-valued functions $\gamma_k(t) \in [0, 1]$ (probabilities)

$$0 \le \gamma_k(t) \le 1, \quad \forall t : \sum_{k=1}^K \gamma_k(t) = 1$$

• *Horenko I.:* Finite Element Approach to Clustering of Multidimensional Time Series, *SIAM J. Sci. Comp. 32(1), 62-83 (2010)*

Optimization: solving the problem

$$\begin{split} [\bar{\Theta},\bar{\Gamma}] &= \arg\min_{\Theta,\Gamma} \sum_{t=1}^{T} \sum_{k=1}^{K} \gamma_{k}(t) \cdot \|x_{t} - \mu(t,\theta_{k})\|^{2} + \varepsilon^{2} \sum_{k=1}^{K} \sum_{t=1}^{T-1} (\gamma_{k}(t+1) - \gamma_{k}(t))^{2} \\ &=: L(\Theta,\Gamma) \\ \text{s.t. } 0 \leq \gamma_{k}(t) \leq 1, \ \forall t : \sum_{k=1}^{K} \gamma_{k}(t) = 1 \\ \end{split}$$

$$\begin{aligned} \text{while } \|L(\Theta_{it},\Gamma_{it}) - L(\Theta_{it-1},\Gamma_{it-1})\| > \varepsilon_{L} \\ \text{solve } \Theta_{it} = \arg\min_{\Theta} L(\Theta,\Gamma_{it}) \quad (with \ fixed \ \Gamma_{it}) \leq 1, \ \forall t : \sum_{k=1}^{K} \gamma_{k}(t) = 1 \\ \text{solve } \Theta_{it} = \arg\min_{\Theta} L(\Theta_{it},\Gamma) \quad (with \ fixed \ \Theta_{it}) \leq 1, \ \forall t : \sum_{k=1}^{K} \gamma_{k}(t) = 1 \\ \text{solve } \Theta_{it} = \arg\min_{\Theta} L(\Theta_{it},\Gamma) \quad (with \ fixed \ \Theta_{it}) \leq 1, \ \forall t : \sum_{k=1}^{K} \gamma_{k}(t) = 1 \\ \text{solve } \Theta_{it} = \arg\min_{\Theta} L(\Theta_{it},\Gamma) \quad (with \ fixed \ \Theta_{it}) \leq 1, \ \forall t : \sum_{k=1}^{K} \gamma_{k}(t) = 1 \\ \text{solve } \Theta_{it} = \arg\min_{\Theta} L(\Theta_{it},\Gamma) \quad (with \ fixed \ \Theta_{it}) \leq 1, \ \forall t : \sum_{k=1}^{K} \gamma_{k}(t) = 1 \\ \text{solve } \Theta_{it} = \arg\min_{\Theta} L(\Theta_{it},\Gamma) \quad (with \ fixed \ \Theta_{it}) \leq 1, \ \forall t : \sum_{k=1}^{K} \gamma_{k}(t) = 1 \\ \text{solve } \Theta_{it} = \arg\min_{\Theta} L(\Theta_{it},\Gamma) \quad (with \ fixed \ \Theta_{it}) \leq 1, \ \forall t : \sum_{k=1}^{K} \gamma_{k}(t) = 1 \\ \text{solve } \Theta_{it} = \arg\min_{\Theta} L(\Theta_{it},\Gamma) \quad (with \ fixed \ \Theta_{it}) \leq 1, \ \forall t : \sum_{k=1}^{K} \gamma_{k}(t) = 1 \\ \text{solve } \Theta_{it} = \arg\min_{\Theta} L(\Theta_{it},\Gamma) \quad (with \ fixed \ \Theta_{it}) \leq 1, \ \forall t : \sum_{k=1}^{K} \gamma_{k}(t) = 1 \\ \text{solve } \Theta_{it} = \arg\min_{\Theta} L(\Theta_{it},\Gamma) \quad (with \ fixed \ \Theta_{it}) \leq 1, \ \forall t \in 1$$

Spectral Projected Gradient method for QP

Given cost function $f : \mathbb{R}^n \to \mathbb{R}$, initial approximation $x^0 \in \Omega$, projection onto feasible set $P_{\Omega}(x)$, parameters $m \in \mathbb{N}, \tau \in (0, 1)$, safeguarding parameters $\sigma_1, \sigma_2 \in \mathbb{R} : 0 < \sigma_1 < \sigma_2 < 1$, precision $\varepsilon > 0$, and initial step-size $\alpha_0 > 0$.

k := 0 $a^0 := Ax^0 - b$ $f^0 := 1/2 \langle q^0 - b, x^0 \rangle$ for $k = 0, 1, \ldots$ $d^k := P_{\Omega}(x^k - \alpha_k q^k) - x^k$ compute matrix-vector multiplication Ad^k compute multiple dot-product $\langle d^k, \{d^k, Ad^k, g^k\} \rangle$ if $\sqrt{\langle d^k, d^k \rangle} < \varepsilon$ then stop. $f_{\max} := \max\{f(x^{k-j}) : 0 \le j \le \min\{k, m-1\}\} \\ \xi := (f_{\max} - f^k) / \langle d^k, Ad^k \rangle$ $\hat{\bar{\beta}} := -\langle g^k, d^k \rangle / \langle d^k, A d^k \rangle \\ \hat{\beta} := \tau \bar{\beta} + \sqrt{\tau^2 \bar{\beta}^2 + 2\xi}$ choose $\beta_k \in \langle \sigma_1, \min\{\sigma_2, \hat{\beta}\} \rangle$ $x^{k+1} := x^k + \beta_k d^k$ $q^{k+1} := q^k + \beta_k A d^k$ $f^{k+1} := f^k + \beta_k \langle d^k, q^k \rangle + \frac{1}{2} \beta_k^2 \langle d^k, Ad^k \rangle$ $\alpha_{k+1} := \langle d^k, d^k \rangle / \langle d^k, A d^k \rangle$ k := k + 1endwhile



Return approximation of solution x^k .

• Birgin E.G., Martínez J.M., Raydan M.: Nonmonotone spectral projected gradient methods on convex sets, (2000)

• Birgin E.G., Raydan M., Martínez J.M.: Spectral Projected Gradient Methods: Review and Perspectives, (2014)

• *Pospíšil L.:* Development of Algorithms for Solving Minimizing Problems with Convex Quadratic Function on Special Convex Sets and Applications, *PhD thesis, supervised by Z. Dostál (2015)*

Feasible set – separable simplexes

$$0 \le \gamma_k(t) \le 1, \quad \forall t : \sum_{k=1}^K \gamma_k(t) = 1$$



T separable simplexes (hypertriangles) = fully parallel projection



Regularization parameter

$$[\bar{\Theta}, \bar{\Gamma}] = \arg\min_{\Theta, \Gamma} \sum_{t=1}^{T} \sum_{k=1}^{K} \gamma_k(t) \cdot \|x_t - \mu(t, \theta_k)\|^2 + \varepsilon^2 \sum_{k=1}^{K} \sum_{t=1}^{T-1} (\gamma_k(t+1) - \gamma_k(t))^2$$
$$\varepsilon^2 = 10$$





 $\varepsilon^2 = 500$





Regularization parameter: error curve

$$[\bar{\Theta}, \bar{\Gamma}] = \arg\min_{\Theta, \Gamma} \sum_{t=1}^{T} \sum_{k=1}^{K} \gamma_k(t) \cdot \|x_t - \mu(t, \theta_k)\|^2 + \varepsilon^2 \sum_{k=1}^{K} \sum_{t=1}^{T-1} (\gamma_k(t+1) - \gamma_k(t))^2$$



Regularization parameter: error curve





Regularization parameter: L-curve



• Hansen P.C. and D. P. O'Leary D.P.: The use of the L-curve in the regularization of discrete ill-posed problems, SIAM J. Sci. Comp. 14, (1993)

Regularization parameter: L-curve



Possible extension: spatial regularization



• Gerber S., Horenko I.: Improving Clustering by Imposing Network Information, Sciences Advances (AAAS), 1(7):e1500163, (2015)

Possible extension: general non-stationary models

$$[\bar{\Theta}, \bar{\Gamma}] = \arg\min_{\Theta, \Gamma} \sum_{t=1}^{T} \sum_{k=1}^{K} \gamma_k(t) \cdot \|x_t - \mu(t, \theta_k)\|^2 + \varepsilon^2 \sum_{k=1}^{K} \sum_{t=1}^{T-1} (\gamma_k(t+1) - \gamma_k(t))^2$$

Piecewise linear regression

$$x(t) = \begin{cases} a_{0,1} + a_{1,1}t + \varepsilon(t) & \text{for } x \in \mathcal{T}_1 \\ a_{0,2} + a_{1,2}t + \varepsilon(t) & \text{for } x \in \mathcal{T}_2 \\ a_{0,3} + a_{1,3}t + \varepsilon(t) & \text{for } x \in \mathcal{T}_3 \end{cases}$$

$$error = \|x(t) - \hat{x}(t)\| = 307.1922$$

$$AIC_K = 2 \cdot 6 - 2 \cdot \log(307.1922) = 0.545$$

To compare: Polynomial regression

$$x(t) = a_0 + a_1 t + a_2 t^2 + ... a_P t^P + \varepsilon(t)$$

 $error = ||x(t) - \hat{x}(t)|| = 309.4182$
 $AIC_P = 2 \cdot 51 - 2 \cdot \log(309.4182) = 90.5306$



t

100 200 300 400 500 600 700 800 900 1000

Possible extension: general non-stationary models

$$[\bar{\Theta},\bar{\Gamma}] = \arg\min_{\Theta,\Gamma} \sum_{t=1}^{T} \sum_{k=1}^{K} \gamma_k(t) \cdot \|x_t - \mu(t,\theta_k)\|^2 + \varepsilon^2 \sum_{k=1}^{K} \sum_{t=1}^{T-1} (\gamma_k(t+1) - \gamma_k(t))^2$$

20/05/99

11/05/01

13/05/03

TV-Entropy (regime-based non-stationary Entropy)

$$\max_{f_t(x),\,\forall t} \left\{ \mathbb{E}\left[H\left[f_t(x)\right] \right] = \mathbb{E}\left[-\int_{\mathcal{X}} f_t(x) \ln f_t(x) \, dx \right] \right\}$$

s.t.
$$\int_{\mathcal{X}} x^j f_t(x) \, dx = \mu_j(t) \quad \forall j \in \{0,\dots,k\},$$

TV-Entropy identifies memoryless models that are simpler and better then the state-ofthe-art for all of the financial benchmark data considered.

 Marchenko G., Gagliardini P., Horenko I.: Towards a Computationally Tractable Maximum Entropy Principle for Nonstationary Financial Time Series, SIAM J. Finan. Math., 9(4), 1249–1285, (2018)



06/05/05

03/05/07

29/04/09

25/04/11





K-means?

$$[S^*, \Gamma^*] = \arg\min_{S, \Gamma} \sum_{t=1}^T \sum_{k=1}^K \Gamma_{k, t} ||x_t - S_k||_2^2$$

 $\Gamma_{k,t}$ - probability for the system to be in particular state S_k at the instance t(K-means: = probability of x_t to belong to k-th cluster)



K-means?

$$[S^*, \Gamma^*] = \arg\min_{S, \Gamma} \sum_{t=1}^T \sum_{k=1}^K \Gamma_{k, t} ||x_t - S_k||_2^2$$

 $\Gamma_{k,t} \text{ - probability for the system to be in particular state } S_k \text{ at the instance } t$ $(K\text{-means: = probability of } x_t \text{ to belong to } k\text{-th cluster}) - \bigcup_{k=1}^{K} (K\text{-means: = probability of observing error } \|x_t - S_k\|_2^2)$ $x_t = \sum_{k=1}^{K} \Gamma_{k,t} S_k$



<u>K-means?</u> $[S^*, \Gamma^*] = \arg\min_{S, \Gamma} \sum_{t=1}^T \sum_{k=1}^K \Gamma_{k,t} ||x_t - S_k||_2^2$

$$\Gamma_{k,t} - \text{probability for the system to be in particular state } S_k \text{ at the instance } t$$

$$(K-\text{means:} = \text{probability of } x_t \text{ to belong to } k-\text{th cluster}) - (K \text{ means:} = \text{probability of observing error } \|x_t - S_k\|_2^2) - (K \text{ means:} = \text{probability of observing error } \|x_t - S_k\|_2^2) - (K \text{ means:} = \text{probability of observing error } \|x_t - S_k\|_2^2) - (K \text{ means:} = \text{probability of observing error } \|x_t - S_k\|_2^2) - (K \text{ means:} = \text{probability of observing error } \|x_t - S_k\|_2^2) - (K \text{ means:} = \text{probability of observing error } \|x_t - S_k\|_2^2) - (K \text{ means:} = \text{probability of observing error } \|x_t - S_k\|_2^2) - (K \text{ means:} = \text{probability of observing error } \|x_t - S_k\|_2^2) - (K \text{ means:} = \text{probability of observing error } \|x_t - S_k\|_2^2) - (K \text{ means:} = \text{probability of observing error } \|x_t - S_k\|_2^2) - (K \text{ means:} = \text{probability of observing error } \|x_t - S_k\|_2^2) - (K \text{ means:} = \text{probability of observing error } \|x_t - S_k\|_2^2) - (K \text{ means:} = \text{probability of observing error } \|x_t - S_k\|_2^2) - (K \text{ means:} = \text{means:} = \text{probability of observing error } \|x_t - S_k\|_2^2) - (K \text{ means:} = \text{means:} = \text$$

$$\begin{split} [S^*, \Gamma^*] &= \arg\min_{S, \Gamma} \sum_{t=1}^T \|x_t - \sum_{k=1}^K \Gamma_{k,t} S_k\|_2^2 \qquad \text{s.t.} \quad \forall t : \sum_{k=1}^K \Gamma_{k,t} = 1, \ \forall k : \Gamma_{k,t} \in [0, 1] \\ \\ \text{Jensen inequality:} &\leq \sum_{t=1}^T \sum_{k=1}^K \Gamma_{k,t} \|x_t - S_k\|_2^2 \\ \\ \Gamma_{k,t} \in \{0, 1\} : \qquad = \sum_{t=1}^T \sum_{k=1}^K \Gamma_{k,t} \|x_t - S_k\|_2^2 \end{split}$$
(K-means is suboptimal)

$$[S^*, \Gamma^*] = \arg\min_{S, \Gamma} \sum_{t=1}^T \|x_t - \sum_{k=1}^K \Gamma_{k, t} S_k\|_2^2$$

s.t.
$$\forall t : \sum_{k=1}^{K} \Gamma_{k,t} = 1, \ \forall k : \Gamma_{k,t} \in [0,1]$$

General SPA:

$$[S^*, \Gamma^*] := \arg\min_{\Gamma \in \Omega_{\Gamma}} L(S, \Gamma)$$
$$L(S, \Gamma) := \sum_{t=1}^T \operatorname{dist}_S(X(t), \Gamma(t)) + \varepsilon_S^2 \Phi_S(S) + \varepsilon_{\Gamma}^2 \Phi_{\Gamma}(\Gamma)$$

Set a feasible initial approximation $\Gamma^0 \in \Omega_{\Gamma}$ while $||L(S^k, \Gamma^k) - L(S^{k-1}, \Gamma^{k-1})|| \ge \varepsilon$ solve $S^k = \arg\min_S L(S, \Gamma^{k-1})$ (with fixed Γ^{k-1}) solve $\Gamma^k = \arg\min_{\Gamma \in \Omega_{\Gamma}} L(S^k, \Gamma)$ (with fixed S^k) k = k + 1endwhile

Return an approximation of the data representation vectors S^k and an approximation of cluster affiliation probability vectors Γ^k .



$$[S^*, \Gamma^*] = \arg\min_{S, \Gamma} \sum_{t=1}^T \|x_t - \sum_{k=1}^K \Gamma_{k, t} S_k\|_2^2$$
$$= \|X - S\Gamma\|_F^2$$

s.t.
$$\forall t : \sum_{k=1}^{K} \Gamma_{k,t} = 1, \ \forall k : \Gamma_{k,t} \in [0,1]$$

 $S^* = X\Gamma^T \left(\Gamma\Gamma^T\right)^+ + \alpha^T R^T, \text{ with parameter } \alpha \in \mathbb{R}^{r,n},$

Introducing regularization:

SPA in action

Exact Law of the Total Probability in a matrix notation:



SPA in action

 $\pi_{\rm Y}(t)$ $\pi_{x}(t)$ Λ $P[Y^{t}=y_1]$ $P[Y^t = y_1 | X^t = x_n]$ $P[X^{t}=x_1]$ $P[Y^t = y_1 | X^t = x_1]$ $P[Y^t=y_1|X^t=x_2]$ $P[Y^{t=y_2}]$ $P[Y^t = y_2 | X^t = x_1]$ $P[Y^t = y_2 | X^t = x_n]$ $P[X^{t}=x_{2}]$ $P[Y^t=y_2|X^t=x_2]$. . . = ... ••• $P[Y^{t}=y_{m}]$ $P[Y^t=y_m|X^t=x_1] \quad P[Y^t=y_m|X^t=x_2]$ $P[Y^t=y_m|X^t=x_n]$ $P[X^{t=x_n}]$. . .

 $\Gamma_{k,t}$ - probability for the system to be in particular state S_k at the instance t

$$\Gamma_{k,t} := P(x_t = S_k)$$

$$\forall t: \Gamma_{:,t}^Y = \Lambda \Gamma_{:,t}^X$$

Markov model: Y(t) := X(t+1)

SPA in action

 $\pi_{x}(t)$ $\pi_{\rm Y}(t)$ Λ $P[X^{t}=x_1]$ $P[Y^{t}=y_1]$ $P[Y^t=y_1|X^t=x_1]$ $P[Y^t=y_1|X^t=x_2]$ $P[Y^t=y_1|X^t=x_n]$ $P[Y^{t}=y_2]$ $P[X^{t}=x_{2}]$ $P[Y^t=y_2|X^t=x_1]$ $P[Y^t=y_2|X^t=x_2]$ $P[Y^t=y_2|X^t=x_n]$. . . = ••• . . . $P[Y^{t}=y_{m}]$ $P[Y^t=y_m|X^t=x_1]$ $P[Y^t=y_m|X^t=x_2]$ $P[Y^t=y_m|X^t=x_n]$ $P[X^{t}=x_{n}]$. . .

 $\Gamma_{k,t}$ - probability for the system to be in particular state S_k at the instance t

$$\Gamma_{k,t} := P(x_t = S_k)$$

Markov model: Y(t) := X(t+1)



Exact Law of the Total Probability in a matrix notation:

Results

- X(t) is continuous (and real valued) set of collected 32 image features
- index t denotes patients and goes from 1 to 569
- Y(t) is binary: 'benign' or 'malignant'

- X(t) contains genetic expression levels for 25'000 genes
- index t goes from 1 to 300 (there are 300 single cell probes),
- Y(t) is a label denoting one of the 11 cell types (e.g., 'blood cell', 'glia 25 cell', etc.)

a) breast cancer diagnostics (569 patients, 32 features, WDBC data)

b) single cell human mRNA classification (25`000 genes, 11 cell types, 300 cell samples)



a) Lorenz-96 1D turbulence model (weakly-chaotic regime)

b) Lorenz-96 1D turbulence model (strongly-chaotic regime)



c) surface temperature dynamics over Europe (1979-2010, 20x30 grid ECMWF resimulation data)



SPA in action: compressing the pipeline



 $\Gamma_{k,t}$ - probability for the system to be in particular state S_k at the instance t





SPA in action: compressing the pipeline

$$[S_{x}^{*}, \Gamma_{x}^{*}, S_{y}^{*}, \Gamma_{y}^{*}] = \arg \min_{\Gamma_{x}, \Gamma_{y} \in \Omega_{\Gamma}} ||Y - S_{y}\Gamma_{y}||_{F}^{2} + \varepsilon ||X - S_{x}\Gamma_{x}||_{F}^{2}$$

$$= \left\| \begin{bmatrix} Y \\ \varepsilon X \end{bmatrix} - \begin{bmatrix} S_{y}\Gamma_{y} \\ \varepsilon S_{x}\Gamma_{x} \end{bmatrix} \right\|_{F}^{2} = \left\| \begin{bmatrix} Y \\ \varepsilon X \end{bmatrix} - \begin{bmatrix} S_{y}\Lambda \\ \varepsilon S_{x} \end{bmatrix} \Gamma_{x} \right\|_{F}^{2}$$

$$= \left\| \hat{X} - \hat{S}\hat{\Gamma} \right\|_{F}^{2}$$

$$Y(t) \in \mathbb{R}^{m}, t = 1, \dots, T$$

$$S_{Y}^{K} \in \mathbb{R}^{m}, k = 1, \dots, K_{Y}$$

$$||Y - S^{Y}\Gamma^{Y}||_{F} \rightarrow \min$$

$$\forall t : \Gamma_{:,t}^{Y} = \Lambda\Gamma_{:,t}^{X}$$

$$X(t) \in \mathbb{R}^{n}, t = 1, \dots, T$$

$$S_{k}^{X} \in \mathbb{R}^{n}, k = 1, \dots, K_{X}$$

$$||Y - S^{Y}\Gamma^{Y}||_{F} \rightarrow \min$$

$$A \text{ Common Machine Learning Pipeline:}$$

$$discretization$$

$$(dimension reduction) \longrightarrow model training (and validation)$$

- Horenko I.: Finite Element Approach to Clustering of Multidimensional Time Series, SIAM J. Sci. Comp. 32(1), 62-83 (2010)
- Metzner P., Putzig L., Horenko I. : Analysis of persistent non-stationary time series and applications. Communications in Applied Mathematics and Computational Science, 7(2):175-229, (2012)
- Gerber S., Horenko I.: Improving Clustering by Imposing Network Information, Sciences Advances (AAAS), 1(7):e1500163, (2015)
- Pospisil L., Gagliardini P., Sawyer W., Horenko I.: On a scalable nonparametric denoising of time series signals. Communications in Applied Mathematics and Computational Science, 13(1), (2018)
- Marchenko G., Gagliardini P., Horenko I.: Towards a Computationally Tractable Maximum Entropy Principle for Nonstationary Financial Time Series, SIAM J. Finan. Math., 9(4), 1249–1285, (2018)
- Gerber S., Pospíšil L., Navandar M., Horenko I.: *Low-cost scalable discretization, prediction and feature selection for complex systems*, almost published in Science Advances, (2019), http://www.biorxiv.org/content/10.1101/720441v1

... thank you for your attention!