



# Time Series Forecasting Using Kernel Regression Methods

**José Daniel de Alencar Santos (IFCE)**  
**Allan Kelvin Mendes de Sales (IFCE)**  
**Guilherme de Alencar Barreto (UFC)**

Graduate Program in Teleinformatics Engineering (PPGETI)  
Federal University of Ceará (UFC), Fortaleza, Ceará, Brazil  
jdaniel@ifce.edu.br, kelvin@ifce.edu.br, gbarreto@ufc.br

Fortaleza, October 2022

# Outline

## ① Introduction

Related Publications

Kernel-Based Regression Models

General and Specific Objectives

## ② The LSSVR Model

## ③ The KOLS Model

## ④ Computational Experiments

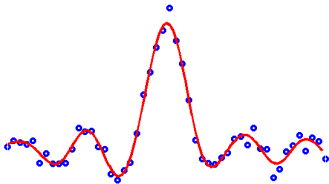
## ⑤ Conclusions

## Related Publications

- J. D. A. Santos & G. A. Barreto (2018). Novel sparse LSSVR models in primal weight space for robust system identification with outliers. **Journal of Process Control**, vol. 67, p. 129-140.
- J. D. A. Santos & G. A. Barreto (2017). An outlier-robust kernel RLS algorithm for nonlinear system identification. **Nonlinear Dynamics**, vol. 90, p. 1707-1726.
- A. K. M. Sales & G. A. Barreto (2022). A novel method for sparsification of kernel adaptive filters. **IEEE Transactions on Neural Networks and Learning Systems**, submitted.

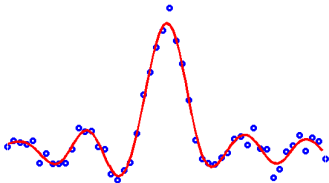
# Kernel-Based Regression Models

- In several application domains, a nonlinear model has to be fitted to data samples generated by nonlinear phenomena.
- This a very demanding task because there are many alternatives out there, such as neural networks, polinomial model, Gaussian process (GP) models, and so on.
- A natural (or would it be unnatural?) alternative involves kernel based models, such as support vector regression.



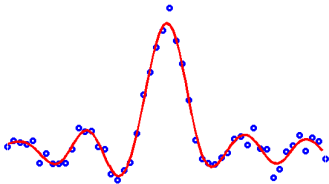
# Kernel-Based Regression Models

- In several application domains, a nonlinear model has to be fitted to data samples generated by nonlinear phenomena.
- This a very demanding task because there are many alternatives out there, such as neural networks, polinomial model, Gaussian process (GP) models, and so on.
- A natural (or would it be unnatural?) alternative involves kernel based models, such as support vector regression.

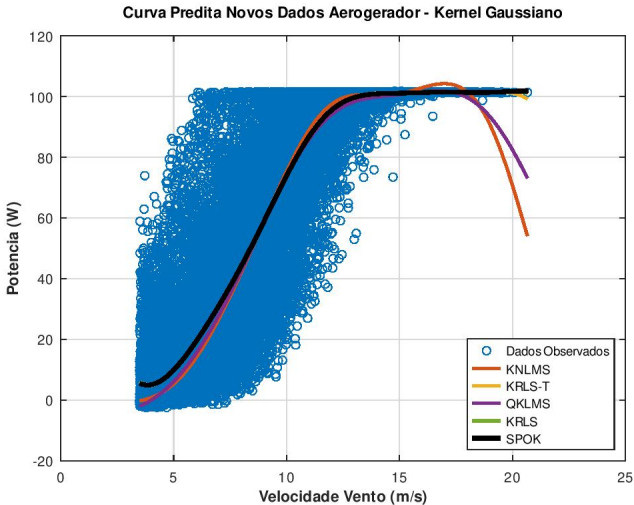


# Kernel-Based Regression Models

- In several application domains, a nonlinear model has to be fitted to data samples generated by nonlinear phenomena.
- This a very demanding task because there are many alternatives out there, such as neural networks, polinomial model, Gaussian process (GP) models, and so on.
- A natural (or would it be unnatural?) alternative involves kernel based models, such as support vector regression.

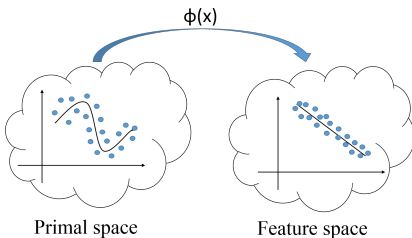


# Kernel-Based Regression Models



# Kernel-Based Regression Models

- The least squares support vector regression (**LSSVR**) and the kernel ordinary least squares (KOLS) <sup>1</sup> models are very competitive tools for nonlinear regression tasks.
- It is easy to understand their backgrounds from the basic concepts of linear regression and least squares estimation.
- Kernel models build a linear model in the feature (rkhs) space using an unknown nonlinear mapping  $\phi$ .

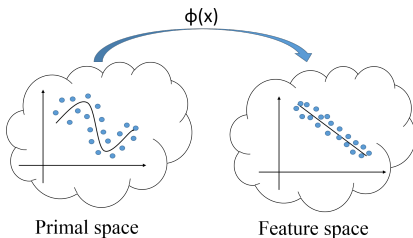


<sup>1</sup>Suykens, J. A. K. et al. *Least Squares Support Vector Machines*. New Jersey, USA: World Scientific, 2002.



# Kernel-Based Regression Models

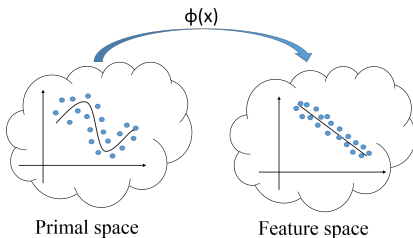
- The least squares support vector regression (**LSSVR**) and the kernel ordinary least squares (KOLS)<sup>1</sup> models are very competitive tools for nonlinear regression tasks.
- It is easy to understand their backgrounds from the basic concepts of linear regression and least squares estimation.
- Kernel models build a linear model in the feature (rkhs) space using an unknown nonlinear mapping  $\phi$ .



<sup>1</sup>Suykens, J. A. K. et al. *Least Squares Support Vector Machines*. New Jersey, USA: World Scientific, 2002.

# Kernel-Based Regression Models

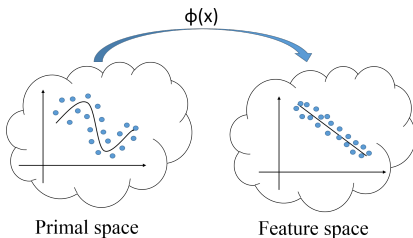
- The least squares support vector regression (**LSSVR**) and the kernel ordinary least squares (KOLS)<sup>1</sup> models are very competitive tools for nonlinear regression tasks.
- It is easy to understand their backgrounds from the basic concepts of linear regression and least squares estimation.
- Kernel models build a linear model in the feature (rkhs) space using an unknown nonlinear mapping  $\phi$ .



<sup>1</sup>Suykens, J. A. K. et al. *Least Squares Support Vector Machines*. New Jersey, USA: World Scientific, 2002.

# Kernel-Based Regression Models

- The least squares support vector regression (**LSSVR**) and the kernel ordinary least squares (KOLS) <sup>1</sup> models are very competitive tools for nonlinear regression tasks.
- It is easy to understand their backgrounds from the basic concepts of linear regression and least squares estimation.
- Kernel models build a linear model in the feature (rkhs) space using an unknown nonlinear mapping  $\phi$ .



<sup>1</sup>Suykens, J. A. K. et al. *Least Squares Support Vector Machines*. New Jersey, USA: World Scientific, 2002.

# General and Specific Objectives

## General Objective

The overall objective of this talk is to introduce kernel-based nonlinear regression models and their applications to time series forecasting.

## Specific Objectives

- 1 To describe the basics of the LSSVR model.
- 2 To describe the basics of the kernel ordinary least squares (KOLS).
- 3 To describe the basics of the kernel regularized least squares (KRELS).
- 4 To discuss Octave/Matlab codes of the aforementioned models.
- 5 To present some applications on time series forecasting.

# General and Specific Objectives

## General Objective

The overall objective of this talk is to introduce kernel-based nonlinear regression models and their applications to time series forecasting.

## Specific Objectives

- ① To describe the basics of the LSSVR model.
- ② To describe the basics of the kernel ordinary least squares (KOLS).
- ③ To describe the basics of the kernel regularized least squares (KRELS).
- ④ To discuss Octave/Matlab codes of the aforementioned models.
- ⑤ To present some applications on time series forecasting.

# General and Specific Objectives

## General Objective

The overall objective of this talk is to introduce kernel-based nonlinear regression models and their applications to time series forecasting.

## Specific Objectives

- 1 To describe the basics of the LSSVR model.
- 2 To describe the basics of the kernel ordinary least squares (KOLS).
- 3 To describe the basics of the kernel regularized least squares (KRELS).
- 4 To discuss Octave/Matlab codes of the aforementioned models.
- 5 To present some applications on time series forecasting.

# General and Specific Objectives

## General Objective

The overall objective of this talk is to introduce kernel-based nonlinear regression models and their applications to time series forecasting.

## Specific Objectives

- 1 To describe the basics of the LSSVR model.
- 2 To describe the basics of the kernel ordinary least squares (KOLS).
- 3 To describe the basics of the kernel regularized least squares (KRELS).
- 4 To discuss Octave/Matlab codes of the aforementioned models.
- 5 To present some applications on time series forecasting.

# General and Specific Objectives

## General Objective

The overall objective of this talk is to introduce kernel-based nonlinear regression models and their applications to time series forecasting.

## Specific Objectives

- 1 To describe the basics of the LSSVR model.
- 2 To describe the basics of the kernel ordinary least squares (KOLS).
- 3 To describe the basics of the kernel regularized least squares (KRELS).
- 4 To discuss Octave/Matlab codes of the aforementioned models.
- 5 To present some applications on time series forecasting.



# General and Specific Objectives

## General Objective

The overall objective of this talk is to introduce kernel-based nonlinear regression models and their applications to time series forecasting.

## Specific Objectives

- 1 To describe the basics of the LSSVR model.
- 2 To describe the basics of the kernel ordinary least squares (KOLS).
- 3 To describe the basics of the kernel regularized least squares (KRELS).
- 4 To discuss Octave/Matlab codes of the aforementioned models.
- 5 To present some applications on time series forecasting.

# General and Specific Objectives

## General Objective

The overall objective of this talk is to introduce kernel-based nonlinear regression models and their applications to time series forecasting.

## Specific Objectives

- 1 To describe the basics of the LSSVR model.
- 2 To describe the basics of the kernel ordinary least squares (KOLS).
- 3 To describe the basics of the kernel regularized least squares (KRELS).
- 4 To discuss Octave/Matlab codes of the aforementioned models.
- 5 To present some applications on time series forecasting.

# Outline

## ① Introduction

Related Publications

Kernel-Based Regression Models

General and Specific Objectives

## ② The LSSVR Model

## ③ The KOLS Model

## ④ Computational Experiments

## ⑤ Conclusions

# LSSVR Primal Problem

- Given an estimation (a.k.a. training) dataset  $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$ , with  $\mathbf{x}_n \in \mathbb{R}^d$  and  $y_n \in \mathbb{R}$ , the kernel regression problem is given by

$$f(\mathbf{x}_n) = \mathbf{w}^\top \phi(\mathbf{x}_n) + b, \quad (1)$$

where  $\mathbf{w} \in \mathbb{R}^{d_h}$  is the unknown parameter vector,  $b \in \mathbb{R}$  is a bias and  $\phi(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^{d_h}$  is a nonlinear map into the feature space.

- The primal optimization problem of the LSSVR model is given by

$$\min_{\mathbf{w}, b, \mathbf{e}} J_p(\mathbf{w}, \mathbf{e}) = \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{smoothness}} + \underbrace{\gamma \frac{1}{2} \sum_{n=1}^N e_n^2}_{\text{training errors}}, \quad (2)$$

$$\text{subject to } \{ y_n = \mathbf{w}^\top \phi(\mathbf{x}_n) + b + e_n, \text{ for } n = 1, \dots, N, \quad (3)$$

where  $\gamma$  is the regularization parameter and  $e_n$  is the  $n$ -th error.

# LSSVR Primal Problem

- Given an estimation (a.k.a. training) dataset  $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$ , with  $\mathbf{x}_n \in \mathbb{R}^d$  and  $y_n \in \mathbb{R}$ , the kernel regression problem is given by

$$f(\mathbf{x}_n) = \mathbf{w}^\top \phi(\mathbf{x}_n) + b, \quad (1)$$

where  $\mathbf{w} \in \mathbb{R}^{d_h}$  is the unknown parameter vector,  $b \in \mathbb{R}$  is a bias and  $\phi(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^{d_h}$  is a nonlinear map into the feature space.

- The primal optimization problem of the LSSVR model is given by

$$\min_{\mathbf{w}, b, \mathbf{e}} J_p(\mathbf{w}, \mathbf{e}) = \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{smoothness}} + \underbrace{\gamma \frac{1}{2} \sum_{n=1}^N e_n^2}_{\text{training errors}}, \quad (2)$$

$$\text{subject to } \{ y_n = \mathbf{w}^\top \phi(\mathbf{x}_n) + b + e_n, \text{ for } n = 1, \dots, N, \quad (3)$$

where  $\gamma$  is the regularization parameter and  $e_n$  is the  $n$ -th error.

# LSSVR Dual Problem

- The LSSVR dual problem is obtained by building the Lagrangian, as

$$\mathcal{L}(\mathbf{w}, b, \mathbf{e}, \boldsymbol{\alpha}) := \frac{1}{2} \|\mathbf{w}\|^2 + \gamma \frac{1}{2} \sum_{n=1}^N e_n^2 - \sum_{n=1}^N \alpha_n (\mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n) + b + e_n - y_n), \quad (4)$$

where  $\alpha_n$  are the Lagrange multipliers.

- From the optimality conditions<sup>2</sup>, one gets

$$\underbrace{\begin{bmatrix} 0 & \mathbf{1}_N^\top \\ \mathbf{1}_N & \mathbf{K} + \gamma^{-1} \mathbf{I}_N \end{bmatrix}}_{\boldsymbol{\Omega}} \underbrace{\begin{bmatrix} b \\ \boldsymbol{\alpha} \end{bmatrix}}_{\boldsymbol{\alpha}_o} = \underbrace{\begin{bmatrix} 0 \\ \mathbf{y} \end{bmatrix}}_{\mathbf{y}_o}. \quad (5)$$

- $\mathbf{K} \in \mathbb{R}^{N \times N}$  is the **kernel matrix**, whose entries are

$$K_{i,j} := k(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\phi}(\mathbf{x}_i)^\top \boldsymbol{\phi}(\mathbf{x}_j), \text{ for } i, j = 1, \dots, N, \quad (6)$$

where  $k(\cdot, \cdot)$  is the chosen kernel function.

<sup>2</sup>FLETCHER, R. Practical methods of optimization. New Jersey, USA: John Wiley & Sons, 2013.

# LSSVR Dual Problem

- The LSSVR dual problem is obtained by building the Lagrangian, as

$$\mathcal{L}(\mathbf{w}, b, \mathbf{e}, \boldsymbol{\alpha}) := \frac{1}{2} \|\mathbf{w}\|^2 + \gamma \frac{1}{2} \sum_{n=1}^N e_n^2 - \sum_{n=1}^N \alpha_n (\mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n) + b + e_n - y_n), \quad (4)$$

where  $\alpha_n$  are the Lagrange multipliers.

- From the optimality conditions<sup>2</sup>, one gets

$$\underbrace{\begin{bmatrix} 0 & \mathbf{1}_N^\top \\ \mathbf{1}_N & \mathbf{K} + \gamma^{-1} \mathbf{I}_N \end{bmatrix}}_{\boldsymbol{\Omega}} \underbrace{\begin{bmatrix} b \\ \boldsymbol{\alpha} \end{bmatrix}}_{\boldsymbol{\alpha}_o} = \underbrace{\begin{bmatrix} 0 \\ \mathbf{y} \end{bmatrix}}_{\mathbf{y}_o}. \quad (5)$$

- $\mathbf{K} \in \mathbb{R}^{N \times N}$  is the **kernel matrix**, whose entries are

$$K_{i,j} := k(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\phi}(\mathbf{x}_i)^\top \boldsymbol{\phi}(\mathbf{x}_j), \text{ for } i, j = 1, \dots, N, \quad (6)$$

where  $k(\cdot, \cdot)$  is the chosen kernel function.

<sup>2</sup>FLETCHER, R. **Practical methods of optimization**. New Jersey, USA: John Wiley & Sons, 2013.

# LSSVR Dual Problem

- The LSSVR dual problem is obtained by building the Lagrangian, as

$$\mathcal{L}(\mathbf{w}, b, \mathbf{e}, \boldsymbol{\alpha}) := \frac{1}{2} \|\mathbf{w}\|^2 + \gamma \frac{1}{2} \sum_{n=1}^N e_n^2 - \sum_{n=1}^N \alpha_n (\mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n) + b + e_n - y_n), \quad (4)$$

where  $\alpha_n$  are the Lagrange multipliers.

- From the optimality conditions<sup>2</sup>, one gets

$$\underbrace{\begin{bmatrix} 0 & \mathbf{1}_N^\top \\ \mathbf{1}_N & \mathbf{K} + \gamma^{-1} \mathbf{I}_N \end{bmatrix}}_{\boldsymbol{\Omega}} \underbrace{\begin{bmatrix} b \\ \boldsymbol{\alpha} \end{bmatrix}}_{\boldsymbol{\alpha}_o} = \underbrace{\begin{bmatrix} 0 \\ \mathbf{y} \end{bmatrix}}_{\mathbf{y}_o}. \quad (5)$$

- $\mathbf{K} \in \mathbb{R}^{N \times N}$  is the **kernel matrix**, whose entries are

$$K_{i,j} := k(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\phi}(\mathbf{x}_i)^\top \boldsymbol{\phi}(\mathbf{x}_j), \text{ for } i, j = 1, \dots, N, \quad (6)$$

where  $k(\cdot, \cdot)$  is the chosen kernel function.

<sup>2</sup>FLETCHER, R. **Practical methods of optimization**. New Jersey, USA: John Wiley & Sons, 2013.



# The Kernel Matrix

- The kernel matrix  $\mathbf{K} = [K_{i,j}]_{N \times N}$  is defined as

$$\mathbf{K} = \begin{bmatrix} \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}_1) & \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}_N) \\ \phi(\mathbf{x}_2)^\top \phi(\mathbf{x}_1) & \phi(\mathbf{x}_2)^\top \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_2)^\top \phi(\mathbf{x}_N) \\ \vdots & \vdots & \cdots & \vdots \\ \phi(\mathbf{x}_N)^\top \phi(\mathbf{x}_1) & \phi(\mathbf{x}_N)^\top \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_N)^\top \phi(\mathbf{x}_N) \end{bmatrix}_{N \times N} \quad (7)$$

$$= \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \cdots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}_{N \times N} \quad (8)$$

# The Kernel Matrix

- The kernel matrix  $\mathbf{K}$  must be positive-definite.
- It is a Gram matrix; that is, a matrix of dot products.
- It is symmetric. Thus, its computation can be optimized for speed.
- There are lots of kernel functions that can be used.

# The Kernel Matrix

- The kernel matrix  $\mathbf{K}$  must be positive-definite.
- It is a Gram matrix; that is, a matrix of dot products.
- It is symmetric. Thus, its computation can be optimized for speed.
- There are lots of kernel functions that can be used.

# The Kernel Matrix

- The kernel matrix  $\mathbf{K}$  must be positive-definite.
- It is a Gram matrix; that is, a matrix of dot products.
- It is symmetric. Thus, its computation can be optimized for speed.
- There are lots of kernel functions that can be used.

# The Kernel Matrix

- The kernel matrix  $\mathbf{K}$  must be positive-definite.
- It is a Gram matrix; that is, a matrix of dot products.
- It is symmetric. Thus, its computation can be optimized for speed.
- There are lots of kernel functions that can be used.

# The Kernel Matrix

- The linear kernel function:  $k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$ .
- The Gaussian (or RBF) kernel function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\lambda^2}\right), \quad (9)$$

where  $\lambda > 0$  is the width of the function.

- The polynomial kernel function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = (c + \mathbf{x}_i^\top \mathbf{x}_j)^d, \quad (10)$$

where  $c$  is a constant and  $d \geq 1$  is the order of the polynomial.

- The sigmoidal (or MLP) kernel function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \tanh(a\mathbf{x}_i^\top \mathbf{x}_j + r), \quad (11)$$

where  $a$  (slope) and  $r$  (bias) are constants.

# The Kernel Matrix

- The linear kernel function:  $k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$ .
- The Gaussian (or RBF) kernel function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\lambda^2}\right), \quad (9)$$

where  $\lambda > 0$  is the width of the function.

- The polynomial kernel function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = (c + \mathbf{x}_i^\top \mathbf{x}_j)^d, \quad (10)$$

where  $c$  is a constant and  $d \geq 1$  is the order of the polynomial.

- The sigmoidal (or MLP) kernel function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \tanh(a\mathbf{x}_i^\top \mathbf{x}_j + r), \quad (11)$$

where  $a$  (slope) and  $r$  (bias) are constants.

# The Kernel Matrix

- The linear kernel function:  $k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$ .
- The Gaussian (or RBF) kernel function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\lambda^2}\right), \quad (9)$$

where  $\lambda > 0$  is the width of the function.

- The polynomial kernel function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = (c + \mathbf{x}_i^\top \mathbf{x}_j)^d, \quad (10)$$

where  $c$  is a constant and  $d \geq 1$  is the order of the polynomial.

- The sigmoidal (or MLP) kernel function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \tanh(a\mathbf{x}_i^\top \mathbf{x}_j + r), \quad (11)$$

where  $a$  (slope) and  $r$  (bias) are constants.



# The Kernel Matrix

- The linear kernel function:  $k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$ .
- The Gaussian (or RBF) kernel function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\lambda^2}\right), \quad (9)$$

where  $\lambda > 0$  is the width of the function.

- The polynomial kernel function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = (c + \mathbf{x}_i^\top \mathbf{x}_j)^d, \quad (10)$$

where  $c$  is a constant and  $d \geq 1$  is the order of the polynomial.

- The sigmoidal (or MLP) kernel function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \tanh(a\mathbf{x}_i^\top \mathbf{x}_j + r), \quad (11)$$

where  $a$  (slope) and  $r$  (bias) are constants.

# LSSVR Dual Problem

- Kernel-based predictors:

$$\underbrace{f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b,}_{\text{primal space}} \quad (12)$$

or

$$\begin{aligned} f(\mathbf{x}) &= \sum_{n=1}^N \alpha_n k(\mathbf{x}, \mathbf{x}_n) + b, \\ &= \underbrace{\alpha^\top \mathbf{K}(\mathbf{x}, \mathcal{D}) + b.}_{\text{dual space}} \end{aligned} \quad (13)$$

# LSSVR Dual Problem

- Kernel-based predictors:

$$\underbrace{f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b,}_{\text{primal space}} \quad (12)$$

or

$$\begin{aligned} f(\mathbf{x}) &= \sum_{n=1}^N \alpha_n k(\mathbf{x}, \mathbf{x}_n) + b, \\ &= \underbrace{\alpha^\top \mathbf{K}(\mathbf{x}, \mathcal{D}) + b.}_{\text{dual space}} \end{aligned} \quad (13)$$

# LSSVR Dual Problem

- Kernel-based predictors:

$$\underbrace{f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b,}_{\text{primal space}} \quad (12)$$

or

$$\begin{aligned} f(\mathbf{x}) &= \sum_{n=1}^N \alpha_n k(\mathbf{x}, \mathbf{x}_n) + b, \\ &= \underbrace{\boldsymbol{\alpha}^\top \mathbf{K}(\mathbf{x}, \mathcal{D}) + b.}_{\text{dual space}} \end{aligned} \quad (13)$$

# Outline

## ① Introduction

Related Publications

Kernel-Based Regression Models

General and Specific Objectives

## ② The LSSVR Model

## ③ The KOLS Model

## ④ Computational Experiments

## ⑤ Conclusions

# The KOLS Models

- The KRLS model ( $f(\mathbf{x}_i) = \boldsymbol{\phi}^\top(\mathbf{x}_i)\mathbf{w}$ ) minimizes the cost function

$$J(\mathbf{w}) = \sum_{i=1}^N (y_i - f(\mathbf{x}_i))^2 = \|\mathbf{y} - \boldsymbol{\Phi}^\top \mathbf{w}\|^2, \quad (14)$$

which can be rewritten as

$$J(\boldsymbol{\alpha}) = \|\mathbf{y} - \mathbf{K}_t \boldsymbol{\alpha}\|^2. \quad (15)$$

where  $\mathbf{K} = \boldsymbol{\Phi}^\top \boldsymbol{\Phi}$  is the kernel matrix built using the  $N$  training samples.

- Theoretically, the KOLS model solution can be given by  $\boldsymbol{\alpha} = \mathbf{K}^{-1}\mathbf{y}$ , by simply setting  $b = 0$  and  $\gamma \rightarrow \infty$  in Eq. (5).

# The KOLS Models

- The KRLS model ( $f(\mathbf{x}_i) = \boldsymbol{\phi}^\top(\mathbf{x}_i)\mathbf{w}$ ) minimizes the cost function

$$J(\mathbf{w}) = \sum_{i=1}^N (y_i - f(\mathbf{x}_i))^2 = \|\mathbf{y} - \boldsymbol{\Phi}^\top \mathbf{w}\|^2, \quad (14)$$

which can be rewritten as

$$J(\boldsymbol{\alpha}) = \|\mathbf{y} - \mathbf{K}_t \boldsymbol{\alpha}\|^2. \quad (15)$$

where  $\mathbf{K} = \boldsymbol{\Phi}^\top \boldsymbol{\Phi}$  is the kernel matrix built using the  $N$  training samples.

- Theoretically, the KOLS model solution can be given by  $\boldsymbol{\alpha} = \mathbf{K}^{-1}\mathbf{y}$ , by simply setting  $b = 0$  and  $\gamma \rightarrow \infty$  in Eq. (5).

# Sparsifying the KOLS Model

- As mentioned before, the size of the kernel matrix and, hence, the dimension of the vector of parameters  $\alpha$  is determined by the number of training samples  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ .
- To avoid inversion of huge kernel matrices, we must use some kind of sparsification method to select a subset of relevant sample pairs to compose a dictionary.
- In the ALD<sup>3</sup> criterion, when a new incoming sample  $\mathbf{x}_t$  is available, one must verify if  $\phi(\mathbf{x}_t)$  is ALD on the **dictionary**  $\mathcal{D}_{t-1}^{sv} = \{\tilde{\mathbf{x}}_j\}_{j=1}^{m_t-1}$ .
- One should estimate  $\mathbf{a} = [a_1, \dots, a_{m_t-1}]^\top$  satisfying the ALD criterion

$$\delta_t \triangleq \min_{\mathbf{a}} \left\| \sum_{m=1}^{m_t-1} a_m \phi(\tilde{\mathbf{x}}_m) - \phi(\mathbf{x}_t) \right\|^2 \leq \nu, \quad (16)$$

where  $\nu$  is the sparsity level (hyper-)parameter.

<sup>3</sup>ENGEL, Y.; MANNOR, S.; MEIR, R. **Sparse online greedy support vector regression**. In: ELOMAA, T.; MANNILA, H.; TOIVONEN, H. (Ed.). Proceedings of the 13th European Conference on Machine Learning (ECML'2002). Helsinki, Finland, 2002. p. 84-96.



# Sparsifying the KOLS Model

- As mentioned before, the size of the kernel matrix and, hence, the dimension of the vector of parameters  $\alpha$  is determined by the number of training samples  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ .
- To avoid inversion of huge kernel matrices, we must use some kind of sparsification method to select a subset of relevant sample pairs to compose a dictionary.
- In the ALD<sup>3</sup> criterion, when a new incoming sample  $\mathbf{x}_t$  is available, one must verify if  $\phi(\mathbf{x}_t)$  is ALD on the **dictionary**  $\mathcal{D}_{t-1}^{sv} = \{\tilde{\mathbf{x}}_j\}_{j=1}^{m_{t-1}}$ .
- One should estimate  $\mathbf{a} = [a_1, \dots, a_{m_{t-1}}]^\top$  satisfying the ALD criterion

$$\delta_t \triangleq \min_{\mathbf{a}} \left\| \sum_{m=1}^{m_{t-1}} a_m \phi(\tilde{\mathbf{x}}_m) - \phi(\mathbf{x}_t) \right\|^2 \leq \nu, \quad (16)$$

where  $\nu$  is the sparsity level (hyper-)parameter.

<sup>3</sup>ENGEL, Y.; MANNOR, S.; MEIR, R. **Sparse online greedy support vector regression**. In: ELOMAA, T.; MANNILA, H.; TOIVONEN, H. (Ed.). Proceedings of the 13th European Conference on Machine Learning (ECML'2002). Helsinki, Finland, 2002. p. 84-96.

# Sparsifying the KOLS Model

- As mentioned before, the size of the kernel matrix and, hence, the dimension of the vector of parameters  $\alpha$  is determined by the number of training samples  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ .
- To avoid inversion of huge kernel matrices, we must use some kind of sparsification method to select a subset of relevant sample pairs to compose a dictionary.
- In the ALD<sup>3</sup> criterion, when a new incoming sample  $\mathbf{x}_t$  is available, one must verify if  $\phi(\mathbf{x}_t)$  is ALD on the **dictionary**  $\mathcal{D}_{t-1}^{sv} = \{\tilde{\mathbf{x}}_j\}_{j=1}^{m_{t-1}}$ .
- One should estimate  $\mathbf{a} = [a_1, \dots, a_{m_{t-1}}]^\top$  satisfying the ALD criterion

$$\delta_t \triangleq \min_{\mathbf{a}} \left\| \sum_{m=1}^{m_{t-1}} a_m \phi(\tilde{\mathbf{x}}_m) - \phi(\mathbf{x}_t) \right\|^2 \leq \nu, \quad (16)$$

where  $\nu$  is the sparsity level (hyper-)parameter.

<sup>3</sup>ENGEL, Y.; MANNOR, S.; MEIR, R. **Sparse online greedy support vector regression**. In: ELOMAA, T.; MANNILA, H.; TOIVONEN, H. (Ed.). Proceedings of the 13th European Conference on Machine Learning (ECML'2002). Helsinki, Finland, 2002. p. 84-96.

# Sparsifying the KOLS Model

- As mentioned before, the size of the kernel matrix and, hence, the dimension of the vector of parameters  $\alpha$  is determined by the number of training samples  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ .
- To avoid inversion of huge kernel matrices, we must use some kind of sparsification method to select a subset of relevant sample pairs to compose a dictionary.
- In the ALD<sup>3</sup> criterion, when a new incoming sample  $\mathbf{x}_t$  is available, one must verify if  $\phi(\mathbf{x}_t)$  is ALD on the **dictionary**  $\mathcal{D}_{t-1}^{sv} = \{\tilde{\mathbf{x}}_j\}_{j=1}^{m_{t-1}}$ .
- One should estimate  $\mathbf{a} = [a_1, \dots, a_{m_{t-1}}]^\top$  satisfying the ALD criterion

$$\delta_t \triangleq \min_{\mathbf{a}} \left\| \sum_{m=1}^{m_{t-1}} a_m \phi(\tilde{\mathbf{x}}_m) - \phi(\mathbf{x}_t) \right\|^2 \leq \nu, \quad (16)$$

where  $\nu$  is the sparsity level (hyper-)parameter.

<sup>3</sup>ENGEL, Y.; MANNOR, S.; MEIR, R. **Sparse online greedy support vector regression**. In: ELOMAA, T.; MANNILA, H.; TOIVONEN, H. (Ed.). Proceedings of the 13th European Conference on Machine Learning (ECML'2002). Helsinki, Finland, 2002. p. 84-96.

## Sparsifying the KOLS Model

- In terms of the kernel matrix, the Eq. (16) can be rewritten as

$$\mathbf{a}_t = \tilde{\mathbf{K}}_{t-1}^{-1} \tilde{\mathbf{k}}_{t-1}(\mathbf{x}_t), \quad \text{and} \quad \delta_t = k_{tt} - \tilde{\mathbf{k}}_{t-1}(\mathbf{x}_t)^\top \mathbf{a}_t \leq \nu, \quad (17)$$

where  $\nu$  is the sparsity level parameter. The 1st element in the dictionary is chosen at random.

- Once a pass of the ALD criterion is concluded, we can rewrite the problem as

$$J(\tilde{\alpha}) = \|\tilde{\mathbf{y}} - \tilde{\mathbf{A}}\tilde{\mathbf{K}}\tilde{\alpha}\|^2, \quad (18)$$

where  $\tilde{\mathbf{A}} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_{m_t}]^\top \in \mathbb{R}^{t \times m_t}$  and  $\tilde{\alpha} \in \mathbb{R}^{m_t}$ .

- Then, computing  $\partial J(\tilde{\alpha}) / \partial \tilde{\alpha} = \mathbf{0}$ , one gets

$$\tilde{\alpha} = \tilde{\mathbf{K}}^{-1} \tilde{\mathbf{P}} \tilde{\mathbf{A}}^\top \tilde{\mathbf{y}}, \quad (19)$$

where  $\tilde{\mathbf{P}} = (\tilde{\mathbf{A}}^\top \tilde{\mathbf{A}})^{-1}$ .

## Sparsifying the KOLS Model

- In terms of the kernel matrix, the Eq. (16) can be rewritten as

$$\mathbf{a}_t = \tilde{\mathbf{K}}_{t-1}^{-1} \tilde{\mathbf{k}}_{t-1}(\mathbf{x}_t), \quad \text{and} \quad \delta_t = k_{tt} - \tilde{\mathbf{k}}_{t-1}(\mathbf{x}_t)^\top \mathbf{a}_t \leq \nu, \quad (17)$$

where  $\nu$  is the sparsity level parameter. The 1st element in the dictionary is chosen at random.

- Once a pass of the ALD criterion is concluded, we can rewrite the problem as

$$J(\tilde{\boldsymbol{\alpha}}) = \|\tilde{\mathbf{y}} - \tilde{\mathbf{A}}\tilde{\mathbf{K}}\tilde{\boldsymbol{\alpha}}\|^2, \quad (18)$$

where  $\tilde{\mathbf{A}} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_{m_t}]^\top \in \mathbb{R}^{t \times m_t}$  and  $\tilde{\boldsymbol{\alpha}} \in \mathbb{R}^{m_t}$ .

- Then, computing  $\partial J(\tilde{\boldsymbol{\alpha}})/\partial \tilde{\boldsymbol{\alpha}} = \mathbf{0}$ , one gets

$$\tilde{\boldsymbol{\alpha}} = \tilde{\mathbf{K}}^{-1} \tilde{\mathbf{P}} \tilde{\mathbf{A}}^\top \tilde{\mathbf{y}}, \quad (19)$$

where  $\tilde{\mathbf{P}} = (\tilde{\mathbf{A}}^\top \tilde{\mathbf{A}})^{-1}$ .

## Sparsifying the KOLS Model

- In terms of the kernel matrix, the Eq. (16) can be rewritten as

$$\mathbf{a}_t = \tilde{\mathbf{K}}_{t-1}^{-1} \tilde{\mathbf{k}}_{t-1}(\mathbf{x}_t), \quad \text{and} \quad \delta_t = k_{tt} - \tilde{\mathbf{k}}_{t-1}(\mathbf{x}_t)^\top \mathbf{a}_t \leq \nu, \quad (17)$$

where  $\nu$  is the sparsity level parameter. The 1st element in the dictionary is chosen at random.

- Once a pass of the ALD criterion is concluded, we can rewrite the problem as

$$J(\tilde{\boldsymbol{\alpha}}) = \|\tilde{\mathbf{y}} - \tilde{\mathbf{A}}\tilde{\mathbf{K}}\tilde{\boldsymbol{\alpha}}\|^2, \quad (18)$$

where  $\tilde{\mathbf{A}} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_{m_t}]^\top \in \mathbb{R}^{t \times m_t}$  and  $\tilde{\boldsymbol{\alpha}} \in \mathbb{R}^{m_t}$ .

- Then, computing  $\partial J(\tilde{\boldsymbol{\alpha}})/\partial \tilde{\boldsymbol{\alpha}} = \mathbf{0}$ , one gets

$$\tilde{\boldsymbol{\alpha}} = \tilde{\mathbf{K}}^{-1} \tilde{\mathbf{P}} \tilde{\mathbf{A}}^\top \tilde{\mathbf{y}}, \quad (19)$$

where  $\tilde{\mathbf{P}} = (\tilde{\mathbf{A}}^\top \tilde{\mathbf{A}})^{-1}$ .

# Outline

## ① Introduction

Related Publications

Kernel-Based Regression Models

General and Specific Objectives

## ② The LSSVR Model

## ③ The KOLS Model

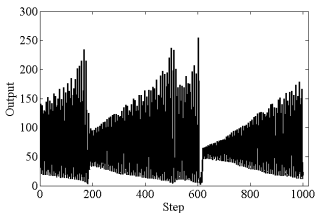
## ④ Computational Experiments

## ⑤ Conclusions

# Computational Experiments - Laser Time Series

Dataset	Task	Prediction Type	$N$	$N'$	$\hat{L}_u$	$\hat{L}_y$
Laser	Time series prediction	Free Simulation	1,000	100/500	-	50

- The **chaotic laser time series**<sup>4 5</sup> is a benchmarking dataset with a total of 10,093 samples.
- Scenarios of test: (i) using the next 100 samples;  
(ii) using the next 500 samples.

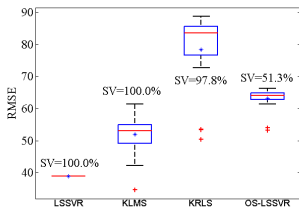


<sup>4</sup> GERSHENFELD, N. A.; WEIGEND, A. S. The future of time series. In: WEIGEND, A. S.; GERSHENFELD, N. A. (Ed.). **Times Series Prediction: Forecasting the Future and Understanding the Past**. Reading, MA: Addison-Wesley, 1993.

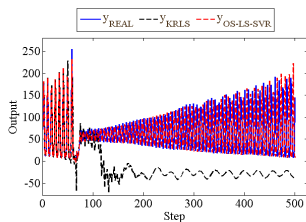
<sup>5</sup> Available for download at [www-psych.stanford.edu/~sim\\$andreas/Time-Series/SantaFe.html](http://www-psych.stanford.edu/~sim$andreas/Time-Series/SantaFe.html).



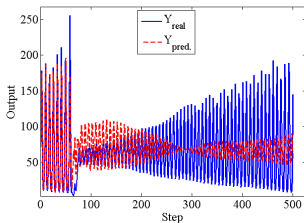
# Computational Experiments - Laser Time Series



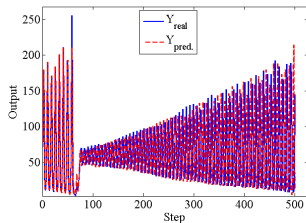
(a) RMSE values.



(b) KRLS and OS-LSSVR.



(c) LSSVR.



(d) OS-LSSVR.

# Computational Experiments - Monthly Rainfall

- Average monthly rainfall at the seashore of Fortaleza from 1983 to 2021.
- Training (1983-2016), Validation (2017-2020), Testing (2021).
- Validation for model selection and Testing for actual prediction.
- Validation and Testing in free simulation mode (recursive prediction).
- Initial 8 regressor values taken from 1982.

# Computational Experiments - Monthly Rainfall

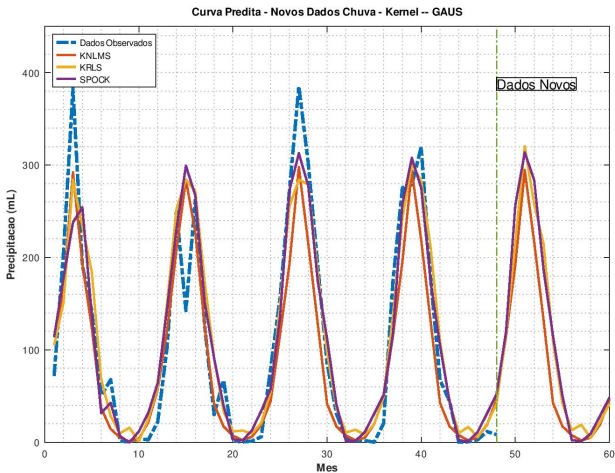
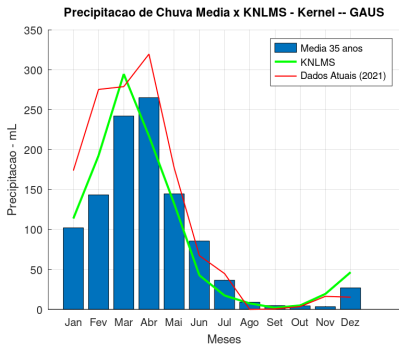
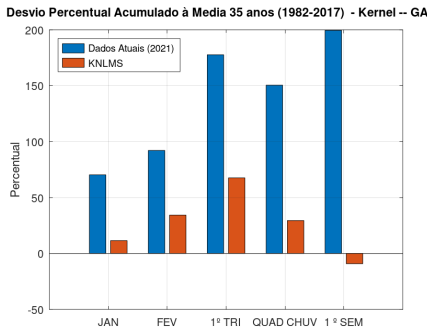


Figure 1: Predicted time series in free simulation mode.

# Computational Experiments - Monthly Rainfall



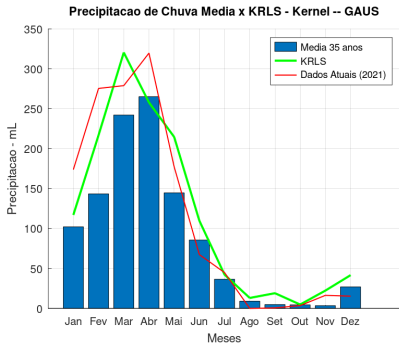
(a) Predicted Rainfall.



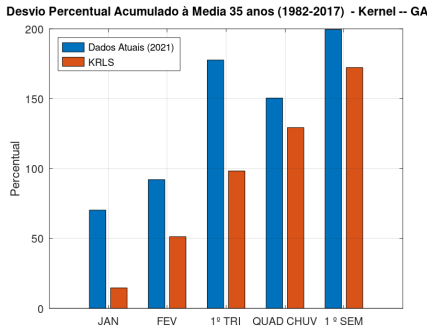
(b) Relative Error.

Figure 2: Results for the KNLMS model.

# Computational Experiments - Monthly Rainfall



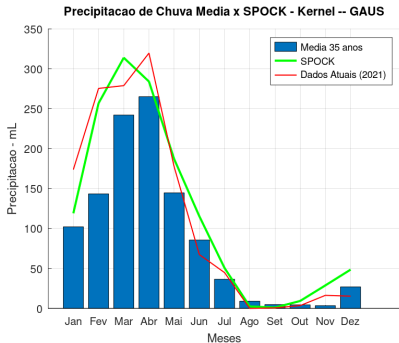
(a) Predicted Rainfall.



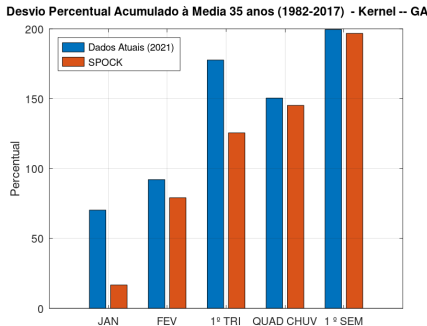
(b) Relative Error.

Figure 3: Results for the KRLS model.

# Computational Experiments - Monthly Rainfall



(a) Predicted Rainfall.



(b) Relative Error.

Figure 4: Results for the SPOCK model.

# Outline

## ① Introduction

Related Publications

Kernel-Based Regression Models

General and Specific Objectives

## ② The LSSVR Model

## ③ The KOLS Model

## ④ Computational Experiments

## ⑤ Conclusions

# Conclusions

- We have introduced some kernel-based nonlinear regression models.
- These models were successfully applied to time series forecasting.
- They can be made more compact by means of sparsification techniques, such as the ALD.
- They can be easily extended to online learning.
- They can be extended in order to handle outlier in the data.



# Conclusions

- We have introduced some kernel-based nonlinear regression models.
- These models were successfully applied to time series forecasting.
- They can be made more compact by means of sparsification techniques, such as the ALD.
- They can be easily extended to online learning.
- They can be extended in order to handle outlier in the data.

# Conclusions

- We have introduced some kernel-based nonlinear regression models.
- These models were successfully applied to time series forecasting.
- They can be made more compact by means of sparsification techniques, such as the ALD.
- They can be easily extended to online learning.
- They can be extended in order to handle outlier in the data.

# Conclusions

- We have introduced some kernel-based nonlinear regression models.
- These models were successfully applied to time series forecasting.
- They can be made more compact by means of sparsification techniques, such as the ALD.
- They can be easily extended to online learning.
- They can be extended in order to handle outlier in the data.

# Conclusions

- We have introduced some kernel-based nonlinear regression models.
- These models were successfully applied to time series forecasting.
- They can be made more compact by means of sparsification techniques, such as the ALD.
- They can be easily extended to online learning.
- They can be extended in order to handle outlier in the data.

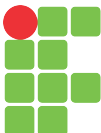
# Conclusions

- We have introduced some kernel-based nonlinear regression models.
- These models were successfully applied to time series forecasting.
- They can be made more compact by means of sparsification techniques, such as the ALD.
- They can be easily extended to online learning.
- They can be extended in order to handle outlier in the data.

# Acknowledgements



UNIVERSIDADE  
FEDERAL DO CEARÁ



INSTITUTO FEDERAL DE  
EDUCAÇÃO, CIÊNCIA E TECNOLOGIA  
CEARÁ



Signal and Information  
Processing for Data Analysis  
and Learning Systems

**THANK YOU!!!**