

Nonparametric Reduced-Rank Regression with Tensor-Structured Response

Guillaume Rabusseau Hachem Kadri François Denis

Journée "Tenseurs et estimation de matrices de covariance"

November 27, 2015

Overview

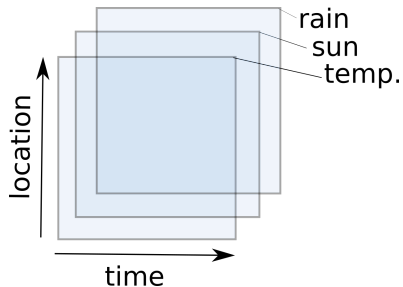
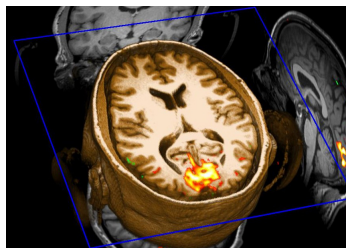
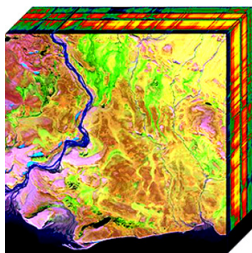
- 1 Introduction
- 2 Regression and Tensors
- 3 Nonparametric Tensor-valued Reduced Rank Regression
- 4 Experiments
- 5 Conclusion and Perspectives

Outline

- 1 Introduction
- 2 Regression and Tensors
- 3 Nonparametric Tensor-valued Reduced Rank Regression
- 4 Experiments
- 5 Conclusion and Perspectives

Introduction

- Data with tensor structure: EEG, hyperspectral images, videos, ...



Introduction

- Data with tensor structure: EEG, hyperspectral images, videos, ...
- Tensor learning
 - ▶ Tensor decomposition and latent variable models
 - ▶ Source separation
 - ▶ Tensor completion, ...

Introduction

- Data with tensor structure: EEG, hyperspectral images, videos, ...
- Tensor learning
 - ▶ Tensor decomposition and latent variable models
 - ▶ Source separation
 - ▶ Tensor completion, ...
 - ▶ Regression with tensor input

- ★ Multilinear model:

$$f(\mathcal{X}) = \langle \mathcal{W}, \mathcal{X} \rangle$$

- ★ Multi-View/Factorisation Machines:

$$f(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}) = \langle \mathcal{W}, \mathbf{x}^{(1)} \otimes \mathbf{x}^{(2)} \otimes \mathbf{x}^{(3)} \rangle$$

- ★ Learning Tensors in RKHS with Multilinear Spectral Penalties (Signoretto et al., 2013)

Introduction

- Data with tensor structure: EEG, hyperspectral images, videos, ...
- Tensor learning
 - ▶ Tensor decomposition and latent variable models
 - ▶ Source separation
 - ▶ Tensor completion, ...
 - ▶ Regression with tensor input

- ★ Multilinear model:

$$f(\mathcal{X}) = \langle \mathcal{W}, \mathcal{X} \rangle$$

- ★ Multi-View/Factorisation Machines:

$$f(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}) = \langle \mathcal{W}, \mathbf{x}^{(1)} \otimes \mathbf{x}^{(2)} \otimes \mathbf{x}^{(3)} \rangle$$

- ★ Learning Tensors in RKHS with Multilinear Spectral Penalties (Signoretto et al., 2013)
- ▶ Regression with tensor output
 - ★ Multilinear models (multitask and spatio-temporal forecasting)
 - ★ Nonlinear / Nonparametric?

Outline

- 1 Introduction
- 2 Regression and Tensors
- 3 Nonparametric Tensor-valued Reduced Rank Regression
- 4 Experiments
- 5 Conclusion and Perspectives

Tensors

- **Tensor** $\mathcal{T} \in \mathbb{R}^{d_1} \otimes \dots \otimes \mathbb{R}^{d_p} \simeq$ Multi-array $(\mathcal{T}_{i_1 \dots i_k}) \in \mathbb{R}^{d_1 \times \dots \times d_p}$.
 - ▶ $p = 3$: higher order tensor $\mathcal{T}_{i_1 i_2 i_3} \in \mathbb{R}$ for $i_1 \in [d_1], i_2 \in [d_2], i_3 \in [d_3]$.

Tensors

- **Tensor** $\mathcal{J} \in \mathbb{R}^{d_1} \otimes \dots \otimes \mathbb{R}^{d_p} \simeq$ Multi-array $(\mathcal{J}_{i_1 \dots i_k}) \in \mathbb{R}^{d_1 \times \dots \times d_p}$.
 - ▶ $p = 3$: higher order tensor $\mathcal{J}_{i_1 i_2 i_3} \in \mathbb{R}$ for $i_1 \in [d_1], i_2 \in [d_2], i_3 \in [d_3]$.
- The **CP-rank** of \mathcal{J} is the smallest integer R such that

$$\mathcal{J} = \sum_{r=1}^R \mathbf{v}_r^1 \otimes \dots \otimes \mathbf{v}_r^p \quad \text{with } \mathbf{v}_r^1 \in \mathbb{R}^{d_1}, \dots, \mathbf{v}_r^p \in \mathbb{R}^{d_p} .$$

- ▶ We write $\mathcal{J} = \llbracket \mathbf{V}^1, \dots, \mathbf{V}^p \rrbracket$ where $\mathbf{V}^i = (\mathbf{v}_1^i \dots \mathbf{v}_R^i) \in \mathbb{R}^{d_i \times R}$.

Tensors

- **Tensor** $\mathcal{T} \in \mathbb{R}^{d_1} \otimes \dots \otimes \mathbb{R}^{d_p} \simeq$ Multi-array $(\mathcal{T}_{i_1 \dots i_k}) \in \mathbb{R}^{d_1 \times \dots \times d_p}$.
 - ▶ $p = 3$: higher order tensor $\mathcal{T}_{i_1 i_2 i_3} \in \mathbb{R}$ for $i_1 \in [d_1], i_2 \in [d_2], i_3 \in [d_3]$.
- The **CP-rank** of \mathcal{T} is the smallest integer R such that

$$\mathcal{T} = \sum_{r=1}^R \mathbf{v}_r^1 \otimes \dots \otimes \mathbf{v}_r^p \quad \text{with } \mathbf{v}_r^1 \in \mathbb{R}^{d_1}, \dots, \mathbf{v}_r^p \in \mathbb{R}^{d_p} .$$

- ▶ We write $\mathcal{T} = \llbracket \mathbf{V}^1, \dots, \mathbf{V}^p \rrbracket$ where $\mathbf{V}^i = (\mathbf{v}_1^i \dots \mathbf{v}_R^i) \in \mathbb{R}^{d_i \times R}$.
- **Multiplication** between tensors. Let $\mathcal{A} \in R^{m_1 \times m_2 \times n_1 \times n_2}$ and $\mathcal{B} \in R^{n_1 \times n_2 \times p_1 \times p_2}$, we define $\mathcal{A}\mathcal{B} \in \mathbb{R}^{m_1 \times m_2 \times p_1 \times p_2}$ by

$$(\mathcal{A}\mathcal{B})_{i_1 i_2 k_1 k_2} = \sum_{j_1 j_2} \mathcal{A}_{i_1 i_2 j_1 j_2} \mathcal{B}_{j_1 j_2 k_1 k_2}$$

⇒ Composition of multilinear maps...

Multivariate Regression

Learn $f : \mathbb{R}^d \rightarrow \mathbb{R}^p$ from samples $\{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}_{n=1}^N$ where $\mathbf{y}^{(n)} \simeq f(\mathbf{x}^{(n)})$.

- Linear model: $f(\mathbf{x}) = \mathbf{W}^\top \mathbf{x}$ ($\mathbf{W} \in \mathbb{R}^{d \times p}$)

Multivariate Regression

Learn $f : \mathbb{R}^d \rightarrow \mathbb{R}^p$ from samples $\{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}_{n=1}^N$ where $\mathbf{y}^{(n)} \simeq f(\mathbf{x}^{(n)})$.

- Linear model: $f(\mathbf{x}) = \mathbf{W}^\top \mathbf{x}$ ($\mathbf{W} \in \mathbb{R}^{d \times p}$)
- Ordinary Least Squares

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W} \in \mathbb{R}^{d \times p}} \|\mathbf{X}\mathbf{W} - \mathbf{Y}\|_F^2 \quad (\mathbf{X} \in \mathbb{R}^{N \times d}, \mathbf{Y} \in \mathbb{R}^{N \times p})$$

⇒ Equivalent to perform p independent linear regressions!
How can we capture linear dependencies in the output?

Multivariate Regression

Learn $f : \mathbb{R}^d \rightarrow \mathbb{R}^p$ from samples $\{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}_{n=1}^N$ where $\mathbf{y}^{(n)} \simeq f(\mathbf{x}^{(n)})$.

- Linear model: $f(\mathbf{x}) = \mathbf{W}^\top \mathbf{x}$ ($\mathbf{W} \in \mathbb{R}^{d \times p}$)
- Ordinary Least Squares

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W} \in \mathbb{R}^{d \times p}} \|\mathbf{X}\mathbf{W} - \mathbf{Y}\|_F^2 \quad (\mathbf{X} \in \mathbb{R}^{N \times d}, \mathbf{Y} \in \mathbb{R}^{N \times p})$$

⇒ Equivalent to perform p independent linear regressions!
How can we capture linear dependencies in the output?

- Reduced Rank Regression (Izenman, 1975)

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W} \in \mathbb{R}^{d \times p}} \|\mathbf{X}\mathbf{W} - \mathbf{Y}\|_F^2 \quad \text{s.t. } \text{rank}(\mathbf{W}) \leq R$$

Tensor-valued Regression

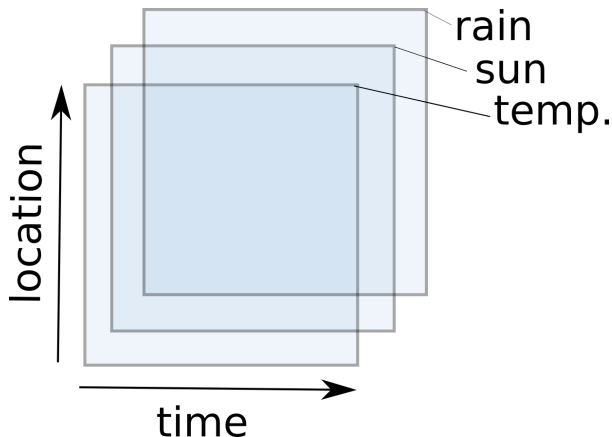
Learn $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d_1 \times \dots \times d_p}$ from $\{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}_{n=1}^N$ where $\mathbf{y}^{(n)} \simeq f(\mathbf{x}^{(n)})$.

Tensor-valued Regression

Learn $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d_1 \times \dots \times d_p}$ from $\{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}_{n=1}^N$ where $\mathbf{y}^{(n)} \simeq f(\mathbf{x}^{(n)})$.

- Spatio-Temporal Forecasting (Bahadori et al., 2014)

$$f(\mathbf{x}) \in \mathbb{R}^{(\text{Times}) \times (\text{Locations}) \times (\text{Variables})}$$



Tensor-valued Regression

Learn $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d_1 \times \dots \times d_p}$ from $\{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}_{n=1}^N$ where $\mathbf{y}^{(n)} \simeq f(\mathbf{x}^{(n)})$.

- Multilinear Multitask Learning (Romera-Paredes et al., 2013)

$$f(\mathbf{x}) \in \mathbb{R}^{(\text{Restaurant Critics}) \times (\text{Evaluation Criteria})}$$

| Rest. 1 | Critic 1 | Critic 2 | Critic 3 |
|---------|----------|----------|----------|
| host | 5 | 3 | 6 |
| food | 7 | 8 | 6.5 |
| price | 5 | 6.5 | 4 |

| Rest. 2 | Critic 1 | Critic 2 | Critic 3 |
|---------|----------|----------|----------|
| host | 7 | 8 | 6 |
| food | 8.5 | 9 | 9 |
| price | 8 | 9.5 | 7 |

...

Tensor-valued Regression

Learn $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d_1 \times \dots \times d_p}$ from $\{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}_{n=1}^N$ where $\mathbf{y}^{(n)} \simeq f(\mathbf{x}^{(n)})$.

- Linear model: $f(\mathbf{x}) = \mathbf{x}\mathcal{W}$ ($\mathcal{W} \in \mathbb{R}^{d \times d_1 \times \dots \times d_p}$)

Tensor-valued Regression

Learn $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d_1 \times \dots \times d_p}$ from $\{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}_{n=1}^N$ where $\mathbf{y}^{(n)} \simeq f(\mathbf{x}^{(n)})$.

- Linear model: $f(\mathbf{x}) = \mathbf{x}\mathbf{W}$ ($\mathbf{W} \in \mathbb{R}^{d \times d_1 \times \dots \times d_p}$)
- Vectorisation of the outputs \rightarrow Reduced Rank Regression

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W} \in \mathbb{R}^{d \times d_1 d_2 \dots d_p}} \|\mathbf{X}\mathbf{W} - \mathbf{Y}\|_F^2 \quad \text{s.t. } \text{rank}(\mathbf{W}) \leq R$$

where $\mathbf{Y}_{n,:} = \text{vec}(\mathbf{y}^{(n)})^\top$

\Rightarrow Tensor structure of the output is lost!

Tensor-valued Regression

Learn $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d_1 \times \dots \times d_p}$ from $\{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}_{n=1}^N$ where $\mathbf{y}^{(n)} \simeq f(\mathbf{x}^{(n)})$.

- Linear model: $f(\mathbf{x}) = \mathbf{x}\mathcal{W}$ ($\mathcal{W} \in \mathbb{R}^{d \times d_1 \times \dots \times d_p}$)
- Vectorisation of the outputs \rightarrow Reduced Rank Regression

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W} \in \mathbb{R}^{d \times d_1 d_2 \dots d_p}} \|\mathbf{X}\mathbf{W} - \mathbf{Y}\|_F^2 \quad \text{s.t. } \text{rank}(\mathbf{W}) \leq R$$

where $\mathbf{Y}_{n,:} = \text{vec}(\mathbf{y}^{(n)})^\top$

\Rightarrow Tensor structure of the output is lost!

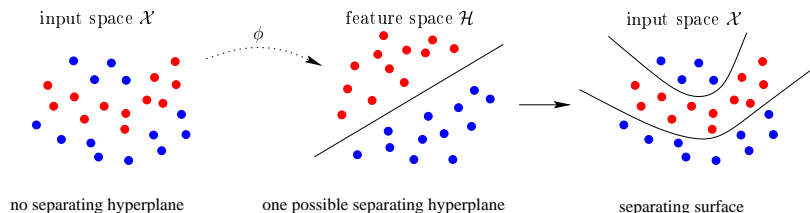
- **Reduced Rank Regression for Tensor Structured Response**

$$\hat{\mathcal{W}} = \arg \min_{\mathcal{W} \in \mathbb{R}^{d \times d_1 \times d_2 \dots \times d_p}} \|\mathbf{X}\mathcal{W} - \mathbf{y}\|_F^2 \quad \text{s.t. } \text{rank}(\mathcal{W}) \leq R$$

Outline

- 1 Introduction
- 2 Regression and Tensors
- 3 Nonparametric Tensor-valued Reduced Rank Regression**
- 4 Experiments
- 5 Conclusion and Perspectives

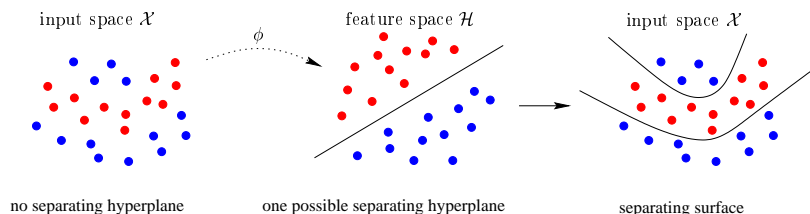
Scalar-valued Kernels



Nonparametric learning of $f : \mathcal{X} \rightarrow \mathcal{Y} = \mathbb{R}$.

- **Scalar-valued** kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ positive semidefinite.
- Learn f in the RKHS $\mathcal{H}_k = \overline{\text{span}\{k(\mathbf{z}, \cdot) \mid \mathbf{z} \in \mathcal{X}\}}$

Scalar-valued Kernels



Nonparametric learning of $f : \mathcal{X} \rightarrow \mathcal{Y} = \mathbb{R}$.

- **Scalar-valued** kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ positive semidefinite.
- Learn f in the RKHS $\mathcal{H}_k = \overline{\text{span}\{k(\mathbf{z}, \cdot) \mid \mathbf{z} \in \mathcal{X}\}}$
- Representer theorem

$$\hat{f} = \arg \min_{f \in \mathcal{H}_k} \sum_{n=1}^N (f(\mathbf{x}^{(n)}) - y^{(n)})^2 + \gamma \|f\|_{\mathcal{H}_k}^2 \implies \hat{f} = \sum_{n=1}^N \alpha_n k(\mathbf{x}^{(n)}, \cdot)$$

for some $\alpha_1, \dots, \alpha_N \in \mathbb{R}$.

Matrix-valued Kernels (Micchelli and Pontil, 2005)

Nonparametric learning of $f : \mathcal{X} \rightarrow \mathcal{Y} = \mathbb{R}^p$.

- **Matrix-valued kernel** $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y}) \simeq \mathbb{R}^{p \times p}$ positive semidefinite.
- Learn f in the RKHS $\mathcal{H}_K = \overline{\text{span}\{K(\mathbf{z}, \cdot)\mathbf{c} \mid \mathbf{z} \in \mathcal{X}, \mathbf{c} \in \mathcal{Y}\}}$

Matrix-valued Kernels (Micchelli and Pontil, 2005)

Nonparametric learning of $f : \mathcal{X} \rightarrow \mathcal{Y} = \mathbb{R}^p$.

- **Matrix-valued kernel** $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y}) \simeq \mathbb{R}^{p \times p}$ positive semidefinite.
- Learn f in the RKHS $\mathcal{H}_K = \overline{\text{span}\{K(\mathbf{z}, \cdot)\mathbf{c} \mid \mathbf{z} \in \mathcal{X}, \mathbf{c} \in \mathcal{Y}\}}$
- Representer theorem

$$\hat{f} = \arg \min_{f \in \mathcal{H}_K} \sum_{n=1}^N (f(\mathbf{x}^{(n)}) - \mathbf{y}^{(n)})^2 + \gamma \|f\|_{\mathcal{H}_K}^2 \implies \hat{f} = \sum_{n=1}^N K(\mathbf{x}^{(n)}, \cdot) \mathbf{c}^{(n)}$$

for some $\mathbf{c}^{(1)}, \dots, \mathbf{c}^{(N)} \in \mathbb{R}^p$.

Matrix-valued Kernels (Micchelli and Pontil, 2005)

Nonparametric learning of $f : \mathcal{X} \rightarrow \mathcal{Y} = \mathbb{R}^p$.

- **Matrix-valued kernel** $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y}) \simeq \mathbb{R}^{p \times p}$ positive semidefinite.
- Learn f in the RKHS $\mathcal{H}_K = \overline{\text{span}\{K(\mathbf{z}, \cdot)\mathbf{c} \mid \mathbf{z} \in \mathcal{X}, \mathbf{c} \in \mathcal{Y}\}}$
- Representer theorem

$$\hat{f} = \arg \min_{f \in \mathcal{H}_K} \sum_{n=1}^N (f(\mathbf{x}^{(n)}) - \mathbf{y}^{(n)})^2 + \gamma \|f\|_{\mathcal{H}_K}^2 \implies \hat{f} = \sum_{n=1}^N K(\mathbf{x}^{(n)}, \cdot) \mathbf{c}^{(n)}$$

for some $\mathbf{c}^{(1)}, \dots, \mathbf{c}^{(N)} \in \mathbb{R}^p$.

- **Separable** kernel $K(\cdot, \cdot) = k(\cdot, \cdot)\mathbf{T}$ for some scalar kernel k and some positive semidefinite $\mathbf{T} \in \mathbb{R}^{p \times p}$

$$\arg \min_{\mathbf{C} \in \mathbb{R}^{N \times p}} \|\mathbf{KCT} - \mathbf{Y}\|_F^2 + \gamma \langle \mathbf{KCT}, \mathbf{C} \rangle$$

Tensor-valued Kernels

Nonparametric learning of $f : \mathcal{X} \rightarrow \mathcal{Y} = \mathbb{R}^{d_1 \times \dots \times d_p}$.

- **Tensor-valued kernel** $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y}) \simeq \mathbb{R}^{d_1 \times \dots \times d_p \times d_1 \times \dots \times d_p}$
positive semidefinite.

Tensor-valued Kernels

Nonparametric learning of $f : \mathcal{X} \rightarrow \mathcal{Y} = \mathbb{R}^{d_1 \times \dots \times d_p}$.

- **Tensor-valued kernel** $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y}) \simeq \mathbb{R}^{d_1 \times \dots \times d_p \times d_1 \times \dots \times d_p}$ positive semidefinite.
- Representer theorem

$$\hat{f} = \arg \min_{f \in \mathcal{H}_K} \sum_{n=1}^N (f(\mathbf{x}^{(n)}) - \mathbf{y}^{(n)})^2 + \gamma \|f\|_{\mathcal{H}_K}^2 \implies \hat{f} = \sum_{n=1}^N K(\mathbf{x}^{(n)}, \cdot) \mathbf{c}^{(n)}$$

for some $\mathbf{c}^{(1)}, \dots, \mathbf{c}^{(N)} \in \mathbb{R}^{d_1 \times \dots \times d_p}$.

Tensor-valued Kernels

Nonparametric learning of $f : \mathcal{X} \rightarrow \mathcal{Y} = \mathbb{R}^{d_1 \times \dots \times d_p}$.

- **Tensor-valued kernel** $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y}) \simeq \mathbb{R}^{d_1 \times \dots \times d_p \times d_1 \times \dots \times d_p}$ positive semidefinite.
- Representer theorem

$$\hat{f} = \arg \min_{f \in \mathcal{H}_K} \sum_{n=1}^N (f(\mathbf{x}^{(n)}) - \mathbf{y}^{(n)})^2 + \gamma \|f\|_{\mathcal{H}_K}^2 \implies \hat{f} = \sum_{n=1}^N K(\mathbf{x}^{(n)}, \cdot) \mathbf{c}^{(n)}$$

for some $\mathbf{c}^{(1)}, \dots, \mathbf{c}^{(N)} \in \mathbb{R}^{d_1 \times \dots \times d_p}$.

- **Separable** kernel $K(\cdot, \cdot) = k(\cdot, \cdot) \mathcal{J}$ for some scalar kernel k and some positive semidefinite $\mathcal{J} \in \mathbb{R}^{d_1 \times \dots \times d_p \times d_1 \times \dots \times d_p}$

$$\arg \min_{\mathbf{c} \in \mathbb{R}^{N \times d_1 \times \dots \times d_p}} \|\mathbf{K} \mathbf{c} \mathcal{J} - \mathbf{y}\|_F^2 + \gamma \langle \mathbf{K} \mathbf{c} \mathcal{J}, \mathbf{c} \rangle$$

Nonparametric Tensor-valued Reduced Rank Regression

Separable kernel $K(\cdot, \cdot) = k(\cdot, \cdot)\mathcal{T}$

$$\arg \min_{\mathcal{C} \in \mathbb{R}^{N \times d_1 \times \dots \times d_p}} \|\mathbf{K}\mathcal{C}\mathcal{T} - \mathbf{y}\|_F^2 + \gamma \langle \mathbf{K}\mathcal{C}\mathcal{T}, \mathcal{C} \rangle$$

\Rightarrow Does not take the tensor structure of the output into account!

We add a low CP-rank constraint on the tensor $\mathcal{C}\mathcal{T}$:

$$\arg \min_{\mathcal{C} \in \mathbb{R}^{N \times d_1 \times \dots \times d_p}} \|\mathbf{K}\mathcal{C}\mathcal{T} - \mathbf{y}\|_F^2 + \gamma \langle \mathbf{K}\mathcal{C}\mathcal{T}, \mathcal{C} \rangle \quad \text{s.t. } \text{rank}(\mathcal{C}\mathcal{T}) \leq R$$

\Rightarrow Back to the linear case when $k(\cdot, \cdot)$ is the linear kernel and $\mathcal{T} = \mathcal{J}$

- TVK-RRR-I: Special case $\mathcal{T} = \mathcal{J}$.

$$\arg \min_{\mathbf{C} \in \mathbb{R}^{N \times d_1 \times \dots \times d_p}} \|\mathbf{K}\mathbf{C} - \mathbf{y}\|_F^2 + \gamma \langle \mathbf{K}\mathbf{C}, \mathbf{C} \rangle \quad \text{s.t. } \text{rank}(\mathbf{C}) \leq R$$

⇒ Alternating Least Squares ($\mathbf{C} = [\mathbf{U}^0, \mathbf{U}^1, \dots, \mathbf{U}^p]$)

Algorithm 1 TVK-RRR-I

Require: Gram matrix $\mathbf{K} \in \mathbb{R}^N \times \mathbb{R}^N$, output tensor $\mathcal{Y} \in \mathbb{R}^N \otimes \mathcal{Y}$, regularization parameter R .

Ensure: $\mathcal{C} \in \mathbb{R}^N \otimes \mathcal{Y}$

Randomly initialize $\mathbf{U}^i \in \mathbb{R}^{d_i \times R}$ for $1 \leq i \leq p + 1$

repeat

$$\mathbf{U}^1 \leftarrow (\mathbf{K}(\mathbf{K} + \gamma \mathbf{I}))^+ (\mathbf{K}\mathcal{Y})_{(1)} \mathbf{U}_{\odot 1} (\mathbf{U}_{\odot 1}^\top \mathbf{U}_{\odot 1})^+$$

for $i = 2, \dots, p + 1$ **do**

$$\tilde{\mathbf{U}} \leftarrow (\mathbf{K}(\mathbf{K} + \gamma \mathbf{I})) \mathbf{U}^1, \mathbf{U}^2, \dots, \mathbf{U}^{p+1}$$

$$\mathbf{U}^i \leftarrow (\mathbf{K}\mathcal{Y})_{(i)} \mathbf{U}_{\odot i} (\tilde{\mathbf{U}}_{\odot i}^\top \mathbf{U}_{\odot i})^+$$

end for

$$\mathcal{C} \leftarrow \llbracket \mathbf{U}^1, \dots, \mathbf{U}^{p+1} \rrbracket$$

until convergence of \mathcal{C}

Learning Algorithms

- TVK-RRR-T: Jointly learn \mathcal{C} and \mathcal{T} .

$$\begin{aligned} \arg \min_{\substack{\mathcal{C} \in \mathbb{R}^{N \times d_1 \times \dots \times d_p}, \\ \mathcal{T} \in \mathbb{R}^{d_1 \times \dots \times d_p \times d_1 \times \dots \times d_p}}} \|\mathbf{K}\mathcal{C}\mathcal{T} - \mathbf{y}\|_F^2 + \gamma \langle \mathbf{K}\mathcal{C}\mathcal{T}, \mathcal{C} \rangle \quad \text{s.t.} \quad & \text{rank}(\mathcal{T}) \leq R, \\ & \mathcal{T} \text{ is } p.s.d. \end{aligned}$$

⇒ Alternating minimisation (\mathcal{C} and $\mathcal{T} = \llbracket \mathbf{V}^1, \dots, \mathbf{V}^p, \mathbf{V}^1, \dots, \mathbf{V}^p \rrbracket$) using gradient descent

Algorithm 2 TVK-RRR-T

Require: Gram matrix $\mathbf{K} \in \mathbb{R}^N \times \mathbb{R}^N$, output tensor $\mathcal{Y} \in \mathbb{R}^N \otimes \mathcal{Y}$, regularization parameter R .

Ensure: $\mathcal{C} \in \mathbb{R}^N \otimes \mathcal{Y}$, $\mathcal{T} \in \mathcal{Y} \otimes \mathcal{Y}$ a p.s.d. tensor.

Randomly initialize $\mathbf{V}^i \in \mathbb{R}^{d_i \times R}$ for $1 \leq i \leq p$

repeat

$\mathcal{T} \leftarrow \llbracket \mathbf{V}^1, \dots, \mathbf{V}^p, \mathbf{V}^1, \dots, \mathbf{V}^p \rrbracket$

$\mathcal{C} \leftarrow$ solution of the Sylvester equation $\mathbf{K}\mathcal{C}\mathcal{T} + \gamma\mathcal{C} - \mathcal{Y} = \mathbf{0}$

for $i = 1, \dots, p$ **do**

$\mathbf{V}^i \leftarrow \arg \min_{\mathbf{V}^i} \frac{1}{2} \|\mathbf{K}\mathcal{C}\mathcal{T} - \mathcal{Y}\|^2 + \frac{\gamma}{2} \langle \mathbf{K}\mathcal{C}\mathcal{T}, \mathcal{C} \rangle$ (solved by gradient descent).

$\mathcal{T} \leftarrow \llbracket \mathbf{V}^1, \dots, \mathbf{V}^p, \mathbf{V}^1, \dots, \mathbf{V}^p \rrbracket$

end for

until convergence of \mathcal{C} and \mathcal{T}

Outline

- 1 Introduction
- 2 Regression and Tensors
- 3 Nonparametric Tensor-valued Reduced Rank Regression
- 4 Experiments**
- 5 Conclusion and Perspectives

Synthetic Data

Generate data from the relation

$$\mathbf{y}^{(n)} = f(\mathbf{x}^{(n)}) + \boldsymbol{\xi}$$

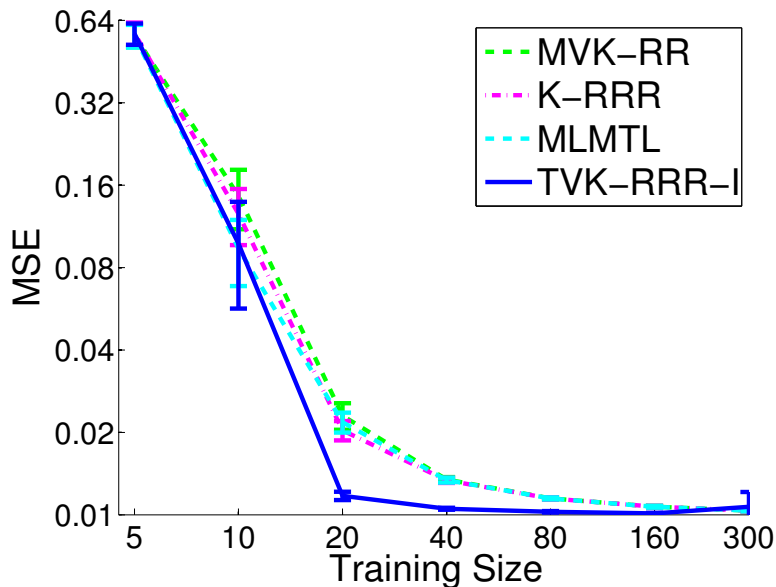
where

$$f : \mathbb{R}^{10} \rightarrow \mathbb{R}^{5 \times 5 \times 5}$$

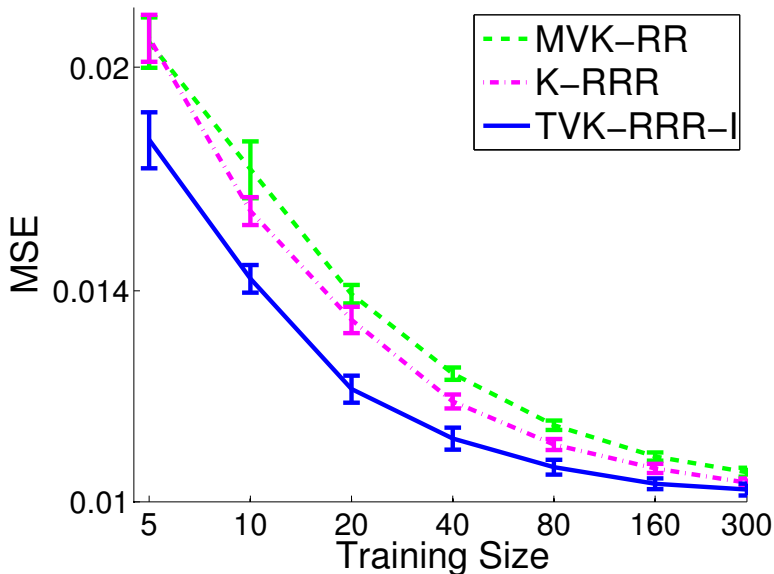
$$\mathbf{x} \mapsto \sum_{i=1}^{100} k(\mathbf{b}_i, \mathbf{x}) \mathcal{T} \mathcal{C}^{(i)}$$

- Combination of 100 basis functions $k(\mathbf{b}_i, \cdot) \mathcal{T} \mathcal{C}^{(i)}$ (drawn at random)
- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a scalar valued kernel
- $\mathcal{C} \in \mathbb{R}^{100 \times 5 \times 5 \times 5}$ and $\mathcal{T} \in \mathbb{R}^{5 \times 5 \times 5 \times 5 \times 5}$ are (low-rank) tensors drawn at random.
- $\mathbf{x}^{(n)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- $\xi_{i_1 i_2 i_3} \sim \mathcal{N}(0, 0.1)$

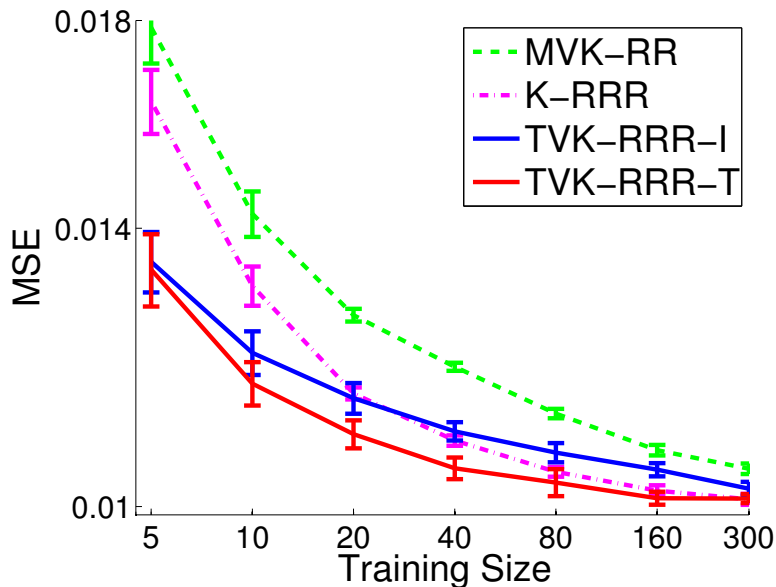
Linear kernel, $\mathcal{T} = \mathcal{J}$, $\text{rank}(\mathcal{C}) = 10$



RBF kernel, $\mathcal{T} = \mathcal{J}$, $\text{rank}(\mathcal{C}) = 10$



RBF kernel, $\text{rank}(\mathcal{T}) = 5$, $\text{rank}(\mathcal{C}) = 100$



Real Data: Meteo Forecasting

- Data from the meteorological office of the UK¹
- Monthly average measurements of 5 variables in 16 stations across the UK from 1960 to 2000
- Predict the 5 variables in the 16 stations from values in the 5 preceding months
- Training set and validation set size: 24 months

¹<http://www.metoffice.gov.uk/public/weather/climate-historic/>

Real Data: Meteo Forecasting

| <i>5-to-1</i> | MLMTL | MVK-RR-I | K-RRR | TVK-RRR-I | TVK-RRR-T |
|-----------------|--------------------|--------------------|--------------------|---------------------------|---------------------------|
| linear | 0.3511 \pm 0.018 | 0.3701 \pm 0.024 | 0.3514 \pm 0.021 | 0.3443 \pm 0.019 | 0.3446 \pm 0.019 |
| RBF | - | 0.3501 \pm 0.012 | 0.3439 \pm 0.014 | 0.3316 \pm 0.011 | 0.3315 \pm 0.011 |
| poly | - | 0.4219 \pm 0.019 | 0.4112 \pm 0.017 | 0.4032 \pm 0.019 | 0.4033 \pm 0.018 |
| <i>time (s)</i> | 14.470 | 0.028 | 0.029 | 12.682 | 56.464 |

| <i>5-to-4</i> | MLMTL | MVK-RR-I | K-RRR | TVK-RRR-I | TVK-RRR-T |
|-----------------|--------------------|--------------------|---------------------------|--------------------|---------------------------|
| linear | 0.3751 \pm 0.009 | 0.3855 \pm 0.008 | 0.3641 \pm 0.007 | 0.3589 \pm 0.008 | 0.3587 \pm 0.008 |
| RBF | - | 0.3630 \pm 0.008 | 0.3595 \pm 0.007 | 0.3492 \pm 0.007 | 0.3487 \pm 0.008 |
| poly | - | 0.4121 \pm 0.017 | 0.3970 \pm 0.015 | 0.3984 \pm 0.017 | 0.3981 \pm 0.014 |
| <i>time (s)</i> | 61.497 | 0.043 | 0.038 | 25.151 | 262.284 |

Table : MSE (mean \pm std) and average running time for the forecasting task: (top) 1 month prediction, (bottom) 4 months prediction.

Outline

- 1 Introduction
- 2 Regression and Tensors
- 3 Nonparametric Tensor-valued Reduced Rank Regression
- 4 Experiments
- 5 Conclusion and Perspectives

Conclusion and Perspectives

- Nonlinear / Nonparametric extension of low rank tensor learning techniques
- Introduction of Tensor-valued Kernels

- Other form of low-rankness (e.g. Tucker)
- Real-world data which could benefit from nonlinear models

Conclusion and Perspectives

- Nonlinear / Nonparametric extension of low rank tensor learning techniques
- Introduction of Tensor-valued Kernels

- Other form of low-rankness (e.g. Tucker)
- Real-world data which could benefit from nonlinear models

Thank you for your attention.

- Bahadori, M. T., Yu, Q. R., and Liu, Y. (2014). Fast multivariate spatio-temporal analysis via low rank tensor learning. In *NIPS*.
- Izenman, A. J. (1975). Reduced-rank regression for the multivariate linear model. *Journal of Multivariate Analysis*, 5(2):248–264.
- Micchelli, C. A. and Pontil, M. (2005). On learning vector-valued functions. *Neural Computation*, 17:177–204.
- Romera-Paredes, B., Aung, M. H., Bianchi-Berthouze, N., and Pontil, M. (2013). Multilinear multitask learning. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, pages 1444–1452.
- Signoretto, M., De Lathauwer, L., and Suykens, J. K. (2013). Learning tensors in reproducing kernel hilbert spaces with multilinear spectral penalties. *arXiv preprint arXiv:1310.4977*.