

Kernel based identification of systems with multiple outputs using nuclear norm regularization

Tillmann Falck, Bart De Moor and Johan A. K. Suykens

KU Leuven, ESAT-SCD and iMinds Future Health Department

International Workshop on Advances in Regularization,
Optimization, Kernel Methods and Support Vector Machines:
theory and applications

July 8 - 10, 2013; Leuven, Belgium

Outline

Introduction

Model derivation

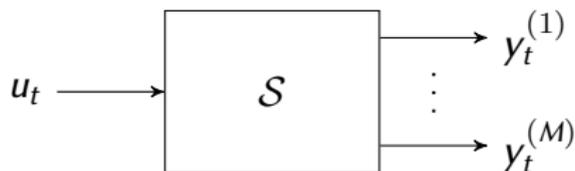
Example

Optimization aspects

Conclusions and outlook

Nonlinear system identification with multiple outputs

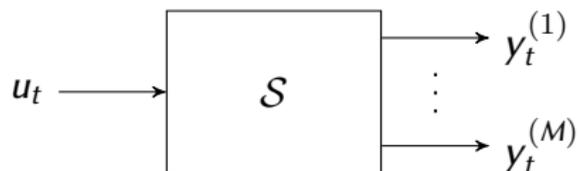
Overview



- ▶ Given input-output data $(u_t, y_t)_{t=1}$ with $y_t = [y_t^{(1)}, \dots, y_t^{(M)}]^T$
- ▶ Estimate a mathematical model for \mathcal{S}
- ▶ Often solved via regression:
 $(y_{t-1}, \dots, y_{t-p}, u_t, \dots, u_{t-Q}) \rightarrow y_t$
- ▶ Kernel based methods and support vector techniques in particular quite successful
- ▶ Applications: load/demand forecasting, virtual sensors, ...

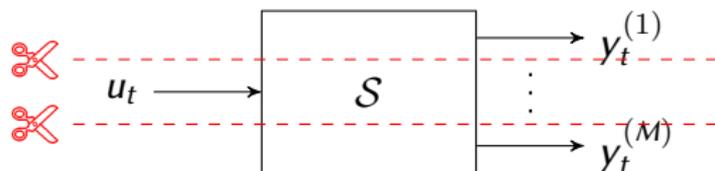
Nonlinear system identification with multiple outputs

Traditional versus proposed approach



Nonlinear system identification with multiple outputs

Traditional versus proposed approach



Traditional approach



Nonlinear system identification with multiple outputs

Traditional versus proposed approach

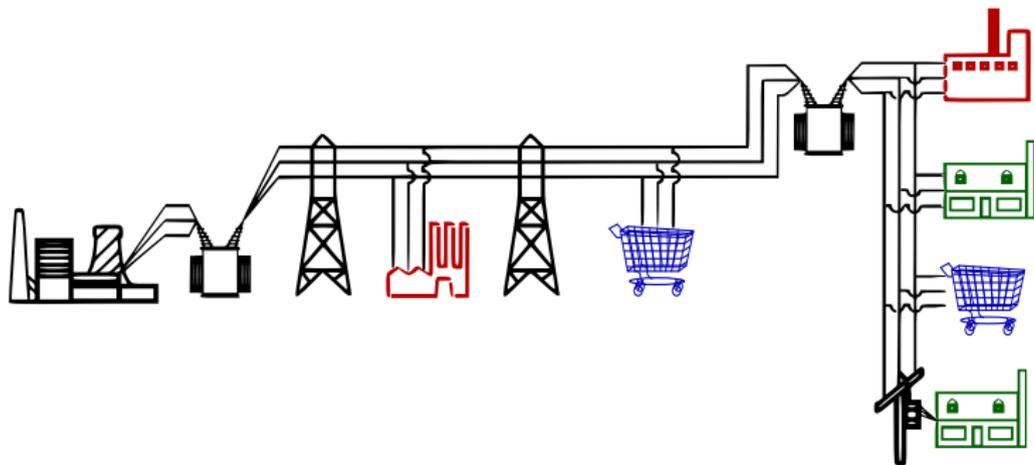


Traditional approach



Proposed approach

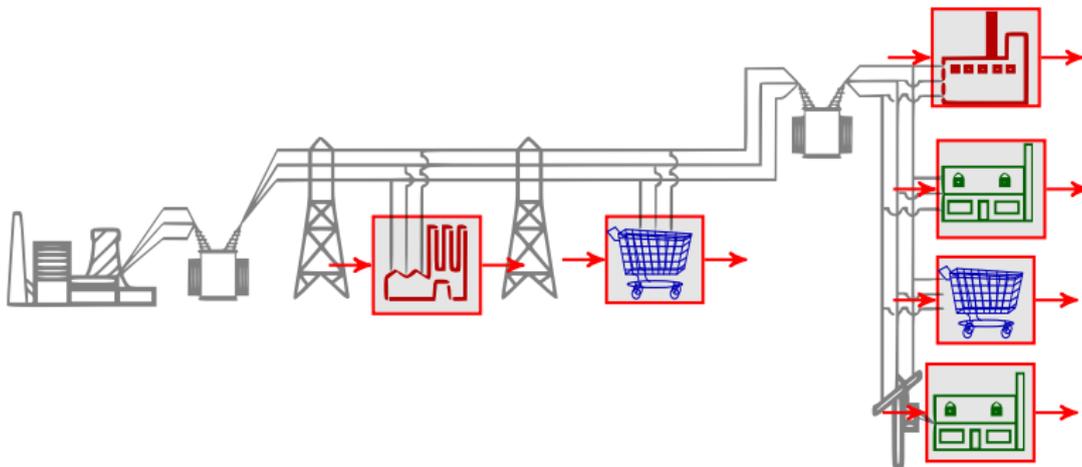
Example for multiple output system



Based on public domain image by United States Department of Energy

C. Alzate, M. Espinoza, B. De Moor, and J.A.K. Suykens (2009). "Identifying customer profiles in power load time series using spectral clustering". In: Proc. of the 19th International Conference on Artificial Neural Networks.

Example for multiple output system



Based on public domain image by United States Department of Energy

C. Alzate, M. Espinoza, B. De Moor, and J.A.K. Suykens (2009). "Identifying customer profiles in power load time series using spectral clustering". In: Proc. of the 19th International Conference on Artificial Neural Networks.

Key contributions and challenges

Contributions

- ▶ Kernel based model for nonlinear systems with multiple related outputs
- ▶ New primal-dual derivation of kernel based model with nuclear norm regularization

Challenges

- ▶ Finding a kernel based problem formulation
- ▶ Connecting dual, kernel based solution, to original model
- ▶ Numerical solution

Nuclear or trace norm

$$\|\mathbf{W}\|_*$$

Basic properties

- ▶ matrix norm
- ▶ sum of singular values
- ▶ induces sparsity
 - ▶ can be interpreted as ℓ_1 -norm of singular values
 - ▶ promotes low-rank solutions

As regularization term in a support vector model

- ▶ columns w_i model parameters, i.e. $y_t^{(i)} = w_i^T \varphi(x_t) + b$
- ▶ promotes relations between models in feature space

M. Fazel (2002). "Matrix Rank Minimization with Applications". PhD thesis, Stanford.

Model formulation in a primal-dual setting

Primal model

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{b}, \mathbf{e}_t} \quad & \eta \|\mathbf{W}\|_* + \sum_{t=1}^N \mathbf{e}_t^T \mathbf{e}_t \\ \text{subject to} \quad & y_t^{(i)} = \mathbf{w}_i^T \boldsymbol{\varphi}(\mathbf{x}_t) + b_i + e_t^{(i)}, \\ & t = 1, \dots, N, i = 1, \dots, M \end{aligned}$$

Derivation in primal-dual setting

- Write down Lagrangian for primal problem
- Take derivatives with respect to optimization variables
- Formulate KKT conditions for optimality
- Write down dual optimization problem
- Substitute dual solution into primal model

A. Argyriou, C.A. Micchelli, and M.A. Pontil (2009). "When Is There a Representer Theorem? Vector Versus Matrix Regularizers". In: Journal of Machine Learning Research.

Model formulation in a primal-dual setting

Primal model

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{b}, \mathbf{e}_t} \quad & \eta \|\mathbf{W}\|_* + \sum_{t=1}^N \mathbf{e}_t^T \mathbf{e}_t \\ \text{subject to} \quad & y_t^{(i)} = \mathbf{w}_i^T \boldsymbol{\varphi}(\mathbf{x}_t) + b_i + e_t^{(i)}, \\ & t = 1, \dots, N, i = 1, \dots, M \end{aligned}$$

Derivation in primal-dual setting

- (a) Write down Lagrangian for primal problem
- (b) Take derivatives with respect to optimization variables
- (c) Formulate KKT conditions for optimality
- (d) Write down dual optimization problem
- (e) Substitute dual solution into primal model

A. Argyriou, C.A. Micchelli, and M.A. Pontil (2009). "When Is There a Representer Theorem? Vector Versus Matrix Regularizers". In: Journal of Machine Learning Research.

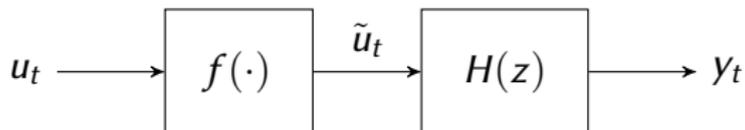
Advantage of primal-dual approach

Straightforward to incorporate additional structure

Example: simple constraints

- ▶ $w_1^T \varphi(x) = -w_2^T \varphi(x)$
- ▶ $w_i^T \varphi(x) \geq y_0$

Example: Hammerstein systems



- ▶ Static nonlinearity $f(x) = w^T \varphi(x) + c$
- ▶ Linear dynamical system $H : \hat{y}_t = \sum_{p=1}^P a_p y_{t-p} + \sum_{q=0}^Q b_q \tilde{u}_{t-q}$
- ▶ Approximate joint model

$$\hat{y}_t = \sum_{p=1}^P a_p y_{t-p} + \sum_{q=0}^Q w_q^T \varphi(u_{t-q}) + \tilde{c}$$

J.A.K. Suykens, C. Alzate, and K. Pelckmans (2010). "Primal and dual model representations in kernel-based learning". In: Statistics Surveys.

Derivation of kernel based model

From Lagrangian to KKT conditions

Lagrangian

$$\mathcal{L} = \eta \|\mathbf{W}\|_* + \sum_{t=1}^N \mathbf{e}_t^T \mathbf{e}_t - \sum_{t=1}^N \boldsymbol{\alpha}^T (\mathbf{W}^T \boldsymbol{\varphi}(\mathbf{x}_t) + \mathbf{b} + \mathbf{e}_t - \mathbf{y}_t)$$

Nuclear norm is not differentiable!

Possible reformulations

- ▶ Dual norm, where $\|\cdot\|_2$ is dual norm of $\|\cdot\|_*$

$$\|\mathbf{W}\|_* = \max_{\|\mathbf{C}\|_2 \leq 1} \text{tr}(\mathbf{C}^T \mathbf{W})$$

- ▶ Conic duality, where \mathcal{K} is convex cone $\{(\mathbf{X}, s) \mid \|\mathbf{X}\|_* \leq s\}$

$$\|\mathbf{W}\|_* = t \text{ where } (\mathbf{W}, t) \in \mathcal{K}$$

Derivation of kernel based model

Kernel based optimization problem

$$\begin{aligned} \max_A \quad & \text{tr}(A^T Y) - \frac{1}{2} \text{tr}(A^T A) \\ \text{subject to} \quad & A^T \Omega A \preceq \eta I_M, \\ & A^T \mathbf{1}_N = \mathbf{0}_M \end{aligned}$$

- ▶ $A = [\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_N]^T, Y = [y_1, \dots, y_N]^T$
- ▶ $\Omega_{ij} = K(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\varphi}(\mathbf{x}_i)^T \boldsymbol{\varphi}(\mathbf{x}_j)$

KKT condition for W

$$C = \sum_{t=1}^N \boldsymbol{\varphi}(\mathbf{x}_t) \boldsymbol{\alpha}_t^T$$

- ▶ $W, C \in \mathbb{R}^{n_h \times M}, \boldsymbol{\varphi}(\mathbf{x}) \in \mathbb{R}^{n_h}$ and $\boldsymbol{\alpha}_t \in \mathbb{R}^M$
- ▶ No expansion of primal variables W!

Formulation of model in terms of dual solution

Characterization of solution set

$$\begin{aligned}\{\mathbf{W} : \text{tr}(\mathbf{W}^T \mathbf{C}) = \xi, \|\mathbf{W}\|_* = \xi\} \\ = \{\mathbf{U}_\eta \mathbf{H}_\eta \mathbf{V}_\eta^T : \text{tr}(\mathbf{H}_\eta) = \xi, 0 \preceq \mathbf{H}_\eta \in \mathbb{R}^{r \times r}\}\end{aligned}$$

- ▶ Here \mathbf{C} and ξ are fixed
- ▶ $\mathbf{C} = \mathbf{\Phi} \mathbf{A}^T$
- ▶ $\mathbf{U}_\eta, \mathbf{V}_\eta$ contain left and right singular vectors corresponding to largest singular value of \mathbf{C} respectively

Connecting \mathbf{W} and \mathbf{A}

$$\begin{aligned}\text{find } & \mathbf{H}_\eta \\ \text{subject to } & \mathbf{H}_\eta \succeq 0, \text{tr}(\mathbf{H}_\eta) = \xi \\ & \mathbf{y}^{(i)} = [\mathbf{\Omega}_{i,1} \boldsymbol{\alpha}_1, \dots, \mathbf{\Omega}_{i,M} \boldsymbol{\alpha}_M] \mathbf{V}_\eta \mathbf{H}_\eta \mathbf{V}_\eta^T \boldsymbol{\epsilon}_i \\ & + \mathbf{b}_i \mathbf{1}_{N_i} + \boldsymbol{\alpha}_i, \quad i = 1, \dots, M\end{aligned}$$

Model representation

Primal model representation

$$\hat{y}^{(i)} = f(z) = \mathbf{w}_i^T \boldsymbol{\varphi}(z) + b_i$$

Dual model representation

$$\hat{y}^{(i)} = f(z) = \sum_{j=1}^M Q_{ji} \sum_{t=1}^N A_{tj} K(\mathbf{x}_t, z) + b_i$$

with $Q = V_{\eta} H_{\eta} V_{\eta}^T$

One dimensional SVM

$$\hat{y} = f(z) = \sum_{t=1}^N \alpha_t K(\mathbf{x}_t, z) + b$$

Example

Description

Toy example

- ▶ number of data: training 50, validation 100, test 150
- ▶ number of outputs: 20
- ▶ number of independent contributions: 3
- ▶ data generation: $Y = W_0^T \Phi + \text{noise}$
- ▶ $\Phi \in \mathbb{R}^{50 \times 50}$, $W_0 = \sum_{i=1}^3 g_i r_i^T = GR^T$, $G \in \mathbb{R}^{50 \times 3}$, $R \in \mathbb{R}^{20 \times 3}$

Evaluated models

MIMO proposed nuclear norm regularized model

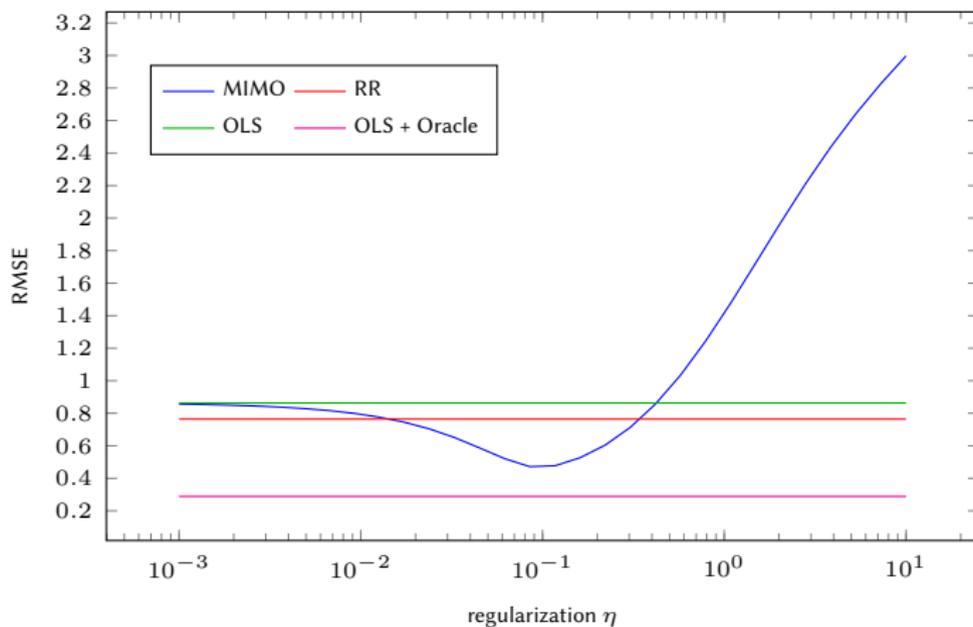
RR ridge regression model (LS-SVM model in primal with given feature map) with independent LS-SVM models for each output

OLS ordinary least squares model

OLS + oracle OLS given the true structure of problem

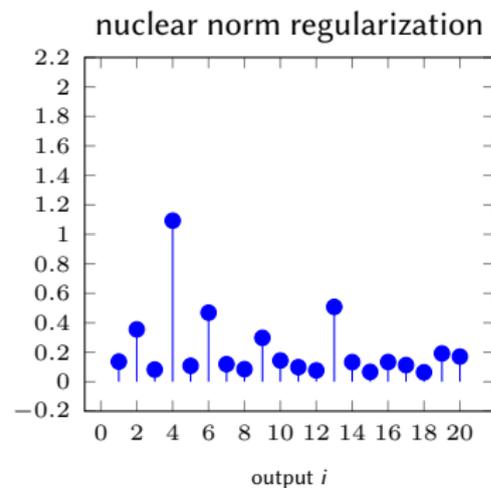
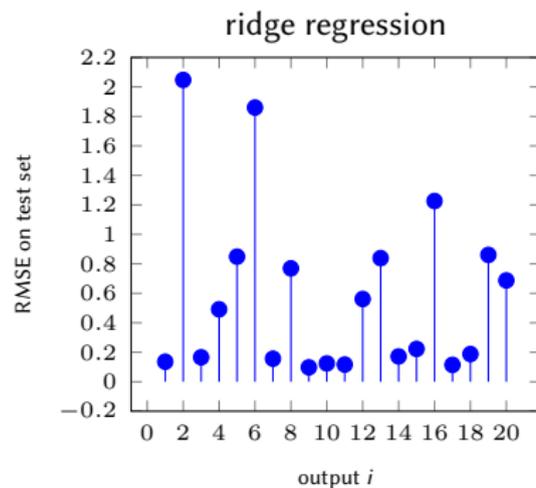
Example

Cross validation



Example

Results



Comparison of optimization problems

Primal formulation

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{b}, \mathbf{e}_t} \quad & \eta \|\mathbf{W}\|_* + \sum_{t=1}^N \mathbf{e}_t^T \mathbf{e}_t \\ \text{subject to} \quad & \mathbf{y}_t^{(i)} = \mathbf{w}_i^T \boldsymbol{\varphi}(\mathbf{x}_t) + b_i + \mathbf{e}_t^{(i)} \end{aligned}$$

Dual formulation

Optimization problem

$$\begin{aligned} \max_{\mathbf{A}} \quad & \text{tr}(\mathbf{A}^T \mathbf{Y}) - \frac{1}{2} \text{tr}(\mathbf{A}^T \mathbf{A}) \\ \text{subject to} \quad & \mathbf{A}^T \boldsymbol{\Omega} \mathbf{A} \preceq \eta \mathbf{I}_M, \\ & \mathbf{A}^T \mathbf{1}_N = \mathbf{0}_M \end{aligned}$$

Reconstruction of model

$$\begin{aligned} \text{find} \quad & \mathbf{H}_\eta \\ \text{subject to} \quad & \mathbf{H}_\eta \succeq 0, \text{tr}(\mathbf{H}_\eta) = \xi \\ & \mathbf{y}^{(i)} = [\boldsymbol{\Omega}_{i,1} \boldsymbol{\alpha}_1, \dots, \boldsymbol{\Omega}_{i,M} \boldsymbol{\alpha}_M] \\ & \quad \cdot \mathbf{V}_\eta \mathbf{H}_\eta \mathbf{V}_\eta^T \boldsymbol{\epsilon}_i + b_i \mathbf{1}_{N_i} + \boldsymbol{\alpha}_i \end{aligned}$$

Dimensionalities of optimization problems

	One input	L inputs
Primal		
Feature map	$n_h \times N$	$L \cdot (n_h \times N)$
Unknown W	$n_h \times M$	$L \cdot (n_h \times M)$
Dual		
Kernel matrix	$N \times N$	$(L \cdot N) \times (L \cdot N)$
Unknown A	$N \times M$	$(L \cdot N) \times M$

- ▶ N : number of data
- ▶ M : number of outputs
- ▶ n_h : dimension of feature map

Solution strategies

SDP formulation

- o Problems are SDP representable
- o Can be solved with general purpose SDP solvers
- + Small implementation effort
- + High accuracy
- High runtime costs, memory & times \Rightarrow Limited to small problem sizes

First order techniques

- o (Accelerated) gradient projection can be applied
- o Higher implementation effort
- + Structure can be exploited
- + Scale to larger problem sizes
- Lower accuracy (crucial for reconstruction of dual model)

Conclusions

- ▶ Proposed novel identification for nonlinear systems with multiple outputs
 - ▶ Exploits relations between output variables
 - ▶ Illustrated improvement on small toy example
- ▶ Presented derivation of a nuclear norm regularized model in primal-dual setting
 - ▶ Allows straightforward integration of additional information

Challenges and Outlook

- ▶ Application on real world datasets
- ▶ Numerical solution on larger datasets
 - ▶ New algorithms are already being developed
 - ▶ Computational power increases exponentially
- ▶ Same primal-dual derivation can be applied to other nonquadratic regularization schemes
- ▶ Conjecture
 - ▶ Promising applications in system identification
 - ▶ Kernel based models can be used for many applications besides regression, these might also benefit from advanced regularization schemes