

---

# Tensorized LSSVMs for Multitask Regression

---

Jiani Liu<sup>1</sup>, Qinghua Tao<sup>2</sup>, Ce Zhu<sup>1\*</sup>, Yipeng Liu<sup>1</sup>, Johan A.K. Suykens<sup>2</sup>

<sup>1</sup> School of Information and Communication Engineering, University of Electronic Science and Technology of China, 610054 Chengdu, China

<sup>2</sup> ESAT-STADIUS, KU Leuven, 3001 Heverlee, Belgium

## ABSTRACT

Multitask learning (MTL) can utilize the relatedness between multiple tasks for performance improvement. The advent of multimodal data allows tasks to be referenced by multiple indices. High-order tensors are capable of providing efficient representations for such tasks, while preserving structural task-relations. In this paper, a new MTL method is proposed by leveraging low-rank tensor analysis and constructing tensorized Least Squares Support Vector Machines, namely the tLSSVM-MTL, where multilinear modelling and its nonlinear extensions can be flexibly exerted. We employ a high-order tensor for all the weights with each mode relating to an index and factorize it with CP decomposition, assigning a shared factor for all tasks and retaining task-specific latent factors along each index. Then an alternating algorithm is derived for the nonconvex optimization, where each resulting subproblem is solved by a linear system. Experimental results demonstrate promising performances of our tLSSVM-MTL.

**Keywords** Multitask learning, tensor regression, CP decomposition, LSSVM, shared factor

## 1 Introduction

Multitask learning (MTL) lies on the exploitation of the coupling information across different tasks, so as to benefit the parameter estimation for each individual task [1, 3, 20]. MTL has been widely applied in many fields, such as social sciences [5, 6, 19], medical diagnosis [7, 14], etc. Various MTL methods have been developed and shown promising performance for related tasks. Among them, support vector machines (SVMs) get great success [4]. Specifically, based on the minimization of regularization functionals, the regularized MTL is proposed in [5, 6] with kernels including a task-coupling parameter. An MTL method based on SVM+, as an extension of SVM, is developed in [8] and compared with standard SVMs in [7] and regularized MTL in [14]. Moreover, the least squares SVM (LSSVM) [15] is also generalized for MTL [19], where the inequality constraints in SVMs are modified into equality ones and a linear system is solved in dual instead of the typical quadratic programming. These SVM-based MTL methods were all applied with the typical vector/matrix expressions.

Tensors, a natural extension for vectors and matrices, provide a more effective way to preserve multimodal information and describe complex dependencies [9, 10]. Different usages of tensor representations have been successfully applied to MTL [13, 17, 18, 21–25]. For instance, motivated by the multidimensional input, [25] proposed to factorize the weight tensor for each task into a sparse task-specific part and a low rank shared part. In [18], it formulates the input as a tensor and extracts its spatial and temporal latent factors, based on which a prediction model is built. It is also

---

\*Corresponding Author. This research is partially supported by the National Natural Science Foundation of China (NSFC) under Grant U19A2052, Grant 62020106011 and Grant 62171088. Johan A.K. Suykens and Qinghua Tao acknowledge the supports from iBOF project Tensor Tools for Taming the Curse (3E221427), European Research Council (ERC) Advanced Grant E-DUALITY (787960), KU Leuven Grant CoE PFV/10/002, Grant FWO GOA4917N, EU H2020 ICT-48 Network TAILOR, and Leuven.AI Institute.

Email: jianiliu@std.uestc.edu.cn, qinghua.tao@esat.kuleuven.be, eczhu@uestc.edu.cn, yipengliu@uestc.edu.cn, johan.suykens@esat.kuleuven.be

intriguing to encode the projection matrices of all classifiers into tensors and apply tensor nuclear norm constraints for task relations [21, 23].

The aforementioned works are all set with a single index for the involved tasks. In practice, tasks can be referenced by multiple indices with physical meanings. Taking a multimodal data task for example, restaurant recommendations consider different aspects of rating (e.g., food and service) and customers. It naturally leads to  $T_1 \times T_2$  tasks spanned by two indices, and thus a single index fails to preserve such information. Therefore, [13] considered tasks with multiple indices and imposed low Tucker rank regularization over the stacked coefficient tensor to explore task relations. In [13], the applied Tucker decomposition can suffer from a dimensionality curse if the tensor order increases. For rank minimization, a convex relaxation is used to handle the whole weight tensor in each iteration and thereby can be problematic for large-scale data. Two variants were later developed in [17, 24] with different convex relaxations for Tucker rank minimization. Though nonconvex optimization was also considered in [13], it required adjusting several ranks within Tucker, making the tuning procedures rather complicated. Besides, they all considered multilinear modelling, while nonlinearity is highly desirable for well describing complex data and tasks.

In this paper, we develop a tensorized MTL method for regression by leveraging LSSVMs, namely the tLSSVM-MTL, which constructs a high-order weight tensor on LSSVMs and indexes the tasks along different modes into groups by multiple indices. Unlike [13, 17, 24], we factorize the constructed tensor into CP forms since the factors are easy to explain from subspace perspective, and enable all tasks to share a common latent factor and meanwhile retain task-specific factors. In our method, both linear and nonlinear feature maps (or kernels) can be flexibly employed. For optimization, an alternating minimization strategy is proposed with each subproblem solved by a linear system in the dual. Numerical experiments show advantageous performances of our tLSSVM-MTL over matrix-based and existing tensor-based MTL methods.

The next section gives some premieres. Section 3 presents the modelling and optimization for our tLSSVM-MTL. Experimental results and conclusions are in Sections 4 and 5.

## 2 Preliminaries

Scalars, vectors, matrices, and tensors are represented as  $x$ ,  $\mathbf{x}$ ,  $\mathbf{X}$ , and  $\mathcal{X}$ , respectively. For clarity, we denote the row and the column in a matrix  $\mathbf{X}$  as  $\mathbf{X}[i, :]^T = \mathbf{x}_{i,:}$  and  $\mathbf{X}[:, j] = \mathbf{x}_{:,j}$ .

**CP decomposition** [2, 12] Given a tensor  $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ , CP decomposition factorizes the tensor into a summation of several rank-one components as  $\mathcal{X} = \sum_{k=1}^K \mathbf{u}_k^1 \circ \dots \circ \mathbf{u}_k^N$ , where  $K$  is the CP rank indicating the smallest number of rank-one components required in this representation. We represent the CP decomposition as  $\mathcal{X} = \llbracket \mathbf{U}^1, \dots, \mathbf{U}^N \rrbracket$  with  $\mathbf{U}^n = [\mathbf{u}_1^n, \dots, \mathbf{u}_K^n]$  for  $n = 1, \dots, N$ .

**LSSVM** LSSVM [15] is a variant of SVMs [4] by forming equality constraints. For regression with data  $\{\mathbf{x}_i, y_i\}_{i=1}^m$ , the primal problem of LSSVM is given as:

$$\begin{aligned} \min_{\mathbf{w}, b, \mathbf{e}} \quad & J(\mathbf{w}, b, \mathbf{e}) = \frac{C}{2} \sum_{i=1}^m (e_i)^2 + \frac{1}{2} \mathbf{w}^\top \mathbf{w} \\ \text{s. t.} \quad & \mathbf{w}^\top \phi(\mathbf{x}_i) + b = y_i - e_i, \end{aligned}$$

where  $\phi: \mathbb{R}^d \mapsto \mathbb{R}^{d_h}$  is the feature mapping function,  $\mathbf{w} \in \mathbb{R}^{d_h}$  and  $b \in \mathbb{R}$  are the modelling coefficients,  $e_i$  denotes the point-wise regression error, and  $C > 0$  is the regularization hyperparameter. In LSSVMs, the Lagrangian dual problem gives a linear system, instead of the quadratic programming in classic SVMs, making certain problems more tractable.

## 3 Tensorized LSSVMs for MTL

### 3.1 Tensorized Modelling

Assuming  $T$  tasks are involved with data  $\{\mathbf{x}_i^t \in \mathbb{R}^{d_t}, y_i^t \in \mathbb{R}\}_{i=1}^{m_t}$ ,  $T$  sets of parameters  $\{\mathbf{w}_t, b^t\}_{t=1}^T$  are thereby required for predictions in MTL. Here we focus on homogeneous attributes with  $d_t = d$ . Thus, the complete weight matrix is  $\mathbf{W} = [\mathbf{w}_1; \dots; \mathbf{w}_T]$ . Instead of using a single index for these  $T$  tasks, multiple indices for an efficient and structured representation can be considered to construct a higher-order tensor [13]. In this paper, the weight tensor is constructed as  $\mathcal{W} \in \mathbb{R}^{d_h \times T_1 \times \dots \times T_N}$  and we factorize it into CP form for the structural relatedness across different tasks, such that:

$$\mathcal{W} = \sum_{k=1}^K \mathbf{l}_{:,k} \circ \mathbf{u}_{:,k}^1 \circ \dots \circ \mathbf{u}_{:,k}^N = \llbracket \mathbf{L}, \mathbf{U}^1, \dots, \mathbf{U}^N \rrbracket, \quad (1)$$

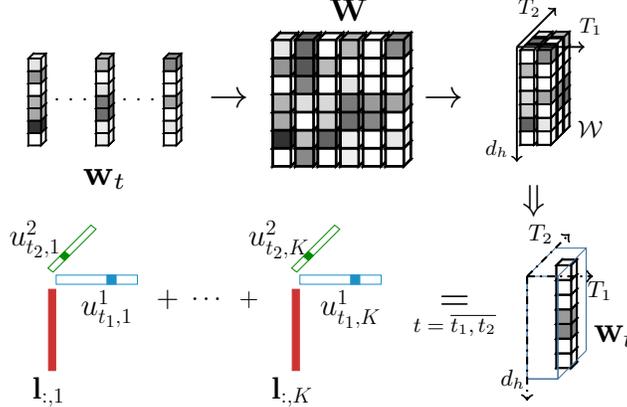


Figure 1: An illustration on our tensorized representations.

where  $\mathbf{L} = [\mathbf{l}_{:,1}; \dots; \mathbf{l}_{:,K}] \in \mathbb{R}^{d_h \times K}$  is the shared factor exploiting coupling information across tasks,  $\mathbf{U}^n = [\mathbf{u}_{:,1}^n; \dots; \mathbf{u}_{:,K}^n] \in \mathbb{R}^{T_n \times K}$  corresponds to the  $n$ -th index with  $\mathbf{u}_{:,k}^n = [u_{1,k}^n, \dots, u_{T_n,k}^n]^\top$ . The task-specific coefficient is thus formulated as:

$$\mathbf{w}_t = \sum_{k=1}^K \mathbf{l}_{:,k} \cdot u_{t_1,k}^1 \dots u_{t_N,k}^N. \quad (2)$$

Each task is now spanned by  $N$  indices, i.e.,  $t = \overline{t_1, \dots, t_N}$  with  $t_n = 1, \dots, T_n$ ,  $n = 1, \dots, N$ , so that the total number of tasks is calculated by  $T = \prod_{n=1}^N T_n$ . Fig. 1 gives a graphical illustration for a third-order case.

It is explicit that  $\{\mathbf{l}_{:,1}, \dots, \mathbf{l}_{:,K}\}$  learns the coupling information across tasks and is always involved in the prediction for each task. In contrast, the variation of  $\mathbf{u}_{t_n,:}^n$  affects a certain group of tasks relating to the index  $t_n$ . For instance, for  $n = 1$ ,  $t_1 = 1$ , the updating of  $\mathbf{u}_{1,:}^1$  affects tasks in  $\{t = \overline{1, \dots, t_N} | t_1 = 1, \dots, T_1, l \neq 1\}$ . In other words, the correlations between tasks can be explored by splitting them into different modes (indices) with a high-order tensor, enabling structural captures of dependencies from multiple modes than using a single mode. In this way, CP rank  $K$  indicates the number of latent shared features  $\mathbf{l}_{:,k}$  in this representation. With the imposed low CP rank, the learned coefficients can be more compact in gaining informative modelling.

Then, our tensorized LSSVM for MTL regression, i.e., tLSSVM-MTL, is constructed in the primal form as:

$$\begin{aligned} \min_{\mathbf{L}, \mathbf{U}^n, b^t, e_i^t} \quad & \frac{C}{2} \sum_{t=1}^T \sum_{i=1}^{m_t} (e_i^t)^2 + \frac{1}{2} \text{tr} \mathbf{L} \mathbf{L}^\top + \frac{1}{2} \sum_{n=1}^N \text{tr} \mathbf{U}^n \mathbf{U}^{n\top} \\ \text{s.t.} \quad & (\sum_{k=1}^K (\mathbf{l}_{:,k} \cdot u_{t_1,k}^1 \dots u_{t_N,k}^N))^\top \phi(\mathbf{x}_i^t) + b^t \\ & = y_i^t - e_i^t, \quad t = \overline{t_1, \dots, t_N}. \end{aligned} \quad (3)$$

With the constructed tensor and the deployed factorization, our proposed tLSSVM-MTL successfully extends the existing LSSVMs to deal with multitasks referenced by multiple indices; the low CP rank factorization enables to explicitly attain the shared factor  $\mathbf{L}$  seeking for common information and these  $\mathbf{U}^n$  maintaining task-specific information, which together boosts the overall performance of all tasks.

### 3.2 Optimization Algorithm

In (3), the product operations between the shared  $\mathbf{L}$  and the task-specific  $\mathbf{U}^1, \dots, \mathbf{U}^N$  result in nonconvexity, but can be decoupled by block coordinate descents. We thus design an alternating updating strategy to optimize each factor iteratively, where each subproblem successfully degenerates to be convex by solving a linear system with Lagrangian duality.

1) **Step  $\mathbf{L}, b^t, e_i^t$  with fixed  $\mathbf{U}^n$ .** The primal problem with respect to  $\mathbf{L}, b^t, e_i^t$  is given by

$$\begin{aligned} \min_{\mathbf{L}, b^t, e_i^t} \quad & \frac{C}{2} \sum_{t=1}^T \sum_{i=1}^{m_t} (e_i^t)^2 + \frac{1}{2} \text{tr}(\mathbf{L} \mathbf{L}^\top) \\ \text{s.t.} \quad & (\sum_{k=1}^K (\mathbf{l}_{:,k} \cdot u_{t,k}))^\top \phi(\mathbf{x}_i^t) + b^t = y_i^t - e_i^t, \end{aligned}$$

where  $u_{t,k} \triangleq u_{t_1,k}^1 \cdots u_{t_N,k}^N$  for  $t = \overline{t_1, \dots, t_N}$ ,  $t_n = 1, \dots, T_n$ . With dual variables  $\alpha_i^t \in \mathbb{R}$  corresponding to each equality constraint, the Lagrangian function is obtained as

$$\mathcal{L}(\mathbf{L}, b^t, e_i^t) = \frac{C}{2} \sum_{t=1}^T \sum_{i=1}^{m_t} (e_i^t)^2 + \frac{1}{2} \text{tr}(\mathbf{L}\mathbf{L}^\top) - \sum_{t=1}^T \sum_{i=1}^{m_t} \alpha_i^t ((\mathbf{L}\mathbf{u}_t)^\top \phi(\mathbf{x}_i^t) + b^t - y_i^t + e_i^t),$$

with  $\mathbf{u}_t \triangleq [u_{t,1}, \dots, u_{t,K}]^\top \in \mathbb{R}^K$ . Then, stationary point conditions are obtained as

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{L}} = 0 &\implies \mathbf{L} = \sum_{t=1}^T \sum_{i=1}^{m_t} \alpha_i^t \phi(\mathbf{x}_i^t) \mathbf{u}_t^\top, \\ \frac{\partial \mathcal{L}}{\partial \mathbf{b}} = 0 &\implies \mathbf{A}^\top \boldsymbol{\alpha} = 0, \quad \mathbf{b} = [b^1, \dots, b^T]^\top, \\ \frac{\partial \mathcal{L}}{\partial \mathbf{e}} = 0 &\implies C\mathbf{e} = \boldsymbol{\alpha}, \\ \frac{\partial \mathcal{L}}{\partial \boldsymbol{\alpha}} = 0 &\implies \Phi \mathbf{w} + \mathbf{A}\mathbf{b} = \mathbf{y} - \mathbf{e}. \end{aligned}$$

where  $\mathbf{A} = \text{blockdiag}(\mathbf{1}_{m_1}, \dots, \mathbf{1}_{m_T}) \in \mathbb{R}^{m \times T}$ ,  $\mathbf{w} = [(\mathbf{L}\mathbf{u}_1)^\top, \dots, (\mathbf{L}\mathbf{u}_T)^\top]^\top \in \mathbb{R}^{Td_h}$ , the task-specific feature mapping matrix  $\Phi^t = [\phi(x_1^t), \dots, \phi(x_{m_t}^t)]^\top \in \mathbb{R}^{m_t \times d_h}$  and  $\Phi = \text{blockdiag}(\Phi^1, \dots, \Phi^T) \in \mathbb{R}^{m \times Td_h}$  for all  $T$  tasks. All outputs, regression errors, and dual variables are denoted as  $\mathbf{y} = [y_1^1, y_2^1, \dots, y_{m_T}^T]^\top \in \mathbb{R}^m$ ,  $\mathbf{e} = [e_1^1, e_2^1, \dots, e_{m_T}^T]^\top \in \mathbb{R}^m$ , and  $\boldsymbol{\alpha} = [\alpha_1^1, \alpha_2^1, \dots, \alpha_{m_T}^T]^\top \in \mathbb{R}^m$ , respectively.

By eliminating  $\mathbf{L}$  and  $e_i^t$ , a linear system is attained as:

$$\left[ \begin{array}{c|c} \mathbf{0}_{T \times T} & \mathbf{A}^\top \\ \hline \mathbf{A} & \mathbf{Q} + \frac{1}{C} \mathbf{I}_{m \times m} \end{array} \right] \begin{bmatrix} \mathbf{b} \\ \boldsymbol{\alpha} \end{bmatrix} = \begin{bmatrix} \mathbf{0}_T \\ \mathbf{y} \end{bmatrix}, \quad (4)$$

where  $\mathbf{Q} \in \mathbb{R}^{m \times m}$  is computed by the components in tensor  $\mathcal{W}$  and the kernel function  $k: \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathbb{R}$  induced by  $\phi(\cdot)$ , such that  $\mathbf{Q}(j, j') = \langle \mathbf{u}_t, \mathbf{u}_q \rangle k(\mathbf{x}_i^t, \mathbf{x}_p^q)$ ,  $j = \sum_{r=1}^{t-1} m_r + i$ ,  $j' = \sum_{r=1}^{q-1} m_r + p$ ,  $i = 1, \dots, m_t$ ,  $p = 1, \dots, m_q$  with  $i, p$  indexing the samples in the involved tasks  $t$  and  $q$ , respectively. With the solution of dual variables (4), i.e.,  $\tilde{\boldsymbol{\alpha}}$ , we can get the updated  $\mathbf{L} = \sum_{t=1}^T \sum_{i=1}^{m_t} \tilde{\alpha}_i^t \phi(\mathbf{x}_i^t) \mathbf{u}_t^\top$ .

2) **Step  $\mathbf{U}^n, b^t, e_i^t$  with fixed  $\mathbf{L}$ .** With fixed  $\mathbf{L}$ , we alternate to optimize  $\mathbf{U}^n, b^t, e_i^t$ . The corresponding primal problem is:

$$\begin{aligned} \min_{\mathbf{u}_{t_n, \cdot}^n, b^t, e_i^t} \quad & \frac{C}{2} \sum_{t \in \mathcal{S}_{t_n}} \sum_{i=1}^{m_t} (e_i^t)^2 + \frac{1}{2} \|\mathbf{u}_{t_n, \cdot}^n\|_2^2 \\ \text{s. t.} \quad & \mathbf{u}_{t_n, \cdot}^n{}^\top \mathbf{z}_i^t + b^t = y_i^t - e_i^t, \end{aligned}$$

where  $\mathbf{z}_i^t$  is calculated by  $\mathbf{L}^\top \phi(\mathbf{x}_i^t) \odot \mathbf{u}_{t_1, \cdot}^1 \odot \cdots \odot \mathbf{u}_{t_{n-1}, \cdot}^{n-1} \odot \mathbf{u}_{t_{n+1}, \cdot}^{n+1} \odot \cdots \odot \mathbf{u}_{t_N, \cdot}^N \in \mathbb{R}^K$ , the involved tasks  $t$  is contained in the index set  $\mathcal{S}_{t_n} = \{t_1, \dots, t_N | t_l = 1, \dots, T_l, l = 1, \dots, N, l \neq n\}$  with cardinality  $|\mathcal{S}_{t_n}| = \prod_{l, l \neq n} T_l$ . With dual variables  $\boldsymbol{\lambda}_{t_n}$ , we have the Lagrangian function:

$$\mathcal{L}(\mathbf{u}_{t_n, \cdot}^n, b^t, e_i^t) = \frac{C}{2} \sum_{t \in \mathcal{S}_{t_n}} \sum_{i=1}^{m_t} (e_i^t)^2 + \frac{1}{2} \|\mathbf{u}_{t_n, \cdot}^n\|_2^2 - \sum_{t \in \mathcal{S}_{t_n}} \sum_{i=1}^{m_t} \lambda_i^t \left( (\mathbf{u}_{t_n, \cdot}^n{}^\top \mathbf{z}_i^t + b^t) - y_i^t + e_i^t \right),$$

where  $\boldsymbol{\lambda}_{t_n} = \{\lambda_i^t | t \in \mathcal{S}_{t_n}, i = 1, \dots, m_t\} \in \mathbb{R}^{M_{t_n}}$  corresponds to the involved constraints in optimizing  $\mathbf{u}_{t_n, \cdot}^n$ .

Similarly, by deriving the stationary conditions and eliminating  $\mathbf{u}_{t_n, \cdot}^n$  and  $e_i^t$  therein, we get the linear system:

$$\left[ \begin{array}{c|c} \mathbf{0}_{|\mathcal{S}_{t_n}| \times |\mathcal{S}_{t_n}|} & \mathbf{A}_{t_n}^\top \\ \hline \mathbf{A}_{t_n} & \mathbf{Q}_{t_n} + \frac{1}{C} \mathbf{I}_{M_{t_n}} \end{array} \right] \begin{bmatrix} \mathbf{b}_{t_n} \\ \boldsymbol{\lambda}_{t_n} \end{bmatrix} = \begin{bmatrix} \mathbf{0}_{|\mathcal{S}_{t_n}|} \\ \mathbf{y}_{t_n} \end{bmatrix}, \quad (5)$$

where  $\mathbf{A}_{t_n} = \text{blockdiag}(\mathbf{1}_{m_t}) \in \mathbb{R}^{M_{t_n} \times |\mathcal{S}_{t_n}|}$  with  $t \in \mathcal{S}_{t_n}$ , and  $\mathbf{y}_{t_n}, \boldsymbol{\alpha}_{t_n}, \mathbf{b}_{t_n} \in \mathbb{R}^{M_{t_n}}$  are vectors collecting  $y_i^t, \alpha_i^t$ , and  $b_i^t$  involved in the equality constraints, respectively. Here, the matrix  $\mathbf{Q}_{t_n} \in \mathbb{R}^{M_{t_n} \times M_{t_n}}$  is computed by  $\mathbf{Q}_{t_n}(j, j') = \langle \mathbf{z}_i^t, \mathbf{z}_p^q \rangle$ , where  $t, q \in \mathcal{S}_{t_n}, i = 1, \dots, m_t, p = 1, \dots, m_q$ .

The proposed alternating algorithm gives the final solutions after convergence. In this paper, we set the convergence condition for factors  $\mathbf{U}^n$ , such that  $\sum_n \|\mathbf{U}_{k+1}^n - \mathbf{U}_k^n\|_F^2 / \|\mathbf{U}_k^n\|_F^2 < 10^{-5}$ . After optimization, the prediction for any given input  $\mathbf{x}$  of the  $t$ -th task is obtained either with

- the expression 1) using explicit feature map  $\phi(\cdot)$ :

$$f_t(\mathbf{x}) = (\mathbf{L}\mathbf{u}_t)^\top \phi(\mathbf{x}) + b^t \quad (6)$$

- the expression 2) using kernel function  $k(\cdot, \cdot)$ :

$$f_t(\mathbf{x}) = \sum_{p=1}^{m_q} \sum_{q=1}^T \lambda_p^q k(\mathbf{x}, \mathbf{x}_p^q) \langle \mathbf{u}_t, \mathbf{u}_q \rangle + b^t. \quad (7)$$

Note that expression 1) is the primal representation, while expression 2) is not strictly the dual representation, due to the existence of parameters  $\mathbf{u}_t, \mathbf{u}_q$  in the primal. This is because the optimization algorithm alternates to update different factors of the tensor and the resulting Lagrangian dual forms correspond to each subproblem during iterations, not to the original nonconvex problem (3). Nonetheless, the problem can be efficiently resolved by sets of linear systems, and both expressions 1) and 2) consider correlations across tasks and task-specific information.

## 4 Numerical Experiments

We evaluate the performance of the proposed method on both synthetic and real-world data. Root mean square error (RMSE),  $Q^2$ , and the correlation of the predicted  $\hat{\mathbf{y}}$  and the ground-truth  $\mathbf{y}$  are measured, where  $Q^2$  is defined as  $1 - \|\mathbf{y} - \hat{\mathbf{y}}\|_{\mathbb{F}}^2 / \|\mathbf{y}\|_{\mathbb{F}}^2$  and each iterative method is repeated 10 times for an average. Except for RMSE, a higher metric value indicates a better result. There are three hyperparameters to be tuned in our tLSSVM-MTL, i.e.,  $K, C$ , and the kernel function, and the hyperparameters in the compared methods are also tuned, where 5-fold cross-validation is used.

### 1) Simulated data

The simulated dataset is generated as: 1) the coefficient tensor via the CP form  $\mathcal{W} = \llbracket \mathbf{L}, \mathbf{U}^1, \dots, \mathbf{U}^N \rrbracket$ , where each entry is randomly generated from  $\mathcal{N}(0, 1)$ ; 2)  $\mathbf{x}_i^t, b^t$  and noise  $e_i^t$  from distribution  $\mathcal{N}(0, 1)$ ; 3) the response  $\mathbf{y}^t = \mathbf{y}^t + \sigma \mathbf{e}^t$  consisting of  $\mathbf{y}^t = \mathbf{X}^t \sum_{k=1}^K \mathbf{1}_k \cdot u_{t_1,k}^1 \dots u_{t_N,k}^N + b^t \mathbf{1}_{m_t}$  and  $\mathbf{e}^t$  given by the signal-to-noise ratio (SNR). We set  $d = 100, N = 3, T_1 = 3, T_2 = 4, T_3 = 5$  with  $T = 60$  tasks,  $K = 3$ , and 60 training samples and 20 test samples for each task.

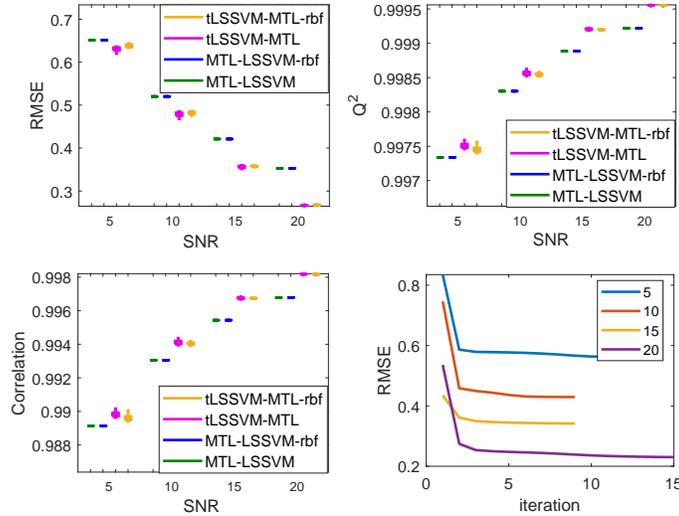


Figure 2: Performance on simulated data with different SNRs.

This experiment mainly aims to validate the efficacy of our tensorized tLSSVM-MTL and optimization results of the proposed algorithm; thus, the MTL-LSSVM counterpart is compared. Fig. 2 presents the performance evaluations on simulated data with different SNR levels, showing that the proposed tLSSVM-MTL consistently provides more accurate predictions on varied SNRs, and its advantage is slightly better with larger SNRs. Additionally, we plot the RMSE during the iterative updates in our method, where RMSE sharply decreases and then converges to a small error. The results of this experiment verify the effectiveness of the proposed method.

### 2) Real-world Data

Three datasets for MTL are employed: Restaurant & Consumer [16], Student performance <sup>2</sup>, and Comprehensive Climate (CCDS). Restaurant & Consumer Dataset contains the rating scores of 138 consumers to different restaurants in 3 aspects, leading to  $138 \times 3$  regression tasks. Student performance Dataset contains student grades in 3 periods and other attributes like sex, and age, where we build  $3 \times 2$  regression tasks by separating the data according to sex and grade period. Comprehensive Climate Dataset (CCDS) gives monthly climate records of 17 variables in North America from 1990 to 2001 [11], where we select 5 locations and construct  $5 \times 17$  regression tasks. MTL-LSSVM [19] and two tensor-based methods, i.e., Convex and Nonconvex Multilinear MTL (MLMTL-C and MLMTL-NC) [13], are compared.

Restaurant & Consumer				
Metric	RMSE	$Q^2$	Correlation	CPU Time
MTL-LSSVM	0.65	41.83%	62.54%	0.45
MTL-LSSVM-rbf	0.65	41.90%	62.55%	0.51
MLMTL-C	0.65	40.42%	61.31%	<b>0.45</b>
MLMTL-NC	0.74	18.61%	56.12%	41.10
tLSSVM-MTL	0.61	45.41%	67.03%	22.86
tLSSVM-MTL-rbf	<b>0.59</b>	<b>49.13%</b>	<b>69.54%</b>	19.36
Student Performance				
Metric	RMSE	$Q^2$	Correlation	CPU Time
MTL-LSSVM	2.99	93.55%	44.66%	<b>0.03</b>
MTL-LSSVM-rbf	2.49	95.56%	67.49%	0.04
MLMTL-C	3.11	93.03%	36.45%	3.21
MLMTL-NC	3.34	91.96%	21.51%	19.10
tLSSVM-MTL	2.99	93.54%	45.79%	0.72
tLSSVM-MTL-rbf	<b>2.44</b>	<b>95.73%</b>	<b>68.59%</b>	0.41
CCDS				
Metric	RMSE	$Q^2$	Correlation	CPU Time
MTL-LSSVM	0.79	29.71%	55.50%	<b>1.08</b>
MTL-LSSVM-rbf	0.70	46.70%	68.36%	1.50
MLMTL-C	0.76	34.56%	58.79%	5.31
MLMTL-NC	0.83	24.04%	50.02%	29.44
tLSSVM-MTL	0.78	32.64%	58.03%	24.07
tLSSVM-MTL-rbf	<b>0.65</b>	<b>54.50%</b>	<b>74.49%</b>	22.01

Table 1: Performance comparison on real-world datasets.

Table 1 presents the prediction results by MTL-LSSVM, MLMTL-C, MLMTL-NC, and the proposed tLSSVM-MTL with both linear and RBF kernels, where the best results are in bold. The results show that our proposed method substantially improves the prediction accuracy in terms of all considered metrics. Our advantages appear more prominent for Restaurant & Consumer and CCDS datasets with RBF kernels, particularly on  $Q^2$  and Correlation metrics which achieve significant improvements. In fact, these two datasets contain larger numbers of tasks, i.e.,  $T = 414$  and  $T = 35$ , and the used multiple indices are endowed with specific meanings in prior to their real-world applications, thereby enabling our model to well learn the underlying structural information.

In Table 1, we also compare the CPU time. We can see that the existing matrix-based MTL-LSSVM and MLMTL-C run faster, due to their convexity benefiting a simple optimization. When comparing with the nonconvex tensor-based MLMTL-NC, our method is more efficient, particularly for the Student Performance dataset, still showing the promising potentials of our tensorized model and the designed iterative updates. Nevertheless, more efficient computations can be expected with further investigations.

## 5 Conclusion

In this paper, we proposed a novel method for MTL regression, which can be regarded as a tensorized generalization and also a multimodal extension of multitask LSSVMs. The proposed method considers multitasks with different indices in the constructed coefficient tensor, which is factorized with low CP rank into a common factor and task-specific factors. In the proposed method, both multilinear and nonlinearity can be flexibly modelled either through feature mappings or kernel functions. In optimization, an alternating strategy is derived to update these factors by solving linear programming subproblems with Lagrangian duality. Experimental results on simulated and real-world data show our great potentials over the compared relevant methods. In future, different tensorization techniques and faster computations are promising to be extended to wider ranges of tasks.

<sup>2</sup><https://archive.ics.uci.edu/ml/datasets/Student+Performance>

## References

- [1] A. Argyriou, T. Evgeniou, and M. Pontil. Convex multi-task feature learning. *Machine Learning*, 73(3):243–272, 2008.
- [2] J D. Carroll and J. Chang. Analysis of individual differences in multidimensional scaling via an n-way generalization of “Eckart-Young” decomposition. *Psychometrika*, 35(3):283–319, 1970.
- [3] R. Caruana. Multitask learning. *Machine Learning*, 28(1):41–75, 1997.
- [4] C. Cortes and V. Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995.
- [5] T. Evgeniou, C. A. Micchelli, and M. Pontil. Learning multiple tasks with kernel methods. *Journal of Machine Learning Research*, 6:615–637, dec 2005.
- [6] T. Evgeniou and M. Pontil. Regularized multi-task learning. In *the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 109–117, 2004.
- [7] L. Liang, F. Cai, and V. Cherkassky. Predictive learning with structured (grouped) data. *Neural Networks*, 22(5-6):766–773, 2009.
- [8] L. Liang and V. Cherkassky. Connection between svm+ and multi-task learning. In *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*, pages 2048–2054. IEEE, 2008.
- [9] J. Liu, C. Zhu, Z. Long, and Y. Liu. Tensor regression. *Foundations and Trends® in Machine Learning*, 14(4):379–565, 2021.
- [10] Y. Liu, J. Liu, Z. Long, and C. Zhu. *Tensor computation for data analysis*. Springer, 2022.
- [11] A. C Lozano, H. Li, A. Niculescu-Mizil, Y. Liu, C. Perlich, J. Hosking, and N. Abe. Spatial-temporal causal modeling for climate change attribution. In *the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 587–596, 2009.
- [12] A. Phan, P. Tichavský, K. Sobolev, K. Sozykin, D. Ermilov, and A. Cichocki. Canonical polyadic tensor decomposition with low-rank factor matrices. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4690–4694. IEEE, 2021.
- [13] B. Romera-Paredes, H. Aung, N. Bianchi-Berthouze, and M. Pontil. Multilinear multitask learning. In *the International Conference on Machine Learning*, pages 1444–1452, 2013.
- [14] H. Shiao and V. Cherkassky. Implementation and comparison of svm-based multi-task learning methods. In *The 2012 International Joint Conference on Neural Networks (IJCNN)*, pages 1–7. IEEE, 2012.
- [15] J. AK Suykens and J. Vandewalle. Least squares support vector machine classifiers. *Neural Processing Letters*, 9(3):293–300, 1999.
- [16] B. Vargas-Govea, G. González-Serna, and R. Ponce-Medellín. Effects of relevant contextual features in the performance of a restaurant recommender system. *ACM RecSys*, 11(592):56, 2011.
- [17] K. Wimalawarne, M. Sugiyama, and R. Tomioka. Multitask learning meets tensor factorization: task imputation via convex optimization. *Advances in Neural Information Processing Systems*, 27, 2014.
- [18] J. Xu, J. Zhou, P. Tan, X. Liu, and L. Luo. Spatio-temporal multi-task learning via tensor decomposition. *IEEE Transactions on Knowledge and Data Engineering*, 33(6):2764–2775, 2019.
- [19] S. Xu, X. An, X. Qiao, and L. Zhu. Multi-task least-squares support vector machines. *Multimedia Tools and Applications*, 71(2):699–715, 2014.
- [20] Z. Xu and K. Kersting. Multi-task learning with task relations. In *the IEEE International Conference on Data Mining*, pages 884–893, 2011.
- [21] Y. Yang and T. Hospedales. Deep multi-task representation learning: A tensor factorisation approach. In *the International Conference on Learning Representations*, 2017.
- [22] Y. Zhang, Y. Zhang, and W. Wang. Deep multi-task learning via generalized tensor trace norm. *arXiv preprint arXiv:2002.04799*, 2020.
- [23] Z. Zhang, Y. Xie, W. Zhang, Y. Tang, and Q. Tian. Tensor multi-task learning for person re-identification. *IEEE Transactions on Image Processing*, 29:2463–2477, 2019.
- [24] Q. Zhao, X. Rui, Z. Han, and D. Meng. Multilinear multitask learning by rank-product regularization. *IEEE Transactions on Neural Networks and Learning Systems*, 31(4):1336–1350, 2020.
- [25] Q. Zheng, Y. Wang, and P. Heng. Multitask feature learning meets robust tensor decomposition for EEG classification. *IEEE Transactions on Cybernetics*, 51(4):2242–2252, 2019.