



Generalized Block-Diagonal Structure Pursuit Learning Soft Latent Task Assignment against Negative Transfer

Zhiyong Yang, Qianqian Xu, Yangbangyan Jiang, Xiaochun Cao, and Qingming Huang
IIE, CAS; ICT, CAS; UCAS; BDKM, CAS; PCL.



Overview

To avoid negative transfer, we propose a novel MTL method with the following contribution:

- ✓ disentangled latent task assignments with a block-diagonal constraint
- ✓ a novel regularizer for generalized block-diagonal structure pursuit
- ✓ explicit connection with optimal transport
- ✓ theoretical guarantee for structural recovery.

Latent Task Representation: A Probabilist View

- ➔ Given T tasks, the per-task parameter W is defined as $W = [W^{(1)}, \dots, W^{(T)}] \in \mathbb{R}^{d \times T}$.
- ➔ To model the relationship among the tasks o_1, \dots, o_T , we assume that $W^{(i)}$ could be represented as a linear combination of latent tasks l_1, \dots, l_k , with $W = LS$.
- ➔ From a probabilist perspective, we regard $S_{i,j}$ as $\mathbb{P}(l = i | o = j)$, namely the possibility of choosing l_i to represent o_j .
- ➔ Another variable of interest is the joint distribution matrix $S_{i,j}^\dagger = \mathbb{P}(l = i, o = j)$. Given marginal distribution \mathbf{a}, \mathbf{b} for l, o , we must have $S_{i,j}^\dagger \mathbf{1}_T = \mathbf{a}$, $S_{i,j}^\dagger \mathbf{1}_k = \mathbf{b}$.
- ➔ In order to avoid suppress, we hope the possibility to assign l_i to o_j is nonzero if and only if (i, j) belongs to the same group. This leads to a block-diagonal structure of S^\dagger up to row and column permutations.

Block-diagonal Structure Pursuit

Auxiliary Bipartite Graph. $A_{l \cup o} = \begin{bmatrix} 0 & S^\dagger \\ S^{\dagger \top} & 0 \end{bmatrix}$, $\Delta(S^\dagger) = \text{diag}(A_{l \cup o} \mathbf{1}) - A_{l \cup o}$. The following shows that squeezing the bottom K eigenvalues of $\Delta(S^\dagger)$ Then leverages the block diagonal property.

Theorem

If $S^\dagger \in \Pi(\mathbf{a}, \mathbf{b})$, $\chi_S = K$ holds if and only if $\dim(\text{Null}(\Delta(S^\dagger))) = K$, i.e., the 0 eigenvalue of $\Delta(S^\dagger)$ has multiplicity K . Moreover, denote $\mathcal{A}^{(i)}$ as the set of latent and output tasks belonging to the i -th block of S , the eigenspace of 0 is spanned by $\mathbf{e}_{\mathcal{A}^{(1)}}, \mathbf{e}_{\mathcal{A}^{(2)}}, \dots, \mathbf{e}_{\mathcal{A}^{(K)}}$, where $\mathbf{e}_{\mathcal{A}^{(i)}} \in \mathbb{R}^{(k+T) \times 1}$, $[\mathbf{e}_{\mathcal{A}^{(i)}}]_j = 1$ if $j \in \mathcal{A}^{(i)}$, otherwise $[\mathbf{e}_{\mathcal{A}^{(i)}}]_j = 0$.

From the variational property of eigenvalues, we reach the regularizer $\Omega(S^\dagger) = \inf \{ \langle \Delta(S^\dagger), U \rangle : U \in \mathcal{M} \}$, where

$$\mathcal{M} = \{U : U \in \mathbb{S}^N, I \succeq U \succeq 0, \text{tr}(U) = K\}$$

Objective Function

Exact problem: Loss + Reg. of L + Structural Reg.

$$\min \tilde{\mathcal{J}} + \Omega_1 + \Omega_2 \quad s.t. \quad S^\dagger \in \Pi(\mathbf{a}, \mathbf{b}), U \in \mathcal{M}, S = TS^\dagger. \quad (1)$$

Inexact problem: Exact + Variable Splitting Between S and TS^\dagger .

$$\min \tilde{\mathcal{J}} + \Omega_1 + \Omega_2 + \Omega_3 \quad s.t. \quad S^\dagger \in \Pi(\mathbf{a}, \mathbf{b}), U \in \mathcal{M}. \quad (2)$$

$$\Omega_1 = \alpha_1 \cdot \|L\|_F^2/2, \Omega_2 = \alpha_3 \cdot \langle \Delta(S^\dagger), U \rangle, \Omega_3 = \alpha_2 \cdot d(S, TS^\dagger)/2$$

Optimization

We present an alternative optimization method to solve.

L, S subproblem. strongly convex, could be solved from off-the-shelf tools.

U subproblem. $U = V_K V_K^\top$, where V_K denotes eigenvectors associated with the smallest k eigenvalues of $\Delta(S^\dagger)$. Define f_i from $V_K = [f_1, \dots, f_{k+T}]^\top$. f_i has a strong grouping power when $\chi_{S^\dagger} = K$.

S^\dagger subroutine: With U updated with $U = V_K V_K^\top$, we reformulate the subproblem as a regularized optimal transport problem:

Proposition (Regularized OT Reformulation)

Reformulation. The S^\dagger subproblem could be reformulated as:

$$\min_{S^\dagger \in \Pi(\mathbf{a}, \mathbf{b})} \frac{\vartheta}{2} \|S^\dagger - \bar{S}\|_F^2 + \langle \mathcal{D}, S^\dagger \rangle \quad (\text{Primal})$$

Connection with OT. Under mild conditions, we have: $0 \leq \mathcal{J}_{\text{REG}} - \mathcal{J}_{\text{OT}} = O(\vartheta/T)$.

Dual Solution The dual problem of (Primal) could be solved from:

$$\argmin_{h, g} \frac{1}{2\vartheta} \cdot \|(h \oplus g - \mathcal{D} + \vartheta \bar{S})_+\|_F^2 - \langle h, \mathbf{a} \rangle - \langle g, \mathbf{b} \rangle,$$

$$\text{with } S^{\dagger*} = \left[\frac{h^* \oplus g^* - \mathcal{D}}{\vartheta} + \bar{S} \right]_+.$$

Negative transfers are punished via a large transportation cost, while positive transfers within group is encouraged with an almost zero cost.

Theoretical Analysis

The hypothesis space \mathcal{H} :

$$\left\{ \left\{ \hat{Y}^{(i)}(X_i^{(t)}) = (LS^{(i)})^\top X_i^{(t)} \right\}_{ti} : \|L\|_F^2 \leq \xi_1, \right.$$

$$\left. d(S, TS^\dagger) \leq \xi_2, \langle \Delta(S^\dagger), U \rangle \leq \xi_3, S^\dagger \in \Pi(\mathbf{a}, \mathbf{b}), U \in \mathcal{M} \right\}$$

$$\xi_1 = 2\mathcal{J}_0/\alpha_1, \xi_2 = 2\mathcal{J}_0/\alpha_2, \xi_3 = 2\mathcal{J}_0/\alpha_3$$

Theorem (Performance Bounds)

Generalization Bound. Picking $\xi_1 = O(1/T)$, $\xi_2 = O(1/T)$, if ℓ is Lipschitz continuous, we have: $|\mathcal{R}(L, S) - \hat{\mathcal{R}}(L, S)| = O_P((nT)^{-1/2})$.

Implicit Spectrum Bound. Picking $\xi_3 = (1/T)^{-3/2}$, $\xi_2 = (1/T^2)$, we have $\sum_{i=N-K+1}^N \lambda_i(\Delta(S)) = O(T^{-1/2})$.

Structure Recovery Bound. Assume that $k \leq T$, for all S^\dagger obtained from the space \mathcal{H} such that $\lambda_{K+1}(\Delta(S^\dagger)) > \lambda_K(\Delta(S^\dagger)) > 0$, there is a co-partition, such that:

$$\|S^{\dagger \text{supp}^c}\|_1 = O\left(k^{1/2}/(T \cdot \lambda_{K+1}(\Delta(S^\dagger)))\right),$$

$$\frac{|\text{supp}^c|}{kT} = O\left(k^{-1/2}/(T \cdot \lambda_{K+1}(\Delta(S^\dagger)))\right).$$

Simulated Dataset

To test the effectiveness of **GBDSP** we generate a simple simulated annotation dataset with $T = 150$ simulated tasks, where the dataset is produced according to the assumption in our model.

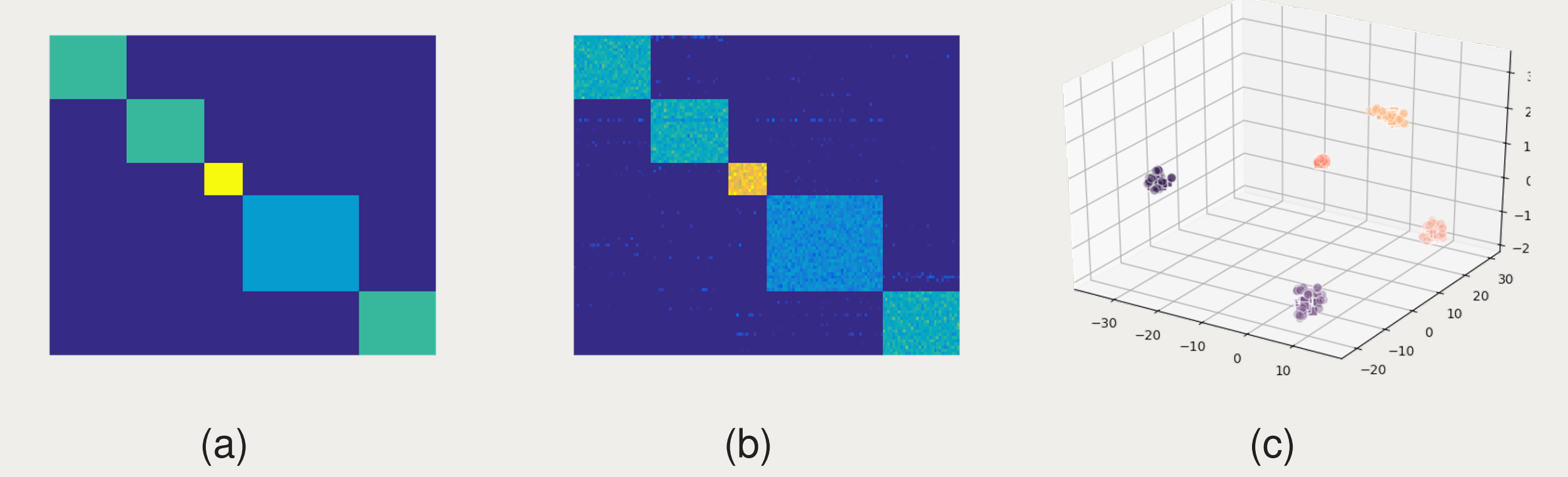


Figure: Visualizations over the Simulated Dataset. (a) shows The true LATM; (b) shows the LATM recovered by **GBDSP** (c) shows the spectral embedding **GBDSP**.

Real World Datasets

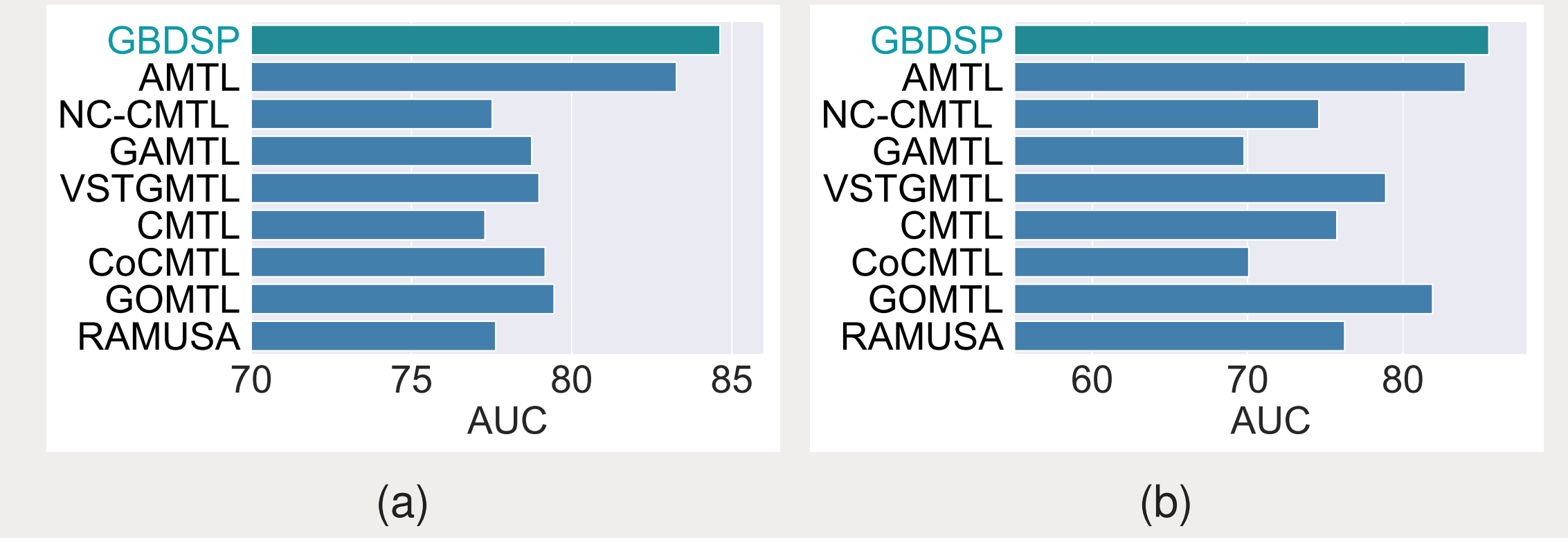


Figure: (a) Performance Comparison Curve over the Simulation Dataset, with varying training data ratio. (b) Performance Comparison over the Sun Dataset.

Table: Performance Comparison over General MTL Datasets (mean \pm std)

Algorithms	AWA2-Attr(\uparrow)	AWA2-Cls(\uparrow)	School(\downarrow)
RAMUSA	88.60 \pm 0.66	93.06 \pm 0.53	10.52 \pm 0.09
GOMTL	89.56 \pm 0.33	88.22 \pm 1.18	10.26\pm0.11
CoCMTL	92.29 \pm 0.35	94.69 \pm 0.73	12.06 \pm 0.09
CMTL	92.95\pm0.35	94.81 \pm 0.70	12.06 \pm 0.09
VSTGMTL	89.31 \pm 0.39	92.03 \pm 0.94	10.17 \pm 0.08
GAMTL	89.39 \pm 0.42	92.55 \pm 0.53	10.50 \pm 0.12
NC-CMTL	92.99\pm0.32	95.10 \pm 0.60	10.53 \pm 0.12
AMTL	92.15 \pm 0.34	95.76\pm0.44	12.15 \pm 0.09
GBDSP	92.73 \pm 0.29	97.86\pm0.22	10.10\pm0.08

Contact Information

Our Group:

<https://qmhuang-ucas.github.io/>

Email me : yangzhiyong@iie.ac.cn

Visit my home page from QR code

