

# Supplementary Material: Structured Sparse Multi-Task Learning with Generalized Group Lasso

## Supplement of Table 1 in the Main Paper

The following table summarizes more structured sparse MTL methods with different settings of  $\mathbf{d}_{g_i}$ , which is a supplement of Table 1 in the main paper.

Table A1: A summary of detailed settings of GenGL for selected MTL methods. For decomposition methods, suppose there are  $h$  components and the  $l$ th component is associated with an operator  $\pi_{l,i}$ . For clarity, when  $h \in \mathbb{N}_+$ , we show  $\Omega(\mathbf{w}_l)$  instead.

| Architecture | Method   | $\Omega(\mathbf{w})$  | $\pi$  | $h$                  |
|--------------|--|---|--|----------------------|
| Shallow      | Lasso<br>[Tibshirani, 1996]  | $\gamma \ \mathbf{w}\ _1$   | $\pi_1 \in \mathcal{P}_F, \pi_2 \in \mathcal{P}_F$   | $h = 1$              |
|              | Group Lasso (GL)<br>[Yuan and Lin, 2006]                                   | $\gamma \ \mathbf{W}\ _{2,1}$   | $\pi_1 \in \mathcal{P}_F, \pi_2 \in \mathcal{P}_C$   | $h = 1$              |
|              | rMTFL<br>[Gong et al., 2012]   | $\gamma_1 \ \mathbf{W}_1\ _{2,1} + \gamma_2 \ \mathbf{W}_2^T\ _{2,1}$<br>$s.t. \quad \mathbf{W} = \mathbf{W}_1 + \mathbf{W}_2$  | $\pi_{1,1} \in \mathcal{P}_F, \pi_{1,2} \in \mathcal{P}_C$<br>$\pi_{2,1} \in \mathcal{P}_C, \pi_{2,2} \in \mathcal{P}_F$   | $h = 2$              |
|              | MeTaG<br>[Han and Zhang, 2015]<br>CCMTL<br>[He et al., 2019]               | $\frac{\gamma}{\phi^t} \sum_{j < k} \ \mathbf{w}_{l,:j} - \mathbf{w}_{l,:k}\ _2$  | $\pi_{l,1} \in \mathcal{P}_C, \pi_{l,2} \in \mathcal{P}_{F^2}$   | $h \in \mathbb{N}_+$ |
|              | Sparse Network Lasso<br>[Okazaki and Kawano, 2022]<br>[Wang and Sun, 2022] | $\sum_{j < k} r_{j,k} \ \mathbf{w}_{l,j} - \mathbf{w}_{l,k}\ _2 + \ \mathbf{w}\ _1$   | $\pi_1^1 \in \mathcal{P}_C, \pi_2^1 \in \mathcal{P}_{F^2}$<br>$\pi_1^2 \in \mathcal{P}_F, \pi_2^2 \in \mathcal{P}_F$   | $h = 1$              |
|              | SSMTL <sub>m</sub><br>(The proposed model)                                 | $\frac{\gamma_1}{\phi^t} \sum_{j < k} \ \mathbf{w}_{l,:j} - \mathbf{w}_{l,:k}\ _2 + \frac{\gamma_2}{\phi^{2t}} \ \mathbf{W}_l\ _{2,1}$<br>$+ \frac{\gamma_1}{\phi^{2t}} \sum_i \sum_{j < k}  w_{l,ij} - w_{l,ik} ^2 + \frac{\gamma_2}{\phi^t} \ \mathbf{w}_l\ _1$   | $\pi_{l,1}^{(1)} \in \mathcal{P}_C, \pi_{l,2}^{(1)} \in \mathcal{P}_{F^2}$<br>$\pi_{l,1}^{(2)} \in \mathcal{P}_F, \pi_{l,2}^{(2)} \in \mathcal{P}_C$<br>$\pi_{l,1}^{(3)} \in \mathcal{P}_F, \pi_{l,2}^{(3)} \in \mathcal{P}_{F^2}$<br>$\pi_{l,1}^{(4)} \in \mathcal{P}_F, \pi_{l,2}^{(4)} \in \mathcal{P}_F$   | $h \in \mathbb{N}_+$ |
| Deep         | Sparse GL (SGL)<br>[Scardapane et al., 2017]                               | $\frac{\gamma_1}{\phi^t} \sum_k \sum_j^{p_2} \ \mathbf{w}_{l,:jk}\ _2 + \frac{\gamma_2}{\phi^t} \sum_{i,j,k}  w_{l,ijk} $   | $\pi_{l,1}^{(1)} \in \mathcal{P}_C, \pi_{l,2}^{(1)} \in \mathcal{P}_F, \pi_{l,3}^{(1)} \in \mathcal{P}_F$<br>$\pi_{l,1}^{(2)} \in \mathcal{P}_F, \pi_{l,2}^{(2)} \in \mathcal{P}_F, \pi_{l,3}^{(2)} \in \mathcal{P}_F$   | $h \in \mathbb{N}_+$ |
|              | Adaptive GL (AGL)<br>[Dinh and Ho, 2020a]                                  | $\sum_k \sum_j^{p_2} \frac{\gamma_j}{\phi^t} \ \mathbf{w}_{l,:jk}\ _2$  | $\pi_{l,1} \in \mathcal{P}_C, \pi_{l,2} \in \mathcal{P}_F, \pi_{l,3} \in \mathcal{P}_F$  | $h \in \mathbb{N}_+$ |
|              | GL + Adaptive GL (GLAGL)<br>[Dinh and Ho, 2020b]                           | $\frac{\gamma}{\phi^t} \sum_k \sum_j^{p_2} \ \mathbf{w}_{l,:jk}\ _2 + \sum_k \sum_j^{p_2} \frac{\gamma_j}{\phi^t} \ \mathbf{w}_{l,:jk}\ _2$   | $\pi_{l,1}^{(1)} \in \mathcal{P}_C, \pi_{l,2}^{(1)} \in \mathcal{P}_F, \pi_{l,3}^{(1)} \in \mathcal{P}_F$<br>$\pi_{l,1}^{(2)} \in \mathcal{P}_F, \pi_{l,2}^{(2)} \in \mathcal{P}_C, \pi_{l,3}^{(2)} \in \mathcal{P}_F$   | $h \in \mathbb{N}_+$ |
|              | SSMTL <sub>t</sub><br>(The proposed model)                                 | $\frac{\gamma_1}{\phi^t} \sum_{j < k} \ \mathbf{W}_{l,:j} - \mathbf{W}_{l,:k}\ _F + \frac{\gamma_2}{\phi^{2t}} \sum_k \sum_j^{p_2} \ \mathbf{w}_{l,:jk}\ _2$<br>$+ \frac{\gamma_1}{\phi^{2t}} \sum_i \sum_{j < k} \ \mathbf{w}_{l,:ij} - \mathbf{w}_{l,:ik}\ _2 + \frac{\gamma_2}{\phi^t} \sum_{i,j,k}  w_{l,ijk} $ | $\pi_{l,1}^{(1)} \in \mathcal{P}_C, \pi_{l,2}^{(1)} \in \mathcal{P}_C, \pi_{l,3}^{(1)} \in \mathcal{P}_{F^2}$<br>$\pi_{l,1}^{(2)} \in \mathcal{P}_C, \pi_{l,2}^{(2)} \in \mathcal{P}_F, \pi_{l,3}^{(2)} \in \mathcal{P}_F$<br>$\pi_{l,1}^{(3)} \in \mathcal{P}_C, \pi_{l,2}^{(3)} \in \mathcal{P}_F, \pi_{l,3}^{(3)} \in \mathcal{P}_{F^2}$<br>$\pi_{l,1}^{(4)} \in \mathcal{P}_F, \pi_{l,2}^{(4)} \in \mathcal{P}_F, \pi_{l,3}^{(4)} \in \mathcal{P}_F$ | $h \in \mathbb{N}_+$ |

## Summary of different types of linear operators in Sec. 4.1

Table A2: Summary of different types of operators for the  $i$ th dimension. (F: feature-level, T: task-level, E: element-wise, V: vector-wise, Net: net-wise, Neu: neuron-wise, W: weight-wise)

|         |   | $\pi_i \in \mathcal{P}_F$ |                | $\pi_i \in \mathcal{P}_{F^2}$ |                 |
|---------|---|---------------------------|----------------|-------------------------------|-----------------|
|         |   | $1 \leq i \leq N-1$       | $i = N$        | $1 \leq i \leq N-1$           | $i = N$         |
| $N = 2$ | $\forall j \neq i, \pi_j \notin \mathcal{P}_C$        | FE-selection              | TE-selection   | FE-clustering                 | TE-clustering   |
|         | $\exists! j \neq i, \pi_j \in \mathcal{P}_C$          | FV-selection              | TV-selection   | FV-clustering                 | TV-clustering   |
| $N = 3$ | $\forall j \neq i, \pi_j \notin \mathcal{P}_C$        | FW-selection              | TW-selection   | FW-clustering                 | TW-clustering   |
|         | $\exists! j \neq i, \pi_j \in \mathcal{P}_C$          | FNeu-selection            | TNeu-selection | FNeu-clustering               | TNeu-clustering |
|         | $\exists j, k \neq i, \pi_j, \pi_k \in \mathcal{P}_C$ | FNet-selection            | TNet-selection | FNet-clustering               | TNet-clustering |

## Derivations for the Gradient of the Function in (14)

To compute the gradient  $\nabla_{\mathbf{w}} f_{\mu}(\mathbf{w})$ , we introduce the following lemma and provide its proof.

**Lemma 1.** For any  $\mu > 0$ ,  $f_{\mu}(\mathbf{w})$  is convex and differentiable in  $\mathbf{w}$ , and the gradient of  $f_{\mu}(\mathbf{w})$  w.r.t.  $\mathbf{w}$  is

$$\nabla_{\mathbf{w}} f_{\mu}(\mathbf{w}) = \gamma \mathbf{D}^T \beta^*, \quad (1)$$

where  $\beta^*$  is the optimal solution to (11) in the main paper. The optimal  $\beta^*$  is got by a projection operator  $S(\cdot)$ , which projects any vector  $\mathbf{u}$  to the  $l_2$  ball:

$$\beta_{g_1, \dots, g_N}^* = S\left(\frac{\gamma \mathbf{D}_{g_1, \dots, g_N} \mathbf{w}}{\mu}\right). \quad (2)$$

*Proof.* We first introduce that the *Fenchel conjugate*  $\phi^*(\alpha)$  of a function  $\phi(\beta)$  is defined as:

$$\phi^*(\alpha) = \sup_{\beta \in \text{dom}(\phi)} (\beta \alpha^T - \phi(\beta)). \quad (3)$$

Recall that  $q(\beta) = \frac{1}{2} \|\beta\|_2^2$  with  $\text{dom}(q) = \mathcal{B}$  in the main paper, the conjugate of  $q(\cdot)$  at  $\frac{\mathbf{D}\mathbf{w}}{\mu}$  is derived as  $q^*\left(\frac{\mathbf{D}\mathbf{w}}{\mu}\right) = \sup_{\beta \in \mathcal{B}} (\beta^T \frac{\mathbf{D}\mathbf{w}}{\mu} - q(\beta))$ . Hence,  $f_\mu(\mathbf{w})$  is reformulated as:

$$f_\mu(\mathbf{w}) = \max_{\beta \in \mathcal{B}} (\gamma \beta^T \mathbf{D}\mathbf{w} - \mu q(\beta)) = \mu q^*\left(\frac{\mathbf{D}\mathbf{w}}{\mu}\right). \quad (4)$$

Since  $q(\beta)$  is strictly convex, its conjugate is smooth. Therefore,  $f_\mu(\beta)$  is a smooth function. Let  $\phi(\beta, \mathbf{w}) = \beta^T \mathbf{D}\mathbf{w} - \mu q(\beta)$ . Since  $q(\cdot)$  is strongly convex,  $\arg\max_{\beta \in \mathcal{B}} \phi(\beta, \mathbf{w})$  has a unique optimal solution denoted as  $\beta^*$ . According to Danskin's theorem [Mangasarian, 1994],

$$\nabla_{\mathbf{w}} f_\mu(\mathbf{w}) = \nabla_{\mathbf{w}} \phi(\beta^*, \mathbf{w}) = \gamma \mathbf{D}^T \beta^*. \quad (5)$$

Then we calculate the optimal  $\beta^*$ .

$$\begin{aligned} \beta^* &= \arg\max_{\beta \in \mathcal{B}} (\gamma \beta^T \mathbf{D}\mathbf{w} - \mu q(\beta)) \\ &= \arg\max_{\beta \in \mathcal{B}} \sum_{g_1, \dots, g_N} (\gamma \beta_{g_1, \dots, g_N}^T \mathbf{D}_{g_1, \dots, g_N} \mathbf{w} - \frac{\mu}{2} \|\beta_{g_1, \dots, g_N}\|_2^2) \\ &= \arg\min_{\beta \in \mathcal{B}} \sum_{g_1, \dots, g_N} \left\| \beta_{g_1, \dots, g_N} - \frac{\gamma \mathbf{D}_{g_1, \dots, g_N} \mathbf{w}}{\mu} \right\|_2^2. \end{aligned}$$

We can see that  $\beta^*$  is a column concatenation of  $\beta_{g_1, \dots, g_N}^*$ , and each  $\beta_{g_1, \dots, g_N}^*$  can be calculated by:

$$\beta_{g_1, \dots, g_N}^* = \arg\min_{\beta_{g_1, \dots, g_N} : \|\beta_{g_1, \dots, g_N}\|_2 \leq 1} \left\| \beta_{g_1, \dots, g_N} - \frac{\gamma \mathbf{D}_{g_1, \dots, g_N} \mathbf{w}}{\mu} \right\|_2^2. \quad (6)$$

According to the property of the  $l_2$  ball, it can be shown that:

$$\beta_{g_1, \dots, g_N}^* = S\left(\frac{\gamma \mathbf{D}_{g_1, \dots, g_N} \mathbf{w}}{\mu}\right), \quad (7)$$

where

$$S(\mathbf{u}) = \begin{cases} \frac{\mathbf{u}}{\|\mathbf{u}\|_2} & \|\mathbf{u}\|_2 > 1, \\ \mathbf{u} & \|\mathbf{u}\|_2 \leq 1. \end{cases}$$

□

## Optimization Algorithm

In Algorithm 1, we provide the pseudocode of the optimization algorithm discussed in Sec. 5 of the main paper.

---

**Algorithm 1** Optimization algorithm for solving (13) in the main paper

---

**Input:**  $\mathbf{X}, \mathbf{y}, \mathbf{D}, \gamma, \hat{\mathbf{w}}^0$ , desired accuracy  $\varepsilon$ , learning rate  $\eta$

**Output:**  $\mathbf{w}$

- 1: Initialize  $\nabla F(\hat{\mathbf{w}}^0)$ ,  $\alpha = 0$ ,  $\theta_0 = 1$ , smoothness parameter  $\mu = \frac{\varepsilon}{\prod_i |\mathcal{G}_i|}$ .
- 2: **repeat**
- 3:   Compute  $\nabla F(\hat{\mathbf{w}}^k)$  using Eq. (16) in the main paper.
- 4:   Solve the proximal step:

$$\mathbf{w}^{k+1} = \arg\min_{\mathbf{w}} F(\hat{\mathbf{w}}^k) + \langle \mathbf{w} - \hat{\mathbf{w}}^k, \nabla F(\hat{\mathbf{w}}^k) \rangle + \frac{2}{\eta} \|\mathbf{w} - \hat{\mathbf{w}}^k\|_2^2. \quad (8)$$

- 5:   Set  $\theta_{k+1} = \frac{2}{\alpha+3}$ .
  - 6:   Set  $\hat{\mathbf{w}}^{k+1} = \mathbf{w}^{k+1} + \frac{1-\theta_k}{\theta_k} \theta_{k+1} (\mathbf{w}^{k+1} - \hat{\mathbf{w}}^k)$ .
  - 7:   Set  $\alpha = \alpha + 1$ .
  - 8: **until** convergence
-

## Real-World Datasets

To evaluate the proposed SSMTL<sub>m</sub>, we conduct experiments on the following six regression datasets:

- **RF1:** The RF1 dataset is to predict the river network flows for 48 h in the future at 8 sites, and each site contributes 8 attribute variables to facilitate prediction, thus there are a total of 64 variables plus 8 target variables. We randomly select 1000 samples in the RF1 dataset.
- **Isolet:** The Isolet dataset is generated as follows. 150 subjects spoke the name of each letter of the alphabet twice, and they are divided into 5 groups, leading to 5 tasks. We randomly select 1000 samples in the Isolet dataset.
- **Energy:** The Energy dataset performs energy analysis based on different building shapes simulated in Ecotect, aiming to predict 2 real valued responses by using 8 features. We randomly select 600 samples in the Energy dataset.
- **Parkinsons:** The Parkinsons dataset is to predict the disease symptom scores of 42 patients by using 16 bio-medical features, which results in 42 tasks.
- **School:** The School dataset contains the examination information of 15,362 students from 139 schools, where each school is regarded as a task. We are aiming to predict the score of each student.
- **SARCOS:** The SARCOS dataset is based on a inverse dynamic problem, which maps from a input space with 21 dimensions to 7 torques. We randomly select 1000 samples in the SARCOS dataset.

To evaluate the proposed SSMTL<sub>t</sub>, we conduct experiments on the following three classification datasets:

- **SSD:** The SSD dataset contains features extracted from electric current drive signals. The drive has intact and defective components, resulting in 11 different classes with different conditions. We randomly select 2000 samples in the SSD dataset.
- **MNIST:** The MNIST dataset is a large database of handwritten digits with images in 28x28 pixel boxes. We reduce the feature dimension to 154 by Principle Component Analysis (PCA) [Wold et al., 1987], and randomly select 1000 samples in the MNIST dataset.
- **COVER:** The COVER dataset predicts forest cover type in the Roosevelt National Forest of northern Colorado. There are 7 kinds of cover types, leading to 7 tasks. We randomly select 2000 samples in the COVER dataset.

## Detailed Settings of the Networks

For the three datasets used for deep models, we construct three soft-parameter sharing networks respectively. Each network is associated with a number of sub-nets corresponding to the tasks. Details are shown in Table A3.

Table A3: Detailed settings of the networks used in deep models.

| Datasets      | SSD        | MNIST      | COVER      |
|---------------|------------|------------|------------|
| Optimizer     | Adam       |            |            |
| Learning rate | $10^{-3}$  |            |            |
| # of neurons  | 40/40/20/1 | 80/40/20/1 | 50/50/25/1 |
| Batch size    | 128        | 64         | 128        |
| # of sub-nets | 11         | 10         | 7          |

## Evaluation Metrics

For those regression tasks in experiments for MTL in shallow architectures, we adopt the following three metrics: normalized Mean Squared Error (nMSE) with the form of  $\frac{1}{mn_t} \sum_{t=1}^m \frac{\|\mathbf{y}_t - \mathbf{X}_t \mathbf{w}_{:t}\|_2^2}{\|\mathbf{y}_t\|_2^2}$ , Mean Absolute Error (MAE) with the form of  $\frac{1}{mn_t} \sum_{t=1}^m \|\mathbf{y}_t - \mathbf{X}_t \mathbf{w}_{:t}\|_1$  and Explained Variance (EV) [Bakker and Heskes, 2003] with the form of  $1 - \frac{1}{mn_t} \sum_{t=1}^m \frac{\text{var}(\mathbf{y}_t - \mathbf{X}_t \mathbf{w}_{:t})}{\text{var}(\mathbf{y}_t)}$  to evaluate the performance. As for the classification tasks in deep MTL experiments, the two metrics we use are: Accuracy with the form of  $\frac{1}{m} \sum_{t=1}^m \frac{\text{TP}_t + \text{TN}_t}{n_t}$ , where  $\text{TP}_t$  and  $\text{TN}_t$  are the number of True Positive and True Negative samples of the  $t$ -th task, and Area Under ROC-Curve (AUC) [Hand and Till, 2001] with the form of  $\frac{1}{m} \text{AUC}_t$ , where  $\text{AUC}_t$  is the AUC score of the  $t$ -th task.

## Comparison Methods

In experiments for MTL in shallow architectures, we compare  $\text{SSMTL}_m$  with the following six methods:

- **Lasso** [Tibshirani, 1996]: A single task learning method by penalizing each task with  $l_1$ -norm regularization independently, which is selected as the baseline method.
- **GL** [Yuan and Lin, 2006]: MTL with Group Lasso, which penalizes sum of the  $l_2$ -norms of task-common feature groups to learn global features that shared by multiple tasks.
- **rMTFL** [Gong et al., 2012]: Robust multi-task feature Learning method, which simultaneously captures features of related tasks and identifies outlier tasks.
- **MeTaG** [Han and Zhang, 2015]: Multi-Level task grouping method, which decomposes the weight into multiple levels and imposes the  $l_2$ -norm regularization on the pairwise difference among the tasks.
- **GBDSP** [Yang et al., 2019]: By assuming that tasks share a latent group structure, it learns a generalized block-diagonal structure for the latent basis of parameters.
- **KMSV** [Chang et al., 2021]: Based on the low-rank assumption, it minimizes exactly  $k$  minimal singular values to learn the latent group structure.

In experiments for deep MTL, we adopt the same soft-sharing network with totally independent sub-nets for multiple tasks, and compare  $\text{SSMTL}_t$  with the following five deep models:

- **DMTRL** [Yang and Hospedales, 2017]: Deep multi-task representation learning, which aims to obtain a low-rank structure by tensor factorization.
- **SGL** [Scardapane et al., 2017]: It imposes the  $l_1$ -norm and GL together to guarantee both weight-wise sparsity and neuron-wise sparsity of layer parameters.
- **GL+AGL** [Dinh and Ho, 2020b]: It combines GL and Adaptive GL [Dinh and Ho, 2020a], and builds neuron-wise sparse model by consistent feature selection.
- **STG** [Yamada et al., 2020]: It is a robust and deep feature selection method using stochastic gates, which is based on a continuous relaxation of the  $l_0$ -norm.

## Hyperparameter Sensitivity Analysis of $\text{SSMTL}_m$

The sensitivity on  $\gamma_1$ ,  $\gamma_2$  and  $\phi$  of  $\text{SSMTL}_m$  is invested on the RF1 dataset for shallow MTL. In  $\text{SSMTL}_m$ ,  $\gamma_1$  controls the regularization degree of *task-level* operations while  $\gamma_2$  controls the regularization degree of *feature-level* operations, which are selected from  $\{10^{-3}, 10^{-2}, \dots, 10^2, 10^3\}$ , and  $\phi$  alters the impact among components (layers), which is selected from  $\{2, 5, 10, 20, 50, 100\}$ . Fig. A1 shows the result in nMSE of three experiments. Specifically, Fig. A1(a), A1(b) and A1(c) are shown by fixing  $\phi = 10$ ,  $\gamma_2 = 1$  and  $\gamma_1 = 1$  respectively. The result shows that: 1) it is recommended to set  $\gamma_1 < 10^{-1}$  and  $\phi < 10$  on the RF1 dataset; 2)  $\gamma_2$  is not as sensitive as the other parameters.

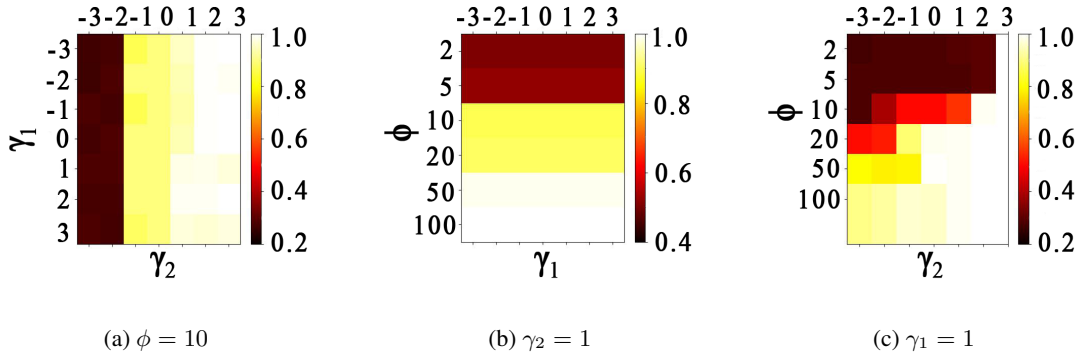


Figure A1: Hyperparameter Sensitivity Analysis on  $\gamma_1$ ,  $\gamma_2$  and  $\phi$  of  $\text{SSMTL}_m$  on the RF1 dataset in shallow MTL experiments. The values of  $\gamma_1$  and  $\gamma_2$  are shown in the logarithmic scale while the value of  $\phi$  is selected from  $\{2, 5, 10, 20, 50, 100\}$ .

## Results in MAE and EV of the Ablation Study for $\text{SSMTL}_m$

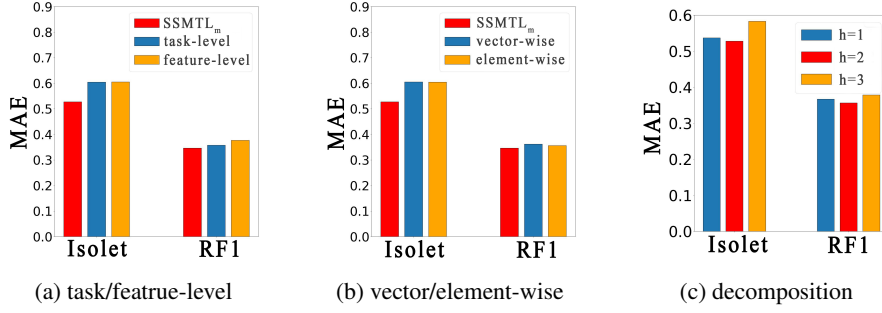


Figure A2: Results in MAE of different operations and model decomposition of  $\text{SSMTL}_m$  on the Isolet and RF1 datasets.

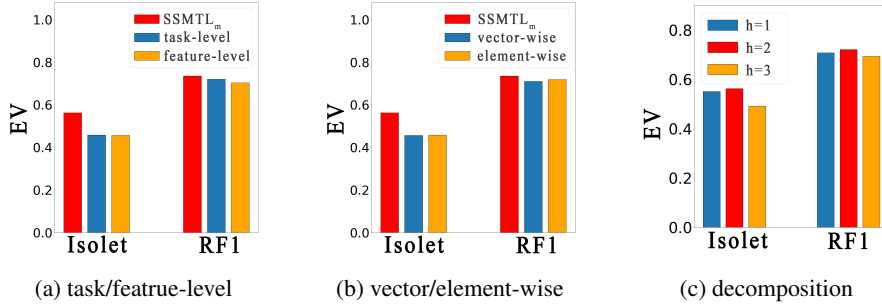


Figure A3: Results in EV of different operations and model decomposition of  $\text{SSMTL}_m$  on the Isolet and RF1 datasets.

## Significant test results

We do paired Wilcoxon tests for statistical test. The results in Table A4 show that  $\text{SSMTL}_m$  significantly outperforms MeTaG and GBDSP on the datasets.

Table A4: Results of paired samples Wilcoxon test at 5% significance level. We show p-values and  $\bullet/\circ$  indicates that  $\text{SSMTL}_m$  is superior/inferior to the comparing method.

| Method                    | Datasets         |                  |                  |                  |
|---------------------------|------------------|------------------|------------------|------------------|
|                           | Parkinsons       | RF1              | Isolet           | SARCOS           |
| MeTaG vs $\text{SSMTL}_m$ | 0.0217 $\bullet$ | 0.0724 $\circ$   | 0.0439 $\bullet$ | 0.0362 $\bullet$ |
| GBDSP vs $\text{SSMTL}_m$ | 0.0278 $\bullet$ | 0.0320 $\bullet$ | 0.0199 $\bullet$ | 0.0225 $\bullet$ |

## References

- [Bakker and Heskes, 2003] Bakker, B. and Heskes, T. M. (2003). Task clustering and gating for bayesian multitask learning. *J. Mach. Learn. Res.*, 4:83–99.
- [Chang et al., 2021] Chang, W., Nie, F., Wang, R., and Li, X. (2021). New tight relaxations of rank minimization for multi-task learning. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pages 2910–2914.
- [Dinh and Ho, 2020a] Dinh, V. and Ho, L. S. T. (2020a). Consistent feature selection for neural networks via adaptive group lasso. *arXiv preprint arXiv:2006.00334*.
- [Dinh and Ho, 2020b] Dinh, V. C. and Ho, L. S. (2020b). Consistent feature selection for analytic deep neural networks. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems*, volume 33, pages 2420–2431.
- [Gong et al., 2012] Gong, P., Ye, J., and Zhang, C. (2012). Robust multi-task feature learning. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '12*, page 895–903, New York, NY, USA.

- [Han and Zhang, 2015] Han, L. and Zhang, Y. (2015). Learning multi-level task groups in multi-task learning. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, AAAI’15, page 2638–2644. AAAI Press.
- [Hand and Till, 2001] Hand, D. J. and Till, R. J. (2001). A simple generalisation of the area under the roc curve for multiple class classification problems. *Mach. Learn.*, 45(2):171–186.
- [He et al., 2019] He, X., Alesiani, F., and Shaker, A. (2019). Efficient and scalable multi-task regression on massive number of tasks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33:3763–3770.
- [Mangasarian, 1994] Mangasarian, O. L. (1994). *Nonlinear programming*. SIAM.
- [Okazaki and Kawano, 2022] Okazaki, A. and Kawano, S. (2022). Multi-task learning for compositional data via sparse network lasso. *Entropy*, 24(12).
- [Scardapane et al., 2017] Scardapane, S., Comminiello, D., Hussain, A., and Uncini, A. (2017). Group sparse regularization for deep neural networks. *Neurocomput.*, 241(C):81–89.
- [Tibshirani, 1996] Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the royal statistical society series b-methodological*, 58:267–288.
- [Wang and Sun, 2022] Wang, J. and Sun, L. (2022). Multi-task personalized learning with sparse network lasso. In Raedt, L. D., editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 3516–3522.
- [Wold et al., 1987] Wold, S., Esbensen, K., and Geladi, P. (1987). Principal component analysis. *Chemometrics and Intelligent Laboratory Systems*, 2(1-3):37–52.
- [Yamada et al., 2020] Yamada, Y., Lindenbaum, O., Negahban, S., and Kluger, Y. (2020). Feature selection using stochastic gates. In *International Conference on Machine Learning*, pages 10648–10659. PMLR.
- [Yang and Hospedales, 2017] Yang, Y. and Hospedales, T. M. (2017). Deep multi-task representation learning: A tensor factorisation approach. In *International Conference on Learning Representations*.
- [Yang et al., 2019] Yang, Z., Xu, Q., Jiang, Y., Cao, X., and Huang, Q. (2019). Generalized block-diagonal structure pursuit: Learning soft latent task assignment against negative transfer. *Advances in Neural Information Processing Systems*, 32.
- [Yuan and Lin, 2006] Yuan, M. and Lin, Y. (2006). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society Series B*, 68:49–67.