Supplementary Material of the Paper: Convex Discriminative Multitask Clustering

Xiao-Lei Zhang, Member, IEEE

Abstract

In this supplementary material, we first give the detailed derivation of equations (6) and (11) in the main paper in Section 1 and Section 2 respectively. Then, we present the optimization procedure of DMTC in detail in Section 3. In the rest of the supplement, we will report the experimental results in detail in four evaluation metrics – Accuracy (ACC), Area Under the ROC Curves (AUC), Adjusted Rand Index (ARI), and Normalized Mutual Information (NMI).

1 DERIVATION OF EQUATION (6) IN THE MAIN PAPER

We can get the Lagrange of the equation (5) in the main paper as follows:

$$\mathcal{L}(\mathbf{W}_{c};\xi_{j,c}^{i};\alpha_{j,c}^{i}) = \frac{\lambda_{1}}{2}\operatorname{tr}(\mathbf{W}_{c}^{T}\mathbf{W}_{c}) + \frac{\lambda_{2}}{2}\operatorname{tr}(\mathbf{W}_{c}^{T}\mathbf{D}^{-1}\mathbf{W}_{c}) + \sum_{i=1}^{m}\frac{1}{n_{i}}\sum_{j=1}^{n_{i}}(\xi_{j,c}^{i})^{2} + \sum_{i=1}^{m}\sum_{j=1}^{n_{i}}\alpha_{j,c}^{i}\left(\bar{y}_{j,c}^{i} - \mathbf{w}_{i,c}^{T}\mathbf{x}_{j}^{i} - \xi_{j,c}^{i}\right) \\ = \frac{1}{2}\operatorname{tr}(\mathbf{W}_{c}^{T}\mathbf{D}^{-1}(\lambda_{1}\mathbf{D} + \lambda_{2}\mathbf{I}_{d})\mathbf{W}_{c}) - \sum_{i=1}^{m}\sum_{j=1}^{n_{i}}\alpha_{j,c}^{i}\mathbf{w}_{i,c}^{T}\mathbf{x}_{j}^{i} + \sum_{i=1}^{m}\sum_{j=1}^{n_{i}}\alpha_{j,c}^{i}\bar{y}_{j,c}^{i} + \sum_{i=1}^{m}\sum_{j=1}^{n_{i}}\left(\frac{1}{n_{i}}(\xi_{j,c}^{i})^{2} - \alpha_{j,c}^{i}\xi_{j,c}^{i}\right) \tag{1}$$

where $\alpha_{j,c}^i$ is the Lagrangian variable, and \mathbf{I}_d is a *d*-dimensional identity matrix. Calculating the partial derivatives of problem (1) over \mathbf{W}_c and $\xi_{j,c}^i$ and letting the derivatives equal to 0 can get the following two equations:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_c} = \mathbf{D}^{-1} (\lambda_1 \mathbf{D} + \lambda_2 \mathbf{I}_d) \mathbf{W}_{i,c} - \sum_{i=1}^m \sum_{j=1}^{n_i} \alpha_{j,c}^i \mathbf{x}_j^i \mathbf{e}_i^T = 0$$
(2)

$$\frac{\partial \mathcal{L}}{\partial \xi_{j,c}^{i}} = \frac{2}{n_{i}} \xi_{j,c}^{i} - \alpha_{j,c}^{i} = 0$$
(3)

From the above equations, we can get:

$$\mathbf{W}_{c} = \sum_{i=1}^{m} \sum_{j=1}^{n_{i}} \alpha_{j,c}^{i} \mathbf{D} (\lambda_{1} \mathbf{D} + \lambda_{2} \mathbf{I}_{d})^{-1} \mathbf{x}_{j}^{i} \mathbf{e}_{i}^{T}$$

$$\tag{4}$$

$$\xi_{j,c}^{i} = \frac{n_{i}}{2} \alpha_{j,c}^{i} \tag{5}$$

Substituting (4) and (5) to (1) can derive the following maximization problem

$$\max_{\boldsymbol{\alpha}_{c}} \sum_{i=1}^{m} \sum_{j=1}^{n_{i}} \alpha_{j,c}^{i} \bar{y}_{j,c}^{i} - \frac{1}{2} \sum_{i_{1}=1}^{m} \sum_{i_{2}=1}^{m} \sum_{j_{1}=1}^{n_{i}} \sum_{j_{2}=1}^{n_{i}} \alpha_{j_{1},c}^{i_{1}} \mathbf{x}_{j_{1}}^{i_{1}}^{T} \mathbf{D} \left(\lambda_{1} \mathbf{D} + \lambda_{2} \mathbf{I}_{d}\right)^{-1} \mathbf{x}_{j_{2}}^{i_{2}} \alpha_{j_{2},c}^{i_{2}} \left\langle \mathbf{e}_{i_{1}}, \mathbf{e}_{i_{2}} \right\rangle - \frac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{n_{i}} \frac{n_{i}}{2} (\alpha_{j,c}^{i})^{2} \quad (6)$$

After making the denotations that $\boldsymbol{\alpha}_c = [\alpha_{1,c}^1, \dots, \alpha_{n_m,c}^m]^T$, $K_F(\mathbf{x}_{j_1}^{i_1}, \mathbf{x}_{j_2}^{i_2}) = \mathbf{x}_{j_1}^{i_1}^T \mathbf{D}(\lambda_1 \mathbf{D} + \lambda_2 \mathbf{I}_d)^{-1} \mathbf{x}_{j_2}^{i_2} \langle \mathbf{e}_{i_1}, \mathbf{e}_{i_2} \rangle$, and $\boldsymbol{\Lambda}$ is the diagonal matrix whose diagonal element equals to n_i if the corresponding observation belongs to the *i*-th task, we can write problem (6) briefly as follows:

$$\max_{\boldsymbol{\alpha}_{c}} \sum_{i=1}^{m} \sum_{j=1}^{n_{i}} \alpha_{j,c}^{i} \bar{y}_{j,c}^{i} - \frac{1}{2} \boldsymbol{\alpha}_{c}^{T} \left(\mathbf{K}_{\mathrm{F}} + \frac{1}{2} \boldsymbol{\Lambda} \right) \boldsymbol{\alpha}_{c}$$
(7)

 Xiao-Lei Zhang was with Tsinghua National Laboratory for Information Science and Technology, Department of Electronic Engineering, Tsinghua University, Beijing, China, 100084.
 E-mail: huoshan6@126.com

This work was finished when the author was a visiting scholar with the Department of Computer Science and Engineering, Ohio State University, OH, USA, 43210.

2 DERIVATION OF EQUATION (11) IN THE MAIN PAPER

Again, we get the Lagrange of the problem in the braces of equation (10) in the main paper as follows:

$$\mathcal{L}(\{\theta_i\}_{i=1}^m) = \sum_{i=1}^m \theta_i - \frac{1}{2} \sum_{c=1}^C \boldsymbol{\alpha}_c^T \widetilde{\mathbf{K}}_{\mathbf{F}} \boldsymbol{\alpha}_c - \sum_{i=1}^m \sum_{k: \mathbf{Y}_k^i \in \mathcal{B}^i} \mu_k^i \left(\theta_i - \sum_{c=1}^C \sum_{j=1}^{n_i} \alpha_{j,c}^i y_{j,c}^j \right)$$
(8)

where $\mu_k^i \ge 0$ is the Lagrange variable. Calculating the partial derivative of problem (8) with respect to θ and letting the derivative equals to 0 can get the following equation:

$$\frac{\partial \mathcal{L}}{\partial \theta_i} = 1 - \sum_{k: \mathbf{Y}_k^i \in \mathcal{B}^i} \mu_k^i = 0 \tag{9}$$

We get

$$\sum_{k:\mathbf{Y}_{k}^{i}\in\mathcal{B}^{i}}\mu_{k}^{i}=1$$
(10)

Substituting (10) to (8) can derive equation (11) in the main paper.

3 OPTIMIZATION PROCEDURE

In this section, we will present the optimization procedure in detail. The readers can directly use the content of this section to replace Section 6 of the main paper.

We are to solve DMTFC (equation (10) in the main paper) and DMTRC (equation (16) in the main paper) in a uniform framework. To facilitate the mathematical representation, we write them as the following uniform objective:

$$\max_{\{\boldsymbol{\alpha}_{c}\}_{c=1}^{C}} \min_{\mathbf{Z} \in \mathcal{Z}} \left\{ \max_{\{\theta_{i}\}_{i=1}^{m}} \sum_{i=1}^{m} \theta_{i} - \frac{1}{2} \sum_{c=1}^{C} \boldsymbol{\alpha}_{c}^{T} \widetilde{\mathbf{K}} \boldsymbol{\alpha}_{c} \right.$$

$$\mathrm{s.t.} \theta_{i} \leq \sum_{c=1}^{C} \sum_{j=1}^{n_{i}} \alpha_{j,c}^{i} y_{j,c}^{j}, \forall i = 1, \dots, m, \forall k : \mathbf{Y}_{k}^{i} \in \mathcal{B}^{i} \right\}.$$

$$(11)$$

where **Z** stands for **D** or Ω , \mathcal{Z} stands for \mathcal{D} or \mathcal{A} , and \widetilde{K} stands for \widetilde{K}_F or \widetilde{K}_R .

The solution framework is an alternating method. First, it decomposes the unsupervised problem (11) to a serial supervised multiclass MTL problem by the cutting-plane algorithm (CPA) [1] and the extended level method (ELM) [2], [3], where the decomposition algorithm can be seen as a multitask extension of the SVR-M3C algorithm [4]. Then, it solves each supervised multiclass MTL problem in an alternating way, which decomposes the multiclass MTL to a serial supervised single-task regression problems eventually. Note that the difference of the optimization procedure between DMTFC and DMTRC only appears in the supervised learning in Section 3.3.

3.1 Optimizing (11) Via Cutting Plane Algorithm

Because the number of the constraints in problem (11) is exponential large with respect to n, directly optimizing (11) is impossible when the data set contains over dozens of examples. Hence, we adopt the cutting-plane algorithm [1] to solve problem (11) approximately.

We present the key idea of the cutting-plane algorithm as follows. Generally, given a constrained optimization problem, the cutting plane algorithm alternates the following two steps until the objective value converges. The first step is to solve a reduced problem of the constrained problem, i.e. a problem that contains only a part of the constraints. The second step is to find the most violated constraint of the reduced problem, and add it to the constraint set so as to form a new reduced problem for the next iteration. It has been proved that the number of the cutting-plane iterations is upper bounded by $O(1/\epsilon)$ [5], where ϵ is a user defined cutting-plane solution precision.

For problem (11), the cutting-plane algorithm iterates the following two steps:

a) Solving the following reduced problem of problem (11):

$$\max_{\{\boldsymbol{\alpha}_{c}\}_{c=1}^{C}} \min_{\mathbf{Z} \in \mathcal{Z}} \left\{ \max_{\{\theta_{i}\}_{i=1}^{m}} \sum_{i=1}^{m} \theta_{i} - \frac{1}{2} \sum_{c=1}^{C} \boldsymbol{\alpha}_{c}^{T} \widetilde{\mathbf{K}} \boldsymbol{\alpha}_{c} \right.$$

$$\mathrm{s.t.} \theta_{i} \leq \sum_{c=1}^{C} \sum_{j=1}^{n_{i}} \alpha_{j,c}^{i} y_{j,c}^{i}, \forall i = 1, \dots, m, \forall k : \mathbf{Y}_{k}^{i} \in \mathcal{Y}^{i} \right\}.$$

$$(12)$$

where \mathcal{Y}^i is the reduced constraint subset of \mathcal{B}^i . Problem (12) is equivalent to

$$\max_{\{\boldsymbol{\alpha}_{c}\}_{c=1}^{C}} \min_{\mathbf{Z}\in\mathcal{Z}} \max_{\{\theta_{i}\}_{i=1}^{m}} \sum_{i=1}^{m} \theta_{i} - \frac{1}{2} \sum_{c=1}^{C} \boldsymbol{\alpha}_{c}^{T} \widetilde{\mathbf{K}} \boldsymbol{\alpha}_{c}$$

$$\text{s.t.} \theta_{i} \leq \sum_{c=1}^{C} \sum_{j=1}^{n_{i}} \alpha_{j,c}^{i} y_{j,c}^{i}, \forall i = 1, \dots, m, \forall k : \mathbf{Y}_{k}^{i} \in \mathcal{Y}^{i} \bigg\}.$$

$$(13)$$

where $\mathcal{M}_{\mathcal{Y}}^{i} = \left\{ \boldsymbol{\mu}^{i} | 0 \leq \boldsymbol{\mu}_{k}^{i} \leq 1, \sum_{k=1}^{|\mathcal{Y}^{i}|} \boldsymbol{\mu}_{k}^{i} = 1 \right\}$ with $|\mathcal{Y}^{i}|$ denoted as the size of \mathcal{Y}^{i} . Here, we leave this complicated problem to Section 3.2.

b) Calculating the most violated constraint $\left\{\mathbf{Y}_{|\mathcal{Y}^i|+1}^i\right\}_{i=1}^m$ by solving the following problem and adding $\mathbf{Y}_{|\mathcal{Y}^i|+1}^i$ to \mathcal{Y}^i .

$$\min_{\left\{\mathbf{Y}_{|\mathcal{Y}^{i}|+1}^{i}\right\}_{i=1}^{m}} \sum_{i=1}^{m} \sum_{c=1}^{C} \sum_{j=1}^{n_{i}} \alpha_{j,c}^{i} \bar{y}_{|\mathcal{Y}^{i}|+1,j,c}^{i}.$$
(14)

We can observe that the subitems $\sum_{c} \sum_{j} \alpha_{j,c}^{i} \bar{y}_{|\mathcal{Y}^{i}|+1,j,c}^{i}$ are mutually independent with respect to *i*. Therefore, optimizing problem (14) is equivalent to optimizing the summation of the following problems:

$$\max_{\mathbf{Y}_{|\mathcal{Y}^{i}|+1}^{i}} \sum_{c=1}^{C} \sum_{j=1}^{n_{i}} \alpha_{j,c}^{i} \bar{y}_{|\mathcal{Y}^{i}|+1,j,c}^{i} , \ \forall i = 1, \dots, m.$$
(15)

Although the above problem is a binary integer matrix optimization problem, it can be solved in time $\mathcal{O}(\sum_{i=1}^{m} Cn_i \log(Cn_i))$ thanks to the constraints of **Y**^{*i*} in equation (3) of the main paper. See [4, Algorithm 6] for the efficient algorithm.

3.2 Optimizing (13) Via Extended Level Method

Like the full problem (11), the cutting-plane subproblem (13) also has an equivalent form:

$$\max_{\{\boldsymbol{\alpha}_{c}\}_{c=1}^{C}} \min_{\mathbf{Z} \in \mathcal{Z}} \min_{\{\boldsymbol{\mu}^{i} \in \mathcal{M}_{\mathcal{Y}}^{i}\}_{i=1}^{m}} \sum_{i=1}^{m} \sum_{c=1}^{C} \sum_{j=1}^{n_{i}} \alpha_{j,c}^{i} \sum_{k: \mathbf{Y}_{k}^{i} \in \mathcal{Y}^{i}} \mu_{k}^{i} \bar{y}_{k,j,c}^{i} - \frac{1}{2} \sum_{c=1}^{C} \boldsymbol{\alpha}_{c}^{T} \widetilde{\mathbf{K}} \boldsymbol{\alpha}_{c}$$
(16)

Problem (13) is a concave-convex optimization problem that is convex on μ and **Z** and concave on α . We will optimize it via the efficient Extended Level Method (ELM) [2], [3].

We present the key idea of the ELM algorithm as follows. ELM tries to solve the concave-convex optimization problem $\min_a \max_b f(a, b)$ that is convex on a and concave on b by iteratively constructing tighter upper and lower bounds for the optimal objective value $f(a^*; b^*)$, where (a^*, b^*) denotes the optimal solution. Specifically, it iterates the following two steps. The first step is to construct the lower bound $\underline{f}_s = \min_a \max_{1 \le r \le s} f(a; b_r)$ and the upper bound $\overline{f}^s = \min_{1 \le r \le s} f(a_r; b_r)$ of f, where r and s denotes the indices of the iterations (i.e. solutions) and $\max_{1 \le r \le s} f(a; b_r)$ is also a cutting-plane model. The second step is to first get a_{s+1} by solving the following optimization problem

$$\min_{a_{s+1}} \|a_{s+1} - a_s\|^2$$
s.t. $f(a_{s+1}; b_r) \le \tau \overline{f}_s + (1 - \tau) \underline{f}_s, \forall r = 1, \dots, s,$

$$(17)$$

and then get b_{s+1} by solving $\max_{b_{s+1}} f(a_{s+1}, b_{s+1})$, where τ is a user defined constant. equation (17) performs like a regularizer that prevents a_{s+1} far from a_s .

For problem (16), because optimizing μ and \mathbf{Z} jointly is difficult, setting $a = {\mu, \mathbf{Z}}$ is improper. We propose to set $a = \mu$ and optimize \mathbf{Z} and α jointly. It is easy to prove the correctness of this new optimization strategy. The proof is similar with the proof of [2, Theorem 1]. Another very important issue is that to make the cutting-plane algorithm presented in Section 3.1 converges, for problem (16) at the *S*-th cutting-plane iteration, we should inherit all previous S-1 ELM models to initialize the upper and lower bounds of the problem, otherwise, the cutting-plane algorithm will fail. The proof is the same as [4, Theorem 3].

With the aforementioned two key points, the ELM algorithm for problem (13) is presented as follows. Suppose we are currently at the *S*-th cutting-plane subproblem. That is to say, we are to solve the *S*-th problem (13). Suppose solving the *R*-th cutting-plane subproblem yields T_R ELM solutions, denoted as $\left\{ \left\{ \alpha_{c,r}^R \right\}_{c=1}^C, \mathbf{Z}_r^R, \left\{ \boldsymbol{\mu}_r^{i,R} \right\}_{i=1}^m \right\}_{r=1}^{T_R}$, where $R = 1, \ldots, S - 1$. Suppose the constraint set of the *i*-th task at the *R*-th cutting plane iteration, denoted as $\mathcal{Y}^{i,R}$, contains $|\mathcal{Y}^{i,R}|$ constraints $(|\mathcal{Y}^{i,R}| \leq R)$.

Initialization of ELM. All previous S - 1 ELM models should be inherited by adding $\boldsymbol{\mu}_r^{i,R}$ with $|\mathcal{Y}^{i,S}| - |\mathcal{Y}^{i,R}|$ zeros: $\boldsymbol{\mu}_r^{i,R'} = \left[\left(\boldsymbol{\mu}_r^{i,R} \right)^T, \mathbf{0}_{|\mathcal{Y}^{i,S}| - |\mathcal{Y}^{i,R}|}^T \right]^T, \forall i = 1, \dots, m, \forall r = 1, \dots, T_R$. Without lose of generality, we further denote all inherited ELM solutions as $\left\{ \{ \boldsymbol{\alpha}_c^r \}_{c=1}^C, \mathbf{Z}_r, \{ \boldsymbol{\mu}^{r,i} \}_{i=1}^m \}_{r=1}^s = \left\{ \left\{ \left\{ \boldsymbol{\alpha}_{c,r}^R \right\}_{c=1}^C, \mathbf{Z}_r^R, \left\{ \boldsymbol{\mu}_r^{i,R'} \right\}_{i=1}^m \right\}_{r=1}^{r_R} \right\}_{R=1}^{S-1}, \forall i = 1, \dots, m, \forall r =$

ELM. The ELM algorithm for (13) iterates the following steps:

a) Constructing the lower bound \underline{h}_s and upper bound \overline{h}_s by

$$\underline{h}_{s} = \min_{\left\{\boldsymbol{\mu}^{i} \in \mathcal{M}_{\mathcal{Y}}^{i}\right\}_{i=1}^{m}} \max_{1 \leq r \leq s} \sum_{i=1}^{m} \sum_{c=1}^{C} \sum_{j=1}^{n_{i}} \alpha_{r,j,c}^{i} \sum_{k: \mathbf{Y}_{k}^{i} \in \mathcal{Y}^{i}} \mu_{k}^{i} \bar{y}_{k,j,c}^{i} - \frac{1}{2} \sum_{c=1}^{C} \boldsymbol{\alpha}_{r,c}^{T} \widetilde{\mathbf{K}}_{r} \boldsymbol{\alpha}_{r,c},$$
(18)

$$\overline{h}_{s}^{i} = \min_{1 \leq r \leq s} \sum_{i=1}^{m} \sum_{c=1}^{C} \sum_{j=1}^{n_{i}} \alpha_{r,j,c}^{i} \sum_{k: \mathbf{Y}_{k}^{i} \in \mathcal{Y}^{i}} \mu_{r,k}^{i} \overline{y}_{k,j,c}^{i} - \frac{1}{2} \sum_{c=1}^{C} \alpha_{r,c}^{T} \widetilde{\mathbf{K}}_{r} \alpha_{r,c}$$

$$\tag{19}$$

where \mathbf{K}_r is a function of \mathbf{Z}_r . With \mathbf{Z}_r fixed, the tasks are mutually independent, hence we can get the bounds of each task separately:

$$\underline{h}_{s}^{i} = \min_{\boldsymbol{\mu}^{i} \in \mathcal{M}_{\mathcal{Y}}^{i}} \max_{1 \leq r \leq s} \sum_{c=1}^{C} \sum_{j=1}^{n_{i}} \alpha_{r,j,c}^{i} \sum_{k: \mathbf{Y}_{k}^{i} \in \mathcal{Y}^{i}} \mu_{k}^{i} \bar{y}_{k,j,c}^{i}, \tag{20}$$

$$\overline{h}_{s}^{i} = \min_{1 \le r \le s} \sum_{c=1}^{C} \sum_{j=1}^{n_{i}} \alpha_{r,j,c}^{i} \sum_{k: \mathbf{Y}_{k}^{i} \in \mathcal{Y}^{i}} \mu_{r,k}^{i} \overline{y}_{k,j,c}^{i}$$

$$\tag{21}$$

b) Get $\{\mu_{s+1}^i\}_{i=1}^m$ by

$$\min_{\{\boldsymbol{\mu}_{s+1}^{i}\in\mathcal{M}_{\mathcal{V}}^{i}\}_{i=1}^{m}} \sum_{i=1}^{m} \|\boldsymbol{\mu}_{s+1}^{i}-\boldsymbol{\mu}_{s}^{i}\|^{2} \qquad (22)$$
s.t.
$$\sum_{i=1}^{m} \sum_{c=1}^{C} \sum_{j=1}^{n_{i}} \alpha_{r,j,c}^{i} \sum_{k:\mathbf{Y}_{k}^{i}\in\mathcal{Y}^{i}} \mu_{s+1,k}^{i} \bar{y}_{k,j,c}^{i} - \frac{1}{2} \sum_{c=1}^{C} \boldsymbol{\alpha}_{r,c}^{T} \tilde{\mathbf{K}}_{r} \boldsymbol{\alpha}_{r,c} \leq \tau \bar{h}_{s} + (1-\tau) \underline{h}_{s}, \quad \forall r = 1, \dots, s.$$

Similar with Step a), we replace problem (22) with the summation of the following problems

$$\min_{\boldsymbol{\mu}_{s+1}^{i} \in \mathcal{M}_{\mathcal{Y}}^{i}} \|\boldsymbol{\mu}_{s+1}^{i} - \boldsymbol{\mu}_{s}^{i}\|^{2}$$
s.t.
$$\sum_{i=1}^{m} \sum_{c=1}^{C} \sum_{j=1}^{n_{i}} \alpha_{r,j,c}^{i} \sum_{k: \mathbf{Y}_{k}^{i} \in \mathcal{Y}^{i}} \mu_{s+1,k}^{i} \bar{y}_{k,j,c}^{i} \leq \tau \bar{h}_{s}^{i} + (1-\tau) \underline{h}_{s}^{i}, \quad \forall r = 1, \dots, s.$$
(23)

c) Get $(\alpha_{s+1}, \mathbf{Z}_{s+1})$ by

$$\min_{\mathbf{Z}_{s+1}\in\mathcal{Z}}\max_{\{\boldsymbol{\alpha}_{s+1,c}\}_{c=1}^{C}}\sum_{i=1}^{m}\sum_{c=1}^{C}\sum_{j=1}^{n_{i}}\alpha_{s+1,j,c}^{i}\sum_{k:\mathbf{Y}_{k}^{i}\in\mathcal{Y}^{i}}\mu_{s+1,k}^{i}\bar{y}_{k,j,c}^{i}-\frac{1}{2}\sum_{c=1}^{C}\alpha_{s+1,c}^{T}\widetilde{\mathbf{K}}_{s+1}\alpha_{s+1,c}.$$
(24)

Here, we leave problem (24) to Section 3.3.

3.3 Optimizing (24) Via the Alternating Method

Problem (24) is a supervised multiclass MTL problem. We adopt an alternating method that is similar with [6] for it.

Specifically, for each problem (24), the method iterates the following two steps:

a) Given fixed **Z**, we aim to solve

$$\max_{\{\boldsymbol{\alpha}_{c}\}_{c=1}^{C}} \sum_{i=1}^{m} \sum_{c=1}^{C} \sum_{j=1}^{n_{i}} \alpha_{j,c}^{i} \sum_{k: \mathbf{Y}_{k}^{i} \in \mathcal{Y}^{i}} \mu_{k}^{i} \bar{y}_{k,j,c}^{i} - \frac{1}{2} \sum_{c=1}^{C} \boldsymbol{\alpha}_{c}^{T} \widetilde{\mathbf{K}} \boldsymbol{\alpha}_{c}$$
$$= \sum_{c=1}^{C} \max_{\boldsymbol{\alpha}_{c}} \sum_{i=1}^{m} \sum_{j=1}^{n_{i}} \alpha_{j,c}^{i} \sum_{k: \mathbf{Y}_{k}^{i} \in \mathcal{Y}^{i}} \mu_{k}^{i} \bar{y}_{k,j,c}^{i} - \frac{1}{2} \boldsymbol{\alpha}_{c}^{T} \widetilde{\mathbf{K}} \boldsymbol{\alpha}_{c}$$
(25)

When Z is fixed, the subitems in the right side of equation (25) are mutually independently. Hence, we solve each item

$$\max_{\boldsymbol{\alpha}_{c}} \sum_{i=1}^{m} \sum_{j=1}^{n_{i}} \alpha_{j,c}^{i} \sum_{\boldsymbol{k}: \mathbf{Y}_{k}^{i} \in \mathcal{Y}^{i}} \mu_{k}^{i} \bar{y}_{k,j,c}^{i} - \frac{1}{2} \boldsymbol{\alpha}_{c}^{T} \widetilde{\mathbf{K}} \boldsymbol{\alpha}_{c}, \quad \forall c = 1, \dots, C$$
(26)

independently, which is a supervised regression problem.

b) Given fixed $\{\alpha_c\}_{c=1}^C$, we aim to optimize

$$\min_{\mathbf{Z}\in\mathcal{Z}}\sum_{i=1}^{m}\sum_{c=1}^{C}\sum_{j=1}^{n_{i}}\alpha_{j,c}^{i}\sum_{k:\mathbf{Y}_{k}^{i}\in\mathcal{Y}^{i}}\mu_{k}^{i}\bar{y}_{k,j,c}^{i}-\frac{1}{2}\sum_{c=1}^{C}\boldsymbol{\alpha}_{c}^{T}\widetilde{\mathbf{K}}\boldsymbol{\alpha}_{c}.$$
(27)

For solving problems (26) and (27), DMTFC and DMTRC should be considered separately as follows:

Specifying (26) and (27) as a part of DMTFC: We replace **Z** and \mathcal{Z} by **D** and \mathcal{D} respectively in the equations. For (26), the multitask kernel $\widetilde{\mathbf{K}}$ should be specified by equation (7) in the main paper. For (27), we can get the closed solution of **D** as $\mathbf{D} = \frac{\left(\sum_{c=1}^{C} \mathbf{W}_{c} \mathbf{W}_{c}^{T}\right)^{\frac{1}{2}}}{2}$ where **W** is defined in (8) in the main paper. The derivation is analogous

solution of **D** as $\mathbf{D} = \frac{\left(\sum_{c=1}^{C} \mathbf{W}_{c} \mathbf{W}_{c}^{T}\right)^{\frac{1}{2}}}{\operatorname{tr}\left(\left(\sum_{c=1}^{C} \mathbf{W}_{c} \mathbf{W}_{c}^{T}\right)^{\frac{1}{2}}\right)}$ where \mathbf{W}_{c} is defined in (8) in the main paper. The derivation is analogous

to [7, Appendix 1].

Specifying (26) and (27) as a part of DMTRC: We replace \mathbf{Z} and \mathcal{Z} by $\boldsymbol{\Omega}$ and \mathcal{A} respectively in the equations. For (26), $\widetilde{\mathbf{K}}$ should be specified by equation (17) in the main paper. For (27), we can get the closed solution of $\boldsymbol{\Omega}$ as $\boldsymbol{\Omega} = \frac{\left(\sum_{c=1}^{C} \mathbf{W}_{c}^{T} \mathbf{W}_{c}\right)^{\frac{1}{2}}}{\operatorname{tr}\left(\left(\sum_{c=1}^{C} \mathbf{W}_{c}^{T} \mathbf{W}_{c}\right)^{\frac{1}{2}}\right)}$ where \mathbf{W}_{c} is defined in (18) in the main paper. The derivation is analogous to [6, equation 12].

13].

3.4 Overview of the Algorithm

Observing that the relaxed DMTFC (equation (10) in the main paper) and DMTRC (equation (16) in the main paper) are quite similar, we propose a unified convex optimization objective (11), which can be solved alternatively by combining several existing efficient algorithms [1]–[4], [8], [9]. The optimization procedure is summarized as Algorithm 1.

4 SUPPLEMENT TO THE RESULTS OF THE PENDIGITS DATASET

In this section, we will first give the average performance of all tasks in the metrics of Accuracy (ACC), Adjusted Rand Index (ARI), and Normalized Mutual Information (NMI) respectively. Then, we will give the performance of each individual task in different jobs in the metric of NMI.

4.1 Average Performance of All Tasks in Different Metrics

Figs. 1, 2, and 3 show the average performance on the Pendigits dataset in the metrics of NMI, ARI, and ACC respectively.



Fig. 1. Comparison of the average NMI of all tasks on the Pendigits dataset. The error bar denotes the standard deviation.

Algorithm 1 Discriminative Multitask Clustering.

Input: The unlabeled observations from *m* tasks $\{\{\mathbf{x}_{j}^{i}\}_{j=1}^{n_{i}}\}_{i=1}^{m}$, class number *P*, regularization parameters λ_{1} and λ_{2} , and the parameters that control the class balance $\{l_{i}\}_{i=1}^{m}$

Output: ŷ

Initialization: Arbitrary initial constraint label matrices $\{\mathbf{Y}_1^i\}_{i=1}^m$ that satisfy the constraints in equation (3) in the main paper, initial constraint sets $\mathcal{Y}^i \leftarrow \{\mathbf{Y}_1^i\}$ and $\mu_1^i \leftarrow 1$ with $i = 1, \dots, m$

- 1: repeat
- 2: Inherit the ELM solutions in the previous cutting-plane iterations
- 3: repeat
- 4: for task $i = 1, \ldots, m$ do
- 5: Construct the lower bound via (20) and the upper bound via (21)
- 6: Update μ^i by solving (23)
- 7: end for
- 8: repeat
- 9: for class $c = 1, \ldots, C$ do
- 10: Update α_c by solving (26)
- 11: end for
- 12: Update **Z** by solving (27)
- 13: **until** the objective value converges
- 14: **until** the objective value converges
- 15: **for** task i = 1, ..., m **do**
- 16: Calculate the most violated constraint $\mathbf{Y}_{|\mathcal{C}^i|+1}^i$ of the cutting-plane algorithm by solving (15)
- 17: $\mathcal{Y}^i \leftarrow \mathcal{Y}^i \cup \left\{ \mathbf{Y}^i_{|\mathcal{Y}^i|+1} \right\}$
- 18: end for

19: **until** the objective value converges or $\{\mathcal{Y}^i\}_{i=1}^m$ are unchanged

- 20: /*Prediction*/
- 21: for task i = 1, ..., m do
- 22: **for** $j = 1, ..., n_i$ **do**
- 23: $\hat{y}_{j}^{i} \leftarrow \arg \max_{p} \sum_{u=1}^{m} \sum_{v=1}^{n_{u}} \alpha_{v,p}^{u} K_{MT}(\mathbf{x}_{v}^{u}, \mathbf{x}_{j}^{i})$, where the kernel function K_{MT} is defined in (7) in the main paper for DMTFC and in (17) in the main paper for DMTRC
- 24: end for
- 25: **end for**



Fig. 2. Comparison of the average ARI of all tasks on the Pendigits dataset.



Fig. 3. Comparison of the average ACC of all tasks on the Pendigits dataset.

4.2 Performance of Each Individual Task in the Metric of NMI

Figs. 4, 5, and 6 show the NMIs of the clustering algorithms with respect to each individual task in Jobs 1, 2, and 3 respectively.

In this subsection, we should pay particular attention to the standard deviations of the NMIs on each individual task, since the standard deviations represent the stabilities of the clustering algorithms partially.



Fig. 4. NMI comparison with respect to each individual task in Job 1 on the Pendigits dataset.







Fig. 6. NMI comparison with respect to each individual task in Job 3 on the Pendigits dataset.

5 SUPPLEMENT TO THE RESULTS OF THE MULTI-DOMAIN NEWSGROUPS DATASET

In this section, we will first give the average performance of all tasks in the metrics of NMI, ARI, and ACC respectively. Then, we will give the performance of each individual task with 5%, 10%, 20% and 40% data of the 20-newsgroups dataset in the metric of NMI.

5.1 Average Performance of All Tasks in Different Metrics

Figs. 7, 8, and 9 show the average performance on the 20-newsgroups dataset in the metrics of NMI, ARI, and ACC respectively.



Fig. 7. Comparison of the average NMI of all tasks on the 20-Newgroups dataset.



Fig. 8. Comparison of the average ARI of all tasks on the 20-Newgroups dataset.





5.2 Performance of Each Individual Task in the Metric of NMI

Figs. 10, 11, 12, and 13 show the NMIs of the clustering algorithms with 5%, 10%, 20% and 40% data of the 20-newsgroups dataset in the metric of NMI.

In this subsection, we should pay particular attention to the standard deviations of the NMIs on each individual task, since the standard deviations represent the stabilities of the clustering algorithms partially. From the figures, we can see that when more data is available, the stability of DMTRC is improved greatly.



Fig. 10. NMI comparison with respect to each individual task on 5% data of the 20-Newgroups dataset.



Fig. 11. NMI comparison with respect to each individual task on 10% data of the 20-Newgroups dataset.



Fig. 12. NMI comparison with respect to each individual task on 20% data of the 20-Newgroups dataset.



Fig. 13. NMI comparison with respect to each individual task on 40% data of the 20-Newgroups dataset.

6 SUPPLEMENT TO THE RESULTS OF THE MULTI-DOMAIN SENTIMENT DATASET

In this section, we will first give the average performance of all tasks in the metrics of NMI, ARI, ACC and the Area under ROC Curve (AUC) respectively. Then, we will give the performance of each individual task with 10%, 30%, and 50% data of the sentiment dataset in the metric of NMI. At last, we will show the CPU time of the methods.

6.1 Average Performance of All Tasks in Different Metrics

Figs. 14, 15, 16 and 17 show the average performance on the sentiment dataset in the metrics of NMI, ARI, ACC and AUC respectively.



Fig. 14. Comparison of the average NMI of all tasks on the sentiment dataset.



Fig. 15. Comparison of the average ARI of all tasks on the sentiment dataset.



Fig. 16. Comparison of the average ACC of all tasks on the sentiment dataset.



Fig. 17. Comparison of the average AUC of all tasks on the sentiment dataset.

6.2 Performance of Each Individual Task in the Metric of NMI

Figs. 18, 19, and 20 show the NMIs of the clustering algorithms with 10%, 30%, and 50% data of the 20-newsgroups dataset in the metric of NMI.

From the figures, we can see that the proposed clustering methods are much more robust than KM and KKM.



Fig. 18. NMI comparison with respect to each individual task on 10% data of the sentiment dataset.



Fig. 19. NMI comparison with respect to each individual task on 30% data of the sentiment dataset.



Fig. 20. NMI comparison with respect to each individual task on 50% data of the sentiment dataset.

REFERENCES

- J. E. Kelley, "The cutting-plane method for solving convex programs," J. Soc. Ind. Appl. Math., vol. 8, no. 4, pp. 703–712, 1960.
 Z. Xu, R. Jin, I. King, and M. R. Lyu, "An extended level method for efficient multiple kernel learning," in Proc. Adv. Neural Inform. Process. Syst., vol. 21, 2009, pp. 1825-1832.
- [3] H. Yang, Z. Xu, J. Ye, I. King, and M. R. Lyu, "Efficient sparse generalized multiple kernel learning," IEEE Trans. Neural Netw., no. 3, pp. 433-446, 2011.
- [4] X. L. Zhang and J. Wu, "Linearithmic time sparse and convex maximum margin clustering," IEEE Trans. Syst., Man, Cybern. B, Cybern., vol. 1, no. 99, pp. 1–24, 2012.
- C. H. Teo, A. Smola, S. V. N. Vishwanathan, and Q. V. Le, "A scalable modular convex solver for regularized risk minimization," in Proc. [5] 13th ACM SIGKDD Int. Conf. Knowl. Disc. Data Min., 2007, pp. 727-736.
- [6] Y. Zhang and D. Y. Yeung, "A convex formulation for learning task relationships in multi-task learning," in Proc. 26th Conf. Uncertainty Artif. Intell., 2010, pp. 733–742. [7] A. Argyriou, T. Evgeniou, and M. Pontil, "Convex multi-task feature learning," Mach. Learn., vol. 73, no. 3, pp. 243–272, 2008.
- [8] Y. F. Li, I. W. Tsang, J. T. Kwok, and Z. H. Zhou, "Tighter and convex maximum margin clustering," in Proc. 12th Int. Conf. Artif. Intell. Statist., Clearwater Beach, FL, 2009, pp. 344-351.
- S. P. Boyd and L. Vandenberghe, Convex Optimization. Cambridge Univ. Pr., 2004. [9]