# Multi-source Domain Adaptation

kongyiji

January 2021

## 1 Problem Formulation

$$\min \frac{\{\text{domain shift }\}_{\text{marginal}} + \{\text{domain shift}\}_{\text{conditional}} + \{\text{dist}\}_{\text{intra}}}{\{\text{dist}\}_{\text{inter}}} \tag{1}$$

### 1.1 Ditribution Matching and Landmark Selection

The model considers both marginal and conditional distribution discrepancy, denoted as $E_{\text{MG}}$ and $E_{\text{CD}}$, which can be formulated as follows:

$$\min_{A,B} E_{\text{MG}}\left(\alpha, \beta, X_{\text{s}}, X_{\text{t}}, A, B\right) + E_{\text{CD}}\left(\alpha, \beta, X_{\text{s}}, X_{\text{t}}, A, B\right)$$
$$\text{s.t. } \{\alpha_{ui}^c, \beta_i^c\} \in [0,1], \frac{\alpha_{\text{u}}^{c\text{T}} \mathbf{1}_{n_{\text{s}}^{\text{u}}}}{n_{\text{s}}^{\text{uc}}} = \delta_{\text{s}}^{\text{u}}, \frac{\beta^{c\text{T}} \mathbf{1}_{n_{\text{t}}^c}}{n_{\text{t}}^c} = \delta_{\text{t}} \tag{2}$$

where $\alpha_{\text{u}} = \left[\alpha_{\text{u}}^1; \cdots; \alpha_{\text{u}}^c; \cdots; \alpha_{\text{u}}^C\right] \in R^{n_{\text{s}}^{\text{u}}}$ are the weights of samples in source domain $X_{\text{s}}^{\text{u}} \in R^{d_{\text{s}}^{\text{u}} \times n_{\text{s}}^{\text{u}}}, \alpha = [\alpha_1; \cdots; \alpha_{\text{u}}] \in R^{n_s}$ are the weights of samples in source domain $X_{\text{s}} = \left[X_{\text{s}}^1; X_{\text{s}}^2; \ldots; X_{\text{s}}^{N_{\text{s}}}\right] \in R^{d_s \times n_s}$. $\beta = \left[\beta^1; \cdots; \beta^c; \cdots; \beta^C\right] \in R^{n_t}$ are the weights of data in the source domain and the target domain, respectively, $\alpha_{\text{u}}^c = \left[\alpha_{\text{u}1}^c; \cdots; \alpha_{\text{u}n_{\text{s}}^{\text{uc}}}^c\right], \beta^c = \left[\beta_1^c; \cdots; \beta_{n^c}^c\right], \mathbf{1}_{n_{\text{s}}^{\text{uc}}} \in R^{n_{\text{s}}^{\text{uc}}}$ and $\mathbf{1}_{n_{\text{t}}^c} \in R^{n_{\text{t}}^c}$ are column vectors with all ones. $\delta_{\text{s}}^{\text{u}}, \delta_t \in [0,1]$ controls the ratio of landmarks in the whole source or target domain samples. The constraints on $\alpha$ and $\beta$ keep them from trivial solutions such as one-hot vectors that only align one sample from source and one sample from target. Then, $E_{\text{MG}}$ and $E_{\text{CD}}$ in Eq.(2) further can be calculated by:

$$E_{\text{MG}} = \sum_{u=1}^{N_{\text{s}}} \left\| \frac{1}{\delta_{\text{s}}^{\text{u}} n_{\text{s}}^{\text{u}}} \sum_{i=1}^{n_{\text{s}}^{\text{u}}} \alpha_{\text{u}i} A^{\text{T}} x_{\text{s}}^{\text{u}i} - \frac{1}{\delta_{\text{t}} n_{\text{t}}} \sum_{j=1}^{n_{\text{t}}} \beta_{\text{j}} B^{\text{T}} x_{\text{t}}^j \right\|^2$$

$$= \sum_{u=1}^{N_{\text{s}}} \left\| \frac{1}{\delta_{\text{s}}^{\text{u}} n_{\text{s}}^{\text{u}}} A^{\text{T}} \left[\ \text{x}_{\text{u}1} \text{x}_{\text{u}2} \cdots \text{x}_{n_{\text{s}}^{\text{u}}}\ \right]_{1 \times n_{\text{s}}^{\text{u}}} \begin{bmatrix} \alpha_{\text{u}1} \\ \alpha_{\text{u}2} \\ \vdots \\ \alpha_{n_{\text{s}}^{\text{u}}} \end{bmatrix}_{n_{\text{s}}^{\text{u}} \times 1} - \frac{1}{\delta_t n_t} B^{\text{T}} \left[\ \text{x}_1 \text{x}_2 \cdots \text{x}_{n_t}\ \right]_{1 \times n_t} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_{n_t} \end{bmatrix}_{n_t \times 1} \right\|$$

$$= \sum_{u=1}^{N_{\text{s}}} \text{tr} \left( \frac{1}{\delta_{\text{s}}^{\text{u}2} n_{\text{s}}^{\text{u}2}} A^{\text{T}} X_{\text{s}}^{\text{u}} \alpha_{\text{u}} \left(A^{\text{T}} X_{\text{s}}^{\text{u}} \alpha_{\text{u}}\right)^{\text{T}} + \frac{1}{\delta_t^2 n_t^2} B^{\text{T}} X_t \beta \left(B^{\text{T}} X_t \beta\right)^{\text{T}} - \frac{1}{\delta_{\text{s}}^{\text{u}} \delta_t n_{\text{s}}^{\text{u}} n_t} A^{\text{T}} X_{\text{s}}^{\text{u}} \alpha_{\text{u}} \left(B^{\text{T}} X_t \beta\right)^{\text{T}} \right.$$

$$\left. - \frac{1}{\delta_{\text{s}}^{\text{u}} \delta_t n_{\text{s}}^{\text{u}} n_t} B^{\text{T}} X_t \beta \left(A^{\text{T}} X_{\text{s}}^{\text{u}} \alpha_{\text{u}}\right)^{\text{T}} \right)$$

$$=\sum_{u=1}^{N_s} \text{tr}\left(\frac{1}{\delta_s^{u2} n_s^{u2}} A^T X_s^u \alpha_u \alpha_u^T X_s^{uT} A + \frac{1}{\delta_t n_t^2} B^T X_t \beta \beta^T X_t^T B - \frac{1}{\delta_s^u \delta_t n_s^u n_t} A^T X_s^u \alpha_u \beta^T X_t^T - \frac{1}{\delta_s^u \delta_t n_s^u n_t} B^T X_t \beta \alpha_u^T X_s^{uT}\right)$$

$$= \text{tr}\left(A^T \left(\sum_{u=1}^{N_s} X_s^u H_{sm}^u X_s^{uT}\right) A + B^T \left(\sum_{u=1}^{N_s} X_t H_{tm} X_t^T\right) B - A^T \left(\sum_{u=1}^{N_s} X_s^u H_{stm}^u X_t^T\right) B$$

$$-B^T \left(\sum_{u=1}^{N_s} X_t H_{stm}^u{}^T X_s^{uT}\right) A\right)$$

$$\text{E}_{CD} = \sum_{u=1}^{N_s} \sum_{c=1}^{C} \left\| \frac{1}{\delta_s^{uc} n_s^{uc}} \sum_{i=1}^{n_s^{uc}} \alpha_{ui}^c A^T x_s^{i,c} - \frac{1}{\delta_t^c n_t^c} \sum_{j=1}^{n_t^c} \beta_i^c B^T x_t^{j,c} \right\|^2$$

$$=\sum_{c=1}^{C} \text{tr}\left(A^T \left(\sum_{u=1}^{N_s} X_s^u H_{sc}^u X_s^{uT}\right) A + B^T \left(\sum_{u=1}^{N_s} X_t H_{tc} X_t^T\right) B - A^T \left(\sum_{u=1}^{N_s} X_s^u H_{stc}^u X_t^T\right) B$$

$$-B^T \left(\sum_{u=1}^{N_s} X_t H_{stc}^u{}^T X_s^{uT}\right) A\right)$$

where

$$H_{sm}^u = \frac{1}{\delta_{us}^2 n_s^{u2}} \alpha_u \cdot \alpha_u^T, H_{tm} \quad = \frac{1}{\delta_t^2 n_t^2} \beta \cdot \beta^T, H_{stm}^u = \frac{1}{\delta_{us} \delta_t n_s^u n_t} \alpha_u \cdot \beta^T,$$

$$H_{sc}^u{}^c = \frac{1}{\delta_{us}^2 n_s^{uc2}} \alpha_u^c \cdot \alpha_u^{c^T}, H_{tc}^c = \frac{1}{\delta_{us}^2 n_t^{c2}} \beta^c \cdot \beta^{c^T}, H_{stc}^u{}^c = \frac{1}{\delta_{us} \delta_t n_s^{uc} n_t^c} \alpha_u^c \cdot \beta^{c^T},$$

After some algebra operations, Eq(2) can be written as the following equivalent equation:

$$A^T M_{ss} A + B^T M_{tt} B - A^T M_{st} B - B^T M_{ts} A$$

$$=A^T \left(\sum_{u=1}^{N_s} X_s^u (H_{sm}^u + H_{sc}^u) X_s^{uT}\right) A + B^T \left(\sum_{u=1}^{N_s} X_t (H_{tm} + H_{tc}) X_t^T\right) B$$

$$- A^T \left(\sum_{u=1}^{N_s} X_s^u (H_{stm}^u + H_{stc}^u) X_t^T\right) B - B^T \left(\sum_{u=1}^{N_s} X_t (H_{tsm}^u + H_{tsc}^u)^T X_s^{uT}\right) A \tag{3}$$

$$M_{ss} = \left(\sum_{u=1}^{N_s} X_s^u (H_{sm}^u + H_{sc}^u) X_s^{uT}\right)$$

$$M_{tt} = \left(\sum_{u=1}^{N_s} X_t (H_{tm} + H_{tc}) X_t^T\right)$$

$$M_{st} = -\left(\sum_{u=1}^{N_s} X_s^u (H_{stm}^u + H_{stc}^u) X_t^T\right) \tag{4}$$

$$M_{ts} = M_{st}^T$$

At last, the above equation(3), can be further transformed to its matrix form as follows:

$$\text{Tr}\left(\begin{bmatrix} A^T & B^T \end{bmatrix} \begin{bmatrix} M_{ss} & M_{st} \\ M_{ts} & M_{tt} \end{bmatrix} \begin{bmatrix} A \\ B \end{bmatrix}\right) \tag{5}$$

## 1.2 Structure Preservation

(a) Construct the intrinsic weight matrix $W_{\mathrm{w}}$ : For each sample $x$, connect the nearest neighbor pair $v$ and $x$ if $v$ has the same label information with $x$. (b) Construct the penalty weight matrix $W_{\mathrm{b}}$ : For each domain, connect the $k$ -nearest vertex pairs where samples in each pair belong to different classes.

$$
\min \sum_{u=1}^{\mathrm{N_s}} \frac{\mathrm{Tr}\left(A^{\mathrm{T}} X_{\mathrm{s}}^{\mathrm{u}} L_{\mathrm{b}}^{\mathrm{us}} X_{\mathrm{s}}^{\mathrm{uT}} A\right)}{\mathrm{Tr}\left(A^{\mathrm{T}} X_{\mathrm{s}}^{\mathrm{u}} L_{\mathrm{w}}^{\mathrm{us}} X_{\mathrm{s}}^{\mathrm{uT}} A\right)} = \min \frac{\mathrm{Tr}\left(A^{\mathrm{T}} S_{\mathrm{w}}^{\mathrm{s}} A\right)}{\mathrm{Tr}\left(A^{\mathrm{T}} S_{\mathrm{b}}^{\mathrm{s}} A\right)}
$$
$$
\min \frac{\mathrm{Tr}\left(B^{\mathrm{T}} X_{\mathrm{t}} L_{\mathrm{w}}^{\mathrm{t}} X_{\mathrm{t}}^{\mathrm{T}} B\right)}{\mathrm{Tr}\left(B^{\mathrm{T}} X_{\mathrm{t}} L_{\mathrm{b}}^{\mathrm{t}} X_{\mathrm{t}}^{\mathrm{T}} B\right)} = \min \frac{\mathrm{Tr}\left(B^{\mathrm{T}} S_{\mathrm{w}}^{\mathrm{t}} B\right)}{\mathrm{Tr}\left(B^{\mathrm{T}} S_{\mathrm{b}}^{\mathrm{t}} B\right)}
\tag{6}
$$

where

$$
S_{\mathrm{b}}^{\mathrm{s}} = \sum_{u=1}^{\mathrm{N_s}} X_{\mathrm{s}}^{\mathrm{u}} L_{\mathrm{b}}^{\mathrm{us}} X_{\mathrm{s}}^{\mathrm{uT}}, \quad S_{\mathrm{w}}^{\mathrm{s}} = \sum_{u=1}^{\mathrm{N_s}} X_{\mathrm{s}}^{\mathrm{u}} L_{\mathrm{w}}^{\mathrm{us}} X_{\mathrm{s}}^{\mathrm{uT}}
$$
$$
S_{\mathrm{b}}^{\mathrm{t}} = X_{\mathrm{t}} L_{\mathrm{b}}^{\mathrm{t}} X_{\mathrm{t}}^{\mathrm{T}}, \quad S_{\mathrm{w}}^{\mathrm{t}} = X_{\mathrm{t}} L_{\mathrm{w}}^{\mathrm{t}} X_{\mathrm{t}}^{\mathrm{T}}
$$

$$
\min_{A,B} \frac{\mathrm{Tr}\left( \begin{bmatrix} A^{\mathrm{T}} B^{\mathrm{T}} \end{bmatrix} \begin{bmatrix} M_{\mathrm{ss}} + \gamma S_{\mathrm{w}}^{\mathrm{s}} & M_{\mathrm{st}} \\ M_{\mathrm{ts}} & M_{\mathrm{tt}} + \gamma S_{\mathrm{w}}^{\mathrm{t}} + \mu I \end{bmatrix} \begin{bmatrix} A \\ B \end{bmatrix} \right)}{\mathrm{Tr}\left( \begin{bmatrix} A^{\mathrm{T}} & B^{\mathrm{T}} \end{bmatrix} \begin{bmatrix} \gamma S_{\mathrm{b}}^{\mathrm{s}} & \mathbf{0} \\ \mathbf{0} & \gamma S_{\mathrm{b}}^{\mathrm{t}} + \mu S_{\mathrm{h}}^{\mathrm{t}} \end{bmatrix} \begin{bmatrix} A \\ B \end{bmatrix} \right)}
\tag{7}
$$

where

$$
S_{\mathrm{h}}^{\mathrm{t}} = X_{\mathrm{t}} \left( I_{\mathrm{t}} - \frac{1}{n_{\mathrm{t}}} \mathbf{1}_{n_{\mathrm{t}}} \mathbf{1}_{n_{\mathrm{t}}}^{\mathrm{T}} \right) X_{\mathrm{t}}^{\mathrm{T}}
$$

is the covariance matrix of the target domain, to avoid projecting features into irrelevant dimensions, we encourage the variances of target domain is maximized in the respective subspaces. $\mathrm{Tr}\left(B^{T} B\right)$ constraint is further imposed small to control the scale of B. $\gamma$ and $\mu$ are trade-off parameters for the locality preserving term and the target variance term, respectively.

# 2 Problem Optimization

1) Optimizing the space mappings $A$ and $B$ : To optimize Eq. (7), we write $[A; B]$ as $P$. Thus, the objective function can be rewritten as:

$$
\min_{P} \frac{\mathrm{Tr}\left( P^{\mathrm{T}} \begin{bmatrix} M_{\mathrm{ss}} + \gamma S_{\mathrm{w}}^{\mathrm{s}} & M_{\mathrm{st}} \\ M_{\mathrm{ts}} & M_{\mathrm{tt}} + \gamma S_{\mathrm{w}}^{\mathrm{t}} + \mu I \end{bmatrix} P \right)}{\mathrm{Tr}\left( P^{\mathrm{T}} \begin{bmatrix} \gamma S_{\mathrm{b}}^{\mathrm{s}} & \mathbf{0} \\ \mathbf{0} & \gamma_{\mathrm{b}}^{\mathrm{t}} + \mu S_{\mathrm{h}}^{\mathrm{t}} \end{bmatrix} P \right)}
\tag{8}
$$

We can reformulate Eq. (8) as:

$$
\max_{P} \mathrm{Tr}\left( P^{\mathrm{T}} \begin{bmatrix} \gamma S_{\mathrm{b}}^{\mathrm{s}} & \mathbf{0} \\ \mathbf{0} & \gamma S_{\mathrm{b}}^{\mathrm{t}} + \mu S_{\mathrm{h}}^{\mathrm{t}} \end{bmatrix} P \right)
$$
$$
\text{s.t. } \mathrm{Tr}\left( P^{\mathrm{T}} \begin{bmatrix} M_{\mathrm{ss}} + \gamma S_{\mathrm{w}}^{\mathrm{s}} & M_{\mathrm{st}} \\ M_{\mathrm{ts}} & M_{\mathrm{tt}} + \gamma S_{\mathrm{w}}^{\mathrm{t}} + \mu I \end{bmatrix} P \right) = 1
\tag{9}
$$

3

According to the constrained optimization theory, we introduce a Lagrange multiplier $\Phi$, and get the Lagrange function for Eq. (19) as follows:

$$
\begin{aligned}
\mathcal{L} = {} & \mathrm{Tr}\left( P^{\mathrm{T}} \begin{bmatrix} \gamma S_{\mathrm{b}}^{\mathrm{s}} & 0 \\ 0 & \gamma S_{\mathrm{b}}^{\mathrm{u}} + \mu S_{\mathrm{h}}^{\mathrm{u}} \end{bmatrix} P \right) \\
& - \mathrm{Tr}\left( \left( P^{\mathrm{T}} \begin{bmatrix} M_{\mathrm{ss}} + \gamma S_{\mathrm{w}}^{\mathrm{s}} & M_{\mathrm{su}} \\ M_{\mathrm{s}}^{\mathrm{u}} & M_{\mathrm{uu}} + \gamma S_{\mathrm{w}}^{\mathrm{u}} + \mu I \end{bmatrix} P - I \right) \Phi \right)
\end{aligned}
\tag{10}
$$

where $\Phi = \mathrm{diag}\left(\phi_1, \cdots, \phi_d\right)$ and $\left(\phi_1, \cdots, \phi_d\right)$ are the $d$ largest eigenvalues of the following eigendecomposition problem:

$$
\begin{bmatrix} \gamma S_{\mathrm{b}}^{\mathrm{s}} & \mathbf{0} \\ \mathbf{0} & \gamma S_{\mathrm{b}}^{\mathrm{u}} + \mu S_{\mathrm{h}}^{\mathrm{u}} \end{bmatrix} P = \begin{bmatrix} M_{\mathrm{ss}} + \gamma S_{\mathrm{w}}^{\mathrm{s}} & M_{\mathrm{su}} \\ M_{\mathrm{s}}^{\mathrm{u}} & M_{\mathrm{uu}} + \gamma S_{\mathrm{w}}^{\mathrm{u}} + \mu I \end{bmatrix} P \Phi
\tag{11}
$$

As a result, $P$ consists of the corresponding $d$ largest eigenvectors of Eq. (11). At last, the subspaces spanned by $A$ and $B$ can be obtained easily once the transformation matrix $P$ is obtained.

2) Optimizing sample weights $\alpha$ and $\beta$ : Regarding $A$ and $B$ as constants, Since $\mathrm{Tr}\left(AB\right) = \mathrm{Tr}\left(A^T B^T\right)$ and $\mathrm{Tr}\left(constant\right) = constant$. The Eq(3) can be formulated as follows:

$$
\begin{aligned}
& \min_{\alpha^u, \beta} \sum_{u=1}^{N_{\mathrm{s}}} \left( \tfrac{1}{2} \alpha_{\mathrm{u}}^{\mathrm{T}} K_{\mathrm{ss}}^{\mathrm{u}} \alpha_{\mathrm{u}} - \tfrac{1}{2} \alpha_{\mathrm{u}}^{\mathrm{T}} K_{\mathrm{st}}^{\mathrm{u}} \beta - \tfrac{1}{2} \beta^{\mathrm{T}} K_{\mathrm{ts}}^{\mathrm{u}} \alpha_{\mathrm{u}} + \tfrac{1}{2} \beta^{\mathrm{T}} K_{\mathrm{tt}} \beta \right) \\
& \text{s.t.} \quad \{\alpha_{ui}^c, \beta_i^c\} \in [0, 1], \frac{\alpha_{\mathrm{u}}^{c\,\mathrm{T}} \mathbf{1}_{n_{\mathrm{s}}^{c\mathrm{u}}}}{n_{\mathrm{s}}^{u c}} = \delta_{\mathrm{s}}^{\mathrm{u}}, \frac{\beta^{c\,\mathrm{T}} \mathbf{1}_{n_{\mathrm{t}}^c}}{n_{\mathrm{t}}^c} = \delta_{\mathrm{t}}
\end{aligned}
\tag{12}
$$

where $(K_{\mathrm{ss}}^{\mathrm{u}})_{i,j}$ in $K_{\mathrm{ss}}^{\mathrm{u}} \in R^{n_{us} \times n_{us}}$ is the coefficient associated with $\left(A^{\mathrm{T}} x_{\mathrm{s}}^{\mathrm{u}i}\right)^{\mathrm{T}} A^{\mathrm{T}} x_{\mathrm{s}}^{\mathrm{u}i}$, $(K_{\mathrm{st}}^{\mathrm{u}})_{i,j}$ in $K_{\mathrm{st}}^{\mathrm{u}} \in R^{n_{\mathrm{s}}^{\mathrm{u}} \times n_{\mathrm{t}}}$ is the coefficient associated with $\left(A^{\mathrm{T}} x_{\mathrm{s}}^{\mathrm{u}i}\right)^{\mathrm{T}} B^{\mathrm{T}} x_{\mathrm{t}}^{j}$, and $(K_{\mathrm{tt}})_{i,j}$ in $K_{\mathrm{tt}} \in R^{n_{\mathrm{t}} \times n_{\mathrm{t}}}$ is the coefficient associated with $\left(B^{\mathrm{T}} x_{\mathrm{t}}^{i}\right)^{\mathrm{T}} B^{\mathrm{T}} x_{\mathrm{t}}^{j}$.

With the above formulation, we can apply Quadratic Programming (QP) solvers to optimize the equivalent problem:

$$
\min_{z_i \in [0,1], V^{\mathrm{T}} \cdot Z = G} \frac{1}{2} Z^{\mathrm{T}} Q Z
\tag{13}
$$

$$
Z = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_u \\ \beta \end{pmatrix}, Q = \begin{bmatrix} K_{\mathrm{ss}}^1 & 0 & \cdots & 0 & -K_{\mathrm{st}}^1 \\ 0 & K_{\mathrm{ss}}^2 & \cdots & 0 & -K_{\mathrm{st}}^2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & K_{\mathrm{ss}}^{\mathrm{u}} & -K_{\mathrm{st}}^{\mathrm{u}} \\ -K_{\mathrm{ts}}^1 & -K_{\mathrm{ts}}^2 & \cdots & -K_{\mathrm{ts}}^{\mathrm{u}} & \sum_{u=1}^{N_s} K_{\mathrm{tt}} \end{bmatrix}
$$

$$
G \in R^{(Ns+1)C \times 1} \text{ with } (G)_c = \begin{cases} \delta_{\mathrm{s}}^{\mathrm{u}} n_{\mathrm{s}}^{uc} & \text{if } (u-1)\mathrm{C} \le c \le \mathrm{uC} \\ \delta_t n_{\mathrm{t}}^c & \text{if } c > \mathrm{Ns} \times \mathrm{C} \end{cases}
$$

$$
V = \begin{bmatrix} V_{1\mathrm{s}} & \mathbf{0}_{n_{2\mathrm{s}} \times C} & \cdots & \mathbf{0}_{n_{1\mathrm{s}} \times C} & \mathbf{0}_{n_{1\mathrm{s}} \times C} \\ \mathbf{0}_{n_{2\mathrm{s}} \times C} & V_{2\mathrm{s}} & \cdots & \mathbf{0}_{n_{2\mathrm{s}} \times C} & \mathbf{0}_{n_{2\mathrm{s}} \times C} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0}_{n_{\mathrm{s}}^{\mathrm{u}} \times C} & \mathbf{0}_{n_{\mathrm{s}}^{\mathrm{u}} \times C} & \cdots & V_{\mathrm{s}}^{\mathrm{u}} & \mathbf{0}_{n_{\mathrm{s}}^{\mathrm{u}} \times C} \\ \mathbf{0}_{n_{\mathrm{t}} \times C} & \mathbf{0}_{n_{\mathrm{t}} \times C} & \cdots & \mathbf{0}_{n_{\mathrm{t}} \times C} & V_{\mathrm{t}} \end{bmatrix} \in R^{(Ns+1)C \times (n_{\mathrm{s}} + n_{\mathrm{t}})} \text{ with}
$$

$$
(V_{\mathrm{s}}^{\mathrm{u}})_{ij} = \begin{cases} 1 & \text{if } x_{\mathrm{s}}^{\mathrm{u}i} \in \text{ class } j \\ 0 & \text{otherwise} \end{cases}, \text{ and}
$$

$$
(V_{\mathrm{t}})_{ij} = \begin{cases} 1 & \text{if } x_{\mathrm{t}}^i \text{ predicted as class } j \\ 0 & \text{otherwise} \end{cases}
\tag{14}
$$

4

# 3   Optimization Procedure

---

**Procedure 1** Multi-site Domain Adaptation via Landmark Selection

---

**Input:** Source and target domain data: $X_{\mathrm{s}}$ and $X_{\mathrm{t}}$; labels for source domain data and pseudo-labels for target domain data: $y_{\mathrm{s}}$ and $\hat{y}_t$ ;Parameters: $\delta_{\mathrm{us}}, \delta_{\mathrm{t}}, d, \mu, \gamma$

**Output:** optimal Predicted labels $y_{\mathrm{u}}$ for target domain unlabeled data

0: Initialize pseudo labels of target domain unlabeled data $\hat{y}_t$ using certain base classifiers with $X_{\mathrm{t}}$ Compute $S_{\mathrm{h}}^{\mathrm{t}}, M_{\mathrm{ss}}, M_{\mathrm{uu}}, M_{\mathrm{st}}, M_{\mathrm{ts}}, S_{\mathrm{h}}^{\mathrm{s}}, S_{\mathrm{W}}^{\mathrm{s}}, S_{\mathrm{h}}^{\mathrm{t}}, S_{\mathrm{W}}^{\mathrm{u}}$;

1: **while** not converge **do**

1:    Solve the generalized eigen-decomposition problem in (11) and select $d$ corresponding eigenvectors of $d$ largest eigenvalues as the transformation $P$, and obtain transformation $A$ and $B$;

1:    Map the original data to respective subspace to get the embeddings: $Z_{\mathrm{s}} = A^{\mathrm{T}} X_{\mathrm{s}}, Z_{\mathrm{t}} = B^{\mathrm{T}} X_{\mathrm{t}}$

1:    Use base classifiers on $Z_{\mathrm{s}}, Z_{\mathrm{t}}, y_{\mathrm{s}}$ to update pseudo labels in target domain $\hat{y}_t$

1:    Update landmark weights $\alpha, \beta$

1:    Update $M_{\mathrm{ss}}, M_{\mathrm{tt}}, M_{\mathrm{st}}, M_{\mathrm{ts}}, S_{\mathrm{b}}^{\mathrm{t}}, S_{\mathrm{w}}^{\mathrm{t}}$

2: **end while**=0

---