

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#) See Section 1.
 - (b) Did you describe the limitations of your work? [\[Yes\]](#) We clearly state the assumptions we use in Section 2.
 - (c) Did you discuss any potential negative societal impacts of your work? [\[No\]](#) We have done a theoretic work that has few societal impacts.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#) See Section 2.
 - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) See supplementary material.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#) See supplemental material.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) See Section 4.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[No\]](#) We think the error bars are not related to the core result of our experiments.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[No\]](#) We use few computation resources in our work.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#) See Section 4.
 - (b) Did you mention the license of the assets? [\[Yes\]](#) See Section 4.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [\[No\]](#) We don’t use new assets.
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [\[Yes\]](#) See Section 4.
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[N/A\]](#)
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#)
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#)
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#)

Appendix

A Loss Function Leading to (2)

To obtain (2), first note that for each $i = 0, 1$, we can define the loss of task i as

$$L_i = \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \hat{P}_{XY}^{(i)}(x, y) \log \frac{1}{Q_{XY}(x, y)}, \quad (19)$$

where Q_{XY} represents the distribution model and $\hat{P}_{XY}^{(i)}$ are the empirical distributions as defined in (1). As a result, training the linear combination of L_0 and L_1 can lead to the convex combination model.

$$\begin{aligned} & \arg \min_{Q_{XY} \in \mathcal{P}} \alpha_0 L_0 + \alpha_1 L_1 \\ &= \arg \min_{Q_{XY} \in \mathcal{P}} \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \left(\alpha_0 \hat{P}_{XY}^{(0)}(x, y) + \alpha_1 \hat{P}_{XY}^{(1)}(x, y) \right) \log \frac{1}{Q_{XY}(x, y)} \\ &= \arg \min_{Q_{XY} \in \mathcal{P}} D \left(\alpha_0 \hat{P}_{XY}^{(0)} + \alpha_1 \hat{P}_{XY}^{(1)} \parallel Q_{XY} \right) \\ &= \alpha_0 \hat{P}_{XY}^{(0)} + \alpha_1 \hat{P}_{XY}^{(1)}. \end{aligned} \quad (20)$$

B Proof of Theorem 2

To compute the testing loss, we first introduce the following lemma.

Lemma 6. Suppose that random vector (X_1, X_2, \dots, X_m) follows the multinomial distribution with the corresponding event probabilities (p_1, p_2, \dots, p_m) and n independent trials, then the variance of random variable X_i is

$$\text{var}(X_i) = np_i(1 - p_i), \quad (21)$$

and the covariance of X_i and X_j is

$$\text{cov}(X_i, X_j) = -np_i p_j. \quad (22)$$

From Lemma 6, for $i = 0, \dots, k$, we have

$$\mathbb{E} \left[\left(\hat{P}_{XY}^{(i)}(x, y) - P_{XY}^{(i)}(x, y) \right)^2 \right] = \frac{1}{n_i} P_{XY}^{(i)}(x, y) (1 - P_{XY}^{(i)}(x, y)). \quad (23)$$

In turn, the testing loss as defined in (3) is

$$\begin{aligned} L_{\text{test}}^{(\alpha_0, \alpha_1)} &= \mathbb{E} \left[\chi^2(P_{XY}^{(0)}, \alpha_0 \hat{P}_{XY}^{(0)} + \alpha_1 \hat{P}_{XY}^{(1)}) \right] \\ &= \mathbb{E} \left[\sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{1}{P_{XY}^{(0)}(x, y)} \left(P_{XY}^{(0)}(x, y) - \alpha_0 \hat{P}_{XY}^{(0)}(x, y) - \alpha_1 \hat{P}_{XY}^{(1)}(x, y) \right)^2 \right] \\ &= \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{\left(P_{XY}^{(0)}(x, y) - \alpha_0 P_{XY}^{(0)}(x, y) - \alpha_1 P_{XY}^{(1)}(x, y) \right)^2}{P_{XY}^{(0)}(x, y)} \\ &\quad + \alpha_0^2 \mathbb{E} \left[\sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{\left(\hat{P}_{XY}^{(0)}(x, y) - P_{XY}^{(0)}(x, y) \right)^2}{P_{XY}^{(0)}(x, y)} \right] \\ &\quad + \alpha_1^2 \mathbb{E} \left[\sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{\left(\hat{P}_{XY}^{(1)}(x, y) - P_{XY}^{(1)}(x, y) \right)^2}{P_{XY}^{(0)}(x, y)} \right] \end{aligned} \quad (24)$$

$$= \alpha_1^2 \chi^2 \left(P_{XY}^{(0)}, P_{XY}^{(1)} \right) + \alpha_0^2 \frac{1}{n_0} V^{(0)} + \alpha_1^2 \frac{1}{n_1} V^{(1)}, \quad (25)$$

where to obtain (24) we have used the facts that $\hat{P}_{XY}^{(0)}$ and $\hat{P}_{XY}^{(1)}$ are independent, $\mathbb{E}[\hat{P}_{XY}^{(0)}] = P_{XY}^{(0)}$, and $\mathbb{E}[\hat{P}_{XY}^{(1)}] = P_{XY}^{(1)}$, and where to obtain (25) we have used (23) and the fact that $n_0 \cdot \hat{P}_{XY}^{(0)}$ and $n_1 \cdot \hat{P}_{XY}^{(1)}$ follows the multinomial distribution.

Similarly, we can obtain Theorem 3.

C Proof of Theorem 4

C.1 Expression of \hat{g}_i as defined in (13)

From

$$\begin{aligned} & \chi_{R_{XY}}^2 \left(\hat{P}_{XY}^{(i)}, P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g})} \right) \\ &= \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{1}{P_X^{(0)}(x) P_Y^{(0)}(y)} \cdot \left(\hat{P}_{XY}^{(i)}(x, y) - P_X^{(0)}(x) P_Y^{(0)}(y) - P_X^{(0)}(x) P_Y^{(0)}(y) \mathbf{f}^T(x) \mathbf{g}(y) \right)^2, \end{aligned}$$

we have, for all $y' \in \mathcal{Y}$,

$$\frac{\partial \chi_{R_{XY}}^2 \left(\hat{P}_{XY}^{(i)}, P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g})} \right)}{\partial \mathbf{g}(y')} \quad (26)$$

$$\begin{aligned} &= -2 \sum_{x \in \mathcal{X}} \left[\hat{P}_{XY}^{(i)}(x, y') - P_X^{(0)}(x) P_Y^{(0)}(y') \right] \mathbf{f}(x) + 2 \sum_{x \in \mathcal{X}} P_X^{(0)}(x) P_Y^{(0)}(y') (\mathbf{f}^T(x) \mathbf{g}(y')) \mathbf{f}(x) \\ &= -2 \sum_{x \in \mathcal{X}} \hat{P}_{XY}^{(i)}(x, y') \mathbf{f}(x) + 2 P_Y^{(0)}(y') \mathbf{\Lambda}_{\mathbf{f}} \mathbf{g}(y') \end{aligned} \quad (27)$$

where to obtain the last equality, we have used the assumption that $\mathbb{E}_{P_X^{(0)}}[\mathbf{f}(X)] = \mathbf{0}$ and the notation that $\mathbf{\Lambda}_{\mathbf{f}} \triangleq \mathbb{E}_{P_X^{(0)}}[\mathbf{f}(X) \mathbf{f}^T(X)]$.

Set the gradient (27) to zero, and we obtain

$$\hat{\mathbf{g}}_i(y) = \frac{1}{P_Y^{(0)}(y)} \mathbf{\Lambda}_{\mathbf{f}}^{-1} \left(\sum_{x \in \mathcal{X}} \hat{P}_{XY}^{(i)}(x, y) \mathbf{f}(x) \right). \quad (28)$$

C.2 Proof of (16)

We first express the testing error (15) as

$$\begin{aligned} L_{\text{test}}(\boldsymbol{\alpha}) &= \mathbb{E} \left[\chi_{R_{XY}}^2 \left(P_X^{(0)}, P_X^{(0)} Q_{Y|X}^{(\boldsymbol{\alpha})} \right) \right] \\ &= \chi_{R_{XY}}^2 \left(P_X^{(0)}, \sum_{i=0}^k \alpha_i P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_i)} \right) + \underbrace{\sum_{i=0}^k \alpha_i^2 \mathbb{E} \left[\chi_{R_{XY}}^2 \left(P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_i)}, P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \hat{\mathbf{g}}_i)} \right) \right]}_{\frac{1}{n_i} \hat{V}^{(i)}}, \end{aligned} \quad (29)$$

where to obtain the last equality we have used the fact that the empirical distributions $\hat{P}_{XY}^{(i)}$ ($i = 0, \dots, k$) are independent and

$$\mathbb{E} \left[\tilde{P}_{Y|X}^{(\mathbf{f}, \hat{\mathbf{g}}_i)} \right] = \mathbb{E} \left[\tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_i)} \right].$$

Next, for the terms in (29), we have

$$\begin{aligned} & \chi_{R_{XY}}^2 \left(P_X^{(0)}, \sum_{i=0}^k \alpha_i P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_i)} \right) \\ &= \chi_{R_{XY}}^2 \left(P_X^{(0)}, P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_0)} \right) + \chi_{R_{XY}}^2 \left(P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_0)}, \sum_{i=0}^k \alpha_i P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_i)} \right), \end{aligned} \quad (30)$$

which comes from the fact that

$$\begin{aligned}
& \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{\left(P_{XY}^{(0)}(x, y) - [P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_0)}](x, y) \right) [P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g})}](x, y)}{P_X^{(0)}(x) P_Y^{(0)}(y)} \\
&= \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \left(P_{XY}^{(0)}(x, y) - [P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_0)}](x, y) \right) \mathbf{f}^T(x) \mathbf{g}(y) \\
&= \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} P_{XY}^{(0)}(x, y) \mathbf{f}^T(x) \mathbf{g}(y) \\
&\quad - \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} P_X^{(0)}(x) P_Y^{(0)}(y) \mathbf{f}^T(x) \frac{1}{P_Y^{(0)}(y)} \mathbf{\Lambda}_{\mathbf{f}}^{-1} \left(\sum_{x \in \mathcal{X}} P_{XY}^{(0)}(x, y) \mathbf{f}(x) \right) \mathbf{f}^T(x) \mathbf{g}(y) \\
&= \mathbb{E}_{P_{XY}^{(0)}} [\mathbf{f}^T(X) \mathbf{g}(Y)] - \mathbb{E}_{P_{XY}^{(0)}} [\mathbf{f}^T(X) \mathbf{g}(Y)] \\
&= 0.
\end{aligned} \tag{31}$$

Moreover, we have

$$\begin{aligned}
& \tilde{V}^{(i)} \\
&= n_i \mathbb{E} \left[\chi_{R_{XY}}^2 \left(P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_i)}, P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \hat{\mathbf{g}}_i)} \right) \right] \\
&= n_i \mathbb{E} \left[\sum_{x \in \mathcal{X}, y \in \mathcal{Y}} P_X^{(0)}(x) P_Y^{(0)}(y) \left(\mathbf{f}^T(x) \frac{1}{P_Y^{(0)}(y)} \mathbf{\Lambda}_{\mathbf{f}}^{-1} \left(\sum_{x \in \mathcal{X}} \left(P_{XY}^{(i)}(x, y) - \hat{P}_{XY}^{(i)}(x, y) \right) \mathbf{f}(x) \right) \right)^2 \right] \\
&= n_i \sum_{y=1}^{|\mathcal{Y}|} \frac{1}{P_Y^{(0)}(y)} \text{tr} \left(\mathbf{\Lambda}_{\mathbf{f}}^{-1} \right. \\
&\quad \cdot \mathbb{E} \left[\left(\sum_{x \in \mathcal{X}} \left(P_{XY}^{(i)}(x, y) - \hat{P}_{XY}^{(i)}(x, y) \right) \mathbf{f}(x) \right) \left(\sum_{x \in \mathcal{X}} \left(P_{XY}^{(i)}(x, y) - \hat{P}_{XY}^{(i)}(x, y) \right) \mathbf{f}(x) \right)^T \right] \Bigg) \\
&= \sum_{y=1}^{|\mathcal{Y}|} \frac{P_Y^{(i)}(y)}{P_Y^{(0)}(y)} \text{tr} \left(\mathbf{\Lambda}_{\mathbf{f}}^{-1} \mathbb{E}_{P_{X|Y=y}^{(i)}} [\mathbf{f}(X) \mathbf{f}^T(X)] \right) - \sum_{y=1}^{|\mathcal{Y}|} \frac{[P_Y^{(i)}(y)]^2}{P_Y^{(0)}(y)} \left\| \mathbf{\Lambda}_{\mathbf{f}}^{-\frac{1}{2}} \mathbb{E}_{P_{X|Y=y}^{(i)}} [\mathbf{f}(X)] \right\|^2.
\end{aligned} \tag{32}$$

Using (29), (30) and (32), we obtain (16) as desired.

D Proof of Proposition 5

Since \mathbf{g}^* as defined in (14) satisfies $\mathbf{g}^* = \sum_{i=0}^k \alpha_i \hat{\mathbf{g}}_i$, the weight \mathbf{g}^* in the convex combination model is the solution of the following optimization problem

$$\mathbf{g}^* = \arg \min_{\mathbf{g}} \chi_{R_{XY}}^2 \left(\sum_{i=0}^k \alpha_i \hat{P}_{XY}^{(i)}, P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g})} \right). \tag{33}$$

Then, we prove that the following two optimization problems have the same solution \mathbf{g}^*

$$\begin{aligned}
\mathbf{g}^* &= \arg \min_{\mathbf{g}} \chi_{R_{XY}}^2 \left(\sum_{i=0}^k \alpha_i \hat{P}_{XY}^{(i)}, P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g})} \right) \\
&= \arg \min_{\mathbf{g}} \sum_{i=0}^k \alpha_i \chi_{R_{XY}}^2 \left(\hat{P}_{XY}^{(i)}, P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g})} \right).
\end{aligned} \tag{34}$$

To obtain (34), we use the fact that the difference

$$\begin{aligned} & \chi_{R_{XY}}^2 \left(\sum_{i=0}^k \alpha_i \hat{P}_{XY}^{(i)}, P_X^{(0)} \hat{P}_{Y|X}^{(\mathbf{f}, \mathbf{g})} \right) - \sum_{i=0}^k \alpha_i \chi_{R_{XY}}^2 \left(\hat{P}_{XY}^{(i)}, P_X^{(0)} \hat{P}_{Y|X}^{(\mathbf{f}, \mathbf{g})} \right) \\ &= \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{\left(\sum_{i=0}^k \alpha_i \hat{P}_{XY}^{(i)}(x, y) \right)^2}{P_X^{(0)}(x) P_Y^{(0)}(y)} - \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{\sum_{i=0}^k \alpha_i \left[\hat{P}_{XY}^{(i)}(x, y) \right]^2}{P_X^{(0)}(x) P_Y^{(0)}(y)} \end{aligned} \quad (35)$$

is irrelevant to \mathbf{g} .

Using (34), we obtain Proposition 5 as desired.

E Details for experiments

E.1 Training Loss

To compute the training loss $L(\alpha, \mathbf{f}, \mathbf{g})$ as defined in (17), we first introduce the following lemma.

Lemma 7 ([48], Proposition 2). *Let $(\mathbf{f}^*, \mathbf{g}^*)$ be the features that minimize the χ^2 -distance loss $\chi_{R_{XY}}^2(\hat{P}_{XY}^{(0)}, \hat{P}_X^{(0)} \hat{P}_{Y|X}^{(\mathbf{f}, \mathbf{g})})$, where $\hat{P}_{Y|X}^{(\mathbf{f}, \mathbf{g})}(x, y) \triangleq \hat{P}_Y(y)(1 + \mathbf{f}^T(x)\mathbf{g}(y))$, for all x, y , and the reference distribution being $\hat{P}_X^{(0)} \hat{P}_Y^{(0)}$. Then, we have*

$$\mathbb{E}_{\hat{P}_X^{(0)}}[\mathbf{f}^*(X)] = \mathbb{E}_{\hat{P}_Y^{(0)}}[\mathbf{g}^*(Y)] = \mathbf{0}, \quad (36)$$

and $(\mathbf{f}^*, \mathbf{g}^*)$ are also the optimal features that maximize the H-score of target samples:

$$H^{(0)}(\mathbf{f}, \mathbf{g}) \triangleq \mathbb{E}_{\hat{P}_{XY}^{(0)}}[\tilde{\mathbf{f}}^T(X)\tilde{\mathbf{g}}(Y)] - \frac{1}{2} \text{tr}(\hat{\Lambda}_{\mathbf{f}} \hat{\Lambda}_{\mathbf{g}}), \quad (37)$$

where $\tilde{\mathbf{f}}(X) \triangleq \mathbf{f}(X) - \mathbb{E}_{\hat{P}_X^{(0)}}[\mathbf{f}(X)]$, $\tilde{\mathbf{g}}(Y) \triangleq \mathbf{g}(Y) - \mathbb{E}_{\hat{P}_Y^{(0)}}[\mathbf{g}(Y)]$, $\hat{\Lambda}_{\mathbf{f}}$ and $\hat{\Lambda}_{\mathbf{g}}$ are the covariance matrices of features on target samples, defined as:

$$\hat{\Lambda}_{\mathbf{f}} \triangleq \mathbb{E}_{\hat{P}_X^{(0)}}[\tilde{\mathbf{f}}(X)\tilde{\mathbf{f}}^T(X)], \quad (38)$$

$$\hat{\Lambda}_{\mathbf{g}} \triangleq \mathbb{E}_{\hat{P}_Y^{(0)}}[\tilde{\mathbf{g}}(Y)\tilde{\mathbf{g}}^T(Y)]. \quad (39)$$

Similarly we can define the H-score for sources. For the task $i = 1, \dots, k$,

$$H^{(i)}(\mathbf{f}, \mathbf{g}) \triangleq \mathbb{E}_{\hat{P}_{XY}^{(i)}}[\tilde{\mathbf{f}}^T(X)\tilde{\mathbf{g}}(Y)] - \frac{1}{2} \text{tr}(\hat{\Lambda}_{\mathbf{f}} \hat{\Lambda}_{\mathbf{g}}). \quad (40)$$

Then, line 4 in Algorithm 1 can be implemented by

$$(\mathbf{f}^*, \mathbf{g}^*) \leftarrow \arg \max_{\mathbf{f}, \mathbf{g}} \sum_{i=0}^k \alpha_i H^{(i)}(\mathbf{f}, \mathbf{g}). \quad (41)$$

E.2 Computation of Testing Loss

In computing the testing loss, after obtaining \mathbf{f}^* from (41), we use the normalization $\tilde{\mathbf{f}}^*(X) \triangleq \mathbf{f}^*(X) - \mathbb{E}_{\hat{P}_X^{(0)}}[\mathbf{f}^*(X)]$ to subtract the sample mean and obtain zero-mean features. Then, the covariance matrix $\hat{\Lambda}_{\mathbf{f}^*}$ of \mathbf{f}^* is computed as

$$\hat{\Lambda}_{\mathbf{f}^*} \triangleq \mathbb{E}_{\hat{P}_X^{(0)}}[\tilde{\mathbf{f}}^*(X)\tilde{\mathbf{f}}^{*T}(X)]. \quad (42)$$

Moreover, for $i = 0, \dots, k$, the $\tilde{V}^{(i)}$ [cf. (32)] in (16) can be estimated as

$$\begin{aligned} & \tilde{V}^{(i)} \\ &= \sum_{y=1}^{|\mathcal{Y}|} \frac{\hat{P}_Y^{(i)}(y)}{\hat{P}_Y^{(0)}(y)} \text{tr} \left(\hat{\Lambda}_{\mathbf{f}^*}^{-1} \mathbb{E}_{\hat{P}_{X|Y=y}^{(i)}}[\tilde{\mathbf{f}}^*(X)\tilde{\mathbf{f}}^{*T}(X)] \right) - \sum_{y=1}^{|\mathcal{Y}|} \frac{[\hat{P}_Y^{(i)}(y)]^2}{\hat{P}_Y^{(0)}(y)} \left\| \hat{\Lambda}_{\mathbf{f}^*}^{-\frac{1}{2}} \mathbb{E}_{\hat{P}_{X|Y=y}^{(i)}}[\tilde{\mathbf{f}}^*(X)] \right\|^2. \end{aligned} \quad (43)$$

The bias term $\chi_{R_{XY}}^2 \left(P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_0)}, \sum_{i=0}^k \alpha_i P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_i)} \right)$ in (16) can be estimated as

$$\chi_{R_{XY}}^2 \left(P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_0)}, \sum_{i=0}^k \alpha_i P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_i)} \right) = \sum_{i=1}^k \sum_{j=1}^k \alpha_i \alpha_j D_{ij}, \quad (44)$$

where

$$D_{ij} \triangleq \text{tr} \left(\hat{\mathbf{\Lambda}}_{\mathbf{f}^*}^{-1} \left(\sum_{y \in \mathcal{Y}} \frac{(h(0, y) - h(i, y))(h(0, y) - h(j, y))^T}{\hat{P}_Y^{(0)}(y)} \right) \right) \quad (45)$$

and

$$h(l, y) \triangleq \hat{P}_Y^{(l)}(y) \mathbb{E}_{\hat{P}_{X|Y=y}^{(l)}} [\tilde{\mathbf{f}}^*(X)], \quad l = 0, \dots, k. \quad (46)$$

Finally, the testing loss (16) can be expressed as

$$L_{\text{test}}^{(\alpha)} = \sum_{i=1}^k \sum_{j=1}^k \alpha_i \alpha_j D_{ij} + \sum_{i=0}^k \frac{\alpha_i^2}{n_i} \tilde{V}^{(i)}. \quad (47)$$

Then, α^* can be computed by solving a non-negative quadratic programming problem.

E.3 Details of Implementations

E.3.1 Multi-Source Transfer Learning

In this experiment, the corresponding labels for the 5 binary classification tasks are as follows: task 0 (airplane and automobile), task 1 (bird and cat), task 2 (deer and dog), task 3 (frog and horse), and task 4 (ship and truck). After training the loss (40) for each source task, if the test accuracy on target samples is less than 50%, we would flip the binary label for this source. Accordingly, in task 3, frog matches automobile and horse matches airplane. In other tasks, the labels of the source match the target in alphabetical order.

Moreover, we normalize all the images for 3 channels under the mean (0.485, 0.456, 0.406) and standard deviation (0.229, 0.224, 0.225). When the sample size of the target training set is 6, to make the training process stable, we use data augmentation to generate 50 samples by random horizontal flips and random crops.

The feature \mathbf{f} is generated of 10 dimensions by GoogLeNet, followed by a fully connection layer (1024 \rightarrow 32) with ReLU activation, and then a fully connection layer (32 \rightarrow 10). Throughout the training process, we use the Adam optimizer with a learning rate of 0.001 and the batch size for each source task is 50. We train the networks in 20 epochs and before each epoch we reshuffle the training samples.

E.3.2 Few-shot Transfer Learning Tasks on Office-31

In this experiment, the feature \mathbf{f} is generated by the pretrained and fixed VGG-16 network, followed by a fully connection layer (4096 \rightarrow 1024) with ReLU activation, and a fully connection layer (1024 \rightarrow 64). Throughout the training process, we use the Adam optimizer with learning rate of 0.0002 and in 100 epochs.

E.3.3 Few-shot Transfer Learning Tasks on Office-Caltech

In this experiment, the feature \mathbf{f} is generated by a fully connection layer (4096 \rightarrow 1024) with ReLU activation followed by a fully connection layer (1024 \rightarrow 10). The inputs of \mathbf{f} are the pretrained DeCAF features. Throughout the training process, we use the Adam optimizer with learning rate of 0.01 and in 100 epochs.

E.4 Instruction for codes

We provide code examples in “code.zip”. In the folder “./code/data”, we provide the features we used. Folder “o31_feature” contains the features of Office-31 dataset, and folder “oc_feature”

contains the features of Office-Caltech dataset. “cifar10_alpha.py” is an example for CIFAR-10 dataset for the case of 6 target samples. “o31_atod_renew.py” is an example for Office-31 dataset for the case of task $A \rightarrow D$, with the details of computing renewed α . “oc_atoc.py” is an example for Office-Caltech dataset for the case of task $A \rightarrow C$.