# Appendix

# A. Proofs

#### A.1. Proof for parameterization gap

**Proof A.1** Recall the feasibility assumption 3 ensures the strong duality such that the primal and dual optimal objectives are equal, which means for  $\forall \ \tilde{\lambda} \in \mathbb{R}^+$ ,  $\tilde{f}_s$ ,  $\tilde{f}_v \in \mathcal{H}$ , the following saddle point condition holds

$$\mathcal{R}(f_s^*, f_v^*, \tilde{\lambda}) \le \max_{\lambda} \min_{f_s, f_v} \mathcal{R}(f_s, f_v, \lambda) = D^* = P^*$$
  
$$= \min_{f_s, f_v} \max_{\lambda} \mathcal{R}(f_s, f_v, \lambda) \le \mathcal{R}(\tilde{f}_s, \tilde{f}_v, \lambda^*)$$
(14)

*Recall the definition of*  $\mathcal{D}^*_{\epsilon}$  *and by the inclusion relation*  $\mathcal{H}_{\theta} \subseteq \mathcal{F}$ *, we can derive the lower bound as* 

$$\mathcal{D}_{\epsilon}^{*}(\gamma) \triangleq \max_{\lambda} \min_{\theta, \phi \in \mathcal{H}} \mathcal{L}(\theta) + \lambda \mathcal{L}_{\mathrm{con}}(\theta, \phi)$$

$$\geq \min_{\theta, \phi \in \mathcal{H}} \mathcal{L}(\theta) + \lambda \mathcal{L}_{\mathrm{con}}(\theta, \phi), \forall \lambda \in \mathbb{R}^{+}$$

$$\geq \min_{f_{s}, f_{v} \in \mathcal{F}} \mathcal{L}(f_{s}) + \lambda \mathcal{L}_{\mathrm{con}}(f_{s}, f_{v}) = \mathcal{P}^{*}$$
(15)

For the upper bound, we add and subtract  $\mathcal{R} \equiv \mathcal{R}(f_s, f_v, \lambda) = \min_{f_s \in \mathcal{F}} \mathcal{L}(f_s) + \min_{f_s, f_v \in \mathcal{F}} \lambda \mathcal{L}_{con}(f_s, f_v)$  from  $\mathcal{D}_{\epsilon}^{\star}(\gamma)$  to get

$$\mathcal{D}_{\epsilon}^{\star}(\gamma) = \max_{\lambda} \min_{\theta, f_s, f_v} \mathcal{L}(\theta) + \lambda \mathcal{L}_{con}(\theta, \phi) + \mathcal{R} - \mathcal{R}$$
$$= \max_{\lambda} \min_{\theta, f_s, f_v} \mathcal{R} + [\mathcal{L}(\theta) - \mathcal{L}(f_s)] + \lambda [\mathcal{L}_{con}(\theta, \phi) - \mathcal{L}_{con}(f_s, f_v)]$$
(16)

With a slight abuse of notation,  $\mathbb{E}$  is short for  $\mathbb{E}_{(\mathbf{x},y)\sim\mathbb{P}(X,Y),\mathbf{\tilde{x}}\sim\mathbb{P}(X)}$ . Then we consider the combination of the second and third term as

$$\begin{bmatrix} 1, \lambda \end{bmatrix} \begin{bmatrix} \mathcal{L}(\theta) - \mathcal{L}(f_{s}) \\ \mathcal{L}_{con}(\theta, \phi) - \mathcal{L}_{con}(f_{s}, f_{v}) \end{bmatrix} \\ \leq (1 + \|\lambda\|_{1}) \left\| \begin{bmatrix} \mathcal{L}(\theta) - \mathcal{L}(f_{s}) \\ \mathcal{L}_{con}(\theta, \phi) - \mathcal{L}_{con}(f_{s}, f_{v}) \end{bmatrix} \right\|_{\infty} \\ = (1 + |\lambda|) \max \left\{ \mathcal{L}(\theta) - \mathcal{L}(f_{s}), \mathcal{L}_{con}(\theta, \phi) - \mathcal{L}_{con}(f_{s}, f_{v}) \right\} \\ = (1 + |\lambda|) \max \left\{ \mathcal{L}(\theta) - \mathcal{L}(f_{s}), \mathcal{L}_{con}(\theta, \phi) - \mathcal{L}_{con}(f_{s}, f_{v}) \right\} \\ = (1 + |\lambda|) \max \left\{ \mathbb{E}[\ell(h_{s}(\mathbf{x};\theta); y) - \ell(f_{s}(\mathbf{x}); y)], \mathbb{E}[d(\mathbf{x}, D(h_{s}(\mathbf{x};\theta), h_{v}(\tilde{\mathbf{x}}; \phi))) - d(\mathbf{x}, D(f_{s}(\mathbf{x}), f_{v}(\tilde{\mathbf{x}})))] \right\} \\ \leq (1 + |\lambda|) \mathbb{E} \left\{ \max[[\ell(h_{s}(\mathbf{x};\theta); y) - \ell(f_{s}(\mathbf{x}); y)], \mathbb{E}[d(\mathbf{x}, D(h_{s}(\mathbf{x};\theta), h_{v}(\tilde{\mathbf{x}}; \phi))) - d(\mathbf{x}, D(f_{s}(\mathbf{x}), f_{v}(\tilde{\mathbf{x}})))] \right\} \\ \leq (1 + |\lambda|) \mathbb{E} \left\{ \max[L_{\ell}|h_{s}(\mathbf{x};\theta) - f_{s}(\mathbf{x})|, L_{d}|D(h_{s}(\mathbf{x};\theta), h_{v}(\tilde{\mathbf{x}}; \phi)) - D(f_{s}(\mathbf{x}), f_{v}(\tilde{\mathbf{x}}))|] \right\} \\ \leq (1 + |\lambda|) \mathbb{E} \max\{L_{\ell}\epsilon_{s}, L_{d}\epsilon_{g}\} \\ = (1 + |\lambda|) \max\{L_{\ell}\epsilon_{s}, L_{d}\epsilon_{g}\}$$

where the first inequality is using Hölder's inequality (18) when  $p = 1, q = \infty$  and the second one is by the convexity of max-norm and Jensen's inequality. The third inequality is applying  $L_{\ell}$  and  $L_d$  lipschitzness on  $\ell$  and d. The fourth one is due to  $\epsilon_s$  and  $\epsilon_g$  parameterization on  $f_s$  and D.

$$\sum_{k=1}^{n} |x_k y_k| \le \left(\sum_{k=1}^{n} |x_k|^p\right)^{1/p} \left(\sum_{k=1}^{n} |y_k|^q\right)^{1/q}$$
(18)

Then Eq. (16) becomes

$$\mathcal{D}_{\epsilon}^{\star}(\gamma) \leq \max_{\lambda} \min_{f_x, f_v} \mathcal{R} + (1 + |\lambda|) \max\{L_{\ell}\epsilon_s, L_d\epsilon_g\} \triangleq D_p^*$$
(19)

In order to use strong duality to bound  $\mathcal{D}_p^*$ , we need to construct primal problem whose optimum is  $\mathcal{P}_p^*$  from  $\mathcal{D}_p^*$ . Therefore, the remaining proof is to show  $\mathcal{D}_p^*$  is the dual problem to a constraint statistical learning problem with perturbation function as  $m = \max \{L_\ell \epsilon_s, L_d \epsilon_g\}$ . Specifically, when  $\lambda > 0$ ,  $\mathcal{D}_p^*$  can be expanded and rearranged as

$$\mathcal{D}_{p}^{*} = \max_{\lambda} \min_{f_{s}, f_{v}} \mathcal{L}(f_{s}) + \lambda \mathcal{L}_{con}(f_{s}, f_{v}) + (1 + |\lambda|)m$$
  
$$= \max_{\lambda} \min_{f_{s}, f_{v}} \ell(f_{s}(\mathbf{x}; \theta); y) + m + \lambda [d(\mathbf{x}, D(f_{s}(\mathbf{x}; \theta), f_{v}(\tilde{\mathbf{x}}; \phi))) - \gamma + m]$$
(20)

By the feasibility assumption 3, Eq. (20) can be considered as the dual problem of the following one:

$$\mathcal{P}_{p}^{\star}(\gamma) \triangleq \min_{f_{s} \in \mathcal{F}} \mathcal{L}\left(f_{s}\right) + m$$

$$s.t. \quad d\left(\mathbf{x}, D\left(f_{s}(\mathbf{x}; \theta), f_{v}(\tilde{\mathbf{x}}; \phi)\right)\right) \leq \gamma - m$$
(21)

Denote  $f_s^*$  and  $f_v^*$  as the primal solutions to  $\mathcal{P}_p^*$ ,  $\lambda^*$  and  $\lambda_p^*$  as the dual variable that give solutions  $\mathcal{D}_p^*$ ,  $\mathcal{P}_p^*$ . With the regularity assumption 2, the saddle point condition shows for  $\forall \tilde{f}_s, \tilde{f}_v \in \mathcal{H}$  and  $\lambda, \lambda_p^* \in \mathbb{R}^+$ , it holds that

$$\mathcal{R}(f_{s,p}^*, f_{v,p}^*, \tilde{\lambda}) + (1 + |\tilde{\lambda}|) \leq \max_{\lambda} \min_{f_s, f_v} \mathcal{R}(f_s, f_v, \lambda) + (1 + |\lambda|)m = \mathcal{D}_p^* = \mathcal{P}_p^*$$
$$= \min_{f_s, f_v} \max_{\lambda} \mathcal{R}(f_s, f_v, \lambda) + (1 + |\lambda|)m$$
$$\leq \mathcal{R}(\tilde{f}_s, \tilde{f}_v, \lambda_p^*) + (1 + |\lambda_p^*|)m$$
(22)

Recall that Eq. (22) holds for  $\forall \tilde{f}_s, \tilde{f}_v \in \mathcal{H}$ . Let  $\tilde{f}_s = f_s^*, \tilde{f}_v = f_v^*$ , then we use it to upper bound  $D_{\epsilon}^*$  as

$$\mathcal{D}_{\epsilon}^{*}(\gamma) \leq \mathcal{P}_{p}^{*} \leq \mathcal{R}(f_{s}^{*}, f_{v}^{*}, \lambda^{*}) + (1 + |\lambda_{p}^{*}|)m = \mathcal{P}^{*}(\gamma) + (1 + |\lambda_{p}^{*}|)m$$

$$\tag{23}$$

The last step is due to Eq. (14).

### A.2. Proof for empirical gap

**Proof A.2** Similar to [16, 60], by KKT conditions and complementary slackness conditions [10] shows

$$\lambda_{\epsilon}^{*}(\mathbb{E}_{\mathbf{x},\tilde{\mathbf{x}}\sim\mathbb{P}(X)}\left[d\left(\mathbf{x}, D\left(f_{s}(\mathbf{x};\theta_{\epsilon}^{*}), f_{v}(\tilde{\mathbf{x}};\phi_{\epsilon}^{*})\right)\right) - \gamma\right] = 0$$
  
$$\lambda_{\epsilon,n}^{*}\left(\sum_{i=1}^{n}\sum_{j\neq i}^{n}d\left(\mathbf{x}_{i}, D\left(f_{s}(\mathbf{x}_{i};\theta_{\epsilon}^{*}), f_{v}(\mathbf{x}_{j};\phi_{\epsilon}^{*})\right) - \gamma\right) = 0$$
(24)

where  $(\theta_{\epsilon}^*, \phi_{\epsilon}^*, \lambda_{\epsilon}^*)$  and  $(\theta_{\epsilon,n}^*, \phi_{\epsilon,n}^*, \lambda_{\epsilon,n}^*)$  are primal-dual pairs for achieving the optimum  $\mathcal{D}_{\epsilon}^*(\gamma)$  and  $\mathcal{D}_{\epsilon,n}^*(\gamma)$ . Eq. (24) implies the constraint-related terms in the objectives to be zero, then consider the remaining term as

$$\mathcal{D}_{\epsilon}^{*}(\gamma) = \mathbb{E}\left[\ell\left(f_{s}(\mathbf{x};\theta);y\right)\right] \triangleq \mathcal{M}(\theta_{\epsilon}^{*})$$
$$\mathcal{D}_{\epsilon,n}^{*}(\gamma) = \sum_{i=1}^{n} \ell\left(f_{s}\left(\mathbf{x}_{i}\right),y_{i}\right) \triangleq \hat{M}(\theta_{\epsilon,n}^{*})$$
(25)

Thus the empirical gap reduces to

$$|\mathcal{D}_{\epsilon}^{*}(\gamma) - \mathcal{D}_{\epsilon,n}^{*}(\gamma)| = |\mathcal{M}(\theta_{\epsilon}^{*}) - \mathcal{\widetilde{M}}(\theta_{\epsilon,n}^{*})|$$
(26)

Using the fact that  $\theta_{\epsilon}^*$  and  $\theta_{\epsilon,n}^*$  are optimal for  $\mathcal{M}(\theta_{\epsilon}^*)$  and  $\mathcal{M}(\theta_{\epsilon,n}^*)$ , the following holds

$$\mathcal{M}(\theta_{\epsilon}^*) - \hat{\mathcal{M}}(\theta_{\epsilon}^*) \le \mathcal{M}(\theta_{\epsilon}^*) - \hat{\mathcal{M}}(\theta_{\epsilon,n}^*) \le \mathcal{M}(\theta_{\epsilon,n}^*) - \hat{\mathcal{M}}(\theta_{\epsilon,n}^*)$$
(27)

Therefore, using the above lower and upper bound, we can bound Eq. (26) as

$$|\mathcal{D}_{\epsilon}^{*}(\gamma) - \mathcal{D}_{\epsilon,n}^{*}(\gamma)| \leq \max\{|\mathcal{M}(\theta_{\epsilon}^{*}) - \hat{\mathcal{M}}(\theta_{\epsilon}^{*})|, |\mathcal{M}(\theta_{\epsilon,n}^{*}) - \hat{\mathcal{M}}(\theta_{\epsilon,n}^{*})|\}$$
(28)

Then we resort to the classical VC-dimension bounds for the above two terms in Eq. (28) as follows

$$|\mathcal{M}(\theta) - \hat{\mathcal{M}}(\theta)| \le 2B\sqrt{\frac{1}{n} \left[1 + \log\left(\frac{4(2n)^{d_{\rm VC}}}{\delta}\right)\right]}$$
(29)

holds with probability  $1 - \delta$  when  $d_{VC}$  is the VC dimension for all  $\theta$ . Combing Eq. (28) and (29) completes the proof.

# A.3. Proof for empirical duality gap

**Proof A.3** Simply combining the results in the above lemmas, i.e. parameterization gap and empirical gap, via applying the triangle inequality completes the proof.

$$\begin{aligned} \mathcal{P}^{\star} - \mathcal{D}_{\varepsilon,n}^{\star}(\gamma) &| = \left| \mathcal{P}^{\star} + \mathcal{D}_{\epsilon}^{\star}(\gamma) - \mathcal{D}_{\epsilon}^{\star}(\gamma) - \mathcal{D}_{\epsilon,n}^{\star}(\gamma) \right| \\ \leq & \left| \mathcal{P}^{\star} - \mathcal{D}_{\varepsilon}^{\star}(\gamma) \right| + \left| \mathcal{D}_{\epsilon}^{\star}(\gamma) - \mathcal{D}_{\epsilon,n}^{\star}(\gamma) \right| \\ \leq & \left( 1 + |\lambda| \right) m + 2B \sqrt{\frac{1}{n} \left[ 1 + \log\left(\frac{4(2n)^{d_{vc}}}{\delta}\right) \right]} \end{aligned}$$

#### **B.** Domain Generalization by Learning on Fictitious Distributions

This section gives a justification for disentanglement from a different perspective by connecting the dots with classical domain adaptation. Specifically, we construct a fictitious distribution to extend it to the DG setting and decompose the target learning objective into empirical learning errors, domain divergence and source domain data diversity. Moreover, we show that learning disentangled representations gives a tighter risk upper bound.

With a slight abuse of notation, let  $\mathcal{H}$  be a hypothesis space and denote  $\hat{\mathcal{D}}$  as the induced distribution over feature space  $\mathcal{Z}$  for every distribution  $\mathcal{D}$  over the raw space. Define  $\mathcal{D}_S^i$  as the source distribution over  $\mathcal{X}$ , which enables a mixture construction of source domains as  $\mathcal{D}_S^{\alpha} = \sum_{i=1}^{N_s} \alpha_i \mathcal{D}_S^i(\cdot)$ . Denote a fictitious distribution  $\mathcal{D}_U^{\alpha} = \sum_{i=1}^{N_s} \alpha_i^* \mathcal{D}_S^i(\cdot)$  as the convex combination of source domains which is the closest to  $\mathcal{D}_U$ , where  $\alpha_1^*, ..., \alpha_{N_S}^* = \arg \min_{\alpha_1, ..., \alpha_{N_s}} d_{\mathcal{H}}(\mathcal{D}_U, \sum_{i=1}^{N_s} \alpha_i \mathcal{D}_S^i(\cdot))$ . The fictitious distribution induces a feature space distribution  $\tilde{\mathcal{D}}_U^{\alpha} = \sum_{i=1}^{N_s} \alpha_i^* \tilde{\mathcal{D}}_S^i(\cdot)$ . The following inequality holds for the risk  $\epsilon_U(h)$  on any unseen target domain  $\mathcal{D}_U$ .

$$\epsilon_{U}(h) \leq \lambda_{\alpha} + \underbrace{\sum_{i=1}^{N_{S}} \alpha_{i} \epsilon_{S,i}(h)}_{(1) \text{ Empirical}} + \underbrace{d_{\mathcal{H}}(\tilde{\mathcal{D}}_{U}^{\alpha}, \tilde{\mathcal{D}}_{S}^{\alpha})}_{(2) \text{ Divergence}} + \underbrace{d_{\mathcal{H}}(\tilde{\mathcal{D}}_{U}, \tilde{\mathcal{D}}_{U}^{\alpha})}_{(3) \text{ Diversity}}$$
(30)

where  $\lambda_{\alpha}$  is the risk of the optimal hypothesis on the mixture source domain  $\mathcal{D}_{S}^{\alpha}$  and  $\mathcal{D}_{U}$ .

We define the symmetric difference hypothesis space as  $\mathcal{H}\Delta\mathcal{H} = \{h(\mathbf{x}) \oplus h'(\mathbf{x}) : h, h' \in \mathcal{H}\}$ , where  $\oplus$  is the XOR operator. Applying [6], we have

$$\begin{aligned} \varepsilon_{U}(h) &\leq \lambda_{U} + \Pr_{\mathcal{D}_{U}}[\mathcal{Z}_{h} \triangle \mathcal{Z}_{h}^{*}] \\ &\leq \lambda_{U} + \Pr_{\mathcal{D}_{S}^{\alpha}}[\mathcal{Z}_{h} \triangle \mathcal{Z}_{h}^{*}] + |\Pr_{\mathcal{D}_{S}^{\alpha}}[\mathcal{Z}_{h} \triangle \mathcal{Z}_{h}^{*}] - \Pr_{\mathcal{D}_{U}}[\mathcal{Z}_{h} \triangle \mathcal{Z}_{h}^{*}]| \\ &\leq \lambda_{U} + \Pr_{\mathcal{D}_{S}^{\alpha}}[\mathcal{Z}_{h} \triangle \mathcal{Z}_{h}^{*}] + d_{\mathcal{H}}(\tilde{D}_{U}, \tilde{\mathcal{D}}_{S}^{\alpha}) \\ &\leq \lambda_{U} + \Pr_{\mathcal{D}_{S}^{\alpha}}[\mathcal{Z}_{h} \triangle \mathcal{Z}_{h}^{*}] + d_{\mathcal{H}}(\tilde{\mathcal{D}}_{U}^{\alpha}, \tilde{\mathcal{D}}_{S}^{\alpha}) + d_{\mathcal{H}}(\tilde{\mathcal{D}}_{U}, \tilde{\mathcal{D}}_{U}^{\alpha}) \\ &\leq \lambda_{u} + \Pr_{\mathcal{D}_{S}^{\alpha}}[\mathcal{Z}_{h} \triangle \mathcal{Z}_{h}^{*}] + d_{\mathcal{H}}(\tilde{\mathcal{D}}_{U}^{\alpha}, \tilde{\mathcal{D}}_{S}^{\alpha}) + d_{\mathcal{H}}(\tilde{\mathcal{D}}_{U}, \tilde{\mathcal{D}}_{U}^{\alpha}) \\ &\leq \lambda_{u} + \sum_{i=1}^{N_{S}} \alpha_{i}\epsilon_{S,i}(h) + \underbrace{d_{\mathcal{H}}(\tilde{\mathcal{D}}_{U}^{\alpha}, \tilde{\mathcal{D}}_{S}^{\alpha})}_{(2) \text{ Divergence}} + \underbrace{d_{\mathcal{H}}(\tilde{\mathcal{D}}_{U}, \tilde{\mathcal{D}}_{U}^{\alpha})}_{(3) \text{ Diversity}} \end{aligned}$$
(31)

The fourth inequality holds because of the triangle inequality. We provide the explanation for our bound in the Eq. (30) and Eq. (31). The second term is the empirical loss for the convex combination of all source domains. The third term corresponds to "To what extent can the convex combination of the source domain approximate the target domain". The minimization of the third term requires diverse data or strong data augmentation, such that the unseen distribution lies within the convex combination of source domains. For the fourth term, the following equation holds for any two distributions  $D'_U, D''_U$ , which are the convex combinations of source domains [12]

$$d_{\mathcal{H}[D'_U, D''_U]} \le \sum_{l=1}^{N_S} \sum_{k=1}^{N_S} \alpha_l \alpha_k d_{\mathcal{H}}[\mathcal{D}_{S,l}, \mathcal{D}_{S,k}]$$
(32)

Such an upper bound will be minimized when  $d_{\mathcal{H}}[\mathcal{D}_{S,l}, \mathcal{D}_{S,k}] = 0, \forall l, k \in \{1, ..., N_S\}$ . Namely projecting the source domain data into a feature space, where the source domain labels are hard to distinguish.

The above (3) Diversity term is also supported by the evidence that compositional generalization and extrapolation can be improved if the training domain data are rich enough [13, 32, 54]. To this end, one can obviously simulate data points with predetermined data augmentation methods such as rotating, cropping, Gaussian blur, color jitter, etc. However, their developments require prior knowledge and domain-specific expertise like translation-invariance on images, which is likely to fail in the unseen domain due to distribution shifts. It motivates learning disentangled representations that are transferable across various domains [21]. Thus we discuss the benefits of disentanglement on the domain generalization gap in the following section. Assume that the semantic and the variation factors are disentangled in the latent space S and V, then the errors [64] on the disentangled source and target domain with a hypothesis h are

$$\epsilon_{S,i}(h) = \epsilon_{S,i}^s(h) + \epsilon_{S,i}^v(h), \\ \epsilon_U(h) = \epsilon_U^s(h) + \epsilon_U^v(h)$$
(33)

Given  $h^* = \arg \min_{h \in \mathcal{H}} \left( \epsilon_{S,i}^s(h), \epsilon_{S,i}^v(h) \right) \forall i \in \{1, \dots, N_S\}$ , since  $\epsilon_U(h) = \epsilon_U^s(h) + \epsilon_U^v(h)$ , combining Eq. (30) and we have

$$\epsilon_U^s(h) + \epsilon_U^v(h) \le \lambda_\alpha + \sum_{i=1}^{N_S} \alpha_i \epsilon_{S,i}(h) + d_{\mathcal{H}}(\tilde{\mathcal{D}}_U^\alpha, \tilde{\mathcal{D}}_S^\alpha) + d_{\mathcal{H}}(\tilde{\mathcal{D}}_U, \tilde{\mathcal{D}}_U^\alpha)$$
(34)

where  $\lambda_{\alpha} = \epsilon_U(h^*) + \sum_{i=1}^{N_S} \alpha_i \epsilon_{S,i}(h^*) = \epsilon_U^s(h^*) + \epsilon_U^v(h^*) + \sum_{i=1}^{N_S} \alpha_i \epsilon_{S,i}^s(h^*) + \sum_{i=1}^{N_S} \alpha_i \epsilon_{S,i}^v(h^*)$ . Then the upper bound for the unseen domain can be further derived as follow,

$$\epsilon_{U}^{s}(h) \leq \lambda_{\alpha} + \sum_{i=1}^{N_{S}} \alpha_{i} \epsilon_{S,i}(h) + d_{\mathcal{H}}(\tilde{\mathcal{D}}_{U}^{\alpha}, \tilde{\mathcal{D}}_{S}^{\alpha}) + d_{\mathcal{H}}(\tilde{\mathcal{D}}_{U}, \tilde{\mathcal{D}}_{U}^{\alpha}) - \epsilon_{U}^{v}(h)$$
(35)

Combining Eq. (30) and Eq. (33), we have

$$\epsilon_{U}^{s}(h) \leq \lambda_{\alpha} + \underbrace{\sum_{i=1}^{N_{S}} \alpha_{i} \epsilon_{S,i}(h)}_{(1) \text{ Empirical}} + \underbrace{d_{\mathcal{H}}(\tilde{\mathcal{D}}_{U}^{\alpha}, \tilde{\mathcal{D}}_{S}^{\alpha})}_{(2) \text{ Divergence}} + \underbrace{d_{\mathcal{H}}(\tilde{\mathcal{D}}_{U}, \tilde{\mathcal{D}}_{U}^{\alpha})}_{(3) \text{ Diversity}} - \underbrace{\epsilon_{U}^{v}(h)}_{(4)}$$
(36)

where (1) denotes the empirical loss over every source domain, (2) means divergence minimization among source domains and (3) encourages the diversity and coverage of the mixture of source domains. The above result shows the benefits of disentanglement over representation spaces. Thus we propose to use two separate encoders representing semantic and variation subspaces respectively in the following section. Such a formulation combined with the disentanglement term (4) in Eq. (36) implies that we should perform ERM over the semantic space only. The above analysis essentially justifies the design of DDG from a classical domain adaptation perspective: optimizing empirical risk over semantic space while promoting diversity and divergence by taking disentanglement as a constraint.

# **C. Experimental Settings**

# **C.1. Other Training Details**

We optimize all models using Adam [40]. For all the detailed hyperparameter settings, please refer to our code which is publicly available at https://github.com/hlzhang109/DDG.

## C.2. Dataset Statistics and Visualization

We show some images of the datasets in Fig. 7 to give an intuitive comparison among these image datasets. One can observe that these images have a diverse set of styles, making it very challenging to transfer knowledge from one to another.



Figure 7. Samples of DG datasets. The training data (a) PACS and (b) VLCS (c) Wilds are shown. PACS has four domains *art* (A), *cartoons* (C), *photos* (P), *sketches* (S). VLCS contains four domains *Caltech101* (C), *LabelMe* (L), *SUN09* (S), *VOC2007* (V). The WILDS dataset includes data from five different medical centers as domains.

# **D.** Additional Experimental Results

**Qualitative comparison with AugMix.** Looking into the Fig. 8, it is harder for the heuristic-based method AugMix to generate diverse samples via interpolation for training compared with DDG as shown in Fig. 3 and Fig. 10b.

More qualitative results via interpolation. Fig. 9 showcases the results of combining the semantic code of one image and the mixture of two variation codes. Results show that the model can generate samples with intermediate variation states.

**More qualitative results via swapping variation and semantic factors.** We showcase the qualitative results of swapping variation and semantic factors with PACS in Fig. 10b, MNIST in Fig. 10a, and WILDS in Fig. 10c. The results demonstrates the strong disentangled capability of DDG. Some interesting observations are DDG learns both intra- (e.g. thickness) and inter-domain (e.g. rotated angle) variations over RotatedMNIST. DDG also maintains semantic information like the color of distinct features across variation-rich data like PACS.



Figure 8. The augmented samples from AugMix [30]. The second and third rows are generated by applying AugMix to the first row.



Figure 9. Interpolation via mixing results on PACS.



(c) Wilds

Figure 10. Qualitative disentanglement results on RotatedMNIST, PACS and Wilds. In every panel, the training data in the first row manifests the semantic factors.

**Numerical comparison with MBDG.** We quantatively compare DDG with MBDG [60] as in Table 2. For a consistent comparison, we run the source code of authors under a test domain validation protocal [27] on PACS with the results as:

	А	С	Р	S	Avg
MBDG	$82.0\pm0.0$	$78.4\pm0.01$	$93.9\pm0.0$	$85.0\pm0.0$	85.8
$MBDG_{Reg}$	$84.9\pm0.0$	$84.9\pm0.0$	$93.9\pm0.0$	$\textbf{85.6} \pm \textbf{0.0}$	87.3
DDG	$\textbf{88.9} \pm \textbf{0.6}$	$\textbf{85.0} \pm \textbf{1.9}$	$\textbf{97.2} \pm \textbf{1.2}$	$84.3\pm0.7$	88.9

Table 2. Numerical results for comparing DDG and MBDG.

Though adopting similar PACCL frameworks and primal-dual algorithms, DDG can consistently outperform MBDG and its variant except over domain S. The primary reason behind the clear performance gain of DDG can be that DDG is better at capturing variations within data via random sampling without relying on domain labels. Specifically, parameterizing  $h_s$ ,  $h_v$ , D

and constrain them based on disentanglement makes the model more robust to both inter- and intra-domain nuisance factors compared to MBDG that only use a pretrained generator to simulate inter-domain variations. Moreover, the performance gain is also partly due to the three major differences between our approach and MBDG we highlight in the Related Work: First, our upper bound of the parameterization gap is tighter under mild conditions, whereas MBDG requires unrealistic assumptions on the distance metric, i.e.,  $d(\cdot, \cdot)$  satisfies Lipschitz-like inequality on both arguments, which is stronger than our normal  $L_d$  Lipschitzness assumption; second, MBDG consumes additional domain labels, which are prohibitively expensive or even infeasible to obtain in safety-critical applications or those containing sensitive demographics; third, DDG enforces invariance constraints via parameterizing semantic and variation encoders, which does not belong to a model-based approach. In contrast, MBDG requires a pre-trained domain transformation model (e.g., CycleGAN) during training, which may result in sub-optimal solutions and parameter inefficiency, while DDG is more flexible by treating this as a design choice.