# Supplementary Material for "STEM: An approach to Multi-source Domain Adaptation with Guarantees"

This is the supplementary material for "STEM: An approach to Multi-source Domain Adaptation with Guarantees". In the first section of this supplementary material, we provide proof for our theory developed, while presenting implementation specification and additional experimental results in the second section.

## **1** Theoretical Developments

**Theorem 1.** If  $\ell$  is a convex function, the following statements hold true:

i)  $\mathcal{L}(h^S, \mathbb{D}^S_{\pi}) \leq \max_{1 \leq k \leq K} \mathcal{L}(h^S_k, \mathbb{D}^S_k).$ 

ii) If each individual expert is an  $\epsilon$ -qualified classifier (i.e.,  $\mathcal{L}(h_k^S, \mathbb{D}_k^S) \leq \epsilon$ ), the multi-source teacher expert  $h^S$  is also an  $\epsilon$ -qualified classifier (i.e.,  $\mathcal{L}(h^S, \mathbb{D}_{\pi}^S) \leq \epsilon$ ).

Proof. i) We have

$$\begin{split} \mathcal{L}\left(h^{S}, \mathbb{D}_{\pi}^{S}\right) &= \int_{\mathcal{X} \times \mathcal{Y}} \ell\left(h^{S}\left(\mathbf{x}, y\right)\right) p_{\pi}^{S}\left(\mathbf{x}, y\right) d\mathbf{x} dy \\ &\leq \int_{\mathcal{X} \times \mathcal{Y}} \ell\left(\sum_{k=1}^{K} \frac{\pi_{k} p_{k}^{S}\left(\mathbf{x}, y\right)}{\sum_{j=1}^{K} \pi_{j} p_{j}^{S}\left(\mathbf{x}, y\right)} h_{k}^{S}\left(\mathbf{x}, y\right)\right) p_{\pi}^{S}\left(\mathbf{x}, y\right) d\mathbf{x} dy \\ &\leq \int_{\mathcal{X} \times \mathcal{Y}} \sum_{k=1}^{K} \frac{\pi_{k} p_{k}^{S}\left(\mathbf{x}, y\right)}{\sum_{j=1}^{K} \pi_{j} p_{j}^{S}\left(\mathbf{x}, y\right)} \ell\left(h_{k}^{S}\left(\mathbf{x}\right), y\right) p_{\pi}^{S}\left(\mathbf{x}, y\right) d\mathbf{x} dy \\ &\leq \sum_{k=1}^{K} \pi_{k} \int_{\mathcal{X} \times \mathcal{Y}} \frac{p_{k}^{S}\left(\mathbf{x}, y\right)}{p_{\pi}^{S}\left(\mathbf{x}, y\right)} \ell\left(h_{k}^{S}\left(\mathbf{x}\right), y\right) p_{\pi}^{S}\left(\mathbf{x}, y\right) d\mathbf{x} dy \\ &= \sum_{k=1}^{K} \pi_{k} \int_{\mathcal{X} \times \mathcal{Y}} \ell\left(h_{k}^{S}\left(\mathbf{x}\right), y\right) p_{k}^{S}\left(\mathbf{x}, y\right) d\mathbf{x} dy \\ &= \sum_{k=1}^{K} \pi_{k} \mathcal{L}\left(h_{k}^{S}, \mathbb{D}_{k}^{S}\right) \leq \max_{1 \leq k \leq K} \mathcal{L}\left(h_{k}^{S}, \mathbb{D}_{k}^{S}\right). \end{split}$$

Note that  $\ell \left( h_k^S \left( \mathbf{x} \right), y \right) := \ell \left( h_k^S \left( \mathbf{x}, y \right) \right)$  where  $h_k^S \left( \mathbf{x} \right) = \left[ h_k^S \left( \mathbf{x}, y \right) \right]_{y=1}^M$  and  $\ell \left( t \right) = -\log \left( t \right)$  for the cross-entropy loss. ii) It is trivial from (i).

As indicated by Theorem 1, the multi-source teacher expert  $h^S$  predicts well data examples sampled from  $\mathcal{D}^S_{\pi}$  which is a mixture of  $\mathcal{D}^S_{1:K}$ . It is natural to ask the question of the factors that influence the performance of  $h^S$  when *predicting on* the target joint distribution  $\mathbb{D}^T$ . To facilitate the following theorem, we define a hybrid joint distribution  $\mathbb{D}^h$  with the density function  $p^h(\mathbf{x}, y)$  consisting of the pairs  $(\mathbf{x}, y)$  in which the data example  $\mathbf{x} \sim \mathbb{P}^S_{\pi} = \sum_{k=1}^K \pi_k \mathbb{P}^S_k$  and the label y is sampled using  $p^T(y \mid \mathbf{x})$ :

$$p^{h}(\mathbf{x}, y) = p^{T}(y \mid \mathbf{x}) \sum_{k=1}^{K} \pi_{k} p_{k}^{S}(\mathbf{x}) = p^{T}(y \mid \mathbf{x}) p_{\pi}^{S}(\mathbf{x}).$$

**Theorem 3.** If  $\ell$  is a convex function and upper-bounded by a positive constant *L*, the general loss  $\mathcal{L}(h^S, \mathbb{D}^T)$  is upperbounded by:

*i)*  $A\left[\max_{k} \mathcal{L}\left(h_{k}^{S}, \mathbb{D}_{k}^{S}\right) + L\max_{k} \mathbb{E}_{\mathbb{P}_{k}^{S}}\left[\left\|\Delta p_{k}\left(y \mid \mathbf{x}\right)\right\|_{1}\right]\right]^{\frac{\alpha-1}{\alpha}}$  where  $A = \exp\left\{R^{\alpha}\left(\mathbb{P}^{T}\|\mathbb{P}_{\pi}^{S}\right)\right\}^{\frac{\alpha-1}{\alpha}} L^{\frac{1}{\alpha}}$  in which  $R^{\alpha}\left(\mathbb{P}^{T}\|\mathbb{P}_{\pi}^{S}\right)$  represents the Rényi divergence between those distributions and  $\Delta p_{k}\left(y \mid \mathbf{x}\right) := \left[\left|p_{k}^{S}\left(y = m \mid \mathbf{x}\right) - p^{T}\left(y = m \mid \mathbf{x}\right)\right|\right]_{m=1}^{M}$  represents the label shift between the labeling assignment mechanisms of an individual source domain and target domain.

ii) 
$$A\left[\epsilon + L \max_{k} \mathbb{E}_{\mathbb{P}_{k}^{S}}\left[\left\|\Delta p_{k}\left(y \mid \mathbf{x}\right)\right\|_{1}\right]\right]^{\frac{\alpha}{\alpha}}$$
 provided that  $\mathcal{L}\left(h_{k}^{S}, \mathbb{D}_{k}^{S}\right) \leq \epsilon, \forall k = 1, ..., K.$ 

*Proof.* i) We have

$$\begin{split} \mathcal{L}\left(h^{S}, \mathbb{D}^{h}\right) &= \int \ell\left(h^{S}\left(\mathbf{x}\right), y\right) p^{h}\left(\mathbf{x}, y\right) d\mathbf{x} dy = \mathcal{L}\left(h^{S}, \mathbb{D}_{\pi}^{S}\right) + \int \ell\left(h^{S}\left(\mathbf{x}\right), y\right) \left[p^{h}\left(\mathbf{x}, y\right) - p_{\pi}^{S}\left(\mathbf{x}, y\right)\right] d\mathbf{x} dy \\ &\leq \mathcal{L}\left(h^{S}, \mathbb{D}_{\pi}^{S}\right) + \int \ell\left(h^{S}\left(\mathbf{x}\right), y\right) |p^{h}\left(\mathbf{x}, y\right) - p_{\pi}^{S}\left(\mathbf{x}, y\right)| d\mathbf{x} dy \\ &= \mathcal{L}\left(h^{S}, \mathbb{D}_{\pi}^{S}\right) + L \int |p^{h}\left(\mathbf{x}, y\right) - p_{\pi}^{S}\left(\mathbf{x}, y\right)| d\mathbf{x} dy \\ &= \mathcal{L}\left(h^{S}, \mathbb{D}_{\pi}^{S}\right) + L \int \sum_{k=1}^{K} \pi_{k} p_{k}^{S}\left(\mathbf{x}\right) |p^{T}\left(y \mid \mathbf{x}\right) - p_{k}^{S}\left(y \mid \mathbf{x}\right)| d\mathbf{x} dy \\ &= \mathcal{L}\left(h^{S}, \mathbb{D}_{\pi}^{S}\right) + L \sum_{k=1}^{K} \pi_{k} \int p_{k}^{S}\left(\mathbf{x}\right) \sum_{y=1}^{M} |p^{T}\left(y \mid \mathbf{x}\right) - p_{k}^{S}\left(y \mid \mathbf{x}\right)| d\mathbf{x} \\ &= \mathcal{L}\left(h^{S}, \mathbb{D}_{\pi}^{S}\right) + L \sum_{k=1}^{K} \pi_{k} \mathbb{E}_{\mathbb{P}_{k}^{S}} \left[ \left\|\Delta p_{k}\left(y \mid \mathbf{x}\right)\right\|_{1} \right] \\ &\leq \max_{k} \mathcal{L}\left(h_{k}^{S}, \mathbb{D}_{k}^{S}\right) + L \max_{k} \mathbb{E}_{\mathbb{P}_{k}^{S}} \left[ \left\|\Delta p_{k}\left(y \mid \mathbf{x}\right)\right\|_{1} \right]. \end{split}$$

Note that we have used  $\mathcal{L}(h^S, \mathbb{D}^S_{\pi}) \leq \max_k \mathcal{L}(h^S_k, \mathbb{D}^S_k)$  (referred to Theorem 1). Finally, we manipulate  $\mathcal{L}(h^S, \mathbb{D}^T)$  as

$$\begin{split} \mathcal{L}\left(h^{S}, \mathbb{D}^{T}\right) &= \int_{\mathcal{X} \times \mathcal{Y}} \ell\left(h^{S}\left(\mathbf{x}\right), y\right) p^{T}\left(\mathbf{x}, y\right) d\mathbf{x} dy \\ &= \int_{\mathcal{X} \times \mathcal{Y}} \frac{p^{T}\left(\mathbf{x}, y\right)}{p^{h}\left(\mathbf{x}, y\right)^{\frac{\alpha-1}{\alpha}}} p^{h}\left(\mathbf{x}, y\right)^{\frac{\alpha-1}{\alpha}} \ell\left(h^{S}\left(\mathbf{x}\right), y\right) d\mathbf{x} dy \end{split}$$

The Holder inequality gives us

$$\mathcal{L}\left(h^{S}, \mathbb{D}^{T}\right) \leq \left[\int_{\mathcal{X} \times \mathcal{Y}} \frac{p^{T}\left(\mathbf{x}, y\right)^{\alpha}}{p^{h}\left(\mathbf{x}, y\right)^{\alpha-1}} d\mathbf{x} dy\right]^{\frac{1}{\alpha}} \left[\int_{\mathcal{X} \times \mathcal{Y}} p^{h}\left(\mathbf{x}, y\right) \ell\left(h^{S}\left(\mathbf{x}\right), y\right)^{\frac{\alpha}{\alpha-1}} d\mathbf{x} dy\right]^{\frac{\alpha-1}{\alpha}}.$$

Referring to the definition of the Rényi divergence and note that  $\ell(h^{S}(\mathbf{x}), y) \leq L$ , we obtain

$$\mathcal{L}\left(h^{S}, \mathbb{D}^{T}\right) \leq \left[\exp\left\{R^{\alpha}\left(\mathbb{D}^{T}\|\mathbb{D}^{h}\right)\right\} \mathcal{L}\left(h^{S}, \mathbb{D}^{h}\right)\right]^{\frac{\alpha-1}{\alpha}} L^{\frac{1}{\alpha}}.$$

We further derive

$$R^{\alpha} \left( \mathbb{D}^{T} \| \mathbb{D}^{h} \right) = \frac{1}{\alpha - 1} \log \int \left[ \frac{p^{T} \left( \mathbf{x}, y \right)}{p^{h} \left( \mathbf{x}, y \right)} \right]^{\alpha - 1} p^{T} \left( \mathbf{x}, y \right) d\mathbf{x} dy$$
$$= \frac{1}{\alpha - 1} \log \left( \int \left[ \frac{p^{T} \left( \mathbf{x} \right)}{p_{\pi}^{S} \left( \mathbf{x} \right)} \right]^{\alpha - 1} \sum_{y = 1}^{M} p^{T} \left( \mathbf{x}, y \right) d\mathbf{x} \right)$$
$$= \frac{1}{\alpha - 1} \log \left( \int \left[ \frac{p^{T} \left( \mathbf{x} \right)}{p_{\pi}^{S} \left( \mathbf{x} \right)} \right]^{\alpha - 1} p^{T} \left( \mathbf{x} \right) d\mathbf{x} \right) = R^{\alpha} \left( \mathbb{P}^{T} \| \mathbb{P}_{\pi}^{S} \right)$$

Finally, we reach the following inequality:

$$\mathcal{L}\left(h^{S}, \mathbb{D}^{T}\right) \leq \left[\exp\left\{R^{\alpha}\left(\mathbb{P}^{T} \|\mathbb{P}_{\pi}^{S}\right)\right\} \mathcal{L}\left(h^{S}, \mathbb{D}^{h}\right)\right]^{\frac{\alpha-1}{\alpha}} L^{\frac{1}{\alpha}}.$$

$$(h^{S}, \mathbb{D}^{S}) \leq \epsilon \forall k = 1 \qquad K$$

ii) It is trivial from (i) and  $\mathcal{L}(h_k^S, \mathbb{D}_k^S) \leq \epsilon, \forall k = 1, ..., K$ .

In what follows, we present how to train the multi-source teacher expert  $h^S$ . Our workaround to train  $h^S$  comes from the following theoretical observation. Assume that we have K distributions  $\mathbb{R}_{1:K}$  with density functions  $r_{1:K}$  (z). We form a joint distribution  $\mathcal{D}$  of a data instance z and label  $t \in \{1, ..., K\}$  by sampling an index  $t \sim \operatorname{Cat}(\pi)$ , sampling  $\mathbf{x} \sim \mathbb{R}_t$ , and collecting  $(\mathbf{z}, t)$  as a sample from  $\mathcal{D}$ . With this setting, we have the following corollary.

**Corollary 2.** If we train a source domain discriminator C to classify samples from the joint distribution D using the crossentropy loss (i.e.,  $CE(\cdot, \cdot)$ ), the optimal source domain discriminator  $C^*$  defined as

$$\mathcal{C}^{*} = argmin_{\mathcal{C}} \mathbb{E}_{(\mathbf{z},t) \sim \mathcal{D}} \left[ CE\left(\mathcal{C}\left(\mathbf{z}\right), t\right) \right]$$

satisfies  $C^{*}(\mathbf{z}) = \left[\frac{\pi_{i}r_{i}(\mathbf{z})}{\sum_{j}\pi_{j}r_{j}(\mathbf{z})}\right]_{i=1}^{K}$ .

Proof. We have

$$\mathbb{E}_{(\mathbf{z},t)\sim\mathcal{D}}\left[CE\left(\mathcal{C}\left(\mathbf{z}\right),t\right)\right] = \sum_{t=1}^{K} \pi_{t} \int CE\left(\mathcal{C}(\mathbf{z}),t\right) r_{t}\left(\mathbf{z}\right) d\mathbf{z}$$
$$= -\int \sum_{t=1}^{K} \log \mathcal{C}\left(\mathbf{z},t\right) \pi_{t} r_{t}\left(\mathbf{z}\right) d\mathbf{z}.$$

Given z, we now find  $C^* = [C_t^*]_{t=1}^K$  subjected to  $\|C^*\|_1 = 1$  and  $C^* \ge 0$  to maximize

$$\max_{\mathcal{C}:\left\|\mathcal{C}\right\|_{1}=1}\sum_{t=1}^{K}\log\mathcal{C}_{t}\pi_{t}r_{t}\left(\mathbf{z}\right).$$

The Lagrange function is as follows:

$$\mathcal{L} = \sum_{t=1}^{K} \log C_t \pi_t r_t \left( \mathbf{z} \right) - \lambda \left( \sum_{t=1}^{K} C_t - 1 \right).$$

Setting the derivatives to 0, we obtain

$$\frac{\partial \mathcal{L}}{\partial \mathcal{C}_{t}} = \frac{\pi_{t} r_{t} \left(\mathbf{z}\right)}{C_{t}} - \lambda = 0, \ t = 1, ..., K.$$

Note that  $\sum_{t=1}^{K} C_t = 1$ , we arrive at

$$\mathcal{C}_{t}^{*} = \frac{\pi_{t} r_{t} \left( \mathbf{z} \right)}{\sum_{j} \pi_{j} r_{j} \left( \mathbf{z} \right)}, \ t = 1, ..., K$$

Finally, we reach the conclusion.

## 2 Implementation Specification and Additional Experimental Results

### 2.1 Data preparation and preprocessing

**Digits-five**. The resolution of digit images is resized to  $32 \times 32$ , and we normalize pixel values to the range of [-1, 1]. **Office-Caltech10 and DomainNet**. We resize the resolution of images to  $224 \times 224$ , and then use those preprocessing images for ResNet-101 [1]. Additionally, for DomainNet, we find that adaptation tasks are much challenging due to the dissimilarity across domains. To reduce the domain gap, we apply the horizontal flip transformation to increase data diversity during the training process.

Architecture	Digit-five	Office-Caltech10	DomainNet
Input size	32  imes 32  imes 3	$224\times224\times3$	$224\times224\times3$
Generator G	instance normalization		transformation
	$3 \times 3$ conv. 64 lReLU	ResNet101	ResNet101
	$3 \times 3$ conv. 64 lReLU	256 dense, ReLU	
	$3 \times 3$ conv. 64 lReLU	dropout, $p = 0.5$	
	$2 \times 2$ max-pool, stride 2	Gaussian noise, $\sigma = 1$	
	dropout, $p = 0.5$		
	Gaussian noise, $\sigma = 1$		
	$3 \times 3$ conv. 64 lReLU		
	$3 \times 3$ conv. 64 lReLU		
	$3 \times 3$ conv. 64 lReLU		
	$2 \times 2$ max-pool, stride 2		
	dropout, $p = 0.5$		
	Gaussian noise, $\sigma = 1$		
Classifier $h_{1:K}^S, h^T$	$3 \times 3$ conv. 64 lReLU	M dense, softmax	M dense, softmax
	$3 \times 3$ conv. 64 lReLU		
	$3 \times 3$ conv. 64 lReLU		
	global average pool		
	M dense, softmax		
Source domain	100 dense, ReLU	K dense, ReLU	K dense, ReLU
discriminator $\mathcal{C}$	K dense, ReLU		

Table 1: Network architecture of STEM. The Leaky ReLU (IReLU) parameter *a* is set to 0.1. *K* and *M* denote the number of source domains and the number of classes respectively.



Figure 1: Test accuracy (%) when tweaking  $\alpha$  and  $\beta$  on " $\rightarrow$ **sv**" (*blue*) and " $\rightarrow$ **sy**" (*red*) tasks.



Figure 2: The comparison of test error of our STEM (red) with the current state-of-the-art LtC-MSDA (blue)

### 2.2 Architecture

The network architecture in detail for each dataset is shown in Table 1. We use a small convolutional neural network for Digitfive since this dataset is easy to be converged. For the two other datasets, we both use Resnet-101 pre-trained on ImageNet as the backbone, which is frozen when training on Office-Caltech10 and fine-tuned on DomainNet. All experiments are run on a computer with an NVIDIA Tesla V100 SXM2 with 16 GB memory.

#### 2.3 Additional ablation study

In section 3.6.5 of the main paper, we update our network by minimizing loss funtion:

$$\sum_{k=1}^{K} \mathcal{L}_{k}^{ie} + \alpha \mathcal{L}^{\mathcal{C}} + \mathcal{L}^{m} + \beta \mathcal{L}^{clus} - \gamma \mathcal{L}^{d}.$$
(1)

These two parameters  $\alpha$  and  $\beta$  help to discriminate source domains in latent space and ensure the clustering assumption [2], respectively. Figure 1 depicts the model performance with a diverse range of  $\alpha$  and  $\beta$  on " $\rightarrow$ **sv**" and " $\rightarrow$ **sv**" tasks. Following the result,  $\beta$  significantly affects both tasks in [0.5, 1] and less sensitive in [0, 0.5], while adjusting  $\alpha$  gives a stable performance in most cases. After searching their values in the range of [0, 1], we observe that our network achieves good performances with  $\alpha = 1.0$  and  $\beta = 0.1$ .



Figure 3: Effect of GAN relevant trade-off parameter  $\gamma$  to the performance on the  $\rightarrow$ **sv** (*blue*) and  $\rightarrow$ **mm** (*red*) tasks.

GAN components (i.e., generator and discriminator) play a vital role in our framework to diminish the data distribution discrepancy of the target domain and mixture of source domains in the latent space. We conduct experiments to study the effect of the parameter  $\gamma$  on the transferring performance. We experiment on the  $\rightarrow$ **sv** and  $\rightarrow$ **mm** task by varying  $\gamma$  in [0, 0.5]. As shown in Figure 3, on the  $\rightarrow$ **mm** task,  $\gamma$  quite significantly influences the performance, while for the case of  $\rightarrow$ **sv** task, it stably affects the performance. This totally depends on how distant the target domain and source domains in the latent space. We empirically find that  $\gamma = 0.1$  works satisfactorily for most of the cases.

#### 2.4 Convergence

We testify the convergence of our STEM with the test error on " $\rightarrow$ sv" and " $\rightarrow$ sv" tasks and compare our proposal with the state-of-the-art method named LtC-MSDA [3]. The comparison is presented in Figure 2. Following the result, our STEM enjoys faster and stable convergence in which the test error is significantly lower than LtC-MSDA method.

## References

- K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 770–778, 2016.
- [2] O. Chapelle and A. Zien. Semi-supervised classification by low density separation. In AISTATS, volume 2005, pages 57–64. Citeseer, 2005.
- [3] H. Wang, M. Xu, B. Ni, and W. Zhang. Learning to combine: Knowledge aggregation for multi-source domain adaptation. In *Computer Vision – ECCV*, 2020.