Supplementary Material: Boosting the Generalization Capability in Cross-Domain Few-Shot Learning via Noise-enhanced Supervised Autoencoder

1. Discriminability Analysis of Deep Features

Below we give the details of the definition of the Interclass correlation (ICC) [1, 2]. Let f be an feature extractor and $\mathcal{D} = \mathcal{D}_1 \cup \mathcal{D}_2 \cup \cdots \cup \mathcal{D}_K$ where $\mathcal{D}_j = \{(x_i, y_i) : y_i = j\}$ be a dataset with K classes. Let $\tilde{f}(x_i) := \frac{f(x_i)}{\|f(x_i)\|_2}$ be the normalized feature extracted by the feature extractor f. Then the center of the images features in jth class is defined as

$$\mu(f|\mathcal{D}_j) = |\mathcal{D}_j|^{-1} \sum_{x_i \in \mathcal{D}_j} \tilde{f}(x_i).$$
(1)

Then the classical intra-class and inter-class variation on the full dataset \mathcal{D} are defined respectively as

$$D_{\text{intra}}(f|\mathcal{D}) = \frac{1}{K} \sum_{k=1}^{K} \left\{ |\mathcal{D}_k|^{-1} \sum_{x_i \in \mathcal{D}_k} \|\tilde{f}(x_i) - \mu(\mathcal{D}_k)\|^2 \right\},\$$
$$D_{\text{inter}}(f|\mathcal{D}) = \frac{1}{K(K-1)} \sum_{k=1}^{K} \sum_{j \neq k} \|\mu(\mathcal{D}_j) - \mu(\mathcal{D}_k)\|^2.$$
(2)

The inter-class variation measures the average pairwise distances of class centers and the intra-class variation measures the within class variation of the image features. Following [2], the intra-class correlation (ICC) is defined as

$$ICC(f|\mathcal{D}) = D_{inter}(f|\mathcal{D})/D_{intra}(f|\mathcal{D}).$$
 (3)

Therefore, the ICC of a feature extractor f on dataset D is larger when the inter-class is larger and the intra-class is smaller. The ICC can therefore measures the discriminability of a feature extractor since a good feature embedding has smaller within-class variation and larger margin across classes.

In our experiment to study the discriminability of the feature extractors, we randomly sample 5 classes and compute the ICC based on the images from these 5 classes. We repeat these procedure for 600 times and use the average ICC as a measure for the discriminability of a feature extractor. The average ICC is computed using the same feature extractor on both the source and the target domains.

2. Model Architecture

In our proposed noise-enhanced supervised autoencoder, we use Conv4 and ResNet10 as the encoder structure and design the corresponding decoders. The decoder can be seen as a mirror mapping of the encoder which consist of deconvolutional blocks, with each block containing 2D transposed convolution operator and ReLU activation, which expand the dimension of the feature map. Before deconvolutional layers, we also add several fully connected layers that transform feature representations from encoder. Detailed architecture and layer specifications of the decoders are shown in Table 1.

Table 1. The architecture and layer specifications of the decoder modules of the Conv4 and ResNet10 based NSAE. Linear represents fully connected layer followed by ReLU activation. Deconv-ReLU represents a ConvTranspose2d-BatchNormalization-ReLU layer. Conv-Sigmoid represents a Conv2d-BatchNormalization-Sigmoid layer.

Module	Specifications							
Conv4	Linear, 1600×512							
	Linear, 512×1600							
	Reshape to $64 \times 5 \times 5$							
	2×2 Deconv-ReLU, 64 filters, stride 2, padding 0							
	2×2 Deconv-ReLU, 64 filters, stride 2, padding 0							
	2×2 Deconv-ReLU, 64 filters, stride 2, padding 0							
	2×2 Deconv-ReLU, 3 filters, stride 2, padding 0							
	3×3 Conv-Sigmoid, 3 filters, stride 1, padding 1							
ResNet10	Linear, 512×512							
	Linear, 512×6272							
	Reshape to $32 \times 14 \times 14$							
	2×2 Deconv-ReLU, 32 filters, stride 2, padding 0							
	2×2 Deconv-ReLU, 32 filters, stride 2, padding 0							
	2×2 Deconv-ReLU, 64 filters, stride 2, padding 0							
	2×2 Deconv-ReLU, 64 filters, stride 2, padding 0							
	3×3 Conv-Sigmoid, 3 filters, stride 1, padding 1							

3. Ablation Study Results

Table 2 gives the detailed experiment result for the 5way 5-shot ablation study on 8 datasets with various model architectures and loss functions. We use four kinds of combinations of the classification loss functions for pretraining and fine-tuning, i.e. CE+CE, BSR+CE, CE+D, and BSR+D. Meanwhile, we respectively test with Conv4

ons.									
Encoder	Method	ISIC	EuroSAT	CropDisease	ChestX	Car	CUB	Plantae	Places
Conv4	CE+CE NSAE(-) NSAE SAE SAE(*)	$\begin{array}{c} 46.04{\pm}0.62\\ 46.64{\pm}0.62\\ \textbf{46.65{\pm}0.63}\\ 46.36{\pm}0.61\\ 44.88{\pm}0.60\\ \end{array}$	$\begin{array}{c} 68.88 {\pm} 0.68 \\ 70.19 {\pm} 0.64 \\ \textbf{70.34 {\pm} 0.65} \\ 68.88 {\pm} 0.65 \\ 70.10 {\pm} 0.68 \end{array}$	83.47 ± 0.67 84.89 ± 0.60 85.22 ± 0.61 83.57 ± 0.62 83.38 ± 0.64	$\begin{array}{c} 24.81 \pm 0.43 \\ 24.86 \pm 0.41 \\ \textbf{25.02 \pm 0.42} \\ 24.88 \pm 0.41 \\ 25.00 \pm 0.41 \end{array}$	38.36 ± 0.58 39.88 ± 0.64 39.90 ± 0.62 37.94 ± 0.59 37.29 ± 0.60	52.94 ± 0.70 55.08 ± 0.70 55.35 ± 0.67 52.74 ± 0.65 53.64 ± 0.70	$\begin{array}{c} 45.55 \pm 0.71 \\ 46.98 \pm 0.69 \\ \textbf{47.18} \pm \textbf{0.75} \\ 44.79 \pm 0.67 \\ 44.60 \pm 0.65 \end{array}$	58.74 ± 0.74 59.54 ± 0.69 59.59 ± 0.69 58.34 ± 0.71 59.12 ± 0.72
	BSR+CE NSAE(-) NSAE SAE SAE(*)	$\begin{array}{r} 48.78 \pm 0.64 \\ 49.32 \pm 0.59 \\ \textbf{49.34} \pm \textbf{0.59} \\ 48.68 \pm 0.63 \\ 47.25 \pm 0.58 \end{array}$	$\begin{array}{c} 69.34 {\pm} 0.68 \\ 71.84 {\pm} 0.66 \\ 72.00 {\pm} 0.65 \\ \textbf{72.17 {\pm} 0.69} \\ 70.52 {\pm} 0.67 \end{array}$	85.88±0.61 86.86±0.59 86.87±0.59 86.23±0.60 85.14±0.62	$\begin{array}{c} 25.41 \pm 0.42 \\ 25.59 \pm 0.42 \\ \textbf{25.62 \pm 0.42} \\ 25.31 \pm 0.41 \\ 25.45 \pm 0.40 \end{array}$	$\begin{array}{r} 42.54{\pm}0.70\\ 42.54{\pm}0.64\\ \textbf{42.56}{\pm}\textbf{0.65}\\ 42.38{\pm}0.68\\ 40.81{\pm}0.63\end{array}$	60.16±0.73 60.00±0.71 60.18±0.72 60.10±0.76 59.54±0.74	50.85±0.78 50.00±0.73 49.48±0.73 48.29±0.74 47.70±0.70	$\begin{array}{c} 62.38 {\pm} 0.77 \\ 63.36 {\pm} 0.72 \\ \textbf{63.40} {\pm} \textbf{0.72} \\ 62.12 {\pm} 0.73 \\ 62.46 {\pm} 0.75 \end{array}$
	CE+D NSAE(-) NSAE SAE SAE(*)	50.54 ± 0.66 50.94 ± 0.63 50.95 ± 0.63 50.62 ± 0.65 49.94 ± 0.66	$76.16 \pm 0.64 77.70 \pm 0.70 77.77 \pm 0.64 74.92 \pm 0.64 77.24 \pm 0.70$	$\begin{array}{c} 89.65 {\pm} 0.55 \\ 90.06 {\pm} 0.52 \\ \textbf{90.11 {\pm} 0.52} \\ 88.11 {\pm} 0.59 \\ 88.31 {\pm} 0.56 \end{array}$	24.07±0.41 24.46±0.40 24.29±0.40 23.52±0.41 24.01±0.40	44.26±0.70 43.96±0.68 44.28±0.72 42.45±0.70 42.12±0.68	$58.61\pm0.8260.00\pm0.8259.90\pm0.7857.04\pm0.8157.15\pm0.82$	52.47±0.74 53.26±0.80 53.36±0.78 51.51±0.80 51.77±0.82	$\begin{array}{c} 61.81 {\pm} 0.74 \\ 62.26 {\pm} 0.73 \\ \textbf{62.42} {\pm} \textbf{0.72} \\ 61.08 {\pm} 0.74 \\ 61.40 {\pm} 0.78 \end{array}$
	BSR+D NSAE(-) NSAE SAE SAE(*)	50.06±0.65 49.95±0.67 49.98±0.67 49.77±0.68 49.46±0.68	$\begin{array}{c} 75.74 {\pm} 0.67 \\ 76.96 {\pm} 0.70 \\ \textbf{77.00 {\pm} 0.69} \\ 75.58 {\pm} 0.69 \\ 76.17 {\pm} 0.70 \end{array}$	87.71 ± 0.56 87.70 ± 0.57 87.71 ± 0.58 87.67 ± 0.57 86.50 ± 0.60	23.66±0.40 23.60±0.41 23.61±0.41 23.35±0.41 23.23±0.39	$\begin{array}{c} 41.11 {\pm} 0.77 \\ 41.05 {\pm} 0.72 \\ \textbf{41.80} {\pm} \textbf{0.72} \\ 41.75 {\pm} 0.75 \\ 40.16 {\pm} 0.70 \end{array}$	58.81 ± 0.81 58.32 ± 0.81 59.42 ± 0.82 58.34 ± 0.81 58.26 ± 0.79	51.35 ± 0.81 51.74 ± 0.84 51.80 ± 0.84 50.92 ± 0.81 50.70 ± 0.83	$\begin{array}{c} 60.51 {\pm} 0.80 \\ 60.34 {\pm} 0.81 \\ \textbf{60.92 {\pm} 0.85} \\ 60.25 {\pm} 0.81 \\ 60.86 {\pm} 0.84 \end{array}$
ResNet10	CE+CE NSAE(-) NSAE SAE SAE(*)	51.28 ± 0.62 53.52 ± 0.62 54.05 ± 0.63 52.28 ± 0.63 52.11 ± 0.65	82.51 ± 0.58 83.83 ± 0.56 83.96 ± 0.57 83.78 ± 0.55 83.50 ± 0.55	92.45±0.45 93.14±0.47 93.14±0.47 93.01±0.42 93.05±0.47	$\begin{array}{c} 26.50 {\pm} 0.43 \\ 26.69 {\pm} 0.44 \\ \textbf{27.10} {\pm} \textbf{0.44} \\ 26.05 {\pm} 0.45 \\ 26.37 {\pm} 0.45 \end{array}$	52.08 ± 0.72 53.49 ± 0.72 54.91 ± 0.74 53.54 ± 0.71 54.26 ± 0.70	64.14±0.77 67.60±0.73 68.51±0.76 64.27±0.75 66.62±0.75	59.27±0.70 59.70±0.74 59.80±0.74 59.87±0.73 59.62±0.75	$\begin{array}{c} 70.06 {\pm} 0.74 \\ 70.74 {\pm} 0.71 \\ \textbf{71.84 {\pm} 0.72} \\ 70.82 {\pm} 0.72 \\ 71.40 {\pm} 0.67 \end{array}$
	BSR+CE NSAE(-) NSAE SAE SAE(*)	54.42 ± 0.66 55.27 ± 0.62 55.88 ± 0.64 54.48 ± 0.65 54.73 ± 0.68	$\begin{array}{c} 80.89 {\pm} 0.61 \\ 84.19 {\pm} 0.54 \\ \textbf{84.33 {\pm} 0.55} \\ 84.10 {\pm} 0.54 \\ 83.90 {\pm} 0.55 \end{array}$	$\begin{array}{c} 92.17 {\pm} 0.45 \\ 92.92 {\pm} 0.47 \\ \textbf{93.31 {\pm} 0.42} \\ 92.92 {\pm} 0.47 \\ 93.02 {\pm} 0.46 \end{array}$	$\begin{array}{c} 26.84 {\pm} 0.44 \\ 27.23 {\pm} 0.45 \\ \textbf{27.30} {\pm} 0.42 \\ 27.20 {\pm} 0.45 \\ 26.74 {\pm} 0.43 \end{array}$	57.49 ± 0.72 58.35 ± 0.76 58.30 ± 0.75 58.30 ± 0.76 57.60 ± 0.71	69.38 ± 0.76 71.30 ± 0.75 71.92± 0.77 71.30 ± 0.75 71.50 ± 0.75	$\begin{array}{c} 61.07 \pm 0.76 \\ 61.92 \pm 0.76 \\ \textbf{62.18} \pm \textbf{0.77} \\ 61.92 \pm 0.76 \\ 62.20 \pm 0.78 \end{array}$	$71.09 \pm 0.68 \\71.76 \pm 0.74 \\73.17 \pm 0.72 \\71.76 \pm 0.74 \\72.99 \pm 0.67$
	CE+D NSAE(-) NSAE SAE SAE(*)	51.62 ± 0.66 54.31 ± 0.68 54.41 ± 0.63 52.64 ± 0.67 51.37 ± 0.66	$\begin{array}{c} 83.72 {\pm} 0.59 \\ 83.77 {\pm} 0.62 \\ \textbf{83.78} {\pm} \textbf{0.56} \\ 83.13 {\pm} 0.63 \\ 83.04 {\pm} 0.63 \end{array}$	$\begin{array}{c} 93.22 {\pm} 0.41 \\ 93.54 {\pm} 0.40 \\ \textbf{93.65 {\pm} 0.40} \\ 93.44 {\pm} 0.41 \\ 92.53 {\pm} 0.42 \end{array}$	26.23 ± 0.44 26.98 ± 0.44 27.25 ± 0.44 26.34 ± 0.44 26.34 ± 0.44	55.12 ± 0.76 55.67 ± 0.78 55.78 ± 0.73 55.44 ± 0.74 55.00 ± 0.73	$\begin{array}{c} 66.56 {\pm} 0.78 \\ 67.17 {\pm} 0.76 \\ \textbf{67.64 {\pm} 0.76} \\ 65.08 {\pm} 0.76 \\ 65.13 {\pm} 0.81 \end{array}$	59.09 ± 0.76 59.46 ± 0.75 59.74 ± 0.75 59.70 ± 0.78 59.46 ± 0.79	$72.81 \pm 0.73 \\72.90 \pm 0.72 \\73.25 \pm 0.73 \\73.13 \pm 0.71 \\73.20 \pm 0.67$
	BSR+D NSAE(-) NSAE SAE SAE(*)	52.85 ± 0.65 53.74 ± 0.67 54.42 ± 0.64 51.84 ± 0.65 53.08 ± 0.67	$\begin{array}{c} 80.13 \pm 0.65 \\ 82.19 \pm 0.64 \\ \textbf{82.79} \pm \textbf{0.62} \\ 80.02 \pm 0.69 \\ 81.77 \pm 0.64 \end{array}$	$\begin{array}{c} 91.20{\pm}0.48\\ 92.22{\pm}0.47\\ \textbf{92.45}{\pm}\textbf{0.45}\\ 91.95{\pm}0.45\\ 91.63{\pm}0.46\end{array}$	26.80±0.45 26.79±0.45 26.69±0.45 26.52±0.42 26.58±0.45	54.99 ± 0.74 55.90 ± 0.77 55.92 ± 0.72 55.90 ± 0.77 54.87 ± 0.78	$\begin{array}{c} 68.15 {\pm} 0.84 \\ 68.32 {\pm} 0.81 \\ \textbf{68.46 {\pm} 0.82} \\ 66.64 {\pm} 0.79 \\ 67.97 {\pm} 0.83 \end{array}$	58.26 ± 0.77 60.25 ± 0.77 60.40 ± 0.78 59.20 ± 0.80 58.61 ± 0.79	71.97 \pm 0.72 73.28 \pm 0.72 73.33 \pm 0.71 72.48 \pm 0.76 73.20 \pm 0.67

Table 2. Ablation study. The ablation study on the 5-way 5-shot support set on 8 datasets with various model architectures and loss functions.

and ResNet10 as backbone of feature encoder. In the table, CE+CE, BSR+CE, CE+D, and BSR+D denote using single feature extractor with different loss functions combinations. SAE denotes that we use auto-encoder but do not further feed in the reconstructed images for classification during the pre-training. SAE(*) denotes that we double the weight on the classification loss of original images as if the auto-encoder works perfectly that the reconstructed images are identical to original images. NSAE(-) denotes using our proposed pre-training strategy but using one-step fine-tuning.

4. Comparison with Handcrafted Noise

The reconstructed images during the pre-training stage can be viewed as noisy inputs to improve the model generalization capability. Can the model generalization capability be improved if we use images with handcrafted noise instead of reconstructed images? To answer this question, we compare the performance of our proposed method with that when images with handcrafted noise are used as data augmentation during pre-training. In our experiment, we consider the following four kinds of handcrafted noise: Gaussian, salt-pepper, Poisson, and speckle. We use the *skimage* package [3] in python to add handcrafted noise to source images. The parameter values for the noise generation are given in Table 3. We use BSR+CE loss combination and

Table 3. Handcrafted Noise Configuration. The parameters for adding noise to the images.

Noise type	Parameter values
Gaussian salt-pepper Poisson speckle	mode='gaussian', mean=0, var=0.1 mode='s&p', salt_vs_pepper=0.5 mode='poisson' mode='speckle', mean=0, var=0.05
	1 · · · ·

consider the following two settings during pre-training: (a) only use the encoder and feed in both source and hand-

crafted noisy images for classification; (b) add a decoder to (a) with reconstruction loss, though the reconstructed images are not used for classification. The rest of the hyperparameter values are the same as that given in Section 4.1 in the main paper. The results averaged over 8 datasets are shown in Fig. 1. It can be seen from the figures that



Figure 1. Ablation study with handcrafted noisy images. The two horizontal lines are baselines without noisy images.

- 1. regardless of the noise type, using auto-encoder scheme with reconstruction loss helps improve the generalization capability owing to regularization effect from decoder, which shows the advantage of our model on top of simple data augmentation;
- adding handcrafted noise may not improve the accuracy, but our design consistently improves the accuracy and surpasses all results with handcrafted noise.

References

- Bin Liu, Yue Cao, Yutong Lin, Qi Li, Zheng Zhang, Mingsheng Long, and Han Hu. Negative margin matters: Understanding margin in few-shot classification. In *European Conference on Computer Vision*, pages 438–455. Springer, 2020. 1
- [2] Sebastian Mika, Gunnar Ratsch, Jason Weston, Bernhard Scholkopf, and Klaus-Robert Mullers. Fisher discriminant analysis with kernels. In *Neural networks for signal processing IX: Proceedings of the 1999 IEEE signal processing society workshop (cat. no. 98th8468)*, pages 41–48. Ieee, 1999. 1

[3] Stefan Van der Walt, Johannes L Schönberger, Juan Nunez-Iglesias, François Boulogne, Joshua D Warner, Neil Yager, Emmanuelle Gouillart, and Tony Yu. scikitimage: image processing in python. *PeerJ*, 2:e453, 2014. 2