# Self-Supervised Domain Adaptation with Consistency Training

Liang Xiao<sup>1</sup>, Jiaolong Xu<sup>1</sup>, Dawei Zhao<sup>1</sup>, Zhiyu Wang<sup>2</sup>, Li Wang<sup>2</sup>, Yiming Nie<sup>1</sup> and Bin Dai<sup>1,2</sup>

<sup>1</sup>National Innovation Institute of Defense Technology, China (NIIDT) <sup>2</sup>National University of Defense Technology, China (NUDT)

xiaoliang.cs@gmail.com, jiaolong@gmail.com

January 2021

#### Self-supervised Domain Adaptation [Xu, 2019]



## Self-supervised Domain Adaptation [Xu, 2019]



Problem: No guarantee on the compatibility between the pretext task and the main task.

Liang Xiao et al. (NIIDT)

January 2021 2 / 13

Explicitly relate the representation of the transformed (specifically, rotated) image to the label of the main task by maximizing the mutual information:

Assuming  $p(y, \mathbf{x}, \tilde{\mathbf{x}}) = p(\mathbf{x}, \tilde{\mathbf{x}})p(y \mid \mathbf{x})$ 

$$-I(\tilde{\mathbf{x}}; y) = -\mathbb{E}_{p(\mathbf{x}, \tilde{\mathbf{x}}, y)} \left[ \log \frac{p(y \mid \tilde{\mathbf{x}})}{p(y)} \right]$$
(1)

$$= \mathbb{E}_{p(\mathbf{x},\tilde{\mathbf{x}})} \Big[ -\sum_{y} p(y \mid \mathbf{x}) \log \frac{p(y \mid \tilde{\mathbf{x}})}{p(y)} + \sum_{y} p(y \mid \mathbf{x}) \log \frac{p(y \mid \mathbf{x})}{p(y \mid \mathbf{x})} \Big]$$
(2)

$$= \mathbb{E}_{p(\mathbf{x},\tilde{\mathbf{x}})} \Big[ \sum_{y} p(y \mid \mathbf{x}) \log \frac{p(y \mid \mathbf{x})}{p(y \mid \tilde{\mathbf{x}})} - \sum_{y} p(y \mid \mathbf{x}) \log \frac{p(y \mid \mathbf{x})}{p(y)} \Big]$$
(3)

$$= \mathbb{E}_{p(\mathbf{x},\tilde{\mathbf{x}})} \Big[ \mathcal{D}_{\mathsf{KL}} \big( p(y \mid \mathbf{x}) \| p(y \mid \tilde{\mathbf{x}}) \big) - \mathcal{D}_{\mathsf{KL}} \big( p(y \mid \mathbf{x}) \| p(y) \big) \Big]$$
(4)

The first item is known as the Kullback-Leibler consistency loss as considered by VAT [Gidaris, 2018] and UDA [Xie, 2019] for semi-supervised learning:

$$\mathcal{L}_{c}(\boldsymbol{\theta_{e}},\boldsymbol{\theta_{m}}) = \mathbb{E}_{\mathbf{x}^{t} \in D_{t}} \mathbb{E}_{\tilde{\mathbf{x}}^{t} \in \tilde{D}_{t}} \left[ \mathcal{D}_{\mathsf{KL}}(\hat{p}(y^{t} \mid \mathbf{x}^{t}) \| p(y^{t} \mid \tilde{\mathbf{x}^{t}})) \right]$$

The first item is known as the Kullback-Leibler consistency loss as considered by VAT [Gidaris, 2018] and UDA [Xie, 2019] for semi-supervised learning:

$$\mathcal{L}_{c}(\boldsymbol{\theta_{e}},\boldsymbol{\theta_{m}}) = \mathbb{E}_{\mathbf{x}^{t} \in D_{t}} \mathbb{E}_{\tilde{\mathbf{x}}^{t} \in \tilde{D}_{t}} \left[ \mathcal{D}_{\mathsf{KL}}(\hat{p}(y^{t} \mid \mathbf{x}^{t}) \| p(y^{t} \mid \tilde{\mathbf{x}^{t}})) \right]$$

Assuming that p(y) is uniform, then the second term can be simplified as:

$$\mathbb{E}_{\mathbf{x}}[-\sum_{y} p(y \mid \mathbf{x}) \log p(y \mid \mathbf{x})],$$

The first item is known as the Kullback-Leibler consistency loss as considered by VAT [Gidaris, 2018] and UDA [Xie, 2019] for semi-supervised learning:

$$\mathcal{L}_{c}(\boldsymbol{\theta_{e}},\boldsymbol{\theta_{m}}) = \mathbb{E}_{\mathbf{x}^{t} \in D_{t}} \mathbb{E}_{\tilde{\mathbf{x}}^{t} \in \tilde{D}_{t}} \Big[ \mathcal{D}_{\mathsf{KL}}(\hat{p}(y^{t} \mid \mathbf{x}^{t}) \| p(y^{t} \mid \tilde{\mathbf{x}^{t}})) \Big]$$

Assuming that p(y) is uniform, then the second term can be simplified as:

$$\mathbb{E}_{\mathbf{x}}[-\sum_{y} p(y \mid \mathbf{x}) \log p(y \mid \mathbf{x})],$$

which suggests the following entropy minimization loss term:

$$\mathcal{L}_{e}(\boldsymbol{\theta}_{e}, \boldsymbol{\theta}_{m}) = \mathbb{E}_{\mathbf{x}^{t} \in D_{t}} \Big[ -\sum_{y^{t}} p(y^{t} \mid \mathbf{x}^{t}) \log p(y^{t} \mid \mathbf{x}^{t}) \Big],$$

Final objective function:

 $\min_{\theta_{e},\theta_{m},\theta_{p}} \mathcal{L}_{m}(\theta_{e},\theta_{m}) + \lambda_{p}\mathcal{L}_{p}(\theta_{e},\theta_{p}) + \lambda_{c}\mathcal{L}_{c}(\theta_{e},\theta_{m}) + \lambda_{e}\mathcal{L}_{e}(\theta_{e},\theta_{m})$ 



Table: Multi-source Domain Adaptation results on PACS (ResNet-18). Three domains are used as source datasets and the remaining one as target.

Method	Art.	Cartoon	Sketch	Photo	Avg.
SRC	74.7	72.4	60.1	92.9	75.0
Dial	87.3	85.5	66.8	97.0	84.2
DDiscovery	87.7	86.9	69.6	97.0	85.3
CDAN	85.7	88.1	73.1	97.2	86.0
CDAN+E	87.4	89.4	75.3	97.8	87.5
JiGen	84.9	81.1	79.1	97.9	85.7
Jigsaw	84.9	83.9	69.0	93.9	82.9
Rot	88.7	86.4	74.9	98.0	87.0
Ours	90.3	87.4	75.1	97.9	87.7

#### Table: Accuracy (%) on Office-31 dataset (ResNet-50).

Method	$A \rightarrow W$	$D \to W$	W  ightarrow D	$A \rightarrow D$	D  ightarrow A	$W \to A$	Avg.
ResNet-50	68.4±0.2	$96.7{\pm}0.1$	99.3±0.1	68.9±0.2	$62.5{\pm}0.3$	60.7±0.3	76.1
DAN	$80.5 {\pm} 0.4$	97.1±0.2	$99.6{\pm}0.1$	78.6±0.2	$63.6{\pm}0.3$	$62.8 {\pm} 0.2$	80.4
RTN	84.5±0.2	$96.8{\pm}0.1$	$99.4{\pm}0.1$	$77.5\pm0.3$	$66.2{\pm}0.2$	$64.8{\pm}0.3$	81.6
DANN	82.0±0.4	96.9±0.2	99.1±0.1	79.7±0.4	$68.2{\pm}0.4$	67.4±0.5	82.2
ADDA	$86.2{\pm}0.5$	96.2±0.3	98.4±0.3	77.8±0.3	$69.5{\pm}0.4$	$68.9{\pm}0.5$	82.9
JAN	85.4±0.3	97.4±0.2	99.8±0.2	84.7±0.3	$68.6{\pm}0.3$	$70.0{\pm}0.4$	84.3
CDAN	93.1±0.2	98.2±0.2	<b>100.0</b> ±0.0	89.8±0.3	$70.1{\pm}0.4$	68.0±0.4	86.6
CDAN+E	<b>94.1</b> ±0.1	<b>98.6</b> ±0.1	<b>100.0</b> ±0.0	<b>92.9</b> ±0.2	$71.0{\pm}0.3$	$69.3{\pm}0.3$	87.7
Jigsaw	86.9±0.8	<b>98.6</b> ±0.5	<b>100.0</b> ±0.0	82.9±1.0	62.9±1.2	$61.2{\pm}0.7$	82.1
Rot	$90.1{\pm}0.8$	$98.1{\pm}0.3$	<b>100.0</b> ±0.0	$88.6{\pm}0.7$	$65.1{\pm}0.8$	$65.0{\pm}0.6$	84.5
Ours	92.5±0.2	<b>98.7</b> ±0.3	<b>100.0</b> ±0.0	88.6±0.2	69.4±0.4	67.2±0.3	86.1

< A

#### Table: Accuracy (%) on Image-CLEF (ResNet-50).

Method	$I \rightarrow P$	$P \rightarrow I$	$I \rightarrow C$	$C \rightarrow I$	$C \rightarrow P$	$P \rightarrow C$	Avg.
ResNet-50	74.8±0.3	$83.9{\pm}0.1$	91.5±0.3	78.0±0.2	65.5±0.3	$91.2{\pm}0.3$	80.7
DAN	$74.5{\pm}0.4$	82.2±0.2	92.8±0.2	86.3±0.4	$69.2{\pm}0.4$	89.8±0.4	82.5
DANN	$75.0{\pm}0.6$	$86.0 {\pm} 0.3$	96.2±0.4	$87.0{\pm}0.5$	$74.3{\pm}0.5$	$91.5{\pm}0.6$	85.0
JAN	76.8±0.4	88.0±0.2	94.7±0.2	$89.5{\pm}0.3$	$74.2{\pm}0.3$	$91.7{\pm}0.3$	85.8
CDAN	76.7±0.3	90.6±0.3	97.0±0.4	90.5±0.4	<b>74.5</b> ±0.3	93.5±0.4	87.1
CDAN+E	77.7±0.3	90.7±0.2	<b>97.7</b> ±0.3	<b>91.3</b> ±0.3	74.2±0.2	94.3±0.3	87.7
Rot	$77.9{\pm}0.8$	$91.6{\pm}0.3$	$95.6{\pm}0.2$	$86.9{\pm}0.6$	$70.5{\pm}0.7$	$94.8{\pm}0.3$	84.2
Ours	<b>78.6</b> ±0.4	<b>92.5</b> ±0.1	96.1±0.3	88.9±0.2	73.9±0.7	<b>95.9</b> ±0.6	87.7

Image: Image:

#### Table: Accuracy (%) on Office-Home (ResNet-50).

Method	$ar \rightarrow cl$	$ar \rightarrow pr$	$ar \rightarrow rw$	$cl \rightarrow ar$	$cl \rightarrow pr$	$cl \rightarrow rw$	$pr \rightarrow ar$	$pr \rightarrow cl$	$pr \rightarrow rw$	$rw \rightarrow ar$	$rw \rightarrow cl$	$rw \rightarrow pr$	Avg.
ResNet-50	34.9	50.0	58.0	37.4	41.9	46.2	38.5	31.2	60.4	53.9	41.2	59.9	46.1
DAN	43.6	57.0	67.9	45.8	56.5	60.4	44.0	43.6	67.7	63.1	51.5	74.3	56.3
DANN	45.6	59.3	70.1	47.0	58.5	60.9	46.1	43.7	68.5	63.2	51.8	76.8	57.6
JAN	45.9	61.2	68.9	50.4	59.7	61.0	45.8	43.4	70.3	63.9	52.4	76.8	58.3
CDAN	49.0	69.3	74.5	54.4	66.0	68.4	55.6	48.3	75.9	68.4	55.4	80.5	63.8
CDAN+E	50.7	70.6	76.0	57.6	70.0	70.0	57.4	50.9	77.3	70.9	56.7	81.6	65.8
Rot	50.4	67.8	74.6	58.7	66.7	67.4	55.7	52.4	77.5	71.0	59.6	81.2	65.3
Ours	51.7	69.0	75.4	60.4	70.3	70.7	57.7	53.3	78.6	72.2	59.9	81.7	66.7

Image: A match a ma

Table: Convergence of different methods on task  $D \rightarrow A$ .



### Feature Visualization



Figure: The t-SNE visualization of deep features in PACS DA task (art painting is used as target domain). (a)-(e) are feature distribution visualization with category colors. (f)-(j) are feature distribution visualization with domain colors. Red and blue points represent samples of source and target domains, respectively.

- J. Xu, L. Xiao, and A. M. Lopez. Self-supervised domain adaptation for computer vision tasks. *IEEE Access.* vol. 7, 156694 156706, 2019.
- Q. Xie, Z. Dai, E. Hovy, M. Luong, and Q. V. Le. Unsupervised data augmentation for consistency training. *arXiv.* 2019.
  - S. Gidaris, P. Singh, and N. Komodakis. Unsupervised representation learning by predicting image rotations. *ICLR*. 2018.
  - T. Miyato, S. ichi Maeda, S. Ishii, and M. Koyama, Virtual adversarial training: a regularization method for supervised and semi-supervised learning, *IEEE T-PAMI*, 2018

Thank you Code available at: https://github.com/Jiaolong/ss-da-consistency.