

Exploring Adversarially Robust Training for Unsupervised Domain Adaptation ACCV 2022





Shao-Yuan Lo and Vishal M. Patel Johns Hopkins University

Adversarial Examples

$$x_{adv} = x + \delta$$

$$f(\boldsymbol{x}_{adv}) \neq y$$

Adversarial Examples

• Deep networks are **vulnerable** to adversarial examples.



Adversarial Defenses

• Image transformation: Remove perturbations from input images.

 $C(x_{adv}) \neq y.$ $C(T(x_{adv})) = y.$

• Adversarial training (AT): Enhance the robustness of networks itself.

$$\theta^* = \arg\min_{\theta} \mathbb{E}_{(x,y)\sim\mathbb{D}} \left[\max_{\delta\in\mathbb{S}} L(x+\delta,y;\theta) \right]$$

Madry et al. Towards deep learning models resistant to adversarial attacks. ICLR'18.

Unsupervised Domain Adaptation (UDA)

- Scenario: Training (source) data and test (target) data are from different domains (i.e. datasets).
 - Cause accuracy drop due to domain shift.
- Setting: Given a labeled source dataset and an unlabeled target dataset, learn a model for the target domain.



Challenges of AT for UDA

- Conventional AT requires ground-truth labels to generate adversarial examples and train models.
- However, UDA considers the scenario that label information is unavailable to a target domain.

Challenges of AT for UDA

- Can we develop an AT algorithm specifically for the UDA problem?
- How to improve the unlabeled data robustness via AT while learning domain-invariant features for UDA?



Conventional AT on UDA

• Natural Training

 $\mathcal{L}_{CE}(C(x_s), y_s) + \mathcal{L}_{DA}(x_s, x_t)$

• Conventional AT on UDA

 $\mathcal{L}_{CE}(C(\tilde{x}_s), y_s) + \mathcal{L}_{DA}(\tilde{x}_s, x_t)$

• Pseudo Labeling

 $\mathcal{L}_{CE}(C(x_s), y_s) + \mathcal{L}_{CE}(C(\tilde{x}_t), y'_t) + \mathcal{L}_{DA}(x_s, \tilde{x}_t)$

Self-supervised AT

• **Conventional AT**: Generate adversarial examples with ground-truth labels (e.g., L: cross-entropy loss)

 $x^{j+1} = \Pi_{\|\delta\|_p \le \epsilon} \left(x^j + \alpha \cdot sign(\nabla_{x^j} \mathcal{L}(C(x^j), y)) \right)$

• Self-supervised AT: Generate adversarial examples without groundtruth labels (e.g., L: L1 loss, L2 loss, KL divergence loss)

$$x_t^{j+1} = \Pi_{\|\delta\|_p \le \epsilon} \left(x_t^j + \alpha \cdot sign(\bigtriangledown_{x_t^j} \mathcal{L}(C(x_t^j), C(x_t))) \right)$$

Self-supervised AT

• Conventional AT (PGD-AT)

$$\min_{F,C} \mathbb{E}\left[\max_{\|\delta\|_p \le \epsilon} \mathcal{L}\big(C(\tilde{x}), y\big)\right]$$

• Self-supervised AT

 $\min_{F,C} \mathbb{E} \left[\max_{\|\delta\|_p \le \epsilon} \mathcal{L} (C(\tilde{x}_t), C(x_t)) \right]$

• Self-supervised AT on UDA

 $\mathcal{L}_{CE}(C(x_s), y_s) + \mathcal{L}_{KL}(C(\tilde{x}_t), C([x_t]_{sg})) + \mathcal{L}_{DA}(x_s, \tilde{x}_t))$

Self-supervised AT Results

- Dataset: VisDA-2017
- Attacks (white-box): FGSM [Goodfellow et al. 2015]

Training method	Clean	FGSM
Natural Training Conventional AT [26] Pseudo Labeling	$\begin{array}{ c c c c }\hline 73.2\\ 62.9 \ (-10.3)\\ 33.1 \ (-40.1) \end{array}$	$\begin{vmatrix} 21.2 \\ 27.1 \ (+5.9) \\ 27.1 \ (+5.9) \end{vmatrix}$
Self-Supervised AT-L1 Self-Supervised AT-L2 Self-Supervised AT-KL	56.2 (-17.0) 51.3 (-21.9) 67.1 (-6.1)	$\begin{vmatrix} 15.8 & (-5.4) \\ 26.0 & (+4.8) \\ 35.0 & (+13.8) \end{vmatrix}$

On the Effects of Clean and Adversarial Examples in Self-Supervise AT

– SSAT-s-t- \tilde{t} -1:

 $\mathcal{L}_{CE}(C(x_s), y_s) + \mathcal{L}_{KL}(C(\tilde{x}_t), C([x_t]_{sg})) + \mathcal{L}_{DA}(x_s, x_t).$

- SSAT-s-t-t-2:

 $\mathcal{L}_{CE}(C(x_s), y_s) + \mathcal{L}_{KL}(C(\tilde{x}_t), C([x_t]_{sg})) + \mathcal{L}_{DA}(x_s, x_t) + \mathcal{L}_{DA}(x_s, \tilde{x}_t).$

– SSAT-s-ŝ-t-ť-1:

 $\mathcal{L}_{CE}(C(x_s), y_s) + \mathcal{L}_{KL}(C(\tilde{x}_t), C([x_t]_{sg}))$ $+ \mathcal{L}_{CE}(C(\tilde{x}_s), y_s) + \mathcal{L}_{DA}(x_s, x_t) + \mathcal{L}_{DA}(\tilde{x}_s, \tilde{x}_t).$

– SSAT-s- \tilde{s} -t- \tilde{t} -2:

 $\mathcal{L}_{CE}(C(x_s), y_s) + \mathcal{L}_{KL}(C(\tilde{x}_t), C([x_t]_{sg}))$ $+ \mathcal{L}_{CE}(C(\tilde{x}_s), y_s) + \mathcal{L}_{DA}(x_s, \tilde{x}_t) + \mathcal{L}_{DA}(\tilde{x}_s, x_t).$

- SSAT-s-s-'t-t-3:

```
\mathcal{L}_{CE}(C(x_s), y_s) + \mathcal{L}_{KL}(C(\tilde{x}_t), C([x_t]_{sg})) + \mathcal{L}_{CE}(C(\tilde{x}_s), y_s) + \mathcal{L}_{DA}(x_s, x_t) + \mathcal{L}_{DA}(x_s, \tilde{x}_t) + \mathcal{L}_{DA}(\tilde{x}_s, x_t) + \mathcal{L}_{DA}(\tilde{x}_s, \tilde{x}_t).
```

On the Effects of Clean and Adversarial Examples in Self-Supervise AT

- Dataset: VisDA-2017
- Attacks (white-box): FGSM [Goodfellow et al. 2015]

Training method	$ x_s $	\tilde{x}_s	x_t	\tilde{x}_t	(x_s, x_t)	(x_s, \tilde{x}_t)	(\tilde{x}_s, x_t)	$(\tilde{x}_s, \tilde{x}_t)$	Clean	FGSM
Natural Training Conventional AT [26] SS-AT-KL	•	•	•	•	•	•	•		73.2 62.9 67.1	$21.2 \\ 27.1 \\ 35.0$
$\begin{array}{l} \mathrm{SS}\text{-}\mathrm{AT}\text{-}\mathrm{s}\text{-}\mathrm{t}\tilde{\text{-}}1\\ \mathrm{SS}\text{-}\mathrm{AT}\text{-}\mathrm{s}\text{-}\mathrm{t}\tilde{\text{-}}2\\ \mathrm{SS}\text{-}\mathrm{AT}\text{-}\mathrm{s}\text{-}\tilde{\mathrm{s}}\text{-}\mathrm{t}\tilde{\text{-}}1\\ \mathrm{SS}\text{-}\mathrm{AT}\text{-}\mathrm{s}\text{-}\tilde{\mathrm{s}}\text{-}\mathrm{t}\tilde{\text{-}}1\\ \mathrm{SS}\text{-}\mathrm{AT}\text{-}\mathrm{s}\text{-}\tilde{\mathrm{s}}\text{-}\mathrm{t}\tilde{\text{-}}1\\ \mathrm{SS}\text{-}\mathrm{AT}\text{-}\mathrm{s}\text{-}\tilde{\mathrm{s}}\text{-}\mathrm{t}\tilde{\text{-}}1\\ \mathrm{SS}\text{-}\mathrm{AT}\text{-}\mathrm{s}\tilde{\text{-}}\mathrm{s}^{-}\mathrm{t}\tilde{\text{-}}1\\ \end{array}$	• • • •	•	• • •	•	• • •	• • •	•	•	$\begin{vmatrix} 67.3 \\ 73.0 \\ 63.4 \\ 62.8 \\ 61.3 \end{vmatrix}$	$27.5 \\ 39.4 \\ 41.6 \\ 42.3 \\ 41.6$

On the Effects of Batch Normalization in Self-Supervise AT

- Dataset: VisDA-2017
- Attacks (white-box): FGSM [Goodfellow et al. 2015]

Method	Mini-batches	Clean	FGSM
$\begin{array}{l} \text{Batch-st-}\tilde{t}\\ \text{Batch-s-}t\tilde{t}\\ \text{Batch-s-}t\tilde{t}\\ \text{Batch-st}\tilde{t} \end{array}$	$egin{aligned} & [x_s, x_t], [ilde{x}_t] \ & [x_s], [x_t, ilde{x}_t] \ & [x_s], [x_t], [ilde{x}_t] \ & [x_s, x_t, ilde{x}_t] \end{aligned}$	73.0 68.2 68.2 69.0	$39.4 \\ 37.0 \\ 35.5 \\ 41.4$

Results

• Comparison with baselines on multiple datasets and attacks

Dataset	Training method	Clean	FGSM	PGD	MI-FGSM	MultAdv	Black-box
VisDA-2017 [29]	Natural TrainingPGD-AT [26]TRADES [42]ARTUDA (ours)	$73.2 \\ 60.5 \\ 64.0 \\ 65.5$	21.2 34.6 42.1 52.5	0.9 21.3 29.7 44.3	0.5 22.7 31.2 45.0	0.3 7.8 16.4 27.3	58.3 59.1 62.6 65.1
Office-31 D \rightarrow W[31]	Natural Training PGD-AT [26] TRADES [42] ARTUDA (ours)	98.0 95.3 88.4 96.5	52.7 91.8 85.3 95.2	0.9 68.2 66.4 92.5	0.6 66.5 67.0 92.5	0.1 31.4 28.2 77.1	95.0 95.3 88.2 96.5
Office-Home Ar \rightarrow Cl [36]	Natural TrainingPGD-AT [26]TRADES [42]ARTUDA (ours)	$54.5 \\ 42.5 \\ 49.3 \\ 54.0$	26.4 38.8 45.1 49. 5	4.7 36.0 41.6 41.3	2.8 35.8 41.6 39.9	2.0 21.7 22.5 21.6	53.1 43.0 49.4 53.9

Results

• Comparison with baselines on multiple UDA algorithms

UDA algorithm \rightarrow Training method \downarrow	Clean	DANN [8] PGD	Drop	Clean	JAN [24] PGD	Drop	Clean	CDAN [23] PGD	Drop
Natural Training	73.2	0.0	-73.2	64.2	0.0	-64.2	75.1	0.0	-75.1
PGD-AT [26]	60.5	13.3	-47.2	47.7	5.8	-41.9	58.2	11.7	-46.5
TRADES $[42]$	64.0	19.4	-44.6	48.7	8.5	-40.2	64.6	15.7	-48.9
Robust PT [2]	65.8	38.2	-27.6	55.1	32.2	-22.9	68.0	41.7	-26.3
RFA [2]	65.3	34.1	-31.2	63.0	32.8	-30.2	72.0	43.5	-28.5
ARTUDA (ours)	65.5	40.7	-24.8	58.5	34.4	-24.1	68.0	43.6	-24.4

Feature Analysis

• Mean square differences between the features of clean images and the features of adversarial examples.



• t-SNE visualization



Conclusion

- We provide a systematic study into various AT methods that are suitable for UDA.
- We propose ARTUDA, a new AT method specifically designed for UDA. To the best of our knowledge, it is the first AT-based UDA defense method that is robust against white-box attacks.
- Comprehensive experiments show that ARTUDA consistently improves UDA models' adversarial robustness under multiple attacks and datasets.