

In this paper, we give a detailed overview of adversarial training in image classification, which has attracted extensive attention of researchers. We believe that this survey can provide up-to-date findings and developments happening on adversarial training. Notably, we carefully review and analyze adversarial training with a novel taxonomy, uniquely discuss the poor generalization ability from different perspectives, and present future research directions.

## Definition of Adversarial Training

Adversarial training [Madry *et al.*, 2018] is a min-max problem: the inner maximization problem is finding the worst-case samples for the given model, and the outer minimization problem is to train a model robust to adversarial examples. Formally it is defined as:

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \max_{\delta \in B(x,\varepsilon)} \mathcal{L}_{ce}(\theta, x + \delta, y) \right], \quad (1)$$

where  $(x, y) \sim \mathcal{D}$  represents training data sampled from distribution  $\mathcal{D}$  and  $B(x, \varepsilon)$  is the allowed perturbation set, expressed as  $B(x, \varepsilon) := \{x + \delta \in \mathbb{R}^{h \times w \times c} \mid \|\delta\|_p \leq \varepsilon\}$ .

One popular method for solving this optimization problem is using Projected Gradient Descent (PGD) [Madry *et al.*, 2018], which is shown below:

$$x^{t+1} := \text{Proj}_{x+B(x,\varepsilon)} (x^t + \alpha \text{sign}(\nabla_{x^t} \mathcal{L}_{ce}(\theta, x^t, y))), \quad (2)$$

where  $t$  is the current step and  $\alpha$  is the step size. Madry *et al.* (2018) investigated the inner maximization problem from the landscape of adversarial examples and gave both theoretical and empirical proofs of local maxima's tractability with PGD. Thus, PGD-based adversarial training became a critical benchmark and is regarded as the standard way to do adversarial training in practice.

## Taxonomy of Adversarial Training for Adversarial Robustness

A summary of selected adversarial training methods is provided in Table 1.

- **Adversarial Regularization** uses both clean and adversarial data for training.
- **Curriculum-based Adversarial Training** utilizes weak adversarial examples.
- **Ensemble Adversarial Training** augments training data with adversarial examples generated from multiple target models.
- **Adversarial Training with Adaptive  $\epsilon$**  generates proper adversarial examples with the intrinsic adversarial robustness of individual samples taken into consideration.
- **Adversarial Training with Semi/Unsupervised Learning** utilizes larger datasets for adversarial training.
- **Efficient Adversarial Training** promotes the efficiency of adversarial training.
- **Other Variants**

## Generalization Problems in Adversarial Training

### Standard Generalization

Adversarial training is observed that hurts standard accuracy badly [Madry *et al.*, 2018]. **On the one hand**, one popular viewpoint is the trade-off between adversarial robustness and standard accuracy. **On the other hand**, it is confirmed the existence of adversarial examples on the manifold of natural data, adversarial robustness on which is equivalent to generalization.

### Adversarially Robust Generalization

The phenomenon that adversarially trained models do not perform well on adversarially perturbed test data is firstly observed in [Madry *et al.*, 2018]. In other words, there is a large gap between the training accuracy and test accuracy on adversarial data.

### Generalization on Unseen Attacks

Adversarially trained models, which are robust to a specific attack, *e.g.*,  $l_{\infty}$  adversarial examples, can be circumvented easily by different types of attacks, *e.g.*, other  $l_p$  norms, or larger  $\epsilon$ , or different target models [Kang *et al.*, 2019].

## Future Direction

- **Min-Max Optimization in Adversarial Training.** The robustness of adversarially trained models is not guaranteed [Kang *et al.*, 2019].
- **Overfitting in Adversarial Training.** The generalization gap between adversarial training accuracy and testing accuracy is very large.
- **Beyond Adversarial Training.** Though many theories have been proposed for improving adversarial training, it is undeniable that these improvements are less effective than claimed [Pang *et al.*, 2021].
- **Adversarial Training in Other Domains.** Adversarial training has been successfully applied to texts [Miyato *et al.*, 2017], graphs [Dai *et al.*, 2019], audios [Pandey and Wang, 2018], and reinforcement learning [Pattanaik *et al.*, 2017].

Taxonomy	Publication	Model Architecture	Attack	$\epsilon$ Dataset	Accuracy
Adversarial Regularization	[Qin <i>et al.</i> , 2019]	ResNet-152	PGD <sub>50</sub>	4/255 ImageNet	47.00%
	[Zhang <i>et al.</i> , 2019b]	Wide ResNet	CW <sub>10</sub>	0.031/1 CIFAR-10	84.03%
	[Wang <i>et al.</i> , 2020]	ResNet-18	PGD <sub>20</sub>	8/255 CIFAR-10	55.45%
	[Kannan <i>et al.</i> , 2018]	InceptionV3	PGD <sub>10</sub>	16/255 ImageNet	27.90%
	[Mao <i>et al.</i> , 2019]	Wide ResNet	PGD <sub>20</sub>	8/255 CIFAR-10	50.03%
Curriculum	[Zhang <i>et al.</i> , 2020]	Wide ResNet	PGD <sub>20</sub>	16/255 CIFAR-10	49.86%
	[Cai <i>et al.</i> , 2018]	DenseNet-161	PGD <sub>7</sub>	8/255 CIFAR-10	69.27%
	[Wang <i>et al.</i> , 2019]	8-Layer ConvNet	PGD <sub>20</sub>	8/255 CIFAR-10	42.40%
Ensemble	[Pang <i>et al.</i> , 2019]	Wide ResNet	PGD <sub>10</sub>	0.005 CIFAR-100	32.10%
	[Kariyappa and Qureshi, 2019]	ResNet-20	PGD <sub>30</sub>	0.09/1 CIFAR-10	46.30%
	[Yang <i>et al.</i> , 2020]	ResNet-20	PGD <sub>20</sub>	0.01/1 CIFAR-10	52.40%
Adaptive $\epsilon$	[Balaji <i>et al.</i> , 2019]	ResNet-152	PGD <sub>1000</sub>	8/255 ImageNet	59.28%
	[Ding <i>et al.</i> , 2020]	Wide ResNet	PGD <sub>100</sub>	8/255 CIFAR-10	47.18%
	[Cheng <i>et al.</i> , 2020]	Wide ResNet	PGD <sub>20</sub>	8/255 CIFAR-10	73.38%
Semi-Unsupervised	[Alayrac <i>et al.</i> , 2019]	Wide ResNet	FGSM	8/255 CIFAR-10	62.18%
	[Carmon <i>et al.</i> , 2019]	Wide ResNet	PGD <sub>10</sub>	8/255 CIFAR-10	63.10%
	[Zhai <i>et al.</i> , 2019]	Customized ResNet	PGD <sub>7</sub>	8/255 CIFAR-10	42.48%
	[Hendrycks <i>et al.</i> , 2019]	Wide ResNet	PGD <sub>20</sub>	0.3/1 ImageNet	50.40%
Efficient	[Shafahi <i>et al.</i> , 2019]	Wide ResNet	PGD <sub>100</sub>	8/255 CIFAR-10	46.19%
	[Wong <i>et al.</i> , 2020]	ResNet-50	PGD <sub>40</sub>	2/255 ImageNet	43.43%
	[Andriushchenko and Flammarion, 2020]	ResNet-50	PGD <sub>50</sub>	2/255 ImageNet	41.40%
	[Kim <i>et al.</i> , 2021]	PreActResNet-18	FGSM	8/255 CIFAR-10	50.50%
	[S. and Babu, 2020]	Wide ResNet	PGD <sub>40</sub>	8/255 MNIST	88.51%
	[Song <i>et al.</i> , 2019]	Customized ConvNet	PGD <sub>20</sub>	4/255 CIFAR-10	58.10%
	[Vivek and Babu, 2020]	Wide ResNet	PGD <sub>100</sub>	0.3/1 MNIST	90.03%
	[Huang <i>et al.</i> , 2020]	Wide ResNet	PGD <sub>20</sub>	8/255 CIFAR-10	45.80%
	[Zhang <i>et al.</i> , 2019a]	Wide ResNet	PGD <sub>20</sub>	8/255 CIFAR-10	47.98%
	Others	[Dong <i>et al.</i> , 2020]	Wide ResNet	PGD <sub>20</sub>	8/255 CIFAR-100
[Wang and Zhang, 2019]		Wide ResNet	CW <sub>200</sub>	4/255 CIFAR-10	60.30%
[Zhang and Wang, 2019]		Wide ResNet	PGD <sub>20</sub>	8/255 CIFAR-100	47.20%
[Pang <i>et al.</i> , 2020]		Wide ResNet	PGD <sub>500</sub>	8/255 CIFAR-10	60.75%
[Lee <i>et al.</i> , 2020]		PreActResNet-18	PGD <sub>20</sub>	8/255 Tiny ImageNet	20.31%
Benchmark	[Madry <i>et al.</i> , 2018]	ResNet-50	PGD <sub>20</sub>	8/255 CIFAR-10	45.80%

**Table 1.** A summary of experimental results for various adversarial training methods. All the attacks are under  $l_{\infty}$  norm.

## Conclusion

In this paper, we present recent advances of adversarial training methods for adversarial robustness. To our best knowledge, for the first time, we review adversarial training with a novel taxonomy and discuss the generalization problem in adversarial training. We also summarize the benchmarks and provide performance comparisons of different methods. Despite extensive efforts, the vulnerability of deep learning models to adversarial examples has not been completely solved by adversarial training and several open problems remain yet to solve.

## Acknowledgements

This paper is supported by 1) Singapore Ministry of Education Academic Research Fund Tier 1 RG128/18, Tier 1 RG115/19, Tier 1 RT07/19, Tier 1 RT01/19, Tier 1 RG24/20, and Tier 2 MOE2019-T2-1-176, 2) NTU-WASP Joint Project, 3) Singapore NRF National Satellite of Excellence, Design Science and Technology for Secure Critical Infrastructure NSoE DeST-SCI2019-0012, 4) AI Singapore (AISG) 100 Experiments (100E) programme, and 5) NTU Project for Large Vertical Take-Off & Landing (VTOL) Research Platform. Bihan Wen was supported in part by the National Research Foundation (NRF), Singapore, through the Singapore Cybersecurity Consortium (SGCSC) Grant Office, under SGCSC\_Grant\_2019-S01. Qian Wang's work was partially supported by the NSFC under Grants U20B2049 and 61822207. We would like to express our very great appreciation to Prof Bo Li and Mr. Qi Bi for their valuable and constructive suggestions of this research work.