

Increasing-Margin Adversarial (IMA) Training to Improve Adversarial Robustness of Neural Networks: Supplementary Information

Anonymous CVPR submission

Paper ID 5479

This document contains supplementary information for the manuscript ‘Increasing-Margin Adversarial (IMA) Training to Improve Adversarial Robustness of Neural Networks’. Appendix A further discusses the necessity of adversarial robustness. Appendix B discusses the effect of parameter β in the IMA. Appendix C discusses the effect of ε_{max} in the IMA. Appendix D discusses the trade-off between robustness and accuracy. Appendix E further discusses the equilibrium state. Appendix F shows some examples of the clean and noisy images used in the experiment of the main paper. Appendix G shows the training details in the experiments in the main paper. Appendix H shows some additional results, including black-box attack evaluation and robust overfitting situation of IMA.

A. Adversarial Robustness is Essential

It seems that adversarial noises are created by algorithms (e.g. PGD) and therefore it is only a security issue caused by hackers. In fact, random imaging noises could also be “adversarial” noises leading to wrong classifications. For the COVID-19 application, we did an additional test and found out that 2.75% of the noisy samples with uniform white noises on the level of 0.05, can cause the model “ce” to make wrong classifications. 2.75% is not a negligible number for this application. We note that CT imaging noises can be better described by Poisson distribution [10]. However, without the hardware parameters of the CT machine, it is impossible to simulate Poisson noises. Nevertheless, adversarial robustness should be the built-in property of a model for this application, and all of the DNN models in the previous COVID-19 studies [8] should be checked and enhanced for adversarial robustness before deploying those models in clinics and hospitals.

B. the Effect of β in the IMA

We evaluated the effect of β in the IMA method on the Fashion-MNIST dataset, and the other settings are the same as those in the corresponding Section in the main paper (e.g. white-box 100-PGD attack on the test set). The results are reported in Table 1.

noise	0	0.5	1	2	3
$\beta=0.1$	89.46	79.95	67.59	43.29	28.27
$\beta=0.3$	89.00	79.87	67.96	45.82	30.90
$\beta=0.5$	88.98	80.73	69.64	47.62	32.42
$\beta=0.7$	87.78	80.86	70.73	47.64	32.64
$\beta=0.9$	86.82	80.74	72.16	50.96	33.51

Table 1. Results on Fashion-MNIST (L2 norm-defined noise level)

It can be clearly seen that smaller β leads to higher accuracy on clean data (noise level = 0) and larger β leads to higher accuracy on noisy data.

The trade-off between robustness and accuracy is highly nonlinear, as shown in Table 2: a small decrease in accuracy on clean data can result in a large increase in accuracy on noisy data.

noise	0	0.5	1	2	3
$\beta = 0.1$	89.46	79.95	67.59	43.29	28.27
$\beta = 0.5$	88.98	80.73	69.64	47.62	32.42
difference	0.48	0.78	2.05	4.33	4.15

Table 2. Accuracy differences caused by different values of β

The nonlinear trade-off between robustness and accuracy makes it difficult to directly compare two methods, as different methods make different trade-offs between robustness and accuracy. The average accuracy is clearly not a good measure of performance on robustness. Adjusting the parameter β of IMA is not a computationally efficient approach to make a trade-off between robustness and accuracy, because for every possible value of β in the range of 0 to 1, the user has to train a model from scratch through many epochs. In Appendix C, we show that the user of IMA can make such a trade-off much more efficiently by “adjusting” the parameter ε_{max} of IMA.

C. the Effect of ε_{max} in the IMA

In this appendix, we show how ε_{max} in Algorithm 2 affects the performance of IMA, and how to choose its value. ε_{max} is the allowed maximum margin of the samples in IMA. And, it is the maximum noise level for training, which exists in almost every adversarial training method.

C.1. the Effect of ε_{max}

From Table 3, Table 4 and Table 5, we can see that, in general, a larger ε_{max} leads to a higher accuracy on noisy data. However, unlike vanilla adversarial training (see Table 7) and MMA (see Table 6), IMA with a larger ε_{max} does not have a significant decrease in the accuracy on clean data. This is because adversarial samples generated by IMA are of good quality and are not likely to cross over the true decision boundary and mislead the model (note: if an adversarial sample goes across the true decision boundary, then its true class label is changed, which means the class label of this adversarial sample is wrong during training). We can also see that, if ε_{max} is too small, the accuracy on noisy data with large noises decreases, which is because the sample margins are not allowed to expand sufficiently. After the ε_{max} reaches a specific noise level (4.0 for Fashion-MNIST, 1.5 for SVHN), the accuracy on noisy data will no longer change significantly as ε_{max} increases, which is because the sample margins have been expanded sufficiently large. Thus, for IMA, the user only needs to avoid setting ε_{max} to a too-small value. A very large ε_{max} is not likely to significantly harm the accuracy on clean data. As a comparison (see Table 7 and Table 6), for MMA, a large ε_{max} will lead to a significant decrease in accuracy on clean data, which is not good for real applications.

noise	0	0.5	1	2	3
$\varepsilon_{max} = 3.0$	89.40	80.98	70.34	48.07	28.34
$\varepsilon_{max} = 4.0$	88.81	80.23	69.07	46.99	31.47
$\varepsilon_{max} = 5.0$	88.98	80.73	69.64	47.62	32.42
$\varepsilon_{max} = 6.0$	88.90	80.68	69.64	47.74	32.92
$\varepsilon_{max} = 7.0$	88.57	79.55	67.74	46.80	33.27

Table 3. Effect of ε_{max} in IMA on Fashion-MNIST test set (L2 norm-defined noise level, $\Delta\varepsilon = 5/60$).

noise	0	0.1	0.25	0.5	1.0
$\varepsilon_{max} = 1.0$	90.53	86.16	77.76	59.63	27.20
$\varepsilon_{max} = 1.5$	89.75	85.50	77.29	60.34	30.30
$\varepsilon_{max} = 2.0$	89.69	85.01	76.82	60.50	31.90
$\varepsilon_{max} = 2.5$	89.96	85.51	77.50	60.53	31.44
$\varepsilon_{max} = 3.0$	89.97	85.80	77.31	60.49	31.57
$\varepsilon_{max} = 3.5$	90.00	85.57	77.49	60.38	31.37
$\varepsilon_{max} = 4.0$	89.42	85.16	77.07	60.41	31.58
$\varepsilon_{max} = 4.5$	89.79	85.45	77.20	60.31	31.89

Table 4. Effect of ε_{max} in IMA on SVHN test set (L2 norm-defined noise level, $\Delta\varepsilon = 2/60$).

noise	0	0.5	1.0	1.5
$\varepsilon_{max} = 1.0$	89.35	65.30	37.08	17.73
$\varepsilon_{max} = 2.0$	88.21	66.28	40.86	22.50
$\varepsilon_{max} = 3.0$	88.07	66.13	41.07	23.09

Table 5. Effect of ε_{max} in IMA on CIFAR10 test set (L2 norm-defined noise level, $\Delta\varepsilon = 3/100$).

noise	0	0.5	1.0	1.5
$\varepsilon_{max} = 1.0$	88.02	66.19	37.80	15.61
$\varepsilon_{max} = 2.0$	84.22	65.98	46.11	28.56
$\varepsilon_{max} = 3.0$	82.11	64.25	47.61	33.48

Table 6. Effect of ε_{max} in MMA on CIFAR10 test set (L2 norm-defined noise level).

noise	0	0.5	1.0	1.5
$\varepsilon_{max} = 1.0$	83.25	66.69	46.11	26.16
$\varepsilon_{max} = 2.0$	71.05	59.80	47.92	35.39

Table 7. Effect of ε_{max} in vanilla adversarial training on CIFAR10 test set (L2 norm-defined noise level).

C.2. How to Choose ε_{max}

Noisy samples with adversarial noises can be correctly classified by humans but may be incorrectly classified by neural networks, which is basically the definition of adversarial noises. By following this definition, we choose the allowed maximum margin by looking at the noisy samples: if the noise magnitude is too large such that we barely can recognize the objects in the images, then this magnitude is chosen as the allowed maximum margin. This is how we choose the allowed maximum margin in the experiments on the datasets.

Next, we provide a strategy for the user of IMA to refine the choice of ε_{max} by making a trade-off between robustness and accuracy. As shown in Fig. 1, the trade-off between clean-accuracy (i.e., accuracy on clean data) and robustness (i.e., accuracy on noisy data) can be observed during training by using our IMA method.

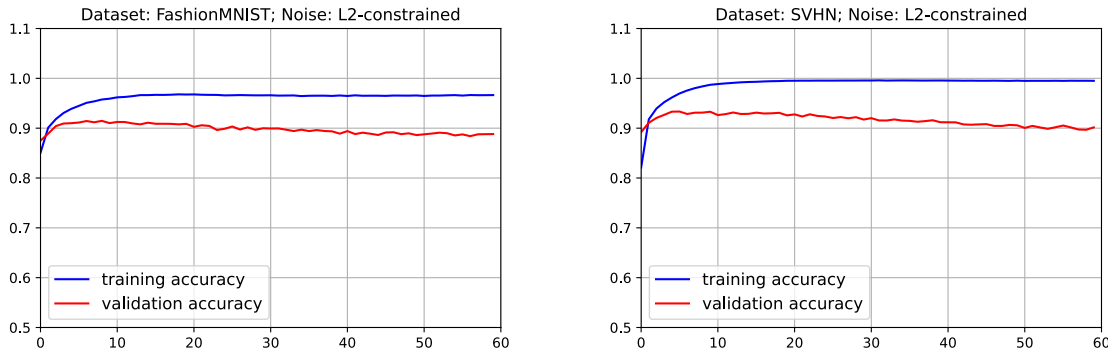


Figure 1. The training and validation curves (accuracy vs epoch) on different datasets obtained by using our IMA method. The accuracy scores are measured on clean data.

During IMA training, the (estimated) margin of a sample x_1 in class-1 is initialized to be close to 0 and it keeps increasing as if the radius of a ball increases during the training process. The sample x_1 is at the center of the ball, and the radius of the ball is the current margin of the sample. When the ball of x_1 collides with the ball of another sample x_2 from a different class-2, then a local decision boundary is formed, and the two balls stop expanding.

If there are enough training data samples, a sample x in some class will eventually meet with its counterpart in another class somewhere in the middle of the two classes, which forms a robust decision boundary. In practice, the amount of training data is never enough to cover the input space, and therefore the margin of a sample could be overestimated because of missing the counterparts that are needed to stop the expansion of the margin. If the margins of many samples are overestimated, then the balls of these samples may penetrate the true decision boundary and cause lower classification accuracy on the validation set of clean samples, **which explains the existence of the trade-off between robustness and clean-accuracy from the perspective of sample margins.**

The above analysis is confirmed by the trend of the training and validation accuracy curves in Fig. 1: after some epochs, the training accuracy curve becomes stable, and the validation accuracy curve starts to decrease, which indicates margin overestimation occurs. **For our IMA method, the cause of margin overestimation is the lack of data in high dimensional space.** In the 2D Moons dataset, there are enough training samples, and therefore, the decision boundary of our IMA method is almost in the “middle” between the two classes, i.e., no margin overestimation.

Using the validation accuracy curve, the user of IMA can choose the allowed maximum margin such that validation accuracy is above a pre-defined threshold that is set by the user to meet some application requirements. Next, we demonstrate this approach on the Fashion-MNIST dataset. We note that during IMA training, the margins of the samples are gradually increasing. Immediately after M epochs, the possible maximum of the training sample margins is $M \times \Delta\epsilon$ where $\Delta\epsilon$ is the margin expansion step size in Algorithm 2. Thus, the maximum margin is $\epsilon_{max} = M \times \Delta\epsilon$ for the model trained for M epochs. Since we do not know the threshold on validation accuracy, which a user may be interested in on the dataset, we selected six models trained by IMA after 10, 20, 30, 40, 50, and 60 epochs, and we evaluated the robustness of the six models on the test set by using the 100-PGD attack. Basically, we re-analyzed the models trained through 60 epochs. The results are reported in Table 8, which reveals the effect of ϵ_{max} on robustness.

noise	0	0.5	1	2	3	0 (Validation)
$\epsilon_{max} = 10\Delta\epsilon$	91.59	76.43	52.74	11.46	0.28	91.43
$\epsilon_{max} = 20\Delta\epsilon$	90.91	80.91	66.72	34.08	8.27	90.87
$\epsilon_{max} = 30\Delta\epsilon$	89.89	81.34	69.85	44.16	20.58	89.93
$\epsilon_{max} = 40\Delta\epsilon$	89.61	81.32	70.15	46.65	26.66	89.38
$\epsilon_{max} = 50\Delta\epsilon$	89.11	81.04	70.17	47.43	29.95	89.08
$\epsilon_{max} = 60\Delta\epsilon$	88.98	80.73	69.64	47.62	32.42	88.61

Table 8. Effect of ϵ_{max} on Fashion-MNIST test set (L2 norm-defined noise level, $\Delta\epsilon = 5/60$).

noise	0	0.1	0.25	0.5	1.0	0 (Validation)
$\varepsilon_{max} = 10\Delta\varepsilon$	93.45	88.12	74.67	45.33	10.84	93.05
$\varepsilon_{max} = 20\Delta\varepsilon$	92.70	88.36	78.80	56.86	20.59	92.83
$\varepsilon_{max} = 30\Delta\varepsilon$	92.34	88.34	79.51	59.92	25.57	91.96
$\varepsilon_{max} = 40\Delta\varepsilon$	91.47	87.33	78.82	60.98	28.45	91.26
$\varepsilon_{max} = 50\Delta\varepsilon$	90.62	86.46	78.10	60.66	29.16	90.80
$\varepsilon_{max} = 60\Delta\varepsilon$	89.69	85.01	76.82	60.50	31.90	90.17

Table 9. Effect of ε_{max} on SVHN test set (L2 norm-defined noise level, $\Delta\varepsilon = 2/60$).

From Table 8, it can be clearly seen that smaller ε_{max} leads to higher accuracy on clean data (noise level = 0), and larger ε_{max} leads to higher accuracy on noisy data. Since the maximum of the sample margins gradually increases during IMA training, the validation accuracy changes gradually (increasing and then decreasing), which makes it easy for the user of our method to choose the trained model with validation accuracy above the user-defined threshold.

We have done similar analyses on the models trained by IMA through 60 epochs on the SVHN dataset. We selected six models trained by IMA after 10, 20, 30, 40, 50, and 60 epochs, and we evaluated the robustness of the six models on the test set by using the 100-PGD attack. The results are reported in Table 9, which reveal the effect of ε_{max} on robustness. By monitoring the validation accuracy curve during training, the user of IMA can choose a trained model such that validation accuracy of the chosen model is above a threshold defined by the user to meet some application requirements.

C.3. COVID-19 application

Next, we plot the curves from the COVID-19 dataset in Fig. 2. The validation accuracy decreased only slightly after 20 epochs, which means the value of the allowed maximum margin is reasonable and the risk of margin overestimation is low.

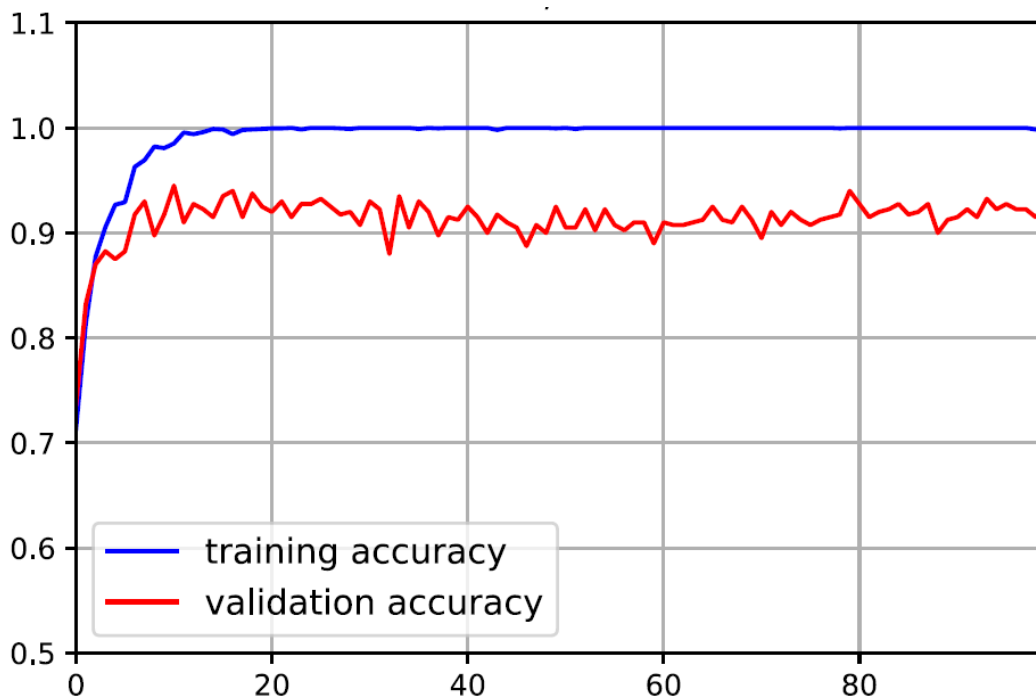


Figure 2. The training and validation curves (accuracy vs epoch) on the COVID-19 CT dataset obtained by using our IMA method. The accuracy scores are measured on clean data.

We note that **using the margin distribution estimated by IMA, we can obtain a good perturbation magnitude ε (10) for the vanilla adversarial training.** When $\varepsilon=10$, vanilla adversarial training achieved the best performance on the

COVID-19 CT dataset, when the noise level is smaller than 5. Please see Fig. 3 which contains the margin distribution estimated by IMA. In practice, we only need robustness against noises within a certain level, because significantly noisy images can only be produced from a malfunctioning CT machine. Thus, vanilla adversarial training (denoted by $\text{adv } \varepsilon$) should be good enough (much less computation cost and time compared to advanced adversarial training methods) as long as its parameter ε is appropriate. In this application, the good value of ε is 10, as revealed by the margin distribution estimated by IMA.

Noise	0	1	3	5	10	30
ce	95.75	9.98	0.0	0.0	0.0	0.0
IMA	91.25	90.25	87.75	83.25	57.00	0.0
MMA	89.00	88.50	84.25	80.00	57.50	0.0
DDN	90.25	88.50	84.75	78.75	59.75	0.0
TRADES	97.25	12.75	10.25	11.00	19.75	48.75
GAIRAT	92.50	88.25	72.50	50.25	12.50	0.0
FAT	90.00	88.00	81.50	64.25	16.00	0.0
adv10	93.50	92.25	89.50	83.50	53.50	0.0
adv20	95.25	25.50	20.00	21.50	22.50	3.00
adv30	85.25	6.25	1.50	2.00	5.00	23.00

Table 10. Classification accuracy on COVID-19 CT image dataset (L2 norm in 100-PGD).

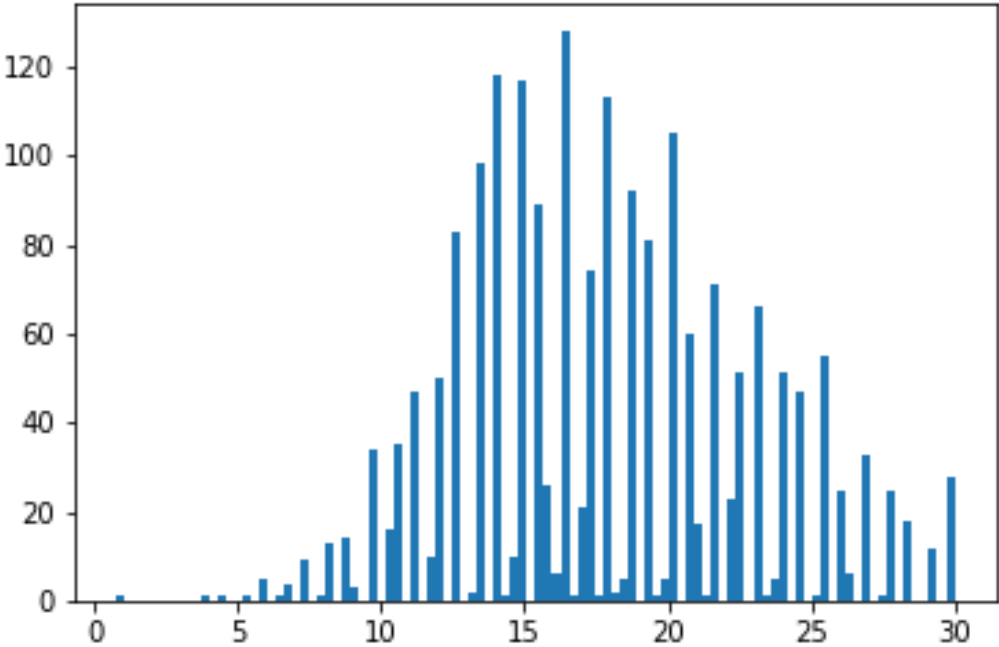


Figure 3. Margin distribution estimated by IMA on training set

For vanilla adversarial training, a straightforward way to find a good ε would be running grid research and evaluating the performance on the validation set. However, the grid-research in a large range (e.g. 0 to 30) is impractical because adversarial training is computationally expensive and time-consuming, compared to standard training with cross-entropy loss and clean

data. The margin distribution estimated by IMA can reveal a good range of the ε , and then grid-research in this small range can be performed to find the optimal ε . After the optimal ε is found for vanilla adversarial training, we could combine other techniques to further improve robustness.

D. Trade-off Between Robustness and Accuracy: Sample Margins' Perspective

In this appendix, we discuss the trade-off between robustness and accuracy from the perspective of sample margins. Let x denote a clean sample (i.e., unaltered by any adversarial attacks). Let $x_{(\delta)}$ denote the noisy sample generated by adding an adversarial noise δ to x , i.e., $x_{(\delta)} = x + \delta$. The vector norm of δ is ε . Let y denote the true class label of x .

For vanilla adversarial training, the class label of the noisy sample $x_{(\delta)}$ is assumed to be the same as the class label y of the clean sample x , no matter how large the noise level ε is. This label assignment can be wrong if ε is very large: large enough such that $x_{(\delta)}$ may “look like” a sample in a different class (not y), which has been shown in numerous studies. Here, “look like” refers to the judgment of a human. However, it is impractical to let a person actually look at every noisy sample to assign the right label to it.

If there are unlimited training samples that cover the input space, the samples can reveal the true class distributions $p(x|y)$. Then, the true decision boundary can be obtained by Bayesian classification using $p(x|y)$ and $p(y)$ of each of the classes. For example, $p(x|y)$ can be simply modeled by a Gaussian mixture model, and $p(y)$ may be assumed to be the same for every class. If we assume Bayesian classification with unlimited data is as good as the judgment of human experts, then Bayesian decision boundary is the most robust against noises. Let $g(x)$ denote the Bayesian classifier (g means ground truth), and it outputs the true label of the input: $y = g(x)$ and $y_{(\delta)} = g(x_{(\delta)})$. Using the Bayesian-optimal decision boundary, the true margin of a clean sample x can be obtained, i.e., the (minimum) distance between x and the decision boundary, which is denoted by $m(x)$.

During adversarial training (vanilla or IMA), a noisy sample is generated, $x_{(\delta)} = x + \delta$ with $\varepsilon = \|\delta\|$, and the class label y is assigned to it. If the noise/perturbation magnitude ε is larger than the true margin $m(x)$, then the assigned class label is wrong for $x_{(\delta)}$, and training the model with $x_{(\delta)}$ may cause its decision boundary to deviate from the Bayesian-optimal decision boundary even if the model $f(x)$ is initialized to be very close to $g(x)$. As a result, the classification accuracy on the test/validation set (clean data) will become lower, but the classification accuracy on the training set (clean data) may not change because the decision boundary has been pushed far away from the clean training sample x . **This phenomenon has been revealed by the training and validation accuracy curves of IMA (see Fig. 1), which explains the well-known trade-off between robustness and accuracy.**

Thus, the success of adversarial training depends on the accurate estimation of the true margin $m(x)$: if $\mathcal{E}(x) = m(x)$, then IMA could be perfect. For the 2D Moons dataset, IMA can indeed find a nearly perfect decision boundary, significantly better than the other methods (see Fig. 6).

However, for a high dimensional dataset (e.g. Fashion-MNIST), there are not enough training samples to cover the input space. As a result, perfect Bayesian classification is not achievable because of inaccurate estimation of $p(x|y)$, and IMA cannot obtain perfect estimations of the true margins because there may not exist the counterparts in other classes (not y) that can stop the expansion of ε -ball (margin estimation) of x in class y . Margin overestimation is revealed by the gradually-decreasing trend of the validation accuracy curve (see Fig. 1). However, the good news for the user of IMA is that the user can choose a good value of ε_{max} such that the validation accuracy is above a pre-defined threshold (e.g. 90% for the COVID-19 application). Also, the user does not need to pre-define the value of ε_{max} : the allowed maximum margin will increase with the number of training epochs (see Appendix C), and the user only needs to monitor the training and validation accuracy curves. Although IMA is not perfect, it brings such a significant **benefit to the user: the ease of “adjusting” ε_{max} to make a trade-off between robustness and accuracy (simply set $\beta=0.5$).**

From IMA, when an equilibrium state is reached, the distributions (i.e., local densities) of noisy samples in different classes are the same along the decision boundary (note: it does not necessarily mean the clean samples in different classes are equally spaced from the decision boundary). This is somewhat analog to Bayesian classification: at the optimal decision boundary, the distributions (densities) of samples in two classes are the same, assuming the classes have equal prior probabilities. From this perspective, the noisy samples, which are generated by IMA, serve as the surrogates of the real samples. Obviously, we cannot claim it is Bayesian classification because noisy samples may not reveal the true distributions. From this perspective, more advanced adversarial training methods may be developed such that the generated samples may reveal the true distributions (i.e., $p(x|y)$); if so, then the resulting decision boundary could be optimal and robust.

Additional Note: The validation accuracy curves of IMA in Fig. 1 in Appendix C increase and then gradually decrease. Please do not confuse this with the common concept of overfitting on clean data. Actually, when the models were trained

with cross-entropy loss and clean data, there is no gradually-decreasing trend in validation accuracy curves, as shown in Fig. 4. Thus, the only explanation of the gradually-decreasing trend in Fig. 1 in Appendix C is margin overestimation.

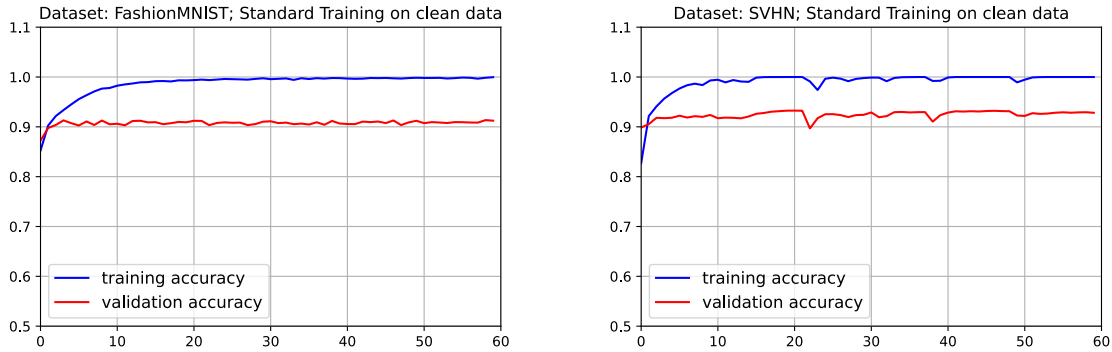


Figure 4. Training and validation accuracy curves (accuracy vs epoch) on the datasets, using standard training with cross-entropy loss and clean data. The accuracy scores are measured on clean data.

E. Further Discussion of the Equilibrium State

E.1. The basic idea of our IMA method

If there are only two classes and the data samples are linearly separable, then linear SVM (support vector machine) will produce a linear decision boundary in the “middle” between the two classes. The decision boundary of SVM is robust against noises: the classification output will not change if a small amount of noise δ is added to x as long as the vector norm of δ is smaller than the margin of x . Here, the margin of x is the (minimum) distance between x and the decision boundary.

In general, the data samples in multiple classes are nonlinearly separable, and the robust decision boundary should be somewhere in the middle between classes, which is the goal that IMA pursues. The result on the Moons dataset shows that IMA can indeed produce a nonlinear decision boundary in the “middle” between classes. We use 2D Moons dataset because it is impossible to directly visualize a nonlinear decision boundary in a high dimensional space.

E.2. The equilibrium state

Our IMA method is a heuristic-based method that is not derived from any theory. We use the equilibrium state analysis to provide a theoretical explanation of the method. We have shown that an equilibrium state can be achieved when the noisy samples have the same spatial distribution on the decision boundary. Here, we will analyze what will happen if the spatial distributions of the noisy samples in different classes are not the same on the current decision boundary. We note that our IMA method will actively generate and put noisy samples on (“close to” due to numerical precision) the current decision boundary of the neural network model, and the training is a dynamic process to adjust the decision boundary of the model. Let’s focus on the following two terms:

$$F_i \triangleq E_{X_n \in c_i \text{ and } X_n \in B_{ij}} = - \int q_i(x) \log(P_i(x)) dx \quad (1)$$

$$F_j \triangleq E_{X_n \in c_j \text{ and } X_n \in B_{ij}} = - \int q_j(x) \log(P_j(x)) dx \quad (2)$$

where $q_i(x)$ and $q_j(x)$ are the distributions (i.e., densities) of the noisy samples on the current decision boundary between the two classes, and $q_i(x)$ and $q_j(x)$ may not be equal to each other. In fact, F_i and F_j can be interpreted as two forces that try to expand the margins of the samples in the two classes against each other. By dividing the decision boundary into small regions (i.e., linear segments), the two integrals can be evaluated in the individual regions. In a region, if $q_i(x) > q_j(x)$ (i.e., more samples in class-i) then the current state is not in equilibrium: after updating the model using these noisy samples, the noisy samples in class-i will be correctly classified and the noisy samples in class-j will be incorrectly-classified (this is a simple result of classification with imbalanced data in the region), which means the decision boundary will shift towards the samples in class-j, and therefore the margins of the corresponding samples in class-i will expand and the margins of the

corresponding samples in class-j will shrink. Thus, the decision boundary may shift locally towards the samples in one of the classes. Obviously, the decision boundary will stop shifting when the local densities of noisy samples in different classes are the same along the decision boundary, i.e., $q_i(x)$ becomes equal to $q_j(x)$, which means an equilibrium state is reached.

F. Examples of Sample Images

F.1. D2 Heart

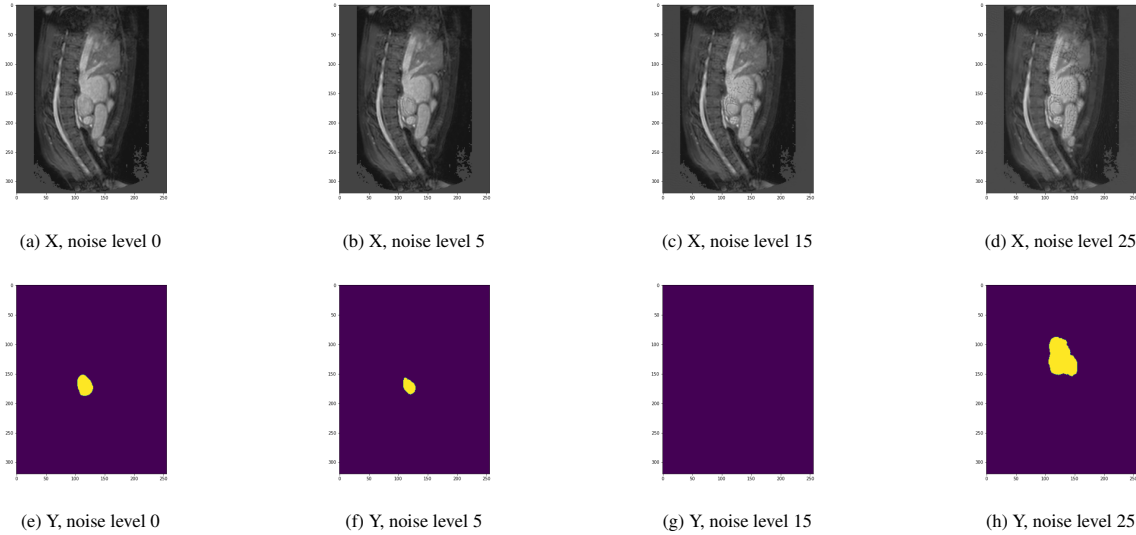


Figure 5. D2 Heart: The first row shows the input of nnUnet under different levels of 100-PGD adversarial noises. The second row shows the corresponding prediction of nnUnet under different levels of 100-PGD adversarial noises.

F.2. D4 Hippocampus

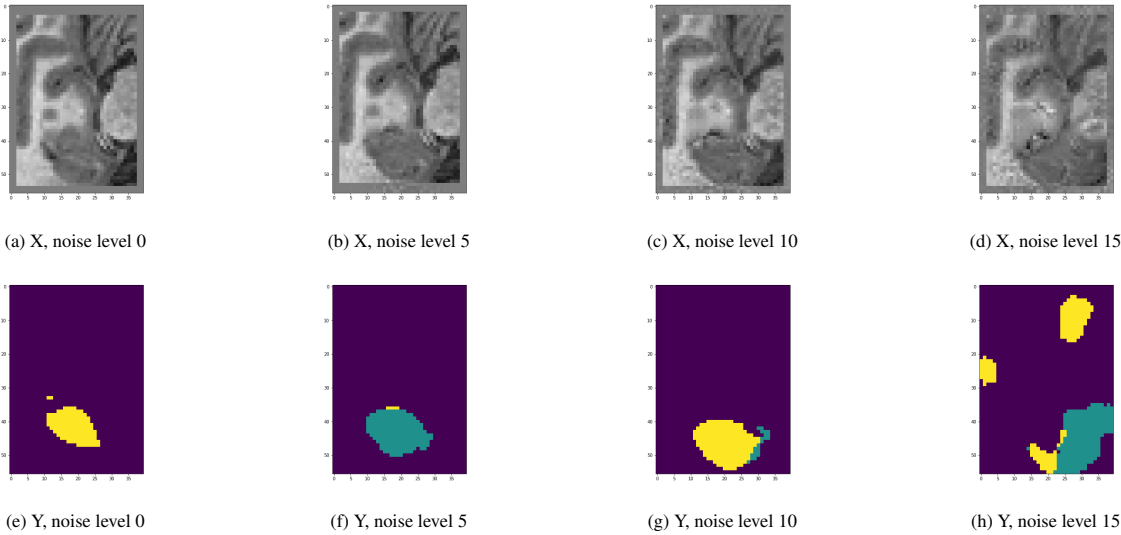


Figure 6. D4 Hippocampus: The first row shows the input of nnUnet under different levels of 100-PGD adversarial noises. The second row shows the corresponding prediction of nnUnet under different levels of 100-PGD adversarial noises.

F.3. D5 Prostate

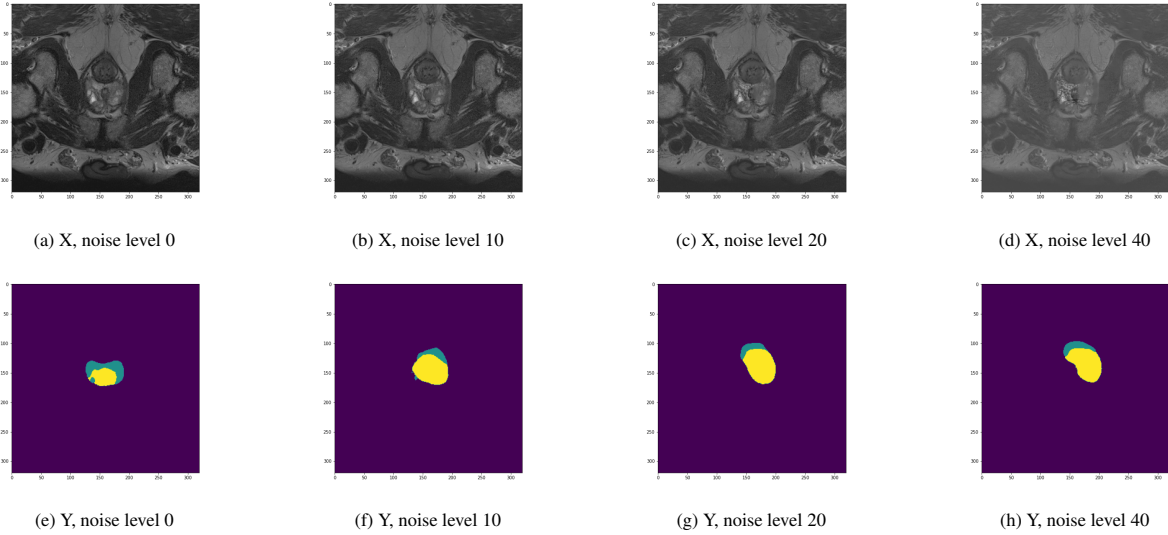


Figure 7. D5 Prostate: The first row shows the input of nnUnet under different levels of 100-PGD adversarial noises. The second row shows the corresponding prediction of nnUnet under different levels of 100-PGD adversarial noises.

F.4. Covid-19 CT Images

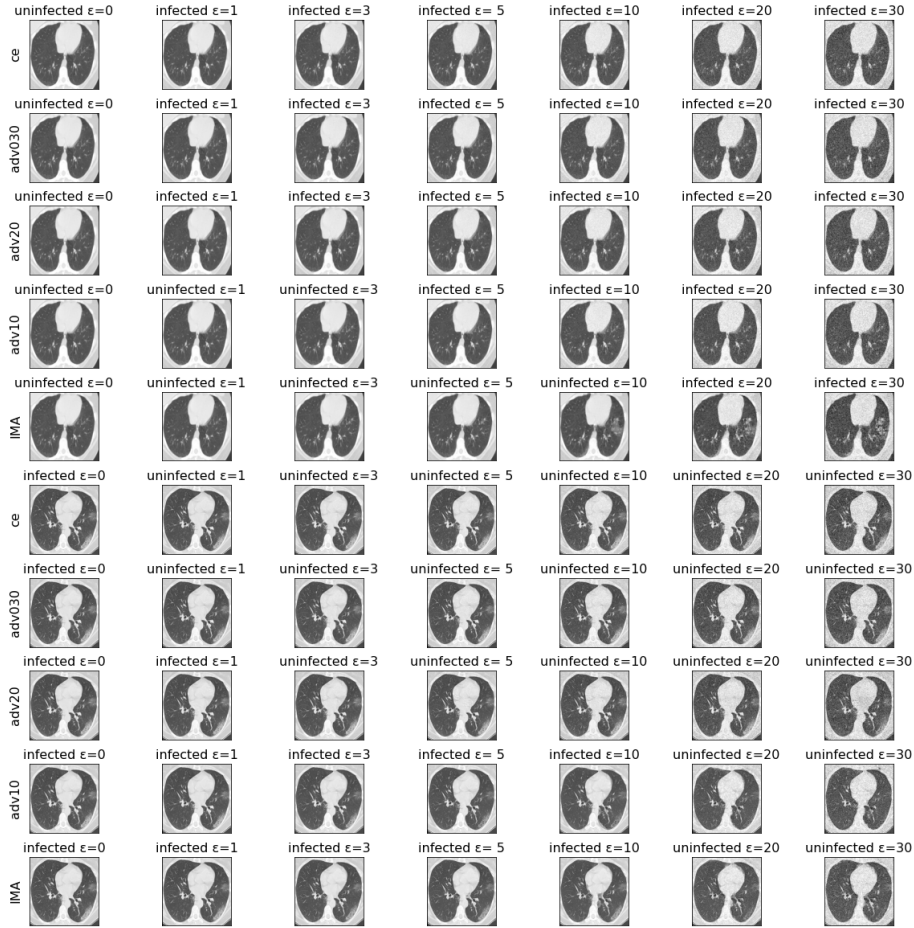


Figure 8. Each row shows a clean image and noisy images associated with a training method. The title of each image shows the predicted class label and the noise level. The clean images are correctly classified.

G. Details in experiments

G.1. Classification Tasks

For the datasets of Moons, Fashion-MNIST, SVHN, CIFAR10 and COVID-19 CT, β in Algorithm 1 is 0.5 because robustness and accuracy are equally important [2]. In Algorithm 3, N_{PGD} is 20, N_{binary} is 10 and there is no repeat. α in PGD is 4. Larger N_{PGD} and N_{binary} may lead to better convergence of the algorithms but more computing time. To make the comparison fair, the number of iterations and step size of TRADES are the same as those of the 20-PGD; and the number of PGD-iterations in MMA is 20. The number of iterations in DDN is 100, as suggested in [7]. Pytorch [5] is used for model implementation. Nvidia V100 and Titan V GPUs are used for model training and testing. All of the adversarial noises are measured in L2 norm.

In the Moons evaluation, for every method, we set the number of training epochs to 30, maximum noise level to 0.3, and batch-size to 128; Adam optimizer was used with default parameters.

In the SVHN evaluation, training details are as follows. For all the methods, batch size was 128. For IMA, the number of training epochs was 60, ε_{max} was 2.0 and Adam optimizer was used with default parameters. For MMA and DDN, the number of training epochs was 60, ε_{max} was 2.0 and Adam optimizer was used with default parameters. For FAT, self-adjusted τ was used as in [11], the number of training epochs was 60 and SGD optimizer was used with parameters the same as those in [11]. For GAIRAT, the number of training epochs was 60, “begin epoch” was 30 (half of the total number of epochs) and SGD optimizer was used with parameters the same as those in [12]. We set ε_{max} 1.0 for FAT and 1.0 for

GAIRAT, because FAT and GAIRAT with 2.0 cannot converge on this network. For TRADES, the number of training epochs was 60, ε_{max} was 2.0 and Adam optimizer was used with default parameters. Besides adversarial training, we did not use any other data augmentation (e.g. crop, etc.).

In the Fashion-FMNIST evaluation, training details are as follows. For all the methods, batch size was 128. For IMA, the number of training epochs was 60, ε_{max} was 5.0 and Adam optimizer was used with default parameters. For MMA and DDN, the number of training epochs was 60, ε_{max} was 5.0 and Adam optimizer was used with default parameters. For FAT, ε_{max} was 5.0, self-adjusted τ was used as in [11], the number of training epochs was 60 and SGD optimizer was used with parameters the same as those in [11]. For GAIRAT, ε_{max} was 5.0, the number of training epochs was 60, “begin epoch” was 30 (half of the total number of epochs) and SGD optimizer was used with parameters the same as those in [12]. For TRADES, the number of training epochs was 60, ε_{max} was 5.0 and Adam optimizer was used with default parameters. Besides adversarial training, we did not use any other data augmentation (e.g. crop, etc.).

In the CIFAR10 evaluation, training details are as follows. The network was WideResNet-28-4 used in [1]. For all the methods, batch size was 128. For IMA, ε_{max} was 1.5 (in the main paper) and Adam optimizer was used with default parameters. Because of the complexity of this dataset, 100 training epochs were used for IMA, starting with a WideResNet-28-4 pre-trained on clean data. For MMA and DDN, we directly used the results from [1]. For FAT, self-adjusted τ was used as in [11], the number of training epochs was 100 and SGD optimizer was used with parameters the same as those in [11]. For GAIRAT, the number of training epochs was 100, “begin epoch” was 50 (half of the total number of epochs) and SGD optimizer was used with parameters the same as those in [12]. For TRADES, the number of training epochs was 100, and Adam optimizer was used with default parameters. The ε_{max} was 1.0 for MMA, DDN, FAT, GAIRAT and TRADES because ε_{max} of 1.5 led to a significant loss in the accuracy on clean data for these methods.

In the COVID-19 CT evaluation, training details are as follows. For all the methods, batch size was 32. For IMA, MMA, TRADES, DDN and “adv ε ”, the number of epochs is 40 and the maximum noise level for training (ε_{max}) is 30. FAT and GAIRAT cannot converge in 40 epochs, so the two methods used 100 training epochs. The maximum noise level was 30 for IMA, MMA, DDN, TRADES, FAT and GAIRAT. Adam optimizer was used with default parameters. For “ce”, weight decay of 0.01 is applied with Adam (i.e., AdamW with default parameters).

G.2. Segmentation Tasks

All of the models in the experiments are trained with 50 epochs, where each epoch is defined as 50 training iterations. The reason for selecting these three datasets, D2, D4 and D5, is that the sizes of these datasets are relatively small so that the models can stably converge within 50 training epochs. All of these three datasets are publically available. The experiments were conducted on Nvidia V100 GPUs. Other training details are the same as [3].

H. Additional Supportive Results

H.1. Robust Overfitting

The three tables, Table 11, Table 12 and Table 13 show that, for each evaluation, our proposed method IMA’s best checkpoint is always the last check point, which means that IMA does not have the robust overfitting problem [4, 6] in our experiment, due to the self-adaptive nature of IMA.

Epoch	9	19	29	39	49	59
noise=0	0.93	0.92	0.92	0.92	0.91	0.90
noise=1.0	0.13	0.22	0.27	0.29	0.31	0.33

Table 11. IMA performance on SVHN validation set under 100-PGD attack with noise level=1.0

Epoch	9	19	29	39	49	59
noise=0	0.91	0.91	0.90	0.90	0.90	0.89
noise=3.0	0	0.08	0.20	0.27	0.31	0.33

Table 12. IMA performance on Fashion-MNIST validation set under 100-PGD attack with noise level=3.0

Epoch	9	19	29	39	49	59	69	79	89	99
noise=0	0.91	0.91	0.91	0.91	0.91	0.90	0.90	0.89	0.89	0.89
noise=1.5	0.0	0.0	0.05	0.09	0.13	0.17	0.17	0.20	0.21	0.21

Table 13. IMA performance on CIFAR10 validation set under 100-PGD attack with noise level=1.5

H.2. Black-box evaluation with SPSA

In this black-box evaluation, SPSA black-box adversarial attack [9] was used. For all the methods, 2048 samples were generated for gradient estimation in SPSA. Because of the high time complexity of SPSA, only a subset from test set was used to evaluate. For SVHN and Fashion-MNIST, the first 10 samples from each class in the test set, a total of 100 samples, formed the new test set for SPSA. For CIFAR10, the first 100 samples, formed the new test set for SPSA. From Table 14, Table 15 and Table 16, IMA had a good resist to black-box adversarial attack, which means there was no gradient obfuscation in IMA [9].

Noise	0.0	0.1	0.25	0.5	1.0
IMA	0.90	0.84	0.78	0.60	0.38
MMA	0.82	0.78	0.75	0.60	0.41

Table 14. Classification accuracy on SVHN (L2 norm in SPSA).

Noise	0.0	0.5	1.0	2.0	3.0
IMA	0.84	0.76	0.66	0.48	0.39
MMA	0.84	0.74	0.65	0.48	0.38

Table 15. Classification accuracy on F-MNIST (L2 norm in SPSA).

Noise	0.0	0.5	1.0	1.5
IMA	0.93	0.64	0.41	0.24
MMA	0.87	0.64	0.4	0.22

Table 16. Classification accuracy on CIFAR10 (L2 norm in SPSA).

References

- [1] Gavin Weiguang Ding, Yash Sharma, Kry Yik Chau Lui, and Ruitong Huang. Mma training: Direct input space margin maximization through adversarial training. In *ICLR*, 2020. 12
- [2] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2015. 11
- [3] Fabian Isensee, Paul F Jaeger, Simon AA Kohl, Jens Petersen, and Klaus H Maier-Hein. nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature methods*, 18(2):203–211, 2021. 12
- [4] Hoki Kim, Woojin Lee, and Jaewook Lee. Understanding catastrophic overfitting in single-step adversarial training. In *AAAI*, 2021. 12
- [5] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 11
- [6] Leslie Rice, Eric Wong, and Zico Kolter. Overfitting in adversarially robust deep learning. In *ICML*, pages 8093–8104. PMLR, 2020. 12

[7] Jérôme Rony, Luiz G Hafemann, Luiz S Oliveira, Ismail Ben Ayed, Robert Sabourin, and Eric Granger. Decoupling direction and norm for efficient gradient-based l2 adversarial attacks and defenses. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4322–4330, 2019. 11

[8] Feng Shi, Jun Wang, Jun Shi, Ziyang Wu, Qian Wang, Zhenyu Tang, Kelei He, Yinghuan Shi, and Dinggang Shen. Review of artificial intelligence techniques in imaging data acquisition, segmentation and diagnosis for covid-19. *IEEE reviews in biomedical engineering*, 2020. 1

[9] Jonathan Uesato, Brendan O’donoghue, Pushmeet Kohli, and Aaron Oord. Adversarial risk and the dangers of evaluating against weak attacks. In *ICML*, pages 5025–5034. PMLR, 2018. 13

[10] Jing Wang, Hongbing Lu, Zhengrong Liang, Daria Eremina, Guangxiang Zhang, Su Wang, John Chen, and James Manzione. An experimental study on the noise properties of x-ray ct sinogram data in radon space. *Physics in Medicine & Biology*, 53(12):3327, 2008. 1

[11] Jingfeng Zhang, Xilie Xu, Bo Han, Gang Niu, Lizhen Cui, Masashi Sugiyama, and Mohan Kankanhalli. Attacks which do not kill training make adversarial learning stronger. In *ICML*, pages 11278–11287. PMLR, 2020. 11, 12

[12] Jingfeng Zhang, Jianing Zhu, Gang Niu, Bo Han, Masashi Sugiyama, and Mohan Kankanhalli. Geometry-aware instance-reweighted adversarial training. In *ICLR*, 2021. 11, 12