

# ECE 661: Homework #5

## Adversarial Attacks and Defenses

Hai Li

ECE Department, Duke University — October 23, 2021

### Objectives

Homework #5 covers the contents of Lectures 16 ~ 18. This assignment starts by implementing several basic gradient-based adversarial attacks and analyzing how the  $\epsilon$  of the attack influences the perceptibility of the noise. Then, you will evaluate the attacks in both the whitebox and blackbox settings, measuring the trade-off between attack success and  $\epsilon$  in each. Finally, you will adversarially train some robust models and measure their ability to defend against such attacks.



**Warning:** You are asked to complete the assignment independently.

This lab has a total of **100** points plus 10 bonus points, yet your final score cannot exceed 100 points. You must submit your report in PDF format and your original codes for the lab questions through **Sakai** before **11:55:00pm, Tuesday, November 16**. You need to submit **two individual files** including (1) *a self-contained report in PDF format* that provides answers to all the Lab questions and clearly demonstrates all your lab results and observations, (2) the *attacks.py* file which contains your attack implementations. **Note that 10 percent of the grade will be deducted for the submissions uploaded in a zip file.**

### 1 True/False Questions (10 pts)

For each question, please provide a short explanation to support your judgment.

**Problem 1.1 (1 pt)** In an evasion attack, the attacker perturbs a subset of training instances which prevents the DNN from learning an accurate model.

**Problem 1.2 (1 pt)** In general, modern defenses not only improve robustness to adversarial attack, but they also improve accuracy on clean data.

**Problem 1.3 (1 pt)** In a backdoor attack, the attacker first injects a specific noise trigger to a subset of data points and sets the corresponding labels to a target class. Then, during deployment, the attacker uses a gradient-based perturbation (e.g., Fast Gradient Sign Method) to fool the model into choosing the target class.

**Problem 1.4 (1 pt)** Outlier exposure is an Out-of-Distribution (OOD) detection technique that uses OOD data during training, unlike the ODIN detector.

**Problem 1.5 (1 pt)** It is likely that an adversarial examples generated on a ResNet-50 model will also fool a VGG-16 model.

**Problem 1.6 (1 pt)** The perturbation direction used by the Fast Gradient Sign Method attack is the direction of steepest ascent on the local loss surface, which is the most efficient direction towards the decision boundary.

**Problem 1.7 (1 pt)** There are two schools of thought for adversarial perturbation optimization: optimization-based and constraint-based, and the PGD attack is an example of a constraint-based attack.

**Problem 1.8 (1 pt)** Analysis shows that the best layer for generating the most transferable feature space attacks is the final convolutional layer, as it is the convolutional layer that has the most effect on the prediction.

**Problem 1.9 (1 pt)** The DVERGE training algorithm promotes a more robust model ensemble, but the individual models within the ensemble still learn non-robust features.

**Problem 1.10 (1 pt)** On a backdoored model, the exact backdoor trigger must be used by the attacker during deployment to cause the proper targeted misclassification.

## 2 Lab 1: Environment Setup and Attack Implementation (20 pts)

In this section, you will train two basic classifier models on the FashionMNIST dataset and implement a few popular adversarial attack methods. The goal is to prepare an “environment” for attacking in the following sections and to understand how the adversarial attack’s  $\epsilon$  value influences the perceptibility of the noise. All code for this set of questions will be in the “Model Training” section of `HWK5_main.ipynb` and in the accompanying `attacks.py` file. Please include all of your results, figures and observations into your PDF report.

### Lab 1 (20 points)

- (a) (4 pts) Train the given *NetA* and *NetB* models on the FashionMNIST dataset. Use the default training parameters and save two checkpoints: “netA\_standard.pt” and “netB\_standard.pt”. What is the final test accuracy of each model? Do both models have the same architecture? (Hint: accuracy should be around 92% for both models).
- (b) (8 pts) Implement the  $L_\infty$ -constrained Projected Gradient Descent (PGD) adversarial attack in the `attacks.py` file. In the report, paste a screenshot of your `PGD_attack` function and describe what each of the input arguments is controlling. Then, using the “Visualize some perturbed samples” cell in `HWK5_main.ipynb`, run your PGD attack using *NetA* as the base classifier and plot some perturbed samples using  $\epsilon$  values in the range  $[0.0, 0.3]$ . At about what  $\epsilon$  does the noise start to become perceptible/noticeable? Do you think that you (or any human) would still be able to correctly predict samples at this  $\epsilon$  value? Finally, to test one important edge case, show that at  $\epsilon = 0$  the computed adversarial example is identical to the original input image.
- (c) (4 pts) Implement the Fast Gradient Sign Method (FGSM) attack in the `attacks.py` file. (Hint: treat the FGSM function as a wrapper around the PGD function). Please include a screenshot of your `FGSM_attack` function in the report. Then, plot some perturbed samples using the same  $\epsilon$  levels from the previous question and comment on the perceptibility of the FGSM noise. Does the FGSM and PGD noise appear visually similar?
- (d) (4 pts) Implement the Fast Gradient Method attack with  $L_2$  constraint in the `attacks.py` file. Please include a screenshot of your `FGM_L2_attack` function in the report. Then, plot some perturbed samples using the same  $\epsilon$  values in the range of  $[0.0, 4.0]$  and comment on the perceptibility of the  $L_2$  constrained noise. How does this noise compare to the  $L_\infty$  constrained FGSM and PGD noise visually? (Note: This attack involves a normalization of the gradient, but since these attack functions take a batch of inputs, the norm must be computed separately for each element of the batch).

### 3 Lab 2: Measuring Attack Success Rate (30 pts)

In this section, you will measure the effectiveness of your PGD and FGSM attacks. Remember, the goal of an adversarial attacker is to perturb the input data such that the classifier outputs a wrong prediction, while the noise is minimally perceptible to a human observer. All code for this set of questions will be in the “Test Attacks” section of `HWK5_main.ipynb` and in the accompanying `attacks.py` file. Please include all of your results, figures and observations into your PDF report.

#### Lab 2 (40 points)

- (a) (2 pts) Briefly describe the difference between a whitebox and blackbox adversarial attacks. Also, what is it called when we generate attacks on one model and input them into another model that has been trained on the same dataset?
- (b) (3 pts) *Random Attack* - To get an attack baseline, we use random uniform perturbations in range  $[-\epsilon, \epsilon]$ . We have implemented this for you in the `attacks.py` file. Test at least eleven  $\epsilon$  values across the range  $[0, 0.1]$  (e.g., `np.linspace(0, 0.1, 11)`) and plot the accuracy vs epsilon curve with y-axis range  $[0, 1]$ . How effective is random noise as an attack? (Note: in the code, whitebox and blackbox accuracy is computed simultaneously)
- (c) (10 pts) *Whitebox Attack* - Using your pre-trained “NetA” as the whitebox model, measure the whitebox classifier’s accuracy versus attack epsilon for both the FGSM and PGD attacks. Test at least eleven  $\epsilon$  values across the range  $[0, 0.1]$  (e.g., `np.linspace(0, 0.1, 11)`) and plot the accuracy vs epsilon curve on the same axes as the curve from part (b). For the PGD attacks, use `perturb_iters = 10` and  $\alpha = 1.85 * (\epsilon / \text{perturb\_iters})$ . Comment on the difference between the attacks. Do either of the attacks induce the equivalent of “random guessing” accuracy? If so, which attack and at what  $\epsilon$  value? (Note: in the code, problem (b) and (c) can be run simultaneously)
- (d) (10 pts) *Blackbox Attack* - Using the pre-trained “NetA” as the whitebox model and the pre-trained “NetB” as the blackbox model, measure the ability of adversarial examples generated on the whitebox model to transfer to the blackbox model. Specifically, measure the blackbox classifier’s accuracy versus attack epsilon for both the FGSM and PGD attacks. Use the same  $\epsilon$  values across the range  $[0, 0.1]$  and plot the blackbox model’s accuracy vs epsilon curve. *Please plot this curve on the same axes as the curves from part (b) and (c).* For the PGD attacks, use `perturb_iters = 10` and  $\alpha = 1.85 * (\epsilon / \text{perturb\_iters})$ . Comment on the difference between the blackbox attacks. Do either of the attacks induce the equivalent of “random guessing” accuracy? If so, which attack and at what  $\epsilon$  value? (Note: in the code, problem (b) and (c) can be run simultaneously)
- (e) (5 pts) Comment on the difference between the attack success rate curves (i.e., the accuracy vs. epsilon curves) for the whitebox and blackbox attacks. How do these compare to effectiveness of the naive uniform random noise attack? Which is the more powerful attack and why? Does this make sense? Also, consider the epsilon level you found to be the “perceptibility threshold” in Lab 1.b. What is the attack success rate at this level and do you find the result somewhat concerning?

## 4 Lab 3: Adversarial Training (40 pts + 10 Bonus)

In this section, you will implement a basic defense called adversarial training (AT). As the name suggests, this involves training the model against adversarial examples. Specifically, we will be using the AT described in <https://arxiv.org/pdf/1706.06083.pdf>, which formulates the training objective as

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

Importantly, the inner maximizer specifies that all of the training data should be adversarially perturbed before updating the network parameters. All code for this set of questions will be in the `HWK5_main.ipynb` file. Please include all of your results, figures and observations into your PDF report.

### Lab 3 (40 points + 10 Bonus)

- (a) (5 pts) Using the given “Model Training” code, adversarially train a “NetA” model using a FGSM attack with  $\epsilon = 0.1$ , and save the model checkpoint as “netA\_advtrain\_fgsm0p1.pt”. What is the final training accuracy of this model on the clean test data? Is the accuracy less than the standard trained model?
- (b) (5 pts) Using the given “Model Training” code, adversarially train a “NetA” model using a PGD attack with  $\epsilon = 0.1$ , `perturb_iters = 4`,  $\alpha = 1.85 * (\epsilon / \text{perturb\_iters})$ , and save the model checkpoint as “netA\_advtrain\_pgd0p1.pt”. What is the final training accuracy of this model on the clean test data? Is the accuracy less than the standard trained model? Are there any noticeable differences in the training convergence between the FGSM-based and PGD-based AT procedures?
- (c) (15 pts) For the model adversarially trained with FGSM (“netA\_advtrain\_fgsm0p1.pt”), compute the accuracy versus attack epsilon curves against both the FGSM and PGD attacks (as whitebox methods only). Use  $\epsilon = [0.0, 0.02, 0.04, \dots, 0.12]$ . Is the model robust to both types of attack? If not, explain and analyze why one attack might be better than another. (Note: you can run this code in the “Test Robust Models” cell of the `HWK5_main.ipynb` notebook).
- (d) (15 pts) For the model adversarially trained with PGD (“netA\_advtrain\_pgd0p1.pt”), compute the accuracy versus attack epsilon curves against both the FGSM and PGD attacks (as whitebox methods only). Use  $\epsilon = [0.0, 0.02, 0.04, \dots, 0.12]$ , `perturb_iters = 10`,  $\alpha = 1.85 * (\epsilon / \text{perturb\_iters})$ . Is the model robust to both types of attack? Explain why or why not and analyze. Can you conclude that one adversarial training method is better than the other? If so, provide an intuitive explanation as to why (this paper may help explain: <https://arxiv.org/pdf/2001.03994.pdf>). (Note: you can run this code in the “Test Robust Models” cell of the `HWK5_main.ipynb` notebook).
- (e) (Bonus 5 pts) Using PGD-based AT, train a few more models with different training  $\epsilon$  values. Is there a trade-off between clean data accuracy and training  $\epsilon$ ? Is there a trade-off between robustness and training  $\epsilon$ ? What happens when the testing PGD’s  $\epsilon$  is larger than the  $\epsilon$  used for training? In the report, provide answers to all of these questions along with evidence (e.g., plots and/or tables) to substantiate your claims.
- (f) (Bonus 5 pts) Plot the saliency maps for a few samples from the FashionMNIST test set as measured on both the standard (non-AT) and PGD-AT models. Do you notice any difference in saliency? What does this difference tell us about the representation that has been learned? (Hint: plotting the gradient w.r.t. the data is often considered a version of saliency, see <https://arxiv.org/pdf/1706.03825.pdf>)