
FREE vs FAST ADVERSARIAL TRAINING

Giuseppe Capaldi

University of Rome "La Sapienza"
capaldi.1699498@studenti.uniroma1.it

Gianluca Capozzi

University of Rome "La Sapienza"
capozzi.1693255@studenti.uniroma1.it

September 24, 2020

ABSTRACT

In this work we propose a comparison between two methods for adversarial training: Free and Fast. The code is available at <https://github.com/not-a-genius/neuralNetworkExam>.

1 Introduction

Adversarial training is a method for training robust deep neural networks, training on adversarial examples as a defense against adversarial attacks. This is typically considered to be more expensive than standard training due to the necessity of constructing adversarial examples via a first-order method like projected gradient descent (PGD). The papers presented show that it is possible to train empirically robust models using methods no more costly than standard training in practice. These methods show results comparable to train a model against PGD but at lower cost in terms of computational time. The goal of adversarial training is to learn a model which is not only accurate on the data (accuracy on train and test data) but also accurate on adversarially perturbed versions of the data (i.e. accuracy on images perturbed using the PGD attack). This leads to particular attention at the results we obtained after adversarial training experiments, aware of the presence of an acceptable compromise in accuracy, given a large increase in robustness due to the tradeoff between robustness and generalization [1, 2, 3].

1.1 Related Works

Our work is based on “Adversarial training for free!” paper [4], that presented an algorithm able to eliminate the overhead cost of generating adversarial examples by recycling the gradient information computed when updating the model parameters. While researching related work in the field we discovered a more recent paper, called “Fast is better than free: Revisiting adversarial training” [5] presenting an already known method called FGSM, previously considered ineffective due to what the paper calls “catastrophic overfitting”. This failure condition is preventable with the use of random initialization points. Moreover the same paper reported a further acceleration in training even for free algorithm thanks to standard techniques for efficient training, including cyclic learning rate and mixed-precision arithmetic. This caught our interest and we decided to try to replicate the reported results and compare them with our implementation, in order to confirm or deny these thesis.

1.2 Dataset and architecture

We decided to choose CIFAR-10 as the dataset on which our experiments have been conducted to adapt to the available hardware (our GPU is an Nvidia RTX 2070 super). The CIFAR-10 dataset consists of 60000 32x32 colour images in 10 classes, with 6000 images per class. There are 50000 training images and 10000 test images. The dataset is divided into five training batches and one test batch, each with 10000 images. The test batch contains exactly 1000 randomly-selected images from each class. The training batches contain the remaining images in random order, but some training batches may contain more images from one class than another. Between them, the training batches contain exactly 5000 images from each class.

Both WRN32 and ResNet50 have been tried as models to train on but their lack of accuracy even at high numbers of epochs suggested us to switch to shallower models, so we used PreActResNet18 which is based on ResNet18.

2 Adversarial Machine Learning

Adversarial machine learning is a machine learning technique that attempts to fool models by supplying deceptive input (called also adversarial example). Adversarial examples exploit the way artificial intelligence algorithms work to disrupt the behavior of artificial intelligence algorithms. In the past few years, adversarial machine learning has become an active area of research as the role of AI continues to grow in many of the applications we use. There's growing concern that vulnerabilities in machine learning systems can be exploited for malicious purposes. The main idea to deal with adversarial examples is represented by adversarial training; it can be traced back to [6], in which models were hardened by producing adversarial examples and injecting them into training data [4].

This technique can be summarized as follows: given a network f_θ parametrized by θ , a dataset (x_i, y_i) , a loss function l and a threat model Δ , the learning problem consists in the following optimization problem:

$$\min_{\theta} \sum_i l(f_\theta(x_i + \delta), y_i)$$

A typical choice for the adversarial model is to take $\Delta = \{\delta : \|\delta\|_\infty \leq \epsilon\}$ for some $\epsilon > 0$.

3 PGD attack

The PGD attack is a white-box attack, which means that the attacker has access to the model's weights, giving to the attacker much more power than a black-box attack (in which the model's weights are not known) as he can specifically craft the attack to fool the chosen model without having to rely on transfer attacks, which often result in human-visible perturbations. The key to understanding the PGD attack is to frame finding an adversarial example as a constrained optimisation problem. PGD attempts to find the perturbation that maximises the loss of a model on a particular input while keeping the size of the perturbation smaller than a specified amount referred to as epsilon. This constraint is usually expressed as the l^2 or l^∞ norm of the perturbation and it is added so the content of the adversarial example is the same as the unperturbed sample. The PGD algorithm can be summarised with the 4 steps below although the attacker is free to apply any optimisation improvements such as momentum, Adam, multiple restarts etc...

1. Start from a random perturbation in the l_p ball around a sample;
2. Take a gradient step in the direction of the greatest loss;
3. Project perturbation back into the l_p ball if necessary;
4. Repeat steps 2-3 until convergence.

In our case, adversarial training consists simply of putting the PGD attack inside the training loop, applying a kind of "data augmentation"; in this case, instead of performing random transformations as a preprocessing step, we create specific perturbations that best fool the model and indeed adversarially trained models do exhibit less overfitting when trained on small datasets.

The PGD attack is based on the simplest version of what is called Fast Gradient Sign Method, used to approximate the inner maximization of Delta. This could be seen as a relatively inaccurate approximation of the inner maximization for l_∞ perturbations, and has the following closed form (as shown in [7]):

$$\delta^* = \epsilon \cdot \text{sign}(\nabla_x l(f(x), y)).$$

A better approximation of the inner maximization is to take multiple, smaller FGSM steps of size alpha. When the iteration leaves the threat model, it is projected back to the set Δ . Multiple restarts within the threat model Δ typically improve the approximation of the inner maximization even further. The combination of all these techniques is the PGD attack. The problem of such technique is that the number of gradient computations here is proportional to $O(MK)$ in a single epoch, where M is the size of the dataset and K is the number of steps taken by the PGD adversary. This is K times greater than standard training so adversarial training is typically K times slower than standard training. The pseudo-code of the PGD attack is reported in Algorithm 1.

Algorithm 1 PGD adversarial training for T epochs, given some radius ϵ , adversarial step size α and K PGD steps and a dataset of size M for a network f_θ

```

1: for  $t = 1 \dots T$  do
2:   for  $i = 1 \dots M$  do
3:     // Perform PGD adversarial attack
4:      $\delta = 0$  // or randomly initialized
5:     for  $j = 1 \dots K$  do
6:        $\delta = \delta + \alpha \cdot \text{sign}(\nabla_\delta l(f_\delta(x_i + \delta), y_i))$ 
7:        $\delta = \max(\min(\delta, \epsilon), -\epsilon)$ 
8:      $\theta = \theta - \nabla_\theta l(f_\theta(x_i + \delta), y_i)$  // Update model weights with some optimizer, e.g. SGD
    
```

The execution of K-PGD in [8] took four days on a Titan X (with model WideResNet and dataset CIFAR-10). This led us to consider the reported results obtained using this adversarial training, without running them directly.

4 Free Adversarial training for free!

Free adversarial training [4] claims to reduce the overhead compared to natural training through the computation of ascent step by reusing the backward pass needed for the descent step. So, the gradient with respect to the network parameters is computed on the backward pass and later on this same backward pass the gradient of the loss with respect to the input image is also computed.

In order to allow multiple adversarial updates on the same image the training is made on the same minibatch m times in a row. Every minibatch is repeated m times before switching to the next minibatch (“minibatch replays”), so the number of epochs can be considered the number of epochs divided by m , such that the overall number of training iterations remains constant. Note that perturbations are not reset between mini-batches. The value of this parameter is shown to be crucial, since increasing m comes with the risk of increasing generalization error so robustness is increased at the cost of validation accuracy on natural images.

The pseudocode of Free is shown in Algorithm 2.

Algorithm 2 “Free” Adversarial Training (Free- m)

Require: Training samples X , perturbation bound ϵ , learning rate τ , hop steps m

```

1: Initialize  $\theta$ 
2:  $\delta \leftarrow 0$ 
3: for  $epoch = 1 \dots N_{ep}/m$  do
4:   for minibatch  $B \in X$  do
5:     for  $i = 1 \dots m$  do
6:       // Update  $\theta$  with stochastic gradient descent
7:        $g_\theta \leftarrow \mathbb{E}_{(x,y) \in B} \nabla_\theta [l(x + \delta, y, \theta)]$ 
8:        $g_{adv} \leftarrow \nabla_x l(x + \delta, y, \theta)$ 
9:        $\theta \leftarrow \theta - \tau g_\theta$ 
10:      // Use gradients calculated for the minimization step to update  $\delta$ 
11:       $\delta \leftarrow +\epsilon \cdot \text{sign}(g_{adv})$ 
12:       $\text{clip}(\delta, -\epsilon, \epsilon)$ 
    
```

It may be possible that catastrophic forgetting happens, e.g. when all the “informative” images of one class are in the first few mini-batches. In this extreme case, we do not see useful examples for most of the epoch, and forgetting may occur. How much does mini-batch replay hurt generalization? To answer the question we followed [4] approach trying different values of m , reaching the same results: dropoff in accuracy for small m , while we see an inversion in PGD and test accuracy for high values of it. On CIFAR10 robustness can be reached with values of $m \leq 10$, our best results have been obtained with $m = 8$.

It is known that adversarially trained classifiers, e.g. PGD, have generative behaviours. When an image is perturbed from one class to another, it adopts features that make it “look” like its adopted class, to a human eye. Free trained models exhibit this benefit and also smooth and flattened loss surface. Indeed, standard adversarially trained classifiers show the property to have a flatten and smooth loss landscape while some defenses work by “masking” the gradients. This techniques makes it difficult to identify adversarial examples using gradient methods, even though adversarial examples remain present. Reference [9] argues that gradient masking adds little security, free training does not operate

by masking gradients and uses a rough loss surface. In the following table we present the results of our experiments conducted using Free method and varying the number of minibatch replays. We can notice that the train time increase as m increases as expected but in the case of $m = 10$ we got values of accuracies lower than $m = 8$. The following table shows the results of our experiments:

| Minibatch-replay (m) | Accuracy using PreActResNet18 | | | | Epoch | Train Time (min) | Avg Epoch Time (sec) |
|----------------------|-------------------------------|--------|--------|--------|-------|------------------|----------------------|
| | Natural Images | PGD-10 | PGD-20 | PGD-50 | | | |
| $m = 2$ | 0.895 | 0.324 | 0.297 | 0.290 | 45 | 22 | 29.4 |
| $m = 4$ | 0.867 | 0.471 | 0.457 | 0.451 | 45 | 44 | 58.2 |
| $m = 8$ | 0.839 | 0.498 | 0.486 | 0.483 | 45 | 88 | 117 |
| $m = 10$ | 0.821 | 0.490 | 0.482 | 0.479 | 45 | 108 | 143.4 |

5 Fast is better than free

| Step-size | Accuracy using PreActResNet18 | | | | Epoch | Train Time (min) | Avg Epoch Time (sec) | Batch size |
|---------------|-------------------------------|--------|--------|--------|-------|------------------|----------------------|------------|
| | Natural Images | PGD-10 | PGD-20 | PGD-50 | | | | |
| $\alpha = 10$ | 0.847 | 0.480 | 0.463 | 0.450 | 45 | 21 | 28.8 | 128 |
| $\alpha = 10$ | 0.850 | 0.454 | 0.438 | 0.434 | 45 | 20 | 25.8 | 256 |
| $\alpha = 16$ | 0.856 | 0 | 0 | 0 | 45 | 21 | 28.2 | 128 |
| $\alpha = 16$ | 0.87 | 0 | 0 | 0 | 45 | 20 | 26.4 | 256 |

6 Conclusions

References

- [1] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. ICLR, 1050:11, 2018.
- [2] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P Xing, Laurent El Ghaoui, and Michael I Jordan. Theoretically principled trade-off between robustness and accuracy. ICML, 2019a.
- [3] Ali Shafahi, W Ronny Huang, Christoph Studer, Soheil Feizi, and Tom Goldstein. Are adversarial examples inevitable? ICLR, 2019a.
- [4] Ali Shafahi, Mahyar Najibi, Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer, Larry S Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! arXiv preprint arXiv:1904.12843, 2019b.
- [5] Eric Wong and Leslie Rice and J. Zico Kolter. Fast is better than free: revisiting adversarial learning, arxiv:2001.03994, 2020.
- [6] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. ICLR, 2015.
- [7] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572, 2014.
- [8] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. ICLR, 2017.
- [9] Logan Engstrom, Andrew Ilyas, and Anish Athalye. Evaluating and understanding the robustness of adversarial logit pairing. arXiv preprint arXiv:1807.10272, 2018.