

CSDS 440 Class Project: Adversarial Machine Learning

Shaochen (Henry) Zhong, sxz517

Minyang Tie, mxt497

Alex Useloff, adu3

Austin Keppers, agk51

David Meshnick, dcm101

Due and submitted on 12/04/2020

Fall 2020, Dr. Ray

Contents

1	Introduction and Significance	3
2	Individual Reports	5
2.1	Shaochen (Henry) Zhong's Individual Report	5
2.1.1	Overview	5
2.1.2	Fast Gradient Sign Method	5
	Algorithm Intuition	5
	Algorithm Implementation	5
	Experiments	5
	Evaluation	5
2.1.3	Hop Skip Jump	6
	Algorithm Intuition	6
	Algorithm Implementation	6
	Experiments	6
	Evaluation	6
2.1.4	Feature Collision	6
	Algorithm Intuition	6
2.2	Minyang Tie's Individual Report	6
2.3	Alex Useloff's Individual Report	6
2.4	Austin Keppers' Individual Report	6
2.5	David Meshnick's Individual Report	8
3	Comparative Study and Discussion	8
3.1	Overview	8
3.1.1	Datasets and Sample Selections	8
3.1.2	ART	9
3.1.3	Metrics	9

3.1.4	Adversarial Rivalry	10
3.2	Attack Algorithms	10
3.2.1	Evasion	11
3.2.2	Poisoning	14
3.2.3	Conclusion	14
3.3	Defense Algorithms	15
3.3.1	Detector	15
3.3.2	Pre-processor	17
3.3.3	Transformer	19
3.3.4	Conclusion	19
4	References	19

1 Introduction and Significance

In the field of machine learning, it is often taken for granted that the testing examples have no malicious intent – as if something is labeled to be a dog, it will indeed look like a dog. However, with the growing popularity of machine learning, robustness against adversarial attacks has been a more and more important metric. Adversarial machine learning is the field studying how to attack a model to make it output incorrect results (e.g., false predictions) in different stage of the model building, and how to make a model more robust against various kinds of attacks.

Here is a classic example of a successful adversarial attack borrowed from Google [?]: by applying the middle perturbation to the original image, a supposedly Labrador Retriever, while still looking like a Labrador Retriever, is now misclassified as a Weimaraner.

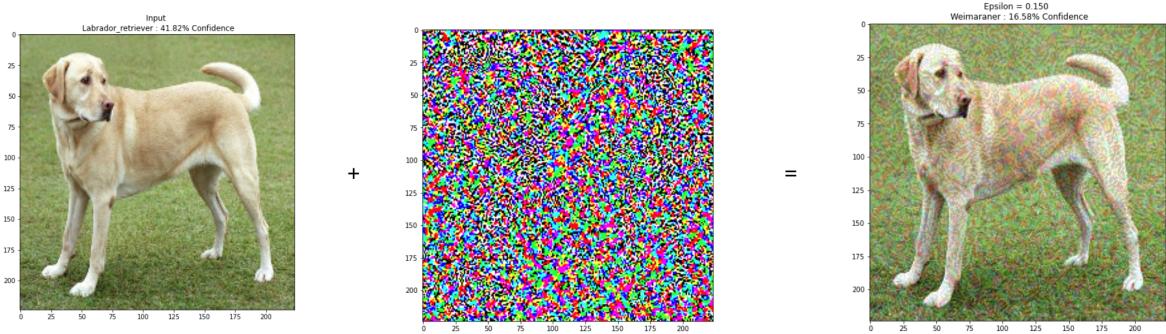


Figure 1: An example of adversarial perturbation resulting successful adversarial attack

One explanation[?] of why adversarial attack works is because there is a natural distinction between how we read and how machine read into a piece of information. When we humans are asked to classify between a *dog* from a *cat*, we know to focus on the ears and snout:

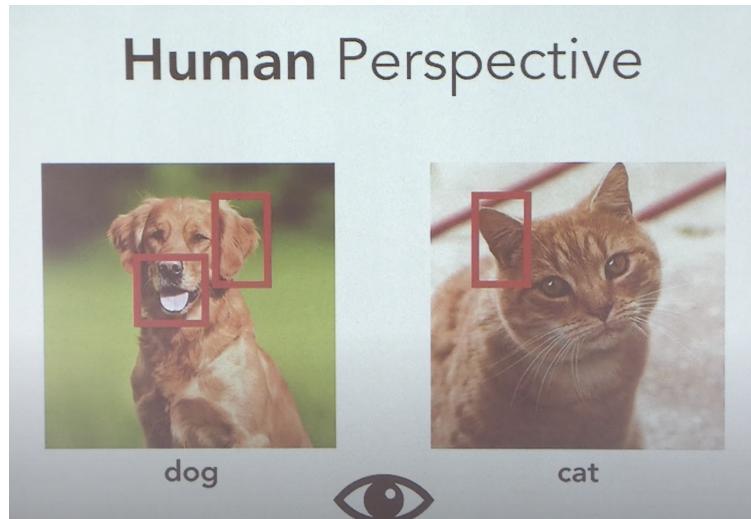


Figure 2: Human perspective

and we considered the perturbation in [Figure 1] to be meaningless because we cannot digest useful information out of it. But that might not be the case to the “eyes” of a machine as it has no knowledge about cats and dogs. To “translate” the machine’s perspective to “human terms,” how machine look at these cat/dog pictures are probably like the following:

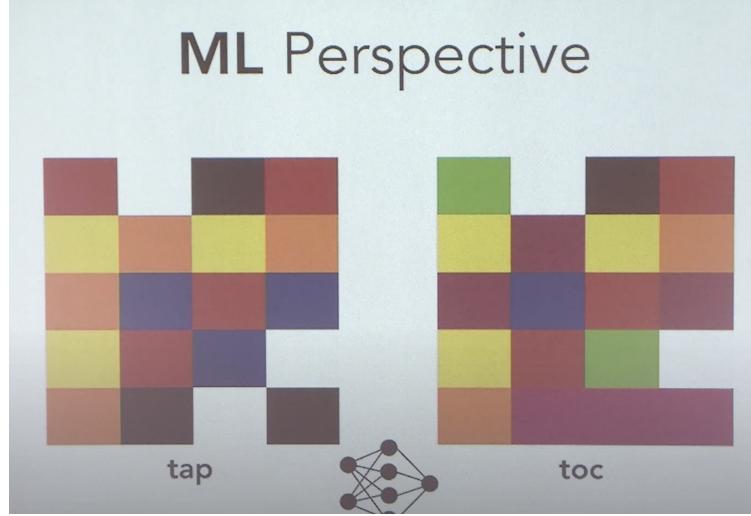


Figure 3: Machine perspective

In fact, a machine might consider the perturbation in [Figure 1] to be more “informative” than the features we care about (ears and snout). And since models are set to maximize accuracy as a general goal, it will utilize both features – both the *robust features* (features that keep being robust after adversarial perturbation), and the *non-robust features* (e.g., some “noise” to human):

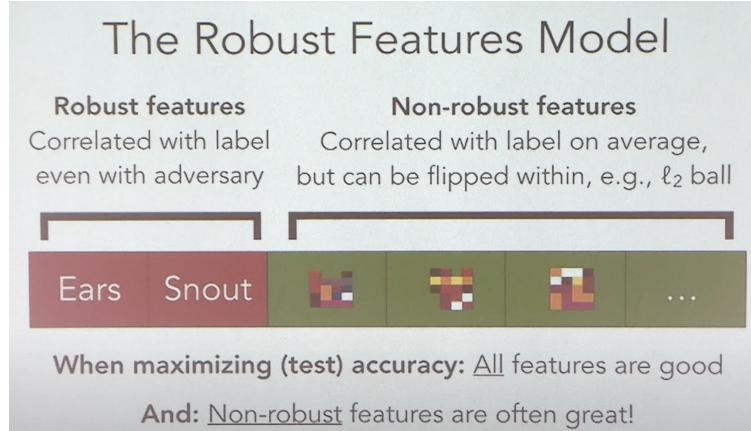


Figure 4: Robust and non-robust features

Thus, if a model’s decision is largely based on these non-robust features, we may apply adversarial perturbations to adjust these features and causing the model to output false results.

For this paper in particular, we will look into *evasion attack* (causing the model to output false result by modifying the input), *poisoning attack* (alter the training set of a model to make it produce inaccurate decision boundaries), and some corresponding defense methods against these attacks.

2 Individual Reports

2.1 Shaochen (Henry) Zhong's Individual Report

2.1.1 Overview

I have wholeheartedly lead and contributed to this group project, below is an itemized list of my contributions. I have inquired Dr. Ray and confirmed these can be considered as “extra works” and maybe give me some grade boost. Thanks :)

- Read 3 papers on algorithms.
- Implemented 2 algorithms (FGSM and Hop Skip Jump) with 3 extensions.
- Piplined attack algorithms to work with 2 datasets, collected almost all (7 algorithms/useable extensions) attack experiments data (except backdoor and one pixel attack) for the comparative evaluations.
- Piplined and collected experiments data for FGSM, Hop Skip Jump, DeepFool attacks (and their useable extensions) against Detector and Spatial Smoothing defenses on 2 datasets.
- Implemented L_2 and L_∞ perturbation budget to aid comparative evaluation.
- Wrote **Introduction and Significance** section.
- Plotted all graphs and charts in **Comparative Study and Discussion**.
- Helped group move forward by making technical decisions, distributing works, setting up deadlines, and facilitating coordination between groupmates.

2.1.2 Fast Gradient Sign Method

Algorithm Intuition FGSM is a *white-box* evasion attack algorithm. Being white-box means the attacker is assumed to have access to the internal of the model: the structure, the parameters... basically each and every details of the model is considered to be known.

In this case, we mostly care about the loss function of the model. The intuition of FGSM is elegant and effective – in short: gradient ascent. Assume we have an benign example x with label y , and we are trying to make it adversarial by letting the model classify x_{adv} to be not y . We apply the following perturbation on x to achieve an x_{adv}

$$x_{\text{adv}} = x + \epsilon \cdot \text{sign}(\nabla_x J(\theta, x, y)) \quad (1)$$

where J is the loss function and θ is the parameters of the model. It is clearly to tell that by finding a proper amount of ϵ , we know x_{adv} will eventually cross the $y - \neg y$

Algorithm Implementation

Experiments

Evaluation

2.1.3 Hop Skip Jump

Algorithm Intuition

Algorithm Implementation

Experiments

Evaluation

2.1.4 Feature Collision

Algorithm Intuition

2.2 Minyang Tie's Individual Report

2.3 Alex Useloff's Individual Report

2.4 Austin Keppers' Individual Report

The main algorithm that I focused on for this report was Defensive Distillation. Defensive distillation is a way to train a deep neural network classifier with the goal of better defending against adversarial examples that was first proposed in the paper *Distillation as a Defense to Adversarial Perturbations against Deep Neural Networks* by Papernot et all. The main idea behind Defensive distillation is that Deep Neural Networks are susceptible to adversarial examples because they make overly confident predictions when a sample has elements of two different classes. Defensive Distillation attempts to deal with this by using a classifier that is trained on soft labels instead of hard labels. Soft labels provide probabilities that an example belongs to each of the classes in a model rather than identifying a specific class the example belongs to.

Distillation is a technique that was created in order to run Deep Neural Networks on devices with lower computational power by reducing the size of the network. The idea is that hard labels can be used to train a large neural network, and then this classifier can be used to create soft labels to train a smaller neural network with the goal of obtaining an accuracy similar to the original larger classifier. Defensive distillation uses a variation of this technique in order to create a deep neural network that can better classify adversarial examples.

The output of the neural network used in distillation must be a Softmax layer with a temperature parameter. The output of a softmax layer is as follows

$$F(X) = \left[\frac{e^{z_i(X)/T}}{\sum_{l=0}^{N-1} e^{z_l(X)/T}} \right]_{i \in 0 \dots N-1}$$

where $Z(X)$ is the output of the previous layer and $F(X)$ is an output that corresponds to a vector with a probability for each class. The temperature value T in the Softmax layer determines how confident the classifier will be about the most probable class. A high temperature value will cause the different class probabilities to be closer together whereas a low temperature will cause the most probable classes to have a probability further apart from the less probable classes. This is because the exponent is divided by the temperature and reducing the exponents value by the same amount for all classes will cause the output values to be closer to each other

In Defensive Distillation a neural network where the last layer is a softmax layer is trained using a dataset with hard labels. The temperature of the Softmax layer is set to a high value (something greater than

1) so that the classifier will output values with larger values for each class. This means that the output will emphasize ambiguities in certain examples that may have similarities to a different class. The examples are then passed through the trained neural network at the same high temperature and the output of the classifier is recorded as a soft label for that example. Whereas in distillation a smaller model is trained using the soft labels, in defensive distillation the goal is a classifier more resilient against adversarial examples rather than a performant one so the second classifier is of the same size as the first one. When this second classifier is trained using the soft labels it the temperature of the Softmax label is set to a high value. When it is then being tested the Temperature is lowered to 1 so that the classifier outputs more confident probabilities.

Theoretical justifications presented in the paper as to why defensive distillation works against adversarial examples include that a higher temperature reduces the model's sensitivity to small perturbations in the input to the classifier, and that defensive distillation improves the generizability of the classifier outside of the training sample.

The researchers were able to use defensive distillation to lower the success rate of adversarial examples from 95.89% to 0.45% on the MNIST dataset and from 87.89% to 5.11% on the CIFAR10 dataset. The distillation did result in a moderate decrease in the accuracy on non-adversarial examples with a 1.28% decrease on the MNIST dataset and a 1.37% decrease on the CFAIR10 dataset. The models performed better against adversarial examples at higher distillation temperatures with the temperature of 100 (the highest that they tested) performing the best. They also show that defensive distillation increases the robustness of the classifier produced.

The second paper I read finds a way to revert adversarial examples into non-adversarial examples before feeding them into a classifier instead of training the classifier to handle adversarial images. This paper was *Defense-GAN: Protecting Classifiers Against Adversarial Attacks Using Generative Models* by Samangouei et all. and its idea is to use Generative Adversarial Networks to help protect Deep Neural Networks against adversarial examples.

Generative Adversarial Networks consist of two networks: D and G . D is a binary classifier for examples x , and G learns how to craft examples with the dimensionality of x from a random vector z . G learns to map z to $G(z)$ similar to x and D then learns how to distinguish between x and $G(z)$. In the paper generative networks were trained with a loss function based on Wasserstein distance which is as follows:

$$\min_G \max_D V_w(D, G) = \mathbf{E}_{x \sim p_{data}(x)}[D(x)] - \mathbf{E}_{x \sim p_z}[D(G(z))]$$

Feeding non-adversarial examples through the generative network should barely effect the examples so long as p_g converges to p_{data} . This means that legitimate example will not be altered by being fed through the network while adversarial example will be altered.

Defense-GAN works by finding an input z^* to the GAN that will match the original input x as closely as possible. $G(z^*)$ is what will then actually be fed into the classifier. The exact expression to be minimized is

$$\min_z \|G(z) - x\|_2^2$$

Gradient descent is then run on this function with random restarts to find z^*

The reason this approach is useful is because it assumes very little about the type of attack or the classifier that is being used. The classifier used to classify the images can also be trained using either the original training set or images generated by the generative network.

The researchers tested different combinations of attacks and defenses on the MNIST dataset as well as the Fashion-MNIST data set. The types of attacks included both black box and white box attacks. Defense-GAN performed well on the MNIST data set on many different types of attacks and classifiers, while other types of defenses tested performed well against only certain types of attacks. The performance from training

the classifier on the original images and the images from the Generative Adversarial Network were both comparable. Increasing the number of random restarts increased the classification accuracy. The number of iterations of gradient descent generally increased the accuracy but against adversarial examples high numbers of iterations eventually decreased the accuracy.

Defense-GAN can also be used to detect the use of adversarial examples. This is because examples with larger perturbations from the original examples will be further away from the image produced by the generative network than unchanged images. Thus the authors propose used the mean squared error of the original image and the image output by the GAN to detect attacks.

Two difficulties the authors note that may need to be considered when deploying defense-GAN are the training of adversarial networks as well as the choice of parameters. They note that there are still challenges in training GANs and the choices of L , the number of gradient descent iterations, and R , the number of random restarts, are both important factors in the effectiveness of defense-GAN.

2.5 David Meshnick's Individual Report

3 Comparative Study and Discussion

In this section, we will present a comparative study between some algorithms (and their comparable extensions) we implemented.

3.1 Overview

3.1.1 Datasets and Sample Selections

For the testing datasets, we opted to use *MINIST* [?] and *CIFAR-10* [?]. In short, *MINIST* is a database of handwritten digits and *CIFAR-10* is a database for tiny images, as demonstrated below in [Figure 5] and [Figure 6].

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9

Figure 5: A selection of *MINIST* database

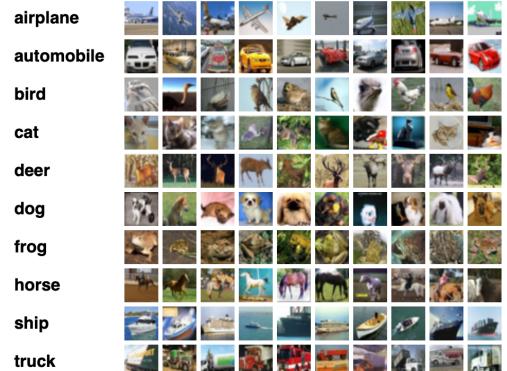


Figure 6: A selection of *CIFAR-10* database

Our decision of using these two datasets are base on the consideration of:

- Both datasets are pre-labeled and have pretrained model available for use. So we may save time on training the victim model(s), and also have some consistent baselines to refer to.

- Both datasets are lightweight in terms of image resolutions, so when evaluating computational-heavy algorithms (e.g. Hope Skip Jump), we are able to get the experiments data either locally or with minimum use of cloud services.
- Having ORC and image classifying tasks (in general) together can be a very fair coverage of the actual applications of adversarial learning.
- *MINIST* specifically has a close-to-pure background color, which is a great for human to evaluate as sometime the perturbation on a colorful picture can be unnoticeable to human eyes.
- The toolbox of choice *adversarial-robustness-toolbox* a.k.a ART have wrapper methods around these two models, and can load their pretrained models as ART's victim classifiers (we have specifically inquired Dr. Ray that it is ok to import this kind of facilities).

We have thought about using a third database like *ImageNet* or *IRIS*, but with the combinations of attacks and defenses algorithms we already have a very heavy experiments workload. So we opted to only use these two. However, considered the image quality of *MINIST* and *CIFAR-10* are rather on the low side, we used some *ImageNet* images to demo our work and concepts.

3.1.2 ART

adversarial-robustness-toolbox[?] a.k.a ART is an IBM-sponsered library that provides tools for necessary adversarial learning experiments. We have utilized ART's facilities on loading dataset, wrapping victim classifiers, and piplining attack and defense algorithms together. We cannot finish this project without this library, so much credit to them.

3.1.3 Metrics

As we are not doing binary classifications, *accuracy* will be our top priority. This is also the case for general goal of adversarial learning, as we either what to attack the model to lower its *accuracy*, or we want to defend from an attack with increased robustness – thus higher *accuracy*.

However, another important aspect of adversarial learning is, in most of the cases, we want the our attack image to be only “adversarial” to a computer model, but not to human eyes. This is first because without such restraint we may simply swap a dog picture with a cat picture and call it a successful attack, which will make the task meaningless. Second, it is because for most of the time we want our attack to be “stealthy” – and a collection of pixels noises is simply not so much of that.

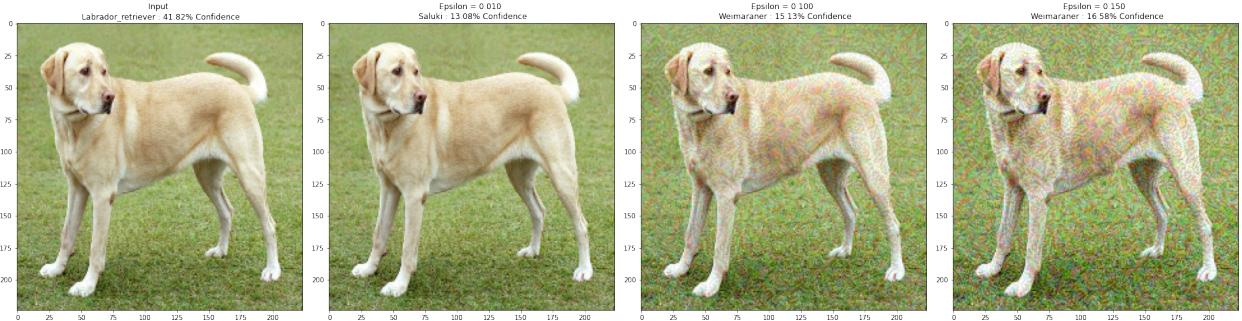


Figure 7: A Labrador Retriever with different levels of perturbations

Thus, keeping the “semantic” of the original image is important in terms of making adversarial examples. However we must have a way to quantify such perturbations – as it can be very hard to tell like shown in [Figure 7] [?] – so that we can numerically conclude and justify the claim if a image has lost its semantic (and by how much).

The metrics we found is *perturbation budget*[?], it is defined as:

$$\epsilon = \|x_{\text{adv}} - x\| \quad (2)$$

where *epsilon* in [Equation 2] is the *perturbation budget*. However, we have different choices on calculating the distance metrics. In this case, we opted to use L_2 and L_∞ as they represent the EUCLID distance and maximum magnitude between pixels – which covered the perturbation of a picture both holistically and “extremely.”

Specifically in practice, we calculate the L_2 and L_∞ of two images channel by channel then add them together. The principle is we will want these numbers to be as small as possible, as long as the adversarial image we created still have an adversarial effect to a model. Note if an algorithm is iterative, a smaller *perturbation budget* also means smaller computation cost – as less iteration were done.

Another thing that maybe worth mentioning is since we used pretrained victim model, we didn’t do anything like N-fold cross validation as we are seeking effects out of a trained model, but not to train one (maybe except Neural Cleanse as it does retain the model). But we do average our trials to make sure our experiment results are reproducible.

3.1.4 Adversarial Rivalry

We have implemented and experimented multiple attack and defense algorithms, but not all of them can be compared together due to various reasons. In our cases, we opted to not have some algorithms compared together, or even not to include some algorithm in our comparative study.

We opted to compare Fast Gradient Sign Method (non-targeted, with one-shot and iterative implementations), Hop Skip Jump, DeepFool (and Dynamic DeepFool – an extension implemented by David) together as they are all the non-targeted evasion algorithm we implemented. We opted to not include Backdoor in this set of comparison as it evasion attacks were done in the assumption of having a trained model, it doesn’t make sense to have poisoning – which can alter the training set of a model – to be part of the comparison. Similarly, the two targeted extensions FGSM is excluded from this comparison as in this context it is the magnitude of decrease of model *accuracy* in general we care about, but not which exact label the model most classifies to.

Likewise, we opted to not compare any other algorithm except Backdoor to against Neural Cleanse, as: Minyang please input here.

Also we are not comparing different defense algorithms against a same attack algorithm. This is not because we are not interested in the performance difference between defense algorithm. It is simply because we have already eliminated Neural Cleanse due to its retrain nature; and for the other two remained defense algorithm Binary Input Detector and Spatial Smoothing, the former one is designed to recognize and flag an adversarial image, where the latter is to increase model robustness so that it can better classify adversarially perturbed images – so they can’t be compared.

3.2 Attack Algorithms

Note I-FGSM represents (non-targeted, one-shot) Iterative Fast Gradient Sign Method; D-DeepFool represents Dynamic DeepFool. The FGSM is running with params being `batch_size = 32`, `eps = 3`, and the I-FGSM is running on the same setting with `eps_step = 0.05` as we have found in Henry’s individual

report with `eps_step > 0.05` we might “overshot” ourself.

Hop Skip Jump a.k.a. HSJ is a very computational-heavy algorithm. When possible, we will use `max_iter=64`, `max_eval=1000`, `init_eval=100` which is close to the authors suggested setting (except the authors suggests `max_eval = 100000` which is impossible to replicate with our resource). But often time we are limited for `max_iter=8`, `max_eval=100`, `init_eval=10`.

Unless specifically addressed, the base line model is a the pretrained model of dataset wrapped in ART’s `KerasClassifier`.

3.2.1 Evasion

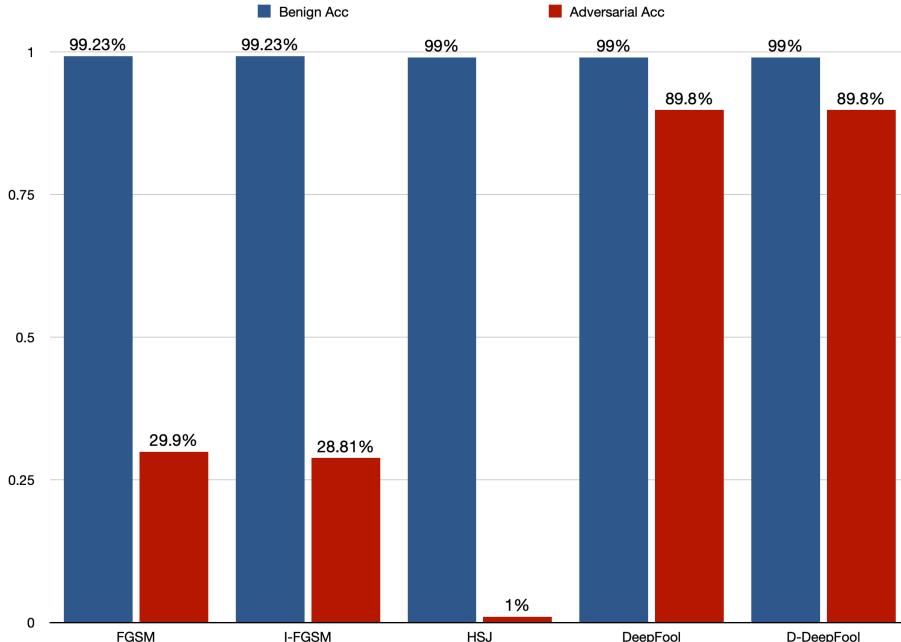


Figure 8: Accuracy comparsion of evasion attack algorithms on *MINIST*

By observing [Figure 8] it can be clearly tell that the presented algorithms are in three different “levels.” With HSJ showing the best performance and DeepFool and D-DeepFool showing identical performance, the only interesting question left is if FGSM is significantly different from I-FGSM.

We then try to find the 95% CI of $E_{\text{FGSM}} - E_{\text{I-FGSM}}$ with a null hypothesis of they have no difference. Note the sample size of tested *MINIST* is 10000.

$$F = 0.299 - 0.2881$$

$$= 0.0109$$

$$V(F) = 0.299(1 - 0.299)/10000 + 0.2881(1 - 0.2881)/10000$$

$$= 0.000041469739$$

$$\Rightarrow \sigma = 0.006439700226$$

$$95\%CI = 0.0109 \pm 1.96 \cdot 0.006439700226 = (-0.001721812443, 0.02352181244)$$

With 0 lies in the 95% CI, we can't reject the null hypothesis and FGSM and I-FGSM are not different in terms of accuracy performance in this experiment. This is consistent to our understanding of the algorithms as I-FGSM is suppose to be the iterative version of FGSM, it performed slightly "better" in number on *adversarial acc* probably just because the adversarial image is generated closer to the decision boundaries of the model and therefore a bit more "confusing" – we will analysize this issue closer with the following budget graph [Figure 9].

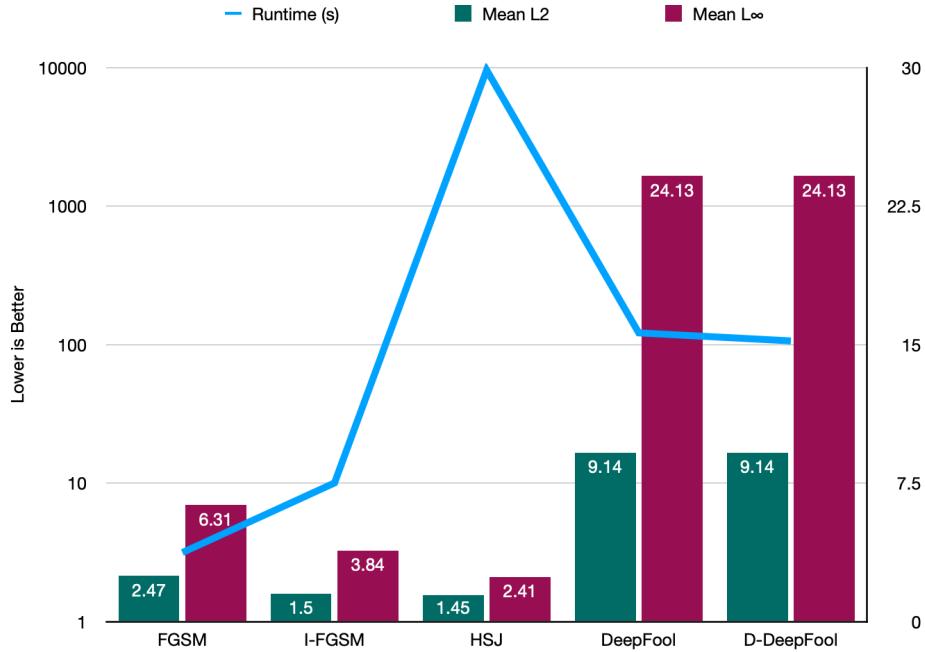


Figure 9: Budget comparsion of evasion attack algorithms on *MINIST*

By investigating the [Figure 9] we may confirm our thinking that I-FGSM generated adversarial examples are less aggressive (and thus closer to the decision boundaries), as it has a lower mean L_2 and mean L_∞ (we can't do 95% CI on this as the sample size has no bearing to the perturbation budget).

In our pervious observation on [Figure 8] we said HSJ has the best performance in terms *adversarial acc*. Now we know that HSJ did such with minimum perturbation budget spent (by having the lowest mean L_2 and mean L_∞ across the board). However, the cost of doing this is very high, the runtime of HSJ is close 100 times of other algorithm. These observations are also consistent to our understanding of HSJ, as it is doing *binary search* until it crosses the decision boundary – so it can perserve a high amount semantic from

the original image, at the cost of spending a lot of time to find such image.

David please explain a bit on DeepFool as why it perturbed a huge amount while performed not so well, and maybe also why its L_∞ is a lot higher than L_2 .

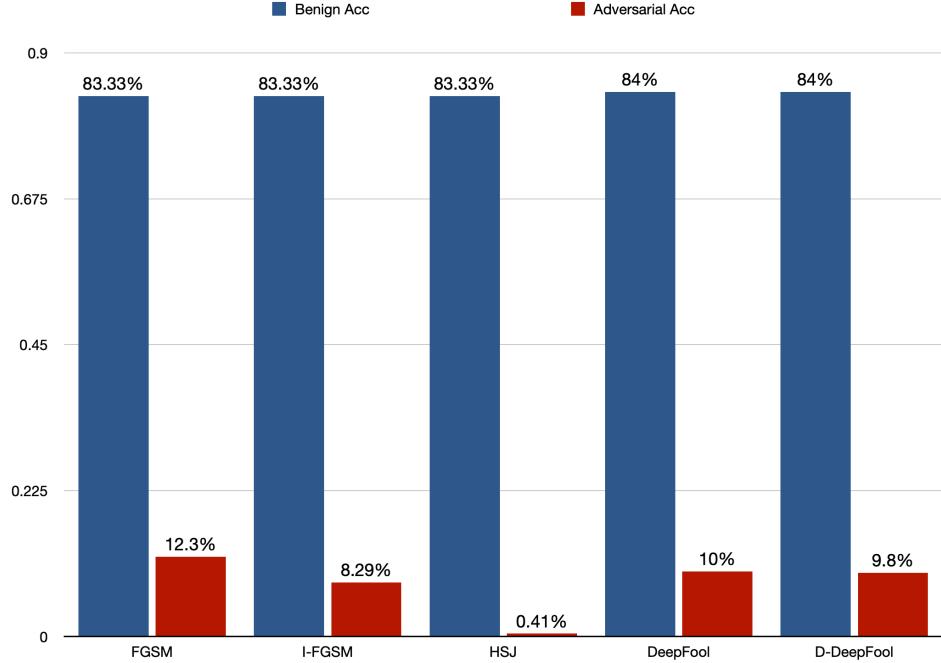


Figure 10: Accuracy comparsion of evasion attack algorithms on *CIFAR-10*

We then repeated the same experiment on 100000 *CIFAR-10* examples and got the result of [Figure 10]. This time we have a different outcome between DeepFool and D-DeepFool, so we will analysis $E_{\text{FGSM}} - E_{\text{I-FGSM}}$, $E_{\text{DeepFool}} - E_{\text{D-DeepFool}}$, and $E_{\text{I-FGSM}} - E_{\text{HSJ}}$ (as they are closer this time). With an aid of a script and a null hypothesis of having no difference, we have:

- 95% CI of $E_{\text{FGSM}} - E_{\text{I-FGSM}}$: $(0.031694853818380074, 0.04850514618161992)$, rejected.
- 95% CI of $E_{\text{DeepFool}} - E_{\text{D-DeepFool}}$: $(-0.006278442326911505, 0.010278442326911509)$, cannot rejected.
- 95% CI of $E_{\text{I-FGSM}} - E_{\text{HSJ}}$: $(0.07325244583218875, 0.08434755416781124)$, rejected.

It is a bit suprised to see the null hypothesis of $E_{\text{FGSM}} - E_{\text{I-FGSM}} = 0$ can be rejected. This is probably because *CIFAR-10* has more label catagories where *MINIST* only has 10 digits, thus decision boundaries between (more) different labels to be closer. The observation of having an overall lower *benign acc* also confirms this assumption. We will look into it again in the following budget analysis [Figure 11].

Also, with $E_{\text{I-FGSM}} - E_{\text{HSJ}} = 0$ rejected, we may say that HSJ indeed performs better than I-FGSM.

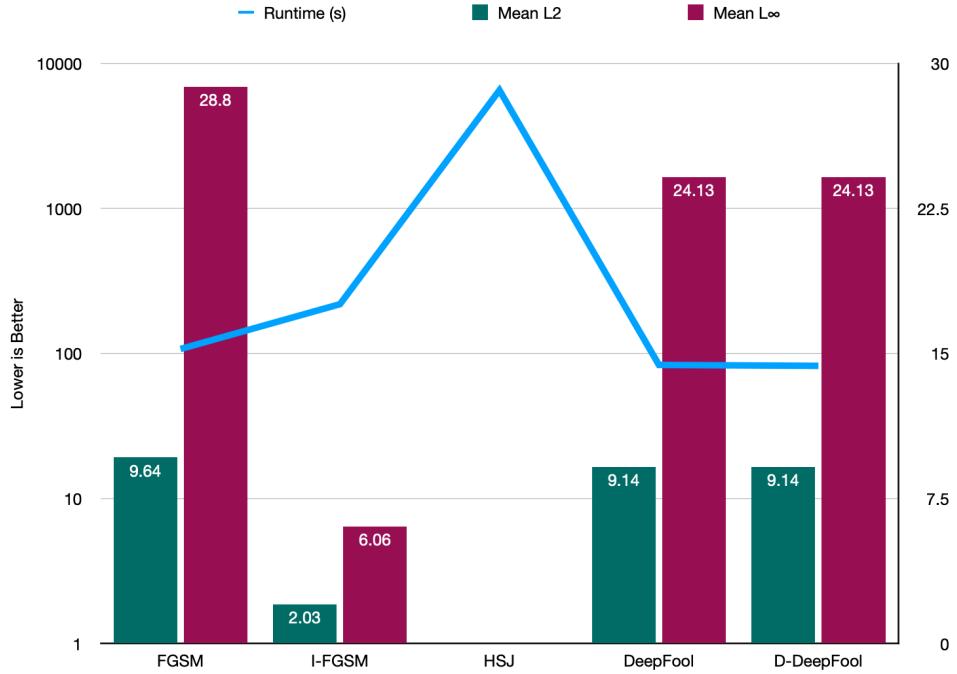


Figure 11: Budget comparsion of evasion attack algorithms on *CIFAR-10*

The budget analysis [Figure 11] confirms our thinking as the perturbation budgets are larger acrossed the board in comparision to *MINIST*. Also all other observation made from the *MINIST* budget graph [Figure 9] are also true in this *CIFAR-10* budget graph.

3.2.2 Poisoning

Please refer to Section 3.3.3 as we have only one poisoning attack algorithm and want to analysize it in combination with its designated defense: Neural Cleanse.

3.2.3 Conclusion

We have made the obersvation of HSJ is having the best *accuracy* performance. With HSJ and I-FGSM being more successful on controlling their perturbation budget – due to their iterative nature – and therefore probably are overall more reliable attacks as they can generate more effective (i.e. more confuse to a modal) adversarial examples while preserve better semantic of the original benign image (i.e. the pertubation is less obvious to human's eyes).

Note we also discovered I-FGSM and HSJ are computationally-heavy (especially the latter), due to their iterative nature and the need of many binary search operations in the case of HSJ. So it is suspected this sort of attacks can be better prevented by implementing flow control mechanism of the model predict() API – as without enough steps of iterations, these two algorithms won't easily find the decision boundary of the model.

DeepFool and D-DeepFool have shown very similar performance across this section of experiments and their performance are considerably lacking both in terms of *adversarial acc* and *perturbation budget*. David, please again summarize a bit here.

3.3 Defense Algorithms

The setup of algorithms and environment remain consistent to Section 3.2, with the exception of we ran HSJ with the params of `max_iter=8`, `max_eval=100`, `init_eval=10` due to resource concern.

Unless specifically addressed, the base line model is a the pretrained model of dataset wrapped in ART's `KerasClassifier`; and the testing sample size of datasets are 500 (per each databset).

3.3.1 Detector

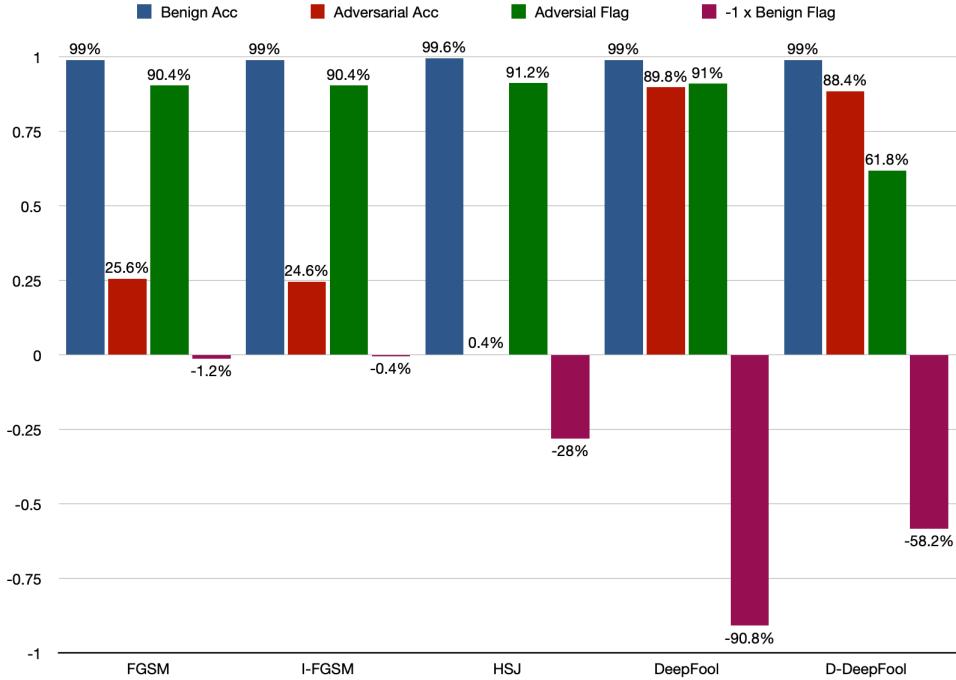


Figure 12: Effectiveness comparision of evasion attack algorithms v. Binary Input Detector on *MINIST*

Binary Input Detector a.k.a BID is algorithm that detects and flags adversarial examples out of a dataset. Since the workflow of our experiment is to get the *benign acc* with original benign examples, then make them into adversarial exmaples to get the *adversarial acc*, then we ask BID to run on both the benign example set and the adversarial example set.

In the adversarial example set, BID should aim for a 100% as every input example of the set is adversarial; vice versa, BID should aim for a 0% in the benign example set as none of the input exmaple are adversarial – we times this benign flagging percentage with an -1 do show that it is a negative impact.

First, we may tell there is much difference on *adversarial flag* except for D-DeepFool, as we have 95% of $E_{\text{HSJ}} - E_{\text{FGSM}}$ (the biggest difference on *adversarial flag* excluding D-DeepFool) to be $(-0.027824596689983806, 0.04382459668998382)$, so we cannot reject the null hypothesis of FGSM, I-FGSM, HSJ, DeepFool having no significant difference on *adversarial flag*.

Thus, the interest is left to *benign flag*, we cannot rejected $E_{\text{FGSM}} - E_{\text{I-FGSM}} = 0$ by having a 95% of $(-0.0030318578671047047, 0.019031857867104707)$. But there are significant differences between

HSJ, DeepFool, D-DeepFool on *benign flag* as we have tested the 95% of $E_{\text{HSJ}} - E_{\text{FGSM}} = 0$ (second smallested difference on *benign flag*) to be $(0.22750277615440784, 0.3084972238455922)$ on *benign flag*. We think this can be better explained by investigating the [Figure 13] graph.

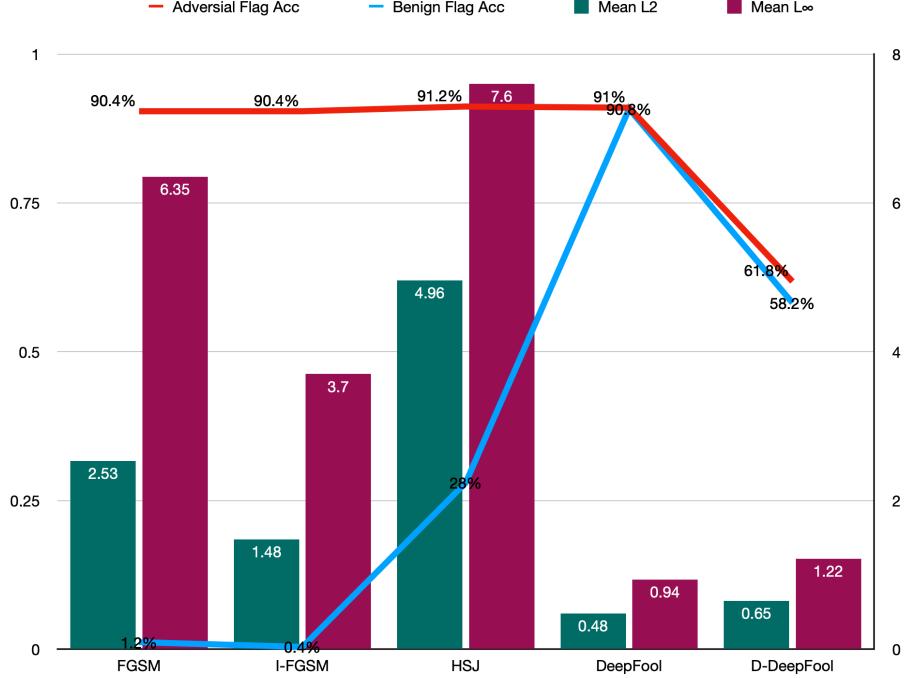


Figure 13: Budget comparision evasion attack algorithms v. Binary Input Detector on *MINIST*

By investigating the [Figure 13] graph, it can be tell there is a clearly correlation between the *perturbation budget* and the *adversarial flag* or *benign flag* (e.g., D-DeepFool). This is in fact very intuitive as less *perturbation budget* means the generated adversarial examples have perserved more semantics or their benign origins, thus making BID hard to distinguish wheather an example is benign or not.

However, it remains unknow why D-DeepFool has a much lesser *benign flag* and *adversarial flag* in comparision to DeepFool. David please have some input here thanks.

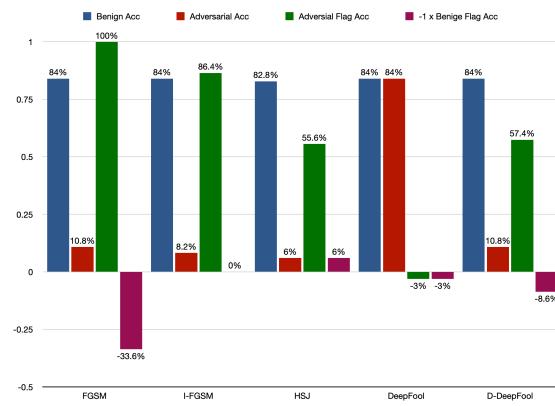


Figure 14: Effectiveness comparision

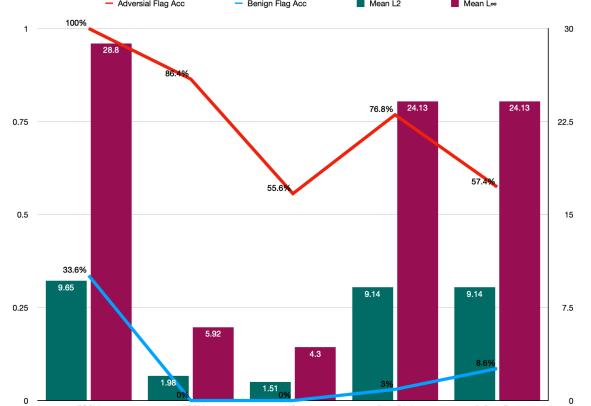


Figure 15: Budget comparision

Evasion attack algorithms v. Binary Input Detector on CIFAR-10

Our discovery continues uphold on the *CIFAR-10* dataset as algorithms with lower *perturbation budget* gets flagged less (regardless adversarial or benign).

However, this time it is I-FGSM and HSJ having the least *perturbation budget*. This is in fact resonable due to their iterative nature. We looked into our experiment and realized it is because HSJ was running on a test set of 250 examples¹, in combinations with `max_iter = 8`, the algorithm might not have enough iterations and random move to detect the decision boundaries of *MINIST*, which can be a lot harder to detect as they all have pure-color backgrounds.

3.3.2 Pre-processor

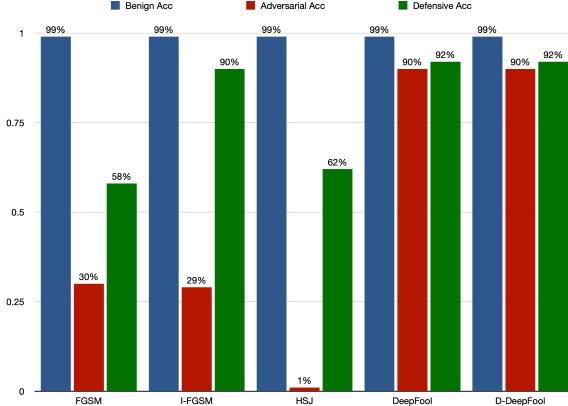


Figure 16: Effectiveness comparision

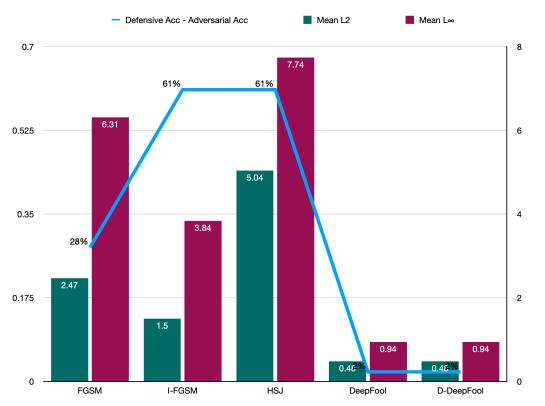


Figure 17: Budget comparision

Evasion attack algorithms v. Spatial Smoothing on MINIST

¹We have to reduce the size as BID needs to generate the adversarial images twices

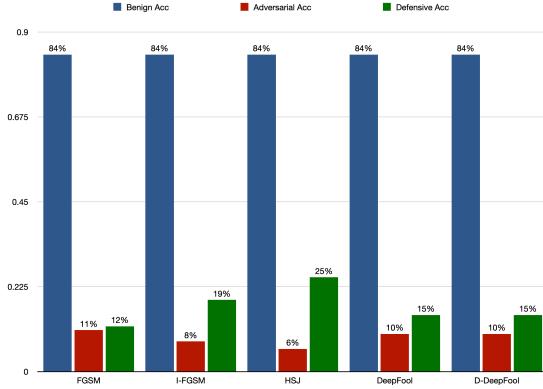


Figure 18: Effectiveness comparision

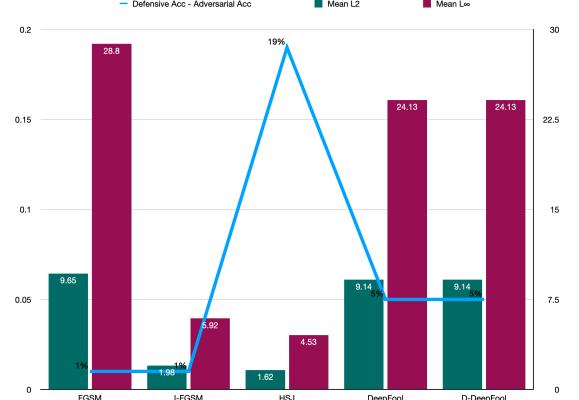


Figure 19: Budget comparision

Evasion attack algorithms v. Spatial Smoothing on *CIFAR-10*

Spatial Smoothing a.k.a. SS is our implemented defense algorithm to increase model’s robustness against adversarial input – as it actually actually help predicting the true label of an adversarial image. We may proudly say that there is a 38% *accuracy* before and after the SS being implemented across our four experimented algorithms (we exclude D-DeepFool for this discussion as it performs identical to standard DeepFool) on *MINIST*.

Note this experiment also confirms the fact that HSJ is probably not “converged” yet with its current setting on *MINIST*, as it has again costed high *perturbation budget*. We also observed the higher the *perturbation budget*, the better the defensive effectiveness – this is again in accordance with our intuition as we as human can also distinguish highly perturbed picture well.

Also note [Figure 13] seems to be suggesting low *perturbation budget* will result in high model effectiveness, this is only semi-true as the overall effectiveness increase on *CIFAR-10* is comparatively low (only +6.5% among the 4 experimented algorithms, where it is +38% in *MINIST*), so it is more of SS being effective against HSJ.

We believe this have something to do with the fact SS look into the activations of neurons of adversarial pattern and either force-zero or oppress them. As HSJ being on the decisions boundaries (implied by the low *perturbation budget*), a HSJ-generated adversarial image might be considered to have adversarial patterns of many kinds and therefore being detected by the defense algorithm.

Minyang, please correct me if I am wrong and add some insight.

3.3.3 Transformer

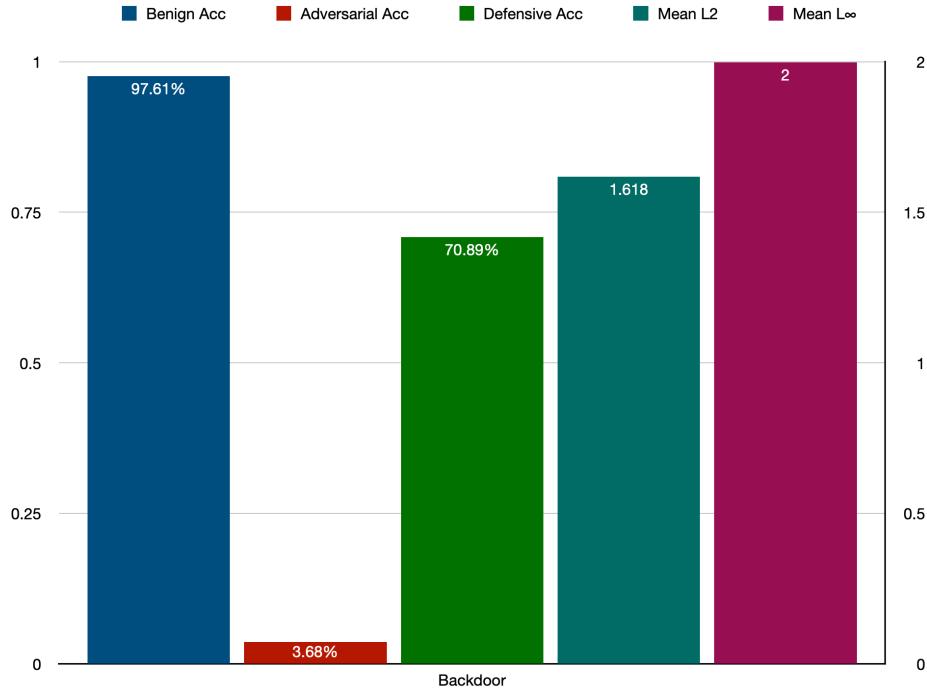


Figure 20: Backdoor v. Neural Cleanse

Minyang please explain why this works.

3.3.4 Conclusion

As each defensive algorithms have their different purposes and properties, it is hard to conclude them generally. But base on our consistent observation, it might be safe to say an adversarial example with less *budget perturbation* is likely likely to be detected by a defensive algorithm – which is consistent to our intuition and the geometrical stucture of feature space and decision boundaries: as if something is the middle of several boundaries, it will be hard to distinguish wheather it is an adversarial example of just an “outliner” of a neighborhood class.

In a parical sense, it might be best

4 References

- [1] Alex Krizhevsky et. al. The cifar-10 dataset. <https://www.cs.toronto.edu/~kriz/cifar.html>.
- [2] Beat Buesser et. al. Adversarial robustness toolbox.
<https://github.com/Trusted-AI/adversarial-robustness-toolbox>.
- [3] Yann LeCun et. al. Minist database. <http://yann.lecun.com/exdb/mnist/>.
- [4] Dan Boneh Florian Tramer. Adversarial training and robustness for multiple perturbations.
<https://github.com/Trusted-AI/adversarial-robustness-toolbox>.

- [5] Google. Adversarial example using fgsm.
https://www.tensorflow.org/tutorials/generative/adversarial_fgsm.
- [6] Aleksander Madry. A new perspective on adversarial perturbations.
<https://simons.berkeley.edu/talks/tbd-57>.