

Adversarial Attacks

Deep Learning

Ali Saffouri

University of Houston

April, 2024



1 Introduction

② Literature Review

③ Methods

4 Results

⑤ Conclusion

⑥ Appendix

1 Introduction

2 Literature Review

③ Methods

4 Results

5 Conclusion

⑥ Appendix

Introduction

Deep Neural Networks are susceptible to adversarial attacks.

Introduction

Deep Neural Networks are susceptible to adversarial attacks.

Adversarial attacks are inputs formed by applying small but intentionally worst-case perturbations to examples from the dataset, such that the perturbed input results in the model outputting an incorrect answer with high confidence.

What is Adversarial Machine Learning?

Adversarial machine learning is a machine learning method that aims to trick machine learning models by providing deceptive input.

What is Adversarial Machine Learning?

Adversarial machine learning is a machine learning method that aims to trick machine learning models by providing deceptive input.

This includes both the generation and detection of adversarial examples, which are inputs specially created to deceive classifiers.

What is Adversarial Machine Learning?

Adversarial machine learning is a machine learning method that aims to trick machine learning models by providing deceptive input.

This includes both the generation and detection of adversarial examples, which are inputs specially created to deceive classifiers.

This method has been extensively explored in areas like spam detection and image classification, where modifications are performed on images to produce incorrect predictions

1 Introduction

② Literature Review

③ Methods

4 Results

5 Conclusion

⑥ Appendix

Literature Review

Explaining and Harnessing Adversarial Examples, by Ian J. Goodfellow, Jonathon Shlens, Christian Szegedy

Literature Review

Explaining and Harnessing Adversarial Examples, by Ian J. Goodfellow, Jonathon Shlens, Christian Szegedy

1. Adversarial examples can be explained as a property of high-dimensional dot products. They are a result of models being too linear, rather than too nonlinear.
 - Previous literature theorized that adversarial examples were due to extreme nonlinearity of deep neural networks, combined with insufficient model averaging and insufficient regularization of the purely supervised learning problem
 - However, Goodfellow et al. demonstrate that linear behavior in high dimension spaces is sufficient for designing a fast method of generating adversarial examples to make adversarial training practical

Literature Review

Explaining and Harnessing Adversarial Examples, by Ian J. Goodfellow, Jonathon Shlens, Christian Szegedy

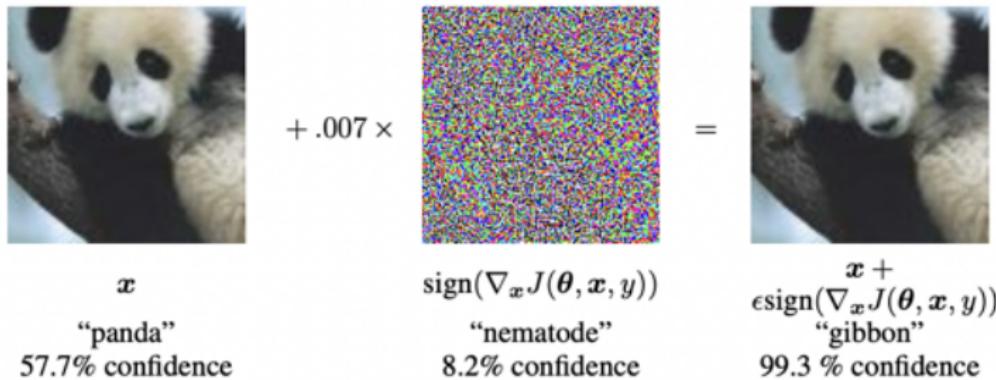
2. Fast Gradient Sign Method

- Let θ be the parameters of a model, x the input to the model, y the targets associated with x (for machine learning tasks that have targets), and $J(\theta, x, y)$ be the cost used to train the neural network. We can linearize the cost function around the current value of θ , obtaining an optimal max-norm constrained perturbation of

$$\eta = \text{sign} (\nabla_x J(\theta, x, y))$$

Literature Review

2. Fast Gradient Sign Method



- By adding an imperceptibly small vector whose elements are equal to the sign of the elements of the gradient of the cost function with respect to the input, we can change GoogLeNet's classification of the image

Literature Review

3. The generalization of adversarial examples across different models can be explained as a result of adversarial perturbations being highly aligned with the weight vectors of a model, and different models learning similar functions when trained to perform the same task.
 4. The direction of perturbation, rather than the specific point in space, matters most.
 5. Models that are easy to optimize are easy to perturb

Literature Review

6. Adversarial Training

- Can result in regularization, even further regularization than drop out
- Linear models lack the capacity to resist adversarial perturbation; only structures with a hidden layer (where the universal approximator theorem applies) should be trained to resist adversarial perturbation.
- By exposing the model to these adversarial examples during training, it learns to be more robust to such perturbations and improves its generalization performance.

Literature Review

7. Broader Implications and Future Directions

- Broader implications besides image classification include natural language processing and reinforcement learning
 - Understanding the vulnerability of machine learning models to adversarial examples is crucial for deploying reliable and trustworthy AI systems in safety-critical applications
 - Challenges remain in scaling it to complex models, as it can be computationally expensive
 - Future research directions include developing more efficient defense mechanisms, understanding the theoretical properties of adversarial examples, and exploring new paradigms for designing robust and interpretable machine learning models

1 Introduction

2 Literature Review

③ Methods

4 Results

⑤ Conclusion

⑥ Appendix

Attacking Pretrained Models

- We perform adversarial attacks on a subset of classes from ImageNet, as opposed to the literature using MNIST, to avoid retraining any models, as the models we obtained were all pretrained on ImageNet. We used the validation set provided by [Imagenette](#).
 - 10 Classes: tench, English springer, cassette player, chain saw, church, French horn, garbage truck, gas pump, golf ball, parachute
 - Imagenette Validation Set Size: 3925
 - We reclassify the classes according to their original ImageNet classification
 - We followed a tutorial on PyTorch applying [FGSM on MNIST](#)
 - We only attack models on images they classify correctly.

Attacking Pretrained Models

Model	Weights	Acc@1	Acc@5	Params	GFLOPS
GoogLeNet	GoogLeNet_Weights.IMAGENET1K_V1	69.778	89.53	6.6M	1.5
VGG11	VGG11_Weights.IMAGENET1K_V1	69.02	88.628	132.9M	7.61
Wide ResNet	Wide_ResNet50_2_Weights.IMAGENET1K_V2	81.602	95.758	68.9M	11.4
ConvNext	ConvNeXt_Tiny_Weights.IMAGENET1K_V1	82.52	96.146	28.6M	4.46

```
def fgsm_attack(image, epsilon, data_grad):
    # Collect the element-wise sign of the data gradient
    sign_data_grad = data_grad.sign()
    # Create the perturbed image by adjusting each pixel of the input image
    perturbed_image = image + epsilon*sign_data_grad
    # Adding clipping to maintain [0,1] range
    perturbed_image = torch.clamp(perturbed_image, 0, 1)
    # Return the perturbed image
    return perturbed_image
```

Retraining VGG11

- To improve the robustness of VGG11 to adversarial attacks, we partially retrain VGG11.
 - ① We generated perturbed images on all **Imagenette**'s training set (size: 9469) by attacking pretrained VGG11 using FGSM on epsilon 0.07
 - ② We froze the VGG11 feature extractor and retrained the classifier on the clean and perturbed images. The classifier architecture was not modified.
 - ③ Early stopping was implemented if the validation loss did not decrease for 3 consecutive epochs.
- We test the robustness by following our previous procedure to attack models

Sample Adversarial Attack Pipeline

Original Image



Image Transformed



Transformed Image Denormalized



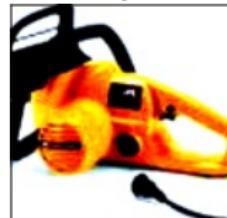
Denormalized Image Perturbed



Noise



Purturbed Image Normalized



Model: GoogLeNet, Epsilon 0.007: chain saw, 0.93 -> muzzle, 0.82

Sample Adversarial Attack

Original Image

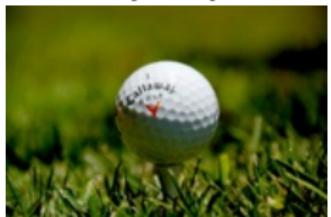


Image Transformed



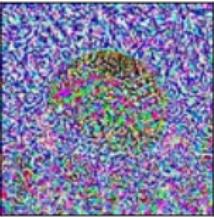
Transformed Image Denormalized



Denormalized Image Perturbed



Noise

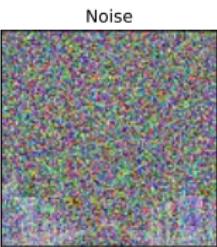
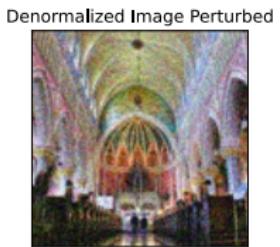
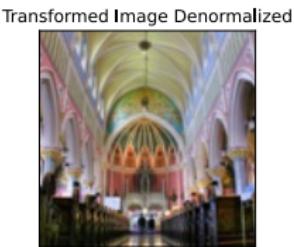
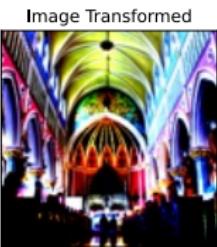


Purturbed Image Normalized



Model: VGG, Epsilon 0.07: golf_ball, 1.00 -> baseball, 0.49

Sample Adversarial Attack



Model: ConvNeXt, Epsilon 0.15: church, 0.86 -> vault, 0.35

Sample Failed Adversarial Attack

Original Image

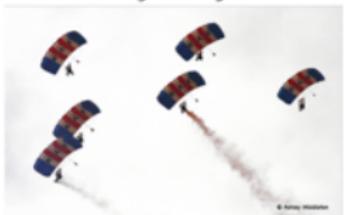


Image Transformed



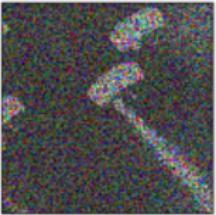
Transformed Image Denormalized



Denormalized Image Perturbed



Noise

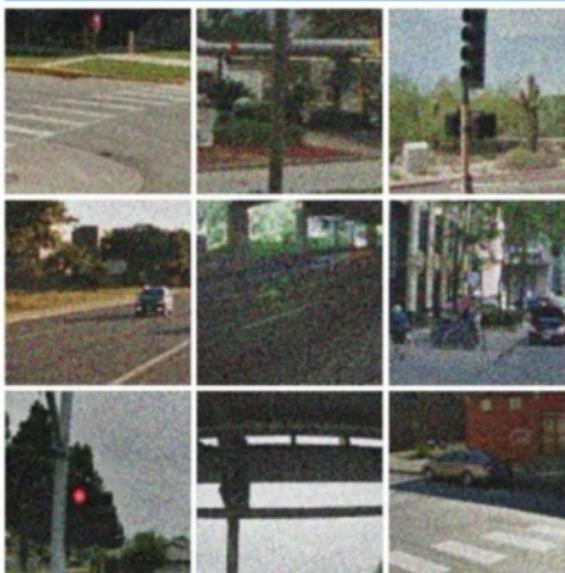


Purturbed Image Normalized



Model: ConvNeXt, Epsilon 0.15: parachute, 0.95 -> parachute, 0.89

Select all images with
cars
Click verify once there are none left



VERIFY

1 Introduction

② Literature Review

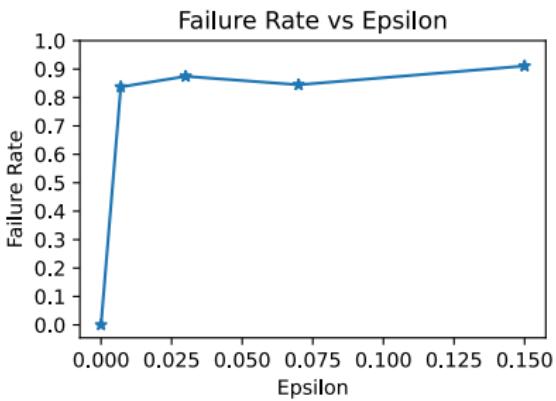
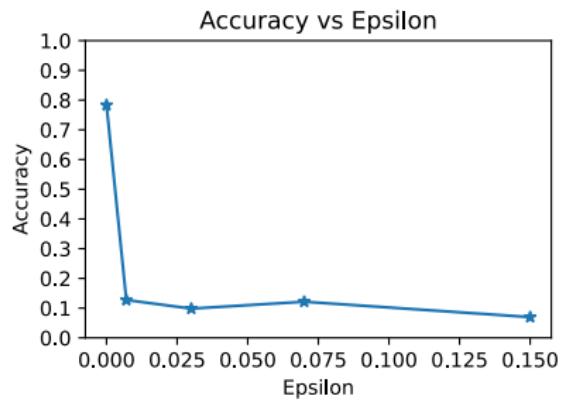
③ Methods

4 Results

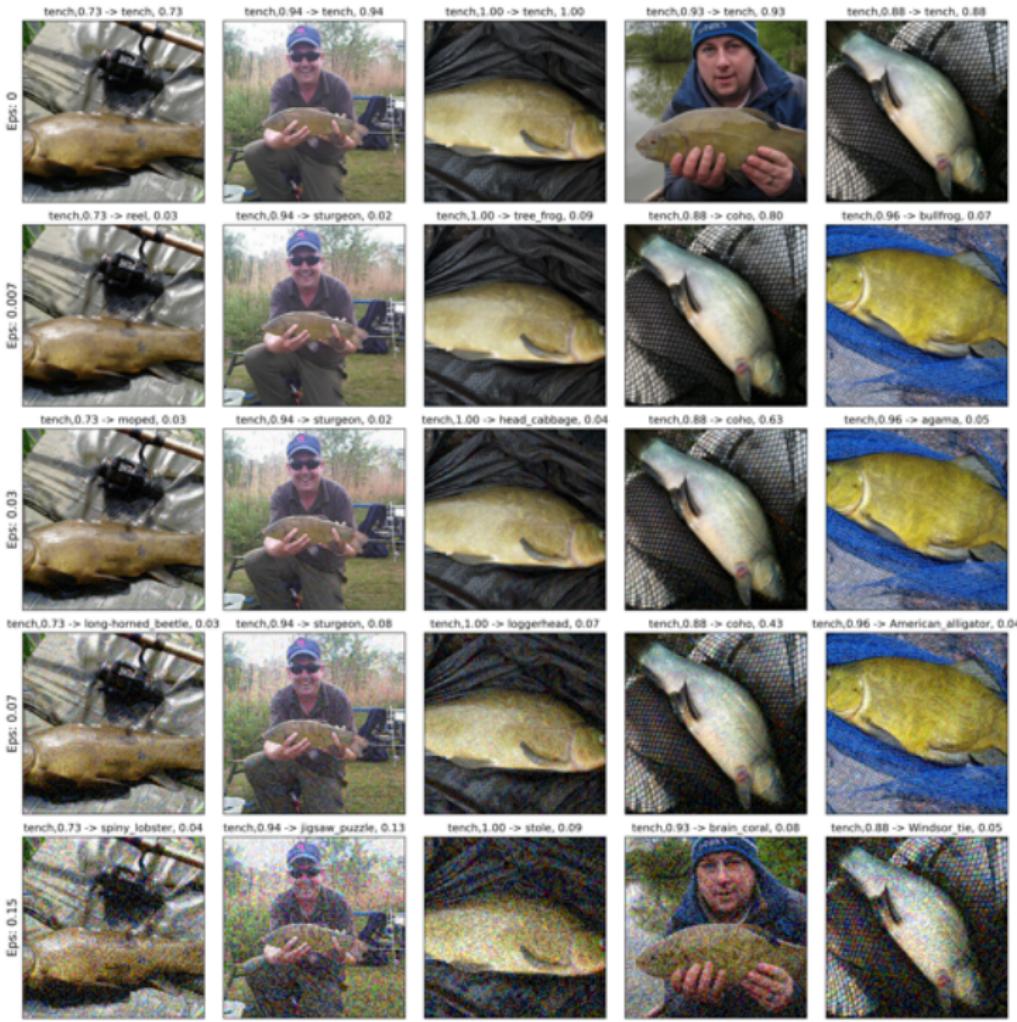
5 Conclusion

⑥ Appendix

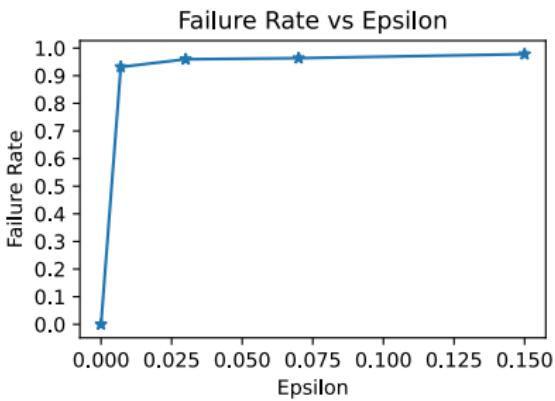
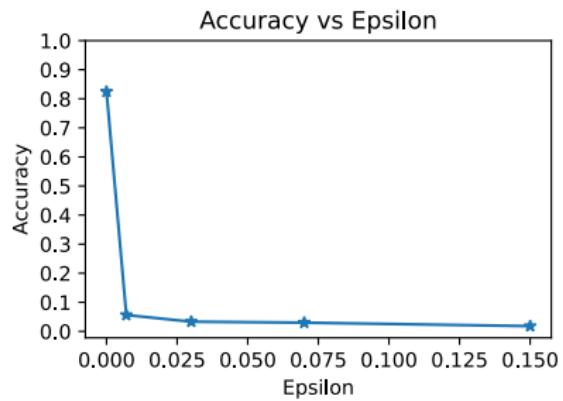
GoogleNet



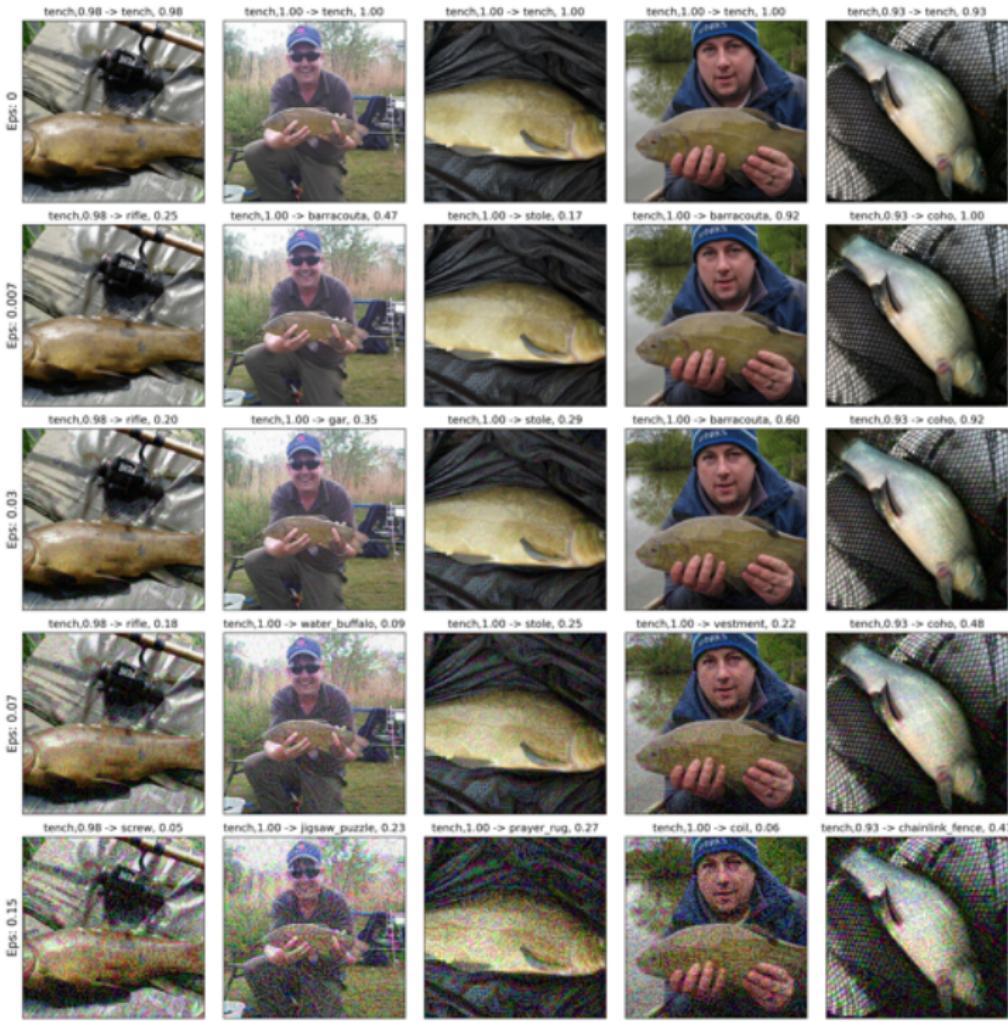
GoogleNet



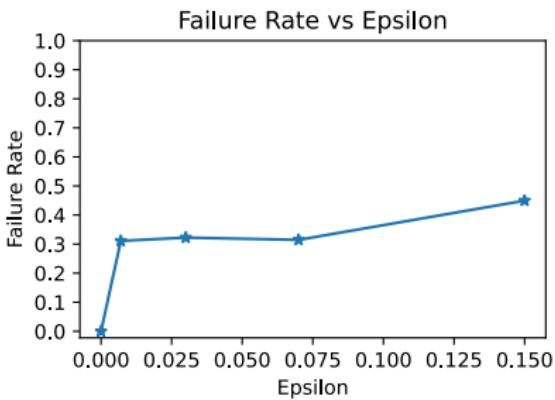
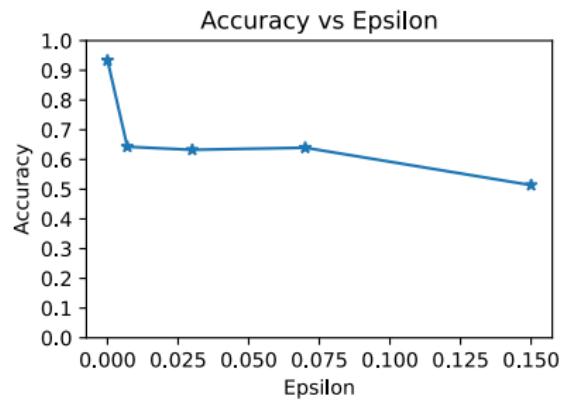
VGG11

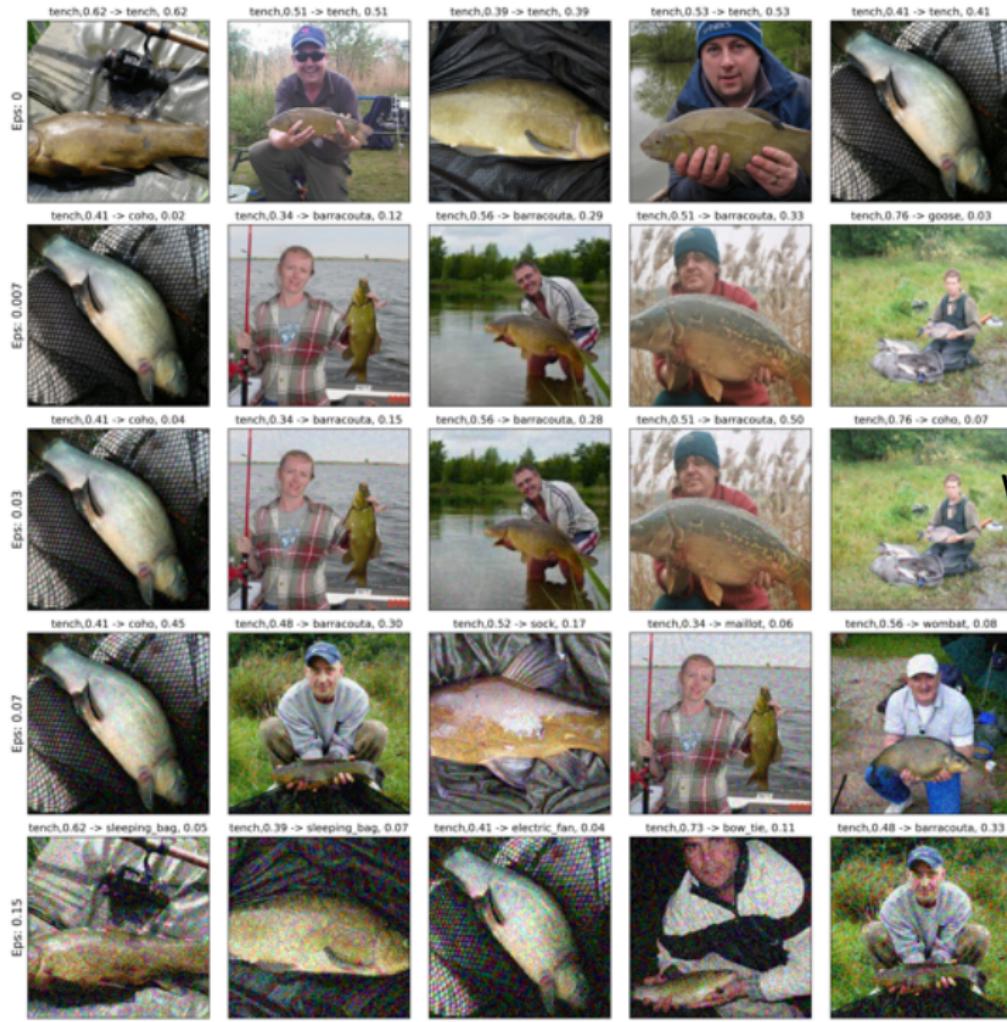


VGG11



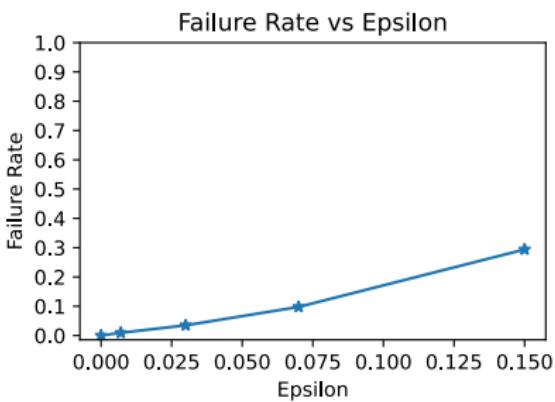
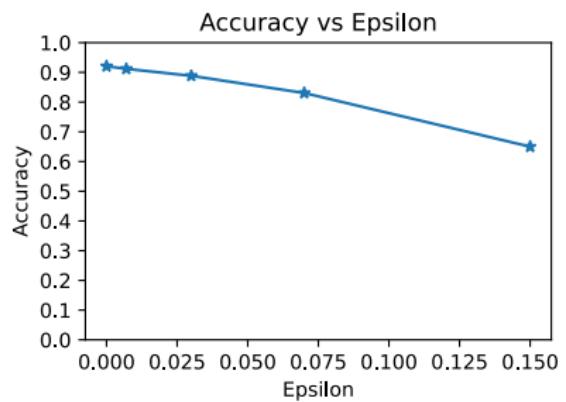
Wide ResNet



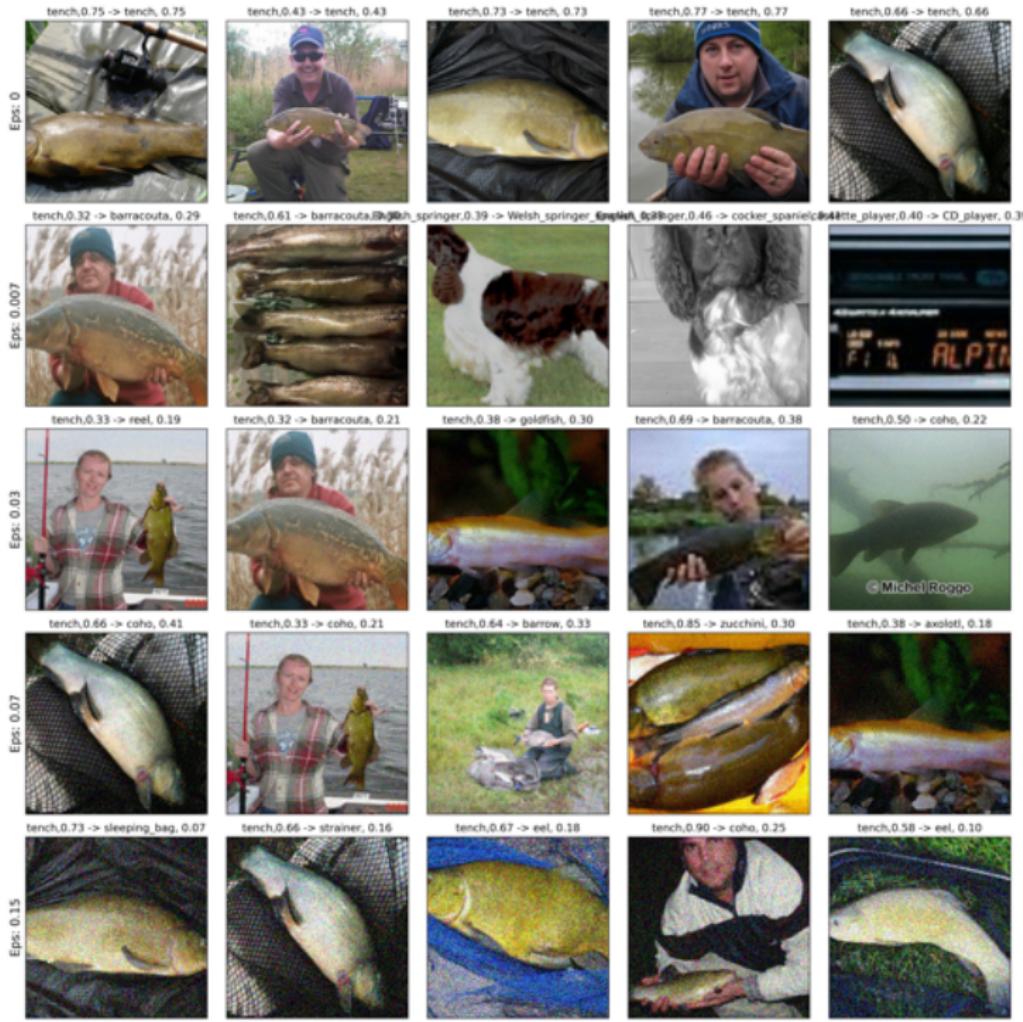


Wide ResNet

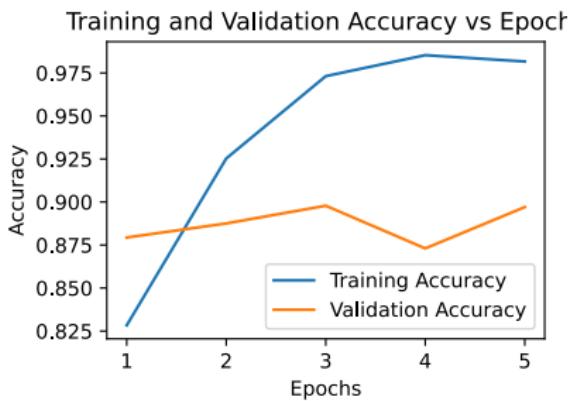
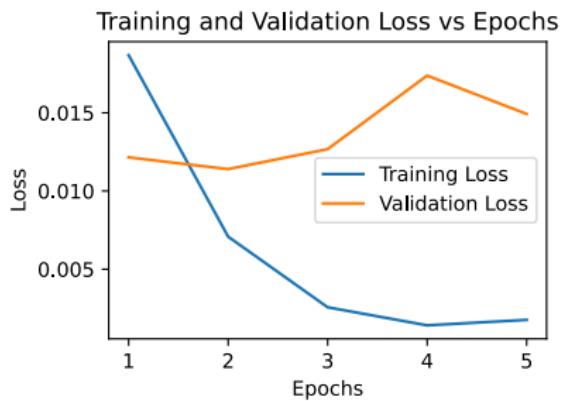
ConvNeXt



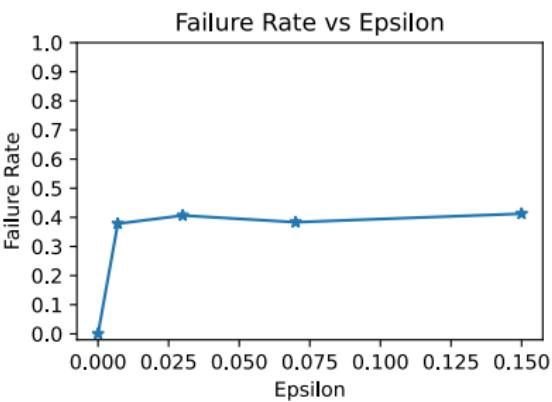
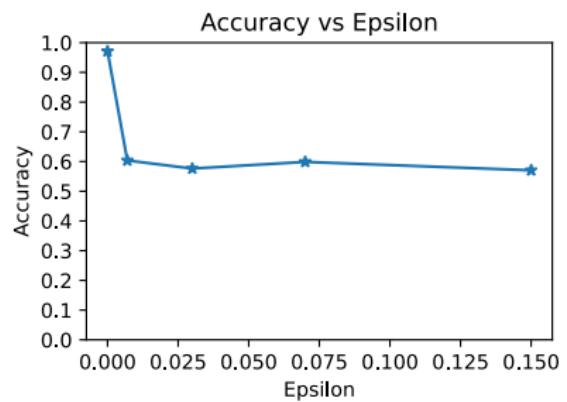
ConvNeXt



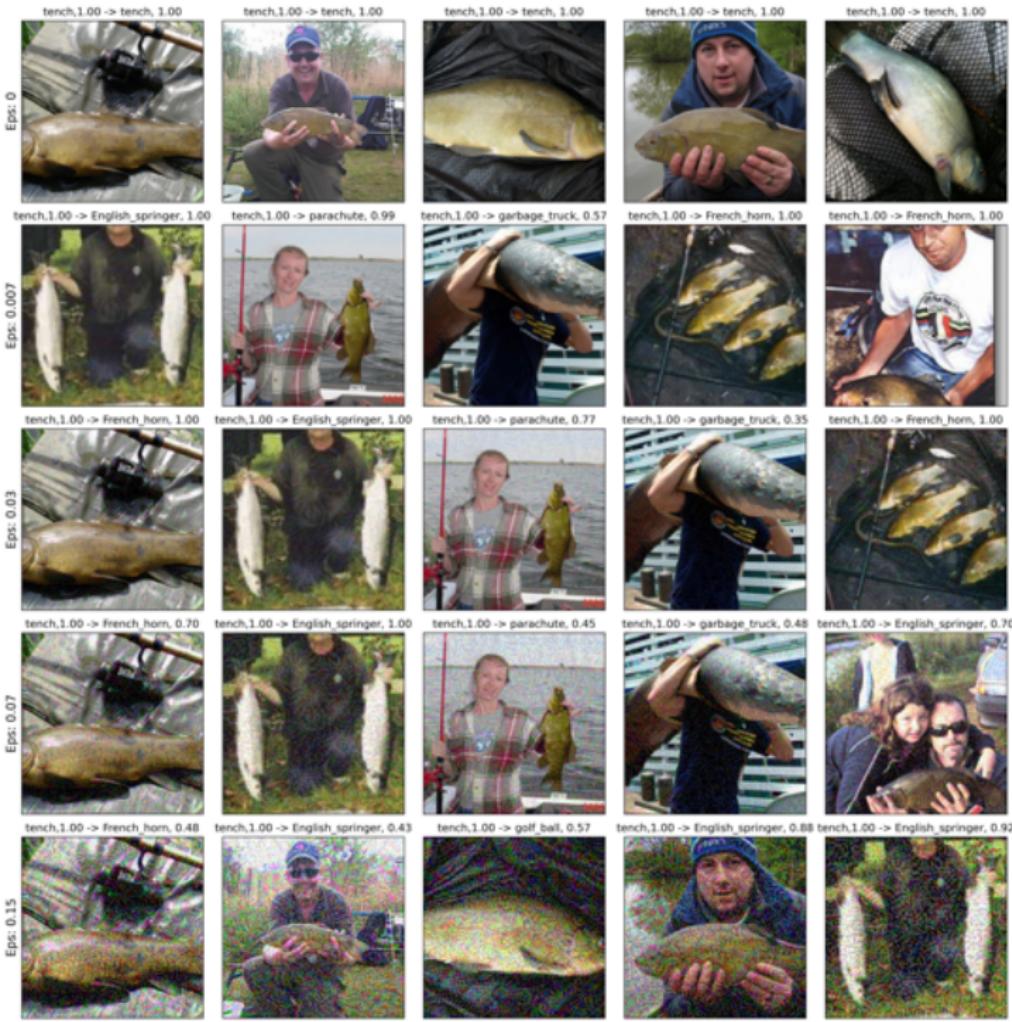
Retraining VGG11



Attacking Retrained VGG11



VGG11



1 Introduction

2 Literature Review

3 Methods

4 Results

5 Conclusion

6 Appendix

Conclusions

- GoogleNet and VGG are not robust to the FGSM adversarial attack. Even with very minimal perturbation that are not noticeable to the human eye.
 - Models more robust to FGSM, are still fooled with sufficient perturbations.
 - Retraining models with adversarial images improves its robustness.

1 Introduction

② Literature Review

③ Methods

4 Results

⑤ Conclusion

⑥ Appendix

Robustness of Retrained VGG11 on New Classes

Does the Retrained VGG11 model retain accuracy on classes not included in the test set using a different subset of ImageNet obtained from [Kaggle](#) validation set that includes 100 ImageNet classes and has a total of 3923 images.

Epsilon: 0.007, Test Accuracy = $29 / 3923 = 0.007$, Failure Rate = $122 / 151 = 0.808$

Transferability of Adversarial Examples

Can adversarial examples generated for VGG11 also deceive the other models?

Yes:

Model Name: GoogLeNet, Test Accuracy = 3213 / 9469 = 0.339

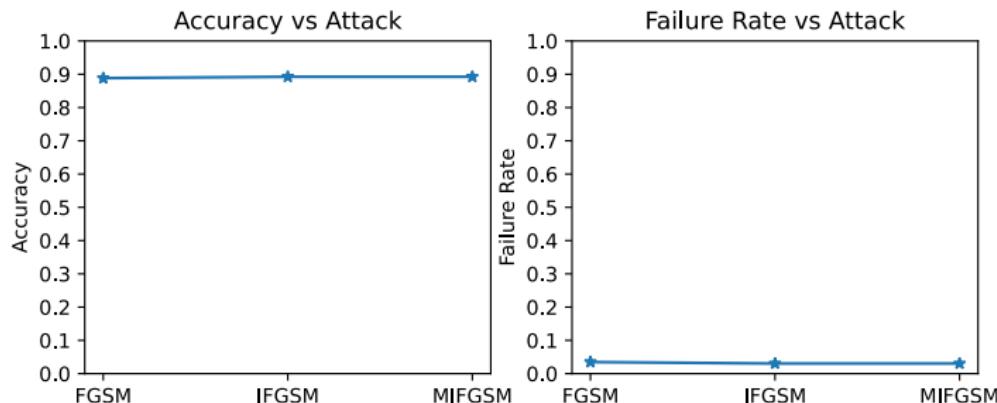
Model Name: VGG, Test Accuracy = 229 / 9469 = 0.024

Model Name: ResNet, Test Accuracy = 4997 / 9469 = 0.528

Model Name: ConvNeXt, Test Accuracy = 4158 / 9469 = 0.439

Different Attacks on ConvNeXt

We tested other attacks on ConvNeXt



Next steps: test on multistep attacks.