

Deep Learning

Chapter 7: Adversarial Examples

Mina Rezaei

Department of Statistics – LMU Munich
Winter Semester 2020



LECTURE OUTLINE

Adversarial Examples

Adversarial Training

Linear Models

Neural Networks

Adversarial Examples : A closer look (Optional)

References

ADVERSARIAL ROBUSTNESS

- It is critical to examine if a trained neural net is robust and reliable.
- **Adversarial robustness** of a model means that a model is robust to (test time) perturbations of its inputs.
- **Adversarial machine learning** studies techniques which attempt to fool machine learning models through malicious input.
- To make a model more robust, we can train our model on adversarially perturbed examples, called **adversarial examples**, derived from the training set.
- This chapter summarizes high-level ideas in adversarial robustness with a particular emphasis on adversarial examples.
- For a deeper dive, [▶ click here.](#)

Adversarial Examples

ADVERSARIAL EXAMPLES

- An adversarial example is an input to a model that is deliberately designed to "fool" the model into misclassifying it.
- The test error of a model is only an indicator of how well the model performs with respect to samples from the data-generating distribution.
- The performance of the same model can be drastically different on samples from a completely different distribution (on the same input space).
- It is possible to make changes to an image that makes a pretrained CNN (for example) output a completely different predicted class even though the change is imperceptible to the human eye.
- These examples suggest that even models that have very good test set performance do not have a deep understanding of the underlying concepts that determine the correct output label.

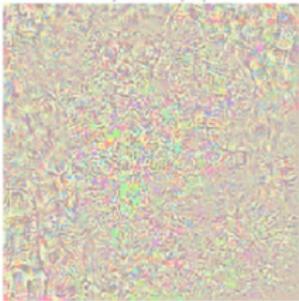
ADVERSARIAL EXAMPLES

- Adversarial examples are *not* unique to deep neural nets. Many other models (such as logistic regression) are also susceptible to them.
- They pose serious security concerns in many areas.
- Example: Fooling autonomous cars into thinking that a stop sign is a 45 km/h sign.
- Example: Evading law enforcement by fooling facial recognition systems into misidentifying individuals.

ADVERSARIAL EXAMPLES



"golden retriever"
80.6 % confidence



added perturbation
(zoomed in by a factor of 10)



"plane"
99 % confidence

Credit: Maxence Prevost

Figure: The difference between the left and the right golden retriever is imperceptible to humans. The last image was classified as a plane by ResNet50 with more than 99% confidence. The only difference between the left and right image are small pixel perturbations, shown in the second picture.

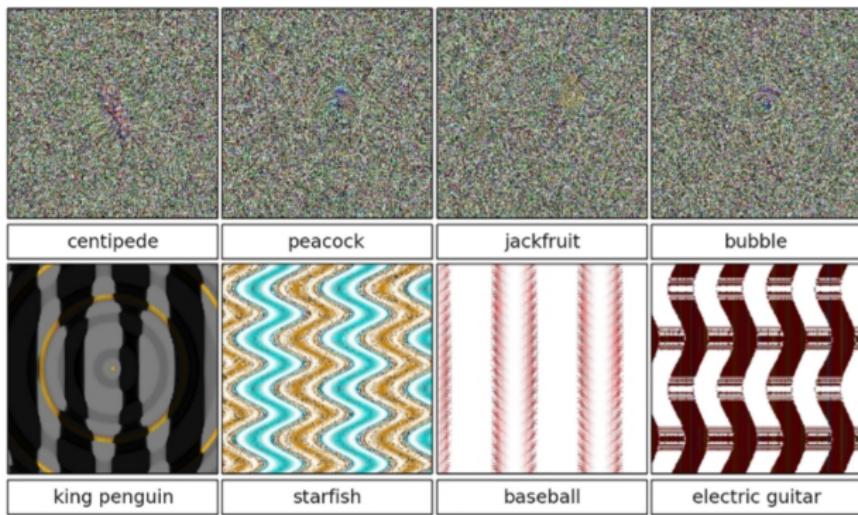
ADVERSARIAL EXAMPLES



Credit: Sharif et al.

Figure: A CNN misidentified each person in the top row (with the funky looking "adversarial" glasses) as the one in the corresponding position in the bottom row. The generated images do often contain some features of the target class.

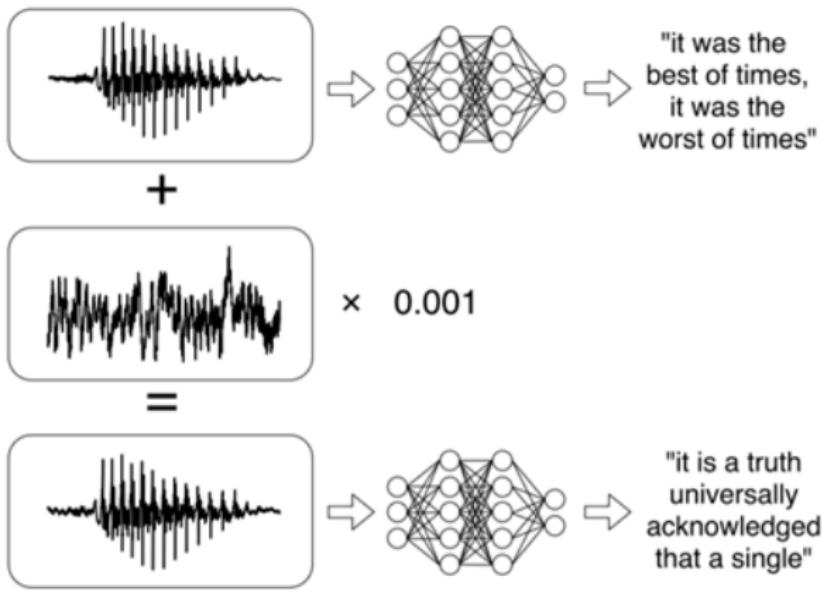
ADVERSARIAL EXAMPLES



Credit: Nguyen et al.

Figure: All 8 images above are unrecognizable to humans but are misclassified by a CNN with higher than 99% confidence. The CNN was trained by Krizhevsky et al. on the ImageNet dataset and consists of five convolutional layers, some of which are followed by max-pooling layers, and three fully-connected layers with a final softmax.

ADVERSARIAL EXAMPLES



Credit: Carlini et al

Figure: It is possible to add a small perturbation to any waveform in order to fool a speech-to-text neural network into transcribing it as any desired target phrase.

CREATION OF ADVERSARIAL EXAMPLES

- In the examples earlier, we saw that adversarial examples can seem recognizable to humans or seem like random noise/patterns.
- In the following, given a datapoint \mathbf{x} , we want to create an adversarial example $\tilde{\mathbf{x}}$ that is very similar to \mathbf{x} .
- Specifically, our goal is to find an input $\tilde{\mathbf{x}}$ close to the datapoint \mathbf{x} such that a pretrained model (which accurately classifies \mathbf{x}), ends up misclassifying $\tilde{\mathbf{x}}$.
- When we train a neural network, we typically want to optimize the parameter θ , so that we minimize the loss.
- By contrast, to find an adversarial example $\tilde{\mathbf{x}}$ that is in the vicinity of \mathbf{x} , we want to optimize the *input* to *maximize* the loss.

CREATION OF ADVERSARIAL EXAMPLES

- To ensure that $\tilde{\mathbf{x}}$ is close to \mathbf{x} , we optimize over the perturbation of \mathbf{x} , denoted as δ , and define a feasible set of perturbations Δ .

$$\arg \max_{\delta \in \Delta} L(y, \hat{f}(\mathbf{x} + \delta | \theta))$$

- A common perturbation set is $\mathcal{B}_\epsilon^\infty$ which is the ϵ -ball measured by $\ell_\infty = \|\cdot\|_\infty$

$$\Delta = \mathcal{B}_\epsilon^\infty(\delta) = \{\delta : \|\delta\|_\infty \leq \epsilon\} \text{ with } \|\delta\|_\infty = \max_i |\delta_i|$$

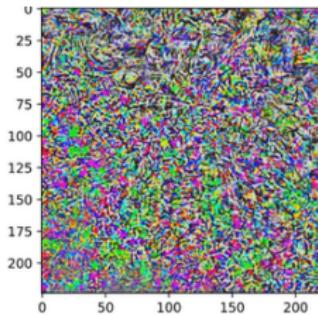
It allows each component of the perturbation δ to lie between $-\epsilon$ and $+\epsilon$.

- In general, Δ can also depend on the input datapoint \mathbf{x} , denoted as $\Delta(\mathbf{x})$.

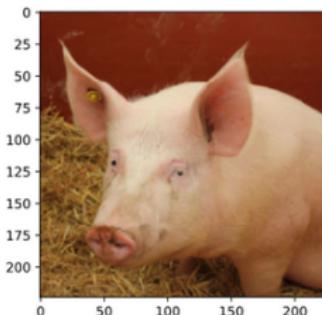
EXAMPLE: RESNET50



"pig"



added perturbation
(heavily zoomed in by a factor of 50)



"wombat"
99.6 % confidence

Credit: Kolter & Madry

Figure: Adversarial example for one datapoint of the ImageNet dataset and pre-trained ResNet50. By adding an imperceptibly small perturbation to the original image, an image was created that looks identical to our original image, but is misclassified.

TARGETED ATTACKS

- It is also possible to generate adversarial examples classified virtually as any desired class. This is known as a **targeted attack**.
- The only difference is that, instead of trying to just maximize the loss of the correct class, we maximize the loss of the correct class while also minimizing the loss of a target class y_{target} .

$$\arg \max_{\delta \in \Delta} (L(y, \hat{f}(\mathbf{x} + \delta | \theta)) - L(y_{target}, \hat{f}(\mathbf{x} + \delta | \theta)))$$

Adversarial Training

ADVERSARIAL TRAINING

- To modify a trained model so that it is more resistant to such attacks, adversarial training can be performed.
- To do so, we minimize the **empirical adversarial risk** which measures the worst-case empirical loss of a model, if we are able to manipulate every input \mathbf{x} in the training data set within the feasible set $\Delta(\mathbf{x})$:

$$\min_{\theta} \mathcal{R}_{adv}(\theta) = \min_{\theta} \frac{1}{N} \sum_{i=1}^N \max_{\delta \in \Delta(\mathbf{x})} L(y^{(i)}, f(\mathbf{x}^{(i)} + \delta | \theta))$$

ADVERSARIAL TRAINING

- To solve the optimization problem, we use SGD over θ . In each SGD step $t \in \{1, 2, \dots\}$ we repeatedly choose a minibatch of size m and repeat the following until a stopping criterion is met:
 - For each $(\mathbf{x}^{(i)}, y^{(i)}), i = 1, \dots, m$, we compute an adversarial example

$$\delta^*(\mathbf{x}^{(i)}) = \arg \max_{\delta \in \Delta(\mathbf{x}^{(i)})} L(y^{(i)}, f(\mathbf{x}^{(i)} + \delta | \theta^{[t]}))$$

- Then we compute the gradient of the empirical adversarial risk given $\delta^* = (\delta^*(\mathbf{x}^{(1)}), \dots, \delta^*(\mathbf{x}^{(m)}))$ and update θ :

$$\theta^{[t+1]} := \theta^{[t]} - \alpha \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L(y^{(i)}, f(\mathbf{x}^{(i)} + \delta^*(\mathbf{x}^{(i)}) | \theta^{[t]}))$$

- The first step is derived from Danskin's theorem, which states that the gradient of the inner function (maximization term) is simply given by the gradient of the function evaluated at its maximum.

Linear Models

LINEAR MODELS

- In case of linear models, the inner maximization problem can be solved exactly. We show this in the case of binary classification using linear models.
- Recall, the hypothesis space for logistic regression consists of models of the form:

$$\mathcal{H} = \left\{ f : \mathbb{R}^p \rightarrow [0, 1] \mid f(\mathbf{x}) = \tau \left(\sum_{j=1}^p \theta_j x_j + \theta_0 \right), \boldsymbol{\theta} \in \mathbb{R}^p, \theta_0 \in \mathbb{R} \right\},$$

where $\tau(z) = (1 + \exp(-z))^{-1}$ is the logistic sigmoid function.

- For class labels $y \in \{+1, -1\}$, the logistic loss is:

$$L(y, f(\mathbf{x} \mid \boldsymbol{\theta})) = \log(1 + \exp(-y(\underbrace{\sum_{j=1}^p \theta_j x_j + \theta_0}_{\boldsymbol{\theta}^T \mathbf{x}}))) \equiv \Psi(y(\boldsymbol{\theta}^T \mathbf{x} + \theta_0))$$

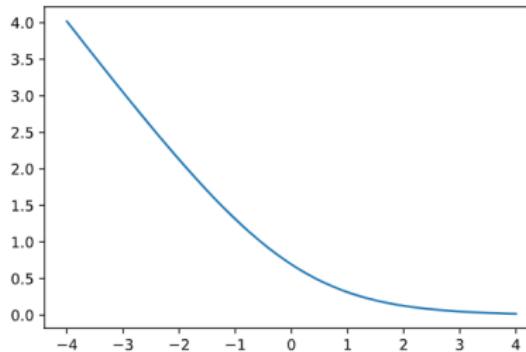
where we define $\Psi(z) = \log(1 + \exp(-z))$.

LINEAR MODELS

- The inner maximization in the adversarial risk, which we saw earlier, can be written as:

$$\max_{\delta \in \Delta(\mathbf{x})} L(y, f(\mathbf{x} + \delta | \theta)) = \max_{\delta \in \Delta(\mathbf{x})} \Psi(y(\theta^T(\mathbf{x} + \delta) + \theta_0))$$

- In this particular case, it is possible to solve the inner maximization exactly.
- First, note that Ψ is a monotonically decreasing function.



Credit: Kolter and Madry

LINEAR MODELS

- Maximizing such a monotonically decreasing function is equivalent to minimizing the argument.
- Therefore

$$\begin{aligned}\max_{\delta \in \Delta(\mathbf{x})} \Psi(y(\theta^T(\mathbf{x} + \delta) + \theta_0)) &= \Psi\left(\min_{\delta \in \Delta(\mathbf{x})} y(\theta^T(\mathbf{x} + \delta) + \theta_0)\right) \\ &= \Psi(y(\theta^T \mathbf{x} + \theta_0) + \min_{\delta \in \Delta(\mathbf{x})} y(\theta^T \delta))\end{aligned}$$

- We have to solve the problem

$$\min_{\delta \in \Delta} y(\theta^T \delta)$$

LINEAR MODELS

- To get a feel for the problem, let us consider the case where $y = +1$ and use $\Delta = \mathcal{B}_\epsilon^\infty$. The latter constrains each element of δ to lie between $-\epsilon$ and $+\epsilon$.
- The quantity $y(\theta^T \delta)$ is then minimized when $\delta_j = -\epsilon$ for $\theta_j \geq 0$ and $\delta_j = \epsilon$ for $\theta_j < 0$.
- For $y = -1$, the signs would be flipped.
- The optimal solution then, is

$$\delta^* = -y\epsilon \cdot \text{sign}(\theta)$$

- Note that the optimal solution does not explicitly depend on \mathbf{x} .

LINEAR MODELS

- The function value achieved by the solution is:

$$y \cdot \theta^T \delta^* = y \cdot \sum_j -y\epsilon \cdot \text{sign}(\theta_j)\theta_j = -y^2\epsilon \sum_j |\theta_j| = -\epsilon \|\theta\|_1$$

- Therefore, we have analytically computed the solution to the inner maximization problem! The solution is:

$$\max_{\delta \in \Delta(\mathbf{x})} \Psi(y(\theta^T(\mathbf{x} + \delta) + \theta_0)) = \Psi(y(\theta^T(\mathbf{x} + \delta)) - \epsilon \|\theta\|_1)$$

- As a result, the adversarial risk, which was a min-max problem, has now been converted to a pure minimization problem:

$$\min_{\theta, \theta_0} \frac{1}{N} \sum_{i=1}^N \Psi \left(y^{(i)} \cdot (\theta^T \mathbf{x}^{(i)} + \theta_0) - \epsilon \|\theta\|_1 \right)$$

- This problem is convex in $\{\theta, \theta_0\}$ and can be solved exactly. An iterative optimizer such as SGD will also approach the global minimum.

MNIST EXAMPLE

- As an example, we look at the MNIST dataset, but this time we perform logistic regression and focus only on the classification of 0s vs. 1s.
- The logistic regression classifier was trained for 10 epochs with SGD on the training set.
- This model obtained a low misclassification rate of 0.0004 on the test set.
- To generate adversarial examples, Δ is defined as $\mathcal{B}_{0.2}^\infty$.
- As we saw earlier, the optimal perturbation δ^* is $-y\epsilon \cdot \text{sign}(\theta)$.

MNIST EXAMPLE

- As δ^* does not directly depend on x , it is the "same" (ignoring the value of label y) across all examples. This is what it looks like:

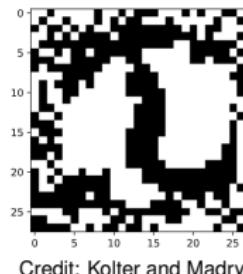
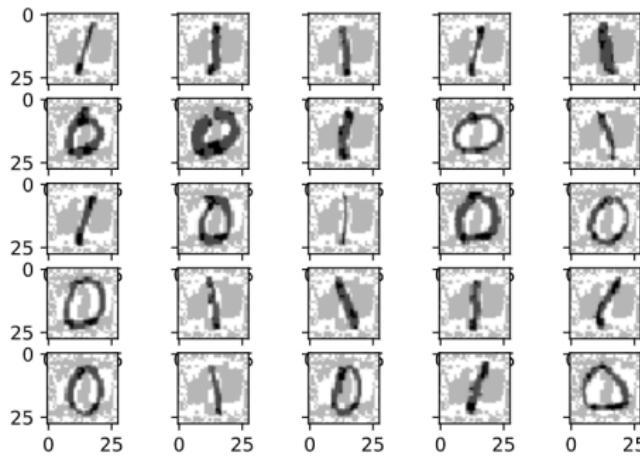


Figure: The optimal perturbation for images that contain 0. For images that contain 1, the signs would be flipped. (The contrast between the black and white pixels is amplified for the sake of visualization.)

- The perturbation (*vaguely*) has a vertical line (like a 1) in black pixels, and a circle (like a 0) in white pixels. Intuition: When a given image is moved (translated) in the black direction, it is more likely to be classified as 1, whereas when moved in the white direction, it is more likely to be classified as 0.

MNIST EXAMPLE



Credit: Kolter and Madry

Figure: Perturbed images from the test set.

- When all the images in test set are perturbed, the misclassification error of the model jumps from 0.0004 to 0.845!
- Interestingly, when the model is trained on similarly perturbed images from the training set (that is, the empirical adversarial risk is minimized), the misclassification error on the perturbed test set drops to 0.025.

Neural Networks

PROJECTED GRADIENT DESCENT

- In contrast to logistic regression, neural networks can have a bumpier, non-convex loss surface.
- As a consequence, Danskin's theorem does not longer hold and the inner optimization problem must be solved approximately.
- One approximation method is projected gradient descent (PGD)¹.
- Let \hat{f} be the pretrained model, \mathbf{x} the input to the model, y the target and $L(y, \hat{f}(\mathbf{x}|\theta))$ the loss function used to train the network.
- In each gradient descent step, the basic PGD algorithm updates δ as:

$$\delta^{[t+1]} = \mathcal{P} \left(\delta^{[t]} + \alpha \nabla_{\delta} L \left(y, \hat{f}(\mathbf{x} + \delta^{[t]} | \theta) \right) \right)$$

where \mathcal{P} denotes the projection onto the ball of interest.

¹ Technically speaking, it is gradient *ascent* since we are maximizing a function, but for the sake of generality, we just refer to the process here as gradient descent.

PROJECTED GRADIENT DESCENT

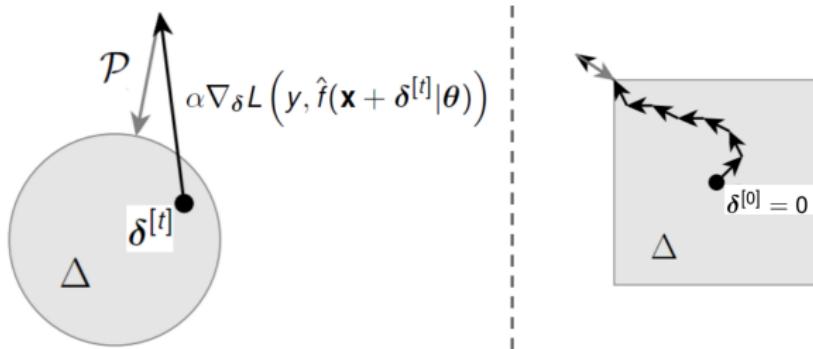
- In essence, $\delta^{[t+1]}$ is obtained by moving from $\delta^{[t]}$ (with a step size α) **in the direction** of the gradient of the loss with respect to δ (evaluated at $\delta^{[t]}$) and then projecting back to Δ .
- In case where the feasible set is $\mathcal{B}_\epsilon^\infty$, the projection of an arbitrary vector \mathbf{z} is simply:

$$\mathcal{P}(\mathbf{z}) = \text{clip}(\mathbf{z}, [-\epsilon, \epsilon])$$

- Therefore, when Δ is $\mathcal{B}_\epsilon^\infty$, one gradient step is then defined as:

$$\delta^{[t+1]} = \text{clip}\left(\delta^{[t]} + \alpha \nabla_\delta L\left(y, \hat{f}(\mathbf{x} + \delta^{[t]} | \boldsymbol{\theta})\right), [-\epsilon, \epsilon]\right)$$

PROJECTED GRADIENT DESCENT



Credit: Kolter and Madry

Figure: Left: $\delta^{[t]} + \alpha \nabla_{\delta} L(y, \hat{f}(\mathbf{x} + \delta^{[t]} | \theta))$ is projected back to the perturbation set Δ using \mathcal{P} . Here, Δ is \mathcal{B}_{ϵ}^2 . Right: Multiple steps of PGD are shown in case of $\Delta = \mathcal{B}_{\epsilon}^{\infty}$; the projection must be only executed in the last step. (Note: A variant of PGD known as *normalized* PGD is shown here. This is why each step has the same size. See Kolter et al. (2019) for details.)

FAST GRADIENT SIGN METHOD

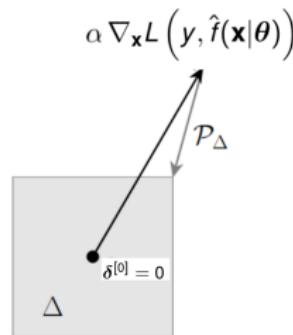
- Fast Gradient Sign Method (FGSM) is a special case of PGD when $\Delta = \mathcal{B}_\epsilon^\infty$ and $\alpha \rightarrow \infty$.
- As before, the projection of an arbitrary vector \mathbf{z} onto $\mathcal{B}_\epsilon^\infty$ is $\mathcal{P}(\mathbf{z}) = \text{clip}(\mathbf{z}, [-\epsilon, \epsilon])$.
- As $\alpha \rightarrow \infty$, the elements of δ are either set to $-\epsilon$ or ϵ depending on the sign of the corresponding component of the gradient.
- Thus, the FGSM algorithm computes δ as

$$\delta = \epsilon \text{sign} \left(\nabla_{\mathbf{x}} L \left(y, \hat{f}(\mathbf{x}|\theta) \right) \right)$$

- Note that for FGSM we only conduct one calculation and not multiple gradient descent steps.
- Furthermore, because we usually initialize $\delta^{[0]}$ as 0

$$\nabla_{\delta} L \left(y, \hat{f}(\mathbf{x} + \delta|\theta) \right) = \nabla_{\mathbf{x}} L \left(y, \hat{f}(\mathbf{x}|\theta) \right)$$

FAST GRADIENT SIGN METHOD



Credit: Kolter and Madry

Figure: δ is obtained by setting each element of $\nabla_{\mathbf{x}} L(y, \hat{f}(\mathbf{x}|\theta))$ to $-\epsilon$ or ϵ depending on its sign. Note that this slightly changes the direction of the step that is taken.

- Implicitly, the (element-wise) *sign* function is simply a way to constrain the size of the "step" that we take in the direction of the gradient. It is basically equal to a projection of the step back on Δ (which is $\mathcal{B}_\epsilon^\infty$).

FAST GRADIENT SIGN METHOD

- Note: The optimal attack against the linear binary classifier we saw in the last section was also FGSM!

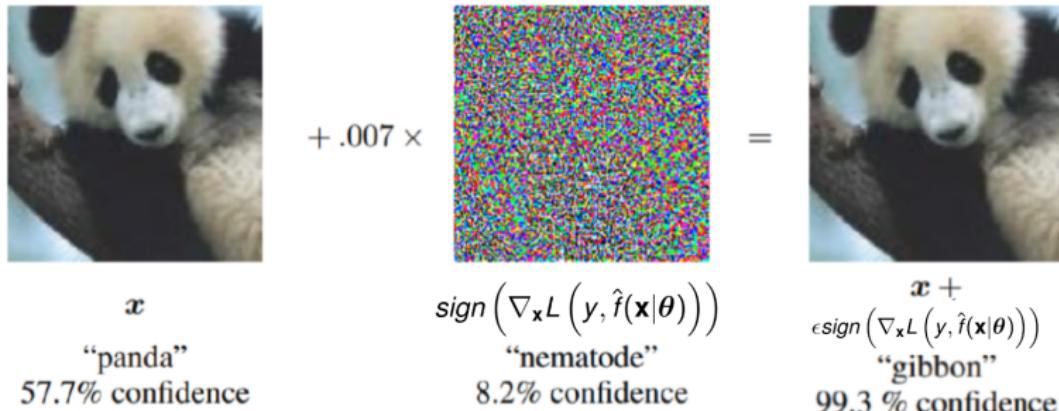


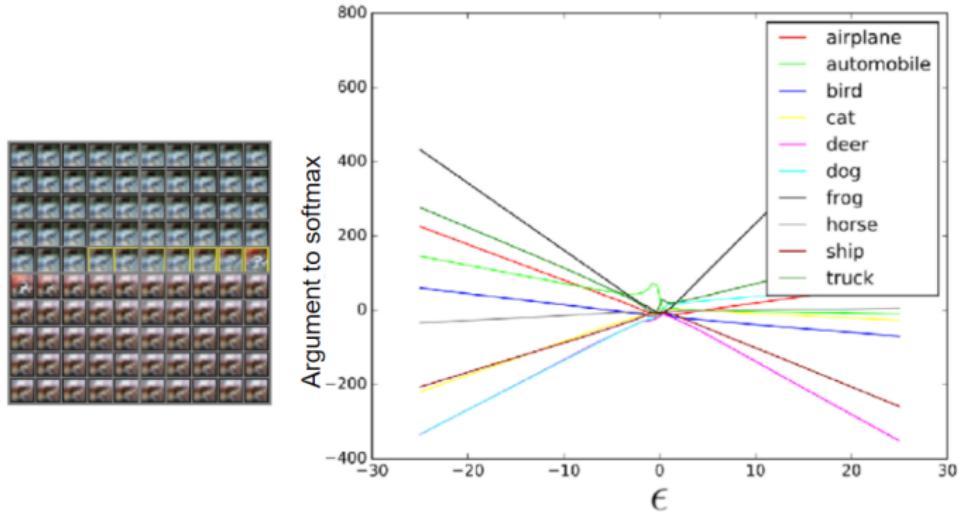
Figure: By adding an imperceptibly small vector whose elements are equal to the sign of the elements of the gradient of the cost function with respect to the input, GoogLeNet's classification of the image was changed from 'panda' to 'gibbon'. In this example, the ϵ is 0.007.

Adversarial Examples : A closer look (Optional)

ADVERSARIAL EXAMPLES AND LINEARITY

- Goodfellow et al. (2014) showed that linear behaviour of models in high dimensional spaces is the reason for the presence of such adversarial examples.
- Neural networks are built from "mostly" linear building blocks. In the case of ReLU activations, the mapping from the input image to the output logits (the inputs to the softmax) is a piece-wise linear function.
- The value of a linear function changes rapidly if it has many inputs. Specifically, if each input is modified by ϵ , a linear function with weights \mathbf{w} can change by as much as $\epsilon \|\mathbf{w}\|_1$, which can be very large if \mathbf{w} is high-dimensional.
- There is a tradeoff between using models that are easy to train due to their linearity and models that can use non-linear effects to become robust to adversarial perturbations.

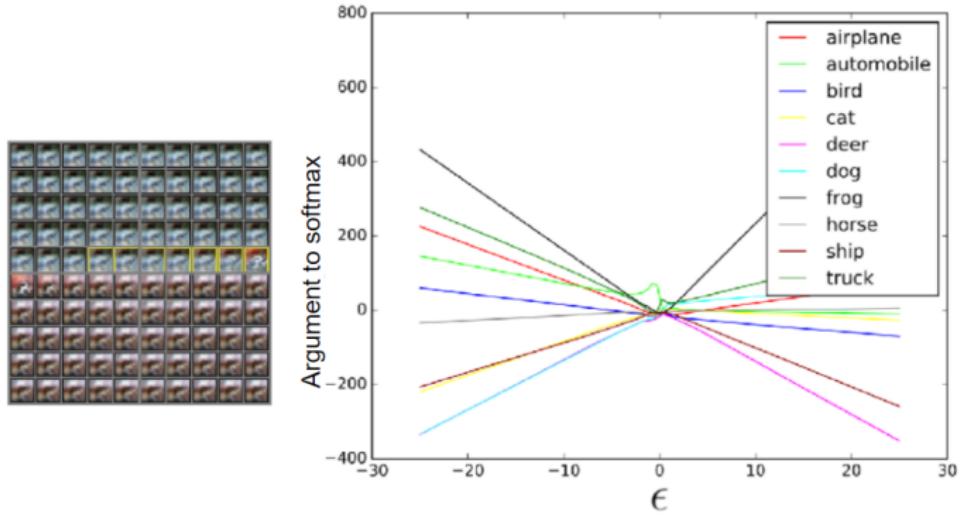
ADVERSARIAL EXAMPLES AND LINEARITY



Credit: Goodfellow

The figure shows the result of moving along a single direction (not necessarily axis-aligned) in the input space of a CNN. We begin with an image of an automobile (somewhere in the end of the fifth row) and move an ϵ amount in this direction (negative ϵ = opposite direction.). The images in the top half are the result of moving in the "negative ϵ " direction and those in the bottom half are the result of moving in the "positive ϵ " direction.

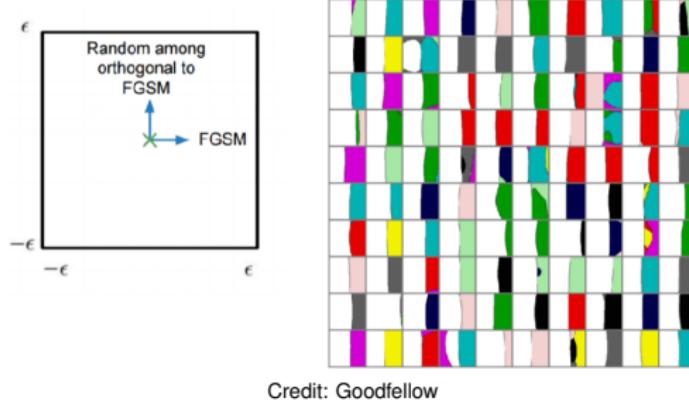
ADVERSARIAL EXAMPLES AND LINEARITY



Credit: Goodfellow

The figure on the right shows the logits (the inputs to the softmax) for each value of ϵ . Each curve is a logit for a specific class. As we move away from the image of the automobile in either direction, the logits for the 'frog' class become extremely high and the images are misclassified by the CNN. The logits for the 'automobile' class are (relatively) high only in the middle of the plot and the CNN correctly classifies these images (yellow boxes on the left).

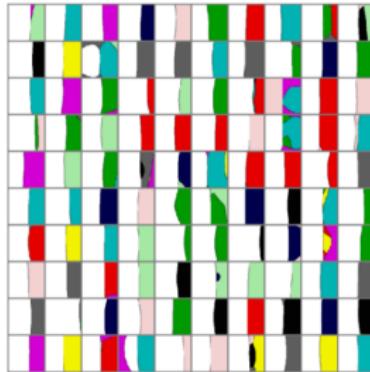
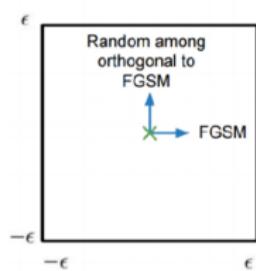
ADVERSARIAL SUBSPACES



Credit: Goodfellow

Each square above represents a 2-dimensional cross section of the input space where the center corresponds to a test example (different squares = different test examples). Moving up or down in a given square indicates moving in a random direction that is orthogonal to the direction of the FGSM. White pixels indicate that the classifier outputs the correct label for the corresponding points and the colored pixels indicate that the classifier misclassifies the corresponding points (different colours = different incorrect classes).

ADVERSARIAL SUBSPACES

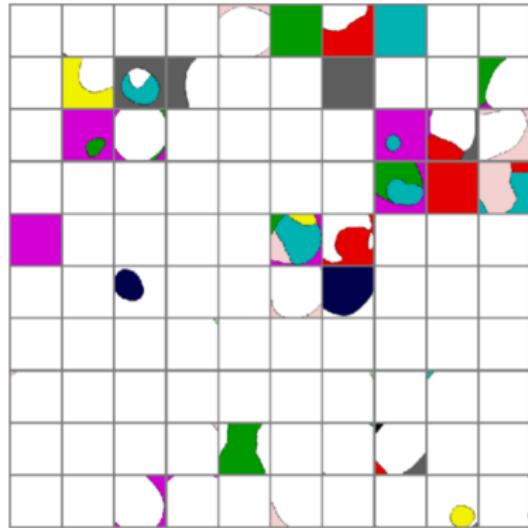
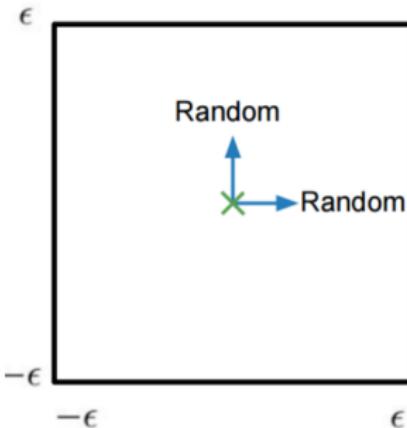


Credit: Goodfellow

The FGSM method identifies a direction such that moving along *any* direction whose unit vector has a large (positive) dot product with the "FGSM vector" results in an adversarial example. Therefore, adversarial examples live in **linear subspaces** of the input space.

ADVERSARIAL SUBSPACES

Adversarial examples
are not noise



Credit:Goodfellow

For a given test example, moving in *completely* random directions (within the ϵ -ball) is unlikely to result in adversarial examples.

PERFORMANCE W.R.T. UNIFORM DISTRIBUTION

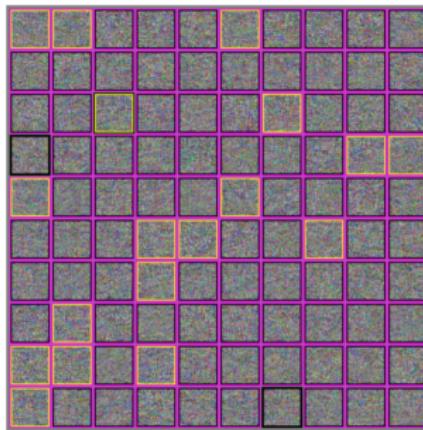


Figure: Each square above corresponds to Gaussian noise that was run through a Cifar-10 classifier.

Instead of measuring the performance of a classifier with respect to the data-generating distribution, if we measure it with respect to a uniform distribution over the whole input space, these classifiers are **wrong almost everywhere**.

PERFORMANCE W.R.T. UNIFORM DISTRIBUTION

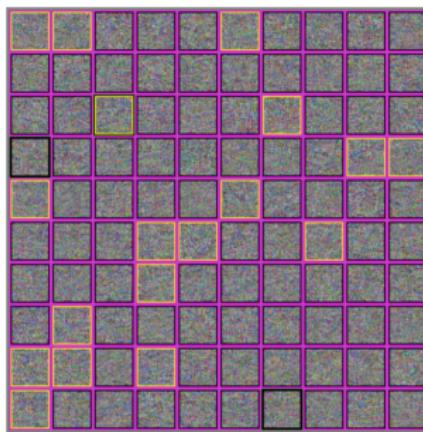


Figure: Each square above corresponds to Gaussian noise that was run through a Cifar-10 classifier.

For all the inputs in the pink boxes, the classifier was reasonably confident (in terms of the softmax values) that the image contained something rather than nothing.

PERFORMANCE W.R.T. UNIFORM DISTRIBUTION

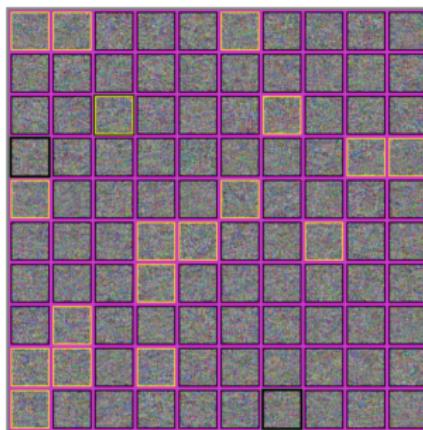


Figure: Each square above corresponds to Gaussian noise that was run through a Cifar-10 classifier.

For the inputs in the yellow boxes, just one step of FGSM was sufficient to convince the model that it was looking at an airplane, specifically.

ADVERSARIAL EXAMPLES: OUTLOOK

- The FGSM method that we have looked at is only one of *many* different algorithms for generating adversarial examples.
- Athalye et al. (2017) 3-D printed a turtle that fooled the network into classifying it as a rifle from most angles.



Credit: Athalye

ADVERSARIAL EXAMPLES: OUTLOOK

- Papernot (2016) discusses ways to fool a classifier even if the model (that is, the network structure and the weights) is unknown. Such methods are called **black-box methods**.
- The library CleverHans can be used to both generate robust adversarial examples and build effective defences against adversarial attacks.

References

REFERENCES

-  Zico Kolter and Aleksander Madry (2019)
Adversarial Robustness - Theory and Practice
<https://adversarial-ml-tutorial.org/>
-  Ian Goodfellow, Yoshua Bengio and Aaron Courville (2016)
Deep Learning
<http://www.deeplearningbook.org/>
-  Ian Goodfellow (2017)
Lecture 16 | Adversarial Examples and Adversarial Training
https://www.youtube.com/watch?v=CIIfsB_EYsVI
-  Ian GoodfellowNicolas PapernotSandy HuangRocky DuanPieter AbbeelJack Clark (2017)
Attacking Machine Learning with Adversarial Examples
<https://openai.com/blog/adversarial-example-research/>

REFERENCES

-  Anh Nguyen, Jason Yosinski and Jeff Clune (2015)
Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images
<https://arxiv.org/abs/1412.1897>
-  Alex Krizhevsky, Ilya Sutskever and Geoffrey E. Hinton (2012)
ImageNet Classification with Deep Convolutional Neural Networks. NIPS.
<https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>
-  Maxence Prevost (2018)
Adversarial ResNet50
<http://arxiv.org/abs/1207.0580>
-  Mahmood Sharif, Sruti Bhagavatula and Lujo Bauer (2016)
Accessorize to a Crime: Real and Stealthy Attacks on State-of-the-Art Face Recognition. Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security.
<https://dl.acm.org/doi/10.1145/2976749.2978392>

REFERENCES

-  Goodfellow, Shlens (2014)
Explaining and Harnessing Adversarial Examples
https://github.com/maxpv/maxpv.github.io/blob/master/notebooks/Adversarial_ResNet50.ipynb
-  Papernot , McDaniel, Goodfellow, Jha, Celik, Swamy (2016)
Practical Black-Box Attacks against Machine Learning
<https://arxiv.org/abs/1602.02697>
-  Athalye , Engstrom, Ilyas, Kwok (2017)
Synthesizing Robust Adversarial Examples
<https://arxiv.org/abs/1707.07397>
-  Tom B. Brown and Catherine Olsson, Research Engineers, Google Brain Team (2018)
Introducing the Unrestricted Adversarial Examples Challenge
<https://ai.googleblog.com/2018/09/introducing-unrestricted-adversarial.html>