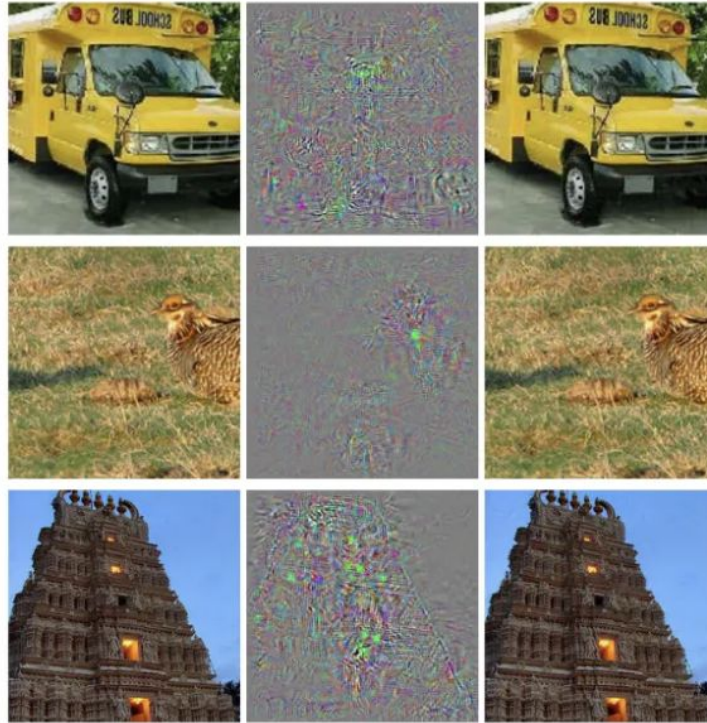


Advanced Machine Learning

Likhith Nayak

Adversarial Examples



Adversarial Examples

Assumption in neural network training:

For a small enough radius r in the vicinity of a given input \mathbf{x} , an input $\mathbf{x} + r$ satisfying $\|r\| < \epsilon$ ($\epsilon > 0$) will predict the correct class

Disproved:

The assumption does not hold true, and it is proved by generating adversarial samples with imperceptibly small perturbations that make the model mispredict classes.

Adversarial Examples

Consider a neural network classifier f that maps an input image $x \in \mathbb{R}^m$ to a set of labels $\{1 \dots k\} — f: \mathbb{R}^m \rightarrow \{1 \dots k\}$

Minimize $\|r\|_2$ subject to:

1. $f(x + r) = l$
2. $x + r \in [0, 1]^m$

Cross-model generalization

	FC10(10^{-4})	FC10(10^{-2})	FC10(1)	FC100-100-10	FC200-200-10	AE400-10	Av. distortion
FC10(10^{-4})	100%	11.7%	22.7%	2%	3.9%	2.7%	0.062
FC10(10^{-2})	87.1%	100%	35.2%	35.9%	27.3%	9.8%	0.1
FC10(1)	71.9%	76.2%	100%	48.1%	47%	34.4%	0.14
FC100-100-10	28.9%	13.7%	21.1%	100%	6.6%	2%	0.058
FC200-200-10	38.2%	14%	23.8%	20.3%	100%	2.7%	0.065
AE400-10	23.4%	16%	24.8%	9.4%	6.6%	100%	0.086
Gaussian noise, stddev=0.1	5.0%	10.1%	18.3%	0%	0%	0.8%	0.1
Gaussian noise, stddev=0.3	15.6%	11.3%	22.7%	5%	4.3%	3.1%	0.3

Cross-training generalization

	FC100-100-10	FC123-456-10	FC100-100-10'
Distorted for FC100-100-10 (av. stddev=0.062)	100%	26.2%	5.9%
Distorted for FC123-456-10 (av. stddev=0.059)	6.25%	100%	5.1%
Distorted for FC100-100-10' (av. stddev=0.058)	8.2%	8.2%	100%
Gaussian noise with stddev=0.06	2.2%	2.6%	2.4%

How unstable are neural networks?

If minor perturbations to the inputs change the model's output to the point that it becomes incorrect, then the model is not stable.

Lipschitz continuity:

Any function where the slope of the line joining any two points on this function is not greater than a real number is called a **Lipschitz continuous function**. This real number is called the **Lipschitz constant**.

$$|f(x_1) - f(x_2)| \leq K|x_1 - x_2|$$

How unstable are neural networks?

The unstability of a neural network can be understood by looking at the Lipschitz constant of each layer $k = 1 \dots K$, as given by:

$$\forall x, r \quad \|\varphi_k(x; W_k) - \varphi_k(x + r; W_k)\| \leq L_k \|r\|$$

where φ_k is the output of the k^{th} **layer**, and L_k is the Lipschitz constant for the k^{th} **layer**. So for the entire network of K layers,

$$\|\varphi(x) - \varphi(x + r)\| \leq L \|r\| \quad \text{where } L = \prod_{k=1}^K L_k$$

How unstable are neural networks?

Layer	Size	Stride	Upper bound
Conv. 1	$3 \times 11 \times 11 \times 96$	4	2.75
Conv. 2	$96 \times 5 \times 5 \times 256$	1	10
Conv. 3	$256 \times 3 \times 3 \times 384$	1	7
Conv. 4	$384 \times 3 \times 3 \times 384$	1	7.5
Conv. 5	$384 \times 3 \times 3 \times 256$	1	11
FC. 1	9216×4096	N/A	3.12
FC. 2	4096×4096	N/A	4
FC. 3	4096×1000	N/A	4

How unstable are neural networks?

Properties of Lipschitz constant for a neural network:

- Small bounds guarantee that no such adversarial examples can appear.
 - This suggests a simple regularization of the parameters, consisting in penalizing each upper Lipschitz bound, which might help improve the generalisation error of the networks
- Large bounds do not automatically translate into existence of adversarial examples
- They don't explain why adversarial examples generalize across different hyperparameters or training sets.

Adversarial Training

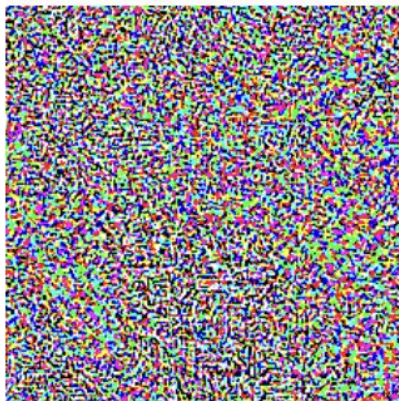


x

“panda”

57.7% confidence

+ .007 ×



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

“nematode”

8.2% confidence

=



$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$

“gibbon”

99.3 % confidence

Adversarial Training - Loss function

- Adversarial examples are a result of models being too linear, rather than too nonlinear.
- Adversarial training can result in regularization; even further regularization than dropout.
- RBF networks are resistant to adversarial examples.

The new loss (or objective) function has a regularizer term denoted by the fast gradient sign method:

$$\tilde{J}(\boldsymbol{\theta}, \mathbf{x}, y) = \alpha J(\boldsymbol{\theta}, \mathbf{x}, y) + (1 - \alpha) J(\boldsymbol{\theta}, \mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} J(\boldsymbol{\theta}, \mathbf{x}, y)))$$

Generative Adversarial Networks (GANs)

