Defense Against Adversarial Attacks

Recall the classical detection methods

- Pre-processing the image with different transformation methods
- Train a network to tell the adversarial instances apart
- Leverage spatial/temporal properties to check the consistency indicating adversarial behaviors
- Map the data to other data manifold by computing meaningful metric to measure differences

Beyond the Min-max Game

- What if we have more knowledge about our learning tasks?
 - Properties of learning tasks and data
 - General understanding about ML models

Characterize Adversarial Examples Based on Spatial Consistency Information for Semantic Segmentation

- Attacks against semantic segmentation
 - State-of-the-art attacks against segmentation: Houdini [NIPS2017], DAG [ICCV 2017]
 - We design diverse adversarial targets: hello kitty, pure color, a real scene, ECCV, color shift, strips of even color of classes
 - Cityscapes and BDD datasets



Benign



Spatial Context Information

- Spatial consistency is a distinct property of image segmentation
- Perturbation at one pixel will potentially affect the prediction of surrounding pixels $\mathcal{H}(m) = -\sum_{j} \mathcal{V}_m[j] \log \mathcal{V}_m[j]$



(a) Benign example

(b) Heatmap of benign image



(c) DAG | Kitty (d) DAG | Pure

(e) Houdini | Kitty (f) Houdini | Pure

Perturbation on single patch may loss its adversarial effect

- Spatial consistency: the consistency of segmentation results for randomly selected patches from an image
- Such spatial consistency information from benign and adversarial instances are distinguishable
- We apply mIOU to compare the segmentation results between patches
 - For each class, Intersection over Union (IOU) is calculated as TP/(TP+FP+FN). Here we calculate the relative mIOU for each pair of patches



Pipeline of spatial consistency based detection for adversarial examples on semantic segmentation

We apply mIOU to evaluate the consistency information for patches from benign and adversarial instances quantitatively

Detection



Spatial Consistency

Adaptive Attack Against Spatial Consistency Based Detection

- Adaptive attack:
 - Assume the attacker has perfect knowledge of #selected patches: K
 - We generate perturbation that the selected k patches can all be mis-segmented to the corresponding regions within adversarial target



Detecting adversarial instances based on spatial consistency information

- Both the spatial consistency based detection and the scaling based baseline achieve promising detection rate on different attacks
- The scaling based baseline fails to detect strong adaptive attacks while the spatial based method can

	Model	mIOU	Detection				Detection Adap			
Method			DAG		Houdini		DAG		Houdini	
			Pure	Kitty	Pure	Kitty	Pure	Kitty	Pure	Kitty
0.5			100%	95%	100%	99%	100%	67%	100%	78%
Scale 3.0	(16.4M)	66.7	100%	100%	100%	100%	100%	0%	97%	0%
(std) 5.0			100%	100%	100%	100%	100%	0%	71%	0%
1			91%	91%	94%	92%	98%	94%	92%	94%
Spatial 5 DRN	66 7	100%	100%	100%	100%	100%	100%	100%	100%	
(K) 10	(16.4M)	00.7	100%	100%	100%	100%	100%	100%	100%	100%
$\left 50 \right $		100%	100%	100%	100%	100%	100%	100%	100%	

Takeaways

- Spatial consistency information can be potentially applied to help distinguish benign and adversarial instances against segmentation models.
- Strong adaptive attacker can hardly succeed when large randomness is incorporated into the model

Adversarial Frames In Videos

Attacks on segmentation



Attacks on pose estimation



Attacks on object detection



Defensing Adversarial behaviors in Videos – Temporal Dependency



Teelr	Attack	Target	Previous	Previous Detection			Detection Adap		
Task	Method	Target	Frames	1	3	5	1	3	5
Semantic Segmentation	Houdini	CVDD	Benign	100%	100%	100%	100%	100%	100%
		CVFK	Adversarial	100%	100%	100%	100%	100%	100%
		Remapping	Benign	100%	100%	100%	100%	100%	100%
			Adversarial	100%	100%	100%	100%	100%	100%
		Stripe	Benign	100%	100%	100%	100%	100%	100%
			Adversarial	100%	100%	100%	99%	100%	100%
	DAG	CVPR	Benign	100%	100%	100%	100%	100%	100%
			Adversarial	100%	100%	100%	100%	100%	100%
		Remapping	Benign	100%	100%	100%	100%	100%	100%
			Adversarial	100%	100%	100%	100%	100%	100%
		Stripe	Benign	100%	100%	100%	100%	100%	100%
			Adversarial	100%	100%	100%	100%	100%	100%
Human Pose Estimation	Houdini	shuffle	Benign	100%	100%	100%	100%	100%	100%
			Adversarial	100%	100%	100%	99%	100%	100%
		Transpose	Benign	100%	100%	100%	98%	100%	100%
			Adversarial	98%	99%	100%	98 %	99%	100%
Object Detection	DAG	all	Benign	100%	100%	100%	100%	100%	100%
			Adversarial	100%	100%	100%	98%	100%	100%
		person	Benign	99%	100%	100 %	100%	100%	100%
			Adversarial	97%	98%	100%	96 %	97%	100%

- The results show that choosing more random patches can improve detection rate while k=5 is enough to achieve AUC 100%
- The spatial consistency based detection is robust against strong adaptive attackers due to the randomness in patch selection

Object Detection



Human pose





Temporal Consistency Based Analysis

• "Yanny" or "Laurel"? – adversarial audio



- Design robust neural networks that are robust to adversarial attacks
- Defense: recover the ground truth instead of just tell adversarial instance apart
- Necessary step: design novel and advanced architectures built on new computational paradigms
- PeerNets:
 - Euclidean convolutions -> graph convolutions
 - Non-local forward propagation: Capture global structure induced by the data graph
 - Design a peer regularization layer

- Peer Regularization layer
- For N images, each image will look for its K nearest neighbors based on cosine similarity
 - For each image, there is a $n \times d$ feature map

$$\tilde{\mathbf{x}}_{p}^{i} = \sum_{k=1}^{K} \alpha_{ij_{k}pq_{k}} \mathbf{x}_{q_{k}}^{j_{k}}, \qquad \alpha_{ij_{k}pq_{k}} = \frac{\text{LeakyReLU}(\exp(a(\mathbf{x}_{p}^{i}, \mathbf{x}_{p_{k}}^{j_{k}})))}{\sum_{k'=1}^{K} \text{LeakyReLU}(\exp(a(\mathbf{x}_{p}^{i}, \mathbf{x}_{p_{k'}}^{j_{k'}})))}$$

- Randomized approximation
- Monte Carlo approximation
 - Select smaller batch and sample the nearest neighbor from each batch $\{l_{m1},\ldots,l_{mN}\}\subset\{1,\ldots,N'\}$

$$\tilde{\mathbf{x}}_{p}^{i} = \frac{1}{M} \sum_{m=1}^{M} \sum_{k=1}^{K} \alpha_{ij_{mk}pq_{mk}} \mathbf{x}_{q_{mk}}^{j_{mk}}$$

• Other optimization method?

- Select M = 1 during training and large M during inference
 - Limitations?



Results for PeerNets



Visualization of perturbation



Takeaways

- Alternate Euclidean Graph convolution to harness information from peers can provide global information
- Can be added to any models as regularized layer -> good principle
- Not affect the benign accuracy -> important
- How to scale up?
- How to consider more peer images instead of pixels?
- Temporal information?

Similar reading

- Countering adversarial images using input transformations
 - Image quilting nearest patches
 - Computationally expensive

Interesting reading

• A simple neural network module for relational reasoning



Interesting reading

- Deformable Convolutional Networks
 - Deformable convolution and deformable RoI pooling
 - Augment the spatial sapling locations with additional offset which can be learned



Towards Deep Learning Models Resistant to Adversarial Attacks

$$\min_{\theta} \rho(\theta), \quad \text{where} \quad \rho(\theta) = \mathbb{E}_{(x,y)\sim\mathcal{D}} \left[\max_{\delta \in \mathcal{S}} L(\theta, x + \delta, y) \right]$$

- Use a natural saddle point (min-max) formulation to capture the notion of security against adversarial attacks in a principled manner.
- The formulation casts both attacks and defenses into a common theoretical framework.
- Motivate projected gradient descent (PGD) as a universal "firstorder adversary".

Madry et al. Towards Deep Learning Models Resistant to Adversarial Attacks, ICLR 2017.

Model Capacity







Towards Deep Learning Models Resistant to Adversarial Attacks



Decision Boundary Based Detection



Decision Boundary Analysis of Adversarial Examples



	False pos.	False	neg.	Accuracy		
Training attack	Benign	OptBrittle	O PT M ARGIN	Our approach	Cao & Gong	
	N					
OptBrittle	1.0%	1.0%	74.1%			
OptMargin	9.6%	0.6%	7.2%	90.4%	10%	
	MNIS	al training	2011/0	2070		
OptBrittle	2.6%	2.0%	39.8%			
OptMargin	10.3%	0.4%	14.5%			
	CII					
OptBrittle	5.3%	3.2%	56.8%			
OptMargin	8.4%	7.4%	5.3%	96.4%	5%	
	CIFAR-		2.12			
OptBrittle	0.0%	2.4%	51.8%			
OptMargin	3.6%	0.0%	1.2%			

Takeaways

- Decision boundaries of DNNs are important towards improving learning robustness
- Isolated islands in the data manifold would lead to harder detected/defensed adversarial behaviors