

On Robustness and Efficiency of Machine Learning Systems

Yao Li

Department of Statistics and Operations Research
UNC Chapel Hill

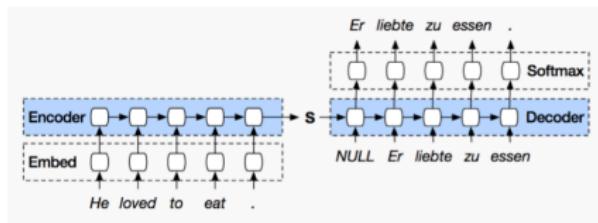
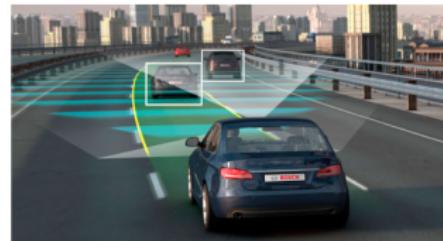
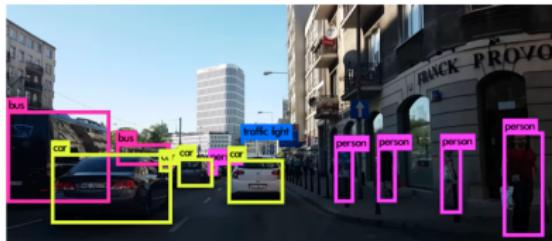


Thomas

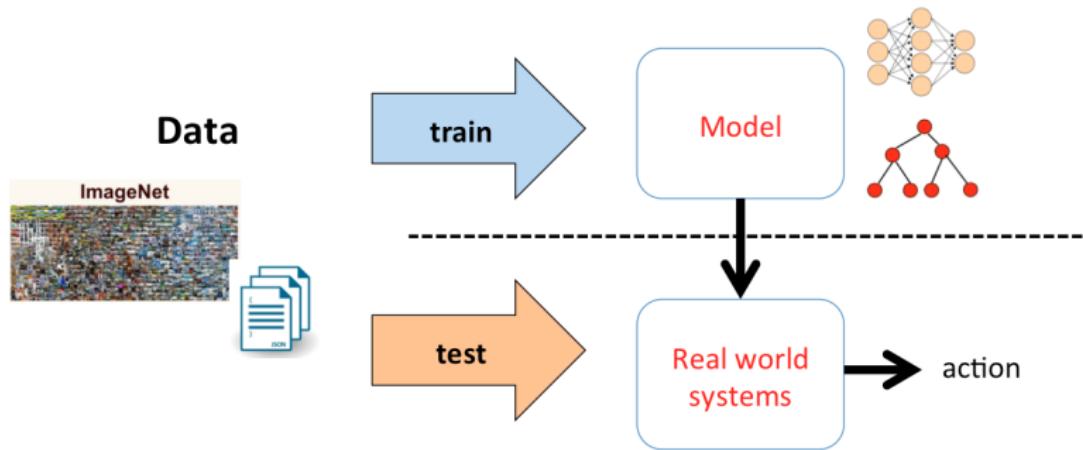


Cho

Machine Learning

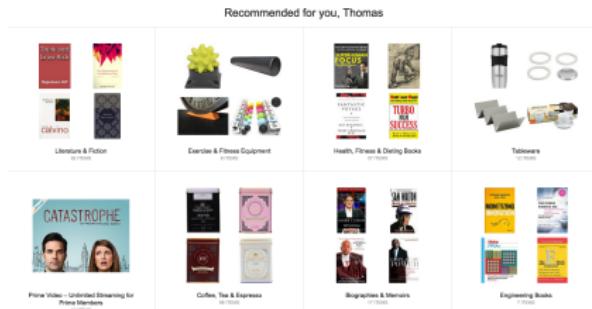


Machine Learning Systems



Training: Scalable and efficient algorithms

- Utilize more information
- Need for Scalabe algorithms



(Recommendation system for e-commerce)



(Game result prediction)

Scalable Demand-Aware Recommendation

- Durable vs. Non-durable goods

YOU MAY ALSO LIKE

New Apple MacBook Pro
(16-inch, 16GB RAM,
512GB Storage) - Space
Gray
★★★★★ 86
\$2,199.00 ✓prime

The Hunger Games:
Catching Fire
★★★★★ 14,708
\$5.99

John Wick Chapter 2
★★★★★ 7,361
\$5.99

Bounty Quick-Size Paper
Towels, White, 8 Family
Rolls = 20 Regular Rolls
★★★★★ 601
\$21.99 ✓prime

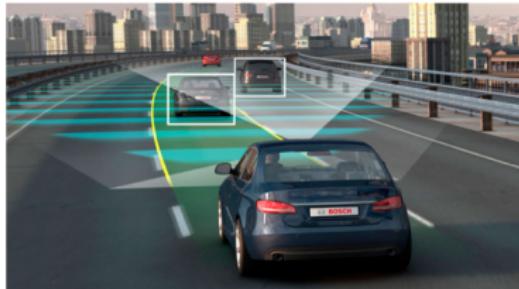
Daewoo FRS-Y22D2T RFS-Y22D2T 20 Cu. Ft. Side
Mounted Silver
Refrigerator, includes ...
★★★★★ 3
\$1,104.71 ✓prime

[ADD TO BAG >](#) [ADD TO BAG >](#) [ADD TO BAG >](#) [ADD TO BAG >](#) [ADD TO BAG >](#)

- Form utility: the item is desired as it is manifested
- Time utility: the item is desired at the given point in time

Prediction: Robustness and Safety

ML systems need to interact with real world



- Robustness and Safety

Today's Outline

1 Robustness and Safety of ML Systems

- Background
- Related Work
- Proposed Defense Methods

2 Efficient Machine Learning Algorithms

- Background
- Previous Work
- Proposed Method

Outline

1 Robustness and Safety of ML Systems

- Background
- Related Work
- Proposed Defense Methods

2 Efficient Machine Learning Algorithms

- Background
- Previous Work
- Proposed Method

Deep Neural Network

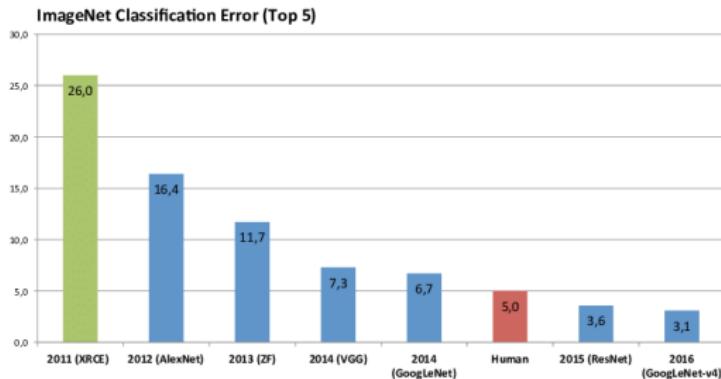
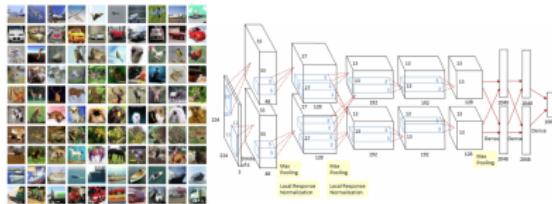
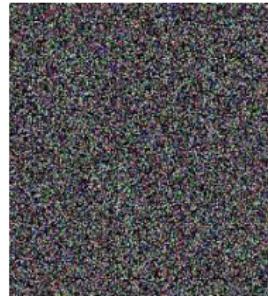


Figure: Winner results of the ImageNet large scale visual recognition challenge (LSVRC) of the past years on the top-5 classification task

Adversarial Examples



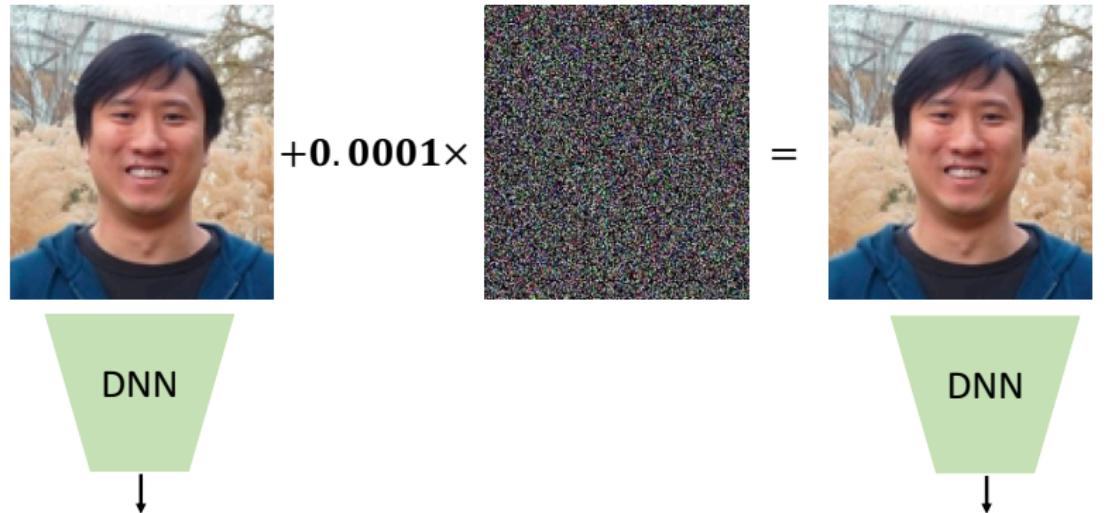
+0.0001×



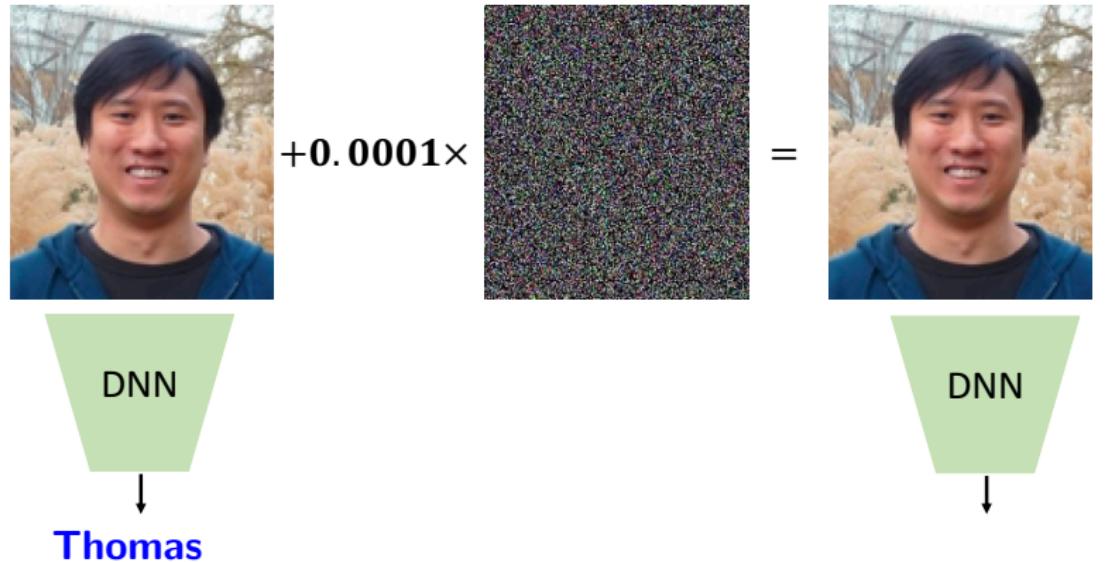
=



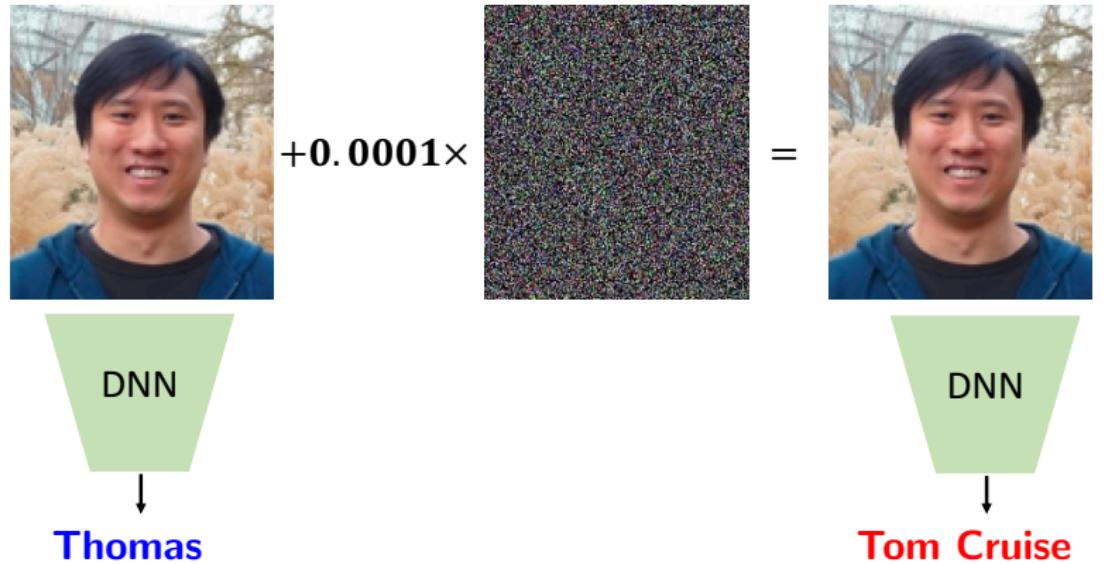
Adversarial Examples



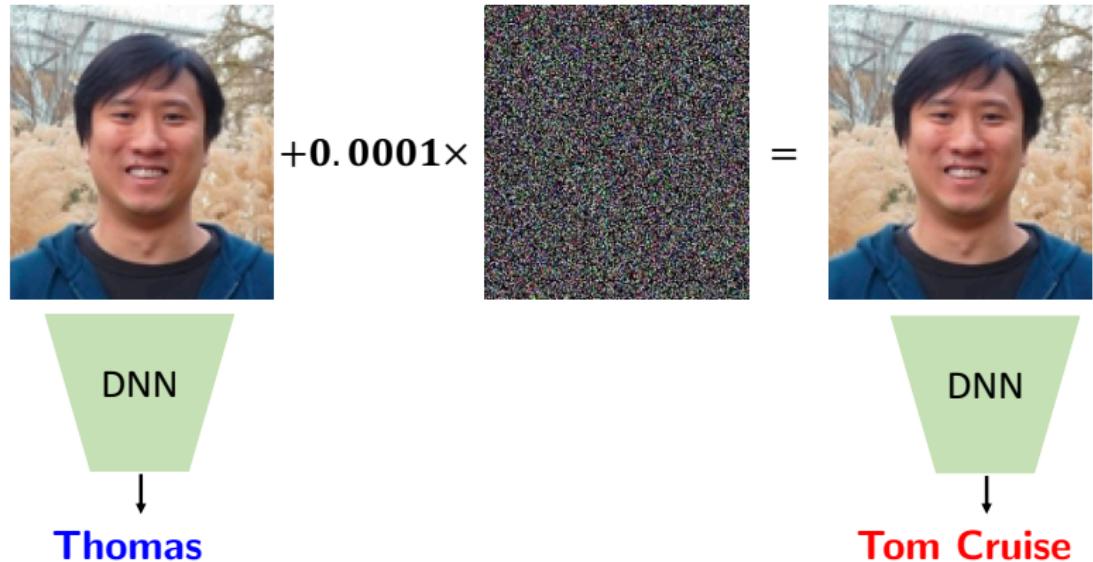
Adversarial Examples



Adversarial Examples



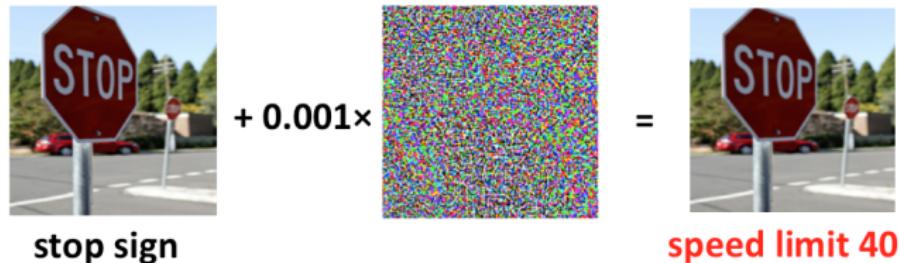
Adversarial Examples



A carefully crafted **adversarial** example can easily fool a deep network

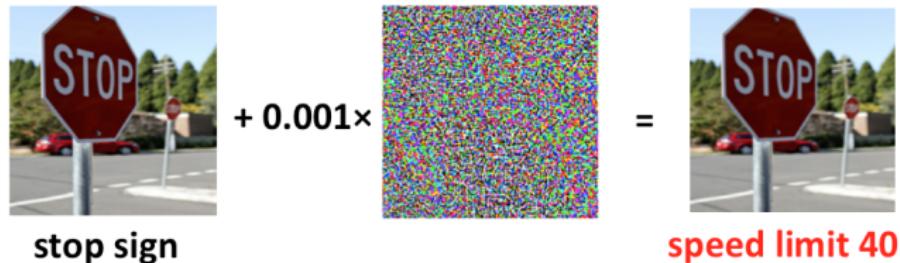
Adversarial Examples

- Robustness is critical in real systems



Adversarial Examples

- Robustness is critical in real systems



Not safe!

Research Questions

Attack: How to craft adversarial example?

Defense: How to **improve** robustness?

Outline

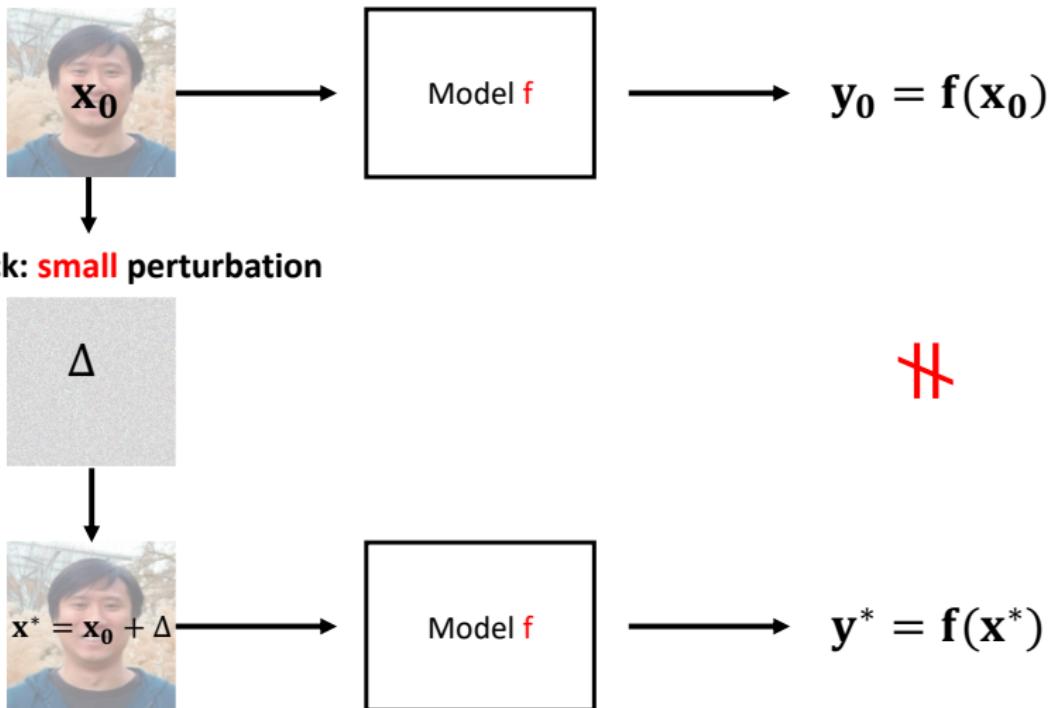
1 Robustness and Safety of ML Systems

- Background
- Related Work
- Proposed Defense Methods

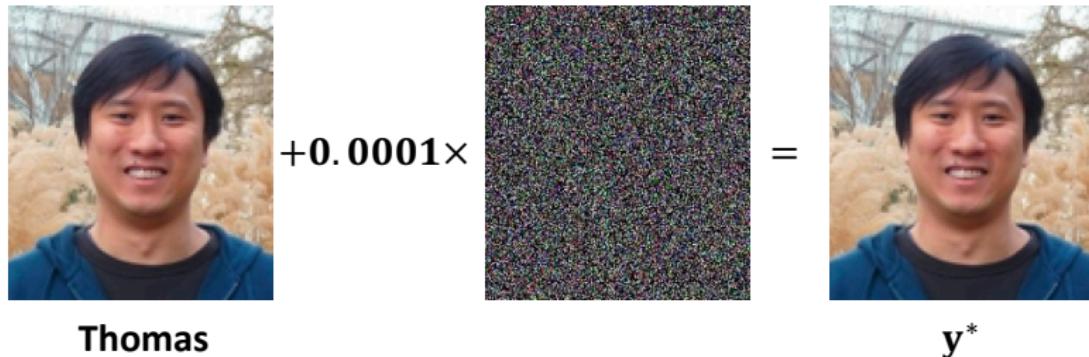
2 Efficient Machine Learning Algorithms

- Background
- Previous Work
- Proposed Method

Notations and Attack Procedure



Type of Attacks



- ① Untargeted attack: $y^* \neq \text{Thomas}$
- ② Targeted attack: For target class $t = \text{Tom Cruise}$, the attacker wants $y^* = \text{Tom Cruise}$

Projected Gradient Descent Attack (PGD)

Attack as an **optimization** problem:

- Given prediction model with fixed parameter θ
- (x_0, y_0) : input image and label

$$\delta = \arg \max_{\delta \in \mathcal{S}} L(\theta, x_0 + \delta, y_0)$$

$$x^* = x_0 + \delta$$

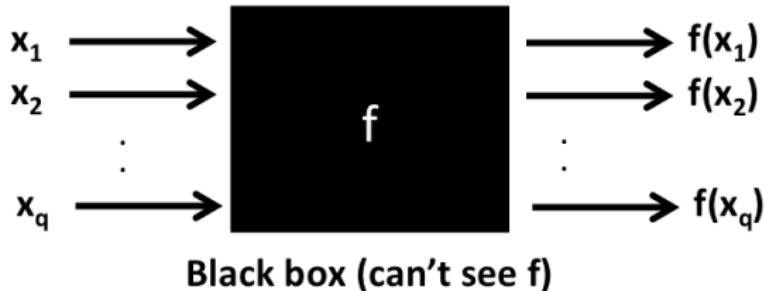
L : loss function training the classifier

δ : adversarial perturbation

$\mathcal{S} \in \mathbb{R}^d$: allowed perturbations, usually chosen to be $\{\delta | \|\delta\|_\infty \leq \epsilon\}$
(Madry et al., 2018)

Adversarial Attack: Other setting

- Black-box setting: Only part of the information is available



- Distortion can be measured by other norms:
(e.g., ℓ_2 , Elastic net, ...)
(Carlini et al., 2017; Chen et al., 2018; Brendel et al. 2017)

Attack vs. Defense

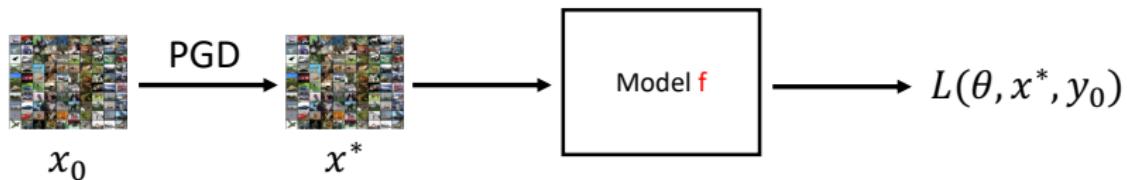
Adversarial Examples Are Not Easily Detected: Bypassing Ten Detection Methods

Nicholas Carlini David Wagner
University of California, Berkeley

Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples

Anish Athalye^{*1} Nicholas Carlini^{*2} David Wagner²

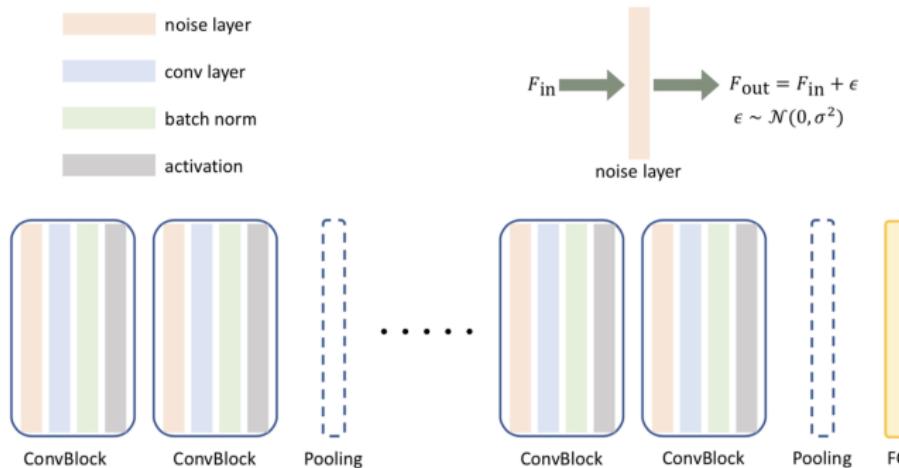
Previous Method 1: Madry's adversarial training



- **Madry's adversarial training:** Madry et al. (2018) proposed to incorporate the adversarial search inside the training process, by solving the following robust optimization problem:

$$\arg \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left\{ \max_{\|\delta\|_\infty \leq \epsilon} L(\theta, x + \delta, y) \right\}$$

Previous Method 2: Random Self-Ensemble



- **Random Self-Ensemble:** Liu et al. proposed a “noise layer”, which fuses output of each layer with Gaussian noise.

Outline

1 Robustness and Safety of ML Systems

- Background
- Related Work
- Proposed Defense Methods

2 Efficient Machine Learning Algorithms

- Background
- Previous Work
- Proposed Method

Proposed Defense Methods

- ① Adversarial Bayesian Neural Network (Adv-BNN)
- ② Embedding Regularized Classifier (ER-Classifier)

Adversarial Bayesian Neural Network (Adv-BNN)

Observation 1

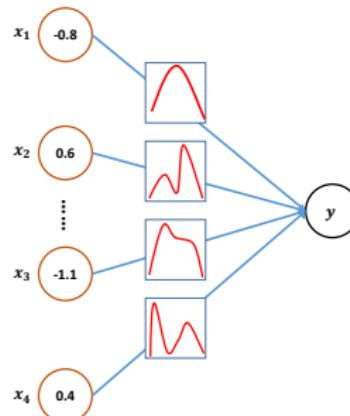
Liu et al., 2018: Adding stochastic components in the neural network to hide gradient information from attacker improves robustness.

Observation 2

Madry et al., 2018: Min-max formulation of training DNN classifier improves the robustness against adversarial examples.

- Our approach: combining adversarial training with Bayesian neural network
 - Bayesian neural network
 - Min-max optimization

Bayesian Neural Network (BNN)



- All weights are represented by probability distributions over possible values
- BNN: a probabilistic model $p(y|\mathbf{x}, \mathbf{w})$

(Blundell et al., 2015)

Adv-BNN: Objective Function

Recall BNN objective function:

$$\arg \max_{\mu, s} \left\{ \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}} \mathbb{E}_{\mathbf{w} \sim q_{\mu, s}} \log p(y_i | \mathbf{x}_i, \mathbf{w}) - \text{KL}(q_{\mu, s}(\mathbf{w}) \| p(\mathbf{w})) \right\}$$

Adv-BNN objective function:

$$\arg \max_{\mu, s} \left\{ \left[\sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}} \min_{\|\delta\|_\infty \leq \epsilon} \mathbb{E}_{\mathbf{w} \sim q_{\mu, s}} \log p(y_i | \mathbf{x}_i + \delta, \mathbf{w}) \right] - \text{KL}(q_{\mu, s}(\mathbf{w}) \| p(\mathbf{w})) \right\} \quad (1)$$

- \mathcal{D} : data distribution
- δ : adversarial distortion
- ϵ : maximum ℓ_∞ distortion
- $p(\mathbf{w})$: prior distribution, $\mathcal{N}(0_d, s_0^2 \mathbf{I}_{d \times d})$.

Embedding Regularized Classifier (ER-Classifier)

Observation 1

Levina et al., 2005: The only reason any methods work in very high dimensions is that, in fact, the data are not truly high-dimensional.

Observation 2

Tanay et al., 2016: Adversarial samples do not lie on the data manifold.

- How to use the two findings to improve robustness of Deep Neural Network?

Deep Classifier: Framework

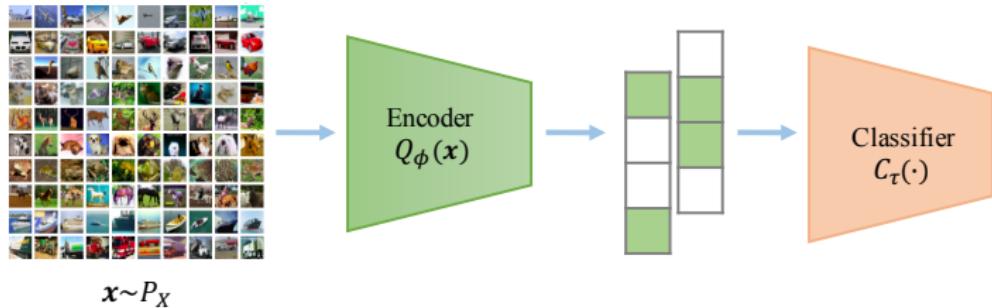


Figure: Deep Classifier

- P_X : Data distribution
- Q_φ : Encoder (neural network)
- C_τ : Classifier (neural network)

ER-Classifier: Framework

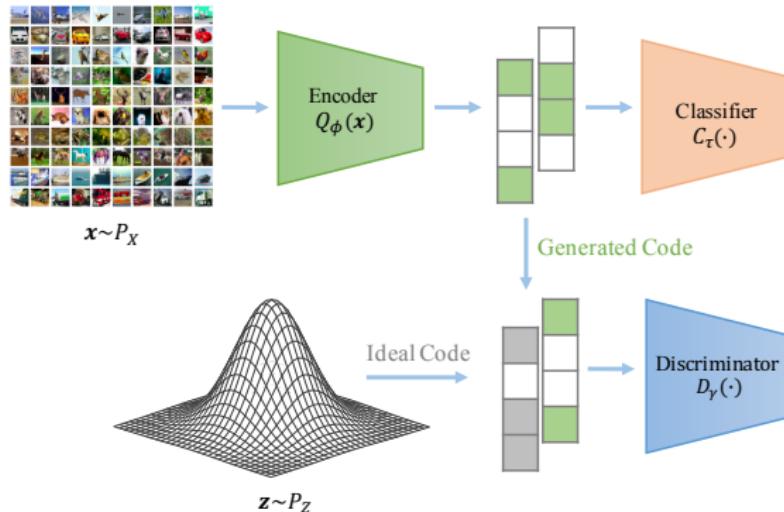


Figure: Embedding Regularized Classifier

- P_Z : Prior distribution
- D_γ : Discriminator (neural network)
(discriminator D used to separate “true” $z \sim P_Z$ and “fake” $\tilde{z} \sim Q_Z$)

Embedding Visualization

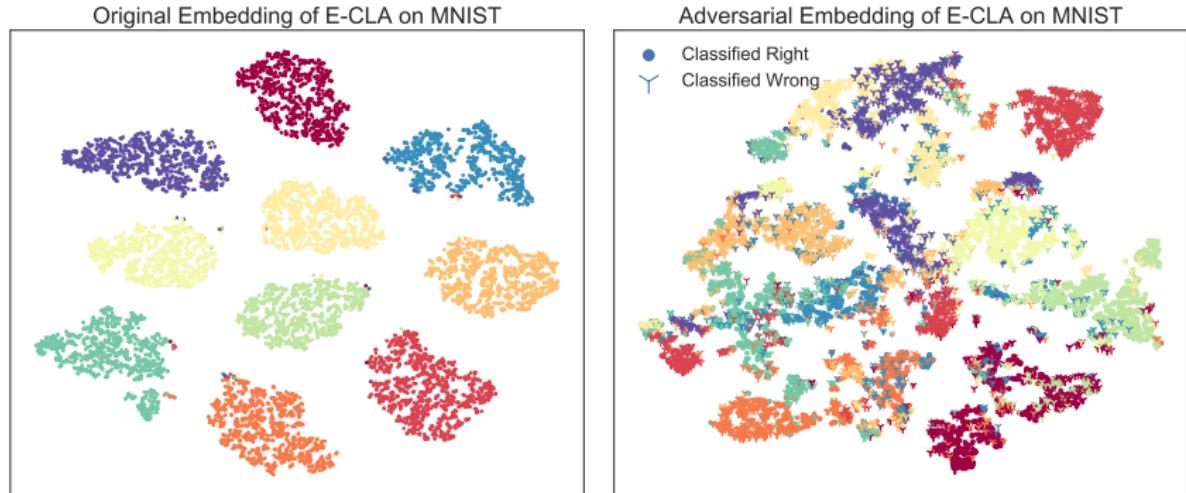


Figure: 2D embeddings for E-CLA on MNIST.

Embedding Visualization

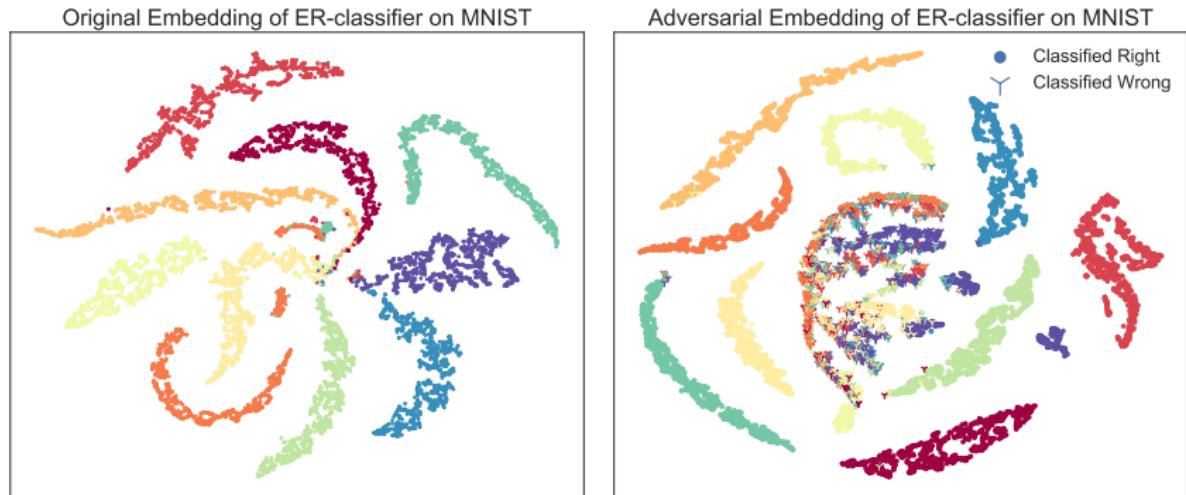


Figure: 2D embeddings for ER-Classifier on MNIST.

Future Work

- Provide uncertainty for label prediction
- Detecting adversarial examples
- Better evaluation methods
- Robustness for statistical model

Outline

1 Robustness and Safety of ML Systems

- Background
- Related Work
- Proposed Defense Methods

2 Efficient Machine Learning Algorithms

- Background
- Previous Work
- Proposed Method

Introduction: Group Comparison



Figure: NBA and Dota 2

- Given the results of **previous** games, how to **predict** the outcome of an unseen game?

Problem Setting

- Number of individuals: n
- Number of observed comparisons: T
- Each game involves two disjoint teams: I_t^+ and I_t^- , each of them is a subset of $\{1, \dots, n\}$
- Possible number of teams: N , where $N = \binom{n}{L}$ and $L = |I_t^+| = |I_t^-|$.
- The outcome of game t : o_t .
 - Measured outcomes (scores): $o_t \in \mathbb{R}$
 - Binary indicator outcomes (wins/losses): $o_t \in \{+1, -1\}$.

Example: only players 2,3,4 are on the winning team, $I_t^+ = \{2, 3, 4\}$

Outline

1 Robustness and Safety of ML Systems

- Background
- Related Work
- Proposed Defense Methods

2 Efficient Machine Learning Algorithms

- Background
- **Previous Work**
- Proposed Method

Assumptions of Previous Work

Assumption 1

The team's score is the sum over team members' scores: $s_t^+ = \sum_{i \in I_t^+} w_i$, where w_i is the ability of player i . The observed outcome is determined by $s_t^+ - s_t^-$.

Logistic Regression Model

$$\min_{\mathbf{w} \in \mathbb{R}^n} \sum_{t=1}^T \ell(\mathbf{w}^T \mathbf{x}_t, o_t) + R(\mathbf{w}), \quad (2)$$

- ℓ is logistic loss.
- \mathbf{w} are the individual scores we want to learn.
- $\mathbf{x}_t \in \mathbb{R}^n$ is the indicator vector, where $(\mathbf{x}_t)_j = 1$ if $j \in I_t^+$, $(\mathbf{x}_t)_j = -1$ if $j \in I_t^-$, and $(\mathbf{x}_t)_j = 0$ for all other elements.

Limits of these Models



- Common weakness: they cannot evaluate the player-interaction effects.
- Solution: include player-interaction effects in the model
- Two questions regarding including player-interaction effects:
 - Can incorporating player-interaction effects improve the prediction accuracy of the model?
 - When n is large, how can our algorithm scale up?

Outline

1 Robustness and Safety of ML Systems

- Background
- Related Work
- Proposed Defense Methods

2 Efficient Machine Learning Algorithms

- Background
- Previous Work
- **Proposed Method**

Assumptions of New Model

Assumption 2

Assume each user has their individual score w_{jj} , and for each pair of users there is a pair-wise score w_{ij} . A team's ability is modeled by

$$s_t^+ = \sum_{i,j \in I_t^+} w_{ij}.$$

Goal: learn the model so that the score s_t^+ is larger than s_t^- for each game.

Basic Model for Higher Order Interactions

- **Basic HOI:**

$$\min_{W \in \mathbb{R}^{n \times n}} \sum_{t=1}^T \ell((\mathbf{e}_t^+)^T W(\mathbf{e}_t^+) - (\mathbf{e}_t^-)^T W(\mathbf{e}_t^-), o_t) + \lambda \|W\|_F^2. \quad (3)$$

where \mathbf{e}_t^+ is the indicator vector for I_t^+

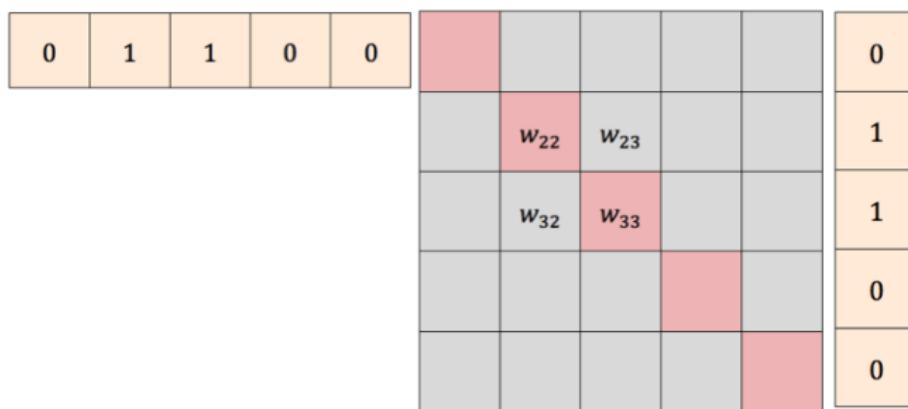


Figure: Example of $(\mathbf{e}_t^+)^T W(\mathbf{e}_t^+)$, and W is the matrix we want to learn

Limits of Basic HOI

- Basic HOI can be reformulated as logistic regression and solved by standard methods.
- Sample complexity: Problem (3) has $\frac{n(n+1)}{2}$ parameters, so the degree of freedom grows with $O(n^2)$.
- Computing: Problem (3) requires $O(n^2)$ memory to store the W matrix. Therefore, a standard solver will be hard to scale beyond 30,000 players.
- Conclusion: When n is very large, Basic HOI doesn't work.

How to handle pair-wise effects when n is large?

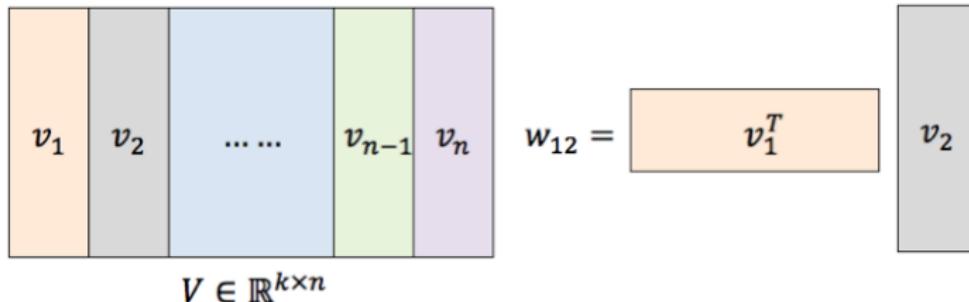
Assumptions of Factorization Model

Assumption 3

We assume a team's score can be written as

$$s_t^+ = \sum_{i \in I_t^+} w_i + \sum_{i \in I_t^+} \sum_{j \in I_t^+} \mathbf{v}_i^T \mathbf{v}_j$$

- Model parameters that have to be estimated are $\mathbf{w} \in \mathbb{R}^n$ and $V \in \mathbb{R}^{k \times n}$, each \mathbf{v}_j is the j -th column of V .
- The hyper-parameter k defines the dimensionality of the factorization.



Factorization Model for Higher Order Interactions

- **Factorization HOI:**

$$\begin{aligned} & \underset{\mathbf{w} \in \mathbb{R}^n, V \in \mathbb{R}^{k \times n}}{\operatorname{argmin}} \sum_{t=1}^T \ell(s_t^+ - s_t^-) + \frac{\lambda_w}{2} \|\mathbf{w}\|_2^2 + \lambda_V \|V\|_F^2 \\ &= \underset{\mathbf{w} \in \mathbb{R}^n, V \in \mathbb{R}^{k \times n}}{\operatorname{argmin}} \sum_{t=1}^T \ell(\mathbf{w}^T (\mathbf{e}_t^+ - \mathbf{e}_t^-) + \sum_{i,j \in I_t^+} \mathbf{v}_i^T \mathbf{v}_j - \sum_{i,j \in I_t^-} \mathbf{v}_i^T \mathbf{v}_j) \quad (4) \\ & \quad + \frac{\lambda_w}{2} \|\mathbf{w}\|_2^2 + \lambda_V \|V\|_F^2, \end{aligned}$$

where λ_w and λ_V are the regularization parameters.

Data Statistics



Figure: Character vs. Character, Player vs. Player

Datasets	HotS Tournament (Character)	HotS Tournament (Player)	HotS Public (Character)	HotS Public (Player)	Dota 2 (Character)	Dota 2 (Player)
Number of Games (T)	9,610	9,610	139,462	139,462	46,459	46,459
Number of Individuals (n)	54	3,470	62	7,251	113	30,452

Table: Dataset Statistics

Methods

- Basic HOI: the proposed basic model using pairwise information.
- Factorization HOI: the proposed latent factor model, which approximates the pairwise interaction by a factor form.
- Trueskill: the state-of-the-art algorithm used in all the online game matching engines.
- Bradley-Terry model: the generalized Bradley-Terry model for group comparison data.
- Logistic Regression (LR): another baseline for individual score model using logistic loss.

Experiment Results

Table: Performance of the proposed algorithm and other algorithms. The numbers are prediction accuracy (%), and “oom” indicates out of memory here.⁴

Datasets	LR	Trueskill (1)	Trueskill (10)	Bradley-Terry	Basic HOI	Factorization HOI
HotS Tournament (C)	59.73	62.90	58.48	59.52	80.59	77.84
HotS Tournament (P)	83.45	80.02	84.50	84.18	83.89	85.17
HotS Public (C)	54.34	53.36	53.06	53.50	54.45	54.75
HotS Public (P)	54.01	53.64	53.87	53.92	53.39	55.76
Dota 2 (C)	61.64	52.50	52.61	61.37	65.34	63.72
Dota 2 (P)	65.98	62.16	64.26	62.72	oom	68.25

Patterns Learned by Factorization HOI

Table: Top-5 pairs and bottom-5 character pairs learned by our model on Heroes of the storm tournament data.

Top 5 pairs	Bottom 5 pairs
(Lunara, Leoric)	(Raynor, Zeratul)
(Kerrigan, Sylvanas)	(Illidan, Thrall)
(Regar, Illidan)	(Sonya, Zeratul)
(Chen, Jaina)	(Muradin, Lt. Morales)
(Thrall, Valla)	(Anub'arak, Illidan)

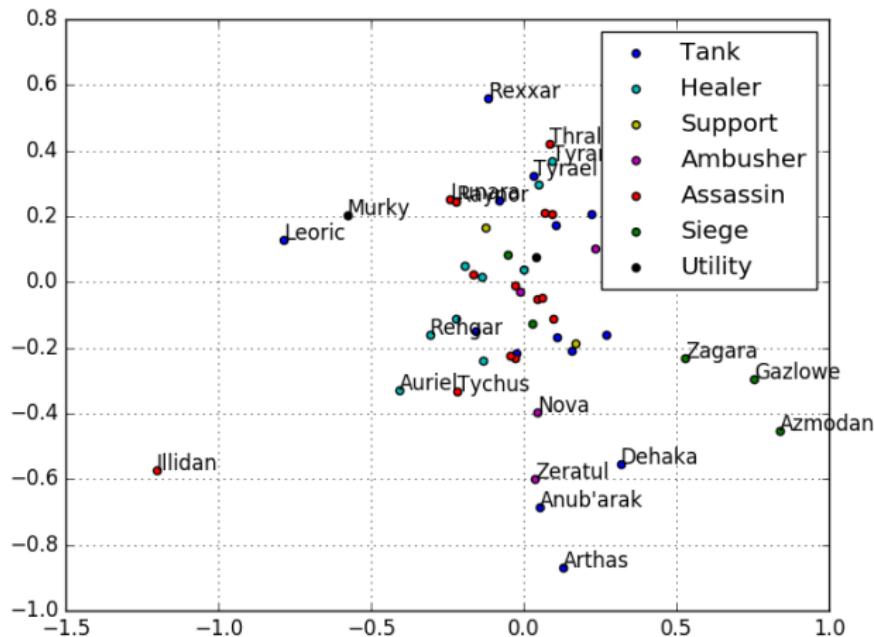


Figure: Projection of interaction features for each character (v_i) to 2-D space. Colors represent the official categorization for these characters.

References

- [1] Blundell, C., Cornebise, J., Kavukcuoglu, K., and Wierstra, D. *Weight uncertainty in neural networks.* ICML, 2015
- [2] Tanay, T., and Griffin, L. *A boundary tilting perspective on the phenomenon of adversarial examples.* ArXiv, 2016
- [3] Athalye, A., Engstrom, L., Ilyas, A., and Kwok, K. *Synthesizing robust adversarial examples.* ICML, 2018.
- [4] Carlini, N., and Wagner, D. *Towards evaluating the robustness of neural networks.* SC, 2017.
- [5] Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. *Towards deep learning models resistant to adversarial attacks.* ICLR, 2018.
- [6] Levina, E., and Bickel, P. J. *Maximum likelihood estimation of intrinsic dimension.* NIPS, 2005.
- [7] Simon-Gabriel, C. J., Ollivier, Y., Bottou, L., Schölkopf, B., and Lopez-Paz, D. *Adversarial vulnerability of neural networks increases with input dimension.* ICML, 2019.

References

- [8] Liu, X., Cheng, M., Zhang, H., and Hsieh, C. J. *Towards robust neural networks via random self-ensemble*. ECCV, 2018.
- [9] Li, Y., Cheng, M., Fujii, K., Hsieh, F., and Hsieh, C. J. *Learning from group comparisons: exploiting higher order interactions*. NIPS, 2018.
- [10] Liu, X., Li, Y., Wu, C., and Hsieh, C. J. *Adv-bnn: Improved adversarial defense through robust bayesian neural network*. ICLR, 2019.
- [11] Yi, J., Hsieh, C. J., Varshney, K. R., Zhang, L., and Li, Y. *Scalable demand-aware recommendation*. NIPS, 2017.
- [12] Li, Y., Min, M. R., Yu, W., Hsieh, C. J., Lee, T., and Kruus, E. *Defending against adversarial examples by regularized deep embedding*. ICML (to be submitted), 2020.
- [13] Li, Y., Yu, W., Min, M. R., Lee, T., Kruus, E., Wang, W., and Hsieh, C. J. *Detecting adversarial examples with regularized deep embedding*. IJCAI (to be submitted), 2020.

Embedding Visualization

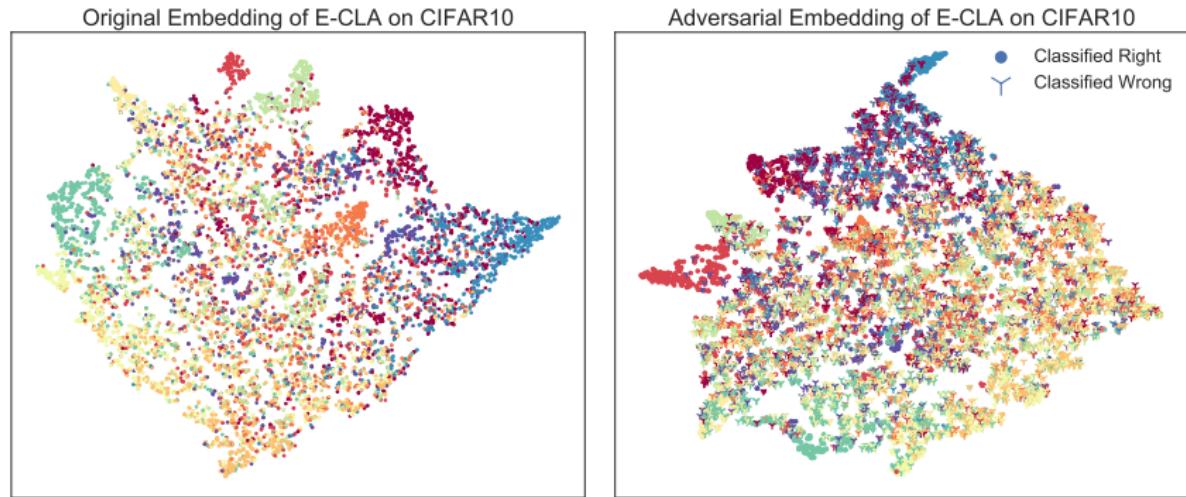
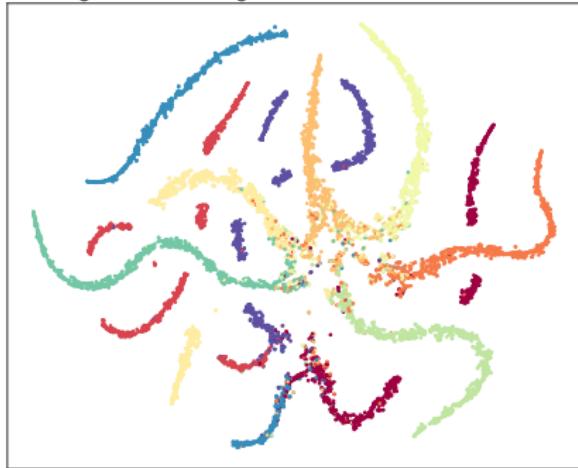


Figure: 2D embeddings for E-CLA on CIFAR10.

Embedding Visualization

Original Embedding of ER-classifier on CIFAR10



Adversarial Embedding of ER-classifier on CIFAR10

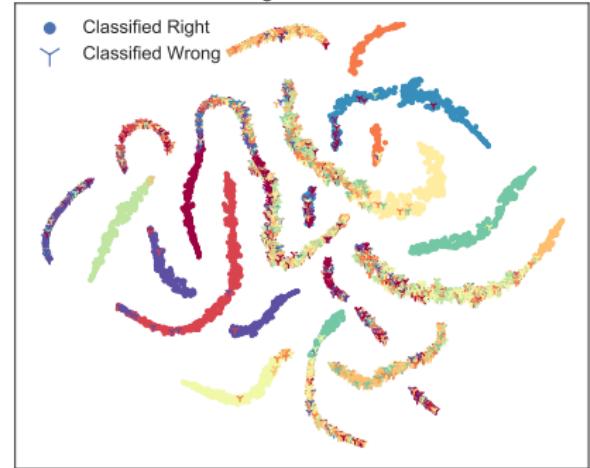


Figure: 2D embeddings for ER-Classifier on CIFAR10.