

Robust Deep Learning: Where Are We?

Ju Sun (Computer Sci. & Eng.)

Oct 20th, 2023

MnRI Colloquium



UNIVERSITY OF MINNESOTA

Driven to DiscoverSM

Success of deep learning (DL) not news anymore

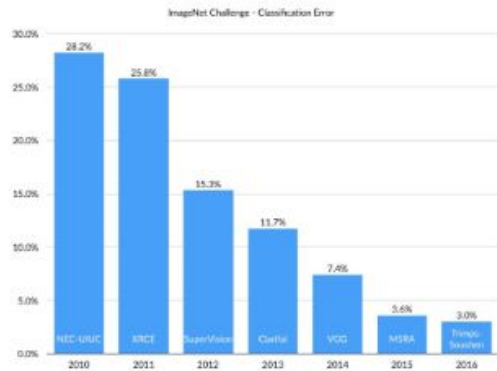


image classification

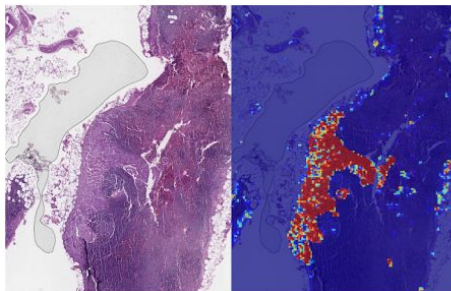


Go game (2017)

Commercial breakthroughs ...



self-driving vehicles credit: wired.com



healthcare credit: Google AI



smart-home devices credit: Amazon



robotics credit: Cornell U.

Robustness issues of DL not news anymore



"panda"

x

δ

FOOLING THE AI

Deep neural networks (DNNs) are brilliant at image recognition — but they can be easily hacked.

These stickers made an artificial-intelligence system read this stop sign as 'speed limit 45'.

Stop



Speed limit 45



Robustness issues across domains/tasks



"panda"

x



[Submitted on 06/01/2018]

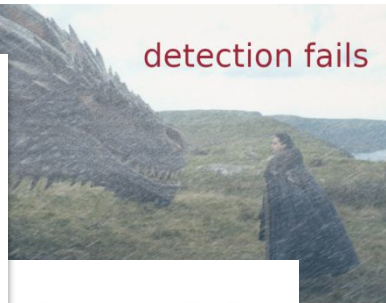
A Multilingual Inputs

Akshay Srinivasan

Adversarial
multilingual
input. Our

man and Hindi. While exact results differ depending on language/datasets, our key findings from these experiments can be summarized as follows:

1. NER models for all three languages are sensitive to adversarial input.
2. Adversarial fine-tuning and re-training could improve the performance of NER models both on original and adversarial test sets, without requiring additional manual labeled data.



detection fails

Adversarial

we performed a
small perturbations in the
(German and Hindi) are

Name entry Recognition

are not very robust to such changes, as indicated by the fluctuations in the overall F1 score as well as in a more fine-grained evaluation. With that knowledge, we further explored whether it is possible to improve the existing NER

Robustness issues across models

Tutorial

Foundational Robustness of Foundation Models

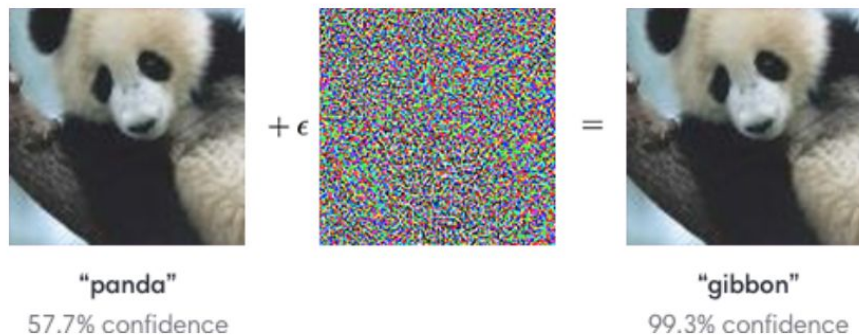
NeurIPS 2022

Abstract

Foundation models adopting the methodology of deep learning with pre-training on large-scale unlabeled data and finetuning with task-specific supervision are becoming a mainstream technique in machine learning. Although foundation models hold many promises in learning general representations and few-shot/zero-shot generalization across domains and data modalities, at the same time they raise unprecedented challenges and considerable risks in robustness and privacy due to the use of the excessive volume of data and complex neural network architectures. This tutorial aims to deliver a Coursera-like online tutorial containing comprehensive lectures, a hands-on and interactive Jupyter/Colab live coding demo, and a panel discussion on different aspects of trustworthiness in foundation models. More information can be found at <https://sites.google.com/view/neurips2022-frfm-tutorial>

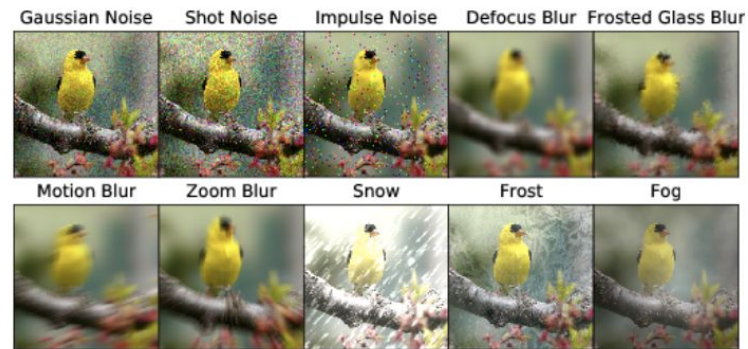
<https://research.ibm.com/publications/foundational-robustness-of-foundation-models>

Two kinds of robustness



credit: openai.com

Adversarial examples

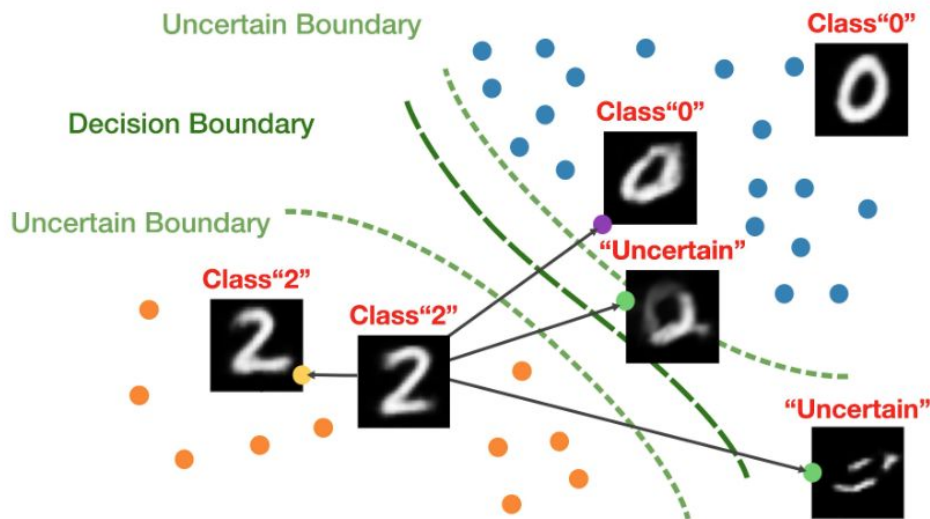


credit: ImageNet-C

Natural corruptions

Other dimensions in trustworthy AI

Trustworthiness: robustness, fairness, explainability, transparency



Boldness

Outline

- **Evaluation of adversarial robustness**

Optimization and Optimizers for Adversarial Robustness <https://arxiv.org/abs/2303.13401>

- **Fundamental challenges in evaluating & achieving robustness**

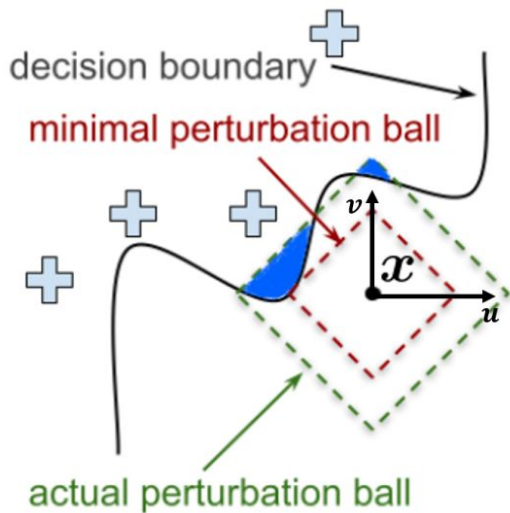
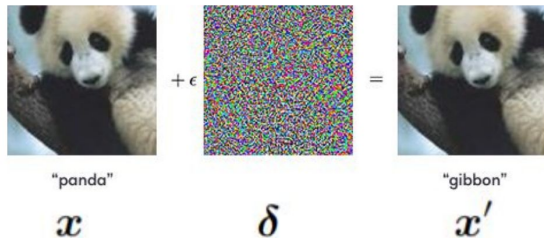
Optimization and Optimizers for Adversarial Robustness <https://arxiv.org/abs/2303.13401>

- **Selective prediction**

Margin As An Effective Confidence Score For Selective Classification Under Distribution Shifts
(Forthcoming)

- **Closing**

Robustness evaluation (RE)



$$\begin{aligned} & \max_{x'} \ell(y, f_{\theta}(x')) && \text{Maximize loss/error function} \\ \text{s. t. } & d(x, x') \leq \epsilon, && \text{Allowable perturbation} \\ & x' \in [0, 1]^n && \text{Valid image} \end{aligned}$$

$$\begin{aligned} & \min_{x'} d(x, x') && \text{Find robustness radius} \\ \text{s. t. } & \max_{i \neq y} f_{\theta}^i(x') \geq f_{\theta}^y(x'), && \text{On the decision boundary} \\ & x' \in [0, 1]^n && \text{Valid image} \end{aligned}$$

Report robust accuracy over an evaluation set

Constrained optimization problems

$$\begin{aligned} & \max_{\mathbf{x}'} \ell(\mathbf{y}, f_{\boldsymbol{\theta}}(\mathbf{x}')) \\ \text{s. t. } & d(\mathbf{x}, \mathbf{x}') \leq \varepsilon, \quad \mathbf{x}' \in [0, 1]^n \end{aligned}$$

$$\begin{aligned} & \min_{\mathbf{x}'} d(\mathbf{x}, \mathbf{x}') \\ \text{s. t. } & \max_{i \neq y} f_{\boldsymbol{\theta}}^i(\mathbf{x}') \geq f_{\boldsymbol{\theta}}^y(\mathbf{x}'), \quad \mathbf{x}' \in [0, 1]^n \end{aligned}$$

Both objective and constraint functions are **nonconvex** in general, e.g., when containing DL models

Projected gradient descent (PGD) for RE

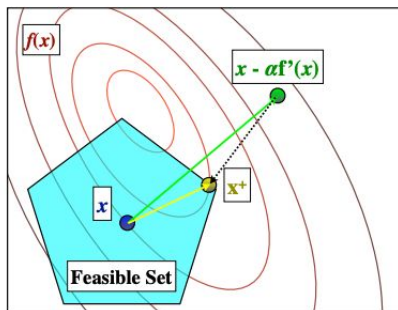
$$\max_{\mathbf{x}'} \ell(\mathbf{y}, f_{\theta}(\mathbf{x}'))$$

$$\text{s. t. } d(\mathbf{x}, \mathbf{x}') \leq \varepsilon, \quad \mathbf{x}' \in [0, 1]^n$$

$$\min_{\mathbf{x} \in \mathcal{Q}} f(\mathbf{x})$$

$$\mathbf{x}_{k+1} = P_{\mathcal{Q}}\left(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)\right)$$

$$P_{\mathcal{Q}}(\mathbf{x}_0) = \arg \min_{\mathbf{x} \in \mathcal{Q}} \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2 \quad \text{Projection operator}$$



Key hyperparameters:

- (1) step size
- (2) iteration number

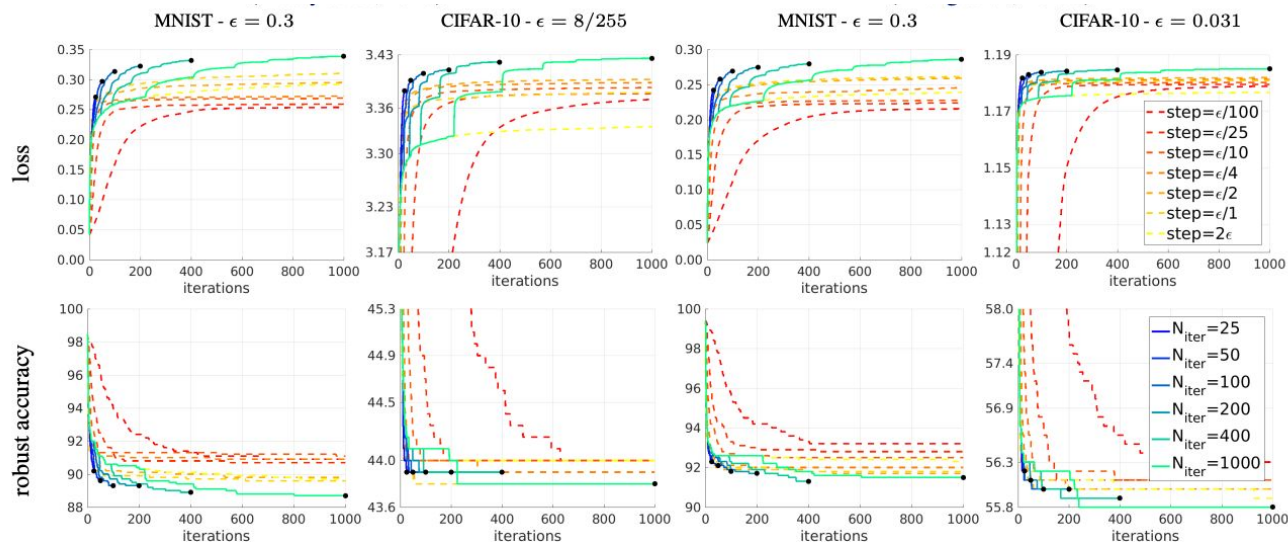
Algorithm 1 APGD

```

1: Input:  $f, S, x^{(0)}, \eta, N_{\text{iter}}, W = \{w_0, \dots, w_n\}$ 
2: Output:  $x_{\max}, f_{\max}$ 
3:  $x^{(1)} \leftarrow P_S(x^{(0)} + \eta \nabla f(x^{(0)}))$ 
4:  $f_{\max} \leftarrow \max\{f(x^{(0)}), f(x^{(1)})\}$ 
5:  $x_{\max} \leftarrow x^{(0)}$  if  $f_{\max} \equiv f(x^{(0)})$  else  $x_{\max} \leftarrow x^{(1)}$ 
6: for  $k = 1$  to  $N_{\text{iter}} - 1$  do
7:    $z^{(k+1)} \leftarrow P_S(x^{(k)} + \eta \nabla f(x^{(k)}))$ 
8:    $x^{(k+1)} \leftarrow P_S(x^{(k)} + \alpha(z^{(k+1)} - x^{(k)}))$ 
        $+ (1 - \alpha)(x^{(k)} - x^{(k-1)})$ 
9:   if  $f(x^{(k+1)}) > f_{\max}$  then
10:      $x_{\max} \leftarrow x^{(k+1)}$  and  $f_{\max} \leftarrow f(x^{(k+1)})$ 
11:   end if
12:   if  $k \in W$  then
13:     if Condition 1 or Condition 2 then
14:        $\eta \leftarrow \eta/2$  and  $x^{(k+1)} \leftarrow x_{\max}$ 
15:     end if
16:   end if
17: end for

```

Problem with projected gradient descent



$$\begin{aligned} & \max_{\mathbf{x}'} \ell(\mathbf{y}, f_{\theta}(\mathbf{x}')) \\ & \text{s.t. } d(\mathbf{x}, \mathbf{x}') \leq \epsilon, \quad \mathbf{x}' \in [0, 1]^n \end{aligned}$$

Tricky to set:
iteration number & step size
i.e., **tricky to decide where to stop**

Penalty methods for complicated d

$$\max_{\mathbf{x}'} \ell(\mathbf{y}, f_{\theta}(\mathbf{x}'))$$

$$\text{s. t. } d(\mathbf{x}, \mathbf{x}') \leq \varepsilon, \quad \mathbf{x}' \in [0, 1]^n$$

$$d(\mathbf{x}, \mathbf{x}') \doteq \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|_2$$

**perceptual
distance**

$$\text{where } \phi(\mathbf{x}) \doteq [\hat{g}_1(\mathbf{x}), \dots, \hat{g}_L(\mathbf{x})]$$

Projection onto the constraint is complicated

Penalty methods

$$\max_{\tilde{\mathbf{x}}} \mathcal{L}(f(\tilde{\mathbf{x}}), y) - \lambda \max\left(0, \|\phi(\tilde{\mathbf{x}}) - \phi(\mathbf{x})\|_2 - \epsilon\right)$$

Solve it for each fixed λ and then increase λ

Algorithm 2 Lagrangian Perceptual Attack (LPA)

```
1: procedure LPA(classifier network  $f(\cdot)$ , LPIPS distance  $d(\cdot, \cdot)$ , input  $\mathbf{x}$ , label  $y$ , bound  $\epsilon$ )
2:    $\lambda \leftarrow 0.01$ 
3:    $\tilde{\mathbf{x}} \leftarrow \mathbf{x} + 0.01 * \mathcal{N}(0, 1)$   $\triangleright$  initialize perturbations with random Gaussian noise
4:   for  $i$  in  $1, \dots, S$  do  $\triangleright$  we use  $S = 5$  iterations to search for the best value of  $\lambda$ 
5:     for  $t$  in  $1, \dots, T$  do  $\triangleright T$  is the number of steps
6:        $\Delta \leftarrow \nabla_{\tilde{\mathbf{x}}} [\mathcal{L}(f(\tilde{\mathbf{x}}), y) - \lambda \max(0, d(\tilde{\mathbf{x}}, \mathbf{x}) - \epsilon)]$   $\triangleright$  take the gradient of (5)
7:        $\hat{\Delta} = \Delta / \|\Delta\|_2$   $\triangleright$  normalize the gradient
8:        $\eta = \epsilon * (0.1)^{t/T}$   $\triangleright$  the step size  $\eta$  decays exponentially
9:        $m \leftarrow d(\tilde{\mathbf{x}}, \tilde{\mathbf{x}} + h\hat{\Delta})/h$   $\triangleright m \approx$  derivative of  $d(\tilde{\mathbf{x}}, \cdot)$  in the direction of  $\hat{\Delta}$ ;  $h = 0.1$ 
10:       $\tilde{\mathbf{x}} \leftarrow \tilde{\mathbf{x}} + (\eta/m)\hat{\Delta}$   $\triangleright$  take a step of size  $\eta$  in LPIPS distance
11:    end for
12:    if  $d(\tilde{\mathbf{x}}, \mathbf{x}) > \epsilon$  then
13:       $\lambda \leftarrow 10\lambda$   $\triangleright$  increase  $\lambda$  if the attack goes outside the bound
14:    end if
15:  end for
16:   $\tilde{\mathbf{x}} \leftarrow \text{PROJECT}(d, \tilde{\mathbf{x}}, \mathbf{x}, \epsilon)$ 
17:  return  $\tilde{\mathbf{x}}$ 
18: end procedure
```

Problem with penalty methods

Method	cross-entropy loss		margin loss	
	Viol. (%) ↓	Att. Succ. (%) ↑	Viol. (%) ↓	Att. Succ. (%) ↑
Fast-LPA	73.8	3.54	41.6	56.8
LPA	0.00	80.5	0.00	97.0
PPGD	5.44	25.5	0.00	38.5
PWCF (ours)	0.62	93.6	0.00	100

$$\begin{aligned} & \max_{\mathbf{x}'} \ell(\mathbf{y}, f_{\theta}(\mathbf{x}')) \\ \text{s. t. } & d(\mathbf{x}, \mathbf{x}') \leq \varepsilon, \quad \mathbf{x}' \in [0, 1]^n \\ & d(\mathbf{x}, \mathbf{x}') \doteq \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|_2 \\ \text{where } & \phi(\mathbf{x}) \doteq [\hat{g}_1(\mathbf{x}), \dots, \hat{g}_L(\mathbf{x})] \end{aligned}$$

LPA, Fast-LPA: penalty methods

PPGD: Projected gradient descent

PWCF, an optimizer with a principled stopping criterion on **stationarity** & **feasibility**

Penalty methods tend to encounter

large constraint violation (i.e., infeasible solution, known in optimization theory) or **suboptimal solution**

Unreliable optimization = Unreliable RE

Issues and answers

projected gradient descent

$$\min_{\mathbf{x} \in \mathcal{Q}} f(\mathbf{x})$$

$$\mathbf{x}_{k+1} = P_{\mathcal{Q}}\left(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)\right)$$

Issue: no principled stopping criterion
/step size rules

penalty methods

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{s. t. } g(\mathbf{x}) \leq 0$$

$$\min_{\mathbf{x}} f(\mathbf{x}) + \lambda \max(0, g(\mathbf{x}))$$

Solved with increasing λ sequence

Issue: infeasible solution

- Feasible & stationary solution **Stationarity and feasibility check:**
KKT condition
- Reasonable speed **Line search & 2nd order methods**
- A hidden problem: nonsmoothness

A principled solver for
constrained, nonconvex,
nonsmooth problems



Nonconvex, nonsmooth, constrained

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \quad \text{s.t. } c_i(\mathbf{x}) \leq 0, \forall i \in \mathcal{I}; \quad c_i(\mathbf{x}) = 0, \forall i \in \mathcal{E}.$$

**Penalty sequential quadratic programming
(P-SQP)**

$$\begin{aligned} \min_{d \in \mathbb{R}^n, s \in \mathbb{R}^p} \quad & \mu(f(x_k) + \nabla f(x_k)^\top d) + e^\top s + \frac{1}{2} d^\top H_k d \\ \text{s.t.} \quad & c(x_k) + \nabla c(x_k)^\top d \leq s, \quad s \geq 0, \end{aligned}$$

Advantage: 2nd order method (BFGS) \rightarrow high-precision solution

Principled line search, stationarity/feasibility check

Ref Curtis, Frank E., Tim Mitchell, and Michael L. Overton. "A BFGS-SQP method for nonsmooth, nonconvex, constrained optimization and its evaluation using relative minimization profiles." Optimization Methods and Software 32.1 (2017): 148-181.

Our PyGranso (and NCVX framework)

<https://ncvx.org/>

GRANVISO + PyTorch

 PyGRANSO

NCVX PyGRANSO
Documentation

🔍 Search the docs ...

Introduction

Installation

Settings

Examples



Home

 NCVX

NCVX Package

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \text{ s.t. } c_i(\mathbf{x}) \leq 0, \forall i \in \mathcal{I}; c_i(\mathbf{x}) = 0, \forall i \in \mathcal{E}$$

**First general-purpose solver for
constrained DL problems**

**NCVX: A General-Purpose Optimization Solver for
Constrained Machine and Deep Learning**

Buyun Liang, Tim Mitchell, Ju Sun

Strategies to speed up PyGranso for RE

$$\begin{aligned} & \max_{\mathbf{x}'} \ell(\mathbf{y}, f_{\theta}(\mathbf{x}')) \\ \text{s. t. } & d(\mathbf{x}, \mathbf{x}') \leq \varepsilon, \quad \mathbf{x}' \in [0, 1]^n \end{aligned}$$

Constraint folding: many constraints into few

$$\begin{aligned} h_j(\mathbf{x}) = 0 & \iff |h_j(\mathbf{x})| \leq 0, \\ c_i(\mathbf{x}) \leq 0 & \iff \max\{c_i(\mathbf{x}), 0\} \leq 0, \\ \mathcal{F}(|h_1(\mathbf{x})|, \dots, |h_i(\mathbf{x})|, \max\{c_1(\mathbf{x}), 0\}, \\ & \dots, \max\{c_j(\mathbf{x}), 0\}) \leq 0, \end{aligned}$$

$$\begin{aligned} & \min_{\mathbf{x}'} d(\mathbf{x}, \mathbf{x}') \\ \text{s. t. } & \max_{i \neq y} f_{\theta}^i(\mathbf{x}') \geq f_{\theta}^y(\mathbf{x}'), \quad \mathbf{x}' \in [0, 1]^n \end{aligned}$$


Two-stage optimization

1. **Stage 1 (selecting the best initialization):** Optimize the problems by PWCF with R different random initialization $\mathbf{x}^{(r,0)}$ for k iterations, where $r = 1, \dots, R$, and collect the final first-stage solution $\mathbf{x}^{(r,k)}$ for each run. Determine the best intermediate result $\mathbf{x}^{(*,k)}$ following [Algorithm 1](#).
2. **Stage 2 (optimization):** Warm start the optimization process with $\mathbf{x}^{*,k}$ until the stopping criterion is met [7](#) (i.e., reaching both the stationarity and feasibility tolerance, or reaching the MaxIter K).

First general-purpose, reliable solver for RE

$$\begin{aligned} & \max_{\mathbf{x}'} \ell(\mathbf{y}, f_{\theta}(\mathbf{x}')) \\ \text{s. t. } & d(\mathbf{x}, \mathbf{x}') \leq \varepsilon, \quad \mathbf{x}' \in [0, 1]^n \end{aligned}$$

Reliability

- SOTA methods 
No stopping criterion (only use maxit); step size scheduler
- PWCF (ours)
Principled line-search criterion and termination condition

$$\begin{aligned} & \min_{\mathbf{x}'} d(\mathbf{x}, \mathbf{x}') \\ \text{s. t. } & \max_{i \neq y} f_{\theta}^i(\mathbf{x}') \geq f_{\theta}^y(\mathbf{x}'), \quad \mathbf{x}' \in [0, 1]^n \end{aligned}$$

Generality

- SOTA methods
Can mostly only handle several lp metrics (l1, l2, linf)
- PWCF (ours)
Any differentiable metrics and both min and max forms
E.g., perceptual distance
$$d(\mathbf{x}, \mathbf{x}') \doteq \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|_2$$

where $\phi(\mathbf{x}) \doteq [\hat{g}_1(\mathbf{x}), \dots, \hat{g}_L(\mathbf{x})]$

A quick example

$$\begin{aligned} & \max_{\mathbf{x}'} \ell(\mathbf{y}, f_{\theta}(\mathbf{x}')) \\ \text{s. t. } & d(\mathbf{x}, \mathbf{x}') \leq \varepsilon, \quad \mathbf{x}' \in [0, 1]^n \end{aligned} \quad \text{where} \quad \begin{aligned} & d(\mathbf{x}, \mathbf{x}') \doteq \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|_2 \\ & \phi(\mathbf{x}) \doteq [\hat{g}_1(\mathbf{x}), \dots, \hat{g}_L(\mathbf{x})] \end{aligned}$$

Method	cross-entropy loss		margin loss	
	Viol. (%) ↓	Att. Succ. (%) ↑	Viol. (%) ↓	Att. Succ. (%) ↑
Fast-LPA	73.8	3.54	41.6	56.8
LPA	0.00	80.5	0.00	97.0
PPGD	5.44	25.5	0.00	38.5
PWCF (ours)	0.62	93.6	0.00	100

PyGranso has enabled much more

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \text{ s.t. } c_i(\mathbf{x}) \leq 0, \forall i \in \mathcal{I}; c_i(\mathbf{x}) = 0, \forall i \in \mathcal{E}$$

First general-purpose solver for constrained DL problems

 **PyGRANSO**
NCVX PyGRANSO
Documentation

🔍 Search the docs ...

Introduction

Installation

Settings

Examples



Home



NCVX Package

NCVX: A General-Purpose Optimization Solver for Constrained Machine and Deep Learning

Buyun Liang, Tim Mitchell, Ju Sun

Topology optimization

$$\min_{\theta, \mathbf{u}} \mathbf{u}^\top \mathbf{K}(\mathbf{G}_\theta(\beta)) \mathbf{u} \quad \text{s.t.} \quad \mathbf{K}(\mathbf{G}_\theta(\beta)) \mathbf{u} = \mathbf{f},$$

$$\sum_{i \in \Omega} [\mathbf{G}_\theta(\beta)]_i = v_0, \quad \mathbf{G}_\theta(\beta) \in \{0, 1\}^n$$

Imbalanced learning

$$\max_{\theta, t} \frac{\sum_{i=1}^N \mathbb{1}\{y_i = +1\} \mathbb{1}\{f_\theta(\mathbf{x}_i) > t\}}{\sum_{i=1}^N \mathbb{1}\{f_\theta(\mathbf{x}_i) > t\}}$$

$$\text{s.t.} \quad \frac{\sum_{i=1}^N \mathbb{1}\{y_i = +1\} \mathbb{1}\{f_\theta(\mathbf{x}_i) > t\}}{\sum_{i=1}^N \mathbb{1}\{y_i = +1\}} \geq \alpha$$

Constrained deep learning for the efficient discovery of stable solid-state materials

PIs: Chris Bartel (CEMS), Ju Sun (CS&E)

Background

Machine/deep learning (MDL) has emerged as a novel tool in material science and engineering (MSE).¹ MDL models in MSE can be broadly categorized as “property prediction models” (PPMs) or “interatomic potentials” (IPs). For the former, the goal is to learn the mapping between material representations and material properties (e.g., formation energy, band gap, etc.). These representations can be compositional,² requiring only the chemical formula (e.g., Al_2O_3), or structural,³ requiring the formula and the 3D arrangement of ions on a periodic lattice (e.g., Al_2O_3 in the corundum structure with specified coordinates for Al and O). IPs make use of a structural representation, but instead of learning to predict a single property, these models learn to predict the energies, forces, and stresses of an arbitrary configuration of ions on a lattice.⁴ Using this learned interatomic model, one can perform a set of tasks and analyses that are usually

Outline

- **Evaluation of adversarial robustness**

Optimization and Optimizers for Adversarial Robustness <https://arxiv.org/abs/2303.13401>

- **Fundamental challenges in evaluating & achieving robustness**

Optimization and Optimizers for Adversarial Robustness <https://arxiv.org/abs/2303.13401>

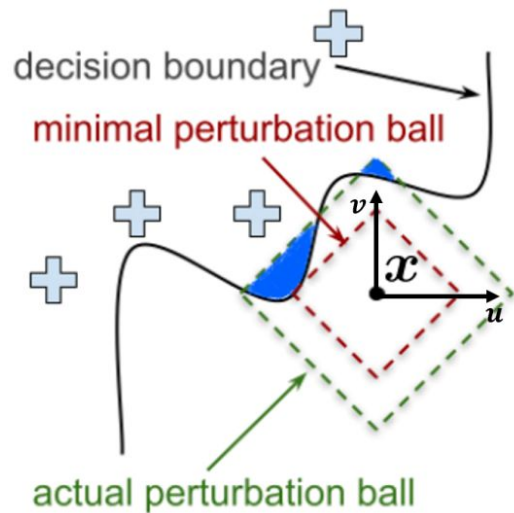
- **Selective prediction**

Margin As An Effective Confidence Score For Selective Classification Under Distribution Shifts
(Forthcoming)

- **Closing**

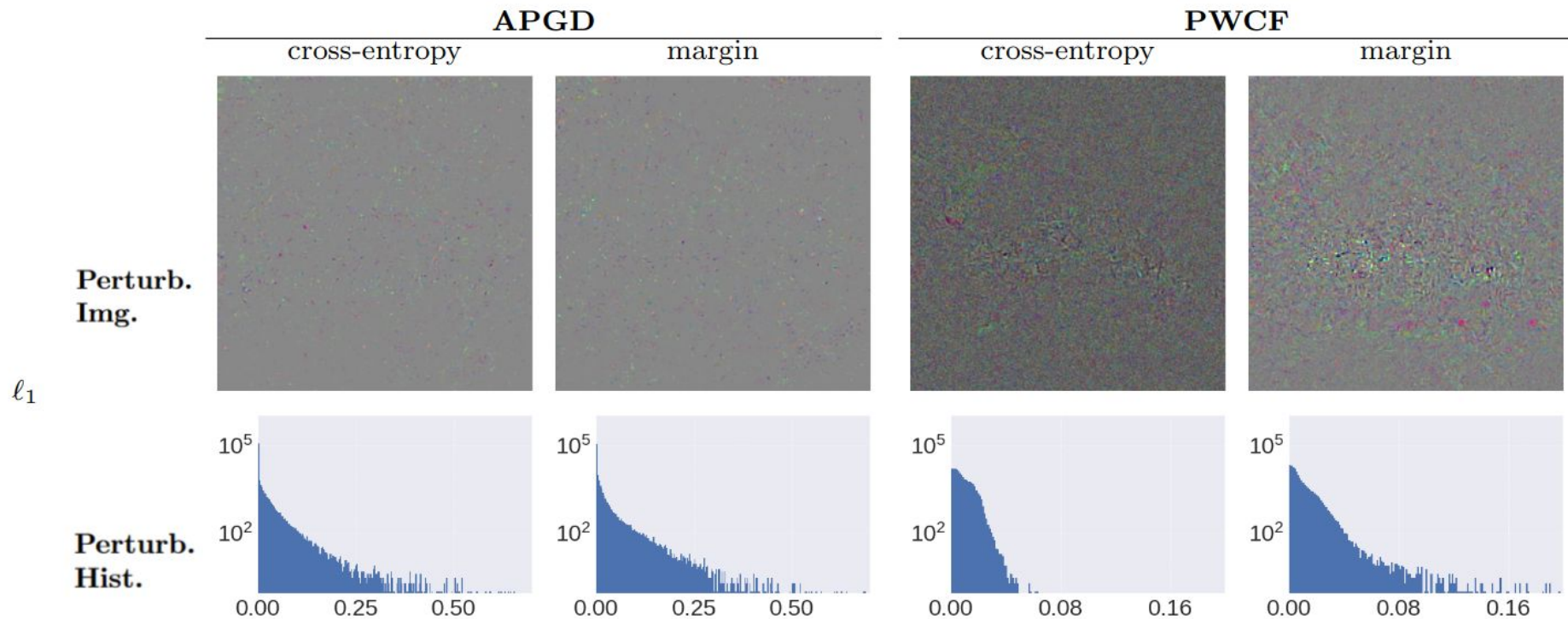
RE tractable even with PW/CF?

$$\begin{aligned} & \max_{\mathbf{x}'} \ell(\mathbf{y}, f_{\theta}(\mathbf{x}')) \\ \text{s.t. } & d(\mathbf{x}, \mathbf{x}') \leq \varepsilon, \quad \mathbf{x}' \in [0, 1]^n \end{aligned}$$



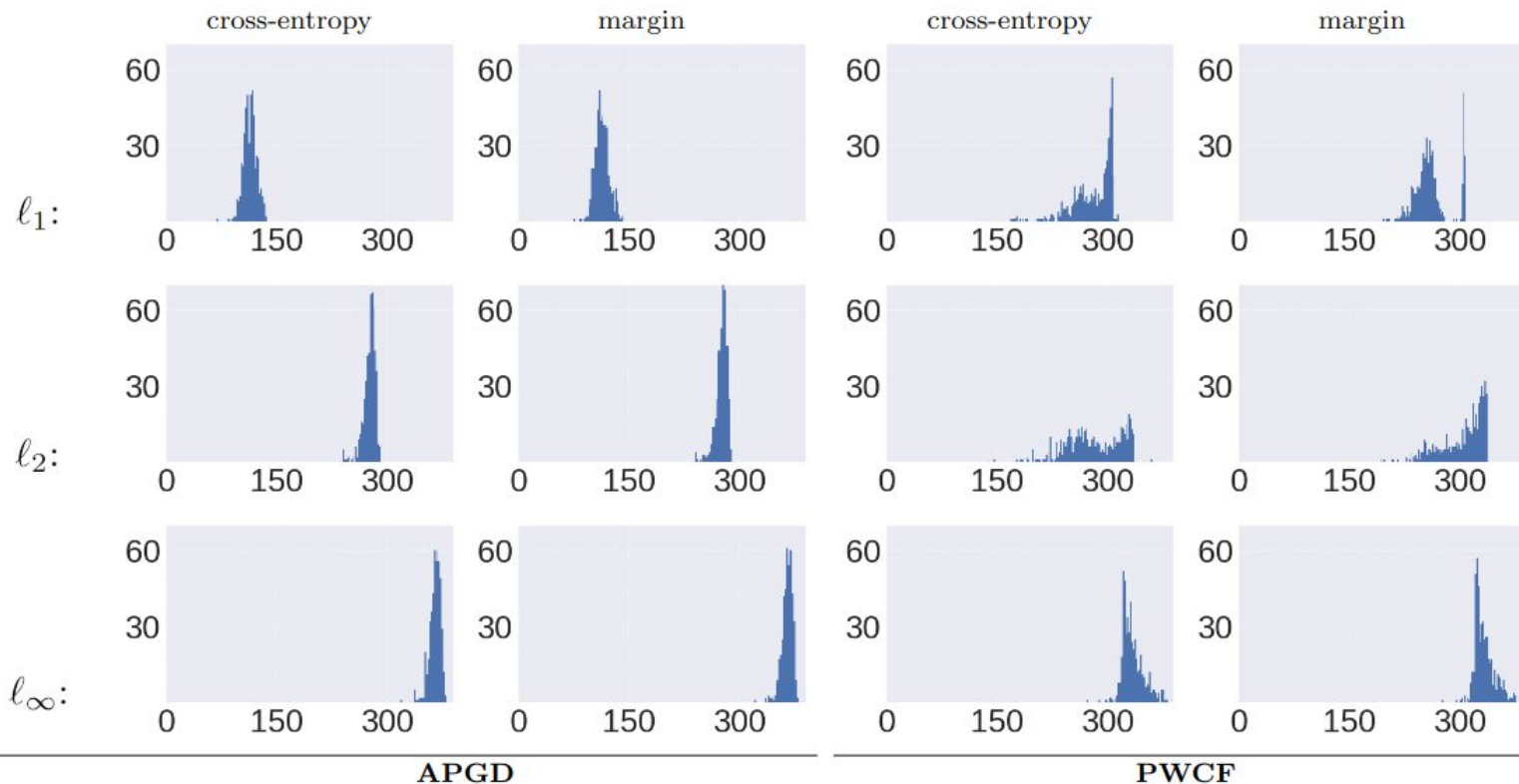
- Assuming 0-1 loss
- Typical **over-specification of ε** means there are potentially infinitely many solutions, with **different patterns**

Is the intuition right?



Is the intuition right?

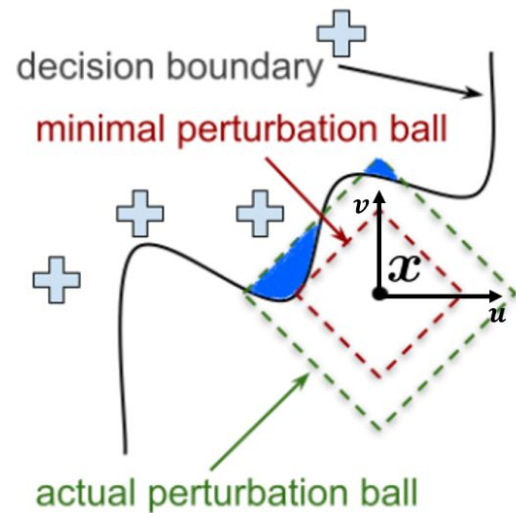
Measured by **sparsity levels** of the perturbations found



Implications - I

$$\begin{aligned} & \max_{\mathbf{x}'} \ell(\mathbf{y}, f_{\theta}(\mathbf{x}')) \\ \text{s. t. } & d(\mathbf{x}, \mathbf{x}') \leq \varepsilon, \quad \mathbf{x}' \in [0, 1]^n \end{aligned}$$

We need to **enumerate** all possible solutions if we want reliable RE using max-form



Take-away: Max-form RE is fundamentally intractable, unless a good ε is set—which is hard

Implications - II

$$\begin{aligned} & \max_{\mathbf{x}'} \ell(\mathbf{y}, f_{\theta}(\mathbf{x}')) \\ \text{s. t. } & \boxed{d(\mathbf{x}, \mathbf{x}') \leq \varepsilon, \quad \mathbf{x}' \in [0, 1]^n} \quad \mathbf{x}' \in \Delta(\mathbf{x}) \end{aligned}$$

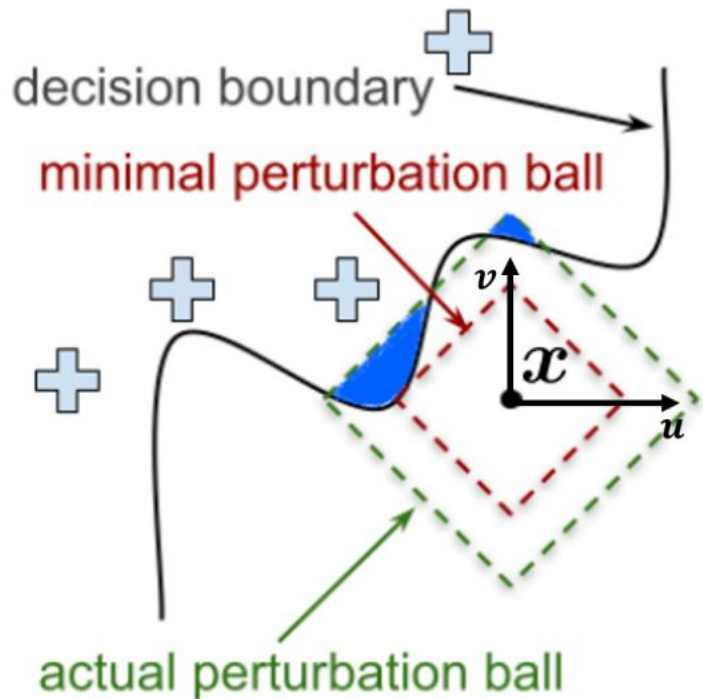
Adversarial training $\min_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} \max_{\mathbf{x}' \in \Delta(\mathbf{x})} \ell(\mathbf{y}, f_{\theta}(\mathbf{x}'))$

i.e., data augmentation with adversarial samples

We need to **enumerate** all possible patterns of adversarial samples if we want to achieve robustness, measured by the same d

Take-away: Adversarial training with the max-form augmentation won't achieve robustness

Any hopes remaining?



$$\begin{aligned} & \max_{x'} \ell(y, f_{\theta}(x')) \\ \text{s. t. } & d(x, x') \leq \varepsilon, \quad x' \in [0, 1]^n \end{aligned}$$

VS

$$\begin{aligned} & \min_{x'} d(x, x') \\ \text{s. t. } & \max_{i \neq y} f_{\theta}^i(x') \geq f_{\theta}^y(x'), \quad x' \in [0, 1]^n \end{aligned}$$

**Take-away: the min-form
(robustness radius) is
more promising**

Outline

- **Evaluation of adversarial robustness**

Optimization and Optimizers for Adversarial Robustness <https://arxiv.org/abs/2303.13401>

- **Fundamental challenges in evaluating & achieving robustness**

Optimization and Optimizers for Adversarial Robustness <https://arxiv.org/abs/2303.13401>

- **Selective prediction**

Margin As An Effective Confidence Score For Selective Classification Under Distribution Shifts
(Forthcoming)

- **Closing**

We have a long way to go



TRUSTWORTHY AI RESEARCH THRUSTS

DARPA experts estimate that research in the following areas will be essential to creating trustworthy technology:

- Foundational theory, to understand the art of the possible, bound the limits of particular system instantiations, and inform guardrails for AI systems in challenging domains such as national security;
- AI engineering, to predictably build systems that work as intended in the real world and not just in the lab; and
- Human-AI teaming, to enable systems to serve as fluent, intuitive, trustworthy teammates to people with various backgrounds.

<https://www.darpa.mil/work-with-us/ai-forward>

Safe Learning-Enabled Systems

PROGRAM SOLICITATION NSF 23-562



National Science Foundation

Directorate for Computer and Information Science and Engineering
Division of Information and Intelligent Systems
Division of Computing and Communication Foundations
Division of Computer and Network Systems



Open Philanthropy Project LLC



Good Ventures Foundation

Full Proposal Deadline(s) (due by 5 p.m. submitter's local time):

May 26, 2023

January 16, 2024

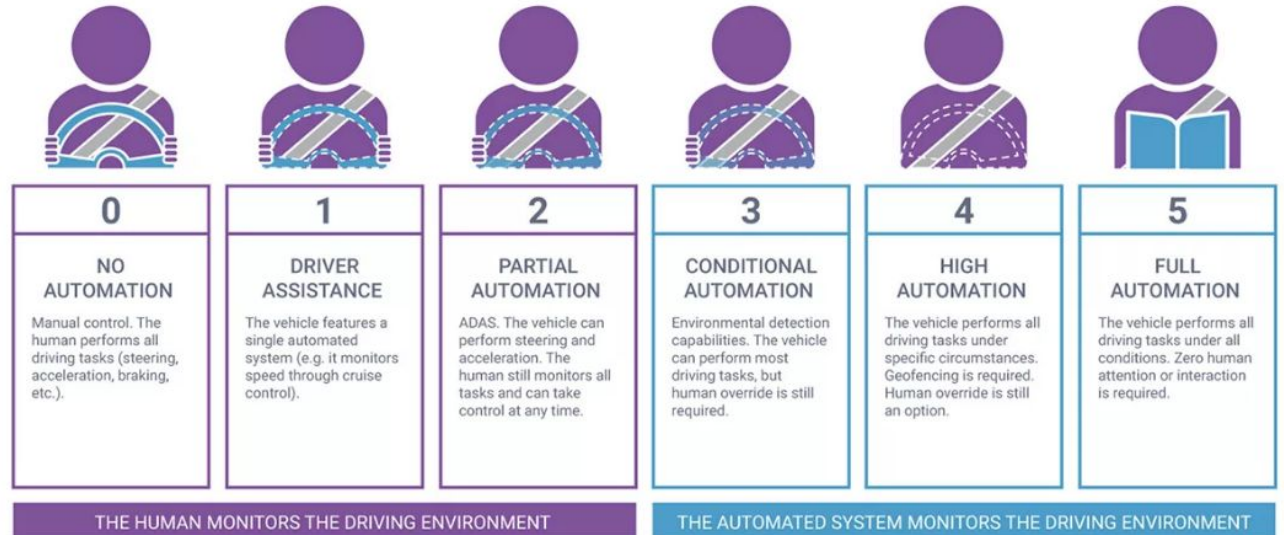
<https://www.nsf.gov/pubs/2023/nsf23562/nsf23562.htm>

Imperfect AI models can still be deployed



SYNOPSIS®

LEVELS OF DRIVING AUTOMATION



A crucial component: allowing AI to restrain itself

predictor $f : \mathcal{X} \rightarrow \mathbb{R}^K$ selector $g : \mathcal{X} \rightarrow \{0, 1\}$

$$(f, g)(\mathbf{x}) \triangleq \begin{cases} f(\mathbf{x}) & \text{if } g(\mathbf{x}) = 1; \\ \text{abstain} & \text{if } g(\mathbf{x}) = 0. \end{cases}$$

No prediction on uncertain samples and defer them to humans

$$g_\gamma(\mathbf{x}) = \mathbb{1}[s(\mathbf{x}) > \gamma]$$

Typically, selection by thresholding prediction confidence

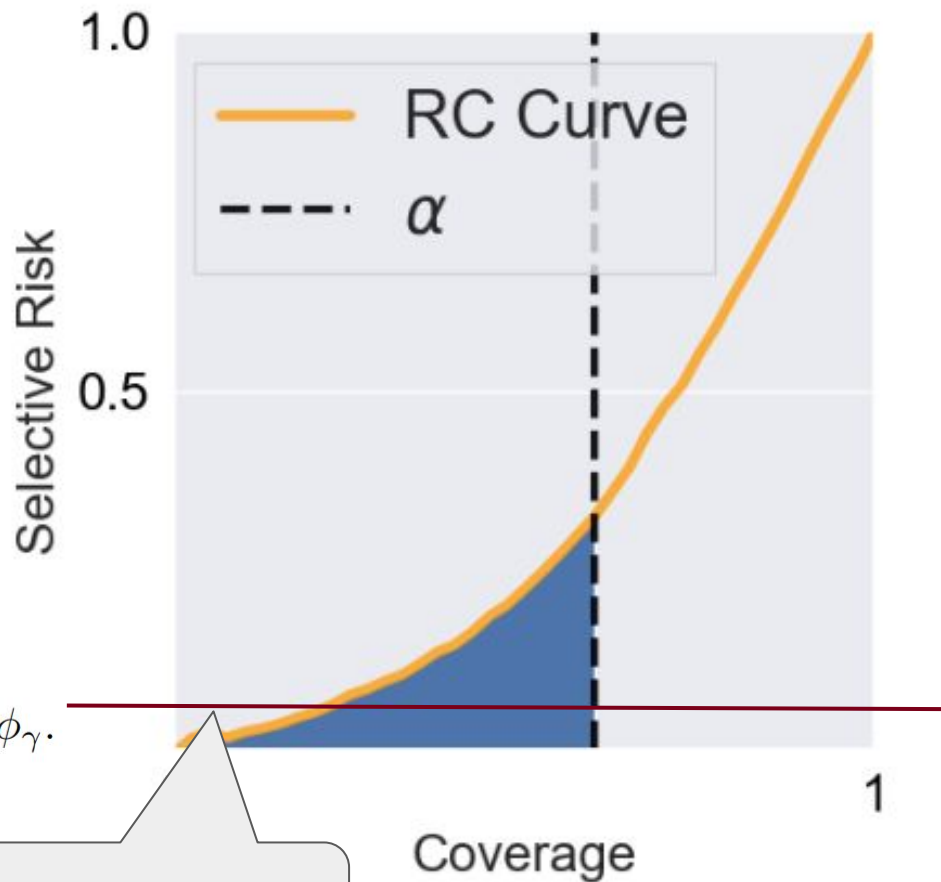
Risk-coverage tradeoff

$$(f, g)(\mathbf{x}) \triangleq \begin{cases} f(\mathbf{x}) & \text{if } g(\mathbf{x}) = 1; \\ \text{abstain} & \text{if } g(\mathbf{x}) = 0. \end{cases}$$

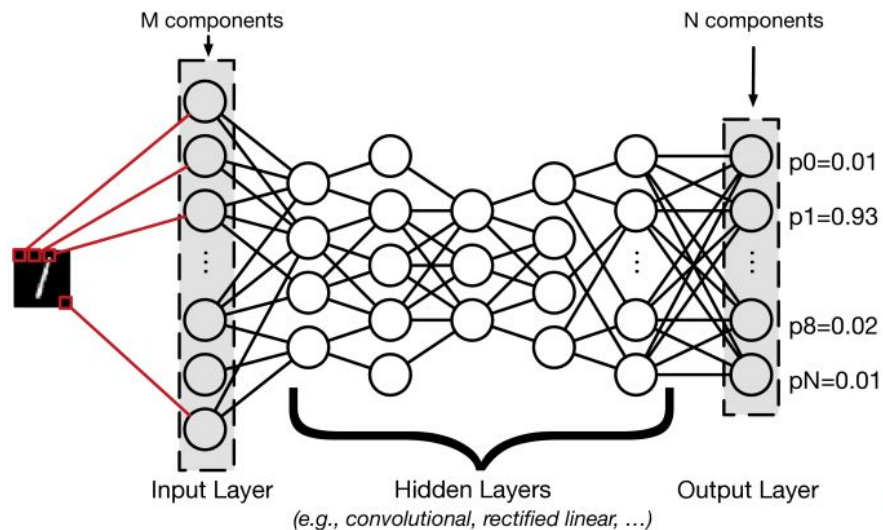
$$g_\gamma(\mathbf{x}) = \mathbb{1}[s(\mathbf{x}) > \gamma]$$

(coverage) $\phi_\gamma = \mathbb{E}_{\mathcal{D}}[g_\gamma(\mathbf{x})]$,

(selection risk) $R_\gamma = \mathbb{E}_{\mathcal{D}}[\ell(f(\mathbf{x}), y)g_\gamma(\mathbf{x})]/\phi_\gamma$.



Which confidence score?



$$\approx p(y = 1|x)$$

$\mathbf{z} \in \mathbb{R}^K$ contains the raw logits (RLs)

$$SR_{\max} \triangleq \max_i \sigma(z^i),$$

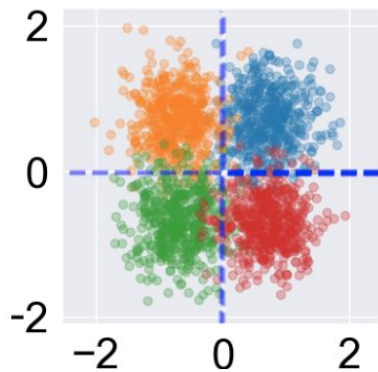
$$SR_{\text{doctor}} \triangleq (\|\sigma(\mathbf{z})\|_2^2 - 1) / \|\sigma(\mathbf{z})\|_2^2 = 1 - \|\sigma(\mathbf{z})\|_1 / \|\sigma(\mathbf{z})\|_2^2,$$

$$SR_{\text{ent}} \triangleq \sum_i \sigma(z^i) \log \sigma(z^i),$$

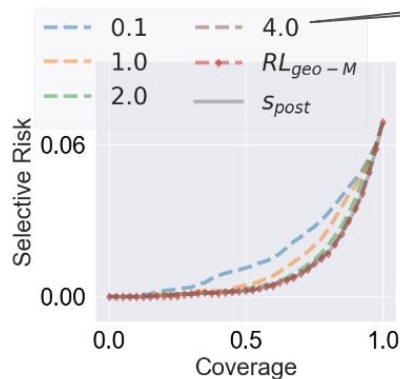
But are they good scores?

$z \in \mathbb{R}^K$ contains the raw logits (RLs)

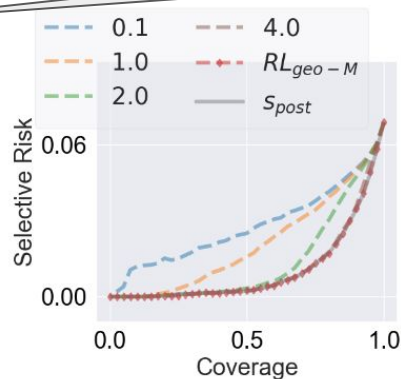
Scale factor applied to RLs



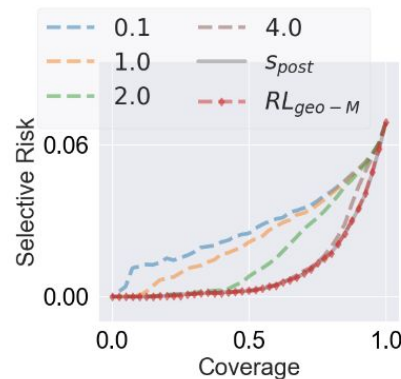
(a) Sample visualization



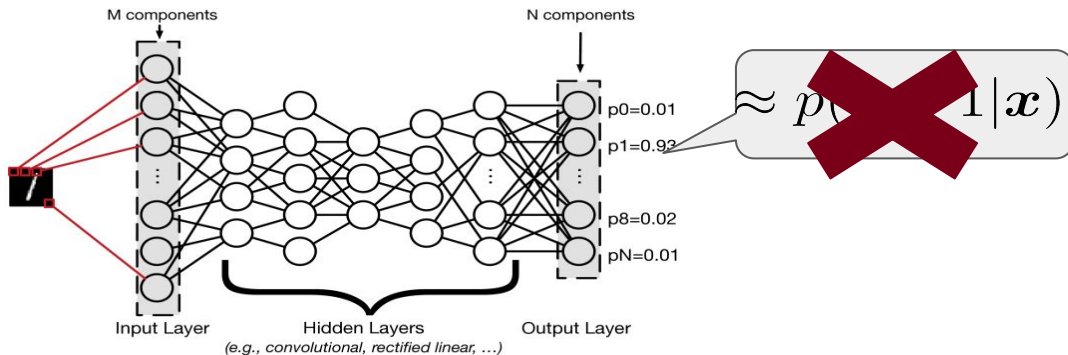
(b) SR_{\max}



(c) SR_{doctor}



(d) SR_{ent}



Calibration: align the outputs with the true posterior probs

Our margin-based scores

Signed dist to the separating hyperplane

Binary SVMs: $f(x) = w^\top x + b$

Geometric margin: $y(w^\top x + b) / \|w\|_2$

Multiclass SVMs: $f(x) = W^\top x + b$

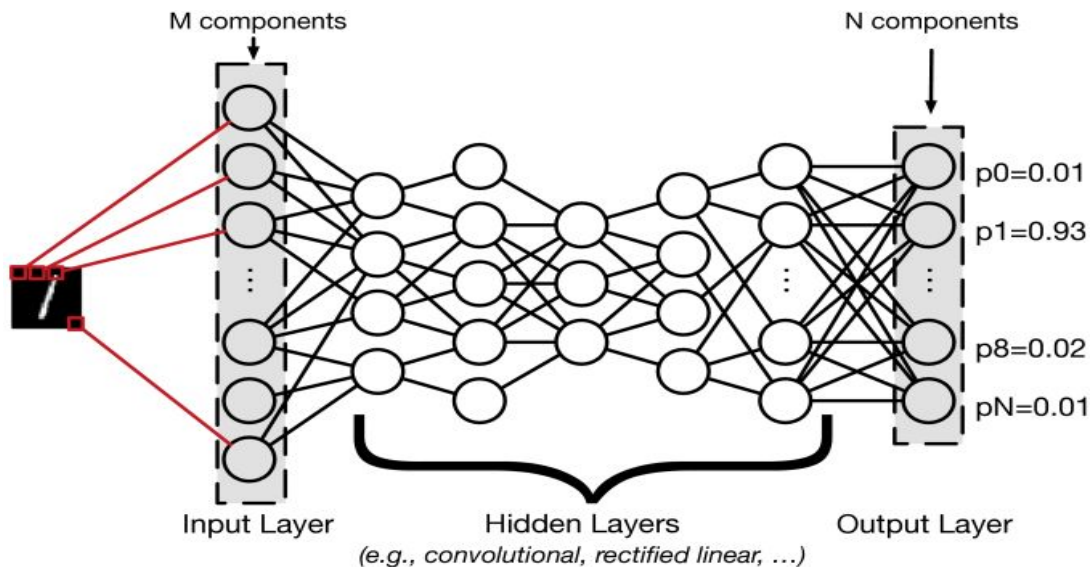
Geometric margin: $\frac{w_{y'}^\top x + b_{y'}}{\|w_{y'}\|_2} - \max_{j \in \{1, \dots, K\} \setminus y'} \frac{w_j^\top x + b_j}{\|w_j\|_2}$

Confidence margin: $(w_{y'}^\top x + b_{y'}) - \max_{i \in \{1, \dots, K\} \setminus y'} (w_i^\top x + b_i)$

Difference of dists between the two nearest hyperplanes

These scores are not affected by the logit scaling

Our margin-based scores



Geometric margin:

$$\frac{\mathbf{w}_{y'}^\top \mathbf{x} + b_{y'}}{\|\mathbf{w}_{y'}\|_2} - \max_{j \in \{1, \dots, K\} \setminus y'} \frac{\mathbf{w}_j^\top \mathbf{x} + b_j}{\|\mathbf{w}_j\|_2}$$

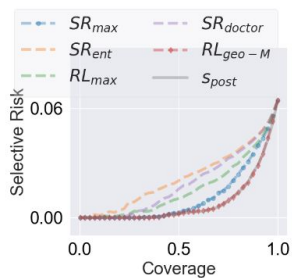
Confidence margin:

$$(\mathbf{w}_{y'}^\top \mathbf{x} + b_{y'}) - \max_{i \in \{1, \dots, K\} \setminus y'} (\mathbf{w}_i^\top \mathbf{x} + b_i)$$

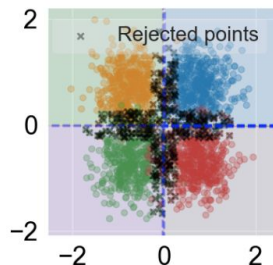
Apply them to the RLs \mathcal{Z}

Benefit: We don't need to worry about the scale of \mathcal{Z}

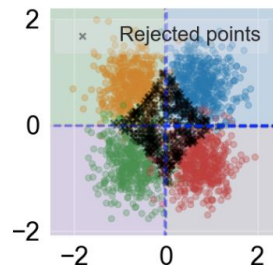
Additional benefit: robustness



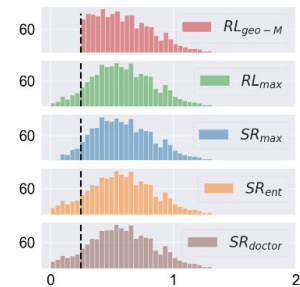
(a-1) RC curves



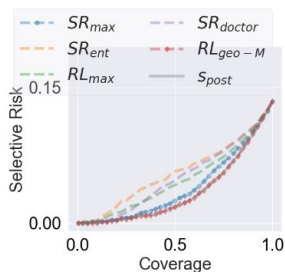
(b-1) RL_{geo-M}



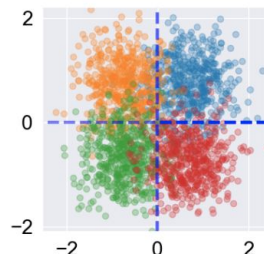
(c-1) SR_{max}



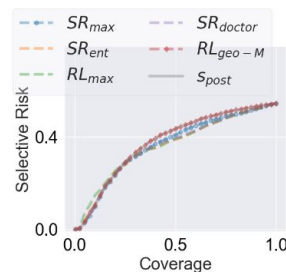
(d-1) Robustness radius



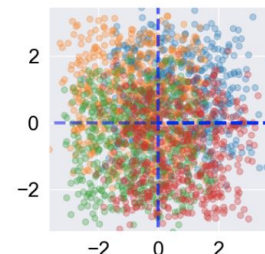
(a-2) RC curves



(b-2) Sample visualization



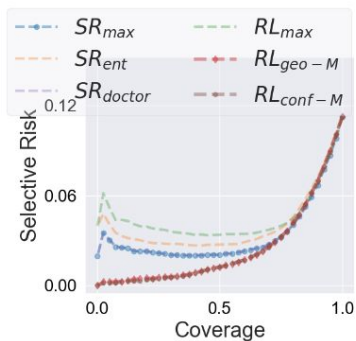
(a-3) RC curves



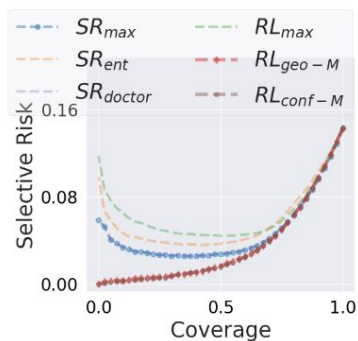
(b-3) Sample visualization

On real data

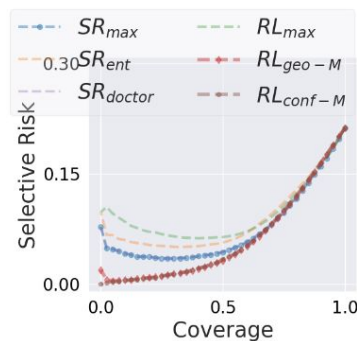
ImageNet vs ImageNet-C



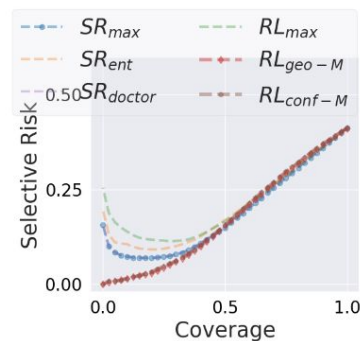
(a) IN (Clean)



(b) Gaussian blur Lv.1



(c) Gaussian blur Lv.3



(d) Gaussian blur Lv.5

	IN (Clean)			Gaussian Blur			Brightness			Fog			Snow		
α	0.1	0.5	1	0.1	0.5	1	0.1	0.5	1	0.1	0.5	1	0.1	0.5	1
$RL_{\text{conf-M}}$	0.16	0.53	2.39	0.37	1.31	<u>6.05</u>	0.21	0.72	3.35	0.14	0.79	4.21	0.17	0.95	4.80
$RL_{\text{geo-M}}$	<u>0.27</u>	<u>0.59</u>	<u>2.43</u>	<u>0.57</u>	<u>1.36</u>	6.04	<u>0.33</u>	<u>0.79</u>	<u>3.39</u>	<u>0.25</u>	<u>0.86</u>	<u>4.22</u>	<u>0.34</u>	<u>1.02</u>	<u>4.81</u>
RL_{max}	5.54	4.05	4.57	9.74	7.38	9.52	7.38	5.17	6.06	7.74	5.77	7.01	9.44	6.44	7.90
SR_{max}	3.19	2.40	3.38	5.02	4.02	7.39	4.07	2.90	4.53	3.92	3.07	5.37	5.35	3.67	6.13
SR_{ent}	4.28	3.13	4.04	6.80	5.63	8.71	5.51	4.01	5.48	5.56	4.37	6.42	7.29	5.07	7.27
SR_{doctor}	3.21	2.38	3.40	5.05	4.05	7.47	4.10	2.93	4.58	3.95	3.10	5.42	5.39	3.71	6.20

Outline

- **Evaluation of adversarial robustness**

Optimization and Optimizers for Adversarial Robustness <https://arxiv.org/abs/2303.13401>

- **Fundamental challenges in evaluating & achieving robustness**

Optimization and Optimizers for Adversarial Robustness <https://arxiv.org/abs/2303.13401>

- **Selective prediction**

Margin As An Effective Confidence Score For Selective Classification Under Distribution Shifts
(Forthcoming)

- **Closing**

Closing

- A long way to go for DL robustness
- Selective prediction crucial for deploying imperfect AI

