

Deep Learning with Nontrivial Constraints

Ju Sun

McKnight Land-Grant Professor
Computer Science & Engineering

Dec 10, 2025

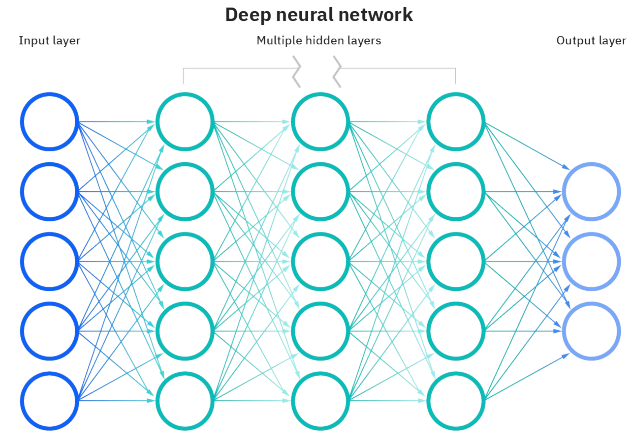
Industrial and Systems Engineering Graduate Seminar

Research in the group



<https://glovex.umn.edu/>

(Machine) **Learning**, (Numerical) **Optimization**, (Computer) **Vision**, healthcare, + X

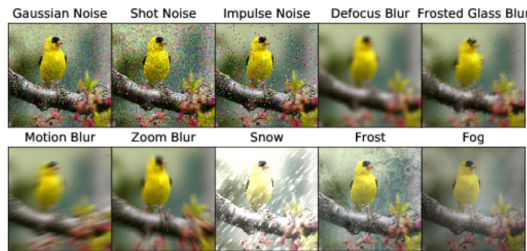
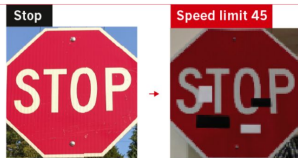


Our research themes

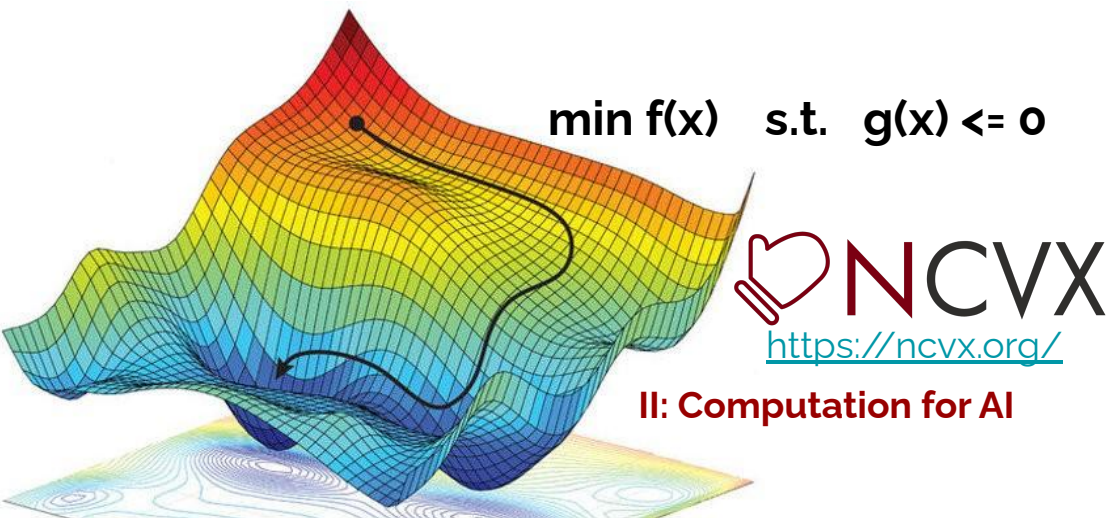
FOOLING THE AI

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

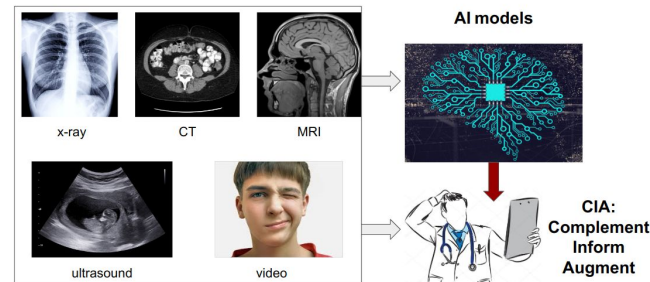
These stickers made an artificial-intelligence system read this stop sign as 'speed limit 45'.



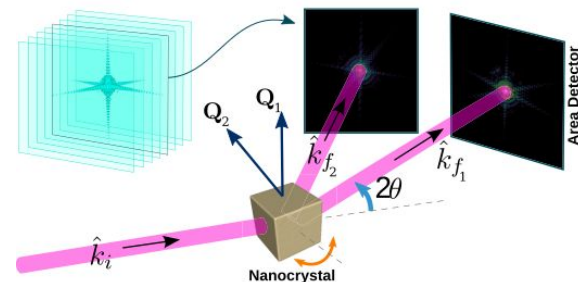
I: Trustworthy AI



II: Computation for AI



III: AI for Healthcare




IV: AI for Science and Engineering

Robustness & Deployability of Deep Learning Models

First general-purpose, reliable solver for robust evaluation

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

Reliability

- SOTA methods  **ROBUSTBENCH**
A standardized benchmark for adversarial robustness
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 $\phi(\mathbf{x}) \doteq [\hat{g}_1(\mathbf{x}), \dots, \hat{g}_L(\mathbf{x})]$

Selective predictions for safety and deployability

predictor $f : \mathcal{X} \rightarrow \mathbb{R}^K$

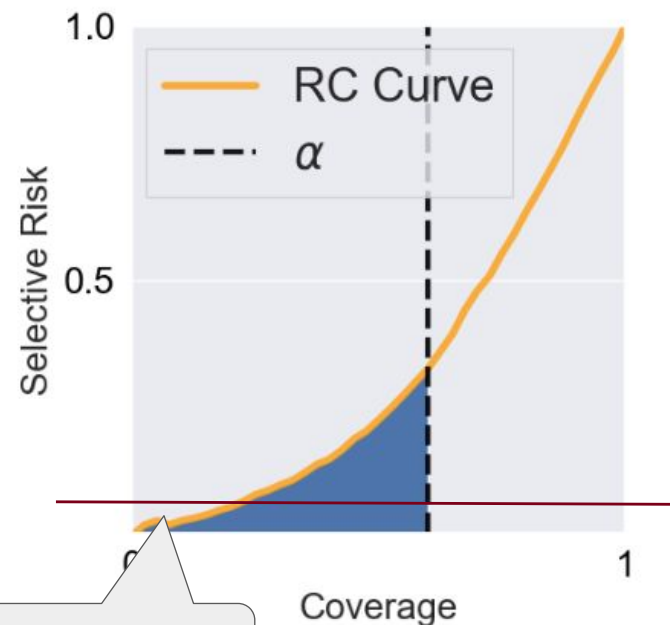
selector $g : \mathcal{X} \rightarrow \{0, 1\}$

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

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



High-stakes corner

Computation for Deep Learning with Nontrivial Constraints



NCVX PyGRANSO
Documentation

Search the docs ...

Introduction

Installation

Settings

Examples

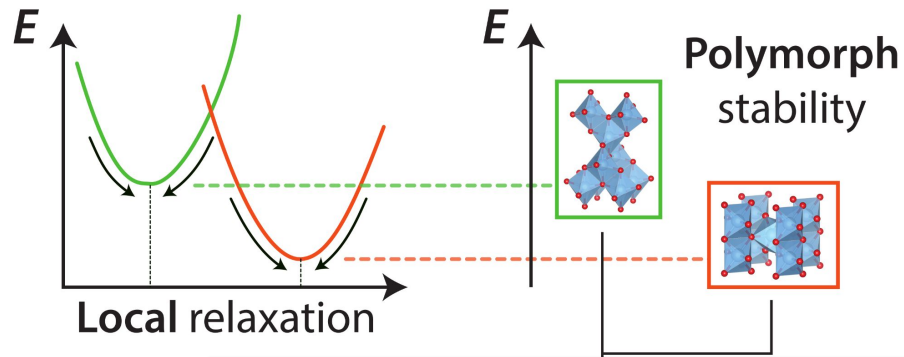
Home

<https://ncvx.org/>



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

- First general-purpose solver for hard-constrained deep learning problems
- Recently updated to be compatible with PyTorch 2.8



Enabling

Robustness evaluation

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

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

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

Stability
respe
phase se

- **(W-CSTR-T) Weak explicit constraints informed by phase transition:** For each composition, the energy/atom of any non-ground-state polymorph is greater than that of the ground-state polymorph
- **(W-CSTR-S) Weak explicit constraints informed by phase separation:** For each chemical space, the energy/atom of any TUS material is above the lower convex envelope in the composition-energy space

Settings

Examples

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{y}_i, f_{\theta}(\mathbf{x}_i)) + \Omega(\theta) \quad \text{s. t. } f_{\theta}(\mathbf{x}_p) \geq f_{\theta}(\mathbf{x}_q) \quad \forall (p, q) \in \mathcal{O}_{\text{W-CSTR-T}},$$

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

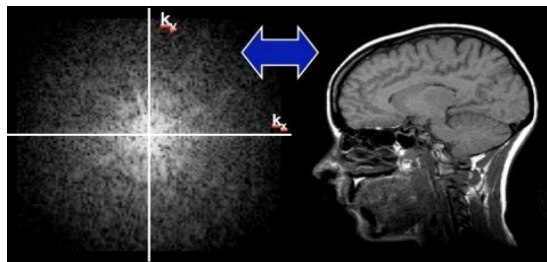
Buyun Liang, Tim Mitchell, Ju Sun

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

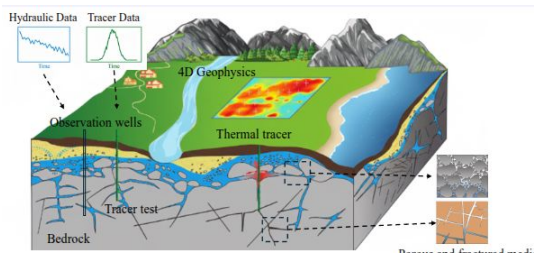
Yash, Le, Zhong, Buyun

Generative Models for Scientific/Engineering Inverse Modeling

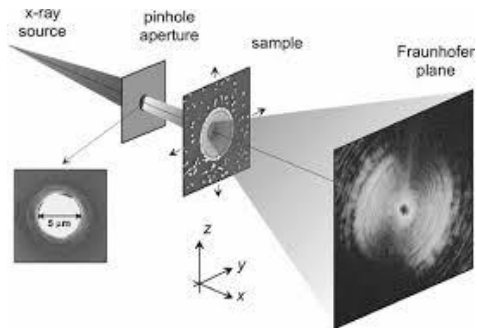
(inverse problems, reconstruction, data assimilation, inverse design/control, conditional generation, estimation ...)



MRI reconstruction



Geophysical inversion



Coherent diffraction imaging



Inverse shape control

Given $\mathbf{y} \approx f(\mathbf{x})$, recover \mathbf{x}

$$\min_{\mathbf{x}} \ell(\mathbf{y}, f(\mathbf{x})) + \Omega(\mathbf{x})$$

Challenges:

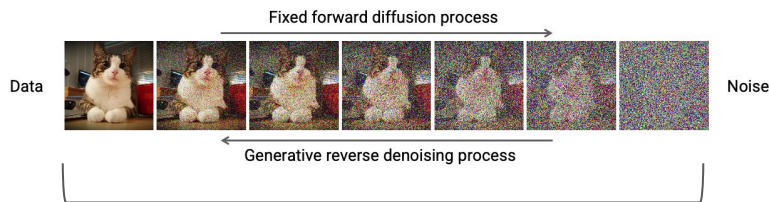
- Linear vs. **nonlinear** f
- Unconstrained vs. **constrained** (e.g., PDE $(x, y) = 0$)
- Explicit vs. **implicit** f

A **plug-in principle** for leveraging **pretrained** deep generative models to solve IM

object-only datasets $\{x_i\}_{i=1,\dots,N}$



Distribution learning via
deep generative models



\mathcal{G}_θ
Plug in pretrained deep
generative priors

$$\min_z \mathcal{L}(z) \doteq \ell(y, \mathcal{A} \circ \mathcal{G}_\theta(z)) + \Omega \circ \mathcal{G}_\theta(z)$$

DMPlug: A Plug-in Method for Solving Inverse Problems with Diffusion Models

Hengkang Wang, Xu Zhang, Taihui Li, Yuxiang Wan, Tiancong Chen, Ju Sun

NeurIPS'24, with **domain-specific** priors

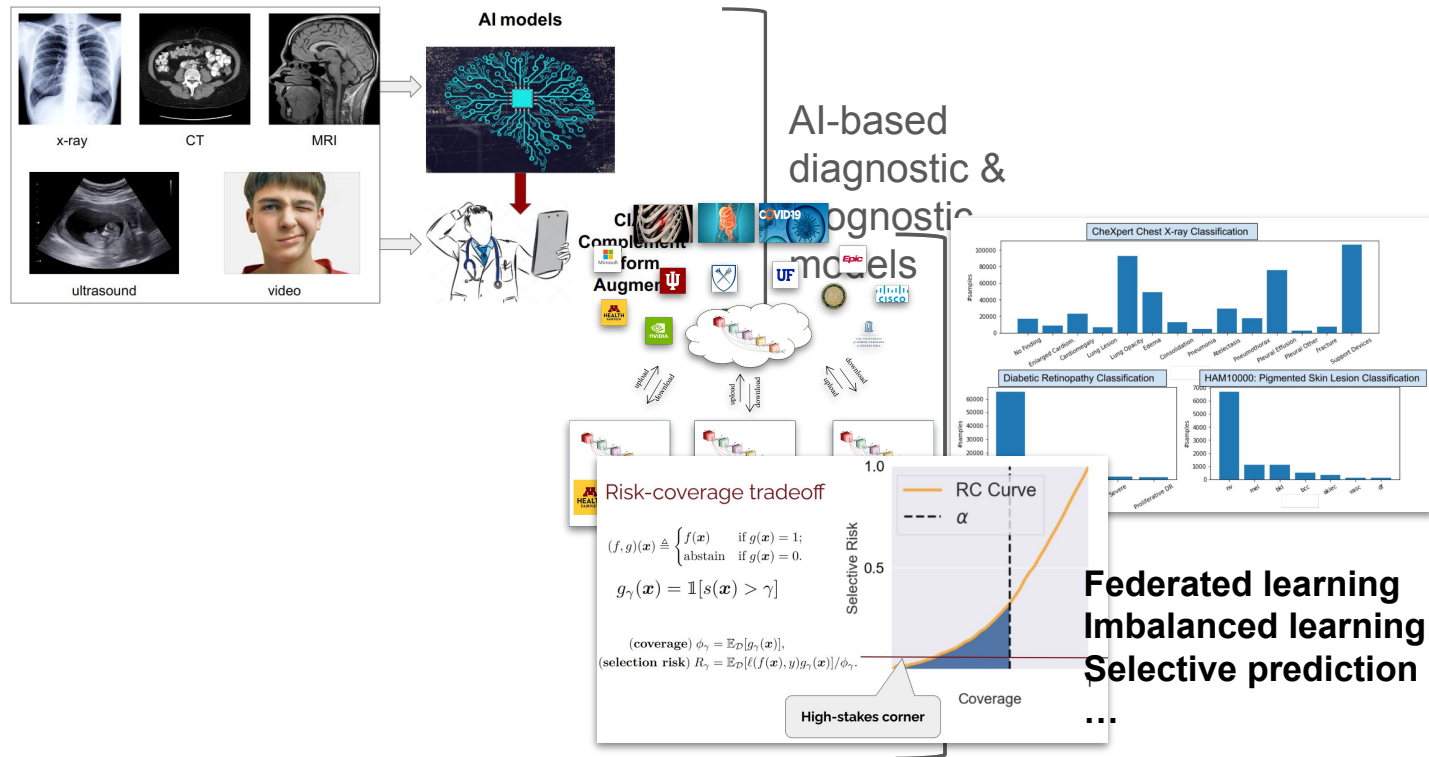
Saving Foundation Flow-Matching Priors for Inverse Problems

Yuxiang Wan, Ryan Devera, Wenjie Zhang, Ju Sun

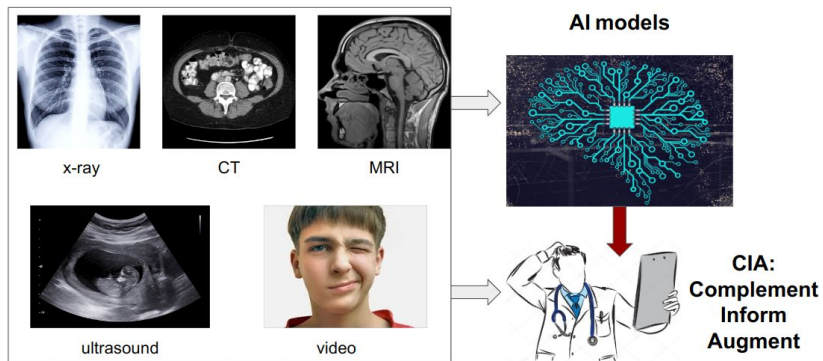
Submitted to ICLR'25, with **foundation**
priors (e.g., stable diffusion models)

Challenge: Foundation generative models are
powerful because **they're not specific**

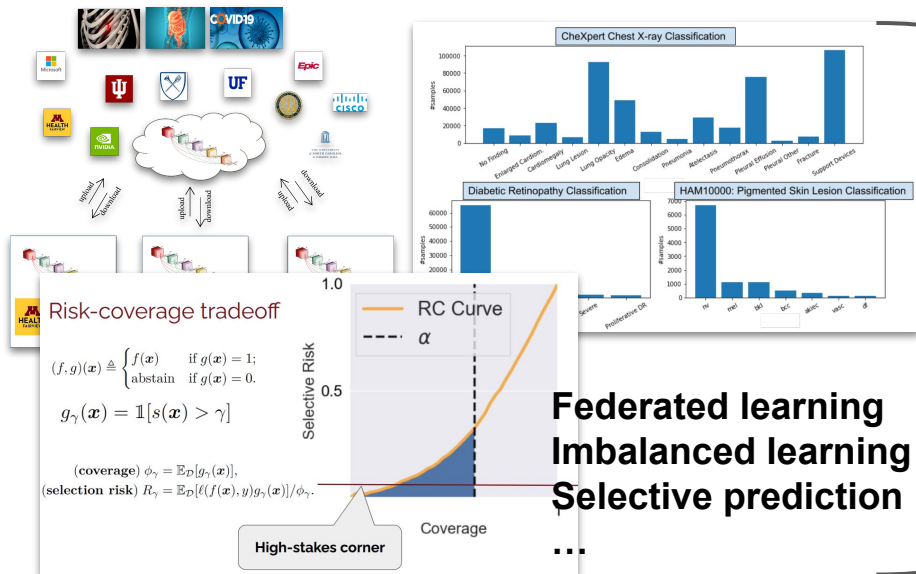
Applied and basic research in AI for Healthcare



Basic research to address data sparsity, data imbalance, and AI safety in healthcare (& beyond)



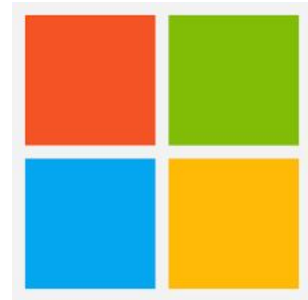
AI-based
diagnostic
models



Basic
research to
address data
sparsity, data
imbalance,
and AI safety
in healthcare
(& beyond)

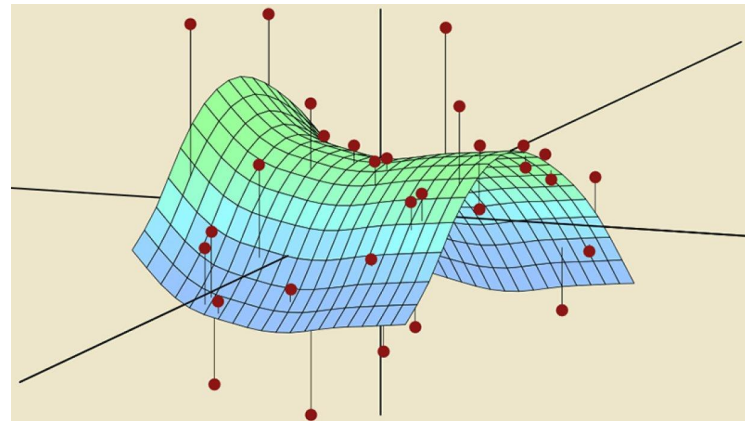
**Supported by
3+ (4th
forthcoming)
active NIH R01
grants**

Thanks to our funders



Practical Optimization for Deep Learning with Nontrivial Constraints

Three fundamental questions in DL



- **Approximation:** is it powerful, i.e., the \mathcal{H} large enough for all possible weights?

- **Optimization:** how to solve

$$\min_{\mathbf{w}'_i, \mathbf{b}'_i} \frac{1}{n} \sum_{i=1}^n \ell[\mathbf{y}_i, \{\text{NN}(\mathbf{w}_1, \dots, \mathbf{w}_k, b_1, \dots, b_k)\}(\mathbf{x}_i)]$$

- **Generalization:** does the learned NN work well on “similar” data?

Isn't it solved?

Base class

CLASS `torch.optim.Optimizer(params, defaults)` [Source]

Base class for all optimizers.

• WARNING

Parameters need to be specified as collections consistent between runs. Examples of objects and iterators over values of dictionaries.

Parameters:

- **params** (*iterable*) – an iterable of `Tensor` objects. Tensors should be optimized.
- **defaults** – (dict): a dict containing default values for the optimizer when a parameter group doesn't specify

Algorithms

Adadelta

Implements Adadelta algorithm.

Adagrad

Implements Adagrad algorithm.

Adamax

Implements Adamax algorithm (a variant of Adam based on infinity norm).

ASGD

Implements Averaged Stochastic Gradient Descent.

L-BFGS

Implements L-BFGS algorithm, heavily inspired by `minFunc`.

NAdam

Implements NAdam algorithm.

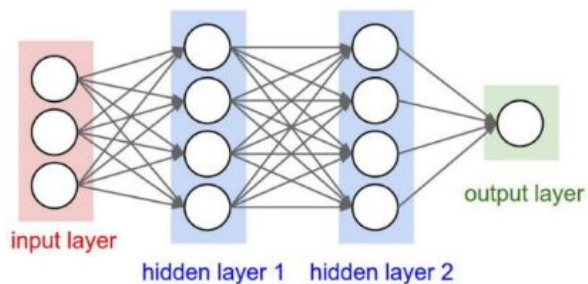
RAdam

Implements RAdam algorithm.

Algorithm

When DL meets constraints

Artificial neural networks



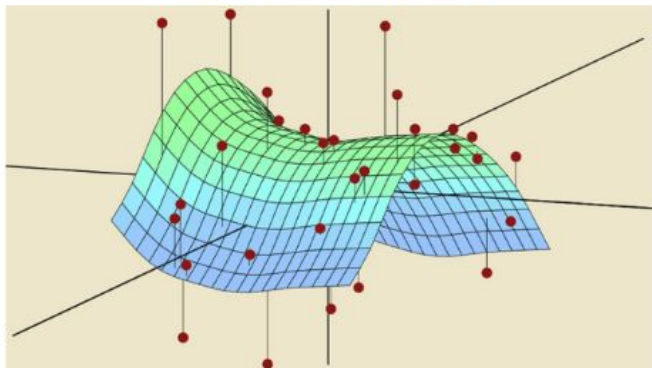
Unconstrained optimization

$$\min_{w'_i s, b'_i s} \frac{1}{n} \sum_{i=1}^n \ell [y_i, \{\text{NN}(w_1, \dots, w_k, b_1, \dots, b_k)\}(x_i)]$$
$$\min_x f(x) \quad \text{“Solved”}$$

Constrained optimization

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

largely “unsolved”



used to approximate nonlinear functions

This talk is about GAPS

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

largely “unsolved”



An imaginary chat between a PhD student working in deep learning (**DLP**) and a PhD student working in optimization (**OP**)

DLP: Man, I've solved a constrained DL problem recently

OP: Oh, that's a hard problem

DLP: Really? I actually did it

OP: How?

DLP: My problem is $\min_x f(x)$, s.t. $g(x) \leq 0$. I put $g(x)$ as a penalty and then call ADAM

OP: Are you sure it works?

DLP: Yes, the performance is improved and my paper is published at ICML

OP: Why don't you try augmented Lagrangian methods?

DLP: No implementation in Pytorch. Is it possible we work out some theory about my method?

OP: I think it's hard. It's not convex

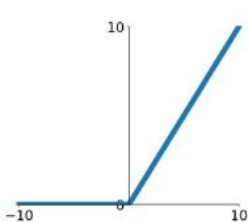
Outline

Constrained deep learning: CDL

- **What, how, and why for CDL**
- No good solvers for CDL yet
- Granso and PyGranso
- PyGranso in action
- Outlook

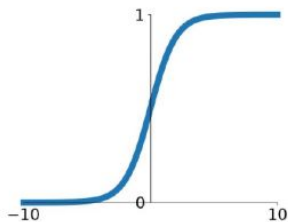
DL with simple constraints

Embedding constraints into DL models



ReLU
(Rectified Linear Unit)

Nonnegativity



Sigmoid

[0, 1]

$$\mathbf{z} \mapsto \left[\frac{e^{z_1}}{\sum_j e^{z_j}}, \dots, \frac{e^{z_p}}{\sum_j e^{z_j}} \right]^T$$

Softmax

Nonnegativity and summed to 1

DL with nontrivial constraints

- Robustness evaluation
- Imbalanced learning
- Topology optimization
- Contrastive learning

More examples: A1

Deep Learning with Nontrivial Constraints: Methods and Applications

Chuan He¹, Ryan Devera¹, Wenjie Zhang¹, Ying Cui², Zhaosong Lu³ and Ju Sun¹

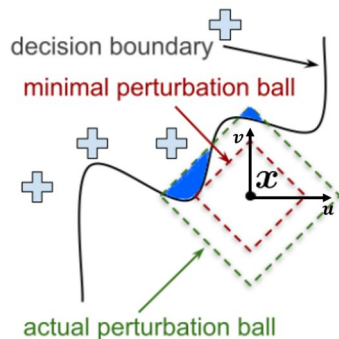
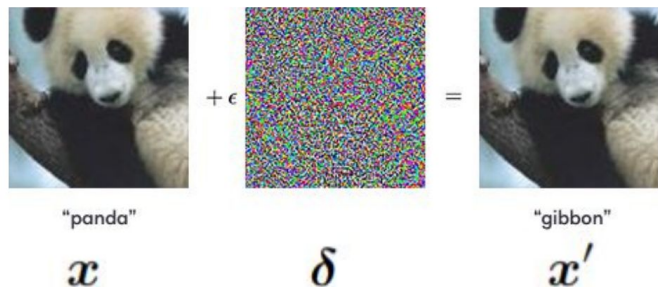
¹Computer Science and Engineering, University of Minnesota

²Industrial Engineering and Operations Research, University of California, Berkeley

³Industrial and Systems Engineering, University of Minnesota

{he000233, dever120, zhan7867}@umn.edu, yingcui@berkeley.edu, {zhaosong, jusun}@umn.edu

Robustness evaluation (RE)



Maximize loss function

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

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

Allowable perturbation

Valid image

Minimize robustness radius

$$\min_{x'} d(x, x')$$

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

Change the predicted class

Valid image

Projected gradient descent (PGD) for RE

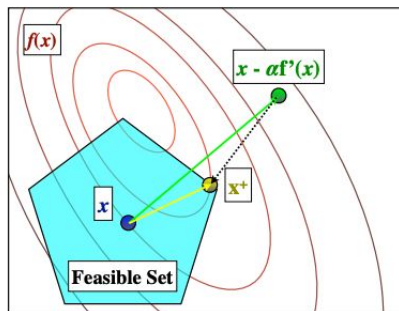
$$\min_{\mathbf{x} \in \mathcal{Q}} f(\mathbf{x})$$

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

Step size

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

Projection operator



Key hyperparameters:

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

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

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

Algorithm 1 APGD

```

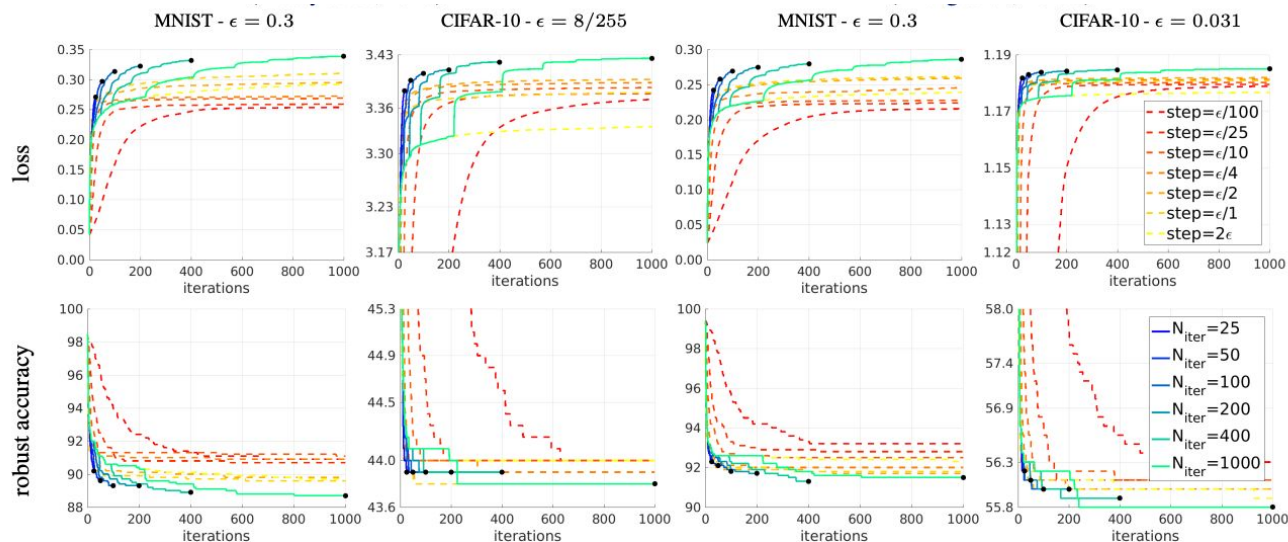
1: Input:  $f, S, x^{(0)}, \eta, N_{\text{iter}}, W = \{w_0, \dots, w_n\}$ 
2: Output:  $x_{\max}, f_{\max}$ 
3:  $x^{(1)} \leftarrow P_S(x^{(0)} + \eta \nabla f(x^{(0)}))$ 
4:  $f_{\max} \leftarrow \max\{f(x^{(0)}), f(x^{(1)})\}$ 
5:  $x_{\max} \leftarrow x^{(0)}$  if  $f_{\max} \equiv f(x^{(0)})$  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:    $x^{(k+1)} \leftarrow P_S\left(x^{(k)} + \alpha(z^{(k+1)} - x^{(k)})\right. \\ \left. + (1 - \alpha)(x^{(k)} - x^{(k-1)})\right)$ 
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

```

Ref https://angms.science/doc/CVX/CVX_PGD.pdf
<https://www.cs.ubc.ca/~schmidtm/Courses/5XX-S20/S5.pdf>

Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. Croce, F., Hein, M., ICML 2020
<https://arxiv.org/pdf/2003.01690.pdf>

Problem with PGD



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

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

RE: penalty methods for complicated d

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

where $d(\mathbf{x}, \mathbf{x}') \doteq \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|_2$ **perceptual distance**
where $\phi(\mathbf{x}) \doteq [\hat{g}_1(\mathbf{x}), \dots, \hat{g}_L(\mathbf{x})]$

Projection onto the constraint is complicated

Penalty methods

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

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

LPA, Fast-LPA: penalty methods

PPGD: Projected gradient descent

Penalty methods tend to encounter

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

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

$$\begin{aligned} & d(\mathbf{x}, \mathbf{x}') \doteq \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|_2 \\ \text{where } & \phi(\mathbf{x}) \doteq [\hat{g}_1(\mathbf{x}), \dots, \hat{g}_L(\mathbf{x})] \end{aligned}$$

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

**Unreliable optimization ==
Unreliable RE**

Outline

Constrained deep learning: CDL

- What, how, and why for CDL
- **No good solvers for CDL yet**
- Granso and PyGranso
- PyGranso in action
- Outlook

DL frameworks



JAX: Autograd and XLA



For unconstrained DL problems

Convex optimization solvers and frameworks



Modeling languages



GUROBI
OPTIMIZATION

**SDPT³ - a M_{ATLAB} software package for
semidefinite-quadratic-linear programming**

[K. C. Toh](#), [R. H. Tütüncü](#), and [M. J. Todd](#).

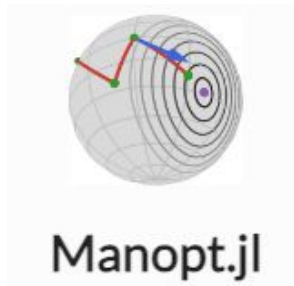
TFOCS: Templates for First-Order Conic Solvers

Solvers

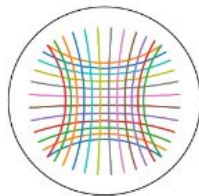
Not for DL, which involves NCVX optimization

Note: Gurobi can handle certain NCVX problems

Manifold optimization



Geomstats



$\mathcal{T}_p \mathcal{G}$
geoopt

McTorch Lib, a manifold optimization library for deep learning

Only for **differentiable manifolds constraints**

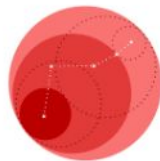
General constrained optimization

KNITRO[®]



IPOPT

Interior-point methods



ensmallen

flexible C++ library for efficient numerical optimization

GENO

Augmented Lagrangian methods



Cooper

TensorFlow Constrained Optimization (TFCO)

Lagrangian-method-based constrained optimization

Specialized ML packages



Problem-specific solvers that **cannot be easily extended** to new formulations

Outline

Constrained deep learning: CDL

- What, how, and why for CDL
- No good solvers for CDL yet
- **Granso and PyGranso**
- PyGranso in action
- Outlook

Issues with typical CDL methods

projected gradient descent

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

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

Issue: no principled stopping criterion/step size rules

Lagrangian method

$$\min_x \max_{\lambda \geq 0} \overset{\blacktriangledown}{f}(x) + \lambda^\top g(x)$$

Idea: alternating minimize x and maximize λ via gradient descent

penalty methods

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

$$\min_x f(x) + \lambda \max(0, g(x))$$

Solved with increasing λ : sequence

Issue: infeasible solution

Issues

- Infeasible solution
- Slow convergence

Want

- **Feasible & stationary solution**
- **Reasonable speed**

Principled answers to these questions

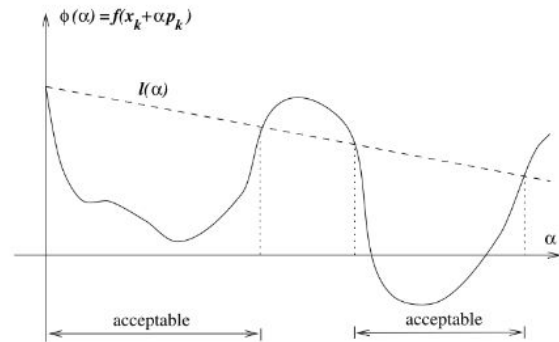
- Feasible & stationary solution

Stationarity and feasibility check: KKT condition

- Reasonable speed

Line search

- A hidden problem: nonsmoothness



Armijo (Sufficient Decrease) Condition

Key algorithm



<http://www.timmitchell.com/software/GRANSO/>

Nonconvex, nonsmooth, constrained

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

Penalty sequential quadratic programming (P-SQP)

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

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.

Algorithm highlights

Steering strategy for the penalty parameter

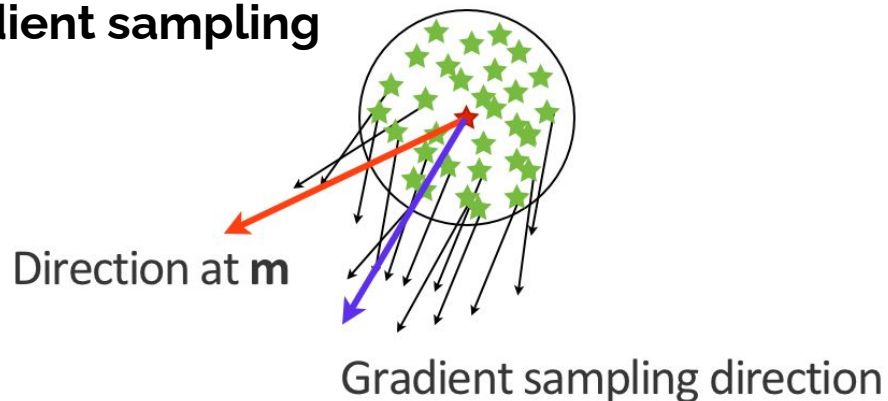
If feasibility improvement is insufficient : $l_\delta(d_k; x_k) < c_v v(x_k)$

Stationarity based on (approximate) gradient sampling

$$G_k := [\nabla f(x^k) \quad \nabla f(x^{k,1}) \quad \dots \quad \nabla f(x^{k,m})]$$

$$\min_{\lambda \in \mathbb{R}^{m+1}} \quad \frac{1}{2} \|G_k \lambda\|_2^2$$

$$\text{s.t. } \mathbf{1}^T \lambda = 1, \quad \lambda \geq 0$$



Key take-away



- Principled stopping criterion and line search, to obtain a **solution with certificate** (stationarity & feasibility check)
- Quasi-newton style method for fast convergence, i.e., **reasonable speed and high-precision solution**

Limitations of GRANSO



```
% Gradient of inner product with respect to A
f_grad      = imag((conj(Bty)*Cx.)/(y'*x));
f_grad      = f_grad(:);

% Gradient of inner product with respect to A
ci_grad     = real((conj(Bty)*Cx.)/(y'*x));
ci_grad     = ci_grad(:);
```

analytical gradients required

```
p      = size(B,2);
m      = size(C,1);
X      = reshape(x,p,m);
```

vector variables only

Lack of Auto-Differentiation

Lack of GPU Support

No native support of tensor variables

⇒ impossible to do deep learning with GRANSO

GRANSO meets PyTorch

GRANSO + PyTorch



NCVX PyGRANSO
Documentation

Search the docs ...

Introduction

Installation

Settings

Examples



Home



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

NCVX Package

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

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

Outline

Constrained deep learning: CDL

- What and how for CDL
- Why CDL
- No good solvers for CDL yet
- Granso and PyGranso
- **PyGranso in action**
- Outlook

Example 1: Support Vector Machine (SVM)

Soft-margin SVM

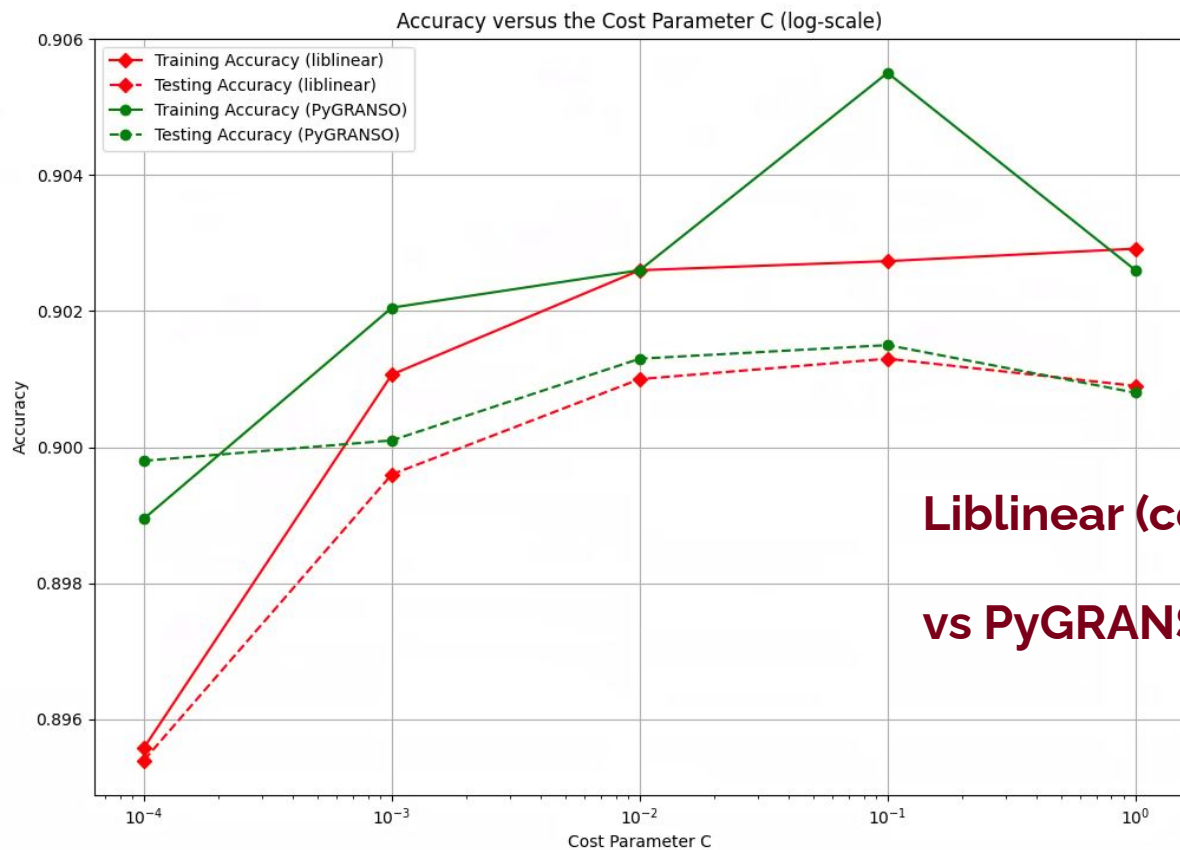
$$\min_{\mathbf{w}, b, \zeta} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \zeta_i$$

s.t. $y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \zeta_i, \zeta_i \geq 0 \quad \forall i = 1, \dots, n$

```
def comb_fn(X_struct):
    # obtain optimization variables
    w = X_struct.w
    b = X_struct.b
    zeta = X_struct.zeta
    # objective function
    f = 0.5*w.T@w + C*torch.sum(zeta)
    # inequality constraints
    ci = pygransoStruct()
    ci.c1 = 1 - zeta - y*(x@w+b)
    ci.c2 = -zeta
    # equality constraint
    ce = None
    return [f,ci,ce]

# specify optimization variables
var_in = {"w": [d,1], "b": [1,1], "zeta": [n,1]}
# pygranso main algorithm
soln = pygranso(var_in,comb_fn)
```

Binary classification (odd vs even digits) on MNIST dataset



**Liblinear (coordinate descent)
vs PyGRANSO**

Example 2: Robustness eval. —min formulation

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

```
def comb_fn(X_struct):
    # obtain optimization variables
    x_prime = X_struct.x_prime
    # objective function
    f = d(x, x_prime)
    # inequality constraints
    ci = pygransoStruct()
    f_theta_all = f_theta(x_prime)
    fy = f_theta_all[:, y] # true class output
    # output except true class
    fi = torch.hstack((f_theta_all[:, :, y], f_theta_all[:, :, y+1:]))
    ci.c1 = fy - torch.max(fi)
    ci.c2 = -x_prime
    ci.c3 = x_prime-1
    # equality constraint
    ce = None
    return [f, ci, ce]
# specify optimization variable (tensor)
var_in = {"x_prime": list(x.shape)}
# pygranso main algorithm
soln = pygranso(var_in, comb_fn)
```

CIFAR10 dataset

Compared with FAB [iterative constraint
linearization + projected gradient]

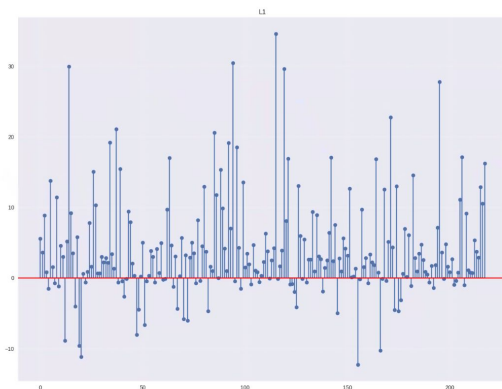
<https://github.com/fra31/auto-attack>

$$\min_{x'} d(x, x')$$

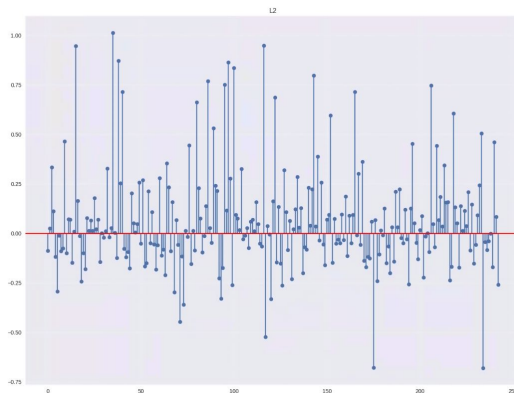
$$\text{s. t. } \max_{\ell \neq c} f_{\theta}^{\ell}(x') \geq f_{\theta}^c(x')$$

$$x' \in [0, 1]^n$$

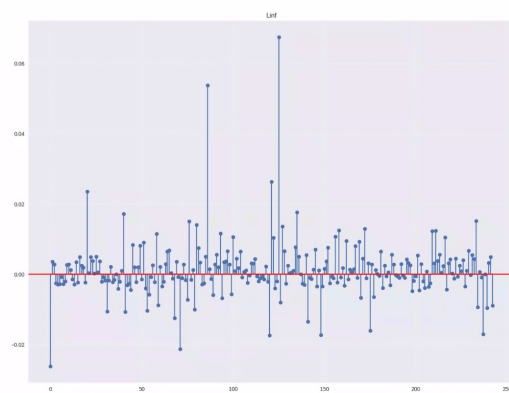
X-axis: image index; Y-axis: PyGRANSO radius - FAB radius



L1 attack



L2 attack



Linf attack

<https://ncvx.org/>

Many
others

NCVX PyGRANSO
Documentation

Search the docs ...

Introduction

Installation

Settings

Examples

- Rosenbrock
- Eigenvalue Optimization
- Dictionary Learning
- Nonlinear Feasibility Problem
- Sphere Manifold
- Trace Optimization
- Robust PCA
- Generalized LASSO
- Logistic Regression
- LeNet5
- Perceptual Attack
- Orthogonal RNN

Highlights



Home



NCVX Package

NCVX (NonConVeX) is a user-friendly and scalable python software package targeting general nonsmooth NCVX problems with nonsmooth constraints. **NCVX** is being developed by **GLOVEX** at the Department of Computer Science & Engineering, University of Minnesota, Twin Cities.

The initial release of **NCVX** contains the solver **PyGRANSO**, a PyTorch-enabled port of **GRANSO** incorporating auto-differentiation, GPU acceleration, tensor input, and support for new QP solvers. As a highlight, **PyGRANSO** can solve general constrained deep learning problems, the first of its kind.



Our current developments

- PyGranso: deterministic SQP solver
- More scalable deterministic solvers (exact penalty, single-loop ALM, etc)
- Algorithm foundations for stochastic objective with stochastic constraints/numerous constraints
- Scalable solvers for stochastic objective with stochastic constraints/numerous constraints

Closing



Deep Learning with Nontrivial Constraints: Methods and Applications

Chuan He¹, Ryan Devera¹, Wenjie Zhang¹, Ying Cui², Zhaosong Lu³ and Ju Sun¹

¹Computer Science and Engineering, University of Minnesota

²Industrial Engineering and Operations Research, University of California, Berkeley

³Industrial and Systems Engineering, University of Minnesota

{he000233, dever120, zhan7867}@umn.edu, yingcui@berkeley.edu, {zhaosong, jusun}@umn.edu



NCVX PyGRANSO
Documentation

Search the docs ...

Introduction

Installation

Settings

Examples

Home



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

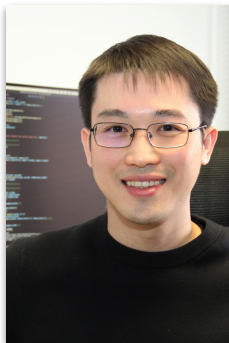
Thanks to all contributors



Prof. Tim Mitchell
(CS, Queens Col.)



Prof. Ying Cui
(IEOR, UC
Berkeley)



Prof. Qizhi He
(CSGE, UMN)



Prof.
Zhaosong Lu
(CSGE, UMN)



Dr. Chuan He
(CS&E, UMN)

Thanks to all contributors



Buyun Liang
(CSE U Penn;
Formally, CS&E, UMN)



Ryan de Vera
(CS&E, UMN)



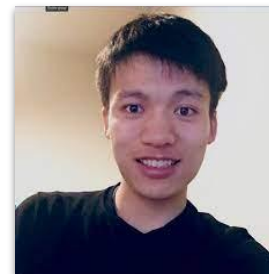
Hengyue Liang
(ECE, UMN)



Wenjie Zhang
(CS&E, UMN)



Yash Travadi
(CS&E, UMN)



Le Peng
(CS&E, UMN)