# Deep Learning with Nontrivial Constraints

**Buyun Liang, Ryan de Vera, and Ju Sun**
Computer Science & Engineering, University of Minnesota
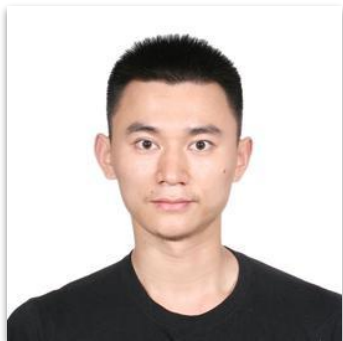**Tim Mitchell**
Queens College, City University of New York

Apr 29, 2023

# Presenters



**Buyun Liang**
**(CS&E, UMN)**



**Ryan de Vera**
**(CS&E, UMN)**



**Prof. Tim Mitchell**
**(CS, Queens C.)**



**Prof. Ju Sun**
**(CS&E, UMN)**

# Contributors

**Hengyue Liang**
**(ECE, UMN)**

**Wenjie Zhang**
**(CS&E, UMN)**

**Yash Travadi**
**(CS&E, UMN)**

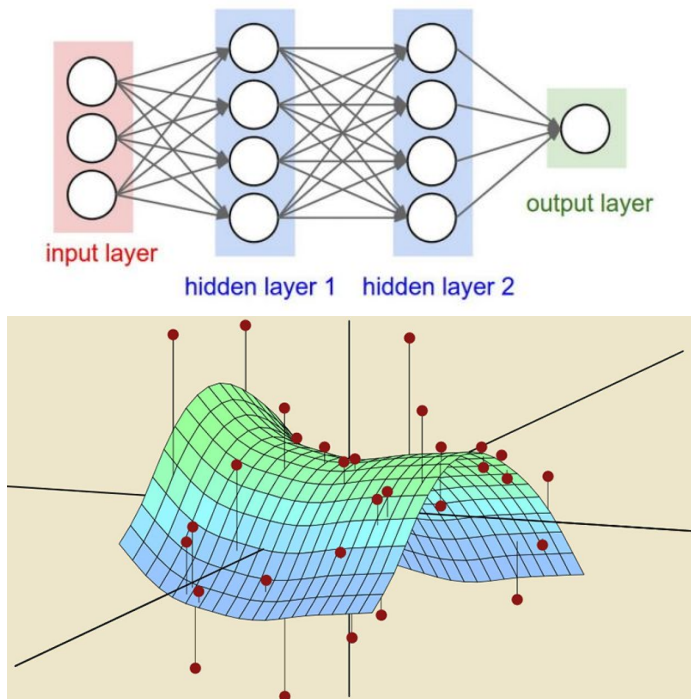**Le Peng**
**(CS&E, UMN)**

**Prof. Ying Cui**
**(ISyE, UMN)**

**Prof. Qizhi He**
**(CSGE, UMN)**

# Deep learning (DL)

## Artificial neural networks



**used to approximate nonlinear functions**

**Typical supervised learning pipeline**

| Step | General view | NN view |
|------|-------------|---------|
| 1 | Gather training set $(\boldsymbol{x}_1, \boldsymbol{y}_1), \ldots, (\boldsymbol{x}_n, \boldsymbol{y}_n)$ | Gather training set $(\boldsymbol{x}_1, \boldsymbol{y}_1), \ldots, (\boldsymbol{x}_n, \boldsymbol{y}_n)$ |
| 2 | Choose a family of functions, e.g., $\mathcal{H}$, so that there is an $f \in \mathcal{H}$ to ensure $\boldsymbol{y}_i \approx f(\boldsymbol{x}_i), \forall i$ | Choose a NN with $k$ neurons, so that there is a group of weights $(\boldsymbol{w}_1, \ldots, \boldsymbol{w}_k, b_1, \ldots, b_k)$ ensuring $\boldsymbol{y}_i \approx \{\mathrm{NN}(\boldsymbol{w}_1, \ldots, \boldsymbol{w}_k, b_1, \ldots, b_k)\}(\boldsymbol{x}_i), \forall i$ |
| 3 | Set up a loss function $\ell$ | Set up a loss function $\ell$ |
| 4 | Find an $f \in \mathcal{H}$ to minimize the average loss $$\frac{1}{n} \sum_{i=1}^{n} \ell\left(\boldsymbol{y}_i, f(\boldsymbol{x}_i)\right)$$ | Find weights $(\boldsymbol{w}_1, \ldots, \boldsymbol{w}_k, b_1, \ldots, b_k)$ to minimize the average loss $$\frac{1}{n} \sum_{i=1}^{n} \ell\left[\boldsymbol{y}_i, \{\mathrm{NN}(\boldsymbol{w}_1, \ldots, \boldsymbol{w}_k, b_1, \ldots, b_k)\}(\boldsymbol{x}_i)\right]$$ |

# Three fundamental questions in DL



– **Approximation**: is it powerful, i.e., the $\mathcal{H}$ large enough for all possible weights? | **Universal approximation theorems** |

– **Optimization**: how to solve

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

– **Generalization**: does the learned NN work well on "similar" data?

# Isn't solved?

## Base class

CLASS   torch.optim.Optimizer(*params*, *defaults*)   [SO

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 t
  Tensors should be optimized.
- **defaults** – (dict): a dict containing d
  when a parameter group doesn't spe

## Algorithms

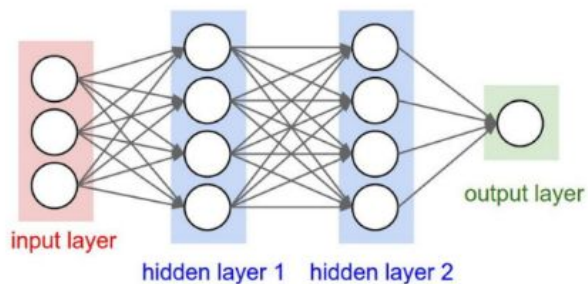| | |
|---|---|
| Adadelta | Implements Adadelta algorithm. |
| Adagrad | Implements Adagrad algorithm. |
| Adamax | Implements Adamax algorithm (a variant of Adam based on infinity norm). |
| ASGD | Implements Averaged Stochastic Gradient Descent. |
| LBFGS | Implements L-BFGS algorithm, heavily inspired by minFunc. |
| NAdam | Implements NAdam algorithm. |
| RAdam | Implements RAdam algorithm. |

# When DL meets constraints



Artificial neural networks

input layer
hidden layer 1    hidden layer 2
output layer

used to approximate nonlinear functions

**Unconstrained optimization**

$$\min_{\boldsymbol{w}_i's, \boldsymbol{b}_i's} \frac{1}{n} \sum_{i=1}^{n} \ell\left[\boldsymbol{y}_i, \{\mathsf{NN}\left(\boldsymbol{w}_1, \ldots, \boldsymbol{w}_k, b_1, \ldots, b_k\right)\}\left(\boldsymbol{x}_i\right)\right]$$

$$\min_{\boldsymbol{x}} f(\boldsymbol{x})$$ **"Solved"**

**Constrained optimization**

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

**largely "unsolved"**

# This tutorial:

## Constrained optimization

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

largely "unsolved"

## how to solve DL problems with constraints





Left: "DL problems with constraints" in DALL-E's mind

# 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

# DL with simple constraints

**Embedding constraints into DL models**



**ReLU**
(Rectified Linear Unit)

**Sigmoid**

$$z \mapsto \left[ \frac{e^{z_1}}{\sum_j e^{z_j}}, \cdots, \frac{e^{z_p}}{\sum_j e^{z_j}} \right]^\top$$

**Softmax**

**Nonnegativity**

**[0, 1]**

**Nonnegativity and summed to 1**

# DL with nontrivial constraints

- **Robustness evaluation**
- Imbalanced learning
- Physics-informed neural networks (PINNs)

# Robustness evaluation (RE)



"panda" $x$   $+\epsilon$   $\delta$   $=$   "gibbon" $x'$

decision boundary

minimal perturbation ball

actual perturbation ball

Maximize loss function

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

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

Allowable perturbation    Valid image

Minimize robustness radius

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

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

Change the predicted class    Valid image

**Ref** Optimization and Optimizers for Adversarial Robustness. Liang, H., Liang, B., Peng, L., Cui, Y., Mitchell, T., & Sun, J. https://arxiv.org/abs/2303.13401

# Projected gradient descent (PGD) for RE

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

**Step size**

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

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

**Projection operator**



**Key hyperparameters:**
(1) step size
(2) iteration number

$$\max_{\boldsymbol{x}'} \ell\left(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x}')\right)$$
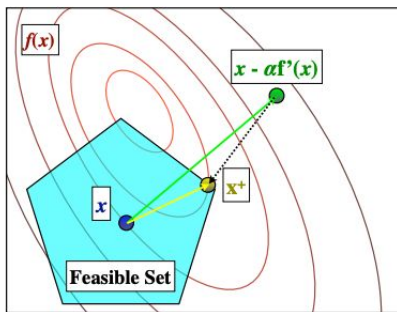
$$\text{s.t. } d\left(\boldsymbol{x}, \boldsymbol{x}'\right) \leq \varepsilon, \quad \boldsymbol{x}' \in [0,1]^n$$
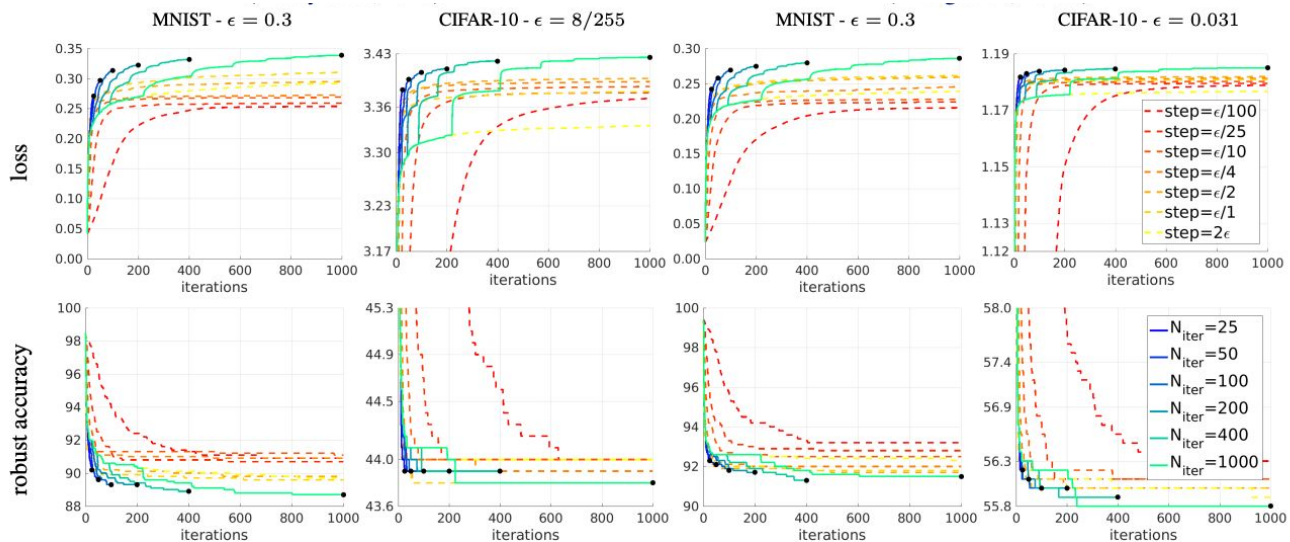
---

**Algorithm 1** APGD

1: **Input:** $f$, $S$, $x^{(0)}$, $\eta$, $N_{\text{iter}}$, $W = \{w_0, \ldots, w_n\}$
2: **Output:** $x_{\max}$, $f_{\max}$
3: $x^{(1)} \leftarrow P_S\left(x^{(0)} + \eta \nabla f(x^{(0)})\right)$
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: $\quad z^{(k+1)} \leftarrow P_S\left(x^{(k)} + \eta \nabla f(x^{(k)})\right)$
8: $\quad x^{(k+1)} \leftarrow P_S\left(x^{(k)} + \alpha(z^{(k+1)} - x^{(k)})\right.$
$\qquad\qquad\qquad \left. +(1-\alpha)(x^{(k)} - x^{(k-1)})\right)$
9: $\quad$ **if** $f(x^{(k+1)}) > f_{\max}$ **then**
10: $\qquad x_{\max} \leftarrow x^{(k+1)}$ and $f_{\max} \leftarrow f(x^{(k+1)})$
11: $\quad$ **end if**
12: $\quad$ **if** $k \in W$ **then**
13: $\qquad$ **if** Condition 1 **or** Condition 2 **then**
14: $\qquad\quad \eta \leftarrow \eta/2$ and $x^{(k+1)} \leftarrow x_{\max}$
15: $\qquad$ **end if**
16: $\quad$ **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 projected gradient descent



MNIST - $\epsilon = 0.3$    CIFAR-10 - $\epsilon = 8/255$    MNIST - $\epsilon = 0.3$    CIFAR-10 - $\epsilon = 0.031$

step=$\epsilon/100$
step=$\epsilon/25$
step=$\epsilon/10$
step=$\epsilon/4$
step=$\epsilon/2$
step=$\epsilon/1$
step=$2\epsilon$

$N_{iter}$=25
$N_{iter}$=50
$N_{iter}$=100
$N_{iter}$=200
$N_{iter}$=400
$N_{iter}$=1000

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

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

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

**Ref** 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

# Robustness evaluation: penalty methods for complicated d (perceptual attack)

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

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

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

$$\text{where} \quad \phi(\boldsymbol{x}) \doteq [\,\widehat{g}_1(\boldsymbol{x}), \ldots, \widehat{g}_L(\boldsymbol{x})\,]$$

**perceptual distance**

**Projection onto the constraint is complicated**

**Penalty methods**

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

Solve it for each fixed $\lambda$ and then increase $\lambda$

---

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

---

**Ref** Perceptual adversarial robustness: Defense against unseen threat models. Laidlaw, C., Singla, S., & Feizi, S. https://arxiv.org/abs/2006.12655

# 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 |

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

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

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

$$\text{where} \quad \phi(\boldsymbol{x}) \doteq [\,\widehat{g}_1(\boldsymbol{x}), \ldots, \widehat{g}_L(\boldsymbol{x})\,]$$

**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**

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

**Ref** Optimization and Optimizers for Adversarial Robustness. Liang, H., Liang, B., Peng, L., Cui, Y., Mitchell, T., & Sun, J. arXiv preprint arXiv:2303.13401.

# Robustness evaluation: quick summary

**Two forms of RE**

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

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

**Two methods for handling constraints**

**projected gradient descent**

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

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

**penalty methods**

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

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

Solved with increasing $\lambda$ sequence

**Issue: no principled stopping criterion/step size rules**

**Issue: infeasible solution**

# DL with nontrivial constraints

- Robustness evaluation
- **Imbalanced learning**
- Physics-informed neural networks (PINNs)

# Imbalanced learning:  background



**Class imbalance in healthcare datasets**

|  | Predicted POS | Predicted NEG |
|---|---|---|
| POS | 70 | 30 |
| NEG | 1000 | 9000 |

**Accuracy:** 9070/10100 = 0.898
**True Positive Rate (Sensitivity, Recall):** 0.7
**True Negative Rate (Specificity):** 0.9
**Balanced Accuracy:** (0.7 + 0.9)/2 = 0.80
**Precision (POS):** 70/1070 = 0.065
**F1 Score:** 2*0.065*0.7/(0.065 + 0.7) = 0.119

Reliable evaluation in  imbalanced learning: **precision needed**!

# Imbalanced learning: direct metric optimization

**Typical learning objective:** $\quad \min\limits_{f \in \mathcal{H}} \mathbb{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}_{\boldsymbol{x},\boldsymbol{y}}} \mathbb{1}\{\boldsymbol{y} \neq f(\boldsymbol{x})\} \qquad$ **accuracy maximization**

$$\text{precision}(f_{\boldsymbol{\theta}}, t) = \frac{\sum_{i=1}^{N} \mathbb{1}\{y_i = +1\}\,\mathbb{1}\{f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) > t\}}{\sum_{i=1}^{N} \mathbb{1}\{f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) > t\}} \qquad \text{recall}(f_{\boldsymbol{\theta}}, t) = \frac{\sum_{i=1}^{N} \mathbb{1}\{y_i = +1\}\,\mathbb{1}\{f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) > t\}}{\sum_{i=1}^{N} \mathbb{1}\{y_i = +1\}}$$

$$F_{\beta}(f_{\boldsymbol{\theta}}, t) = (1 + \beta^2) \frac{\text{precision}(f_{\boldsymbol{\theta}}, t) \cdot \text{recall}(f_{\boldsymbol{\theta}}, t)}{\beta^2 \text{precision}(f_{\boldsymbol{\theta}}, t) + \text{recall}(f_{\boldsymbol{\theta}}, t)}$$

$$\text{AP}(f_{\boldsymbol{\theta}}) = \frac{1}{|\{i : y_i = +1\}|} \sum_{i=1}^{N} \mathbb{1}\{y_i = +1\} \frac{\sum_{s=1}^{N} \mathbb{1}\{y_s = +1\}\,\mathbb{1}\{f_{\boldsymbol{\theta}}(\boldsymbol{x}_s) > f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)\}}{\sum_{s=1}^{N} \mathbb{1}\{f_{\boldsymbol{\theta}}(\boldsymbol{x}_s) > f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)\}}.$$

**fix precision, optimize recall (FPOR):** $\quad \max\limits_{\boldsymbol{\theta}, t} \text{recall}(f_{\boldsymbol{\theta}}, t) \quad$ s.t. $\text{precision}(f_{\boldsymbol{\theta}}, t) \geq \alpha,$

**fix recall, optimize precision (FROP):** $\quad \max\limits_{\boldsymbol{\theta}, t} \text{precision}_t \quad$ s.t. $\text{recall}(f_{\boldsymbol{\theta}}, t) \geq \alpha,$

**optimize $F_{\beta}$ score (OFBS):** $\quad \max\limits_{\boldsymbol{\theta}, t} F_{\beta}(f_{\boldsymbol{\theta}}, t),$

**optimize AP (OAP):** $\quad \max\limits_{\boldsymbol{\theta}} \text{AP}(f_{\boldsymbol{\theta}}).$

# Imbalanced learning: Lagrangian methods

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

$$\min_{\boldsymbol{x}} \max_{\boldsymbol{\lambda} \geq \boldsymbol{0}} \ f(\boldsymbol{x}) + \boldsymbol{\lambda}^{\mathsf{T}} g(\boldsymbol{x})$$

**Idea: alternating minimize $\boldsymbol{x}$ and maximize $\boldsymbol{\lambda}$ via gradient descent**

Reminder on gradient descent

$$\min_{\boldsymbol{x}} \ f(\boldsymbol{x})$$

iteration step :

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - t\nabla f(\boldsymbol{x}_k)$$

$$\max_{\boldsymbol{\theta},t} \ \text{recall}(f_{\boldsymbol{\theta}}, t) \quad \text{s.\,t.} \ \text{precision}(f_{\boldsymbol{\theta}}, t) \geq \alpha,$$

$$\max_{f,b} \ \frac{1}{|Y^+|} tp(f)$$

$$\text{s.t.} \quad tp(f) \geq \alpha(tp(f) + fp(f)).$$

$$f^{(t+1)} = f^{(t)} - \gamma \nabla L(f^{(t)}, \lambda^{(t)})$$
$$\lambda^{(t+1)} = \lambda^{(t)} + \gamma \nabla L(f^{(t+1)}, \lambda^{(t)})$$

where

$$L(f, \lambda) = (1 + \lambda)\mathscr{L}^+(f) + \lambda \frac{\alpha}{1 - \alpha} \mathscr{L}^-(f) - \lambda |Y^+|.$$

Eban, Elad, et al. "Scalable learning of non-decomposable objectives." *Artificial intelligence and statistics*. PMLR, 2017.

# Imbalanced learning: quick summary

**fix precision, optimize recall (FPOR):** $\max\limits_{\boldsymbol{\theta},t} \mathrm{recall}(f_{\boldsymbol{\theta}}, t)$  s.t. $\mathrm{precision}(f_{\boldsymbol{\theta}}, t) \geq \alpha,$

**fix recall, optimize precision (FROP):** $\max\limits_{\boldsymbol{\theta},t} \mathrm{precision}_t$  s.t. $\mathrm{recall}(f_{\boldsymbol{\theta}}, t) \geq \alpha,$

**Lagrangian method**

$$\min_{\boldsymbol{x}} \max_{\boldsymbol{\lambda} \geq \mathbf{0}} \ \check{f}(\boldsymbol{x}) + \boldsymbol{\lambda}^{\mathsf{T}} g(\boldsymbol{x})$$

**Issues**

- Infeasible solution
- Slow convergence

**Idea: alternating minimize $\boldsymbol{x}$ and maximize $\lambda$ via gradient descent**

# DL with nontrivial constraints

- Robustness evaluation
- Imbalanced learning
- **Physics-informed neural networks (PINNs)**

# PINNs: DL for PDEs

## Physics-informed neural networks (PINNs)

$$f\left(\mathbf{x}; \frac{\partial u}{\partial x_1}, \ldots, \frac{\partial u}{\partial x_d}; \frac{\partial^2 u}{\partial x_1 \partial x_1}, \ldots, \frac{\partial^2 u}{\partial x_1 \partial x_d}; \ldots; \boldsymbol{\lambda}\right) = 0, \quad \mathbf{x} \in \Omega,$$

$$\mathcal{B}(u, \mathbf{x}) = 0 \quad \text{on} \quad \partial\Omega,$$

U is represented as a DNN



**Continuous modeling instead of finite-difference for derivatives**

Penalty parameters

$$\mathcal{L}(\boldsymbol{\theta}; \mathcal{T}) = \boxed{w_f} \mathcal{L}_f(\boldsymbol{\theta}; \mathcal{T}_f) + \boxed{w_b} \mathcal{L}_b(\boldsymbol{\theta}; \mathcal{T}_b)$$

$$\mathcal{L}_f(\boldsymbol{\theta}; \mathcal{T}_f) = \frac{1}{|\mathcal{T}_f|} \sum_{\mathbf{x} \in \mathcal{T}_f} \left\| f\left(\mathbf{x}; \frac{\partial \hat{u}}{\partial x_1}, \ldots, \frac{\partial \hat{u}}{\partial x_d}; \frac{\partial^2 \hat{u}}{\partial x_1 \partial x_1}, \ldots, \frac{\partial^2 \hat{u}}{\partial x_1 \partial x_d}; \ldots; \boldsymbol{\lambda}\right) \right\|_2^2$$

$$\mathcal{L}_b(\boldsymbol{\theta}; \mathcal{T}_b) = \frac{1}{|\mathcal{T}_b|} \sum_{\mathbf{x} \in \mathcal{T}_b} \|\mathcal{B}(\hat{u}, \mathbf{x})\|_2^2,$$

**Ref** Liang, Buyun, Tim Mitchell, and Ju Sun. "NCVX: A general-purpose optimization solver for constrained machine and deep learning." *arXiv preprint arXiv:2210.00973* (2022). wiki https://en.wikipedia.org/wiki/Physics-informed_neural_networks

# PINNs: methods

Typical methods

$$\min_{u(\boldsymbol{x})} \mathcal{L}(u(\boldsymbol{x})) \quad \text{s. t.} \begin{cases} f\left(\boldsymbol{x}; \frac{\partial u}{\partial x_1}, \ldots, \frac{\partial u}{\partial x_d}; \frac{\partial^2 u}{\partial x_1 \partial x_1}, \ldots, \frac{\partial^2 u}{\partial x_1 \partial x_d}; \ldots\right) = 0, \quad \forall \boldsymbol{x} \in \Omega \\ \mathcal{B}(u, \boldsymbol{x}) = 0, \quad \forall \boldsymbol{x} \in \partial\Omega \end{cases}$$

- Penalty methods
- Lagrangian methods
- Augmented Lagrangian methods

⟹ Infeasible solution

**+**

- First-order solver

⟹ Low quality solution

# 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

# There's no free lunch!

**Supervised learning as data fitting**



**Typically, #data points we need grow exponentially with respect to dimension (i.e., curse of dimensionality)**

**Knowledge**

**Small data AI**



**Data**

**Building in prior knowledge is _crucial_ for reducing the data complexity e.g., "convolutional" layers**

# AI for science



**Thrust B: How Should Domain Knowledge Be Incorporated into Supervised Machine Learning?**

The central question for this thrust is "which knowledge should be leveraged in SciML, and how should this knowledge be included?" Any answers will naturally depend on the SciML task and computational budgets, thus mirroring standard considerations in traditional scientific computing.

**Hard Constraints.** One research avenue involves incorporation of domain knowledge through imposition of constraints that cannot be violated. These hard constraints could be enforced during training, replacing what typically is an unconstrained optimization problem with a constrained one. In general, such constraints could involve simulations or highly nonlinear functions of the training parameters. Therefore, there is a need to identify particular cases when constraint qualification conditions can be ensured as these conditions are necessary regularity conditions for constrained optimization [57–59]. Although incorporating constraints during training generally makes maximal use of training data, there may be additional opportunities to employ constraints at the time of prediction (e.g., by projecting predictions onto the region induced by the constraints).

**Soft Constraints.** A similar avenue for incorporating domain knowledge involves modifying the objective function (soft constraints) used in training. It is understood that ML loss function selection should be guided by the task and data. Therefore, opportunities exist for developing loss functions that incorporate domain knowledge and analyzing the resulting impact on solvability

**BASIC RESEARCH NEEDS FOR**
**Scientific Machine Learning**
Core Technologies for Artificial Intelligence

Prepared for U.S. Department of Energy Advanced Scientific Computing Research

U.S. DEPARTMENT OF ENERGY

**Ref** https://www.osti.gov/servlets/purl/1478744

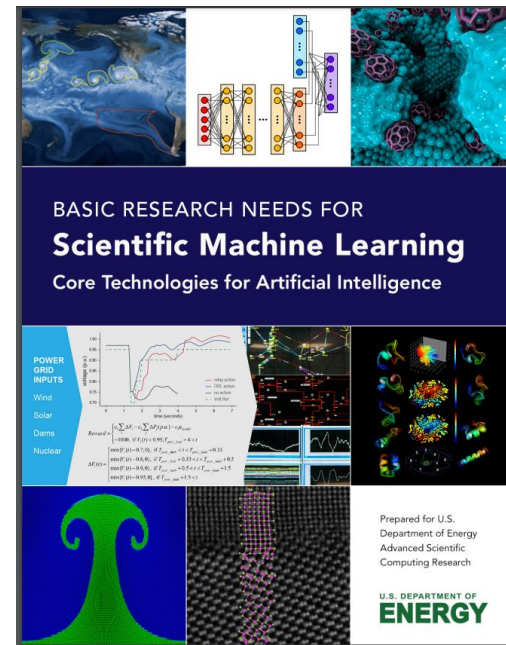**Domain-Aware Scientific Machine Learning**

# AI for science



**Thrust C: How Should the Robustness, Performance, and Quality of Scientific Machine Learning Be Assessed?**

The outcome of an ML process is either a decision (classification) or a prediction. For reliable and credible use of SciML, we need the ability to rigorously quantify ML performance in these outcomes. Performance measurement implies an assessment of quality, as well as a cost measure of computations and/or data preparation and management. Traditional measures of acceptable quality based on statistical cross-validation-type approaches often are heuristic. Measures of prediction quality such as *a priori* and *a posteriori* error estimates for numerical approximations of PDEs [96] (familiar to the finite element modeling community) will be transformative in allowing the development of optimal and reliable ML algorithms for different uses. Such error estimates also will enable SciML processes that allow iterative model improvement. Research establishing quantitative estimates of prediction quality, including effective confidence bounds, will greatly enhance the usefulness of SciML to decision makers and users. Finally, research is needed on algorithms that have proven convergence rates with weak dependence on bad data, especially in situations with a large amount of data of unproven quality or minimal availability of human expertise.

**BASIC RESEARCH NEEDS FOR Scientific Machine Learning**
Core Technologies for Artificial Intelligence

Prepared for U.S. Department of Energy Advanced Scientific Computing Research

U.S. DEPARTMENT OF ENERGY

**Ref** https://www.osti.gov/servlets/purl/1478744

**Robust Scientific Machine Learning**

# 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

# DL frameworks

JAX: Autograd and XLA

For unconstrained DL problems

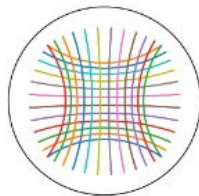# Convex optimization solvers and frameworks



**Modeling languages**

**Solvers**

**Not for DL**, which involves NCVX optimization

Note: Gurobi can handle certain NCVX problems

# Manifold optimization



Manopt.jl



Geomstats



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

geoopt

McTorch Lib, a manifold optimization library for deep learning
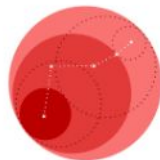
Only for **differentiable manifolds constraints**

# General constrained optimization



**IPOPT**

**Interior-point methods**

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

# Issues with typical CDL methods

## projected gradient descent

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

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

**Issue: no principled stopping criterion/step size rules**

## penalty methods

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

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

Solved with increasing $\lambda$ sequence

**Issue: infeasible solution**

## Lagrangian method

$$\min_{\boldsymbol{x}} \max_{\boldsymbol{\lambda} \ge \mathbf{0}} \ f(\boldsymbol{x}) + \boldsymbol{\lambda}^\mathsf{T} g(\boldsymbol{x})$$

**Idea: alternating minimize $\boldsymbol{x}$ and maximize $\lambda$ via gradient descent**

**Issues**

- Infeasible solution
- Slow convergence

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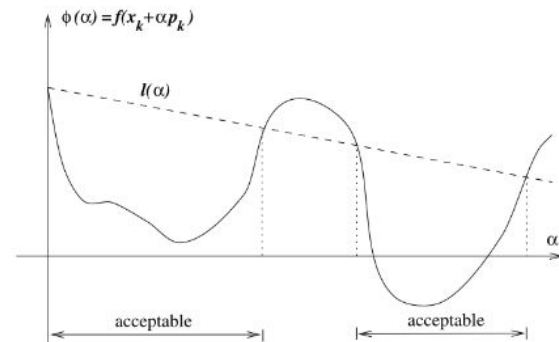- **Reasonable speed**

# Principled answers to these questions

- **Feasible & stationary solution**

  **Stationarity and feasibility check: KKT condition**

- **Reasonable speed**

  **Line search**

- **A hidden problem: nonsmoothness**



$\phi(\alpha) = f(x_k + \alpha p_k)$

$l(\alpha)$

$\alpha$

acceptable        acceptable

Armijo (Sufficient Decrease) Condition

# A principled solver for constrained, nonconvex, nonsmooth problems



**Nonconvex, nonsmooth, constrained**

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

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

$$\min_{d \in \mathbb{R}^n, \, s \in \mathbb{R}^p} \quad \mu(f(x_k) + \nabla f(x_k)^\mathsf{T} d) + e^\mathsf{T} s + \frac{1}{2} d^\mathsf{T} H_k d$$

$$\text{s.t.} \quad c(x_k) + \nabla c(x_k)^\mathsf{T} d \leq s, \quad s \geq 0,$$

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

**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.

# Determining the search direction

**Corresponding dual**

$$\max_{\lambda \in \mathbb{R}^p} \quad \mu f(x_k) + c(x_k)^\mathsf{T}\lambda - \frac{1}{2}(\mu \nabla f(x_k) + \nabla c(x_k)\lambda)^\mathsf{T} H_k^{-1}(\mu \nabla f(x_k) + \nabla c(x_k)\lambda)$$

$$\text{s.t.} \quad 0 \leq \lambda \leq e, \tag{8}$$

**Primal solution (recovered from dual solution): searching direction**

$$d_k = -H_k^{-1}(\mu \nabla f(x_k) + \nabla c(x_k)\lambda_k). \tag{9}$$

**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.

# Adjusting the penalty parameter

**Linear model of constraint violation**

$$l(d; x_k) := \| \max\{c(x_k) + \nabla c(x_k)^\mathsf{T} d, 0\}\|_1$$

**Corresponding reduction**

$$l_\delta(d; x_k) := l(0; x_k) - l(d; x_k)$$
$$= v(x_k) - \| \max\{c(x_k) + \nabla c(x_k)^\mathsf{T} d, 0\}\|_1$$

Advantage: feasibility guarantee

GRANSO

---

**Procedure 1** $[d_k, \mu_{\text{new}}] = $ **sqp_steering_strategy**$(x_k, H_k, \mu)$

**Input:**
    Current iterate $x_k$ and BFGS Hessian approximation $H_k$
    Current value of the penalty parameter $\mu$

**Constants:**
    Values $c_v \in (0, 1)$ and $c_\mu \in (0, 1)$

**Output:**
    Search direction $d_k$
    Penalty parameter $\mu_{\text{new}} \in (0, \mu]$

1: Solve QP (8) using $\mu_{\text{new}} := \mu$ to obtain search direction $d_k$ from (9)
2: **if** $l_\delta(d_k; x_k) < c_v v(x_k)$ **then**
3:     Solve (8) using $\mu = 0$ to obtain reference direction $\tilde{d}_k$ from (9)
4:     **while** $l_\delta(d_k; x_k) < c_v l_\delta(\tilde{d}_k; x_k)$ **do**
5:         $\mu_{\text{new}} := c_\mu \mu_{\text{new}}$
6:         Solve QP (8) using $\mu := \mu_{\text{new}}$ to obtain search direction $d_k$ from (9)
7:     **end while**
8: **end if**

---

**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.

# Estimating the stationarity

*GRANSO*

**Gradient from l most recent iterates**

$$G := [\nabla f(x_{k+1-l}) \cdots \nabla f(x_k)]$$

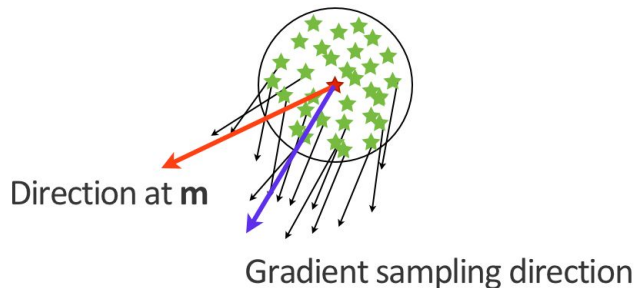$$J_i := [\nabla c_i(x_{k+1-l}) \cdots \nabla c_i(x_k)], \quad i \in \{1, \ldots, p\}$$

**Augmented QP**

$$\max_{\sigma \in \mathbb{R}^l, \lambda \in \mathbb{R}^{pl}} \quad \sum_{i=1}^{p} c_i(x_k) e^{\mathsf{T}} \lambda_i - \frac{1}{2} \begin{bmatrix} \sigma \\ \lambda \end{bmatrix}^{\mathsf{T}} [G, J_1, \ldots, J_p]^{\mathsf{T}} H_k^{-1} [G, J_1, \ldots, J_p] \begin{bmatrix} \sigma \\ \lambda \end{bmatrix}$$

$$\text{s.t.} \quad 0 \le \lambda_i \le e, \quad e^{\mathsf{T}} \sigma = \mu, \quad \sigma \ge 0. \tag{12}$$

**Primal solution: termination condition**

$$d_\diamond = H_k^{-1} [G, J_1, \ldots, J_p] \begin{bmatrix} \sigma \\ \lambda \end{bmatrix}$$

**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.

# Estimating the stationarity

**Augmented QP**

$$\max_{\sigma \in \mathbb{R}^l, \lambda \in \mathbb{R}^{pl}} \quad \sum_{i=1}^{p} c_i(x_k) e^{\mathsf{T}} \lambda_i - \frac{1}{2} \begin{bmatrix} \sigma \\ \lambda \end{bmatrix}^{\mathsf{T}} [G, J_1, \ldots, J_p]^{\mathsf{T}} H_k^{-1} [G, J_1, \ldots, J_p] \begin{bmatrix} \sigma \\ \lambda \end{bmatrix}$$

$$\text{s.t.} \quad 0 \le \lambda_i \le e, \quad e^{\mathsf{T}} \sigma = \mu, \quad \sigma \ge 0. \tag{12}$$

**Stationarity based on (approximate) gradient sampling**

$$G_k := \begin{bmatrix} \nabla f(x^k) & \nabla f(x^{k,1}) & \cdots & \nabla f(x^{k,m}) \end{bmatrix}$$

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

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

Direction at **m**

Gradient sampling direction

Advantage: can handle nonsmoothness

**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.

```
1:  Set H₀ := I and μ := μ₀
```
1: Set $H_0 := I$ and $\mu := \mu_0$
2: Set $\phi(\cdot)$ as the penalty function given in (2) using $f(\cdot)$ and $c(\cdot)$
3: Set $\nabla\phi(\cdot)$ and $v(\cdot)$ as the associated gradient (4) and violation function (3)
4: Evaluate $\phi_0 := \phi(x_0; \mu)$, $\nabla\phi_0 := \nabla\phi(x_0; \mu)$, and $v_0 := v(x_0)$
5: **for** $k = 0, 1, 2, \ldots$ **do**
6: $\quad [d_k, \hat{\mu}] := $ **sqp_steering_strategy**$(x_k, H_k, \mu)$
7: $\quad$ **if** $\hat{\mu} < \mu$ **then**
8: $\quad\quad$ // *Penalty parameter has been lowered by steering; update current iterate*
9: $\quad\quad$ Set $\mu := \hat{\mu}$
10: $\quad\quad$ Reevaluate $\phi_k := \phi(x_k; \mu)$, $\nabla\phi_k := \nabla\phi(x_k; \mu)$, and $v_k := v(x_k)$
11: $\quad$ **end if**
12: $\quad [x_{k+1}, \phi_{k+1}, \nabla\phi_{k+1}, v_{k+1}] := $ **inexact_linesearch**$(x_k, \phi_k, \nabla\phi_k, d_k, \phi(\cdot), \nabla\phi(\cdot))$
13: $\quad$ Compute $d_\diamond$ via (12) and (13)
14: $\quad$ **if** $\|d_\diamond\|_2 < \tau_\diamond$ and $v_{k+1} < \tau_v$ **then**
15: $\quad\quad$ // *Stationarity and feasibility sufficiently attained; terminate successfully*
16: $\quad\quad$ **break**
17: $\quad$ **end if**
18: $\quad$ Set $H_{k+1}$ using BFGS update formula
19: **end for**

### Advantages

- **Reliable step-size rule**

- **Principled stopping criterion**

**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.

# 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**

**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.

# 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

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

# NCVX PyGRANSO: Advantages

**1) Auto-Differentiation**

https://ncvx.org/

Orthogonal Dictionary Learning (ODL)

$$\min_{\boldsymbol{q}\in\mathbb{R}^n} f(\boldsymbol{q}) \doteq \frac{1}{m}\|\boldsymbol{q}^\mathsf{T}\boldsymbol{Y}\|_1, \quad \text{s.t. } \|\boldsymbol{q}\|_2 = 1$$

```
function[f,fg,ci,cig,ce,ceg]=comb_fn(q)
    f = 1/m*norm(q'*Y, 1); % obj
    fg = 1/m*Y*sign(Y'*q); % obj grad
    ci = []; cig = []; % no ineq constr
    ce = q'*q - 1; % eq constr
    ceg = 2*q; % eq constr grad
end
soln = granso(n,comb_fn);
```

Analytical gradients

```
def comb_fn(X_struct):
    q = X_struct.q
    f = 1/m*norm(q.T@Y, p=1) # obj
    ce = pygransoStruct()
    ce.c1 = q.T@q - 1 # eq constr
    return [f,None,ce]
var_in = {"q": [n,1]} # define variable
soln = pygranso(var_in,comb_fn)
```

No Analytical gradients

Demo 1: GRANSO for ODL

Demo 2: PyGRANSO for ODL

**Ref** Buyun Liang, Tim Mitchell, Ju Sun. NCVX: A General-Purpose Optimization Solver for Constrained Machine and Deep Learning. In Neural Information Processing Systems (NeurIPS) Workshop on Optimization for Machine Learning (OPT 2022).

# NCVX PyGRANSO: Advantages

## 2) GPU acceleration for large scale problems

https://ncvx.org/

Orthogonality-constrained RNN

### GPU: ~7.2 s for 100 iter



### CPU: ~17.6 s for 100 iter



**Ref** Buyun Liang, Tim Mitchell, Ju Sun. NCVX: A General-Purpose Optimization Solver for Constrained Machine and Deep Learning. In Neural Information Processing Systems (NeurIPS) Workshop on Optimization for Machine Learning (OPT 2022).

# NCVX PyGRANSO: Advantages

### 3) General Tensor Variables

```
var_in = {"x1": [1], "x2": [1]}
```

*Scalar input*

```
# objective function
f = (8 * abs(x1**2 - x2) + (1 - x1)**2)
```

```
var_in = {"q": [n,1]}
```

*Vector input*

```
# objective function
qtY = q.T @ Y
f = 1/m * torch.norm(qtY, p = 1)
```

```
var_in = {"M": [d1,d2],"S": [d1,d2]}
```

*Matrix inputs*

```
# objective function
f = torch.norm(M, p = 'nuc') + eta * torch.norm(S, p = 1)
```

```
var_in = {"x_tilde": list(inputs.shape)}
```

*Higher order tensor input*

```
adv_inputs = X_struct.x_tilde
epsilon = eps
logits_outputs = model(adv_inputs)
f = -torch.nn.functional.cross_entropy(logits_outputs,labels)
```

**Ref** Buyun Liang, Tim Mitchell, Ju Sun. NCVX: A General-Purpose Optimization Solver for Constrained Machine and Deep Learning. In Neural Information Processing Systems (NeurIPS) Workshop on Optimization for Machine Learning (OPT 2022).

**User-friendly** is our philosophy

# Answering DOE's call

## Thrust C: How Should Domain Knowledge Be Modeled and Represented in Scientific Machine Learning?

An additional opportunity for domain-aware SciML research is in constructing modeling languages and frameworks that facilitate the inclusion of domain knowledge into the training process. Often, modeling languages and frameworks (e.g., [65, 66]) are designed to lower the barrier of entry for users by facilitating rapid and robust problem formulation. Extending the ways that SciML can express and incorporate domain knowledge could have far-reaching implications in much the same way that these tools now are regularly used for implicit features, such as algorithmic differentiation.

BASIC RESEARCH NEEDS FOR
**Scientific Machine Learning**
Core Technologies for Artificial Intelligence

POWER GRID INPUTS
Wind
Solar
Dams
Nuclear

Prepared for U.S.
Department of Energy
Advanced Scientific
Computing Research

U.S. DEPARTMENT OF ENERGY

**Ref** https://www.osti.gov/servlets/purl/1478744

**Robust Scientific Machine Learning**

# Constraint folding



(a) $n$ box constraints    (b) folded constraints

**Reduce # constraints**
- Reduce cost of QP in the SQP

**Equality into non-negative inequality**

$$h_j(\boldsymbol{x}) = 0 \iff |h_j(\boldsymbol{x})| \leq 0$$

**Inequality into nonnegative inequality**

$$c_i(\boldsymbol{x}) \leq 0 \iff \max\{c_i(\boldsymbol{x}), 0\} \leq 0$$

**All non-negative inequalities into one**

$$\mathcal{F}(|h_1(\boldsymbol{x})|, \cdots, |h_i(\boldsymbol{x})|, \max\{c_1(\boldsymbol{x}), 0\},$$
$$\cdots, \max\{c_j(\boldsymbol{x}), 0\}) \leq 0,$$

$\mathcal{F} : \mathbb{R}_+^{i+j} \mapsto \mathbb{R}_+ \ (\mathbb{R}_+ = \{\alpha : \alpha \geq 0\})$   Can be any function satisfying   $\mathcal{F}(\boldsymbol{z}) = 0 \implies \boldsymbol{z} = \boldsymbol{0}$

**Ref** Hengyue Liang, Buyun Liang, Le Peng, Ying Cui, Tim Mitchell, Ju Sun. *Optimization and Optimizers for Adversarial Robustness*. Under review at International Journal of Computer Vision (IJCV).

# 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

# General instruction

PyGRANSO

NCVX PyGRANSO Documentation

🔍 Search the docs ...

Introduction

Installation

Settings ⌄

Examples ⌄

Tips

NCVX Methods ⤢

Citing PyGRANSO

Tutorial Sessions ⌄

NCVX PyGRANSO Forum ⤢

## 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.

Our software announcement paper is available at https://arxiv.org/abs/2210.00973. This paper is accepted by the NeurIPS Workshop on Optimization for Machine Learning (OPT 2022). See our poster for more details.

# SVM: mathematical form

$$\min_{w,b} \frac{\nu}{2}\|w\|^2 + b\nu + \frac{1}{n}\sum_{i=1}^{n} \boxed{\max(0, 1 - (\langle w, x_i \rangle + b))}$$

nonsmoothness

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

$$\text{subject to } y_i(w^T \phi(x_i) + b) \geq 1 - \zeta_i,$$
$$\zeta_i \geq 0, i = 1, \ldots, n$$

SVC constrained version

**Ref** https://scikit-learn.org/stable/modules/sgd.html#online-one-class-svm
https://scikit-learn.org/stable/modules/svm.html#svc
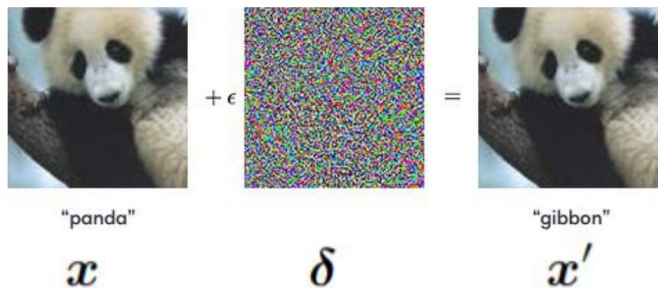
# NCVX PyGRANSO live coding for SVM

# NCVX PyGRANSO quick summary: SVM

**NVCX for unconstrained SVM**

- can handle nonsmoothness
- reliable termination condition (w/o ad-hoc maxiteration)
- line search criterion (w/o step size scheduler)

**NVCX is able to deal with general constrained problem (SVC)**

# Robustness evaluation: mathematical form



"panda"
$x$

$\delta$

"gibbon"
$x'$

Maximize loss function

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

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

Allowable perturbation with radius $\varepsilon$

Valid image constraints

Minimize robustness radius

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

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

Change the predicted class

Valid image constraints

# Robustness evaluation

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

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

**First general-purpose method for evaluating adversarial robustness**

**ROBUSTBENCH**
A standardized benchmark for adversarial robustness

**Reliability**

- SOTA methods
  No stopping criterion (only use maxit); step size scheduler

- PWCF (ours)
  Line search criterion and termination condition

**Generality**

- SOTA methods
  Can mostly only handle standard lp norm (l1,l2,linf)

- PWCF (ours)
  Distance metric beyond standard lp norm (l1,l2,linf). E.g., perceptual distance

$$d(\boldsymbol{x}, \boldsymbol{x}') \doteq \|\phi(\boldsymbol{x}) - \phi(\boldsymbol{x}')\|_2$$
$$\text{where} \quad \phi(\boldsymbol{x}) \doteq [\ \widehat{g}_1(\boldsymbol{x}), \ldots, \widehat{g}_L(\boldsymbol{x})\ ]$$

**Ref** Optimization and Optimizers for Adversarial Robustness. Liang, H., Liang, B., Peng, L., Cui, Y., Mitchell, T., & Sun, J. arXiv preprint arXiv:2303.13401.

# NCVX PyGRANSO live coding for robust evaluation

https://colab.research.google.com/drive/1vO4YCnfhFokyYG7D_ufUy_q-QKrFho48
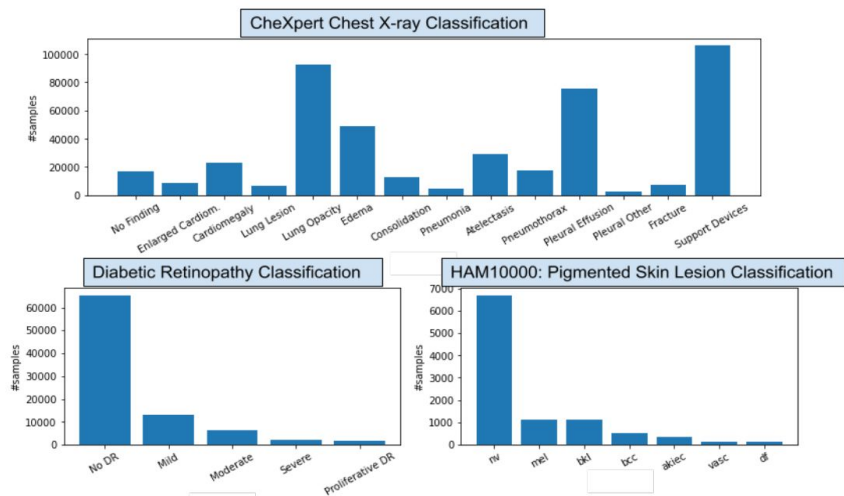
NCVX

# NCVX PyGRANSO quick summary: robustness

**NVCX for robustness evaluation**

- reliable termination condition (w/o ad-hoc maxiteration)
- line search criterion (w/o step size scheduler)

**NVCX is able to deal with general constraints (perceptual attack)**

# Learning with label imbalance



**Imbalance data in healthcare**

| | Predicted POS | Predicted NEG |
|---|---|---|
| POS | 70 | 30 |
| NEG | 1000 | 9000 |

**Accuracy:**      9070/10100 = 0.898
**True Positive Rate (Sensitivity, Recall):**    0.7
**True Negative Rate (Specificity):**    0.9
**Balanced Accuracy:**    (0.7 + 0.9)/2 = 0.80
**Precision (POS):**    70/1070 = 0.065
**F1 Score:**    2*0.065*0.7/(0.065 + 0.7) = 0.119

Reliable imbalance learning: **precision needed**!

# Learning with label imbalance

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

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

$$F_{\beta}(f_{\boldsymbol{\theta}}, t) = (1 + \beta^2) \frac{\text{precision}(f_{\boldsymbol{\theta}}, t) \cdot \text{recall}(f_{\boldsymbol{\theta}}, t)}{\beta^2 \text{precision}(f_{\boldsymbol{\theta}}, t) + \text{recall}(f_{\boldsymbol{\theta}}, t)}$$

One direction: directly optimizing the evaluation metric

**fix precision, optimize recall (FPOR):** $\quad \max_{\boldsymbol{\theta}, t} \text{recall}(f_{\boldsymbol{\theta}}, t) \quad$ s.t. $\text{precision}(f_{\boldsymbol{\theta}}, t) \geq \alpha,$

**fix recall, optimize precision (FROP):** $\quad \max_{\boldsymbol{\theta}, t} \text{precision}_t \quad$ s.t. $\text{recall}(f_{\boldsymbol{\theta}}, t) \geq \alpha,$

**optimize $F_{\beta}$ score (OFBS):** $\quad \max_{\boldsymbol{\theta}, t} F_{\beta}(f_{\boldsymbol{\theta}}, t),$

# NCVX PyGRANSO live coding for imbalance learning

https://colab.research.google.com/drive/1__OeV8OSbpszqPImaYQgwqOQC
XquuACl

# NCVX PyGRANSO quick summary: imbalance learning

**NVCX for robustness evaluation**

- reliable termination condition (w/o ad-hoc maxiteration)
- line search criterion (w/o step size scheduler)

**NVCX is able to deal with general constraints (e.g., precision/recall)**

# Closing

- Constrained DL (CDL) problems are everywhere

- Current methods for solving CDL are suboptimal
    - Projected gradient descent
    - Penalty methods
    - Lagrangian methods

- NCVX modeling framework + PyGranso solver is to address the gap
    - Principled stopping criterion, line search, and quasi-Newton method to obtain high-quality solution with reasonable speed

- Next steps
    - Stochastic optimization
    - Autoscaling