



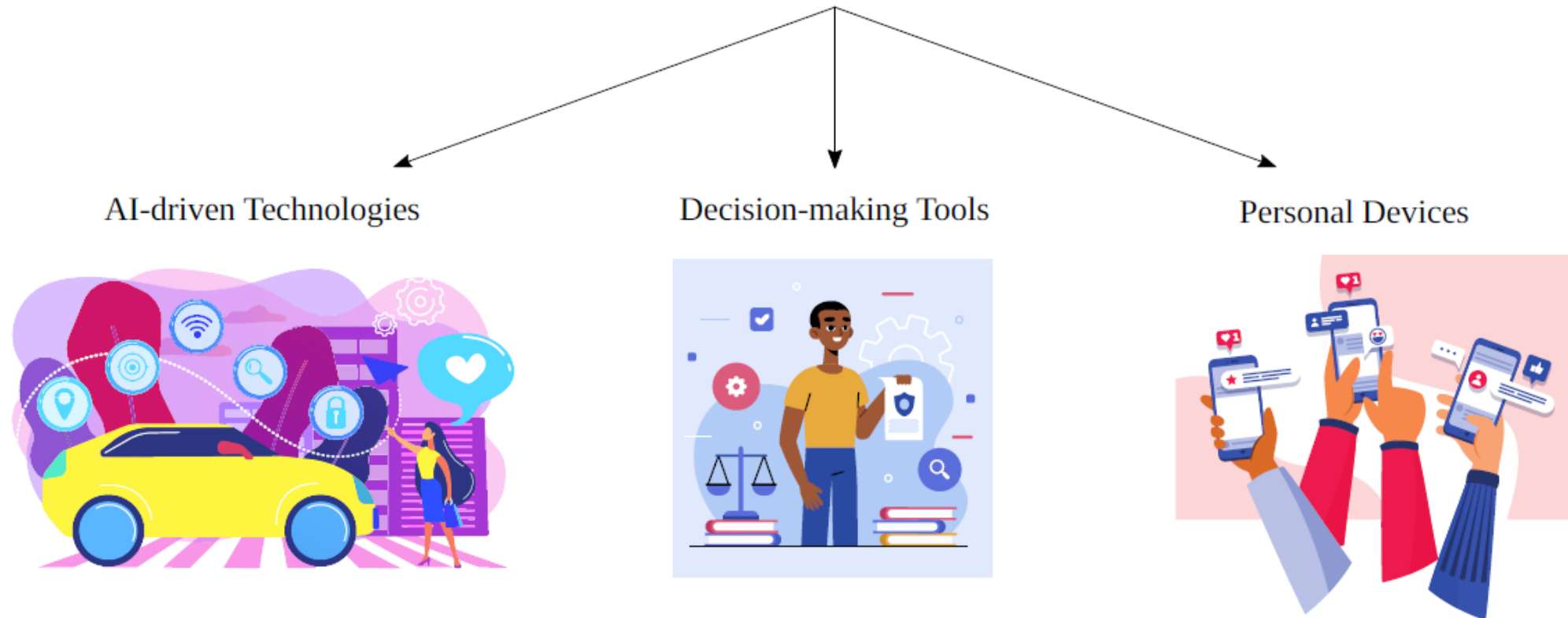
DE LA RECHERCHE À L'INDUSTRIE

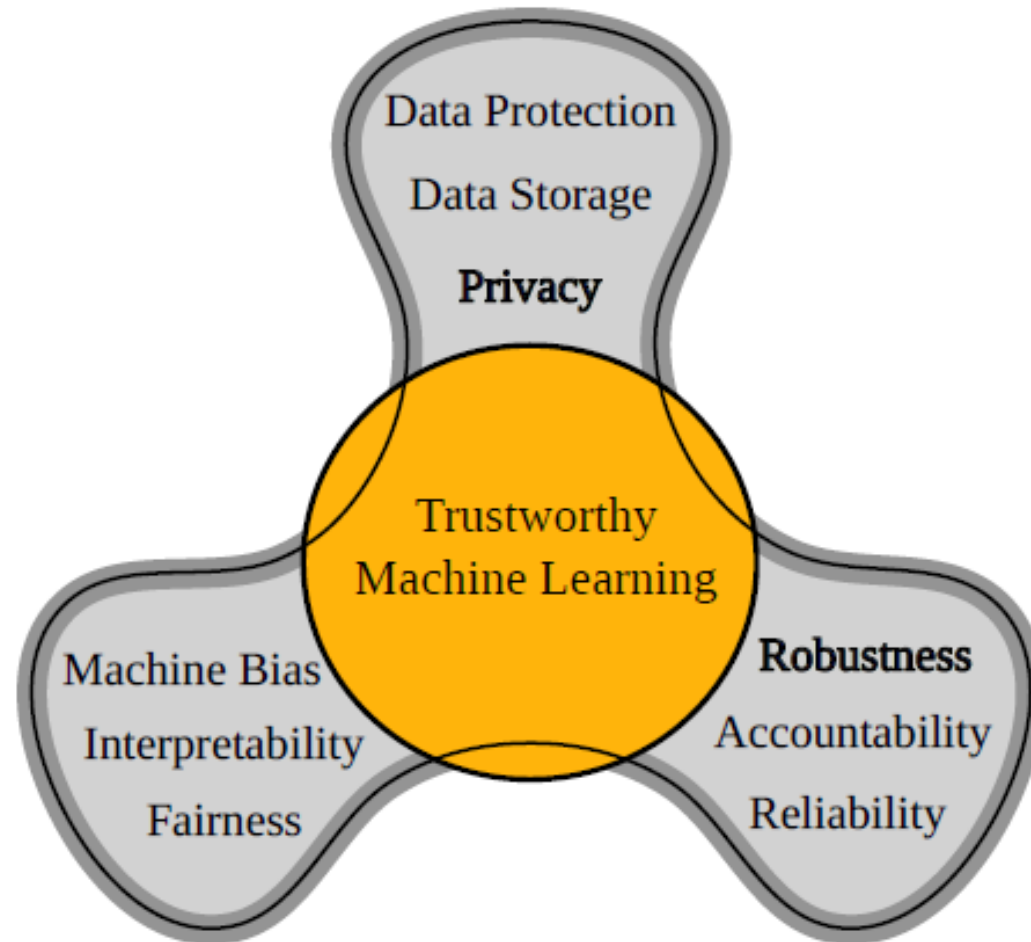
Robustesse des réseaux de neurones face aux attaques adversariales

2022-2023

Aurélien Mayoue & Cédric Gouy-Pailler (cours)

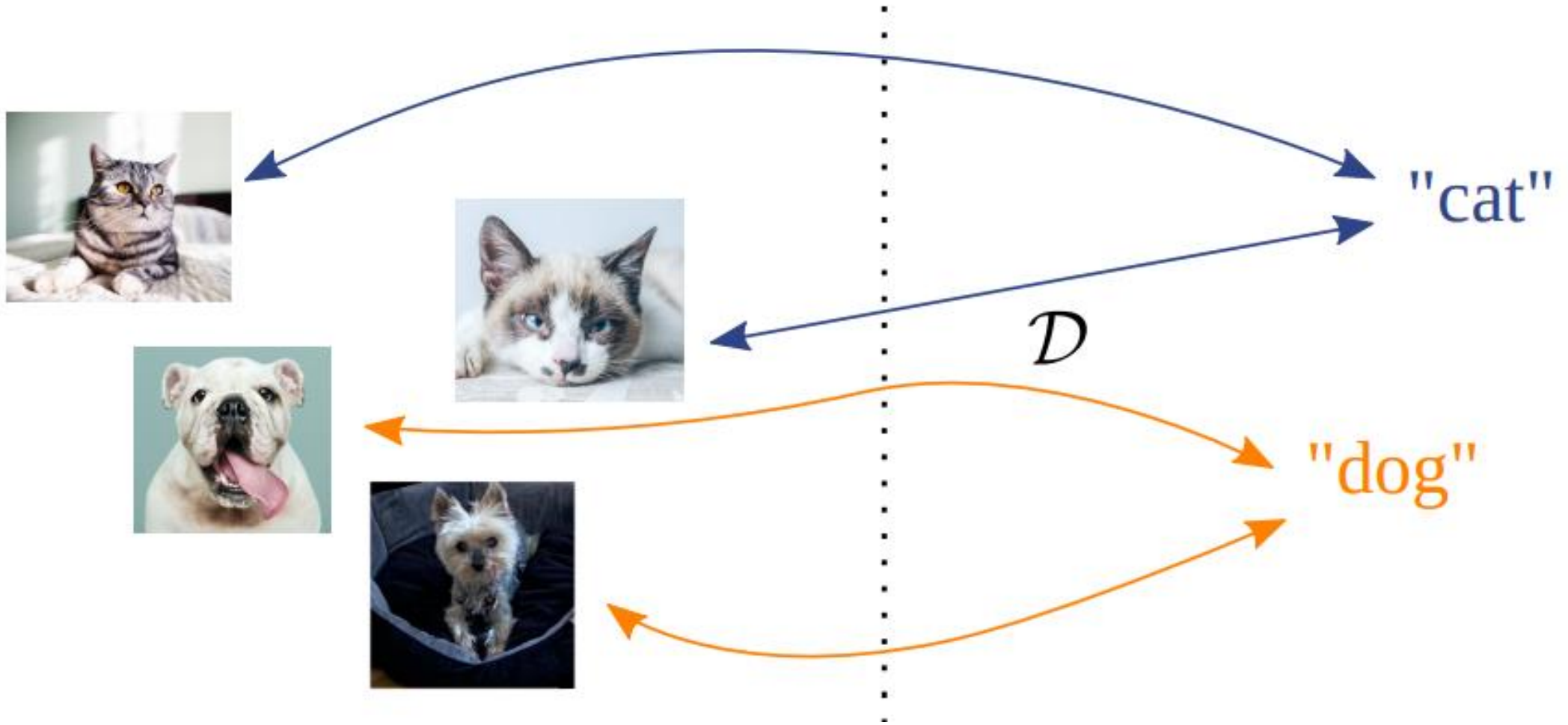
- Machine learning models recently gave **outstanding results** (e.g. vision, NLP)
- Industries and governments are starting to use them in **critical applications**



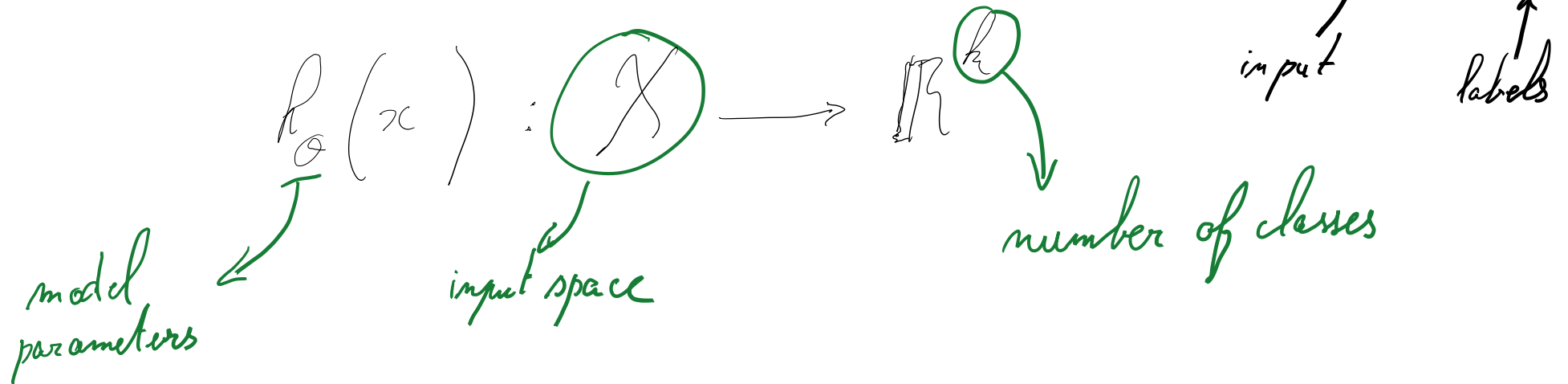


Adversarial attacks





A ground truth distribution \mathcal{D} explains the link between X and Y .



Loss function

$$l : \mathbb{R}^k \times \mathbb{Z}^+ \longrightarrow \mathbb{R}^+$$

$l(h(x), y)$ is called the loss function

Cross-entropy of softmax

softmax

[maximize prob of true class label]

[maximize \log of probability]

$\sigma : \mathbb{R}^k \rightarrow \mathbb{R}^k$
mapping from class logits to probabilities.

$$\sigma(z)_{[i]} = \frac{\exp(z_i)}{\sum_{j=1}^k \exp(z_j)}$$

$$\log(\sigma(h_\theta(x))_{[y]}) = \log\left(\frac{\exp(h_\theta(x)_{[y]})}{\sum_{j=1}^k \exp(h_\theta(x)_{[j]})}\right)$$

Cross-entropy

[minimize negative log of probability]

$$-\log \sigma(h(x)[y]) = \log \left(\sum_{j=1}^k \exp(h_\theta(x)[j]) \right) - h_\theta(x)[y]$$

Loss function

$$l(h_\theta(x), y) = \log \left(\sum_{j=1}^k \exp(h_\theta(x)[j]) \right) - h_\theta(x)[y]$$

Training a classifier

Optimize parameters Θ to minimize the average loss over some training set $\{x_i \in \mathcal{X}; y_i \in \mathbb{R}^+\}_{i=1..m}$

$$\underset{\Theta}{\text{minimize}} \frac{1}{m} \sum_{i=1}^m l(h_{\Theta}(x_i), y_i)$$

\Rightarrow Stochastic gradient descent on some minibatch $\mathcal{B} \subseteq \{1, \dots, m\}$

$$\Theta^{[t+1]} \leftarrow \Theta^{[t]} - \frac{\alpha}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla_{\Theta} l(h_{\Theta}(x_i), y_i)$$

$$\nabla_{\theta} \ell(h_{\theta}(x_i), y_i)$$

Gradient



Backpropagation

Automatic differentiation

Can we compute the gradient of the loss with respect to input x_i ?

Adjust the image to maximize the loss.

Untargeted Attack

$$\text{maximize}_{\delta \in \Delta} \ell(h_{\theta}(x + \delta), y)$$

set of allowable perturbations

$$\text{example } \|\cdot\|_{\infty} : \Delta = \{ \delta : \|\delta\|_{\infty} \leq \epsilon \}$$

$$\|\delta\|_{\infty} = \max_i |\delta_i|$$

Targeted attacks

maximize
 $\delta \in \Delta$

$$\underbrace{\ell(h_{\theta}(x + \delta), y)}_{\text{max the loss of the correct class}} - \underbrace{\ell(h_{\theta}(x + \delta), y_{\text{target}})}_{\text{minimize the loss of the target class}}$$

\Rightarrow

maximize
 $\delta \in \Delta$

$$h_{\theta}(x + \delta)_{[y_{\text{target}}]} - h_{\theta}(x + \delta)_{[y]}$$

More formal notation

$$R(h_\theta) = \mathbb{E}_{x,y \sim \mathcal{D}} [l(h_\theta(x), y)]$$

Risk

$$\hat{R}(h_\theta, D) = \frac{1}{|D|} \sum_{(x,y) \in D} l(h_\theta(x), y)$$

Empirical Risk.

⚠ We are trying to maximize the loss!
 ↳ in the direction of the positive gradient.

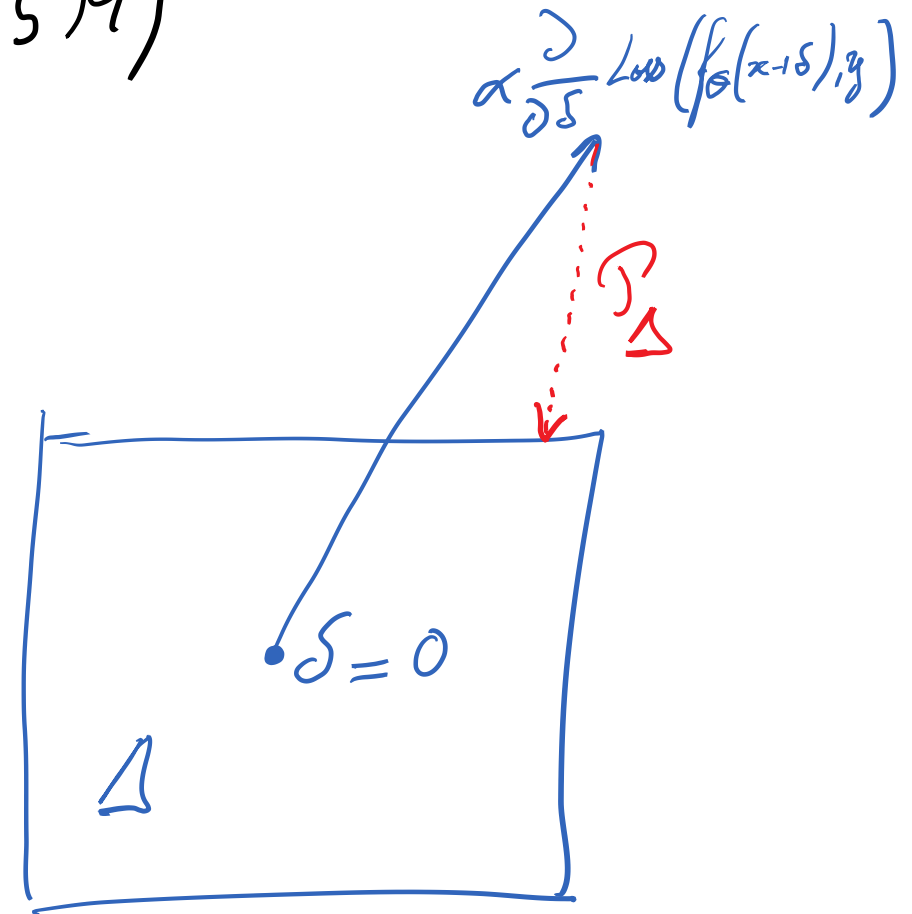
Projected Gradient Descent → move toward desired direction
 → ensure $\delta \in \Delta$

$$\delta := \text{Proj} \left(\delta + \alpha \frac{\partial}{\partial \delta} \text{Loss} (f_{\theta}(x + \delta), y) \right)$$

to be defined depending on kind of attack

Descent → projection

The Fast Gradient Sign Method (FGSM)



Linear
decision

$$h_{\theta}(x) = w^T x + b$$

$$\left| \begin{array}{l} w \in \mathbb{R}^m \\ b \in \mathbb{R} \\ y \in \{-1; +1\} \end{array} \right.$$

Cross-entropy

$$\ell(h_{\theta}(x), y) = \log(1 + \exp(-y \cdot h_{\theta}(x))) = L(y \cdot h_{\theta}(x))$$

Probability

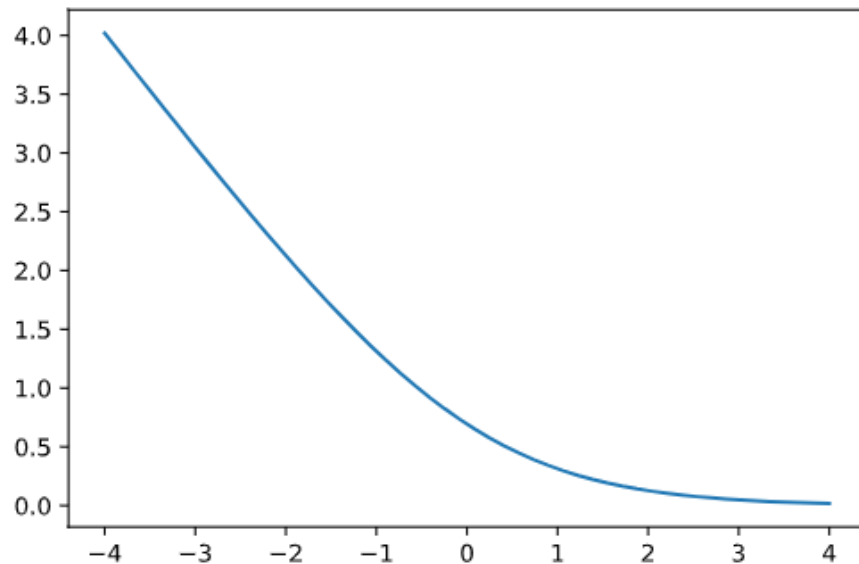
$$p(y = +1 | x) = \frac{1}{1 + \exp(-h_{\theta}(x))}$$

$$L(z) = \log(1 + \exp(-z))$$

$$\underset{\|\delta\| \leq \epsilon}{\text{maximize}} \ell(w^T(x+\delta), y) = \underset{\|\delta\| \leq \epsilon}{\text{maximize}} L(y \cdot (w^T(x+\delta) + b))$$

monotonic \Rightarrow maximize = minimize inside

```
x = np.linspace(-4,4)
plt.plot(x, np.log(1+np.exp(-x)))
```



$$\underset{\|\delta\| \leq \epsilon}{\text{max}}$$

$$= L\left(\min_{\|\delta\| \leq \epsilon} y \cdot (w^T(x+\delta) + b)\right)$$

$$= L\left(y \cdot (w^T x + b) + \min_{\|\delta\| \leq \epsilon} y \cdot w^T \delta\right)$$

$$\epsilon \begin{cases} w_i \geq 0 & \delta_i = -\epsilon \\ w_i < 0 & \delta_i = \epsilon \end{cases}$$

$$\Rightarrow \delta^* = -y \epsilon \cdot \text{sign}(w)$$

$$y = -1$$

$$\begin{cases} w_i \geq 0 & \delta_i = \epsilon \\ w_i < 0 & \delta_i = -\epsilon \end{cases}$$

$$y \cdot w^T \delta = y \cdot \sum_i -y \varepsilon \cdot \text{sign}(w_i) w_i = -y^2 \varepsilon \sum |w_i| \\ = -\varepsilon \|w\|_1$$

$$\max_{\|w\|_1 \leq \varepsilon} L(y \cdot (w^T (x + \delta) + b)) = L(y \cdot (w^T x + b) - \varepsilon \|w\|_1)$$

Analytical solution.

$$z_1 = x$$

$$z_{i+1} = f_i(w_i z_i + b_i) \quad i \in [1; d]$$

$$h_\theta(x) = z_{d+1}$$

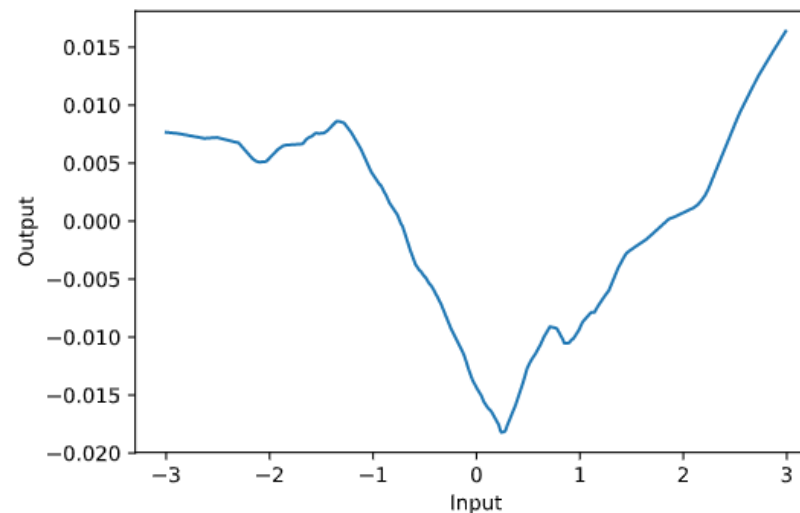
z_i : activations at layer i

$f_i(z) = \max(0, z)$: ReLU operator.

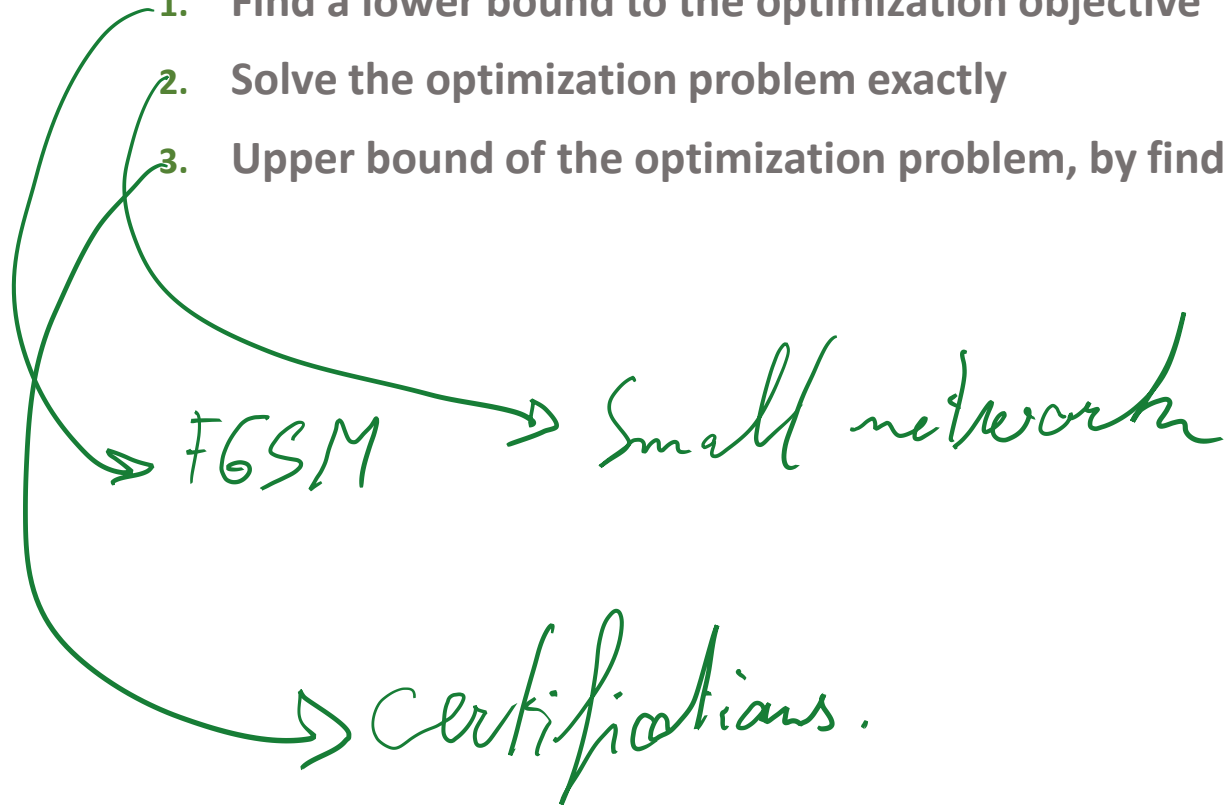
```
import torch
import torch.nn as nn
import torch.optim as optim

torch.manual_seed(0)
model = nn.Sequential(nn.Linear(1,100), nn.ReLU(),
                      nn.Linear(100,100), nn.ReLU(),
                      nn.Linear(100,100), nn.ReLU(),
                      nn.Linear(100,1))
opt = optim.SGD(model.parameters(), lr=1e-2)
for _ in range(100):
    loss = nn.MSELoss()(model(torch.randn(100,1)), torch.randn(100,1))
    opt.zero_grad()
    loss.backward()
    opt.step()

plt.plot(np.arange(-3,3,0.01), model(torch.arange(-3,3,0.01)[:,:None]).detach().numpy())
plt.xlabel("Input")
plt.ylabel("Output")
```



1. Find a lower bound to the optimization objective
2. Solve the optimization problem exactly
3. Upper bound of the optimization problem, by finding convex relaxation of the network structure



Defending against adversarial attacks

Untargeted adversarial attack : SOLVING THE INNER OPTIMIZATION PROBLEM

$$\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x + \delta), y)$$

Key insight: The same process that enabled us to learn the model parameters via gradient descent also allows us to create an adversarial example via gradient descent.

$$\frac{\partial}{\partial \delta} \text{Loss}(f_{\theta}(x + \delta), y)$$

Adversarial attacks: small, imperceptible change of an image maliciously designed to fool the model.

Clean example x



Label: "cat"

+ 0.006 ×

Perturbation δ



=

adversarial example $x + \delta$



Label: "dog"



Label: "cat"

+ 0.006 ×



=

Label: "~~dog~~
cat"

$$\min_{\theta} E_{x,y} [\text{Loss}(f_{\theta}(x), y)] \Rightarrow \min_{\theta} E \left[\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x+\delta), y) \right]$$

Adversarial training: during the learning procedure, replace x with x_{adv} . STATE-OF-THE-ART EMPIRICAL DEFENSE FOR l_2 ATTACKS.

Randomized smoothing: point-wise certification of the classification. STATE-OF-THE-ART CERTIFIED DEFENSE FOR l_2 ATTACKS

Provable defense: provably upper bound inner maximization

Adversarial training objective

$$\min_{\theta} \sum_{x, y \in S} \max_{\delta \in \Delta} \mathcal{L}_{\text{ass}}(f_{\theta}(x + \delta), y)$$

Gradient descent ?? Difficult due to max term inside.

Danskin's theorem :

$$\frac{\partial}{\partial \theta} \max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x+\delta), y) = \frac{\partial}{\partial \theta} \text{Loss}(f_{\theta}(x+\delta^*), y)$$

where

$$\delta^* = \arg \max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x+\delta), y)$$

\Rightarrow we can optimize through the max by just finding its maximizing value !!

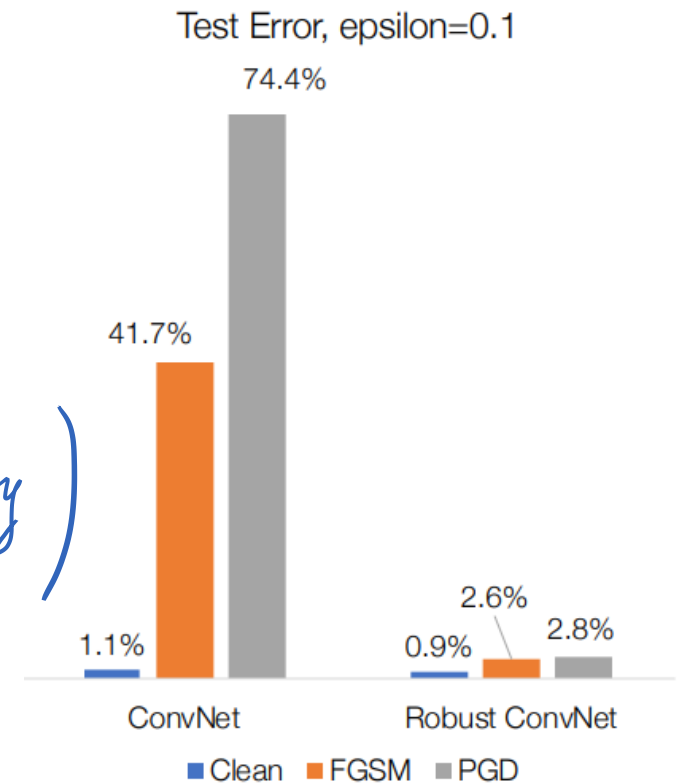
Repeat:

① Select minibatch B

② For each $(x, y) \in B$, compute adversarial example $\delta^*(x)$

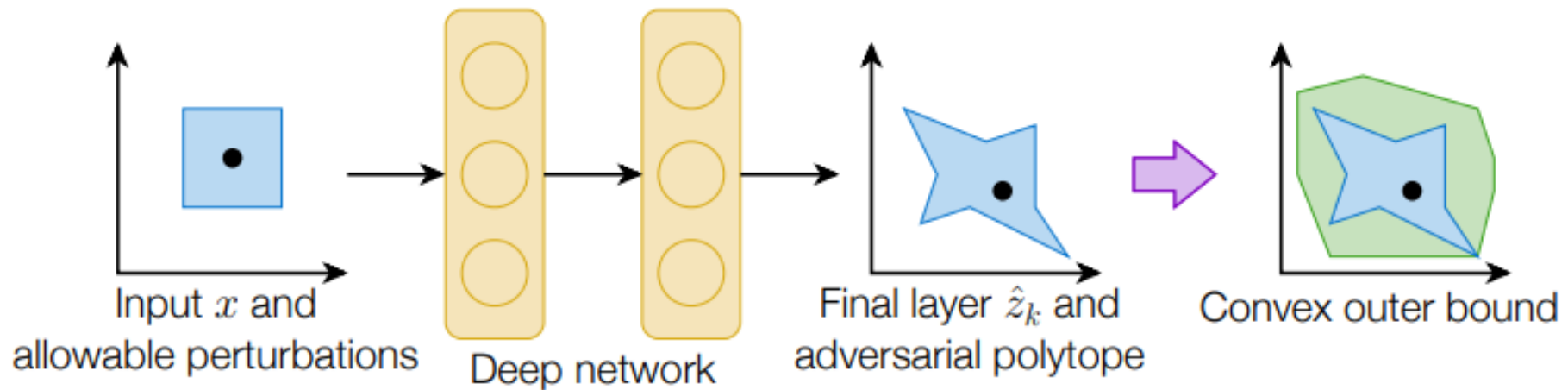
③ Update parameters

$$\theta := \theta - \frac{\alpha}{|B|} \sum_{(x, y) \in B} \frac{\partial}{\partial \theta} \text{Loss}(f_{\theta}(x + \delta^*(x)), y)$$



$$\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x+\delta), y) \leq \max \text{Loss}(f_{\theta}^{\text{rel}}(x+\delta), y)$$

E. Wong & Z. Kolter (2018)



Let's consider a k -layer feedforward ReLU-based neural network, $f_{\theta}: \mathbb{R}^{|\mathbf{x}|} \rightarrow \mathbb{R}^{|\mathbf{y}|}$

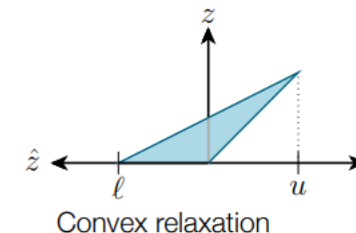
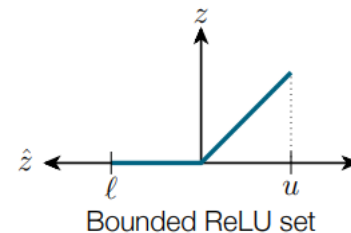
$$\text{s.t.} \quad \hat{z}_{i+1} = W_i z_i + b_i \quad \text{for } i=1, \dots, k-1$$

$$z_i = \max(\hat{z}_i, 0) \quad \text{for } i=2, \dots, k-1$$

Convex outer bound: linear relaxation of the ReLU activations.

Given known lower and upper bounds l, u for the pre-ReLU activations, we can replace ReLU $z = \max\{0; \hat{z}\}$ with their upper convex envelopes.

$$\begin{aligned} z &\geq 0 \\ z &\geq \hat{z} \\ -u\hat{z} + (u-l)z &\leq -ul \end{aligned}$$

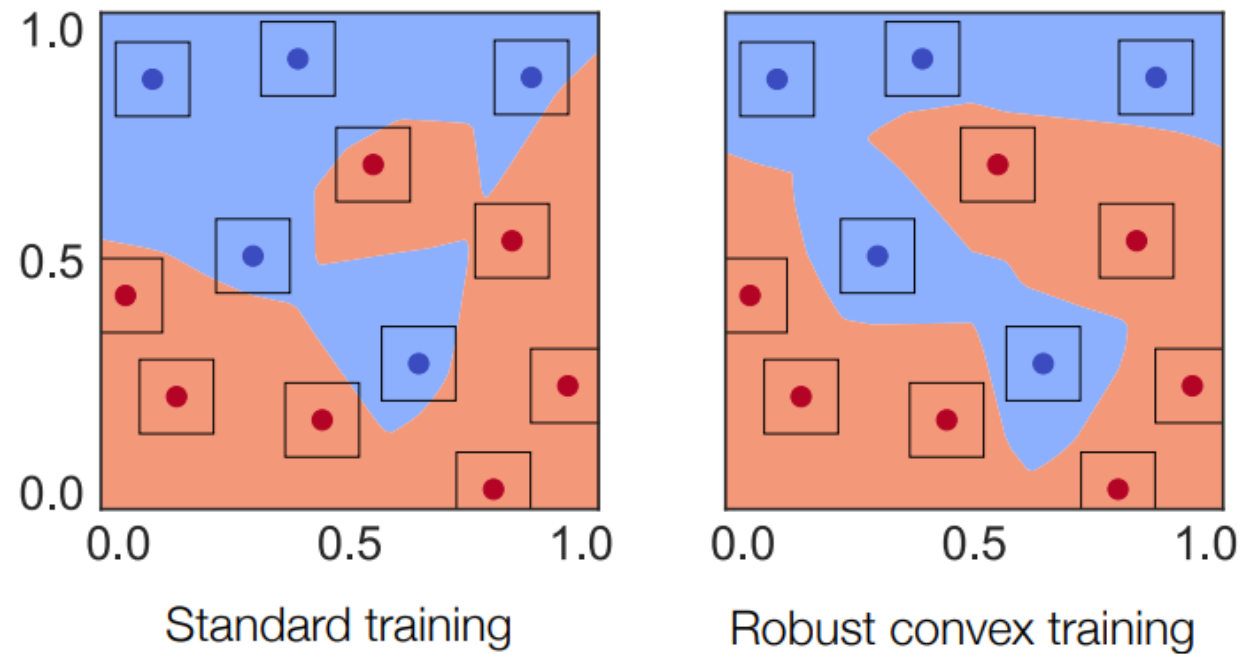


\Rightarrow MAXIMIZATION problem is now a convex linear program!

minimize the computed bound on loss (autodiff)

$$\underset{\Theta}{\text{minimize}} \sum_{i=1}^m l(\mathcal{J}_{\varepsilon, \Theta}(x_i), y_i) \geq \underset{\Theta}{\text{minimize}} \sum_{i=1}^m l\left(\max_{\delta \in \Delta} l\left(\mathcal{J}_{\Theta}(x_i + \delta), y_i\right)\right)$$

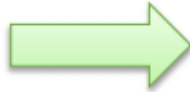
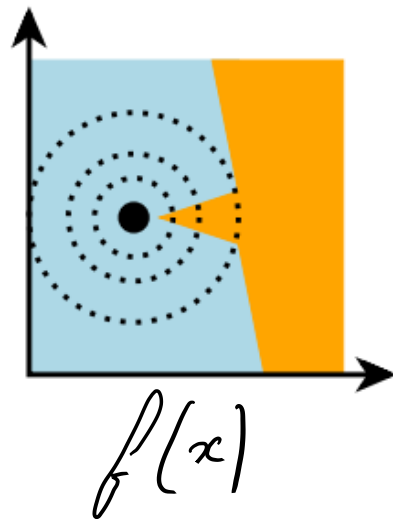
Simple 2D toy problem, 2-100-100-100-2 MLP network, trained with Adam (learning rate = 0.001, no hyperparameter tuning)



Example from Kolter.

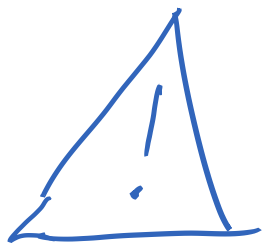
Randomization as a defense

We can smooth the decision region by adding Gaussian noise to the input and picking the majority class of the classifier over this noise.



$$g(x) = \arg \max_y \mathbb{P}_{\epsilon \sim \mathcal{N}(0; \sigma^2 I)} [f(x + \epsilon) = y]$$

Classify a bunch of versions perturbed by random noise. Take the majority vote.



Need to train on noisy images.



Randomized classifiers

$$\Gamma_{\Sigma} = \left\{ m : x \mapsto h(x+z) \mid h : \mathcal{X} \rightarrow \mathbb{R} \text{ continuous and } z \sim \mathcal{N}(0; \Sigma) \right\}$$

Gaussian noise injection: regularize the training phase

Here: test time to get robustness to data manipulation.

References

- ▶ Murphy, K.P. (2021). Probabilistic Machine Learning: An introduction. MIT Press. <https://probml.github.io/pml-book/book1.html>
- ▶ Kolter, Z. and Madry, A. (2018). Adversarial Robustness - Theory and Practice. <https://adversarial-ml-tutorial.org>
- ▶ Li, J. (2019). Course: Robustness in Machine Learning. <https://jerryzli.github.io/robust-ml-fall19.html>