

Mixed Nash Equilibria in the Adversarial Examples Game

¹ Facebook AI Research, Paris, France

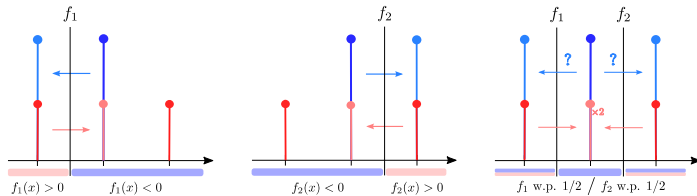
² Miles Team, LAMSADE, Université Paris-Dauphine, Paris, France

³ CREST, ENSAE, Paris, France

⁴ Ecole Polytechnique Fédérale de Lausanne (EPFL), Switzerland

* Equal contribution

A Motivating Example



- On left and in the middle: the classifier is deterministic. The adversarial risk is 3/4.
- On right: The classifier is randomized. The adversarial risk is 1/2.

The best attack is also randomized: if the attacker takes a deterministic decision, the classifier can play a deterministic strategy to counter him.

There exists a **Mixed Nash equilibrium in this game**. Can we generalize it?

Setting

- Classification problem on $X \times Y$. P is a Borel probability distribution over $X \times Y$.
- Θ is a set of classifiers.
- A loss $l : \Theta \times (X \times Y) \rightarrow \mathbb{R}$ (possibly the 0/1 loss).

Adversarial deterministic risk:

$$R_{adv}^e(\theta) := \mathbb{E}_{(x,y) \sim P} \left[\sup_{x' \in \mathcal{X}, d(x,x') \leq \epsilon} l(\theta, (x', y)) \right]$$

Adversarial randomized risk :

$$R_{adv}^e(\mu) := \mathbb{E}_{(x,y) \sim P} \left[\sup_{x' \in \mathcal{X}, d(x,x') \leq \epsilon} \mathbb{E}_{\theta \sim \mu} [l(\theta, (x', y))] \right]$$

Risk minimization problems:

$$V_{rand}^e := \inf_{\mu \in \mathcal{M}_+^1(\Theta)} R_{adv}^e(\mu), \quad V_{det}^e := \inf_{\theta \in \Theta} R_{adv}^e(\theta)$$

Remark: $V_{rand}^0 = V_{det}^0 \leq V_{rand}^e \leq V_{det}^e$

Adversarial distributions/random adversaries:

$$\mathcal{A}_\epsilon(P) := \left\{ Q \mid \exists \gamma, d(x, x') \leq \epsilon, y = y' \text{ } \gamma\text{-a.s., } \Pi_{1\mu}\gamma = P, \Pi_{2\mu}\gamma = Q \right\}$$

Such an attacker can map any point randomly in the ball of radius ϵ . It is a Wasserstein ball for a well chosen cost.

Proposition: For a given classifier μ , the adversarial randomized risk equals:

$$R_{adv}^e(\mu) = \sup_{Q \in \mathcal{A}_\epsilon(P)} \mathbb{E}_{(x',y') \sim Q, \theta \sim \mu} [l(\theta, (x', y'))]$$

The supremum is attained and the optimum might be attained with a deterministic mapping.

Adversarial Examples Game

Primal game: $\inf_{\mu \in \mathcal{M}_+^1(\Theta)} \sup_{Q \in \mathcal{A}_\epsilon(P)} \mathbb{E}_{Q, \mu} [l(\theta, (x, y))]$. **Classifier goal:** be robust to every attacks.

Dual game: $\sup_{Q \in \mathcal{A}_\epsilon(P)} \inf_{\mu \in \mathcal{M}_+^1(\Theta)} \mathbb{E}_{(x,y) \sim Q, \theta \sim \mu} [l(\theta, (x, y))]$. **Attacker goal:** find an attack to fool

any classifier.

Denoting by D^e the value of the dual formulation, we have:

$$D^e \leq V_{rand}^e \leq V_{det}^e$$

Is there always equality of the left terms? Does there exist Nash equilibria in this game?

Theorem: Strong duality always holds in the randomized setting

$$\inf_{\mu \in \mathcal{M}_+^1(\Theta)} \max_{Q \in \mathcal{A}_\epsilon(P)} \mathbb{E}_{\theta \sim \mu, (x,y) \sim Q} [l(\theta, (x, y))] = \max_{Q \in \mathcal{A}_\epsilon(P)} \inf_{\mu \in \mathcal{M}_+^1(\Theta)} \mathbb{E}_{\theta \sim \mu, (x,y) \sim Q} [l(\theta, (x, y))]$$

Interpretations:

- Always exist approximate Mixed Nash Equilibria in the adversarial examples game.
- If the infimum is attained, there exist Mixed Nash Equilibria.

Entropic Regularization and Algorithms

Given empirical distribution $P_n = \sum_{i=1}^n \delta_{(x_i, y_i)}$ and a finite set of classifiers $\{\theta_1, \dots, \theta_L\}$. Can we

learn the optimal randomized classifier over this finite set of classifiers, i.e. optimize the weights of the probability that θ_i appears?

Entropic Relaxation:

$$\inf_{\mu \in \mathcal{M}_+^1(\Theta)} \sum_{i=1}^n \sup_{Q \in \Gamma_{\epsilon, \alpha}} \mathbb{E}_{Q, \mu} [l(\theta, (x, y))] - \alpha \text{KL} \left(Q \left\| \frac{1}{N} \mathbb{U}_{(x,y)} \right. \right) = \inf_{\mu \in \mathcal{M}_+^1(\Theta)} \sum_{i=1}^n \frac{\alpha_i}{N} \log \left(\int \exp \frac{\mathbb{E}_{\mu} [l(\theta, (x, y))]}{\alpha_i} d\mathbb{U}_{(x,y)} \right)$$

- Approximates well the adversarial risk.

- Convex and smooth objective: Algorithm with rate of convergence in $O(T^{-2})$.

Oracle-Based Algorithm:

- Inspired by robust optimization and subgradient methods (Danskin Theorem).
- Rate of convergence of order $O(\delta + T^{-1/2})$ where δ denotes the quality of the gradient estimation.

Experiments

Algorithm 2 Adversarial Training for Mixtures

L : number of models, T : number of iterations,

T_θ : number of updates for the models θ ,

T_λ : number of updates for the mixture λ ,

$\lambda_0 = (\lambda_0^1, \dots, \lambda_0^L)$, $\theta_0 = (\theta_0^1, \dots, \theta_0^L)$

for $t = 1, \dots, T$ **do**

 Let B_t be a batch of data.

if $t \bmod (T_\theta L + 1) \neq 0$ **then**

k sampled uniformly in $\{1, \dots, L\}$

$\tilde{B}_t \leftarrow$ Attack of images in B_t for the model (λ_t, θ_k)

$\theta_k^t \leftarrow$ Update θ_k^{t-1} with \tilde{B}_t for fixed λ_t with a SGD step

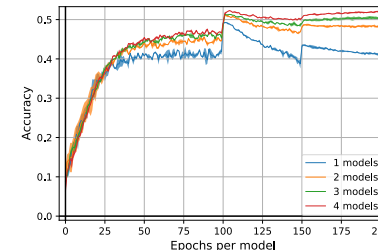
else

$\lambda_t \leftarrow$ Update λ_{t-1} on B_t for fixed θ_t with oracle-based or regularized algorithm with T_λ iterations.

end

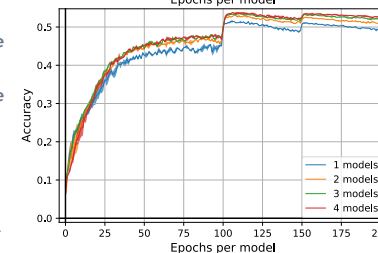
end

Proposed heuristic algorithm for deep learning



Models	Acc.	APGD _{CE}	APGD _{DLR}	Rob. Acc.
1	81.9%	47.6%	47.7%	45.6%
2	81.9%	49.0%	49.6%	47.0%
3	81.7%	49.0%	49.3%	46.9%
4	82.6%	49.7%	49.8%	47.2%

Accuracy under PGD attack on a ResNet18 model for CIFAR10 dataset using Adversarial Training loss



Models	Acc.	APGD _{CE}	APGD _{DLR}	Rob. Acc.
1	79.6%	50.9%	48.9%	48.3%
2	80.3%	52.3%	51.2%	50.2%
3	80.7%	52.8%	51.7%	50.7%
4	80.9%	53.0%	51.8%	50.8%

Accuracy under PGD attack on a ResNet18 model for CIFAR10 dataset using TRADES loss

Take a photo to learn more:

