

An analytical and geometric perspective on adversarial learning.

Nicolás García Trillos

UNIVERSITY OF WISCONSIN - MADISON

• Adversarial training / learning.

$$(AT) \inf_{\theta \in \Theta} \sup_{\tilde{\mu}: D(\mu, \tilde{\mu}) \leq \epsilon} \mathbb{E}_{\tilde{z} \sim \tilde{\mu}} [\ell(\tilde{z}, \theta)]$$

- μ distrib on $z = (x, y)$ $x \in \mathbb{R}^d, y \in \{0, 1\}$
- θ parameter of statistical model
- $\ell(\cdot, \cdot)$
- $D(\cdot, \cdot)$

- Regularized Risk minimization.

$$(R) \inf_{\theta \in \Theta} \mathbb{E}_{z \sim \mu} [\ell(z, \theta)] + \lambda R(\theta)$$

- Classical statistics.
- Inverse Problems
- Graph-Based Learning.

Q: What is the relationship between adversarial training and regularization?

Sometimes a very direct one.
Consider the following setup:

- $\Theta = \mathbb{R}^d$
- $\ell(z, \theta) = (\langle \theta, x \rangle - y)^2$
- $D(\mu, \tilde{\mu}) = W_{C_p}(\mu, \tilde{\mu})$

$$:= \inf_{\pi \in \Pi(\mu, \tilde{\mu})} \iint C_p(z, \tilde{z}) d\pi(z, \tilde{z})$$

where $C_p(z, \tilde{z}) = \begin{cases} \|x - \tilde{x}\|_p & \text{if } y = \tilde{y} \\ +\infty & \text{if } y \neq \tilde{y}. \end{cases}$

Then :

$$\inf_{\theta \in \Theta} \sup_{\tilde{\mu} : D_p(\tilde{\mu}, \mu) \leq \varepsilon} |E_{(x, \tilde{x}) \sim \tilde{\mu}}[\ell(\tilde{z}, \theta)]|$$

=

$$\inf_{\theta \in \Theta} \left(\sqrt{|E_{(x, y) \sim \mu}[\ell(z, \theta)]|} + \sqrt{\varepsilon} \|\theta\|_q \right)^2$$

$$\frac{1}{p} + \frac{1}{q} = 1$$

[Chen et al 20'], [Blanchet et al 19']
[Kuhn et al 19'], ...

Today's talk Based on:

- [NGT, MURRAY 20']:
"Adversarial classification:
necessary conditions
and geometric flows".
- [C. GARCIA TRILLOS, NGT 21']
"On the regularized risk of
distributionally robust learning
over deep neural networks".
- [BUNGERT, NGT, MURRAY, 21']
"The geometry of adversarial
learning in binary classification".

Setup:

- \textcircled{H} : Borel Subsets of \mathbb{R}^d

From now on use $A \subseteq \mathbb{R}^d$ instead of Θ

- $(x, y) \in \mathbb{R}^d \times \{0, 1\}$.

- $D(\mu, \tilde{\mu}) = W_{\infty}(\mu, \tilde{\mu})$

$\hat{d}(z, \tilde{z}) := \begin{cases} \underline{d}(x, \tilde{x}) & \text{if } y = \tilde{y} \\ +\infty & \text{if } y \neq \tilde{y} \end{cases}$

$\left\{ \begin{array}{l} W_{\infty}(\mu, \tilde{\mu}) \\ \inf_{\pi \in \Pi(\mu, \tilde{\mu})} \sup_{\{z, \tilde{z}\} \in \text{supp} \pi} \hat{d}(z, \tilde{z}) \end{array} \right.$

- $Q(z, A) = \begin{cases} 1 & \text{if } \underline{1}_A(x) \neq y \\ 0 & \text{otherwise.} \end{cases}$

Here:

$L_{\varepsilon}(A)$

$$\inf_A \sup_{\tilde{\mu}: W_{\infty}(\mu, \tilde{\mu}) \leq \varepsilon} \mathbb{E}_{\tilde{z} \sim \tilde{\mu}} [Q(\tilde{z}, \Theta)]$$

$$\varepsilon = 0 : A_0^* := \operatorname{argmin}_A L_0(A)$$

$$A_0^* := \left\{ x \in \mathbb{R}^d : \omega_1 p_1(x) \geq \omega_0 p_0(x) \right\}$$

Here:

$$\begin{aligned} \omega_0 &\rightarrow \mu(\mathbb{R}^d \times \{0\}) \\ \omega_1 &\rightarrow \mu(\mathbb{R}^d \times \{1\}) \\ \rho_0 &\rightarrow \mu(dx | y=0) \\ \rho_1 &\rightarrow \mu(dx | y=1) \\ \rho &\rightarrow \mu(dx) \end{aligned}$$

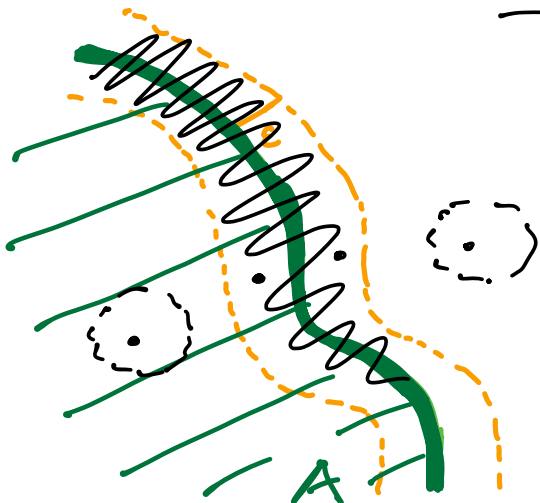
Q: How should the boundary of A_0^* change to track solutions to $\inf_A L_\varepsilon(A)$ as ε grows?

$$\begin{cases} \varepsilon \in (0, \varepsilon_0) \mapsto A_\varepsilon \\ A_0 = A_0^* \end{cases}$$

s.t. $A_\varepsilon \in \underset{A}{\operatorname{argmin}} L_\varepsilon(A)$

Here it is convenient to notice that:

$$\begin{aligned} L_\varepsilon(A) &= \frac{\sup_{\tilde{\mu}: W_{\tilde{\mu}}(\mu, \tilde{\mu}) \leq \varepsilon} \mathbb{E}_{\tilde{x} \sim \tilde{\mu}} [\ell(\tilde{x}, A)]}{=} \\ &= \frac{\mathbb{E}_{(x, y) \sim \mu} \left[\sup_{\tilde{x} \in B_\varepsilon(x)} \ell((\tilde{x}, y), A) \right]}{=} \end{aligned}$$



$$\begin{aligned} &= \int_{\partial A} \rho(x) dx \\ &+ \int_{A^{-\varepsilon}} \omega_0 \rho_0 dx + \int_{(A^\varepsilon)^c} \omega_1 \rho_1 dx \end{aligned}$$

[NGT, MURRAY, 20']

In 1D first: Suppose Bayes classifier

has the form $A_0 = \bigcup_{k=1}^K [a_k(0), b_k(0)]$

Under a "strict crossing" condition

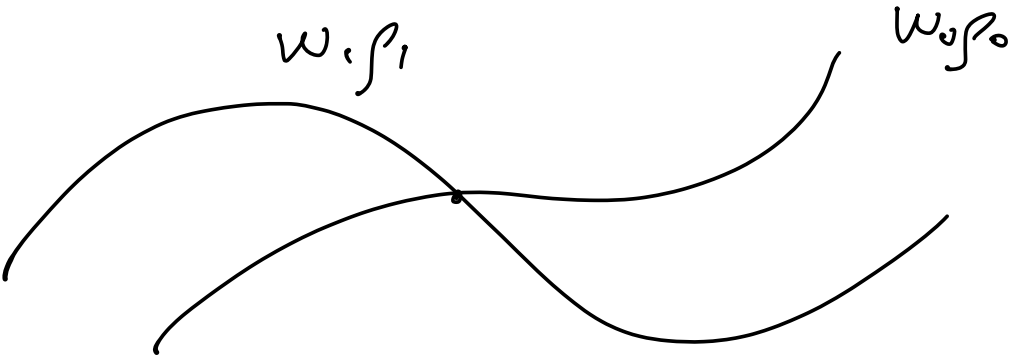
for $w_0 p_0$ and $w_i p_i$, the

following system of ODEs

tracks solutions for all
small enough ε :

$$\begin{cases} \frac{db_k(\varepsilon)}{d\varepsilon} = - \left(\frac{w_0 p_0'(b_k + \varepsilon) + w_i p_i'(b_k - \varepsilon)}{w_0 p_0'(b_k + \varepsilon) - w_i p_i'(b_k - \varepsilon)} \right) \\ \underline{b_k(0)} \end{cases}$$

and similar eqns for a_k .



Comments:

① Connection to Optimal Trespot problem:

- [Bhagoji et al 19']
 - [Pydi + Jog 19']
- } $w_0 = w,$

② 1D setting does not reveal the geometric structure of the general problem...

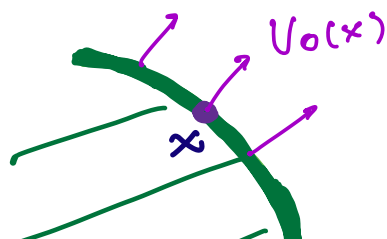
• In $d > 1$:

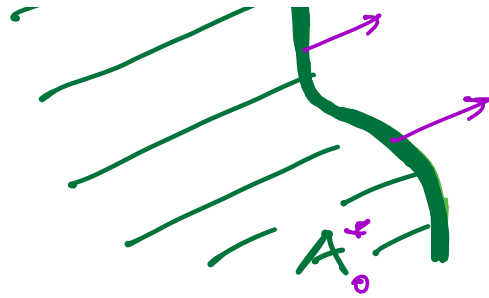
Messy equations (existence?)

but at $\varepsilon = 0$ we can try to

answer how the boundary of A_0

changes infinitesimally:





$$V_0(x) = \left(- \frac{\nabla \rho \cdot \vec{n} + \rho \sum_i \kappa_i}{(\omega_i \nabla \rho_i - \omega_0 \nabla \rho_0) \cdot \vec{n}} \right) \vec{n}$$

SAME infinitesimal change
as if we were tracking

solutions of : (R')

$$\inf_A \left\{ \int_{\partial A} \rho(x) d\mathcal{H}^{d-1}(x) + \varepsilon \text{Per}(A) \right\}$$

where :

$$\text{Per}(A) = \int_{\partial A} \rho(x) d\mathcal{H}^{d-1}(x)$$

Comments :

- ① So Perimeter is connected to the regularization induced by (AT) .
- ② Now, is (R') equivalent to (AT) ?

NO, BUT :

Take :

$$\begin{aligned} & \underline{E_{z \sim \mu} \left[\sup_{\tilde{x} \in B_\varepsilon(x)} Q(\tilde{x}, y, A) \right]} \\ &= E_{z \sim \mu} [Q(x, y, A)] \\ &+ \varepsilon \left(\frac{E_{z \sim \mu} \left[\sup_{\tilde{x} \in B_\varepsilon(x)} Q(\tilde{x}, y, A) \right] - E_{z \sim \mu} [Q(x, y, A)]}{\varepsilon} \right) \end{aligned}$$

Per_ε(A).

Theorem : [Bungert, NGT, MURRAY, 21]

$$(AT) = \inf_A \left\{ E_{z \sim \mu} [Q(x, y, A)] + \varepsilon \text{Per}_\varepsilon(A) \right\}$$

where :

$$(1) \text{Per}_\varepsilon(A) \geq 0 \quad \forall A.$$

② $\text{Per}_\varepsilon(\cdot)$ is submodular:

$$A, B \subseteq \mathbb{R}^d$$

$$\begin{aligned} \text{Per}_\varepsilon(A \cup B) + \text{Per}_\varepsilon(A \cap B) \\ \leq \text{Per}_\varepsilon(A) + \text{Per}_\varepsilon(B) \end{aligned}$$

③ let

$$\text{TV}_\varepsilon(u) := \int_{-\infty}^{\infty} \text{Per}_\varepsilon(\{u \geq t\}) dt$$

$$u: \mathbb{R}^d \rightarrow \mathbb{R}.$$

Then: TV_ε is convex, 1-homogeneous,
and l.s.c w.r.t appropriate
topology.

④ The problem:

$$\min_{f: \mathbb{R}^d \rightarrow [0,1]} \left\{ \mathbb{E}_{(x,y) \sim \mu} [|f(x) - y|] + \varepsilon \text{TV}_\varepsilon(f) \right\}$$

is an exact convex relaxation
of (AL).

Remark:

• $TV_\varepsilon(u) =$

$$\frac{\omega_1}{\varepsilon} \int_{\mathbb{R}^d} (u(x) - \inf_{\tilde{x} \in B_\varepsilon(x)} u(\tilde{x})) (\rho_1(dx))$$

$$+ \frac{\omega_0}{\varepsilon} \int_{\mathbb{R}^d} (\sup_{\tilde{x} \in B_\varepsilon(x)} u(\tilde{x}) - u(x)) (\rho_0(dx))$$

• Notice that TV_ε depends on
the full μ !

$$(Per_\varepsilon(A) = TV_\varepsilon(\mathbb{1}_A) \text{ too})$$

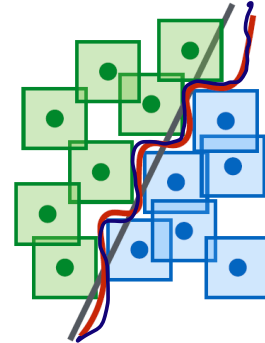
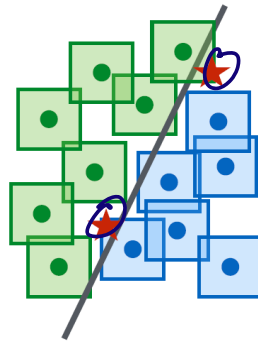
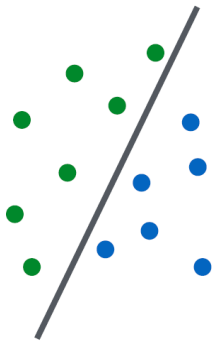
- Take $\omega_0 = \omega_1 = \frac{1}{2}$ and $\rho_0 = \rho_1 = \rho$ in the formula for TV_e . We recover:

Related to:

- [Barchiesi et al 10']
- [Chumbolle et al 12']



Per vs Per_ε



From MADRY et al 19'

Q: What is the connection with

regularizers used in graph Based learning?

Different adversarial model:

- Native chooses $\tilde{x} \sim \mathcal{P}_{B_\epsilon(x)}$
Adversary decides to accept/reject \tilde{x} (with the goal of maximizing their payoff).

(AL')

$$= \inf_{u: \mathbb{R}^d \rightarrow [0,1]} \left\{ \mathbb{E}_{z \sim \mu} [|u(x) - y|] + \epsilon \tilde{TV}_\epsilon(u) \right\}$$

where

$$\tilde{TV}_\epsilon(u) =$$

$$\frac{\omega_1}{\varepsilon} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \frac{\eta_\varepsilon(|x-\tilde{x}|)}{\mathcal{P}(B_\varepsilon(x))} (u(x) - u(\tilde{x}))_+ \rho(d\tilde{x}) \rho(dx)$$

$$+ \frac{\omega_0}{\varepsilon} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \frac{\eta_\varepsilon(|x-\tilde{x}|)}{\mathcal{P}(B_\varepsilon(x))} (u(x) - u(\tilde{x}))_- \rho(d\tilde{x}) \rho_0(dx)$$

Remark: $\omega_0 = \omega_1 = \frac{1}{2}$

$$\rho_0 = \rho_1 = \rho \quad :$$

$$\frac{1}{2\varepsilon} \iint \frac{\eta_\varepsilon(|x-\tilde{x}|)}{\mathcal{P}(B_\varepsilon(x))} |u(x) - u(\tilde{x})| \rho(d\tilde{x}) \rho(dx)$$

$$\text{When } \rho(dx) = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$$

- $\rho(B_\epsilon(x)) = \text{degree of geometric graph}$

$$\frac{1}{2n^2\epsilon} \sum_i \sum_j \frac{\chi_\epsilon(|x_i - x_j|)}{d_\epsilon(x_i)} |u(x_i) - u(x_j)|$$

Thank you for your attention!

Special thanks to:

-NSF Grant: DMS-2005797

-All my collaborators.

