

# Adversarial Training Through the Lens of Optimal Transport

Nicolás García Trillos  
University of Wisconsin-Madison

Kantorovich Initiative seminar series  
February 2023

Based on joint works with:

- Jakwang Kim (Wisc)
- Camilo García Trillos (UCL)
- Matt Jacobs (Purdue)

Based on joint works with:

- Jakwang Kim (joining Kantorovich initiative this Fall!)
- Camilo García Trillos (UCL)
- Matt Jacobs (Purdue)

*Neural networks, although accurate on clean data, are sensitive to adversarial attacks:*

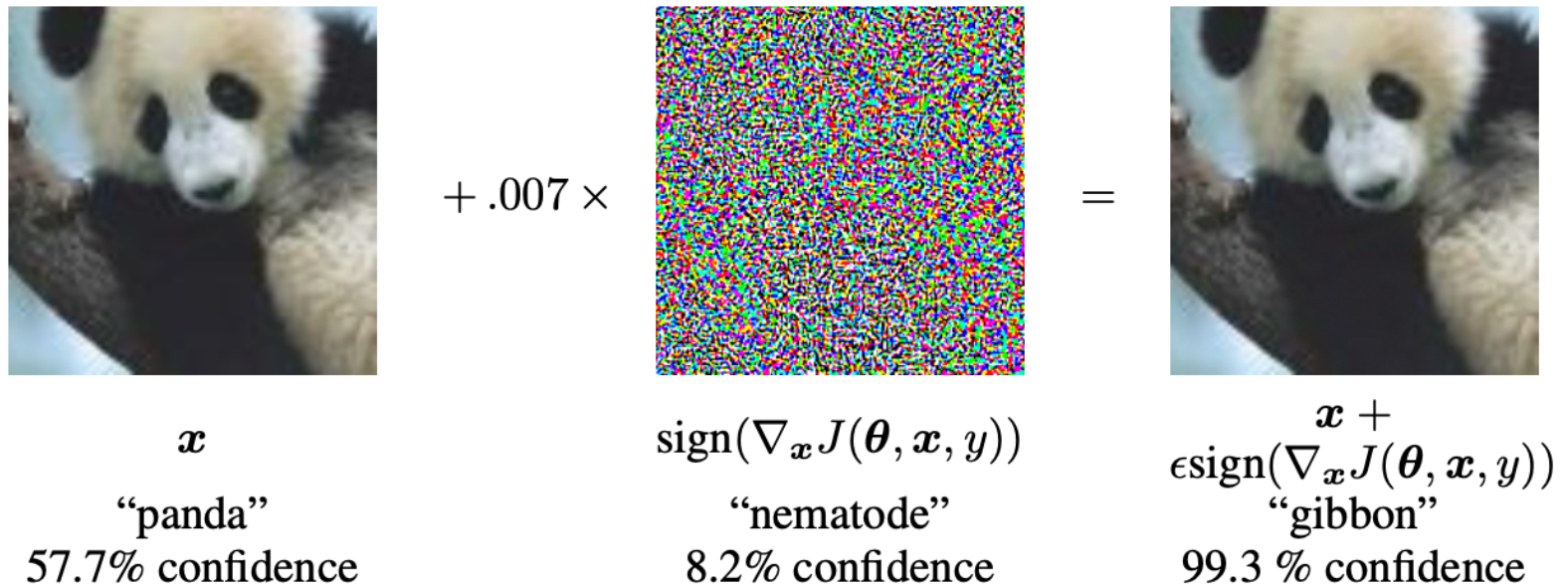


Figure: Picture taken from Goodfellow et al. (2015)

[Szegedy et al. 2014] , [Goodfellow et al. 2015]



**Figure:** An adversarial attack of a clean image in a safety-critical setting. Picture taken from Eykholt et al. (2018)

# Formalization of adversarial training problem

How to train classifiers to be robust to (specific) adversarial attacks?:

$$\min_{\theta \in \Theta} \mathbb{E}_{(x,y) \sim \mu} \left[ \sup_{\tilde{x} \in B_\varepsilon(x)} \ell(\theta, (\tilde{x}, y)) \right]. \quad (\text{AT})$$

[Madry et al 2017]

Compare to unrobust problem:

$$\min_{\theta \in \Theta} \mathbb{E}_{(x,y) \sim \mu} [\ell(\theta, (x, y))].$$

# Formalization of adversarial training problem

How to train classifiers to be robust to (specific) adversarial attacks?:

$$\min_{\theta \in \Theta} \mathbb{E}_{(x,y) \sim \mu} \left[ \sup_{\tilde{x} \in B_\varepsilon(x)} \ell(\theta, (\tilde{x}, y)) \right],$$

or its **distributionally robust optimization (DRO)** version:

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

# Formalization of adversarial training problem

How to train classifiers to be robust to (specific) adversarial attacks?:

$$\min_{\theta \in \Theta} \mathbb{E}_{(x,y) \sim \mu} \left[ \sup_{\tilde{x} \in B_\varepsilon(x)} \ell(\theta, (\tilde{x}, y)) \right],$$

or its distributionally robust optimization (DRO) version:

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

or its explicit penalization version:

$$\inf_{\theta \in \Theta} \sup_{\tilde{\mu}} \mathbb{E}_{\tilde{z} \sim \tilde{\mu}} [\ell(\tilde{z}, \theta)] - C(\mu, \tilde{\mu}).$$



$$\inf_{\theta \in \Theta} \sup_{\tilde{\mu}} \mathbb{E}_{\tilde{z} \sim \tilde{\mu}} [\ell(\tilde{z}, \theta)] - C(\mu, \tilde{\mu}). \quad (\text{AT})$$

- How do we find a solution to this problem?
- Can we find meaningful upper and lower bounds?

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

- How do we find a solution to this problem?
- Can we find meaningful upper and lower bounds?
- How is a problem like (AT) related to regularization methods?

$$\inf_{\theta \in \Theta} \sup_{\tilde{\mu}} \mathbb{E}_{\tilde{z} \sim \tilde{\mu}} [\ell(\tilde{z}, \theta)] - C(\mu, \tilde{\mu}). \quad (\text{AT})$$

- How do we find a solution to this problem?
- Can we find meaningful lower and upper bounds?
- How is a problem like (AT) related to regularization methods?  
**i.e. a problem like:**

$$\inf_{\theta \in \Theta} R(\mu, \theta) + \lambda F(\theta), \quad (\text{Reg})$$

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

- How do we find a solution to this problem?
- Can we find meaningful lower and upper bounds?
- How is a problem like (AT) related to regularization methods?  
**i.e. a problem like:**

$$\inf_{\theta \in \Theta} R(\mu, \theta) + \lambda F(\theta), \quad (\text{Reg})$$

$$\inf_{\theta \in \Theta} \sup_{\tilde{\mu}} \mathbb{E}_{\tilde{z} \sim \tilde{\mu}} [\ell(\tilde{z}, \theta)] - C(\mu, \tilde{\mu}). \quad (\text{AT})$$

What is the **geometry** of:

- Optimal robust classifiers.
- Optimal adversarial attacks.

Instead of the parametric problem

$$\inf_{\theta \in \Theta} \sup_{\tilde{\mu}} \mathbb{E}_{\tilde{z} \sim \tilde{\mu}} [\ell(\tilde{z}, \theta)] - C(\mu, \tilde{\mu}). \quad (1)$$

we'll consider non-parametric problems:

$$\inf_{f \in \mathcal{F}} \sup_{\tilde{\mu}} \mathbb{E}_{\tilde{z} \sim \tilde{\mu}} [\ell(\tilde{z}, f)] - C(\mu, \tilde{\mu}). \quad (2)$$

We'll consider two settings:

- 1 A multilabel classification problem with an agnostic learner.
- 2 A regression problem in a mean field regime.

We'll consider two settings:

- ① A multilabel classification problem with an agnostic learner.
  - Lower bounds for general AT problems.
  - Connections to MOT and (generalized) barycenter problems.
- ② A regression problem in a mean field regime.
  - How to find (approximate) Nash equilibria in mean-field learning settings.

**Overarching goal:** an invitation to look at (AT) from geometric and analytic perspectives.



# 1. A multilabel classification problem with an agnostic learner

# A multilabel classification problem with an agnostic learner

- *Type of data:*  $z = (x, y) \in \mathbb{R}^d \times \{1, \dots, k\}$ ,  $k \geq 2$ .
- *Learner's actions:* measurable  $f = (f_1, \dots, f_k)$  with:  $f_l : \mathbb{R}^d \rightarrow [0, 1]$ , and  $\sum_{l=1}^k f_l = 1$  (**Agnostic learner**).
- *Loss function:*  $\ell(z, f) = \ell((x, y), f) = 1 - f_y(x)$ , i.e. 0-1 loss.

# A multilabel classification problem with an agnostic learner

$$\inf_f \sup_{\tilde{\mu} \in \mathcal{P}(\mathcal{Z})} \left\{ \mathbb{E}_{(\tilde{x}, \tilde{y}) \sim \tilde{\mu}} (\ell(\tilde{z}, f)) - C(\mu, \tilde{\mu}) \right\},$$

where

$$C(\mu, \tilde{\mu}) := \min_{\pi \in \Gamma(\mu, \tilde{\mu})} \int c_{\mathcal{Z}}(z, \tilde{z}) d\pi(z, \tilde{z})$$

for some cost function  $c_{\mathcal{Z}} : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  of the form:

$$c_{\mathcal{Z}}(z, \tilde{z}) = \begin{cases} c(x, \tilde{x}) & \text{if } y = \tilde{y} \\ \infty & \text{if } y \neq \tilde{y}, \end{cases} \quad c : \mathbb{R}^d \times \mathbb{R}^d \rightarrow [0, \infty].$$

# Lower bounds for more general AT problems:

$$\inf_{f \text{ measurable}} \sup_{\tilde{\mu} \in \mathcal{P}(\mathcal{Z})} \left\{ \mathbb{E}_{(\tilde{x}, \tilde{y}) \sim \tilde{\mu}} (\ell(\tilde{z}, f)) - C(\mu, \tilde{\mu}) \right\},$$

is smaller than

$$\inf_{f \in \mathcal{F}'} \sup_{\tilde{\mu} \in \mathcal{P}(\mathcal{Z})} \left\{ \mathbb{E}_{(\tilde{x}, \tilde{y}) \sim \tilde{\mu}} (\ell(\tilde{z}, f)) - C(\mu, \tilde{\mu}) \right\}.$$

# Example of cost function:

$$\inf_f \sup_{\tilde{\mu} \in \mathcal{P}(\mathcal{Z})} \left\{ \mathbb{E}_{(\tilde{x}, \tilde{y}) \sim \tilde{\mu}} (\ell(\tilde{z}, f)) - C(\mu, \tilde{\mu}) \right\},$$

where

$$C(\mu, \tilde{\mu}) := \min_{\pi \in \Gamma(\mu, \tilde{\mu})} \int c_{\mathcal{Z}}(z, \tilde{z}) d\pi(z, \tilde{z})$$

for some cost function  $c_{\mathcal{Z}} : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  of the form:

$$c_{\mathcal{Z}}(z, \tilde{z}) = \begin{cases} c(x, \tilde{x}) & \text{if } y = \tilde{y} \\ \infty & \text{if } y \neq \tilde{y}, \end{cases} \quad c(x, \tilde{x}) = \begin{cases} 0 & \text{if } d(x, \tilde{x}) \leq \varepsilon \\ \infty & \text{if } d(x, \tilde{x}) > \varepsilon \end{cases}$$

# Example of cost function:

When

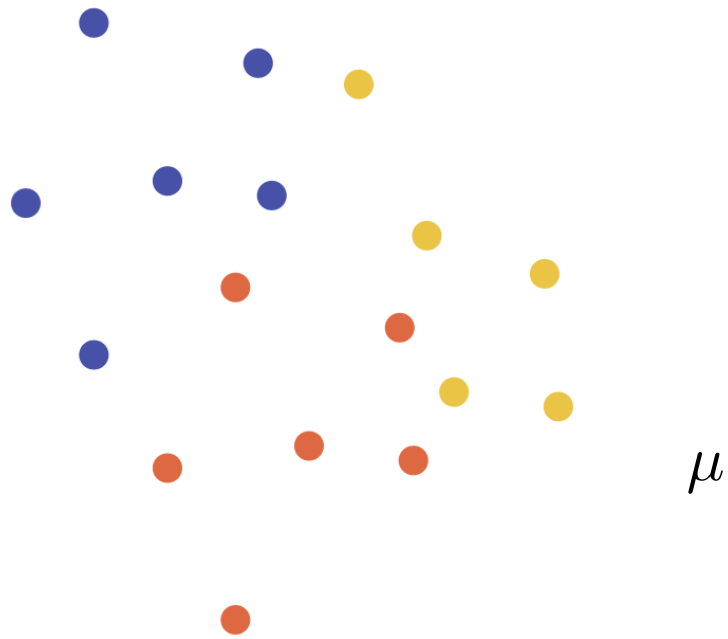
$$c(x, \tilde{x}) = \begin{cases} 0 & \text{if } d(x, \tilde{x}) \leq \varepsilon \\ \infty & \text{if } d(x, \tilde{x}) > \varepsilon \end{cases},$$

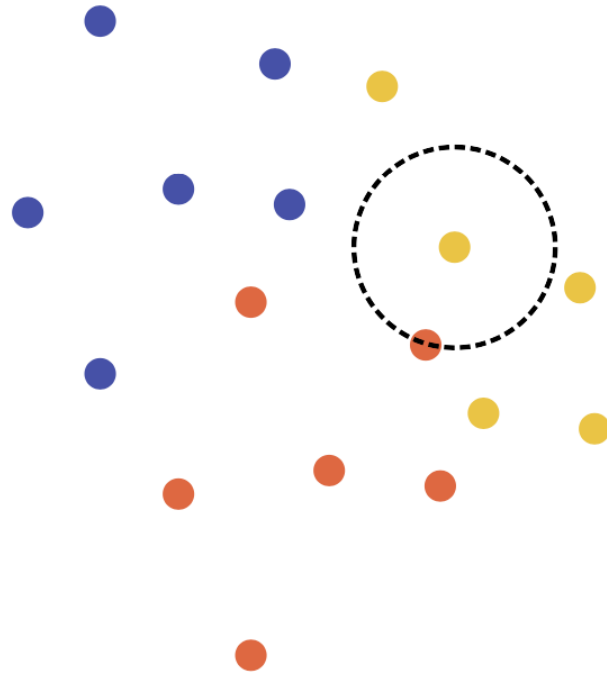
problem

$$\inf_f \sup_{\tilde{\mu} \in \mathcal{P}(\mathcal{Z})} \left\{ \mathbb{E}_{(\tilde{x}, \tilde{y}) \sim \tilde{\mu}} (\ell(\tilde{z}, f)) - C(\mu, \tilde{\mu}) \right\}$$

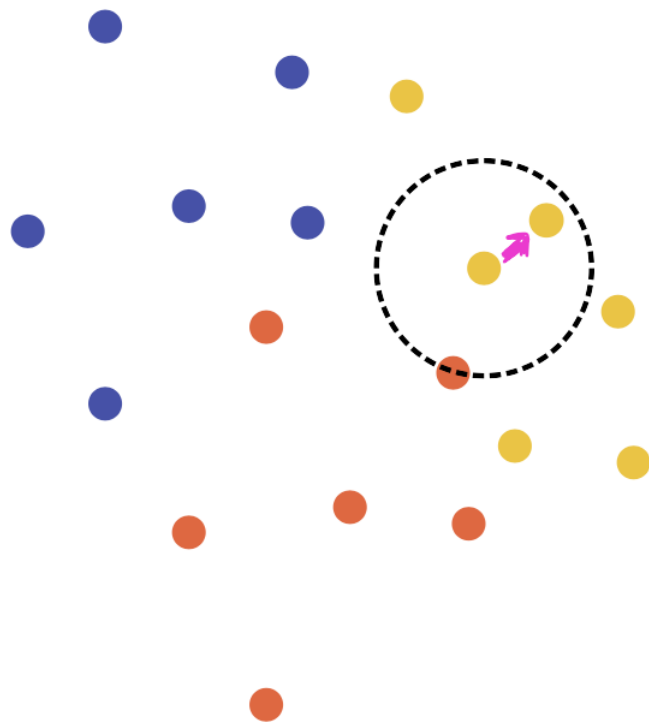
becomes:

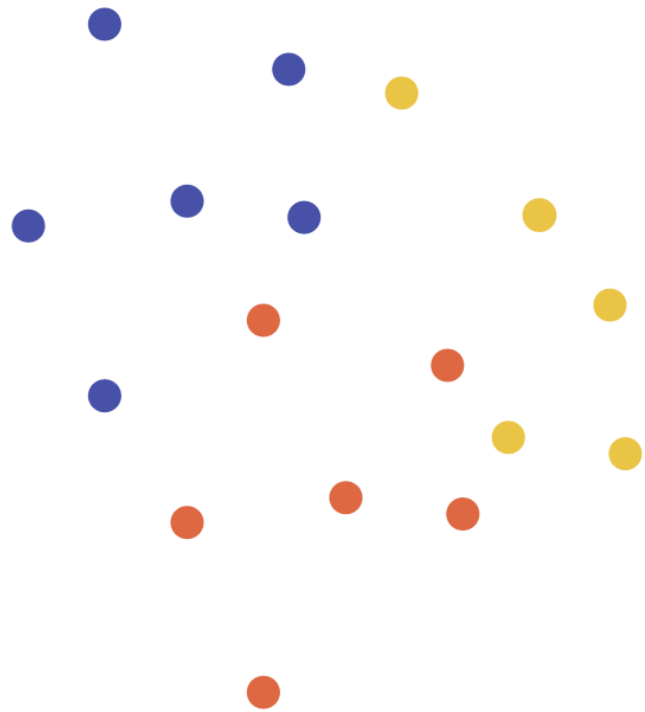
$$\inf_f \mathbb{E}_{(x, y) \sim \mu} \left( \sup_{\tilde{x} \in B_\varepsilon(x)} \ell((\tilde{x}, y), f) \right).$$











$\tilde{\mu}$

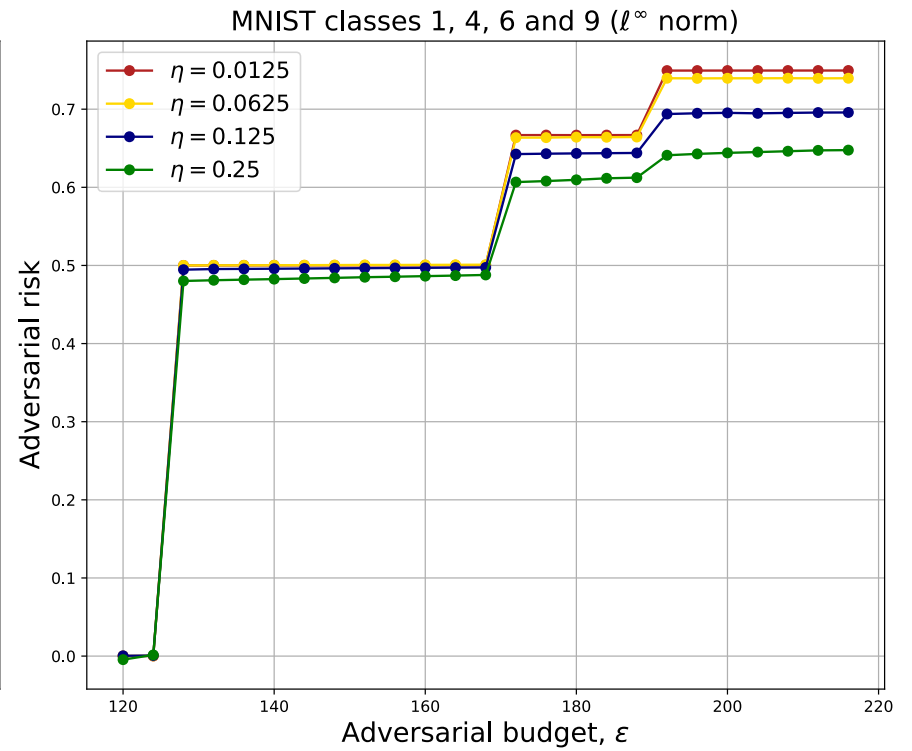
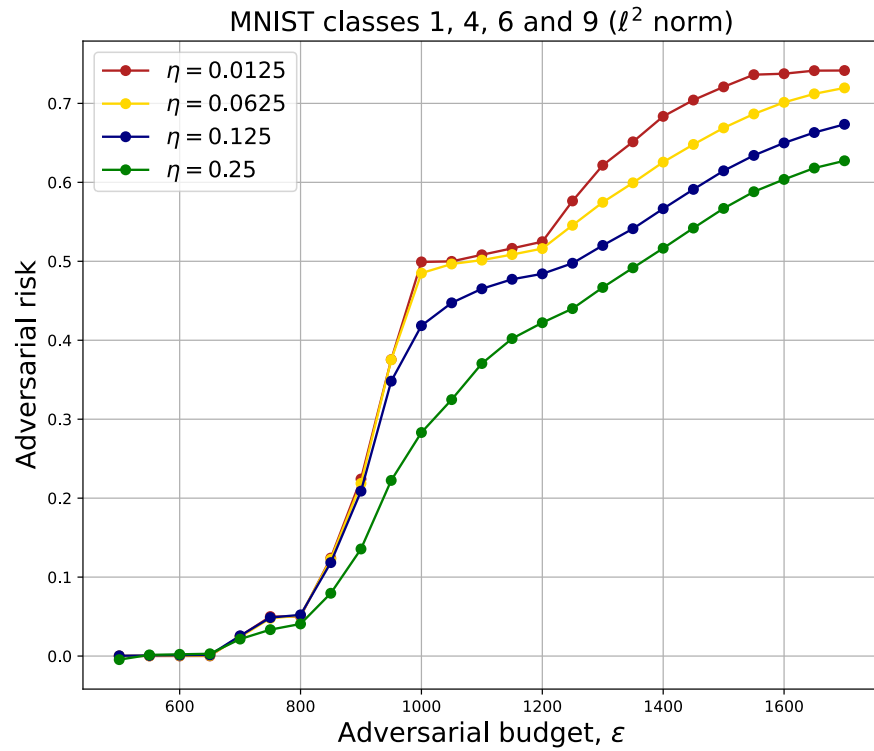
$$\inf_f \sup_{\tilde{\mu} \in \mathcal{P}(\mathcal{Z})} \left\{ \mathbb{E}_{(\tilde{x}, \tilde{y}) \sim \tilde{\mu}} (\ell(\tilde{z}, f)) - C(\mu, \tilde{\mu}) \right\},$$

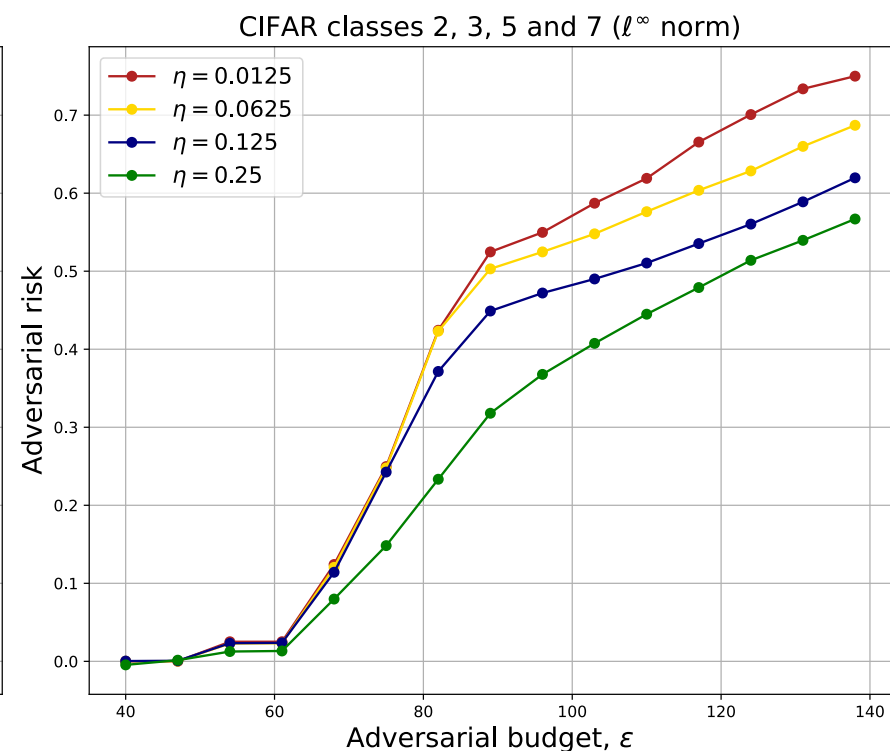
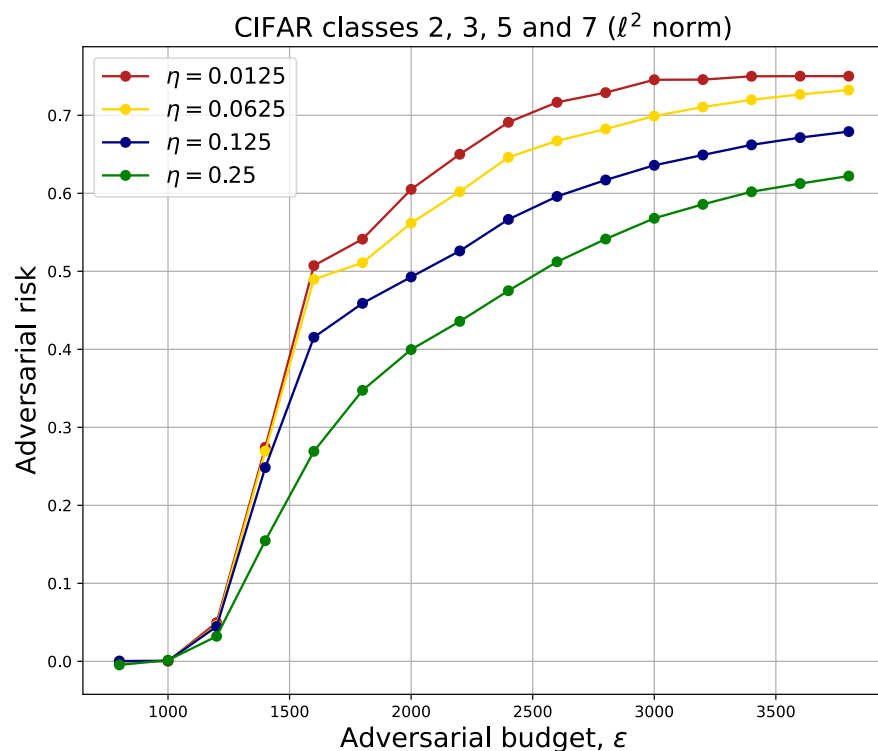
where

$$C(\mu, \tilde{\mu}) := \min_{\pi \in \Gamma(\mu, \tilde{\mu})} \int c_{\mathcal{Z}}(z, \tilde{z}) d\pi(z, \tilde{z})$$

for some cost function  $c_{\mathcal{Z}} : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  of the form:

$$c_{\mathcal{Z}}(z, \tilde{z}) = \begin{cases} c(x, \tilde{x}) & \text{if } y = \tilde{y} \\ \infty & \text{if } y \neq \tilde{y}, \end{cases} \quad c : \mathbb{R}^d \times \mathbb{R}^d \rightarrow [0, \infty].$$





We computed the above using off-the-shelf MOT solvers...

# Multimarginal Optimal Transport (MOT)

$$\inf_{\Gamma(\rho_1, \dots, \rho_K)} \int \mathbf{c}(\xi_1, \dots, \xi_K) d\pi(\xi_1, \dots, \xi_K) \quad (\text{MOT})$$

- Applications in Physics.
- Applications in Economics.
- Machine learning.

# Density functional theory

$$\inf_{\Gamma(\rho_1, \dots, \rho_K)} \int \mathbf{c}(\xi_1, \dots, \xi_K) d\pi(\xi_1, \dots, \xi_K)$$

where

$$\mathbf{c}(\xi_1, \dots, \xi_K) := \sum_{1 \leq i < j \leq K} f(d(x_i, x_j)).$$

**[Seidl 1999], [Gori-Giorgi et al. 2009].**



# Barycenter problems

$$\inf_{\Gamma(\rho_1, \dots, \rho_K)} \int \mathbf{c}(\xi_1, \dots, \xi_K) d\pi(\xi_1, \dots, \xi_K)$$

where

$$\mathbf{c}(\xi_1, \dots, \xi_K) := \inf_{\xi' \in \mathcal{X}} \sum_{i=1}^K c(\xi', \xi_i).$$

# Barycenter problems

$$\inf_{\Gamma(\rho_1, \dots, \rho_K)} \int c(\xi_1, \dots, \xi_K) d\pi(\xi_1, \dots, \xi_K)$$

where

$$c(\xi_1, \dots, \xi_K) := \inf_{\xi' \in \mathcal{X}} \sum_{i=1}^K c(\xi', \xi_i).$$

Equivalent to:

$$\inf_{\rho} \sum_{i=1}^K C(\rho_i, \rho),$$

where

$$C(\rho_i, \rho) := \inf_{\pi \in \Gamma(\rho_i, \rho)} \int c(x, x') d\pi(x, x').$$

**[Agueh and Carlier 2011].**

What is the connection between (AT) and (MOT)?

# What is the connection between (AT) and (MOT)?

- How to find a saddle  $(\tilde{\mu}^*, f^*)$  for the (AT) problem?

# What is the connection between (AT) and (MOT)?

- How to find a saddle  $(\tilde{\mu}^*, f^*)$  for the (AT) problem?  
**Answer:** Solve a certain MOT problem and its dual.

**Theorem [NGT, Jacobs, Kim 22']:** For arbitrary  $k \geq 2$

$$(AT)(\mu) = 1 - \frac{1}{2} \inf_{\pi \in \Pi_k(\mu)} \int_{\mathcal{Z}_*^k} \mathbf{c}(z_1, \dots, z_k) d\pi(z_1, \dots, z_k),$$

for some cost function  $\mathbf{c}$ .

- From  $\pi^*$  can construct  $\tilde{\mu}^*$ .
- $\tilde{\mu}^*$  concentrates on barycenters (w.r.t. cost  $c$ ) of groups of  $k$  or less points in the support of  $\mu_x$ .
- From dual of (MOT) can construct  $f^*$ .

**Theorem [NGT, Jacobs, Kim 22']:** For arbitrary  $k \geq 2$

$$(\text{AT})(\mu) = 1 - \frac{1}{2} \inf_{\pi \in \Pi_k(\mu)} \int_{\mathcal{Z}_*^k} \mathbf{c}(z_1, \dots, z_k) d\pi(z_1, \dots, z_k),$$

for some cost function  $\mathbf{c}$ . A given  $c : \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty]$  induces a  $\mathbf{c}$ .

- From  $\pi^*$  can construct  $\tilde{\mu}^*$ .
- $\tilde{\mu}^*$  concentrates on barycenters (w.r.t. cost  $c$ ) of groups of  $k$  or less points in the support of  $\mu_x$ .
- From dual of (MOT) can construct  $f^*$ .

# Precise MOT problem

Set  $\mathcal{Z}_* := \mathcal{Z} \cup \{\text{ghost}\}$ .

$$\inf_{\pi \in \Pi_K(\mu)} \int_{\mathcal{Z}_*^K} \mathbf{c}(z_1, \dots, z_K) d\pi(z_1, \dots, z_K).$$

- Couplings:

$$\Pi_K(\mu) := \left\{ \pi \in \mathcal{P}(\mathcal{Z}_*^K) : P_{i\#} \pi = \frac{1}{2\mu(\mathcal{Z})} \mu(\cdot \cap \mathcal{Z}) + \frac{1}{2} \delta_{\text{ghost}}, \quad \forall i \right\}.$$

- Cost:

$$\mathbf{c}(z_1, \dots, z_K) := \hat{\mu}_{\vec{z}}(\mathcal{Z}) - AT(\hat{\mu}_{\vec{z}}),$$

where  $\hat{\mu}_{\vec{z}}$  is the positive measure defined as:

$$\hat{\mu}_{\vec{z}} := \frac{1}{K} \sum_{\substack{I \text{ s.t. } z_I \neq \text{ghost}}} \delta_{z_I}.$$



# Toy example

$$\text{Let } c(x, \tilde{x}) = c_\varepsilon(x, \tilde{x}) = \begin{cases} 0 & \text{if } d(x, \tilde{x}) \leq \varepsilon \\ +\infty & \text{otherwise} \end{cases}$$

$$\mu = \omega_1 \delta_{(x_1,1)} + \omega_2 \delta_{(x_2,2)} + \omega_3 \delta_{(x_3,3)}$$

$x_2$   


$x_1$   


$x_3$   


# Case 1:

$x_2$



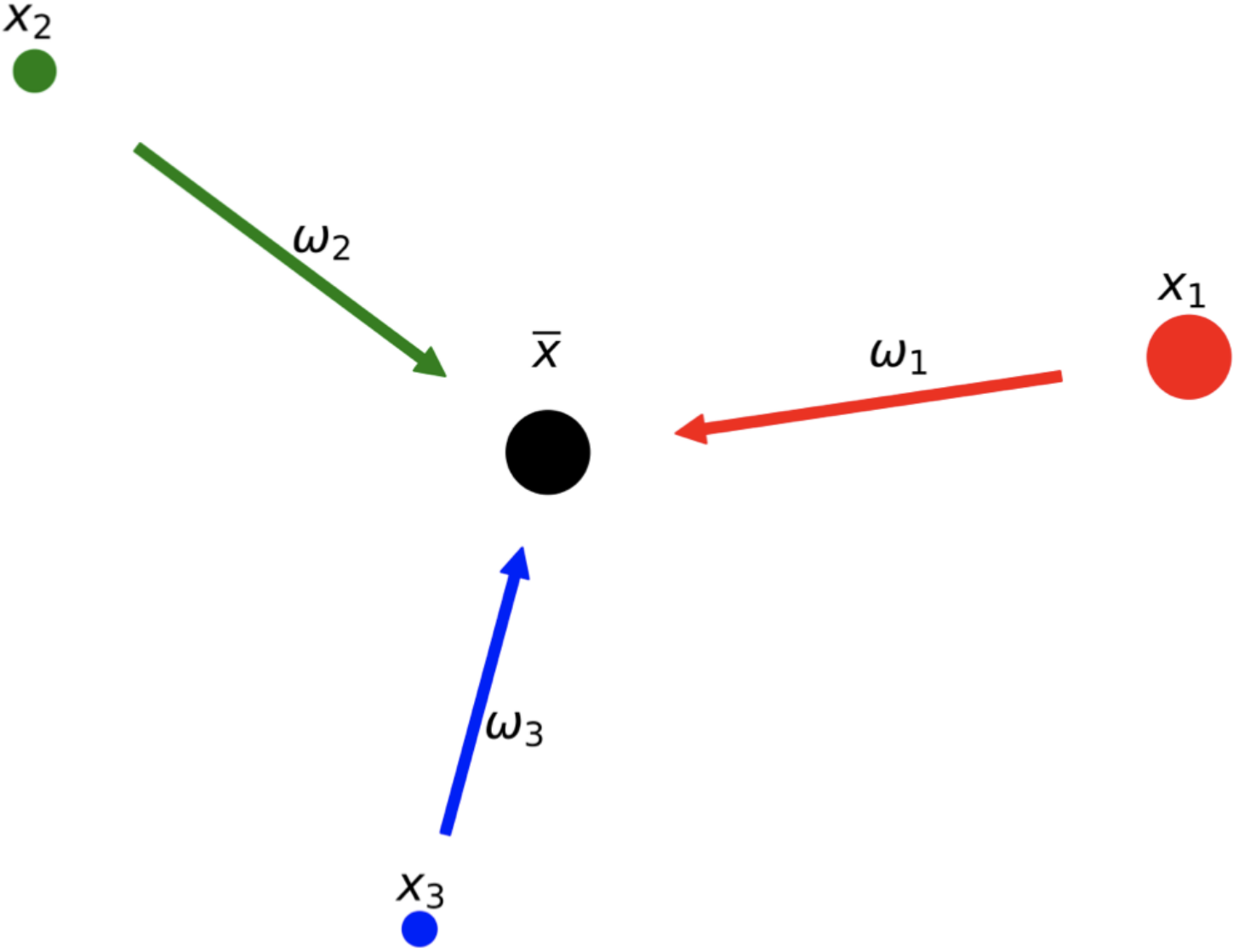
$x_1$



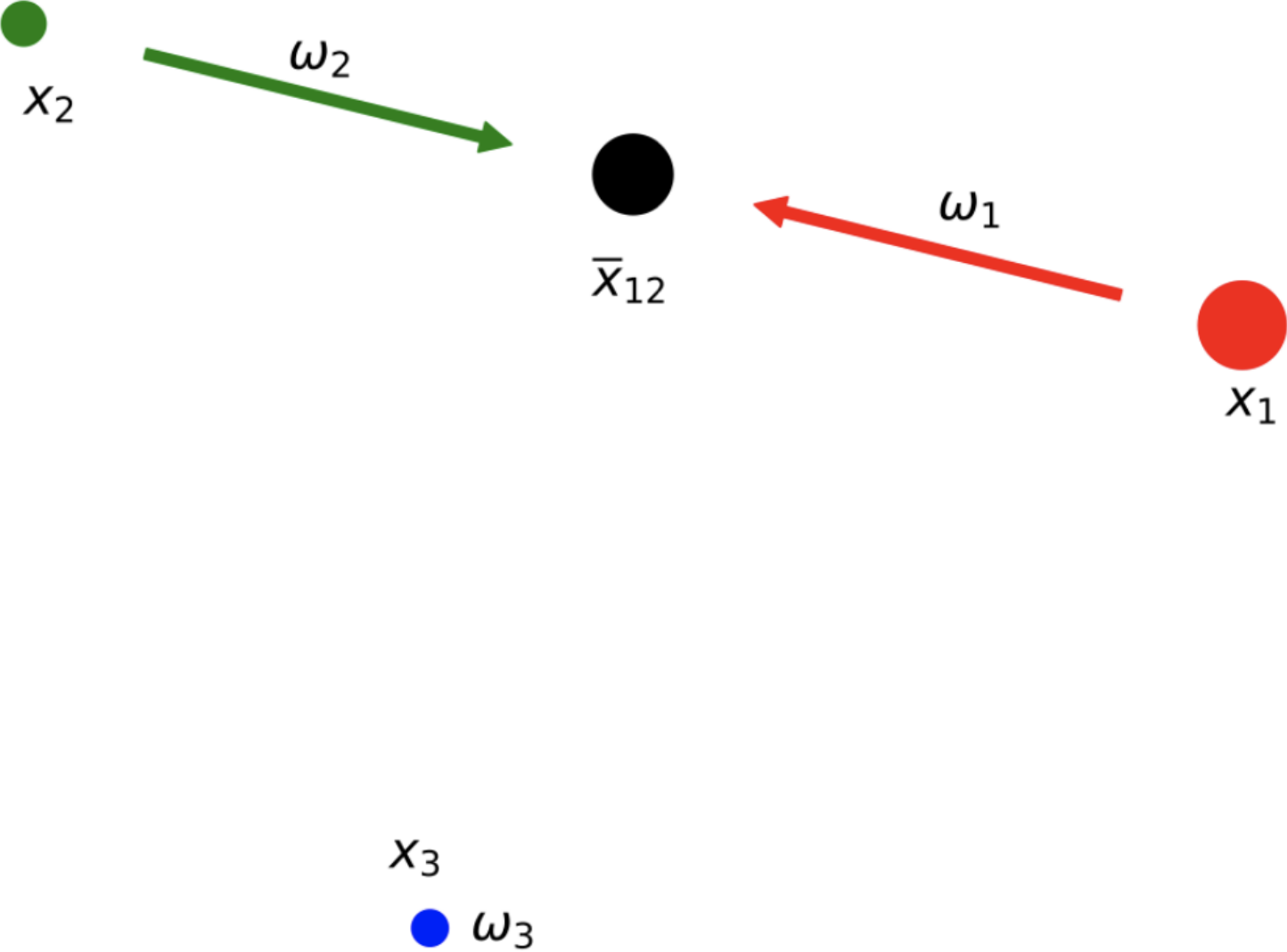
$x_3$



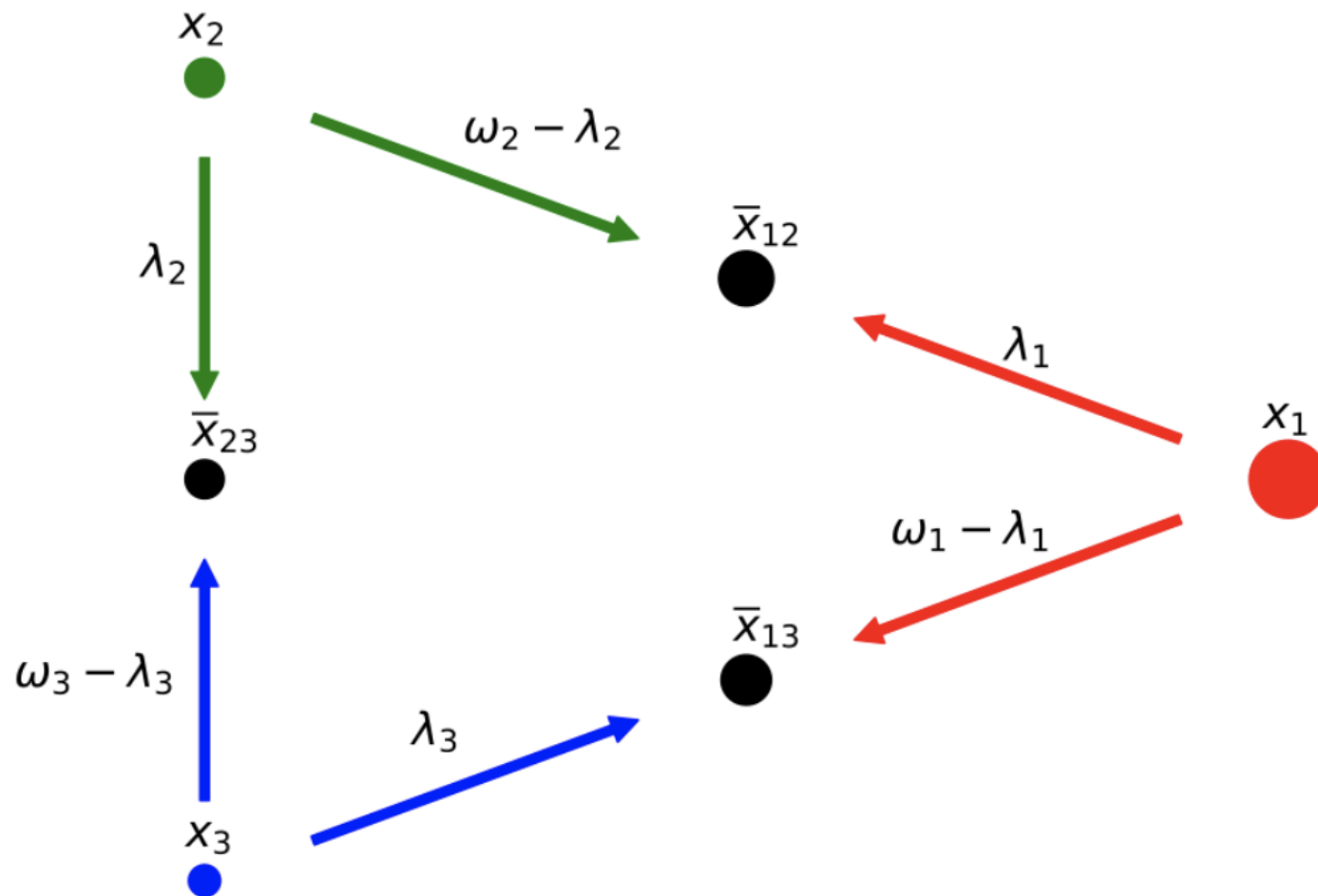
Case 2:



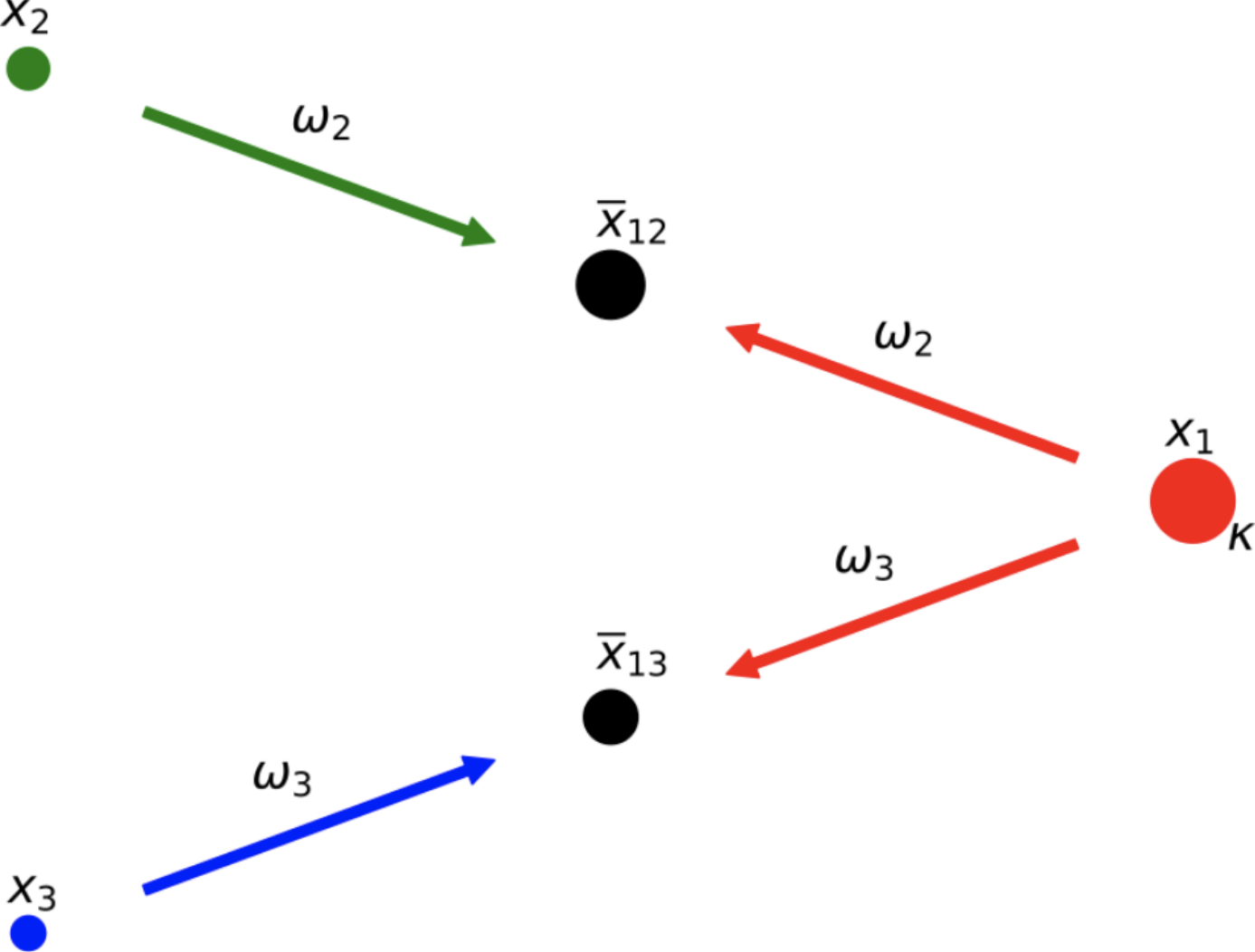
# Case 3:

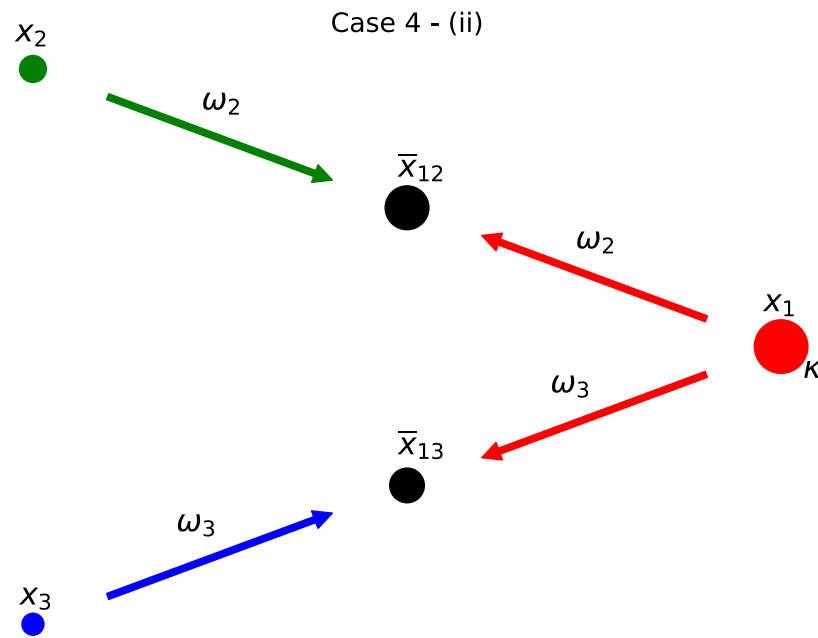
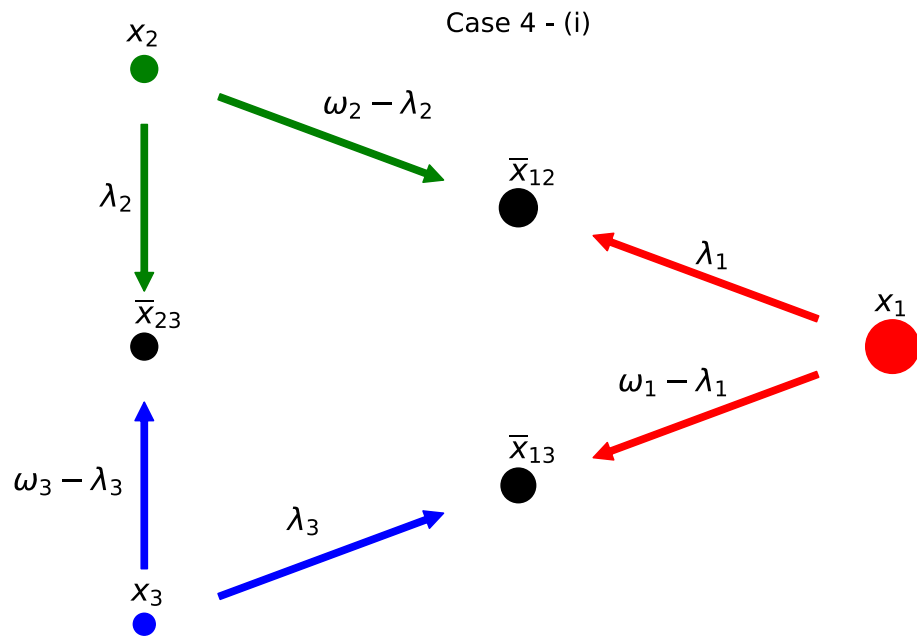
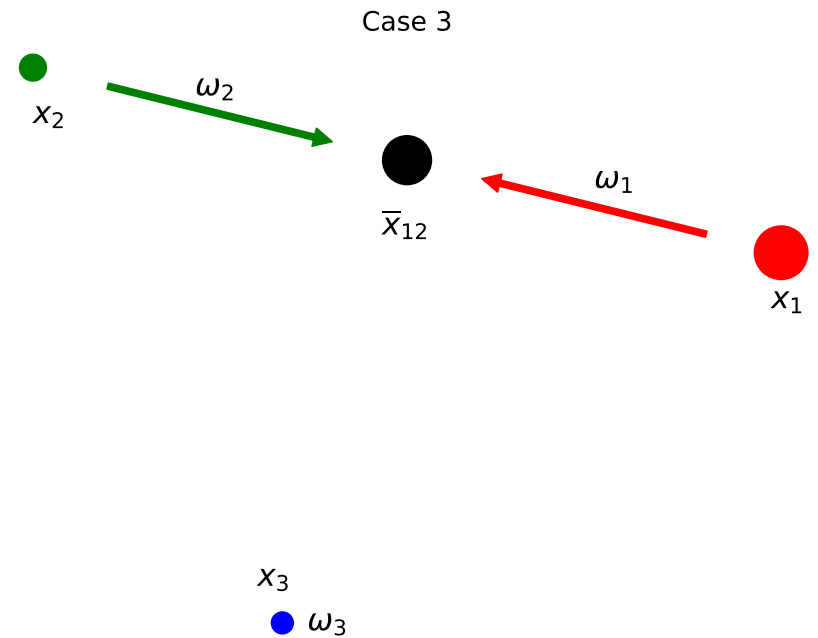
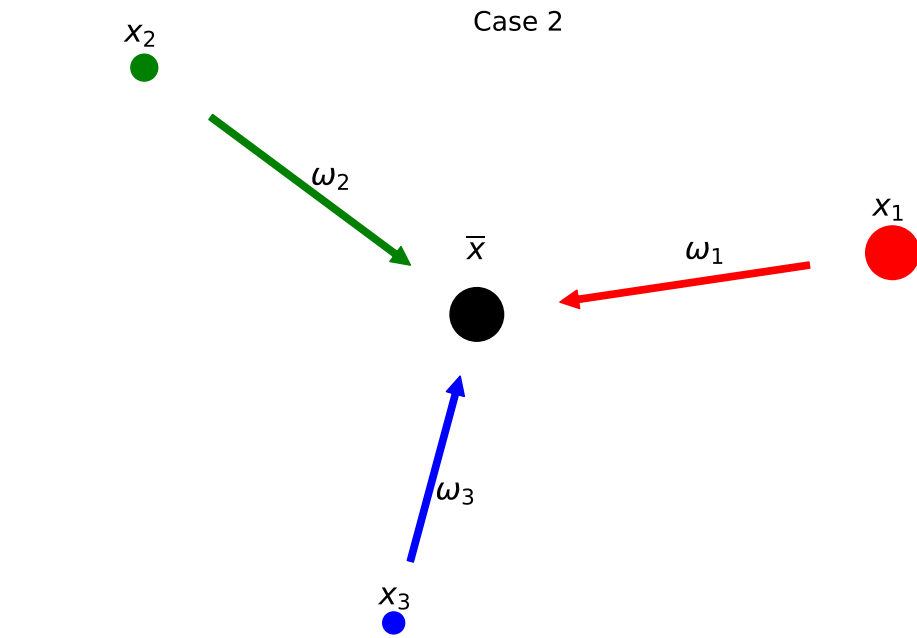


# Case 4 i:



Case 4 ii:





$$(AT)(\mu) = 1 - \frac{1}{2} \inf_{\pi \in \Pi_k(\mu)} \int_{\mathcal{Z}_*^K} \mathbf{c}(z_1, \dots, z_K) d\pi(z_1, \dots, z_K),$$

for cost function  $\mathbf{c}$ :

$$\mathbf{c}(z_1, \dots, z_K) := \hat{\mu}_{\vec{z}}(\mathcal{Z}) - AT(\hat{\mu}_{\vec{z}}),$$

where  $\hat{\mu}_{\vec{z}}$  is the positive measure defined as:

$$\hat{\mu}_{\vec{z}} := \frac{1}{K} \sum_{l \text{ s.t. } z_l \neq \emptyset}^K \delta_{z_l}.$$



## Theorem (NGT, Jacobs, Kim, 2022)

Suppose that  $(\pi^*, \phi^*)$  is a solution pair for the MOT problem and its dual. Define  $f^*$  and  $\tilde{\mu}^*$  according to:

$$f_i^* := \left( \max \left\{ \sum_{j=1}^K \phi_j^*(\cdot, i) + \sum_{j=1}^K \phi_j^*(\emptyset), 0 \right\} \right)^{\bar{c}}$$

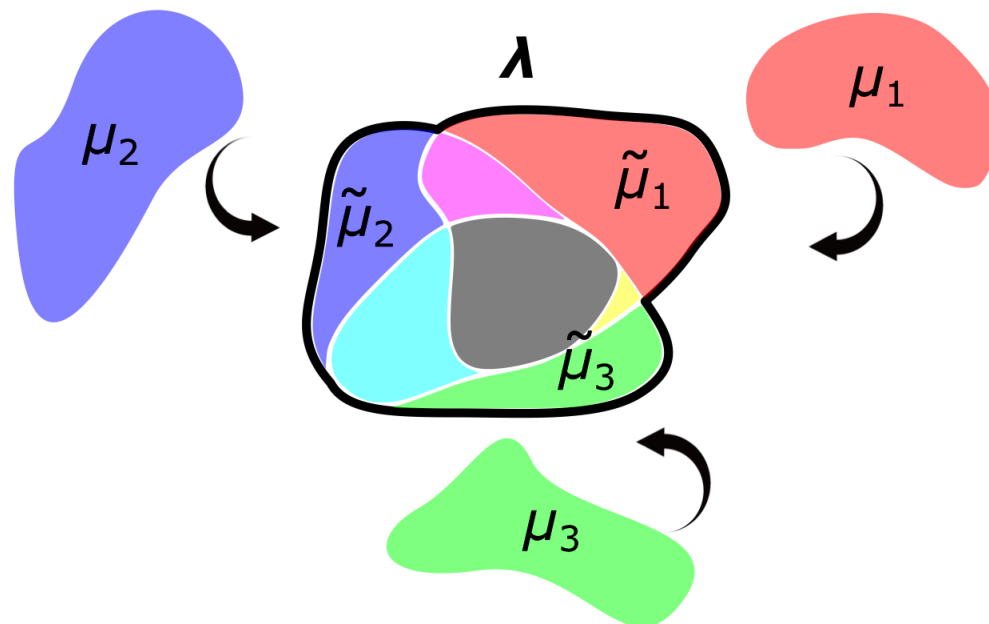
and for any test function  $h$  on  $\mathcal{X}$ ,

$$\int_{\mathcal{X}} h(\tilde{x}) d\tilde{\mu}_i^*(\tilde{x}) := \int_{\mathcal{Z}_*^K} \left\{ \int_{\mathcal{X}} h(\tilde{x}) d\tilde{\mu}_{\vec{z}, i}^*(\tilde{x}) \right\} d\pi^*(\vec{z}),$$

where  $\tilde{\mu}_{\vec{z}, i}^*$  is the  $i$ -th marginal of  $\tilde{\mu}_{\vec{z}}^*$ , an optimal adversarial attack which achieves  $\mathbf{c}(z_1, \dots, z_K)$  given  $\vec{z} = (z_1, \dots, z_K)$ . Then  $(f^*, \tilde{\mu}^*)$  is a saddle for problem (AT).

# Generalized barycenter problems

$$\inf_{\lambda, \tilde{\mu}_1, \dots, \tilde{\mu}_K} \lambda(\mathcal{X}) + \sum_{i=1}^K C(\mu_i, \tilde{\mu}_i) \quad \text{s.t. } \lambda \geq \tilde{\mu}_i \quad \forall i = 1, \dots, K.$$



# Generalized barycenter problems

$$\inf_{\lambda, \tilde{\mu}_1, \dots, \tilde{\mu}_K} \lambda(\mathcal{X}) + \sum_{i=1}^K C(\mu_i, \tilde{\mu}_i) \quad \text{s.t. } \lambda \geq \tilde{\mu}_i \quad \forall i = 1, \dots, K.$$



$(AT)(\mu)$

# Generalized barycenter problems

$$\inf_{\lambda, \tilde{\mu}_1, \dots, \tilde{\mu}_K} \lambda(\mathcal{X}) + \sum_{i=1}^K C(\mu_i, \tilde{\mu}_i) \quad \text{s.t. } \lambda \geq \tilde{\mu}_i \forall i = 1, \dots, K.$$



**(MOT)**

[NGT, Jacobs, Kim 22']

$$\inf_{\rho} \sum_{i=1}^K C(\rho_i, \rho)$$



**(MOT)**

[Agueh and Carlier 2011].

$$(\text{AT})(\mu) = 1 - \frac{1}{2} \inf_{\pi \in \Pi_k(\mu)} \int_{\mathcal{Z}_*^k} \mathbf{c}(z_1, \dots, z_k) d\pi(z_1, \dots, z_k),$$

- Equivalence between (AT) and computational OT!
- Geometric description of optimal adversarial attacks!
- Specific OT algorithms for this problem?
- Generalizations to other loss functions?
- In binary case (i.e.,  $k = 2$ ): [Baghoji, Cullina, Mittal 19'], [Pydi, Jog 20'], [NGT and Murray 20'].

## 2. A regression problem in a mean field regime

# A mean field model of NNs

- $z = (x, y)$ ,  $x \in \mathbb{R}^d$  and  $y \in \mathbb{R}$ .
- $f(x) = f_\nu(x) := \int_{\Theta} ah(b \cdot x) d\nu(a, b)$ , where:  
 $\theta = (a, b) \in \Theta$ ,  $\nu \in \mathcal{P}(\Theta)$ ;  $h$  is non-linearity.
- $\ell(\tilde{z}, f_\nu) := (f_\nu(\tilde{x}) - \tilde{y})^2$ .

- $z = (x, y)$ ,  $x \in \mathbb{R}^d$  and  $y \in \mathbb{R}$ .
- $f(x) = f_\nu(x) := \int_{\Theta} ah(b \cdot x) d\nu(a, b)$ , where:  
 $\theta = (a, b) \in \Theta$ ,  $\nu \in \mathcal{P}(\Theta)$ ;  $h$  is non-linearity.
- $\ell(\tilde{z}, f_\nu) := (f_\nu(\tilde{x}) - \tilde{y})^2$ .
- *(AT) problem:*

$$\inf_{\nu \in \mathcal{P}(\Theta)} \sup_{\tilde{\mu} \in \mathcal{P}(\mathcal{Z})} \left\{ \mathbb{E}_{(\tilde{x}, \tilde{y}) \sim \tilde{\mu}} (\ell(\tilde{z}, f_\nu)) - c_a W_2^2(\mu, \tilde{\mu}) \right\}.$$



Now, problem:

$$\inf_{\nu \in \mathcal{P}(\Theta)} \sup_{\tilde{\mu} \in \mathcal{P}(\mathcal{Z})} \left\{ \mathbb{E}_{(\tilde{x}, \tilde{y}) \sim \tilde{\mu}} (\ell(\tilde{z}, f_\nu)) - c_a W_2^2(\mu, \tilde{\mu}) \right\}$$

is equivalent to

$$\min_{\nu \in \mathcal{P}(\Theta)} \max_{\pi \in \mathcal{P}(\mathcal{Z} \times \mathcal{Z}) \text{ s.t. } \pi_z = \mu} \mathcal{U}(\pi, \nu),$$

where

$$\mathcal{U}(\pi, \nu) := \int_{\mathcal{Z}} \int_{\mathcal{Z}} (f_\nu(\tilde{x}) - \tilde{y})^2 d\pi(z, \tilde{z}) - c_a \int_{\mathcal{Z} \times \mathcal{Z}} |z - \tilde{z}|^2 d\pi(z, \tilde{z}).$$

Target:

$$\min_{\nu \in \mathcal{P}(\Theta)} \max_{\pi \in \mathcal{P}(\mathcal{Z} \times \mathcal{Z}) \text{ s.t. } \pi_Z = \mu} \mathcal{U}(\pi, \nu),$$

Ascent-Descent in spaces of measures:

$$\begin{cases} \partial_t \pi_t &= -\eta_t \operatorname{div}_{z, \tilde{z}}(\pi_t(0, \nabla_{\tilde{z}} \mathcal{U}_\pi)) \\ &+ \kappa_t \pi_t (\mathcal{U}_\pi(z, \tilde{z}) - \int \mathcal{U}_\pi(z, \tilde{z}') d\pi_t(\tilde{z}'|z)) \\ \partial_t \nu_t &= \eta_t \operatorname{div}_\theta(\nu_t \nabla_\theta \mathcal{U}_\nu(\theta)) - \kappa_t \nu_t (\mathcal{U}_\nu(\theta) - \int \mathcal{U}_\nu(\theta') d\nu_t(\theta')), \end{cases}$$

where  $\mathcal{U}_\pi, \mathcal{U}_\nu$  first variations of  $\mathcal{U}$  w.r.t.  $\pi, \nu$ , respectively.

Target:

$$\min_{\nu \in \mathcal{P}(\Theta)} \max_{\pi \in \mathcal{P}(\mathcal{Z} \times \tilde{\mathcal{Z}}) \text{ s.t. } \pi_z = \mu} \mathcal{U}(\pi, \nu),$$

Ascent-Descent in spaces of measures (precisely, projected ascent-descent w.r.t. Wasserstein-Fisher-Rao metric):

$$\begin{cases} \partial_t \pi_t &= -\eta_t \operatorname{div}_{z, \tilde{z}}(\pi_t(0, \nabla_{\tilde{z}} \mathcal{U}_\pi)) \\ &+ \kappa_t \pi_t (\mathcal{U}_\pi(z, \tilde{z}) - \int \mathcal{U}_\pi(z, \tilde{z}') d\pi_t(\tilde{z}'|z)) \\ \partial_t \nu_t &= \eta_t \operatorname{div}_\theta(\nu_t \nabla_\theta \mathcal{U}_\nu(\theta)) - \kappa_t \nu_t (\mathcal{U}_\nu(\theta) - \int \mathcal{U}_\nu(\theta') d\nu_t(\theta')), \end{cases}$$

where  $\mathcal{U}_\pi, \mathcal{U}_\nu$  first variations of  $\mathcal{U}$  w.r.t.  $\pi, \nu$ , respectively.

Particle system approximation:

$$\pi_t^N = \frac{1}{N} \sum_{i=1}^N \omega_t^i \delta_{(Z_t^i, \tilde{Z}_t^i)}, \quad \nu_t^N = \frac{1}{N} \sum_{i=1}^N \alpha_t^i \delta_{\theta_t^i},$$

where:

$$d_t(Z_t^i, \tilde{Z}_t^i) = (0, \eta_t \nabla_{\tilde{z}} \mathcal{U}_\pi(\pi_t^N, \nu_t^N; Z_t^i, \tilde{Z}_t^i))$$

$$d_t \omega_t^i = \kappa_t \omega_t^i \left( \mathcal{U}_\pi(\pi_t^N, \nu_t^N; Z_t^i, \tilde{Z}_t^i) - \int \mathcal{U}_\pi(\pi_t^N, \nu_t^N; Z_t^i, \tilde{z}') d\pi_t^N(\tilde{z}' | Z_t^i) \right)$$

$$d_t \theta_t^i = -\eta_t \nabla_{\theta} \mathcal{U}_\nu(\pi_t^N, \nu_t^N; \theta_t^i)$$

$$d_t \alpha_t^i = -\kappa_t \alpha_t^i \left( \mathcal{U}_\nu(\pi_t^N, \nu_t^N; \theta_t^i) - \int \mathcal{U}_\nu(\pi_t^N, \nu_t^N; \theta') d\nu_t^N(\theta') \right);$$

and given initial condition  $(Z_0^i, \tilde{Z}_0^i, \omega_0^i, \vartheta_0^i, \alpha_0^i)$  (possibly random).

# Part 1: Mean field limit of particle system

## Theorem (C.A. García Trillos, NGT 23')

*Suppose that:*

- $\Theta, \mathcal{Z}$  are bounded subsets of Euclidean space.
- $\nabla U_\pi, \nabla U_\nu$  are Lipschitz.
- Initial conditions  $(Z_0^i, \tilde{Z}_0^i, \omega_0^i, \theta_0^i, \alpha_0^i)$  are well prepared.

*Then, for every fixed  $T > 0$ , we have:*

$$\sup_{t \in [0, T]} W_1(\pi_t^N, \pi_t) \rightarrow 0; \quad \sup_{t \in [0, T]} W_1(\nu_t^N, \nu_t) \rightarrow 0,$$

*as  $N \rightarrow \infty$ , where  $(\pi_t, \nu_t)$  solve Ascent-Descent dynamics PDE.*

Both  $(\pi_t^N, \nu_t^N)$  and  $(\pi_t, \nu_t)$  solve the same equation:

$$\begin{cases} \partial_t \pi_t &= -\eta_t \operatorname{div}_{z, \tilde{z}}(\pi_t(0, \nabla_{\tilde{z}} \mathcal{U}_\pi)) \\ &+ \kappa_t \pi_t (\mathcal{U}_\pi(z, \tilde{z}) - \int \mathcal{U}_\pi(z, \tilde{z}') d\pi_t(\tilde{z}'|z)) \\ \partial_t \nu_t &= \eta_t \operatorname{div}_\theta(\nu_t \nabla_\theta \mathcal{U}_\nu(\theta)) - \kappa_t \nu_t (\mathcal{U}_\nu(\theta) - \int \mathcal{U}_\nu(\theta') d\nu_t(\theta')) , \end{cases}$$

but they differ in their initial conditions  $(\pi_0^N, \nu_0^N)$  and  $(\pi_0, \nu_0)$ .

# Mean field limit of particle system

## Theorem (C.A. García Trillos, NGT 23')

*Suppose that:*

- $\Theta, \mathcal{Z}$  are bounded subsets of Euclidean space.
- $\nabla U_\pi, \nabla U_\nu$  are Lipschitz.
- *Initial conditions  $(Z_0^i, \tilde{Z}_0^i, \omega_0^i, \theta_0^i, \alpha_0^i)$  are well prepared.*

*Then, for every fixed  $T > 0$ , we have:*

$$\sup_{t \in [0, T]} W_1(\pi_t^N, \pi_t) \rightarrow 0; \quad \sup_{t \in [0, T]} W_1(\nu_t^N, \nu_t) \rightarrow 0,$$

*as  $N \rightarrow \infty$ , where  $(\pi_t, \nu_t)$  solve Ascent-Descent dynamics PDE.*

# An example of well prepared initial conditions

Set  $\omega_0^i = \alpha_0^i = 1$  and suppose that, as  $N \rightarrow \infty$ , we have:

$$W_1(\nu_0^N, \nu_0) \rightarrow 0,$$

as well as

$$\inf_{\nu \in \Gamma_{\text{opt}}(\pi_{0,z}^N, \pi_{0,z})} \int W_1(\pi_0^N(\cdot|z'_0), \pi_0(\cdot|z_0)) d\nu(z'_0, z_0) \rightarrow 0$$

(Knothe transport and reminiscent to TLp metric).



# An example of well prepared initial conditions

To satisfy:

$$\inf_{\nu \in \Gamma_{\text{Opt}}(\pi_0^N, \pi_{0,z})} \int W_1(\pi_0^N(\cdot|z'_0), \pi_0(\cdot|z_0)) d\nu(z'_0, z_0) \rightarrow 0,$$

set, for example,

$$\pi_0^N = \frac{1}{nm} \sum_{ij} \delta_{(Z_0^i, \tilde{Z}_0^{ij})},$$

where

- $Z_0^i \sim \pi_{0,z} = \mu, i = 1, \dots, n,$
- $\tilde{Z}_0^{ij} \sim \pi_0(\cdot|Z_0^i), j = 1, \dots, m, i = 1, \dots, n.$

# An example of well prepared initial conditions

## Lemma (C.A. García Trillos, NGT 23')

Let  $A, B$  be two bounded Borel subsets of  $\mathbb{R}^d$  and  $\mathbb{R}^{d'}$ , respectively. Let  $\mu \in \mathcal{P}(A)$ , and let  $u \in A \mapsto \mu_u(\cdot) \in \mathcal{P}(B)$  be a measurable map.

Then, for every sequence  $\{\Upsilon_n\}_{n \in \mathbb{N}} \subseteq \Gamma(\mu, \mu)$  satisfying

$$\lim_{n \rightarrow \infty} \int_{A \times A} |u - u'| d\Upsilon_n(u, u') = 0,$$

we have

$$\lim_{n \rightarrow \infty} \int_{A \times A} W_1(\mu_u, \mu_{u'}) d\Upsilon_n(u, u') = 0.$$

## Part 2: Long time behavior mean field system

### Theorem (C.A. García Trillos, NGT 23')

Fix  $\delta > 0$ . Let  $\pi, \nu$  the solution to descent-ascent dynamics for  $\eta_t, \kappa_t$  appropriately tuned. Define:

$$\bar{\nu}_t = \frac{1}{t} \int_0^t \nu_s ds, \quad \bar{\pi}_t = \frac{1}{t} \int_0^t \pi_s ds.$$

Then, for all large enough  $t$ ,

$$\sup_{\tilde{\pi} \text{ s.t. } \tilde{\pi}_z = \mu} \mathcal{U}(\tilde{\pi}, \bar{\nu}(t)) - \inf_{\tilde{\nu}} \mathcal{U}(\bar{\pi}(t), \tilde{\nu}) \leq \delta.$$

# Long time behavior mean field system

## Theorem (C.A. García Trillos, NGT 23')

Fix  $\delta > 0$ . Let  $\pi, \nu$  the solution to descent-ascent dynamics for  $\eta_t, \kappa_t$  appropriately tuned. Define:

$$\bar{\nu}_t = \frac{1}{t} \int_0^t \nu_s ds, \quad \bar{\pi}_t = \frac{1}{t} \int_0^t \pi_s ds.$$

Then, for all large enough  $t$ ,

$$\sup_{\tilde{\pi} \text{ s.t. } \tilde{\pi}_z = \mu} \mathcal{U}(\tilde{\pi}, \bar{\nu}(t)) - \inf_{\tilde{\nu}} \mathcal{U}(\bar{\pi}(t), \tilde{\nu}) \leq \delta.$$

However, this is under very stringent conditions on initializations (both  $\pi_0, \nu_0$ ). On  $\nu_0$ , these conditions are not so different to those in **Chizat and Bach 17'**, for example.

# The “strongly concave” case

However, roles of  $\pi$  and  $\nu$  are quite different. In the setting:

$$\mathcal{U}(\pi, \nu) = \int_{\mathcal{Z}} \int_{\tilde{\mathcal{Z}}} (f_{\nu}(\tilde{x}) - \tilde{y})^2 d\pi(z, \tilde{z}) - c_a \int_{\mathcal{Z} \times \tilde{\mathcal{Z}}} |z - \tilde{z}|^2 d\pi(z, \tilde{z}),$$

# The “strongly concave” case

However, roles of  $\pi$  and  $\nu$  are quite different. In the setting:

$$\mathcal{U}(\pi, \nu) = \int_{\mathcal{Z}} \int_{\tilde{\mathcal{Z}}} (f_{\nu}(\tilde{x}) - \tilde{y})^2 d\pi(z, \tilde{z}) - c_a \int_{\mathcal{Z} \times \tilde{\mathcal{Z}}} |z - \tilde{z}|^2 d\pi(z, \tilde{z}),$$

if  $c_a$  sufficiently large, then there exists  $\lambda > 0$  such that  
 $\forall \nu \in \mathcal{P}(\Theta), \forall \pi \in \mathcal{P}(\mathcal{Z}^2)$  with  $\pi_z = \mu$ :

$$\int |\nabla_{\tilde{z}} \mathcal{U}_{\pi}(\pi, \nu; z, \tilde{z})|^2 d\pi(z, \tilde{z}) \geq \lambda(m_{\nu}^* - \mathcal{U}(\pi, \nu)), \quad (\text{PL})$$

where  $m_{\nu}^* := \sup_{\tilde{\pi} \text{ s.t. } \tilde{\pi}_z = \mu} \mathcal{U}(\tilde{\pi}, \nu)$ .

## Theorem (C.A. García Trillos, NGT 23')

Fix  $\delta > 0$ . Suppose PL assumption holds. Let  $\pi, \nu$  the solution to (slightly modified) descent-ascent dynamics for  $\eta_t, \kappa_t$  appropriately tuned, and with  $\nu_0$  appropriately initialized and  $\pi_0$  arbitrary.

Define:

$$\bar{\nu}_t = \frac{1}{t} \int_0^t \nu_s ds, \quad \bar{\pi}_t = \frac{1}{t} \int_0^t \pi_s ds.$$

Then, for all large enough  $t$ ,

$$\sup_{\tilde{\pi} \text{ s.t. } \tilde{\pi}_z = \mu} \mathcal{U}(\tilde{\pi}, \bar{\nu}(t)) - \inf_{\tilde{\nu}} \mathcal{U}(\bar{\pi}(t), \tilde{\nu}) \leq \delta.$$

Related work: "Certifying Some Distributional Robustness with Principled Adversarial Training" **Sinha, Namkoong, and Duchi 18'**

$$\min_{\nu \in \mathcal{P}(\Theta)} \max_{\pi \in \mathcal{P}(\mathcal{Z} \times \mathcal{Z}) \text{ s.t. } \pi_{\mathcal{Z}} = \mu} \mathcal{U}(\pi, \nu),$$

- Ascent-descent algorithms.
- Some convergence results.
- Less stringent assumptions.
- Other geometries modeling adversarial costs?
- Other related work:
  - 1 "A mean-field analysis of two-player zero-sum games" **Domingo-Enrich et al 20'**.
  - 2 "An Exponentially Converging Particle Method for the Mixed Nash Equilibrium of Continuous Games" **Chizat and Wang 22'**.



# An analyst's perspective on adversarial training:

- NGT and R. Murray "Adversarial classification: necessary conditions and geometric flows" *Journal of Machine Learning research (JMLR)* 22'.
- C. García Trillos, NGT "On the regularized risk of distributionally robust learning over deep neural networks" *Research in the Mathematical Sciences (RMS)* 22'.
- L. Bungert, NGT, R. Murray "The Geometry of Adversarial Training in Binary Classification" *To appear in Information and Inference: A Journal of the IMA*.
- NGT, M. Jacobs, J. Kim "The multimarginal optimal transport formulation of adversarial multiclass classification" *To appear in JMLR*.
- C. García Trillos, NGT "On adversarial robustness and the use of Wasserstein ascent-descent dynamics to enforce it" <https://arxiv.org/abs/2301.03662> 23'.

# Thank you for your attention!

**Special thanks to:**

- NSF Grants: DMS-2005797 and DMS-2236447
- All my collaborators.

