

Optimal Transport and Deep Learning: Learning from one another

PhD defense

Kilian Fatras

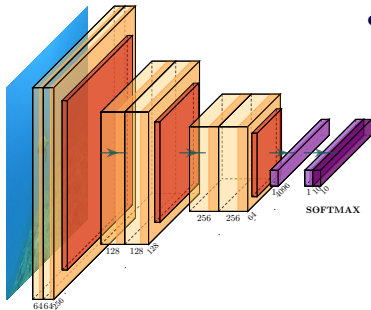
25th November 2021

Obelix and Panama Teams

Supervised by Nicolas Courty and Rémi Flamary

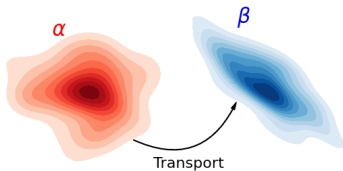


Introduction on deep learning and optimal transport



- Introduction on deep learning

- Neural networks
- Applications
- Probability distributions



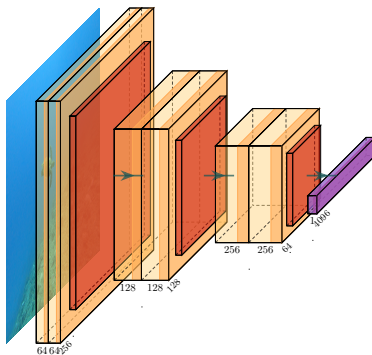
- Introduction on optimal transport

- Definitions
- Properties
- Entropic variant

Neural network illustration

Deep learning is a tool to estimate non-linear complex functions

- Neural networks: many stacked layers and each layer is made of neurons
- Parameters of neural networks: connections between layers
- Different layers: convolutional layers, fully connected layers, ...



Motivating example: Classification

Classification problem: predicting the class of a given image

Motivating example: Classification

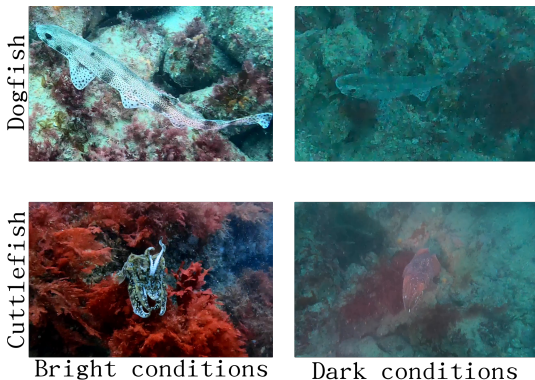
- Find a function f_θ which describes the relationship between the space of images and the space of classes
- f_θ is a **neural network** !

$$f_\theta \left(\begin{array}{c} \text{Image of a clownfish} \end{array} \right) = \begin{pmatrix} 0.1 \\ 0.2 \\ 0.7 \end{pmatrix} \begin{array}{l} \text{clown fish} \\ \text{grouper} \\ \text{turtle} \end{array}$$

- n training samples: $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$
- Goal: minimizing the *empirical risk* with respect to θ

$$\min_{\theta} R(f_\theta) = \min_{\theta} \frac{1}{N} \sum_{i=1}^N L(\mathbf{y}_i, f_\theta(\mathbf{x}_i))$$

Motivating example: Domain adaptation

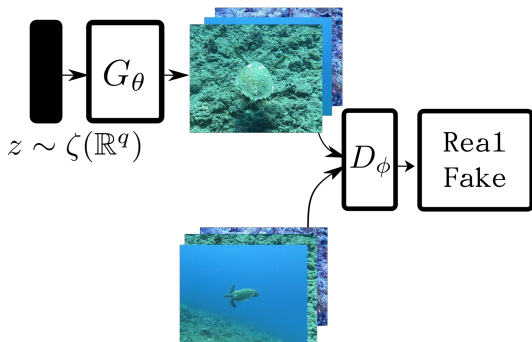


Domain adaptation (DA) setting

- Two domains with same classes, only one with labels
- Goal: classify unlabeled target data with source labeled data
- $\mathbf{x}_i^s, \mathbf{x}_j^t$ have same class $\rightarrow g_\phi(\mathbf{x}_i) \approx g_\phi(\mathbf{x}_j)$ and $\mathbf{y}_i = f_\theta(g_\phi(\mathbf{x}_j))$

Motivating example: Generative adversarial networks

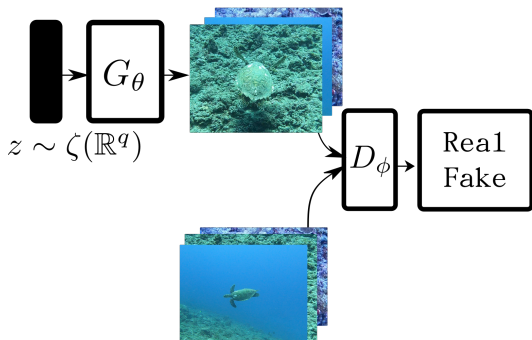
Goal: generating new images



- Generative adversarial networks (GANs) developed in [Goodfellow et al., 2014]
- G_θ tries to fool D_ϕ
- D_ϕ tries to predict if an image is real or not

Motivating example: Generative adversarial networks

Goal: generating new images



- $\alpha \in \mathcal{P}(\mathcal{X}), \zeta \in \mathcal{P}(\mathcal{Z})$ are probability distributions
- Loss: $\min_{\theta} \max_{\phi} \mathbb{E}_{\mathbf{x} \sim \alpha} \log(D_{\phi}(\mathbf{x})) - \mathbb{E}_{\mathbf{z} \sim \zeta} \log(1 - D_{\phi}(G_{\theta}(\mathbf{z})))$

The loss can be reformulated with a Jensen-Shannon divergence between generated and training distributions

Training samples as distributions paradigm

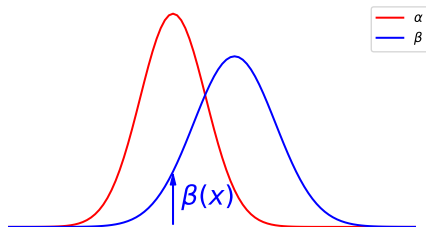
Applications use probability distributions to train neural networks

- Classification: function L takes probability vectors as inputs
- Domain adaptation: align embedding probability distributions from domains
- GANs: distance between generated and training distributions

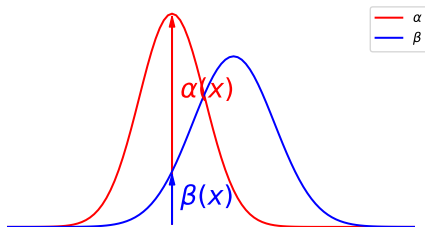
$$\hat{\theta} = \arg \min_{\theta \in \Theta} L(\alpha_n, \beta_{\theta})$$

Goal : Find a suitable function L between probability distributions

Comparing probability distributions



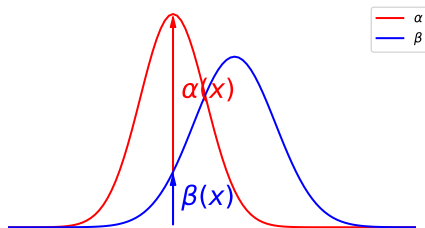
Comparing probability distributions



φ -divergences compare mass ratio point-wise $\alpha(\mathbf{x})/\beta(\mathbf{x})$ ($\beta(\mathbf{x}) > 0$)

$$L_{\varphi}(\alpha|\beta) = \int_{\mathcal{X}} \varphi \left(\frac{d\alpha}{d\beta} \right) d\beta$$

Comparing probability distributions

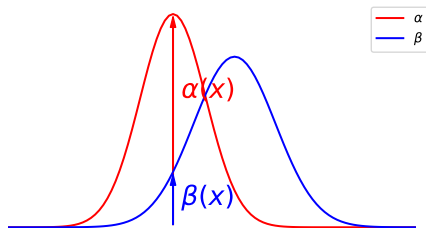


φ -divergences compare mass ratio point-wise $\alpha(\mathbf{x})/\beta(\mathbf{x})$ ($\beta(\mathbf{x}) > 0$)

$$L_{\varphi}(\alpha|\beta) = \int_{\mathcal{X}} \varphi \left(\frac{d\alpha}{d\beta} \right) d\beta$$

- φ -divergences cannot compare Diracs
- fail to capture the geometry
- $\text{KL}(\alpha|\beta_t) = +\infty$

Comparing probability distributions

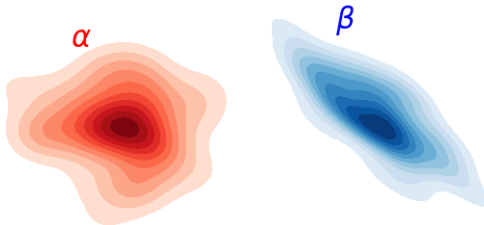


φ -divergences compare mass ratio point-wise $\alpha(\mathbf{x})/\beta(\mathbf{x})$ ($\beta(\mathbf{x}) > 0$)

$$L_{\varphi}(\alpha|\beta) = \int_{\mathcal{X}} \varphi \left(\frac{d\alpha}{d\beta} \right) d\beta$$

- φ -divergences cannot compare Diracs
→ fail to capture the geometry
- $\text{KL}(\alpha|\beta_t) = +\infty$ but $\text{KL}(\alpha|\beta_{\infty}) = 0$

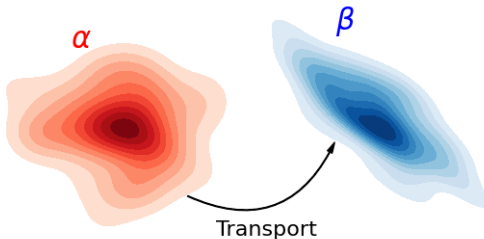
Optimal Transport definition



Ingredients

- Probability distributions $\alpha \in \mathcal{P}(\mathcal{X})$ and $\beta \in \mathcal{P}(\mathcal{Y})$
- A ground cost $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ with \mathcal{X} and \mathcal{Y} metric spaces

Optimal Transport definition



Ingredients

- Probability distributions $\alpha \in \mathcal{P}(\mathcal{X})$ and $\beta \in \mathcal{P}(\mathcal{Y})$
- A ground cost $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ with \mathcal{X} and \mathcal{Y} metric spaces

Optimal Transport definition

Definition (Kantorovich problem [[Kantorovich, 1942](#)])

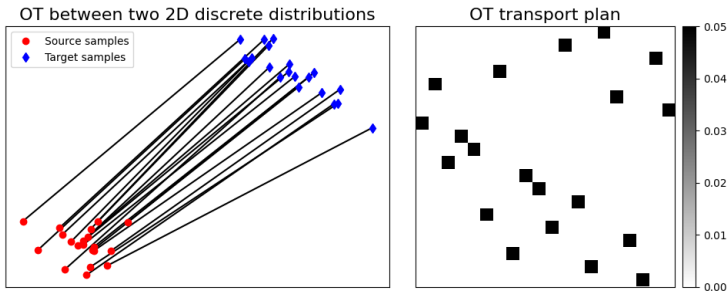
$$\min_{\pi \in U(\alpha, \beta)} \int_{\mathcal{X} \times \mathcal{Y}} c(\mathbf{x}, \mathbf{y}) d\pi(\mathbf{x}, \mathbf{y})$$

with : $U(\alpha, \beta) = \{\pi \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}), \int_{\mathcal{Y}} \pi(\mathbf{x}, \mathbf{y}) d\mathbf{y} = \alpha, \int_{\mathcal{X}} \pi(\mathbf{x}, \mathbf{y}) d\mathbf{x} = \beta\}$

Discrete ingredients

- Discrete distributions $\alpha = \sum_{i=1}^n a_i \delta_{\mathbf{x}_i}$ and $\beta = \sum_{j=1}^n b_j \delta_{\mathbf{y}_j}$
- Cost matrix $C = C(X, Y)$, such that $C_{i,j} = c(\mathbf{x}_i, \mathbf{y}_j)$

Discrete Optimal Transport



For discrete distributions, OT becomes a linear program:

Definition (Discrete Optimal Transport)

$$\text{OT}(\alpha, \beta, C) = \min_{\Pi \in U(\alpha, \beta)} \sum_{i,j} \Pi_{i,j} C_{i,j}$$

$$U(\alpha, \beta) = \left\{ \Pi \in (\mathbb{R}^+)^{n \times n} \mid \Pi \mathbf{1}_n = \mathbf{a}, \Pi^T \mathbf{1}_n = \mathbf{b} \right\}$$

Wasserstein distance

Some properties

- Leverages geometry of sample spaces through C
- A solution always exists (ex. $\pi = \alpha \otimes \beta$)
- $\langle \Pi, C \rangle_F$ is linear in the transport plan and in the cost
- Computational complexity of discrete OT is $\mathcal{O}(n^3 \log(n))$

Definition (Wasserstein distance)

C is a ground metric, then OT cost W_p is a metric for $p \geq 1$ and where

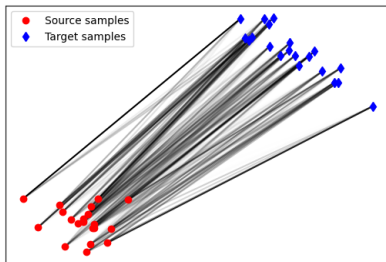
$$W_p(\alpha, \beta, C^p) = \left(\min_{\Pi \in U(\alpha, \beta)} \langle \Pi, C^p \rangle_F \right)^{1/p}$$

Proposition (Kantorovich–Rubinstein duality)

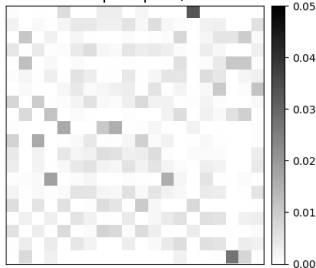
$$W_1(\alpha, \beta, C) = \sup_{f \in \text{Lip}^1(\mathcal{X})} \mathbb{E}_{\mathbf{x} \sim \alpha}[f(\mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim \beta}[f(\mathbf{z})]$$

Entropic Optimal Transport

E-OT between two 2D discrete distributions



E-OT transport plan, $\varepsilon = 0.05$



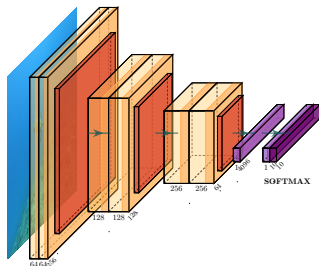
Definition (Entropic Optimal Transport [Cuturi, 2013])

$$\text{OT}^\varepsilon(\alpha, \beta, C) = \min_{\Pi \in U(\alpha, \beta)} \sum_{i,j} \Pi_{i,j} C_{i,j} + \varepsilon \text{KL}(\Pi | \alpha \otimes \beta)$$

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}_+^n, \text{KL}(\mathbf{x} | \mathbf{y}) = \sum_i \mathbf{x}_i \log \left(\frac{\mathbf{x}_i}{\mathbf{y}_i} \right) - \mathbf{x}_i + \mathbf{y}_i$$

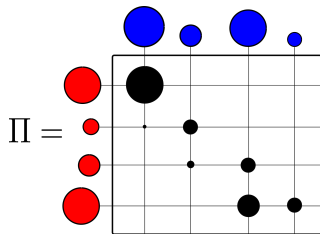
- Functional is strongly convex in the transport plan
- Computational complexity of entropic OT is $\mathcal{O}\left(\frac{n^2}{\varepsilon}\right)$

Summary on neural networks and optimal transport



• Summary on neural networks

- Neural networks are stacked layers of neurons
- Competitive methods on classification, domain adaptation and GANs



• Summary on optimal transport

- + Loss function/distance between distributions of samples
- + Leverages geometry of sample spaces through C
- Cubical computational complexity of discrete OT
- +/- Faster and easy computable entropic variant

- **Optimal transport as a loss function in deep learning**

Geometry on label distributions \leftarrow Optimal transport adversarial regularization for noisy labels

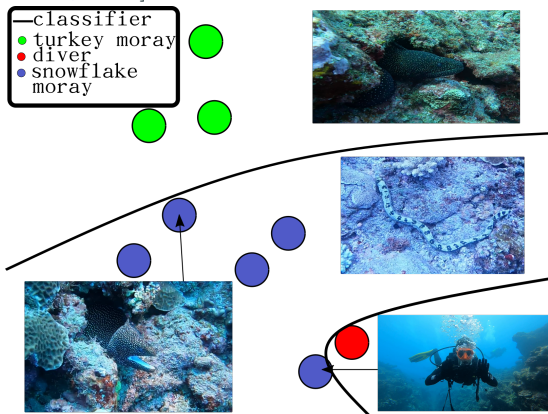
Geometry on sample distributions \leftarrow Optimal transport loss function to generate misclassified data

- **Minibatch optimal transport**

- Minibatch optimal transport formalism
- Unbalanced minibatch optimal transport

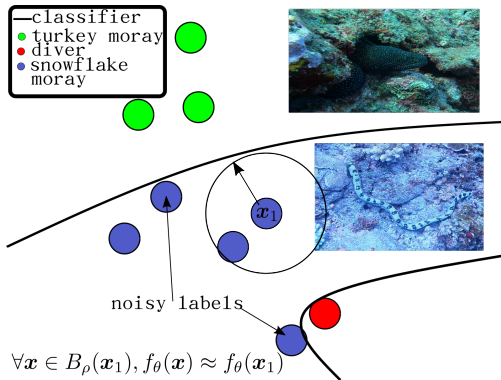
Wasserstein Adversarial Regularization for noisy labels

- Label noise: label does not correspond to the image class
 - Occurs in real-world dataset → neural networks overfitting
- [Zhang et al., 2017].



- Enforce uniformity prediction around the vicinity of samples
- Use optimal transport in the regularization

Robust optimization illustration

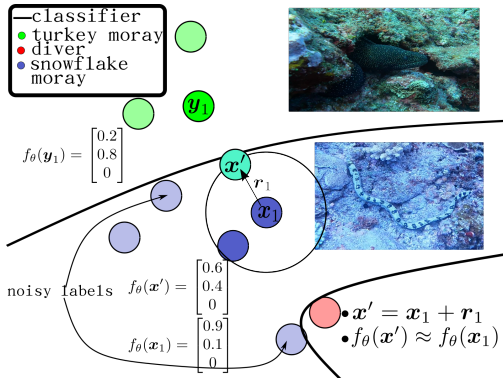


Robust optimization:

$$\operatorname{argmin}_{\theta} \sum_{i=1}^n \max_{\mathbf{x}_i^u \in B_{\rho}(\mathbf{x}_i)} L_{\text{CE}}(f_{\theta}(\mathbf{x}_i^u), \mathbf{y}_i) \quad (1)$$

- Intuition: noisy labels mitigated by uniform prediction
- Cannot rely on labels \rightarrow use prediction

Virtual adversarial training

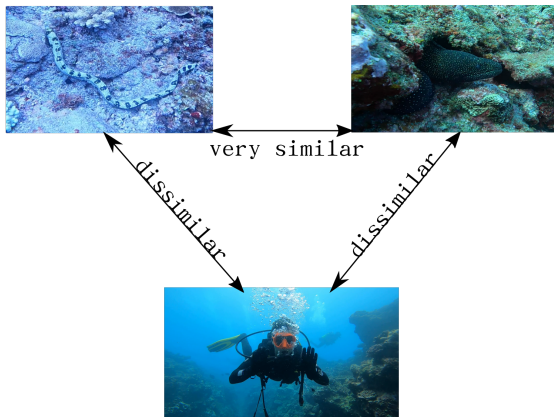


VAT's loss function [Miyato et al., 2018]:

$$\mathcal{L}_{\text{VAT}}((X, Y), f_{\theta}) = \frac{1}{n} \sum_{i=1}^n \underbrace{\text{L}(\mathbf{y}_i, f_{\theta}(\mathbf{x}_i))}_{\text{Supervised loss}} + \beta \underbrace{\text{KL}(\hat{f}_{\theta}(\mathbf{x}_i), f_{\theta}(\mathbf{x}_i + \mathbf{r}_i))}_{\text{regularization term}}$$

$$\text{where } \mathbf{r}_i = \underset{\mathbf{r}, \|\mathbf{r}\| \leq \rho}{\text{argmax}} \text{KL}(\hat{f}_{\theta}(\mathbf{x}_i), f_{\theta}(\mathbf{x}_i + \mathbf{r}))$$

Encoding class similarities



- KL divergence penalizes errors between classes in the same manner
- Replace KL divergence by optimal transport

Wasserstein adversarial regularization

The WAR loss function is:

$$\mathcal{L}_{\text{WAR}}((X, Y), f_{\theta}) = \frac{1}{n} \sum_{i=1}^n \underbrace{L(\mathbf{y}_i, f_{\theta}(\mathbf{x}_i))}_{\text{Supervised loss}} + \beta \underbrace{\text{OT}_{\mathbf{C}}^{\varepsilon}(\hat{f}_{\theta}(\mathbf{x}_i), f_{\theta}(\mathbf{x}_i + \mathbf{r}_i))}_{\text{regularization term}}$$

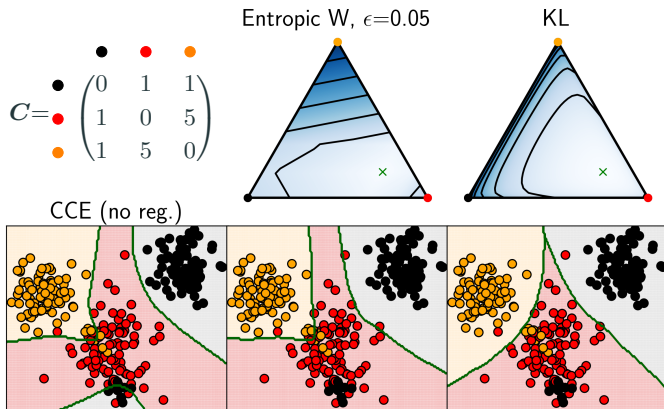
where $\mathbf{r}_i = \underset{\mathbf{r}, \|\mathbf{r}\| \leq \rho}{\operatorname{argmax}} \text{OT}_{\mathbf{C}}^{\varepsilon}(\hat{f}_{\theta}(\mathbf{x}_i), f_{\theta}(\mathbf{x}_i + \mathbf{r}))$

For the ground cost \mathbf{C} , we want:

- High cost between close classes to get complex boundaries
- Small cost between non-similar classes to get smooth boundaries

The OT cost between labels was also studied in [\[Frogner et al., 2015\]](#)

Learning with WAR



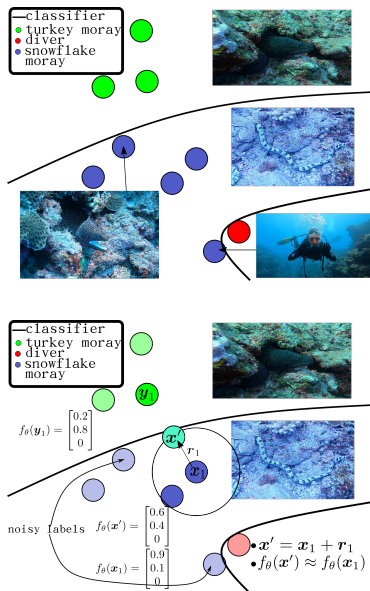
- Regularization geometry for different adversarial regularizations
- (Top) Regularization values on the simplex of class probabilities
- (Down) Classification boundaries for different methods

Experiments

- Test accuracy (%) on Fashion-MNIST, CIFAR-10, and CIFAR-100 datasets
- Varying noise rates (20% and 40%)

| Dataset / | noise | CCE | Bootsoft | CoTeaching | VAT | WAR _C |
|-----------|-------|------------|-------------|--------------|------------|------------------|
| F-MNIST | 20% | 89.02±0.47 | 88.17±0.11 | 91.24±0.06 | 93.10±0.14 | 93.37±0.08 |
| | 40% | 78.85±0.56 | 73.84±0.28 | 86.83±0.10 | 89.74±0.10 | 90.41±0.02 |
| CIFAR-10 | 20% | 85.26±0.09 | 85.35 ± 0.8 | 86.19 ± 0.07 | 88.91±0.09 | 89.12±0.48 |
| | 40% | 76.23±0.15 | 74.32 ± 0.2 | 80.87±0.09 | 81.98±0.25 | 84.55±0.78 |
| CIFAR-100 | 20% | 58.81±0.10 | 58.97±0.08 | 60.90±0.03 | 65.44±0.11 | 62.72±0.16 |
| | 40% | 42.45±0.12 | 41.73±0.08 | 42.73±0.08 | 55.75±0.14 | 58.86±0.21 |

Summary on Wasserstein Adversarial Regularization

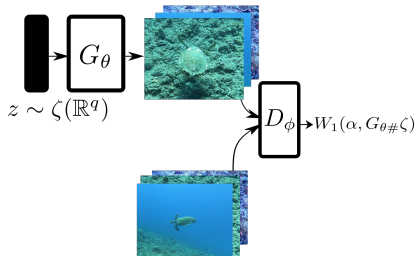
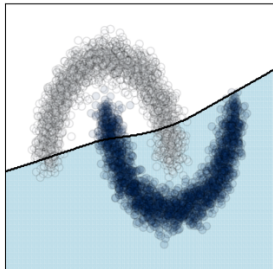


- Noisy labels are corrupted labels and hurt the performances of neural networks
- Promoting uniform classification around inputs mitigate their influence
- Integrate optimal transport to control the uniform classification

Published in [\[Fatras et al., 2021a\]](#)

Generating misclassified data

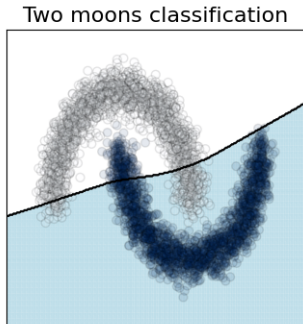
Two moons classification



- Attack the classifier
- Introduction on misclassified examples
- How OT can be used to generate data
- Generating misclassified examples with just the output of the classifier

Generating misclassified data

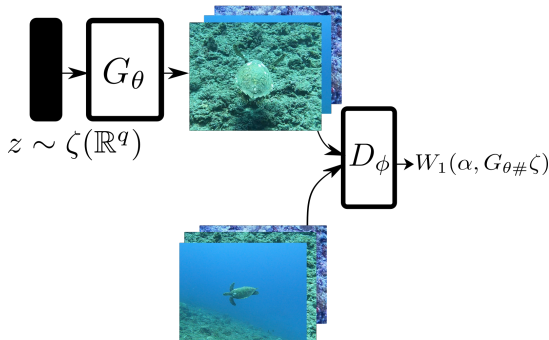
Our objective is to generate samples from the gray class which are classified as blue data



- Generate adversarial examples
- Most adversarial examples generation uses classifier architecture
- Some methods hurt quality of adversarial images

Generating data with Wasserstein GAN

To generate data, we use the *Wasserstein Generative Adversarial Networks* (WGAN) [Arjovsky et al., 2017]



We use the Kantorovich-Rubinstein duality theorem and we minimize:

$$\min_{\theta} W_1(\alpha, G_\theta \# \zeta) = \min_{\theta} \max_{\phi} \mathbb{E}_{\mathbf{x} \sim \alpha} [\mathcal{D}_\phi(\mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim \zeta} [\mathcal{D}_\phi(G_\theta(\mathbf{z}))]$$

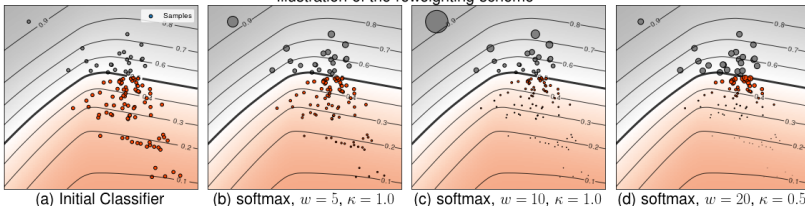
→ \mathcal{D}_ϕ needs to be 1-Lipschitz

Defining a new distribution

Create a new distribution which gives bigger weights to misclassified data, $\frac{1}{n} \sum_{i=1}^n \delta_{x_i} \rightarrow \sum_{i=1}^n a_i \delta_{x_i}, \sum_{i=1}^n a_i = 1$

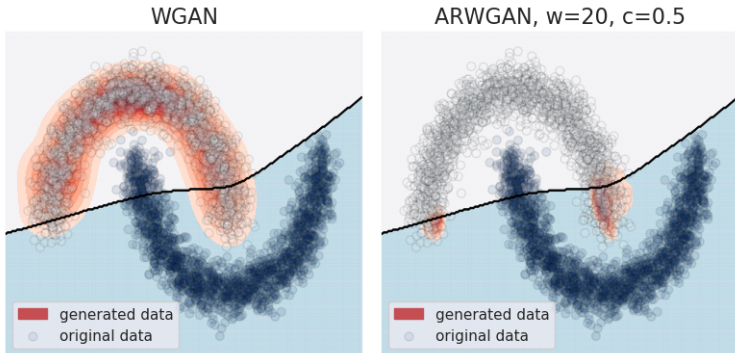
- Hard weighting (only consider misclassified samples)
- Soft weighting (weight depends on how much the sample is misclassified)

Illustration of the reweighting scheme



Generating misclassified data

- Fool classifiers on hyperspectral images
- Transfer misclassified examples to unseen classifiers
- Can modify images to fool classifier
- Can fool state of the art detector Yolo V3

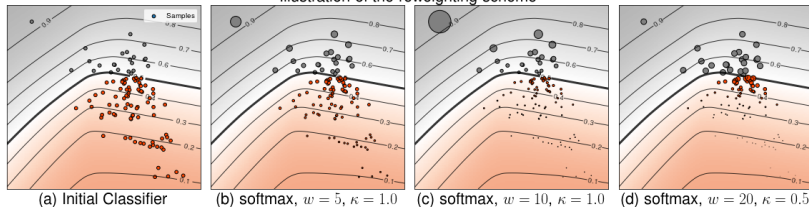


Summary on ARWGAN

- Attacked the classifier
- Used WGAN to generate samples
- Created new empirical distributions
- Applied on remote sensing data

Published in [Burnel et al., 2021, Burnel et al., 2020]

Illustration of the reweighting scheme



- **Optimal transport as a loss function in deep learning**

Geometry on label distributions \leftarrow Optimal transport adversarial regularization for noisy labels

Geometry on sample distributions \leftarrow Optimal transport loss function to generate misclassified data

- **Minibatch optimal transport**

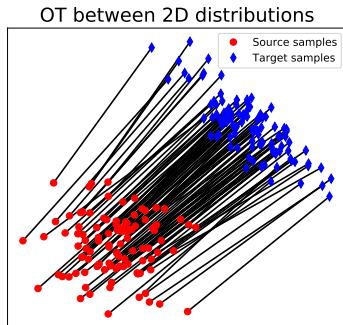
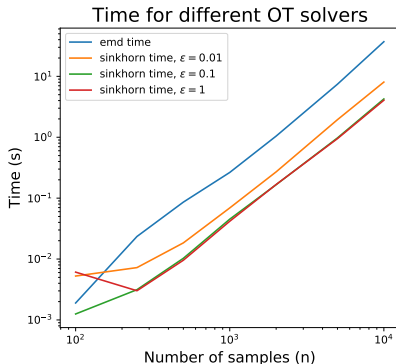
- Minibatch optimal transport formalism
- Unbalanced minibatch optimal transport

Minibatch optimal transport

- Minibatch optimal transport formalism
- Loss properties
- Statistical and optimization properties

Time experiment

Optimal transport can be computed between a lot of samples depending on the application



Limits

Can not be used in Big Data scenario !

Minibatch Optimal Transport definition

Let $m \leq n$, [Damodaran et al., 2018, Genevay et al., 2018] compute optimal transport between minibatch (MBOT) of distributions

Minibatch strategy

- Select m samples without replacement at random in domains
- Compute OT between the minibatches
- Average several MBOT terms \rightarrow complexity $\mathcal{O}(m^3)$

Minibatch Optimal Transport definition

Expectation of minibatches

Computing OT kernel h between minibatches estimates:

$$E_h(\alpha, \beta, C) := \mathbb{E}_{(X,Y) \sim \alpha^{\otimes m} \otimes \beta^{\otimes m}} [h(\mu_m, \mu_m, C(X, Y))]$$

can be any OT variants h (Gromov-Wasserstein distance, Sliced Wasserstein distance, ...)

Estimate minibatch OT distance

Definition (Complete minibatch estimator)

$$\bar{h}^m(X, Y) := \binom{n}{m}^{-2} \sum_{I, J \in \mathcal{P}_m} h(\mu_m, \mu_m, C_{I, J})$$

$$\Pi^m(X, Y) := \binom{n}{m}^{-2} \sum_{I, J \in \mathcal{P}_m} \Pi_{I, J}$$

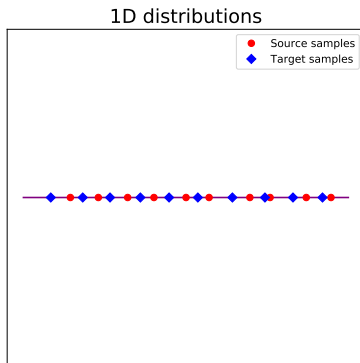
- where \mathcal{P}_m is the set of all m -tuples without replacement
- $\Pi^m(X, Y)$ is an admissible transport plan between the input probability distributions $\Pi \in U(\mu_n, \mu_n)$

Definition (Incomplete minibatch estimator)

$$\tilde{h}_k^m(X, Y) := k^{-1} \sum_{(I, J) \in D_k} h(\mu_m, \mu_m, C_{I, J})$$

where $k > 0$ is an integer and D_k is a set of cardinality k whose elements are minibatches drawn at random

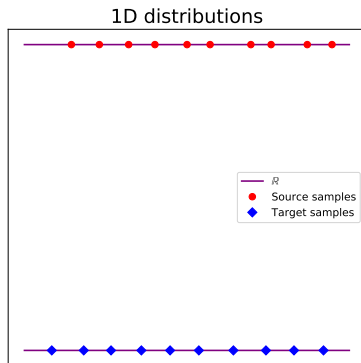
1D OT closed-form



1D closed-form

- Sorted 1D data with uniform weights
- Optimal Transport plan is the identity scaled by a uniform weight

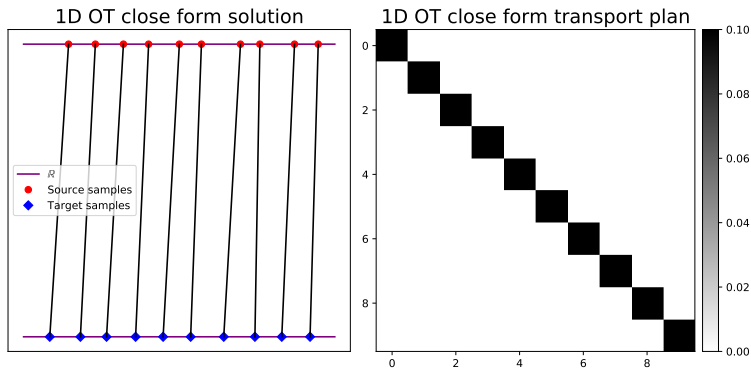
1D OT closed-form



1D closed-form

- Sorted 1D data with uniform weights
- Optimal Transport plan is the identity scaled by a uniform weight

1D OT closed-form



1D closed-form

- Sorted 1D data with uniform weights
- Optimal Transport plan is the identity scaled by a uniform weight

1D Minibatch Optimal Transport closed-form

$$\pi_{j,k} = \frac{1}{m} \binom{n}{m}^{-2} \sum_{i=i_{\min}}^{i_{\max}} \binom{j-1}{i-1} \binom{k-1}{i-1} \binom{n-j}{m-i} \binom{n-k}{m-i}$$

where $i_{\min} = \max(0, m - n + j, m - n + k)$ and $i_{\max} = \min(j, k)$

MBOT on a toy example

- 2 balanced classes in the source and target domains
- Classes are in different clusters

Proposition

We have the following properties:

- \bar{h}^m, \tilde{h}_k^m are unbiased estimators of E_h
- Strictly positive loss: $\bar{h}(\alpha, \alpha) > 0$

Difference with OT

- \triangleleft Minibatch OT is not a metric!
- OT empirical estimator is a biased estimator of OT between continuous measures ($\mathbb{E}_{\alpha_n, \beta_n} W(\alpha_n, \beta_n) \neq W(\alpha, \beta)$)
[Bellemare et al., 2017]

From now on, we suppose that α and β compactly supported and the cost c is at least continuous on \mathcal{X} and \mathcal{Y}

Deviation bounds

How far is our incomplete estimator \tilde{h}_k^m to the expectation over minibatches E_h ?

Theorem (Maximal deviation bound)

Let $\delta \in (0, 1)$, three integers $k \geq 1$ and $m \leq n$ be fixed. Consider two n -tuples $X \sim \alpha^{\otimes n}$ and $Y \sim \beta^{\otimes n}$. With probability at least $1 - \delta$ on the draw of X, Y and D_k we have:

$$|\tilde{h}_k^m(X, Y) - E_h| \leq M \left(\sqrt{\frac{\log(\frac{2}{\delta})}{2 \lfloor \frac{n}{m} \rfloor}} + \sqrt{\frac{2 \log(\frac{2}{\delta})}{k}} \right), \quad (2)$$

where M is an OT upper bound

Data fitting problem

Let discrete samples $(\mathbf{x}_i)_{i=1}^n \in \mathbb{R}^d$ and their empirical distribution α_n

Goal: to fit a parametric model $\theta \mapsto \beta_\theta \in \mathcal{M}_1^+(\mathbb{R}^d)$ to α_n using OT

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \text{OT}_c(\alpha_n, \beta_\theta)$$

It is known as Minimum Wasserstein estimator

Parameters are updated as $\theta_{t+1} = \theta_t + \eta_{\text{lr}} \nabla_{\theta} \text{OT}_c(\alpha_n, \beta_\theta)$

Data fitting problem

Let discrete samples $(\mathbf{x}_i)_{i=1}^n \in \mathbb{R}^d$ and their empirical distribution α_n

Goal: to fit a parametric model $\theta \mapsto \beta_\theta \in \mathcal{M}_1^+(\mathbb{R}^d)$ to α_n using OT

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \tilde{h}_k^m(\alpha_n, \beta_\theta)$$

Replace Wasserstein distance by minibatch OT to use SGD

Parameters are updated as $\theta_{t+1} = \theta_t + \eta_{\text{lr}} \nabla_{\theta} \tilde{h}_k^m(\alpha_n, \beta_\theta)$

Minimum Wasserstein estimator

Minimum MB Wasserstein estimator

Data fitting problem

Let discrete samples $(\mathbf{x}_i)_{i=1}^n \in \mathbb{R}^d$ and their empirical distribution α_n

Goal: to fit a parametric model $\theta \mapsto \beta_\theta \in \mathcal{M}_1^+(\mathbb{R}^d)$ to α_n using OT

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \tilde{h}_k^m(\alpha_n, \beta_\theta)$$

We replace the Wasserstein distance by minibatch OT to use SGD

Parameters are updated as $\theta_{t+1} = \theta_t + \eta_{lr} \nabla_\theta \tilde{h}_k^m(\alpha_n, \beta_\theta)$

Does minimizing this expectation with SGD converge towards the right minimum ?

- \bar{h}^m, \tilde{h}_k^m are unbiased estimators of E_h
- Can we exchange gradients and expectations?

Exchange gradient and expectation

Consider Clarke generalized gradients (defined for locally Lipschitz functions [Clarke, 1990])

Theorem

Let $\hat{X}, \{\hat{Y}_\theta\}_{\theta \in \Theta}$ be two m -tuples of random vectors compactly supported and C^m a \mathbf{C}^1 cost. Under an additional integrability assumption, we have:

$$\partial_\theta \mathbb{E}[h(\mu_m, \mu_m, C^m(\hat{X}, \hat{Y}_\theta))] = \mathbb{E}[\partial_\theta h(\mu_m, \mu_m, C^m(\hat{X}, \hat{Y}_\theta))],$$

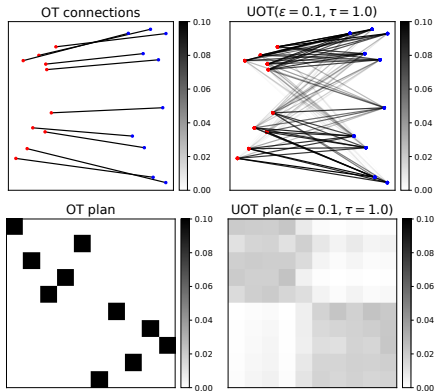
with both expectation being finite. Furthermore the function $\theta \mapsto -\mathbb{E}[h(\mu_m, \mu_m, C^m(\hat{X}, \hat{Y}_\theta))]$ is also Clarke regular.

→ SGD converges almost surely [Davis et al., 2020]

- Minibatch optimal transport formalism
- MBOT is a transport problem but not a distance
- Good statistical and optimization properties

Published in [Fatras et al., 2020b, Fatras et al., 2020a] and submitted in [Fatras et al., 2021c]

Unbalanced minibatch optimal transport



- Limits of MBOT
- Unbalanced minibatch OT
- Domain adaptation experiments

Limits of minibatch OT

- ⚠ Minibatches and marginal constraints create connections between classes !

Limits of Minibatch Optimal Transport

This is due to

- the sampling effect
- the marginal constraints

These force users to use large minibatch size [Damodaran et al., 2018]

Unbalanced Optimal Transport

Definition

Unbalanced optimal transport (UOT) measures the OT cost between probability distributions with relaxed marginals

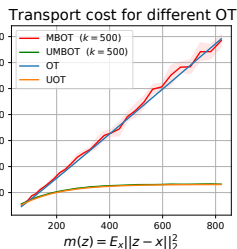
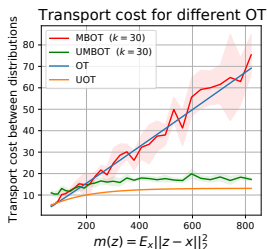
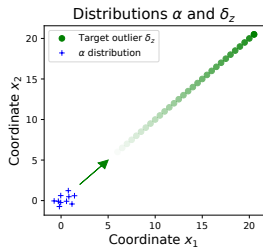
$$\begin{aligned} \text{UOT}^{\tau, \varepsilon}(\alpha, \beta, c) = \min_{\pi \in \mathcal{M}_+(\mathcal{X} \times \mathcal{Y})} & \int c d\pi + \varepsilon \text{KL}(\pi | \alpha \otimes \beta) \\ & + \tau (\text{KL}(\pi_1 \| \alpha) + \text{KL}(\pi_2 \| \beta)), \end{aligned}$$

where π is the transport plan, π_1 and π_2 the plan's marginals, $\tau \geq 0$ is the marginal penalization and $\varepsilon \geq 0$ is the regularization coefficient

Difference with OT

- $\pi \in U(\alpha, \beta) \longrightarrow \pi \in \mathcal{M}_+(\mathcal{X} \times \mathcal{Y})$
- Fixed marginal constraints are replaced by $\text{KL}(\pi_1 \| \alpha)$ penalties

Robust OT

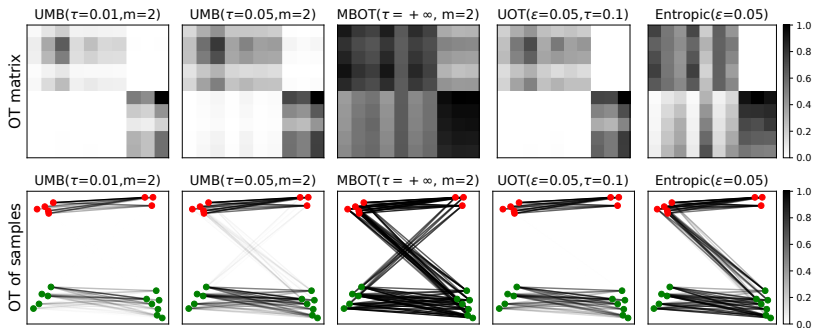


Lemma

Let (α, β) be two probability distributions. For $\zeta \in [0, 1]$, write $\tilde{\alpha} = \zeta\alpha + (1 - \zeta)\delta_z$. Write $m(z) = \int C(z, y)d\beta(y)$.

$$\text{UOT}^{\tau, 0}(\tilde{\alpha}, \beta, C) \lesssim \zeta \text{UOT}^{\tau, 0}(\alpha, \beta, C) + 2\tau(1 - \zeta)(1 - e^{-m(z)/2\tau})$$

Unbalanced minibatch OT plan

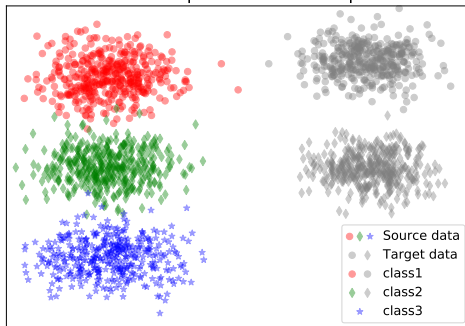


Unbalanced MBOT keeps the same properties as MBOT

- Same loss properties (not a distance but symmetric)
- Same deviation bounds rates
- Same unbiased gradients

Domain adaptation problem

Illustration of partial domain adaptation



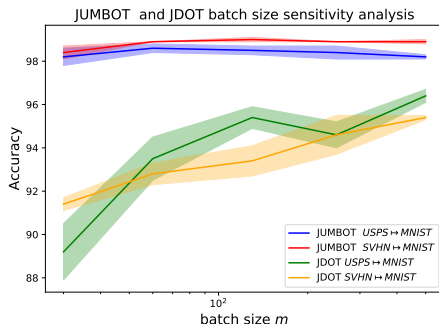
Method

- JUMBOT aligns a joint law between embedded samples and labels like DEEPJDOT
- Use unbalanced minibatch OT instead of minibatch OT
- Can be applied to partial DA unlike DEEPJDOT

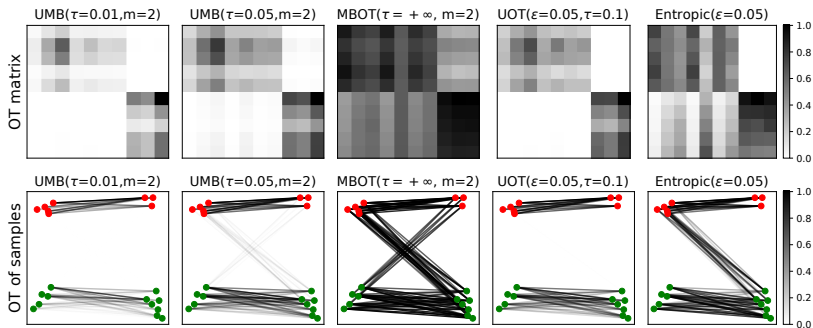
Results, ablation and sensitivity

JUMBOT outperforms state of the art methods on digits (U-M-S), Office Home and VisDA datasets including Partial Office-Home

| Methods | U \rightarrow M | S \rightarrow M |
|-------------------|----------------------------------|----------------------------------|
| DEEPPDOT | 96.4 \pm 0.3 | 95.4 \pm 0.1 |
| ENTROPIC DEEPPDOT | 97.1 \pm 0.3 | 97.6 \pm 0.1 |
| JUMBOT | 98.2 \pm 0.1 | 98.9 \pm 0.1 |

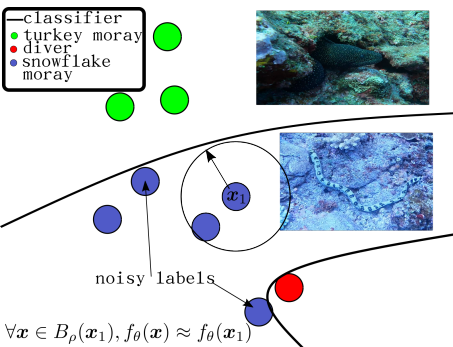


Summary on unbalanced minibatch optimal transport



- Minibatch optimal transport creates non-optimal connections
- Replace OT by unbalanced OT
- Same statistical and optimization properties
- JUMBOT outperforms MBOT methods on DA experiments

Conclusion



- **Optimal transport in deep learning**
 - as an adversarial regularization for noisy labels
 - to generate misclassified data
- **Minibatch optimal transport**
 - transport problem but not a distance
 - good statistical and optimization properties
 - unbalanced optimal transport to mitigate bad connections

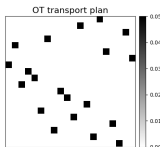
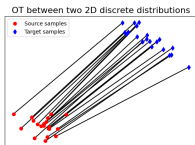
Future work on

- **Optimal transport in deep learning**

- Ground cost for WAR
- Optimal transport for out-of-distribution samples
- Normalizing flow
- To learn weights to make OT robust to outliers

- **Optimal transport**

- Sliced unbalanced OT
- Unbalanced OT when $\tau \rightarrow 0$
- New sampling schemes of MBOT



- Wasserstein adversarial regularization [Fatras et al., 2021a]
- ARWGAN [Burnel et al., 2020, Burnel et al., 2021]
- Minibatch optimal transport formalism
[Fatras et al., 2020b, Fatras et al., 2020a, Fatras et al., 2021c]
- Unbalanced minibatch optimal transport [Fatras et al., 2021b]
- Stochastic optimization [Pedregosa et al., 2019]
- Open source OT library [Flamary et al., 2021]

Thank you for your attention !

- [Arjovsky et al., 2017] Arjovsky, M., Chintala, S., and Bottou, L. (2017).
Wasserstein generative adversarial networks.
In Proceedings of the 34th International Conference on Machine Learning.
- [Bellemare et al., 2017] Bellemare, M. G., Danihelka, I., Dabney, W., Mohamed, S., Lakshminarayanan, B., Hoyer, S., and Munos, R. (2017).
The cramer distance as a solution to biased wasserstein gradients.
CoRR, abs/1705.10743.
- [Burnel et al., 2020] Burnel, J.-C., Fatras, K., and Courty, N. (2020).
Generating natural adversarial hyperspectral examples with a modified wasserstein gan.
In C&ESAR 2020.
- [Burnel et al., 2021] Burnel, J.-C., Fatras, K., Flamary, R., and Courty, N. (2021).
Generating natural adversarial remote sensing images.
IEEE Transactions on Geoscience and Remote Sensing.
- [Clarke, 1990] Clarke, F. H. (1990).
Optimization and nonsmooth analysis.
SIAM.

- [Cléménçon et al., 2016] Cléménçon, S., Colin, I., and Bellet, A. (2016).
Scaling-up empirical risk minimization: Optimization of incomplete u -statistics.
Journal of Machine Learning Research.
- [Courty et al., 2017] Courty, N., Flamary, R., Tuia, D., and Rakotomamonjy, A. (2017).
Optimal transport for domain adaptation.
IEEE Transactions on Pattern Analysis and Machine Intelligence.
- [Cuturi, 2013] Cuturi, M. (2013).
Sinkhorn distances: Lightspeed computation of optimal transport.
In Burges, C. J. C., Bottou, L., Welling, M., Ghahramani, Z., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems 26*, pages 2292–2300. Curran Associates, Inc.
- [Damodaran et al., 2018] Damodaran, B. B., Kellenberger, B., Flamary, R., Tuia, D., and Courty, N. (2018).
DeepJDOT: Deep Joint Distribution Optimal Transport for Unsupervised Domain Adaptation.
In *ECCV 2018 - 15th European Conference on Computer Vision*. Springer.
- [Davis et al., 2020] Davis, D., Drusvyatskiy, D., Kakade, S., and Lee, J. D. (2020).
Stochastic subgradient method converges on tame functions.
Foundations of computational mathematics, 20(1):119–154.

- [Fatras et al., 2021a] Fatras, K., Bushan, B., Lobry, S., Flamary, R., Tuia, D., and Courty, N. (2021a).
Wasserstein adversarial regularization for learning with label noise.
IEEE Transactions on Pattern Analysis and Machine Intelligence, pages 1–1.
- [Fatras et al., 2021b] Fatras, K., Séjourné, T., Courty, N., and Flamary, R. (2021b).
Unbalanced minibatch optimal transport; applications to domain adaptation.
In *Proceedings of the 38th International Conference on Machine Learning*.
- [Fatras et al., 2020a] Fatras, K., Zine, Y., Flamary, R., Gribonval, R., and Courty, N. (2020a).
Divergence wasserstein par lots.
In *52 èmes Journées de Statistiques de la Société Française de Statistique*.
- [Fatras et al., 2020b] Fatras, K., Zine, Y., Flamary, R., Gribonval, R., and Courty, N. (2020b).
Learning with minibatch wasserstein: asymptotic and gradient properties.
In *AISTATS*.
- [Fatras et al., 2021c] Fatras, K., Zine, Y., Majewski, S., Flamary, R., Gribonval, R., and Courty, N. (2021c).
Minibatch optimal transport distances; analysis and applications.
CoRR.

[Ferradans et al., 2013] Ferradans, S., Papadakis, N., Rabin, J., Peyré, G., and Aujol, J.-F. (2013).

Regularized discrete optimal transport.

In Scale Space and Variational Methods in Computer Vision. Springer Berlin Heidelberg.

[Feydy et al., 2019] Feydy, J., Séjourné, T., Vialard, F.-X., Amari, S.-i., Trounev, A., and Peyré, G. (2019).

Interpolating between optimal transport and mmd using sinkhorn divergences.

In Proceedings of Machine Learning Research.

[Flamary et al., 2021] Flamary, R., Courty, N., Gramfort, A., Alaya, M. Z., Boisbunon, A., Chambon, S., Chapel, L., Corenflos, A., Fatras, K., Fournier, N., Gautheron, L., Gayraud, N. T., Janati, H., Rakotomamonjy, A., Redko, I., Rolet, A., Schutz, A., Seguy, V., Sutherland, D. J., Tavenard, R., Tong, A., and Vayer, T. (2021).

Pot: Python optimal transport.

Journal of Machine Learning Research, 22(78):1–8.

- [Frogner et al., 2015] Frogner, C., Zhang, C., Mobahi, H., Araya, M., and Poggio, T. A. (2015).
Learning with a wasserstein loss.
In Cortes, C., Lawrence, N. D., Lee, D. D., Sugiyama, M., and Garnett, R., editors, *Advances in Neural Information Processing Systems 28*, pages 2053–2061. Curran Associates, Inc.
- [Genevay et al., 2018] Genevay, A., Peyre, G., and Cuturi, M. (2018).
Learning generative models with sinkhorn divergences.
In *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*.
- [Goodfellow et al., 2014] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014).
Generative adversarial nets.
In Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N. D., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc.
- [Kantorovich, 1942] Kantorovich, L. V. (1942).
On translation of mass (in russian).
Proceedings of the USSR Academy of Sciences.

- [Majewski et al., 2018] Majewski, S., Miasojedow, B., and Moulines, E. (2018).
Analysis of nonsmooth stochastic approximation: the differential inclusion approach.
arXiv preprint arXiv:1805.01916.
- [Miyato et al., 2018] Miyato, T., Maeda, S., Ishii, S., and Koyama, M. (2018).
Virtual Adversarial Training: A regularization method for supervised and semi-supervised learning.
IEEE Transactions on Pattern Analysis and Machine Intelligence.
- [Pedregosa et al., 2019] Pedregosa, F., Fatras, K., and Casotto, M. (2019).
Proximal splitting meets variance reduction.
In Chaudhuri, K. and Sugiyama, M., editors, *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pages 1–10. PMLR.
- [Peng et al., 2017] Peng, X., Usman, B., Kaushik, N., Hoffman, J., Wang, D., and Saenko, K. (2017).
Visda: The visual domain adaptation challenge.
- [Peyré and Cuturi, 2019] Peyré, G. and Cuturi, M. (2019).
Computational optimal transport.
Foundations and Trends® in Machine Learning, 11(5-6):355–607.

[Venkateswara et al., 2017] Venkateswara, H., Eusebio, J., Chakraborty, S., and Panchanathan, S. (2017).

Deep hashing network for unsupervised domain adaptation.

In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5018–5027.

[Zhang et al., 2017] Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. (2017).

Understanding deep learning requires rethinking generalization.

In *The International Conference on Learning Representations*.

Appendix

Hyperparameters i

| Datasets | η_1 | η_2 | η_3 | τ | ε |
|---------------------|----------|----------|----------|--------|---------------|
| DIGITS | 0.1 | 0.1 | 1 | 1 | 0.1 |
| VisDA | 0.005 | 1 | 1 | 0.3 | 0.01 |
| OFFICE-HOME | 0.01 | 0.5 | 1 | 0.5 | 0.01 |
| PARTIAL OFFICE-HOME | 0.003 | 0.75 | 10 | 0.06 | 0.01 |

Clarke regularity

Definition

A function f is said to be Clarke regular at \boldsymbol{x} provided:

- For all \boldsymbol{v} , the usual one-sided directional derivative $f'(\boldsymbol{x}, \boldsymbol{v})$ exists
- For all \boldsymbol{v} , $f'(\boldsymbol{x}; \boldsymbol{v}) = f^\circ(\boldsymbol{x}; \boldsymbol{v})$

Full echange gradients and expectations theorem

Theorem

Let \mathbf{u} be uniform probability vectors and let \mathbf{X} be a \mathbb{R}^{dm} -valued random variable, and $\{\mathbf{Y}_\theta\}$ a family of \mathbb{R}^{dm} -valued random variables defined on the same probability space, indexed by $\theta \in \Theta$, where $\Theta \subset \mathbb{R}^q$ is open. Assume that $\theta \mapsto \mathbf{Y}_\theta$ is \mathbf{C}^1 . Consider a \mathbf{C}^1 cost C , $h \in \{\text{UOT}^{\tau, \varepsilon}\}$, and assume in addition that the random variables $\mathbf{X}, \{\mathbf{Y}_\theta\}_{\theta \in \Theta}$ are compactly supported. If for all $\theta \in \Theta$ there exists an open neighbourhood U , $\theta \in U \subset \Theta$, and a random variable $K_U : \Omega \rightarrow \mathbb{R}$ with finite expected value, such that

$$\|C(\mathbf{X}(\omega), \mathbf{Y}_{\theta_1}(\omega)) - C(\mathbf{X}(\omega), \mathbf{Y}_{\theta_2}(\omega))\| \leq K_U(\omega) \|\theta_1 - \theta_2\| \quad (3)$$

then we have

$$\partial_\theta \mathbb{E} [h(\mathbf{u}, \mathbf{u}, C(\mathbf{X}, \mathbf{Y}_\theta))] = \mathbb{E} [\partial_\theta h(\mathbf{u}, \mathbf{u}, C(\mathbf{X}, \mathbf{Y}_\theta))] . \quad (4)$$

with both expectation being finite. Furthermore the function $\theta \mapsto -\mathbb{E} [h(\mathbf{u}, \mathbf{u}, C(\mathbf{X}, \mathbf{Y}_\theta))]$ is also Clarke regular.