# **Optimal Transport for Machine learning**

Domain Adaptation and structured data

R. Flamary - Lagrange, OCA, CNRS, Université Côte d'Azur

Meeting ANR MAGA, December 2018, Nancy











N. Courty

A. Rakotomamonjy

D. Tuia

A. Habrard





V. Seguy



B. B. Damodaran



T. Vayer







R. Tavenard

+ ANR OATMIL project members

# Introduction



- Probability measures  $\mu_s$  and  $\mu_t$  on and a cost function  $c: \Omega_s \times \Omega_t \to \mathbb{R}^+$ .
- The Monge formulation [Monge, 1781] aim at finding a mapping  $T: \Omega_s \to \Omega_t$

$$\inf_{T # \boldsymbol{\mu}_{\boldsymbol{s}} = \boldsymbol{\mu}_{\boldsymbol{t}}} \quad \int_{\Omega_{\boldsymbol{s}}} c(\mathbf{x}, T(\mathbf{x})) \boldsymbol{\mu}_{\boldsymbol{s}}(\mathbf{x}) d\mathbf{x}$$
(1)

• Non-convex optimization problem, mapping does not exist in the general case.

# **Optimal transport (Kantorovich formulation)**



The Kantorovich formulation [Kantorovich, 1942] seeks for a probabilistic coupling π ∈ P(Ω<sub>s</sub> × Ω<sub>t</sub>) between Ω<sub>s</sub> and Ω<sub>t</sub>:

$$\pi_0 = \underset{\pi}{\operatorname{argmin}} \int_{\Omega_s \times \Omega_t} c(\mathbf{x}, \mathbf{y}) \pi(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y},$$
(2)

s.t. 
$$\pi \in \Pi = \left\{ \pi \ge 0, \ \int_{\Omega_t} \pi(\mathbf{x}, \mathbf{y}) d\mathbf{y} = \mu_s, \int_{\Omega_s} \pi(\mathbf{x}, \mathbf{y}) d\mathbf{x} = \mu_t \right\}$$

- $\pi$  is a joint probability measure with marginals  $\mu_s$  and  $\mu_t$ .
- Linear Program that always have a solution.

# Wasserstein distance



#### Wasserstein distance

$$W_p^p(\boldsymbol{\mu}_s, \boldsymbol{\mu}_t) = \min_{\boldsymbol{\pi} \in \Pi} \quad \int_{\Omega_s \times \Omega_t} c(\mathbf{x}, \mathbf{y}) \boldsymbol{\pi}(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} = E_{(\mathbf{x}, \mathbf{y}) \sim \boldsymbol{\pi}}[c(\mathbf{x}, \mathbf{y})]$$
(3)

where  $c(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|^p$  is the ground metric.

- A.K.A. Earth Mover's Distance  $(W_1^1)$  [Rubner et al., 2000].
- Do not need the distribution to have overlapping support.
- Subgradients can be computed with the dual variables of the LP.
- Works for continuous and discrete distributions (histograms, empirical).



#### Short history of OT for ML

- Recently introduced to ML (well known in image processing since 2000s).
- Computationnal OT allow numerous applications (regularization).
- Deep learning boost (numerical optimization and GAN).

# Table of content

## Introduction

Optimal transport

Optimal transport and machine learning

## Optimal transport for domain adaptation

Supervised learning and Domain adapation

Optimal Transport for Domain Adaptation (OTDA)

Joint distribution OT for domain adaptation (JDOT)

## **Optimal Transport on structured data**

Gromov-Wasserstein distance for structured data

Structured data as distributions

Fused Gromov-Wasserstein distance

Applications on structured data classification

Optimal transport for domain adaptation

# Supervised learning

Amazon



# Traditional supervised learning

- We want to learn predictor such that  $y\approx f(\mathbf{x}).$
- Actual  $\mathcal{P}(X, Y)$  unknown.
- We have access to training dataset  $(\mathbf{x}_i, y_i)_{i=1,...,n} \ (\widehat{\mathcal{P}}(X, Y)).$
- We choose a loss function  $\mathcal{L}(y,f(\mathbf{x}))$  that measure the discrepancy.

#### **Empirical risk minimization** We week for a predictor f minimizing

$$\min_{f} \left\{ \mathbb{E}_{(\mathbf{x}, y) \sim \widehat{\mathcal{P}}} \mathcal{L}(y, f(\mathbf{x})) = \sum_{j} \mathcal{L}(y_j, f(\mathbf{x}_j)) \right\}$$
(4)

- Well known generalization results for predicting on new data.
- Loss is usually  $\mathcal{L}(y, f(\mathbf{x})) = (y f(\mathbf{x}))^2$  for least square regression and is  $\mathcal{L}(y, f(\mathbf{x})) = \max(0, 1 yf(\mathbf{x}))^2$  for squared Hinge loss SVM.

# **Domain Adaptation problem**



Probability Distribution Functions over the domains

#### Our context

- Classification problem with data coming from different sources (domains).
- Distributions are different but related.

# Unsupervised domain adaptation problem



# Problems

- Labels only available in the **source domain**, and classification is conducted in the **target domain**.
- Classifier trained on the source domain data performs badly in the target domain

# Optimal transport for domain adaptation



#### Assumptions

- $\bullet\,$  There exist a transport in the feature space  ${\bf T}$  between the two domains.
- The transport preserves the conditional distributions:

$$P_s(y|\mathbf{x}_s) = P_t(y|\mathbf{T}(\mathbf{x}_s)).$$

# 3-step strategy [Courty et al., 2016]

- 1. Estimate optimal transport between distributions.
- 2. Transport the training samples with barycentric mapping .
- 3. Learn a classifier on the transported training samples.

# **OT** for domain adaptation : **Step 1**



#### Step 1 : Estimate optimal transport between distributions.

- Choose the ground metric (squared euclidean in our experiments).
- Using regularization allows
  - Large scale and regular OT with entropic regularization [Cuturi, 2013].
  - Class labels in the transport with group lasso [Courty et al., 2016].
- Efficient optimization based on Bregman projections [Benamou et al., 2015] and
  - Majoration minimization for non-convex group lasso.
  - Generalized Conditionnal gradient for general regularization (cvx. lasso, Laplacian).

# OT for domain adaptation : Steps 2 & 3



Step 2 : Transport the training samples onto the target distribution.

- The mass of each source sample is spread onto the target samples (line of  $\pi_0$ ).
- Transport using barycentric mapping [Ferradans et al., 2014a].
- The mapping can be estimated for out of sample prediction [Perrot et al., 2016, Seguy et al., 2017].

#### Step 3 : Learn a classifier on the transported training samples

- Transported sample keep their labels.
- Classic ML problem when samples are well transported.

# Visual adaptation datasets



#### Datasets

- Digit recognition, MNIST VS USPS (10 classes, d=256, 2 dom.).
- Face recognition, PIE Dataset (68 classes, d=1024, 4 dom.).
- Object recognition, Caltech-Office dataset (10 classes, d=800/4096, 4 dom.).

#### Numerical experiments

- Comparison with state of the art on the 3 datasets.
- OT works very well on digits and object recognition.
- Works well on deep features adaptation and extension to semi-supervised DA.  $_{14/40}$

# Optimal transport for domain adaptation



# Discussion

- Works very well in practice for large class of transformation [Courty et al., 2016].
- Can use estimated mapping [Perrot et al., 2016, Seguy et al., 2017].

## $\mathsf{But}$

- Model transformation only in the feature space.
- Requires the same class proportion between domains [Tuia et al., 2015].
- We estimate a  $T : \mathbb{R}^d \to \mathbb{R}^d$  mapping for training a classifier  $f : \mathbb{R}^d \to \mathbb{R}$ .

# **Objectives of JDOT**

- Model the transformation of labels (allow change of proportion/value).
- Learn an optimal target predictor with no labels on target samples.
- Approach theoretically justified.

## Joint distributions and dataset

- Let  $\Omega \in \mathbb{R}^d$  be a feature space of dimension d and  $\mathcal{C}$  the set of labels.
- Let  $\mathcal{P}_s(X,Y) \in \mathcal{P}(\Omega \times C)$  and  $\mathcal{P}_t(X,Y) \in \mathcal{P}(\Omega \times C)$  the source and target joint distribution.
- We have access to an empirical sampling  $\hat{\mathcal{P}}_s = \frac{1}{N_s} \sum_{i=1}^{N_s} \delta_{\mathbf{x}_i^s, \mathbf{y}_i^s}$  of the source distribution defined by  $\mathbf{X}_s = \{\mathbf{x}_i^s\}_{i=1}^{N_s}$  and label information  $\mathbf{Y}_s = \{\mathbf{y}_i^s\}_{i=1}^{N_s}$ .
- but the target domain is defined only by an empirical distribution in the feature space with samples  $\mathbf{X}_t = {\{\mathbf{x}_i^t\}}_{i=1}^{N_t}$ .

# Proxy joint distribution

- Let f be a  $\Omega \to \mathcal{C}$  function from a given class of hypothesis  $\mathcal{H}$ .
- $\bullet\,$  We define the following joint distribution that use f as a proxy of y

$$\mathcal{P}_t^f = (\mathbf{x}, f(\mathbf{x}))_{\mathbf{x} \sim \mu_t} \tag{5}$$

and its empirical counterpart  $\hat{\mathcal{P}}_t^f = \frac{1}{N_t} \sum_{i=1}^{N_t} \delta_{\mathbf{x}_i^t, f(\mathbf{x}_i^t)}$ .

#### Learning with JDOT

We propose to learn the predictor f that minimize :

$$\min_{f} \left\{ W_1(\hat{\mathcal{P}}_s, \hat{\mathcal{P}}_t^f) = \inf_{\boldsymbol{\pi} \in \Pi} \sum_{ij} \mathcal{D}(\mathbf{x}_i^s, \mathbf{y}_i^s; \mathbf{x}_j^t, f(\mathbf{x}_j^t)) \boldsymbol{\pi}_{ij} \right\}$$
(6)

- $\Pi$  is the transport polytope.
- $\mathcal{D}(\mathbf{x}_i^s, \mathbf{y}_i^s; \mathbf{x}_j^t, f(\mathbf{x}_j^t)) = \alpha \|\mathbf{x}_i^s \mathbf{x}_j^t\|^2 + \mathcal{L}(\mathbf{y}_i^s, f(\mathbf{x}_j^t)) \text{ with } \alpha > 0.$
- We search for the predictor f that better align the joint distributions.
- Generalization bound show that expected risk on target is bounded by 6.

$$\min_{f \in \mathcal{H}, \pi \in \Pi} \sum_{i,j} \pi_{i,j} \left( \alpha d(\mathbf{x}_i^s, \mathbf{x}_j^t) + \mathcal{L}(y_i^s, f(\mathbf{x}_j^t)) \right) + \lambda \Omega(f)$$
(7)

## **Optimization procedure**

- $\Omega(f)$  is a regularization for the predictor f
- We propose to use block coordinate descent (BCD)/Gauss Seidel.
- Provably converges to a stationary point of the problem.

#### $\pi$ update for a fixed f

- Classical OT problem.
- Solved by network simplex.
- Regularized OT can be used (add a term to problem (7))

# f update for a fixed $\pi$ $\min_{f \in \mathcal{H}} \quad \sum_{i,j} \pi_{i,j} \mathcal{L}(y_i^s, f(\mathbf{x}_j^t)) + \lambda \Omega(f)$ (8)

- Weighted loss from all source labels.
- $\pi$  performs label propagation.

# **Regression with JDOT**



Least square regression with quadratic regularization For a fixed  $\pi$  the optimization problem is equivalent to

$$\min_{f \in \mathcal{H}} \quad \sum_{j} \frac{1}{n_t} \| \hat{y}_j - f(\mathbf{x}_j^t) \|^2 + \lambda \| f \|^2$$
(9)

- $\hat{y}_j = n_t \sum_j \pi_{i,j} y_i^s$  is a weighted average of the source target values.
- Note that this problem is linear instead of quadratic.
- Can use any solver (linear, kernel ridge, neural network).

# **Classification with JDOT**



#### **Multiclass classification with Hinge loss** For a fixed $\pi$ the optimization problem is equivalent to

$$\min_{f_k \in \mathcal{H}} \sum_{j,k} \hat{P}_{j,k} \mathcal{L}(1, f_k(\mathbf{x}_j^t)) + (1 - \hat{P}_{j,k}) \mathcal{L}(-1, f_k(\mathbf{x}_j^t)) + \lambda \sum_k \|f_k\|^2$$
(10)

- $\hat{\mathbf{P}}$  is the class proportion matrix  $\hat{\mathbf{P}} = \frac{1}{N_t} \boldsymbol{\pi}^\top \mathbf{P}^s$ .
- $\mathbf{P}^s$  and  $\mathbf{Y}^s$  are defined from the source data with One-vs-All strategy as

$$Y_{i,k}^s = \begin{cases} 1 & \text{if } y_i^s = k \\ -1 & \text{else} \end{cases}, \quad P_{i,k}^s = \begin{cases} 1 & \text{if } y_i^s = k \\ 0 & \text{else} \end{cases}$$

with  $k \in 1, \cdots, K$  and K being the number of classes.

# DeepJDOT



- $\bullet\,$  Learn simultaneously the embedding g and the classifier f.
- JDOT performed in the joint embedding/label space.

# DeepJDOT



- Learn simultaneously the embedding g and the classifier f.
- JDOT performed in the joint embedding/label space.
- Use minibatch to estimate OT and update g, f at each iterations.
- Scales to large datasets and estimate a representation for both domains.



- Evaluation of DeepJDOT on visual classification tasks.
- Digit adaptation between MNIST, USPS, SVHN, MNIST-M.
- Home-office [Venkateswara et al., 2017] and VisDA 2017 [Peng et al., 2017] dataset.
- Ablation study : all terms are impportant.
- TSNE projections of embeddings (MNIST $\rightarrow$ MNIST-M).

# **DeepJDOT** in action



- Evaluation of DeepJDOT on visual classification tasks.
- Digit adaptation between MNIST, USPS, SVHN, MNIST-M.
- Home-office [Venkateswara et al., 2017] and VisDA 2017 [Peng et al., 2017] dataset.
- Ablation study : all terms are impportant.
- TSNE projections of embeddings (MNIST $\rightarrow$ MNIST-M).



- Evaluation of DeepJDOT on visual classification tasks.
- Digit adaptation between MNIST, USPS, SVHN, MNIST-M.
- Home-office [Venkateswara et al., 2017] and VisDA 2017 [Peng et al., 2017] dataset.
- Ablation study : all terms are impportant.
- TSNE projections of embeddings (MNIST $\rightarrow$ MNIST-M).

# **DeepJDOT** in action



- Evaluation of DeepJDOT on visual classification tasks.
- Digit adaptation between MNIST, USPS, SVHN, MNIST-M.
- Home-office [Venkateswara et al., 2017] and VisDA 2017 [Peng et al., 2017] dataset.
- Ablation study : all terms are impportant.
- TSNE projections of embeddings (MNIST → MNIST-M).

# **DeepJDOT** in action



- Evaluation of DeepJDOT on visual classification tasks.
- Digit adaptation between MNIST, USPS, SVHN, MNIST-M.
- Home-office [Venkateswara et al., 2017] and VisDA 2017 [Peng et al., 2017] dataset.
- Ablation study : all terms are impportant.
- TSNE projections of embeddings (MNIST → MNIST-M).



## Optimal transport for DA

- Model transformation of the features.
- Conditional distribution preserved.
- Mapping between distributions.
- Learn classifier on the transported samples.

# Joint distribution OT for DA

- Model transformation of the joint distribution.
- General framework for DA.
- Theoretical justification with generalization bound.

**Optimal Transport on structured data** 

## Structured data



#### Structured data

- A structure data is viewed as a combination of features informations linked within each other by some structural information.
- Can be seen as a distribution on a joint feature/structure space.
- Example : labeled graph.

#### Meaningful distances on structured data

- Us both features (labels) and structure (graph).
- Allows for comparison, classification.
- Data science (statistics, means)

# Structured data



## Structured data

- A structure data is viewed as a combination of features informations linked within each other by some structural information.
- Can be seen as a distribution on a joint feature/structure space.
- Example : labeled graph.

#### Meaningful distances on structured data

- Us both features (labels) and structure (graph).
- Allows for comparison, classification.
- Data science (statistics, means)

## Structured data as distributions



Graph data representation

$$\mu = \sum_{i=1}^{n} h_i \delta_{(x_i a_i)}$$

- Nodes are weighted by their mass  $h_i$ .
- Features values  $a_i$  and  $b_j$  can be compared through the common metric
- But no common between the structure points  $x_i$  and  $y_j$ .

# **Optimal Transport for structured data**



Wasserstein distance for structures data

$$\mathcal{W}_p(\boldsymbol{\mu_A}, \boldsymbol{\mu_B}) = \left(\min_{\pi \in \Pi(\boldsymbol{\mu_A}, \boldsymbol{\mu_B})} \sum_{i,j} M_{i,j}^p \pi_{i,j}\right)^{\frac{1}{p}}$$

 $\mu_A = \sum_i h_i \delta_{a_i}$  and  $\mu_B = \sum_j g_j \delta_{b_j}, M_{i,j} = \|a_i - b_j\|$ 

- Wasserstein good for (empirical) distributions, samples as IID.
- OT can encode structure with OT Lp [Thorpe et al., 2017] by extending the feature space but requires the same ambient space.

# Gromov-Wasserstein distance for structured data



Inspired from Gabriel Peyré

GW for structured data [Memoli, 2011]

$$\mathcal{GW}_p(D, D', \boldsymbol{\mu}_{\boldsymbol{X}}, \boldsymbol{\mu}_{\boldsymbol{Y}}) = \left(\min_{\pi \in \Pi(\boldsymbol{\mu}_{\boldsymbol{s}}, \boldsymbol{\mu}_{\boldsymbol{t}})} \sum_{i, j, k, l} |D_{i,k} - D'_{j,l}|^p \pi_{i,j} \pi_{k,l}\right)^{\frac{1}{p}}$$

 $\mu_{X} = \sum_{i} h_i \delta_{x_i}$  and  $\mu_{Y} = \sum_{j} g_j \delta_{y_j}$  and  $D_{i,k} = \|x_i - x_k\|, D_{j,l}' = \|y_j - y_l\|$ 

- Distance over measures with no common ground space.
- Works well on graphs (using distances between nodes) but do not handle labels.
- Invariant to rotations and translation in either spaces.

## Fused Gromov-Wasserstein distance



#### **Fused Gromov Wasserstein distance**

$$\mathcal{FGW}_{p,q,\alpha}(D,D',\boldsymbol{\mu_s},\boldsymbol{\mu_t}) = \left(\min_{\pi \in \Pi(\boldsymbol{\mu_s},\boldsymbol{\mu_t})} \sum_{i,j,k,l} \left( (1-\alpha)M_{i,j}^q + \alpha |D_{i,k} - D_{j,l}'|^q \right)^p \pi_{i,j} \pi_{k,l} \right)^{\frac{1}{p}}$$

 $\mu_s = \sum_{i=1}^n h_i \delta_{x_i,a_i}$  and  $\mu_t = \sum_{j=1}^m g_j \delta_{y_j,b_j}$ 

- Parameters q > 1,  $\forall p \ge 1$ .
- $\alpha \in [0,1]$  is a trade off parameter between structure and features.

$$\mathcal{FGW}_{p,q,\alpha}(D,D',\boldsymbol{\mu_s},\boldsymbol{\mu_t}) = \left(\min_{\pi \in \Pi(\boldsymbol{\mu_s},\boldsymbol{\mu_t})} \sum_{i,j,k,l} \left( (1-\alpha)M_{i,j}^q + \alpha |D_{i,k} - D_{j,l}'|^q \right)^p \pi_{i,j} \pi_{k,l} \right)^{\frac{1}{p}}$$

# Metric properties

- *FGW* defines a metric over structured data with measure and features preserving isometries as invariants.
- $\mathcal{FGW}$  is a metric for q = 1 a semi metric for q > 1,  $\forall p \ge 1$ .
- The distance is nul iff :
  - There exists a Monge map  $T \# \mu_s = \mu_t$ .
  - Structures are equivalent through this Monge map (isometry).
  - Features are equal through this Monge map.

#### Other properties for sontinuous distributions

- Interpolation between  $\mathcal{W}$  ( $\alpha = 0$ ) and  $\mathcal{GW}$  ( $\alpha = 1$ ) distances.
- Geodesic properties (constant speed, unicity).

$$\mathcal{FGW}_{p,q,\alpha}(D,D',\boldsymbol{\mu_s},\boldsymbol{\mu_t}) = \left(\min_{\pi \in \Pi(\boldsymbol{\mu_s},\boldsymbol{\mu_t})} \sum_{i,j,k,l} \left( (1-\alpha)M_{i,j}^q + \alpha |D_{i,k} - D'_{j,l}|^q \right)^p \pi_{i,j} \pi_{k,l} \right)^{\frac{1}{p}}$$

#### Bounds and convergence to finite samples

• The following inequalities hold:

$$\mathcal{FGW}(\mu_s, \mu_t) \ge (1 - \alpha) \mathcal{W}(\mu_A, \mu_B)^q$$
$$\mathcal{FGW}(\mu_s, \mu_t) \ge \alpha \mathcal{GW}(\mu_X, \mu_Y)^q$$

• Bound when  $\mathcal{X} = \mathcal{Y}$ :

$$\mathcal{FGW}(\mu_s,\mu_t)^p \leq 2\mathcal{W}(\mu_s,\mu_t)^p$$

• Convergence of finite samples when  $\mathcal{X} = \mathcal{Y}$  with  $d = Dim(\mathcal{X}) + Dim(\Omega)$ :

$$\mathbb{E}[\mathcal{FGW}(\mu,\mu_n)] = O\left(n^{-\frac{1}{d}}\right)$$

$$\pi^* = \underset{\pi \in \Pi(\mu_s, \mu_t)}{\operatorname{arg\,min}} \quad \operatorname{vec}(\pi)^T Q \operatorname{vec}(\pi) + \operatorname{vec}((1-\alpha)M)^T \operatorname{vec}(\pi) \tag{12}$$

where  $Q = -2\alpha D' \otimes D$ 

# Algorithmic resolution (p = 1)

- Problem is a non-convex Quadratic Program.
- We use Conditional gradient [Ferradans et al., 2014b] with network simplex solver.
- Convergence to a local minima [Lacoste-Julien, 2016].
- With entropic regularization, projected gradient descent [Peyré et al., 2016].

$$\pi^* = \operatorname*{arg\,min}_{\pi \in \Pi(\mu_s, \mu_t)} \operatorname{vec}(\pi)^T Q \operatorname{vec}(\pi) + \operatorname{vec}((1 - \alpha)M)^T \operatorname{vec}(\pi)$$
(12)

## Algorithm 1 Conditional Gradient (CG) for FGW

- 1:  $\pi^{(0)} \leftarrow \mu_X \mu_Y^\top$
- 2: for  $i = 1, \ldots, \operatorname{do}$
- 3:  $G \leftarrow \text{Gradient from Eq. (12) } w.r.t. \pi^{(i-1)}$
- 4:  $\tilde{\pi}^{(i)} \leftarrow \text{Solve OT with ground loss } G$
- 5:  $\tau^{(i)} \leftarrow \text{Line-search for loss with } \tau \in (0,1)$

6: 
$$\pi^{(i)} \leftarrow (1 - \tau^{(i)})\pi^{(i-1)} + \tau^{(i)}\tilde{\pi}^{(i)}$$

7: end for

# Algorithmic resolution (p = 1)

- Problem is a non-convex Quadratic Program.
- We use Conditional gradient [Ferradans et al., 2014b] with network simplex solver.
- Convergence to a local minima [Lacoste-Julien, 2016].
- With entropic regularization, projected gradient descent [Peyré et al., 2016].

# Illustration of FGW distance



#### FGW maps on toy tree

- Uniform weights on the leafs of the tree.
- Structure distance taken as shortest path on the tree.
- Only FGW can encode both features and structures.

Vector attributes	AIDS	BZR	COX2	CUNEIFORM	ENZYMES	PROTEIN	SYNTHETIC
FGW SP FGW SP REGUL FGW WSP FGWDMM SP FGWDMM WSP	99.44+/-0.47 - 99.55+/-0.35 - -	85.12+/-4.15 <b>85.61+/-5.05</b> 84.88+/-4.34 84.39+/-5.48 83.17+/-5.05	$\begin{array}{c} 77.23 + /-4.86 \\ 77.66 + /-4.17 \\ 78.09 + /-3.81 \\ 76.81 + /-4.30 \\ 78.30 + /-3.53 \end{array}$	76.67+/-7.04 - - - -	$\begin{array}{c} 71.00+/-6.76\\ 70.17+/-6.81\\ 69.50+/-7.30\\ 61.67+/-7.19\\ 59.17+/-6.55 \end{array}$	$\begin{array}{c} 74.55+/-2.74\\ 74.64+/-2.99\\ 75.09+/-2.34\\ 75.00+/-2.59\\ 75.09+/-3.03 \end{array}$	100.00+/-0.00 - - - -
HOPPER ALL CV PROPA ALL CV PSCN K=10 PSCN K=5	99.50+/-0.59 98.45+/-1.06 99.80+/-0.24 <b>99.85+/-0.23</b>	$\begin{array}{c} 84.15+/-5.26\\ 79.51+/-5.02\\ 80.00+/-4.47\\ 82.20+/-4.23\end{array}$	<b>79.57+/-3.46</b> 77.66+/-3.95 71.70+/-3.57 71.91+/-3.40	32.59+/-8.73 12.59+/-6.67 25.19+/-7.73 24.81+/-7.23	45.33+/-4.00 71.67+/-5.63 26.67+/-4.77 27.33+/-4.16	$\begin{array}{c} 71.96+/-3.22\\ 61.34+/-4.38\\ 67.95+/-11.28\\ 71.79+/-3.39 \end{array}$	90.67+/-4.67 64.67+/-6.70 <b>100.00+/-0.00</b> <b>100.00+/-0.00</b>

#### **Graph classification**

- Classifiation accuracy on classical graph datasets.
- Comparison with state-of-the-art graph kernel approaches and Graph CNN.
- We use  $\exp(-\gamma \mathcal{FGW})$  as a non-positive kernel for an SVM [Loosli et al., 2016] (FGW).
- Train Wassertsein Distance Measure Machine [Rakotomamonjy et al., 2018] (FGWDMM).

DISCRETE ATTRIBUTES	MUTAG	NCI1	PTC	
FGW RAW SP	83.26+/-10.30	72.82+/-1.46	55.71+/-6.74	
FGW WL H=2 SP	$86.42 \pm -7.81$	$85.82 \pm -1.16$	$63.20 \pm -7.68$	
FGW WL H=2 SP REGUL	84.74 + / - 8.03	-	$63.37 \pm -6.75$	
FGW WL H=4 SP FGW WL H=4 SP REGUL	88.42+/-5.67 86.42+/-8.81	86.42 +/- 1.63	65.31+/-7.90 63.83+/-7.83	WITHOUT ATT
GK $\kappa=3$ PSCN $\kappa=10$ PSCN $\kappa=5$	82.42+/-8.40 83.47+/-10.26 82.05+/_10.80	60.78+/-2.48 70.65+/-2.58 60.85+/-1.70	56.46+/-8.03 58.34+/-7.71	FGW raw sp GK k=3 SP all cv
RW ALL CV SP ALL CV	83.05+/-10.80 79.47+/-8.17 82.95+/-8.19	58.63 + / -2.44 74.26 + / -1.53	55.09+/-7.34 -	
WL ALL CV WL H=2 WL H=4	86.21+/-8.48 86.21+/-8.15 83.68+/-9.13	85.77+/-1.07 81.85+/-2.28 85.13+/-1.61	62.86+/-7.23 61.60+/-8.14 62.17+/-7.80	

WITHOUT ATTRIBUTE	IMDB-B	IMDB-M
FGW RAW SP	63.80+/-3.49	48.00+/-3.22
GK K=3	56.00 + / - 3.61	$41.13 \pm -4.68$
SP all cv	$55.80 \pm -2.93$	$38.93 \pm -5.12$

## Graph classification

- Classifiation accuracy on classical graph datasets.
- Comparison with state-of-the-art graph kernel approaches and Graph CNN.
- We use  $\exp(-\gamma \mathcal{FGW})$  as a non-positive kernel for an SVM [Loosli et al., 2016] (FGW).
- Train Wassertsein Distance Measure Machine [Rakotomamonjy et al., 2018] (FGWDMM).

# FGW barycenter



**FGW barycenter** p = 1, q = 2

- Estimate FGW barycenter using Frechet means.
- Barycenter optimization solved via block coordinate descent (on  $\pi$ , D,  $\{a_i\}_i$ ).
- Can chose to fix the structure (D) or the features  $\{a_i\}_i$  in the barycenter.
- $a_{ii}$ , and D updates are weighted averages using  $\pi$ .



- We select a clean graph, change the number of nodes and add label noise and random connections.
- We compute the barycenter on n = 15 and n = 7 nodes.
- Barycenter graph is obtained through thresholding of the D matrix.



- We select a clean graph, change the number of nodes and add label noise and random connections.
- We compute the barycenter on n = 15 and n = 7 nodes.
- Barycenter graph is obtained through thresholding of the D matrix.



- We select a clean graph, change the number of nodes and add label noise and random connections.
- We compute the barycenter on n = 15 and n = 7 nodes.
- Barycenter graph is obtained through thresholding of the D matrix.



- We select a clean graph, change the number of nodes and add label noise and random connections.
- We compute the barycenter on n = 15 and n = 7 nodes.
- Barycenter graph is obtained through thresholding of the D matrix.



# Time series averaging

- Comparsion with Euclidean, DBA [Petitjean et al., 2011] and Soft-DTW [Cuturi and Blondel, 2017].
- Structure is time position of samples, fetaure value of the signal.
- Temporal position of nodes recovered with a MDS of *D*.
- Barycenter have non-regular sampling.



# Mesh interpolation

- Two meshes (deer and cat).
- Fix structure from cat, estimate barycenter for the positions of the edges.
- Wasserstien ( $\alpha = 0$ ) do not respect the graph (mesh neighborhood).
- FGW conserve the graph, regularized FGW smoothes the surface.



# Mesh interpolation

- Two meshes (deer and cat).
- Fix structure from cat, estimate barycenter for the positions of the edges.
- Wasserstien ( $\alpha = 0$ ) do not respect the graph (mesh neighborhood).
- FGW conserve the graph, regularized FGW smoothes the surface.



# Mesh interpolation

- Two meshes (deer and cat).
- Fix structure from cat, estimate barycenter for the positions of the edges.
- Wasserstien ( $\alpha = 0$ ) do not respect the graph (mesh neighborhood).
- FGW conserve the graph, regularized FGW smoothes the surface.

# FGW for community clustering



#### Graph approximation and comunity clustering

 $\min_{D,\mu} \quad \mathcal{FGW}(D, D_0, \mu, \mu_0)$ 

- Approximate the graph  $(D_0, \mu_0)$  with a small number of nodes.
- OT matrix give the clustering affectation.
- Works for signle and multiple modes in the clusters.

# FGW for community clustering



#### Graph approximation and comunity clustering

 $\min_{D,\mu} \quad \mathcal{FGW}(D, D_0, \mu, \mu_0)$ 

- Approximate the graph  $(D_0, \mu_0)$  with a small number of nodes.
- OT matrix give the clustering affectation.
- Works for signle and multiple modes in the clusters.





# Fused Gromov-Wasserstein distance [Vayer et al., 2018]

- Model structured data as distributions.
- New versatile method for comparing structured data based on Optimal Transport
- Many desirable distance properties
- New notion of barycenter of structured data such as graphs or time series
- Promising applications for signal over graphs and deep learning for structured data

#### What next ?

- Devise efficient optimization shemes for large structures.
- Add interpretability to deep neural networks on graph.

# Thank you

# Python code available on GitHub: https://github.com/rflamary/POT

- OT LP solver, Sinkhorn (stabilized,  $\epsilon$ -scaling, GPU)
- Domain adaptation with OT.
- Barycenters, Wasserstein unmixing.
- Wasserstein Discriminant Analysis.

Python code for JDOT on GitHub: https://github.com/rflamary/JDOT

Papers available on my website: https://remi.flamary.com/

#### Post docs available in:

Nice, Rouen, Rennes (France)



#### Expected loss

The expected loss on a domain D and for a given predictor f is defined as

$$err_D(f) \stackrel{\text{def}}{=} \mathop{\mathbb{E}}_{(\mathbf{x},y)\sim\mathcal{P}_t} \mathcal{L}(y, f(\mathbf{x})).$$

Probabilistic Lipschitzness [Urner et al., 2011, Ben-David et al., 2012] Let  $\phi : \mathbb{R} \to [0, 1]$ . A labeling function  $f : \Omega \to \mathbb{R}$  is  $\phi$ -Lipschitz with respect to a distribution P over  $\Omega$  if for all  $\lambda > 0$ 

$$Pr_{x \sim P}\left[\exists y : \left[|f(x) - f(y)| > \lambda d(x, y)\right]\right] \le \phi(\lambda).$$

#### Probabilistic Transfer Lipschitzness

Let  $\mu_s$  and  $\mu_t$  be respectively the source and target distributions. Let  $\phi : \mathbb{R} \to [0, 1]$ . A labeling function  $f : \Omega \to \mathbb{R}$  and a joint distribution  $\Pi(\mu_s, \mu_t)$  over  $\mu_s$  and  $\mu_t$  are  $\phi$ -Lipschitz transferable if for all  $\lambda > 0$ :

$$Pr_{(\mathbf{x}_1,\mathbf{x}_2)\sim\Pi(\mu_s,\mu_t)}\left[|f(\mathbf{x}_1) - f(\mathbf{x}_2)| > \lambda d(\mathbf{x}_1,\mathbf{x}_2)\right] \le \phi(\lambda).$$

#### Theorem 1

Let f be any labeling function of  $\in \mathcal{H}.$  Let

$$\begin{split} \Pi^* &= \operatorname{argmin}_{\Pi \in \Pi(\mathcal{P}_s, \mathcal{P}_t^f)} \int_{(\Omega \times \mathcal{C})^2} \alpha d(\mathbf{x}_s, \mathbf{x}_t) + \mathcal{L}(y_s, y_t) d\Pi(\mathbf{x}_s, y_s; \mathbf{x}_t, y_t) \text{ and } W_1(\hat{\mathcal{P}}_s, \hat{\mathcal{P}}_t^f) \text{ the} \\ \text{associated 1-Wasserstein distance. Let } f^* \in \mathcal{H} \text{ be a Lipschitz labeling function that verifies the} \\ \phi \text{-probabilistic transfer Lipschitzness (PTL) assumption w.r.t. } \Pi^* \text{ and that minimizes the joint error} \\ err_S(f^*) + err_T(f^*) \text{ w.r.t all PTL functions compatible with } \Pi^*. \text{ We assume the input instances are} \\ \text{bounded s.t. } |f^*(\mathbf{x}_1) - f^*(\mathbf{x}_2)| \leq M \text{ for all } \mathbf{x}_1, \mathbf{x}_2. \text{ Let } \mathcal{L} \text{ be any symmetric loss function, } k\text{-Lipschitz} \\ \text{and satisfying the triangle inequality. Consider a sample of } N_s \text{ labeled source instances drawn from } \mathcal{P}_s \text{ and } \\ N_t \text{ unlabeled instances drawn from } \mu_t, \text{ and then for all } \lambda > 0, \text{ with } \alpha = k\lambda, \text{ we have with probability at least } 1 - \delta \text{ that:} \end{split}$$

$$\begin{aligned} \operatorname{err}_{T}(f) &\leq W_{1}(\hat{\mathcal{P}_{s}}, \hat{\mathcal{P}_{t}^{f}}) + \sqrt{\frac{2}{c'}\log(\frac{2}{\delta})} \left(\frac{1}{\sqrt{N_{S}}} + \frac{1}{\sqrt{N_{T}}}\right) \\ &+ \operatorname{err}_{S}(f^{*}) + \operatorname{err}_{T}(f^{*}) + kM\phi(\lambda). \end{aligned}$$

- First term is JDOT objective function.
- Second term is an empirical sampling bound.
- Last terms are usual in DA [Mansour et al., 2009, Ben-David et al., 2010].

# References i

- Ben-David, S., Blitzer, J., Crammer, K., Kulesza, A., Pereira, F., and Vaughan, J. W. (2010).

A theory of learning from different domains.

Machine Learning, 79(1-2):151–175.

Ben-David, S., Shalev-Shwartz, S., and Urner, R. (2012).

Domain adaptation-can quantity compensate for quality? In *Proc of ISAIM*.

- Benamou, J.-D., Carlier, G., Cuturi, M., Nenna, L., and Peyré, G. (2015). Iterative Bregman projections for regularized transportation problems. *SISC*.
- Courty, N., Flamary, R., Tuia, D., and Rakotomamonjy, A. (2016).
  Optimal transport for domain adaptation.

Pattern Analysis and Machine Intelligence, IEEE Transactions on.

# References ii

Cuturi, M. (2013).

Sinkhorn distances: Lightspeed computation of optimal transportation. In Neural Information Processing Systems (NIPS), pages 2292–2300.

🚺 Cuturi, M. and Blondel, M. (2017).

Soft-DTW: a differentiable loss function for time-series.

volume 70, pages 894–903, International Convention Centre, Sydney, Australia. PMLR.

Damodaran, B. B., Kellenberger, B., Flamary, R., Tuia, D., and Courty, N. (2018).

Deepjdot: Deep joint distribution optimal transport for unsupervised domain adaptation.

🔋 Ferradans, S., Papadakis, N., Peyré, G., and Aujol, J.-F. (2014a).

Regularized discrete optimal transport.

SIAM Journal on Imaging Sciences, 7(3).

# References iii

Ferradans, S., Papadakis, N., Peyré, G., and Aujol, J.-F. (2014b). Regularized discrete optimal transport. SIAM Journal on Imaging Sciences, 7(3):1853–1882.

Kantorovich, L. (1942).

On the translocation of masses.

C.R. (Doklady) Acad. Sci. URSS (N.S.), 37:199-201.

Lacoste-Julien, S. (2016).

**Convergence rate of frank-wolfe for non-convex objectives.** *arXiv preprint arXiv:1607.00345.* 

Loosli, G., Canu, S., and Ong, C. S. (2016).

Learning svm in krein spaces.

IEEE transactions on pattern analysis and machine intelligence, 38(6):1204–1216.

# References iv

Mansour, Y., Mohri, M., and Rostamizadeh, A. (2009). **Domain adaptation: Learning bounds and algorithms.** In *Proc. of COLT* 

Memoli, F. (2011).

**Gromov wasserstein distances and the metric approach to object matching.** *Foundations of Computational Mathematics*, pages 1–71.

```
Monge, G. (1781).
```

Mémoire sur la théorie des déblais et des remblais.

De l'Imprimerie Royale.

Peng, X., Usman, B., Kaushik, N., Hoffman, J., Wang, D., and Saenko, K. (2017).

Visda: The visual domain adaptation challenge.

arXiv preprint arXiv:1710.06924.

# References v

- Perrot, M., Courty, N., Flamary, R., and Habrard, A. (2016).
  Mapping estimation for discrete optimal transport.
  In Neural Information Processing Systems (NIPS).
- Petitjean, F., Ketterlin, A., and Gançarski, P. (2011).

A global averaging method for dynamic time warping, with applications to clustering.

44(3):678-693.

Peyré, G., Cuturi, M., and Solomon, J. (2016).

Gromov-Wasserstein Averaging of Kernel and Distance Matrices. In *ICML 2016*, Proc. 33rd International Conference on Machine Learning, New-York, United States.

Rakotomamonjy, A., Traore, A., Berar, M., Flamary, R., and Courty, N. (2018). Wasserstein Distance Measure Machines. preprint.

# References vi

Rubner, Y., Tomasi, C., and Guibas, L. J. (2000).

The earth mover's distance as a metric for image retrieval.

International journal of computer vision, 40(2):99–121.

Seguy, V., Bhushan Damodaran, B., Flamary, R., Courty, N., Rolet, A., and Blondel, M. (2017).

Large-scale optimal transport and mapping estimation.

Thorpe, M., Park, S., Kolouri, S., Rohde, G. K., and Slepcev, D. (2017).

A transportation lp distance for signal analysis.

Journal of Mathematical Imaging and Vision, 59(2):187–210.

Tuia, D., Flamary, R., Rakotomamonjy, A., and Courty, N. (2015).

Multitemporal classification without new labels: a solution with optimal transport.

In 8th International Workshop on the Analysis of Multitemporal Remote Sensing Images.

- Urner, R., Shalev-Shwartz, S., and Ben-David, S. (2011).
  Access to unlabeled data can speed up prediction time.
  In *Proceedings of ICML*, pages 641–648.
- Vayer, T., Chapel, L., Flamary, R., Tavenard, R., and Courty, N. (2018).
  Fused gromov-wasserstein distance for structured objects: theoretical foundations and mathematical properties.
- Venkateswara, H., Eusebio, J., Chakraborty, S., and Panchanathan, S. (2017).
  Deep hashing network for unsupervised domain adaptation.
  In (IEEE) Conference on Computer Vision and Pattern Recognition (CVPR).