

Université Jean Monnet de Saint-Etienne, Université de Lyon, Laboratoire Hubert Curien, CNRS, UMR5516, F-42000, Saint-Etienne, France. {michael.perrot,amaury.habrard}@univ-st-etienne.fr

METRIC LEARNING: MAIN IDEA

Learning how to compare objects: learn a new space where some constraints are fulfilled, e.g. move closer circles of the same color (class) and keep far away circles of different colors (classes).



Mahalanobis-like Distance (with \mathbf{M} a PSD matrix):

$$d_{\mathbf{M}}(\mathbf{x}, \mathbf{x}') = \sqrt{(\mathbf{x} - \mathbf{x}')^T \mathbf{M}(\mathbf{x} - \mathbf{x}')}$$
$$= \sqrt{(\mathbf{L}^T \mathbf{x} - \mathbf{L}^T \mathbf{x}')^T (\mathbf{L}^T \mathbf{x} - \mathbf{L}^T \mathbf{x}')}$$
$$= \sqrt{\mathbf{x}^T \mathbf{M} \mathbf{x} + {\mathbf{x}'}^T \mathbf{M} \mathbf{x}' - 2\mathbf{x}^T \mathbf{M} \mathbf{x}'}$$

METRIC LEARNING: CLASSICAL APPROACH

Classical approaches in metric learning use two kind of constraints:

- Similarity constraints: the goal is to bring closer similar examples, e.g. examples of the same class.
- Dissimilarity constraints: the goal is to push far away dissimilar examples, e.g. examples of different classes.





It induces a quadratic number of constraints: $\mathcal{O}(n^2)$ (e.g. [JWZ09]).

METRIC LEARNING: VIRTUAL POINTS APPROACH

In our virtual points approach we only use similarity constraints. However instead of pairing examples with each other, we associate each example to a given virtual point.



It induces a linear number of constraints: $\mathcal{O}(n)$.

Regressive Virtual Metric Learning

Michaël Perrot and Amaury Habrard

Objective: Propose a new framework for metric learning using only positive constraints between examples and a priori defined virtual points.

FORMULATION

Let $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \subset (\mathcal{X} \times \mathcal{Y})^n$ be a set of examples. Let $f_{\mathbf{v}} : \mathcal{X} \times \mathcal{Y} \to \mathcal{V}$ where $\mathcal{V} \subseteq \mathbb{R}^{d'}$ be the function which associates each example to a virtual point. We consider the learning set $S_{\mathbf{v}} = \{(\mathbf{x}_i, \mathbf{v}_i)\}_{i=1}^n \subset \mathcal{X} \times \mathcal{V}.$ Let $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$ and $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n)^T$, we learn \mathbf{L} such that $\mathbf{M} = \mathbf{L}\mathbf{L}^T$ with the following optimization problem:

$$\min_{\mathbf{L}} \hat{R}(\mathbf{L}) + \lambda \|\mathbf{L}\|_{\mathcal{F}}^2 = \min_{\mathbf{L}} \frac{1}{n} \|\mathbf{X}\mathbf{L} - \mathbf{V}\|_{\mathcal{F}}^2 + \lambda \|\mathbf{L}\|_{\mathcal{F}}^2.$$
(1)

Using the closed form solution for \mathbf{L} , we get:

$$\mathbf{M} = \mathbf{L}\mathbf{L}^{T} = \mathbf{X}^{T} \left(\mathbf{X}\mathbf{X}^{T} + \lambda n\mathbf{I}\right)^{-1} \mathbf{V}\mathbf{V}^{T} \left(\mathbf{X}\mathbf{X}^{T} + \lambda n\mathbf{I}\right)^{-1} \mathbf{X}.$$
 (2)

Let $K_{\mathbf{X}} = \phi(\mathbf{X})\phi(\mathbf{X})^T$, we can kernelize our approach as:

 $d_{\mathbf{M}_{K}}^{2}(\phi(\mathbf{x}),\phi(\mathbf{x}')) = \phi(\mathbf{x})^{T}\mathbf{M}_{K}\phi(\mathbf{x}) + \phi(\mathbf{x}')^{T}\mathbf{N}_{K}\phi(\mathbf{x})$

with $\mathbf{M}_K = \phi(\mathbf{X})^T (K_{\mathbf{X}} + \lambda n \mathbf{I})^{-1} \mathbf{V} \mathbf{V}^T (K_{\mathbf{X}} + \lambda n \mathbf{I})^{-1} \mathbf{V} \mathbf{V}^T$

THEORETICAL ANALYSIS: BOUNDING THE TRUE RISK OF [JWZ09] BY THE EMPIRICAL RISK OF OUR APPROACH

Theorem 1. Let \mathcal{D} be a distribution over $\|\mathbf{v}\|_2 \leq C_{\mathbf{v}}$ for any $\mathbf{v} \in \mathcal{V}$ and $\|\mathbf{x}\|_2 \leq C_{\mathbf{x}}$

$$\mathbb{E}_{(\mathbf{x}_i, y_i) \sim \mathcal{D}, (\mathbf{x}_j, y_j) \sim \mathcal{D}} \left[y_{ij} (d^2 (\mathbf{L}^T \mathbf{x}_i, \mathbf{L}^T)) \right]$$

with $y_{ij} = 1$ for examples of the same class

EXPERIMENTS

		Baselines						Our approach		
Base	1NN		LMNN		SCML		R	VML-Lin-OT	RVML-Lin-Class	
$Amazon 41.51 \pm 3$		$5.24 65.50 \pm$		2.28 71.		68 ± 1.86	7	71.62 ± 1.34	$\textbf{73.09} \pm \textbf{2.49}$	
Caltech 18.04 ± 2		$20 49.68 \pm 2.$		2.76	52.84 ± 1.61		Ę	52.51 ± 2.41	$\textbf{55.41} \pm \textbf{2.55}^{*}$	
DSLR	$\boxed{\text{DSLR}} \qquad 29.61 \pm 4$		$\boxed{38 \textbf{76.08} \pm 4}$		65.10 ± 9.00		7	74.71 ± 5.27	75.29 ± 5.08	
Isolet	Isolet 88.97		95.8			89.61		91.40	94.61	
Letters	Letters 94.74 ± 0		$\textbf{96.43}~\pm$	0.28^{*}	96.13 ± 0.20		Ģ	90.25 ± 0.60	95.51 ± 0.26	
Splice	Splice 71.17		82.0)2	85.43			84.64	78.44	
Svmguide	1 95.12		95.03		87.38			94.83	85.25	
Webcam	42.90 ± 4	.19	85.81 ± 3.75		$\textbf{90.43} \pm \textbf{2.70}$		8	88.60 ± 3.63	88.60 ± 2.69	
					-					
		Baselines						Our approach		
Base	1NN-KPCA	LMN	IN-KPCA	GBLM	[NN	SCMLLOCA	AL	RVML-RBF-OT	RVML-RBF-Class	
Amazon	20.27 ± 2.42	53.1	16 ± 3.73	$73 65.53 \pm$		69.14 ± 1.74		73.51 ± 0.83	$76.22 \pm 2.09^{*}$	
Caltech	20.82 ± 8.29	29.8	8 ± 10.89	$49.91 \pm$	2.80	50.56 ± 1.62		54.39 ± 1.89	$57.98 \pm 2.22^{*}$	
DSLR	64.90 ± 5.81	73.9	92 ± 7.57 $76.08 \pm$		4.79	62.55 ± 6.94		70.39 ± 4.48	$\textbf{76.67} \pm \textbf{4.57}$	
Isolet	68.70		96.28 96.0		2	91.40		95.96	96.73	
Letter	95.39 ± 0.27	97.1	$17^* \pm 0.18$ 96.51		$= 0.25 \underline{96.63 \pm 0.2}$		26	91.26 ± 0.50	96.09 ± 0.21	
Splice	66.99	8	88.97	82.2	1	87.13		88.51	88.32	
Svmguide1	de1 95.72		95.60		0	87.40		95.67	95.05	
Webcam	73.55 ± 4.57	84.52 ± 3.83		85.81 ± 3.75		88.71 ± 2.83		88.71 ± 4.28	88.92 ± 2.91	

$$\mathbf{M}_{K}\phi(\mathbf{x}') - 2\phi(\mathbf{x})^{T}\mathbf{M}_{K}\phi(\mathbf{x}')$$

$$(\mathbf{X}_{\mathbf{X}} + \lambda n \mathbf{I})^{-1} \phi(\mathbf{X}).$$



Class Based Representation Space Approach The virtual points are defined as \blacklozenge unit vectors of a space of dimension the number of classes, i.e. there is one class associated with one virtual point. Each example is associated with a virtual point using its class.

$\mathcal{X} \times \mathcal{Y}$. Let $\mathcal{V} \subset \mathbb{R}^{d'}$ be a finite set of	virtual points and $f_{\mathbf{v}}$ is	s defined as $f_{\mathbf{v}}(\mathbf{x}_i, y_i) = \mathbf{v}_i$,	$\mathbf{v}_i \in \mathcal{V}$. Let
for any $\mathbf{x} \in \mathcal{X}$. Let $\gamma_1 = 2 \max_{\mathbf{x}_k, \mathbf{x}_l, y_k}$	$a_{kl=1} d^2(\mathbf{v}_k, \mathbf{v}_l) and \gamma_{-1}$	$= \frac{1}{2} \min_{\mathbf{x}_k, \mathbf{x}_l, y_{kl} = -1} d^2(\mathbf{v}_k, \mathbf{v}_k)$	(r_l) , we have:
$[\mathbf{x}_j) - \gamma_{y_{ij}}]_+ \le 8 \left(\hat{R}(\mathbf{L}) + \frac{8C_{\mathbf{v}}^2 C_{\mathbf{x}}^2}{\lambda n} \right) $	$1 + \frac{C_{\mathbf{x}}}{\sqrt{\lambda}} \bigg)^2 + \left(\left(\frac{16C_{\mathbf{x}}^2}{\lambda} + \right)^2 \right)^2 + \left(\frac{16C_{\mathbf{x}}^2}{\lambda} + \right)^2 \right)^2 + \left(\frac{16C_{\mathbf{x}}^2}{\lambda} + \right)^2 + \left(\frac{16C_{\mathbf{x}}^2}{\lambda} + \right)^2 + \left(\frac{16C_{\mathbf{x}}^2}{\lambda} + \right)^2 \right)^2 + \left(\frac{16C_{\mathbf{x}}^2}{\lambda} + \left(\frac{16C_{\mathbf{x}}^2}{\lambda} + \right)^2 + \left(\frac{16C_{\mathbf{x}}^2}{\lambda} + \left(\frac{16C_{\mathbf{x}}^2}{\lambda} + \right)^2 + \left(\frac{16C_{\mathbf{x}}^2}{\lambda} + \left$	$+1 C_{\mathbf{v}}^2 \left(1 + \frac{C_{\mathbf{x}}}{\sqrt{\lambda}}\right)^2 \sqrt{\frac{\ln \frac{1}{2\pi}}{2\pi}}$	$\left(\frac{1}{\delta}{\delta}\right)$.
and-1 otherwise. Note that the marg	ins are expressed w.r.t.	the distances between virtu	al points.





REFERENCES

[BHS15]	Aurélien Bellet, Amaury H & Claypool Publishers, 201
[CFT14]	Nicolas Courty, Rémi Flan ularized optimal transport
[JWZ09]	Rong Jin, Shijun Wang, an Theory and algorithm. In
[KJ11]	Purushottam Kar and Pra- embeddings. In <i>Proc. of N</i>



Optimal Transport Based Approach This is a two steps approach:

- Selection of several landmarks in the training set using a diversity criteria |KJ11|.
- Application of Optimal Transport with regularization [CFT14] between the training set and the landmarks. The points obtained after transport are used as virtual points.



- Iabrard, and Marc Sebban. *Metric Learning*. Morgan
- mary, and Devis Tuia. Domain adaptation with regt. In *Proc.* of ECML, pages 274–289, 2014.
- nd Yang Zhou. Regularized distance metric learning: Proc. of NIPS, pages 862–870, 2009.
- ateek Jain. Similarity-based learning via data driven NIPS, pages 1998–2006, 2011.