

ICML | 2019

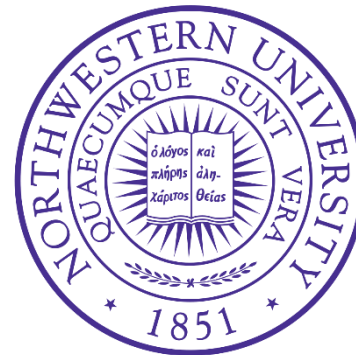
Thirty-sixth International Conference on
Machine Learning

Learning Distance for Sequences by Learning a Ground Metric

Bing Su



Ying Wu

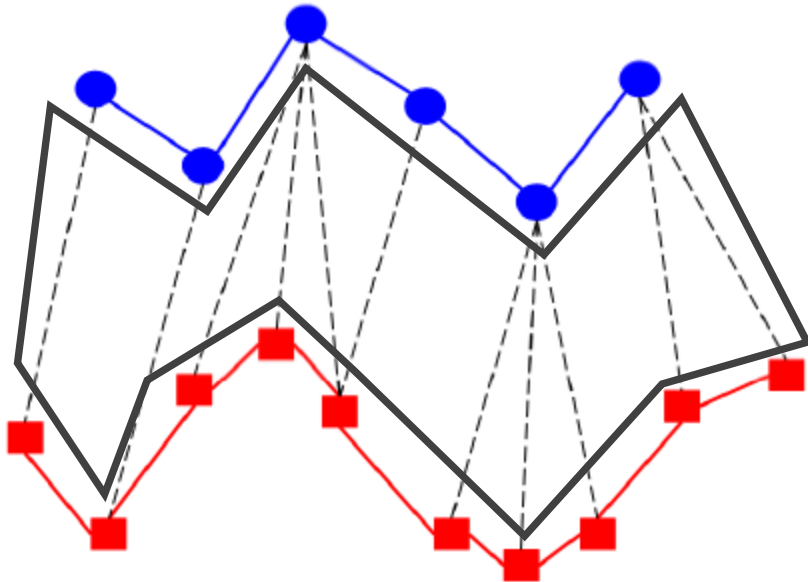


Motivation

- Distance between sequences depends on *temporal alignment* to eliminate the local temporal discrepancies.

Temporal alignment

$$X = [x_1, \dots, x_{L_X}] \in \Omega^{L_X}$$

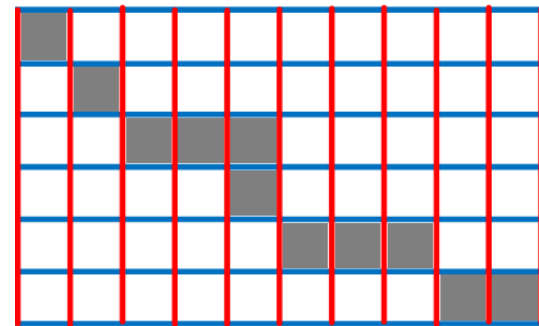


$$Y = [y_1, \dots, y_{L_Y}] \in \Omega^{L_Y}$$

$t_{i,j} = T(i,j)$ indicates whether or the probability of the pair x_i and y_j is aligned.



T



Motivation

- The inference of alignment depends on the *ground metric* between elements in sequences.

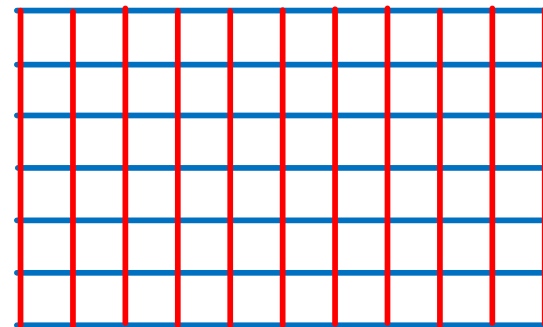
Let Ω be a space, $d(M) : \Omega \times \Omega \rightarrow \mathbb{R}$ be the metric on this space.



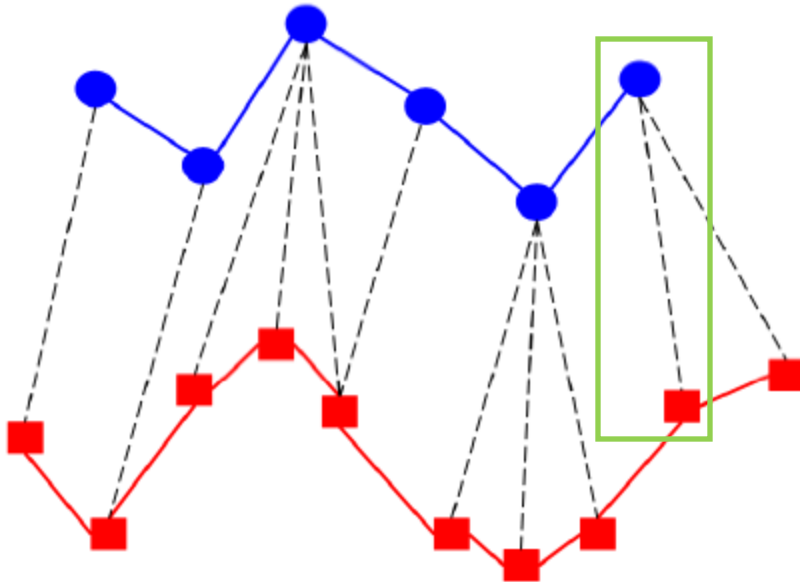
Ground metric



$$D(M) := [d(M, \mathbf{x}_i, \mathbf{y}_j)]_{ij} \in \mathbb{R}^{L_X \times L_Y}$$



$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_{L_X}] \in \Omega^{L_X}$$



$$\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_{L_Y}] \in \Omega^{L_Y}$$

A Unified Perspective

- Distance between two sequences: a general formulation

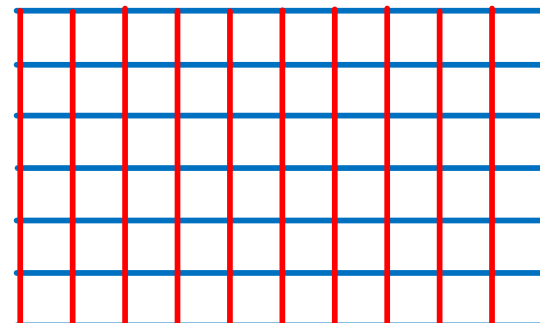
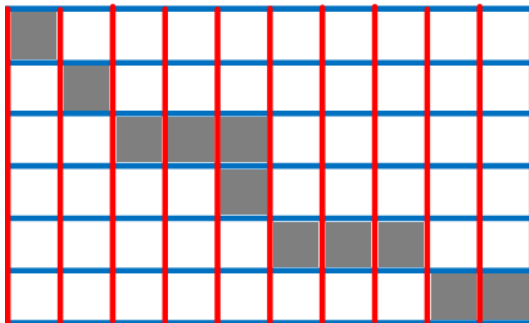
$$g_M(X, Y) = \langle T^*, D(M) \rangle$$

T

$D(M)$

The *temporal alignment* matrix

The *ground metric* matrix of pairwise distances between elements



A Unified Perspective

- T^* is generally inferred by

$$T^* = \underset{T \in \Phi}{\operatorname{arg\,min}} \langle T, D(M) \rangle + \mathcal{R}(T)$$

- Φ is the feasible set of T , which is a subset of $\mathbb{R}^{L_X \times L_Y}$ with some constraints; $\mathcal{R}(T)$ is a regularization term.
- Different distance measures for sequences differ in the constraints imposed to the feasible set, the regularization term, and the optimization method.

A Unified Perspective

- Connection to *dynamic time warping (DTW)*

$$\mathcal{R}(\mathbf{T}) = 0;$$

$$\Phi = \{\mathbf{T} \in \{0, 1\}^{L_X \times L_Y} \mid \mathbf{T}_{1,1} = 1, \mathbf{T}_{L_X, L_Y} = 1; \mathbf{T}\mathbf{1}_{L_Y} > \mathbf{0}_{L_X}, \mathbf{T}^T\mathbf{1}_{L_X} > \mathbf{0}_{L_Y}; \\ \text{if } t_{i,j} = 1, \text{ then } t_{i-1,j+1} = 0, t_{i+1,j-1} = 0, \forall 1 < i < L_X, 1 < j < L_Y\}$$

DTW infers \mathbf{T} via dynamic programming.

- Connection to *order-preserving Wasserstein distance (OPW)*

$$\mathcal{R}(\mathbf{T}) = \lambda_1 I(\mathbf{T}) + \lambda_2 KL(\mathbf{T} \parallel \mathbf{P});$$

$$\Phi = \{\mathbf{T} \in \mathbb{R}_+^{L_X \times L_Y} \mid \mathbf{T}\mathbf{1}_{L_Y} = \frac{1}{L_X}\mathbf{1}_{L_X}, \mathbf{T}^T\mathbf{1}_{L_X} = \frac{1}{L_Y}\mathbf{1}_{L_Y}\}$$

OPW infers \mathbf{T} by the Sinkhorn's matrix scaling algorithm.

Problem

- The distance between sequences is formulated as a function of the ground metric: *meta-distance*
- *Learn meta-distance by learning the ground metric*
- Given a set of N training sequences and the corresponding labels, $\{\mathbf{X}^n, z^n\}_{n=1}^N$ $\mathbf{X}^n = [\mathbf{x}_1, \dots, \mathbf{x}_{L^n}] \in \mathbb{R}^{b \times L^n}$
- Learn a meta-distance $g_M(\mathbf{X}^n, \mathbf{X}^{n'})$ by learning a Mahalanobis distance as the ground metric:

$$d(\mathbf{M}, \mathbf{x}_i, \mathbf{y}_j) = (\mathbf{x}_i - \mathbf{y}_j)^T \mathbf{M} (\mathbf{x}_i - \mathbf{y}_j)$$

- $\mathbf{M} = \mathbf{W}\mathbf{W}^T$, $\mathbf{W} \in \mathbb{R}^{b \times b'}$
- Goal: with the learned \mathbf{W} , the resulting meta-distance

$$g_M(\mathbf{X}^n, \mathbf{X}^{n'}) = g_I(\mathbf{W}^T \mathbf{X}^n, \mathbf{W}^T \mathbf{X}^{n'})$$

better separates sequences from different classes.

Objective

- *Regressive virtual sequence metric learning (RVSML)*
- Associate a virtual sequence $V^n = [v_1, \dots, v_{l^n}] \in \mathbb{R}^{b' \times l^n}$ with each training sequence X^n
- Minimize the meta-distances between the training sequences and their associated virtual sequences

$$\begin{aligned} \min_{\mathbf{W}} \frac{1}{N} \sum_{n=1}^N g_I(\mathbf{W}^T \mathbf{X}^n, V^n) + \beta \|\mathbf{W}\|_{\mathcal{F}}^2 \\ = \frac{1}{N} \sum_{n=1}^N \langle \mathbf{T}^{n*}, D_I^n(\mathbf{W}) \rangle + \beta \|\mathbf{W}\|_{\mathcal{F}}^2 \\ \text{s.t. } \mathbf{T}^{n*} = \arg \min_{\mathbf{T} \in \Phi} \langle \mathbf{T}^n, D_I^n(\mathbf{W}) \rangle + \mathcal{R}(\mathbf{T}^n) \end{aligned}$$

- If $\mathcal{R}(\mathbf{T})$ does not depend on \mathbf{W} , it is equivalent to

$$\min_{\mathbf{W}, \mathbf{T}^n} \frac{1}{N} \sum_{n=1}^N \langle \mathbf{T}^n, D_I^n(\mathbf{W}) \rangle + \beta \|\mathbf{W}\|_{\mathcal{F}}^2 + \mathcal{R}(\mathbf{T}^n)$$

Optimization

$$\min_{\mathbf{W}, \mathbf{T}^n} \frac{1}{N} \sum_{n=1}^N \langle \mathbf{T}^n, \mathbf{D}_I^n(\mathbf{W}) \rangle + \beta \|\mathbf{W}\|_{\mathcal{F}}^2 + \mathcal{R}(\mathbf{T}^n)$$

- Fix \mathbf{T}^n , optimize \mathbf{W} : standard regression, closed form solution

$$\mathbf{W}^* = \mathbf{A}^{-1} \left(\sum_{n=1}^N \sum_{i=1}^{L^n} \sum_{j=1}^{l^n} t_{ij}^n \mathbf{x}_i^n \mathbf{v}_j^{nT} \right) \quad \mathbf{A} = \sum_{n=1}^N \sum_{i=1}^{L^n} \sum_{j=1}^{l^n} t_{ij}^n \mathbf{x}_i^n \mathbf{x}_i^{nT} + \beta N \mathbf{I}$$

- Fix \mathbf{W} , optimize \mathbf{T}^n : standard inference, e.g. DTW, OPW

$$\mathbf{T}^{n*} = \underset{\mathbf{T}^n \in \Phi}{\text{arg min}} \langle \mathbf{T}^n, \mathbf{D}_I^n(\mathbf{W}) \rangle + \mathcal{R}(\mathbf{T}^n)$$

- Guaranteed convergence

Evaluation

- Generating $V^n = f(X^n, z^n) = [e_{(z^n-1)m+1}, \dots, e_{(z^n-1)m+m}]$
- RVSML instantiated by (a) DTW and (b) OPW using the NN classifier with the (a) DTW and (b) OPW distance
- Comparison with other metric learning methods on the ChaLearn and SAD datasets

(a) DTW

Method	MAP	Accuracy
Ori (Su & Hua, 2018)	11.75	61.12
ITML (Davis et al., 2007)	13.46	52.17
LMNN (Weinberger & Saul, 2009)	11.67	63.78
RVML (Perrot & Habrard, 2015)	31.21	83.79
LDMLT (Mei et al., 2014)	21.30	84.37
SWMD (Huang et al., 2016)	14.39	64.45
RVSML	33.83	87.38

(b) OPW

Method	MAP	Accuracy
Ori (Su & Hua, 2018)	12.21	59.38
ITML (Davis et al., 2007)	13.92	64.71
LMNN (Weinberger & Saul, 2009)	12.07	62.83
RVML (Perrot & Habrard, 2015)	30.19	80.66
LDMLT (Mei et al., 2014)	21.56	82.74
SWMD (Huang et al., 2016)	15.36	60.31
RVSML	33.07	83.82

(a) DTW

Method	MAP	Accuracy
Ori (Su & Hua, 2018)	56.58	96.36
ITML (Davis et al., 2007)	51.13	95.55
LMNN (Weinberger & Saul, 2009)	56.25	96.00
SCML (Shi et al., 2014)	47.98	93.27
RVML (Perrot & Habrard, 2015)	57.94	96.59
LDMLT (Mei et al., 2014)	59.54	96.50
SWMD (Huang et al., 2016)	52.44	93.95
RVSML	60.24	96.23

(b) OPW

Method	MAP	Accuracy
Ori (Su & Hua, 2018)	59.77	96.36
ITML (Davis et al., 2007)	54.51	96.36
LMNN (Weinberger & Saul, 2009)	59.33	96.27
SCML (Shi et al., 2014)	50.08	94.50
RVML (Perrot & Habrard, 2015)	60.71	95.77
LDMLT (Mei et al., 2014)	61.07	96.73
SWMD (Huang et al., 2016)	58.00	95.41
RVSML	65.63	97.09

Results

- Comparison with state-of-the-art methods on the MSR Activity3D and MSR Action3D datasets

Method	Accuracy
Actionlet Ensemble (Wang et al., 2012)	85.8%
Moving Pose (Zanfir et al., 2013)	73.8%
COV- $J_{\mathcal{H}}$ -SVM (Harandi et al., 2014)	75.5%
Ker-RP-POL (Wang et al., 2015)	96.9%
Ker-RP-RBF (Wang et al., 2015)	96.3%
Kernelized-COV (Cavazza et al., 2016)	96.3%
Luo et al. (Luo et al., 2017)	86.9%
Ji et al. (Ji et al., 2018)	81.3%
DSSCA SSLM (Shahroudy et al., 2018)	97.5%
RVSML-DTW+Kernelized-COV	96.9%
RVSML-OPW+Kernelized-COV	97.5%

Method	Accuracy
Actionlet Ensemble (Wang et al., 2012)	88.2%
Moving Pose (Zanfir et al., 2013)	91.7%
COV- $J_{\mathcal{H}}$ -SVM (Harandi et al., 2014)	80.4%
Ker-RP-POL (Wang et al., 2015)	96.2%
Ker-RP-RBF (Wang et al., 2015)	96.9%
Kernelized-COV (Cavazza et al., 2016)	96.2%
SCK+DCK (Koniusz et al., 2016)	91.45%
TS-LSTM-GM (Lee et al., 2017)	91.21%
FTP-SVM (Ben Tanfous et al., 2018)	90.01%
Bi-LSTM (Ben Tanfous et al., 2018)	86.18%
RVSML-DTW+Kernelized-COV	82.78%
RVSML-OPW+Kernelized-COV	96.34%
RVSML-DTW+TS-LSTM-GM	93.04%
RVSML-OPW+TS-LSTM-GM	90.48%

- Please visit our poster for more details.
- Thank you very much!