Soft-DTW: A Differentiable Loss Function for Time Series

Marco Cuturi UNIVERSITE PARIS-SACLAY



Mathieu Blondel



Follow pollution levels in real time in your city



From: Plume App

Ground truth (reality)

Follow pollution levels in real time in your city



From: Plume App

Follow pollution levels in real time in your city



How wrong was this prediction?

This depends on the loss function used to train the algorithm.

From: Plume App

Ground truth (reality)

Follow pollution levels in real time in your city

- In this talk we propose to use the celebrated
 Dynamic Time Warping discrepancy as a loss.
- Loss functions should be differentiable. We show that an appropriate smoothing, soft-DTW, helps
- We apply this to **several problems**:
 - Computation of barycenters,
 - Clustering of time series,
- G •• Learning with structured (time series) output him. (reality)

From: Plume App

0. The DTW Geometry

1. Soft-DTW

2. Soft-DTW as a Loss Function

Dynamic Time Warping

A discrepancy function between two **time series of observations** supported **on a metric space**.





Pairwise Distance Matrix



Pairwise Distance Matrix



Pairwise Distance Matrix





































 $Cost = \langle A, \Delta \rangle, A \in \{0, 1\}^{n \times m}$

Min Cost Alignment Matrix?



Min Cost Alignment Matrix?



Set of all valid path matrices: $\mathcal{A}(\boldsymbol{n}, \boldsymbol{m}) \subset \{0, 1\}^{\boldsymbol{n} \times \boldsymbol{m}}$

Best Alignment Matrix



Set of all valid path matrices: $\mathcal{A}(\boldsymbol{n}, \boldsymbol{m}) \subset \{0, 1\}^{\boldsymbol{n} \times \boldsymbol{m}}$

Best Alignment Matrix



Set of all valid path matrices: $\mathcal{A}(\boldsymbol{n}, \boldsymbol{m}) \subset \{0, 1\}^{\boldsymbol{n} \times \boldsymbol{m}}$



 $r_{1,1} = \Delta_{11} \qquad r_{0,j} = r_{i,0} = 0$



 $r_{i,j} = \min(r_{i-1,j-1}, r_{i-1,j}, r_{i,j-1}) + \Delta_{i,j}$



 $r_{i,j} = \min(r_{i-1,j-1}, r_{i-1,j}, r_{i,j-1}) + \Delta_{i,j}$



$$\operatorname{dtw}_0(\boldsymbol{X},\boldsymbol{Y}) = \boldsymbol{r_{n,m}}$$

Optimal Path



 $\mathbf{dtw}_0(\boldsymbol{X},\boldsymbol{Y}) = \boldsymbol{r_{n,m}} = \langle A^\star, \boldsymbol{\Delta} \rangle$

Optimal Path



0. The DTW Geometry

1. Soft-DTW

2. Soft-DTW as a Loss Function







$$\mathbf{dtw}_0(\boldsymbol{X},\boldsymbol{Y}) = \min_{A \in \mathcal{A}(\boldsymbol{n},\boldsymbol{m})} \langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

$$\mathbf{dtw}_0(\boldsymbol{X}, \boldsymbol{Y}) = \min_{A \in \mathcal{A}(\boldsymbol{n}, \boldsymbol{m})} \langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

$$\begin{pmatrix} \frac{\partial \operatorname{dtw}_{0}(\boldsymbol{X},\boldsymbol{Y})}{\partial \boldsymbol{X}} \end{pmatrix}^{T} = \begin{pmatrix} \frac{\partial \operatorname{dtw}_{0}(\boldsymbol{X},\boldsymbol{Y})}{\partial \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}}} & \frac{\partial \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}}}{\partial \boldsymbol{X}} \end{pmatrix}^{T}$$
$$\downarrow$$
$$\mathbb{R}^{\boldsymbol{n} \times \boldsymbol{m}} \to \mathbb{R} \qquad \mathbb{R}^{|\Omega| \times \boldsymbol{n}} \to \mathbb{R}^{\boldsymbol{n} \times \boldsymbol{m}}$$

$$\mathbf{dtw}_0(\boldsymbol{X}, \boldsymbol{Y}) = \min_{A \in \mathcal{A}(\boldsymbol{n}, \boldsymbol{m})} \langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

$$\nabla_{\mathbf{X}} \operatorname{dtw}_{0}(\mathbf{X}, \mathbf{Y}) = \begin{pmatrix} \partial \Delta_{\mathbf{X}\mathbf{Y}} \\ \partial \mathbf{X} \end{pmatrix}^{T} \nabla_{\Delta} \min_{\mathcal{A}(n,m)} \langle \cdot, \Delta_{\mathbf{X}\mathbf{Y}} \rangle$$
$$= A^{\star}$$
$$\mathbb{R}^{n \times m} \to \mathbb{R}^{|\Omega| \times n} \quad \text{iff optimal solution}$$
is unique

$$\mathbf{dtw}_0(\boldsymbol{X}, \boldsymbol{Y}) = \min_{A \in \mathcal{A}(\boldsymbol{n}, \boldsymbol{m})} \langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

- \mathbf{dtw}_0 is piecewise **linear** w.r.t Δ
- if $\Delta_{ij} = \delta(x_i, y_j) = ||x_i y_j||^2$, dtw_0 is piecewise quadratic w.r.t. X.

$$\mathbf{dtw}_0(\boldsymbol{X}, \boldsymbol{Y}) = \min_{A \in \mathcal{A}(\boldsymbol{n}, \boldsymbol{m})} \langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

• If A^* is unique,

$$\nabla_X \operatorname{\mathbf{dtw}}_0(\boldsymbol{X}, \boldsymbol{Y}) = \left(\frac{\partial \boldsymbol{\Delta}(\boldsymbol{X}, \boldsymbol{Y})}{\partial \boldsymbol{X}}\right)^T A^*$$

• dtw_0 has a discontinuous gradient.

Any way to fix this?

$$\mathbf{dtw}_0(\boldsymbol{X}, \boldsymbol{Y}) = \min_{A \in \mathcal{A}(\boldsymbol{n}, \boldsymbol{m})} \langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

Problem: non-differentiability of min operator over finite family of values.

Fix

$$\min^{\gamma}(u_1, \dots, u_n) = \begin{cases} \min_{i \le n} u_i, & \gamma = 0, \\ -\gamma \log \sum_{i=1}^n e^{-u_i/\gamma}, & \gamma > 0. \end{cases}$$

Example softmin of quadratic functions

Example softmin of quadratic functions

Soft-DTW

$$\mathbf{dtw}_0(\boldsymbol{X}, \boldsymbol{Y}) = \min_{A \in \mathcal{A}(\boldsymbol{n}, \boldsymbol{m})} \langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

Fix: Replace min by $\min^{\gamma}, \gamma > 0$

$$\mathbf{dtw}_{\gamma}(\boldsymbol{X},\boldsymbol{Y}) = \min_{A \in \mathcal{A}(\boldsymbol{n},\boldsymbol{m})}^{\gamma} \langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

$$\operatorname{dtw}_{\gamma}(\boldsymbol{X},\boldsymbol{Y}) = -\gamma \log \sum_{A \in \mathcal{A}(\boldsymbol{n},\boldsymbol{m})} e^{-\frac{\langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle}{\gamma}}$$

Soft-DTW

$$\mathbf{dtw}_0(\boldsymbol{X}, \boldsymbol{Y}) = \min_{A \in \mathcal{A}(\boldsymbol{n}, \boldsymbol{m})} \langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

Fix: Replace min by $\min^{\gamma}, \gamma > 0$

$$\mathbf{dtw}_{\gamma}(\boldsymbol{X},\boldsymbol{Y}) = \min_{A \in \mathcal{A}(\boldsymbol{n},\boldsymbol{m})}^{\gamma} \langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

$$\mathbf{dtw}_{\gamma}(\boldsymbol{X}, \boldsymbol{Y}) = -\gamma \log \sum_{A \in \mathcal{A}(\boldsymbol{n}, \boldsymbol{m})} e^{-rac{\langle A, \boldsymbol{\Delta}_{\boldsymbol{X} \boldsymbol{Y}}
angle}{\gamma}}$$

Recursive Computation

$$\mathbf{dtw}_0(\boldsymbol{X}, \boldsymbol{Y}) = \min_{A \in \mathcal{A}(\boldsymbol{n}, \boldsymbol{m})} \langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

$$r_{i,j} = \min(r_{i-1,j-1}, r_{i-1,j}, r_{i,j-1}) + \Delta_{i,j}$$

$$\mathbf{dtw}_{\gamma}(\boldsymbol{X},\boldsymbol{Y}) = \min_{A \in \mathcal{A}(\boldsymbol{n},\boldsymbol{m})}^{\gamma} \langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

$$r_{i,j} = \min^{\gamma} (r_{i-1,j-1}, r_{i-1,j}, r_{i,j-1}) + \Delta_{i,j}$$

$$\mathbf{dtw}_{\gamma}(\boldsymbol{X},\boldsymbol{Y}) = \min_{A \in \mathcal{A}(\boldsymbol{n},\boldsymbol{m})} \langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

$$\operatorname{dtw}_{\gamma}(\boldsymbol{X},\boldsymbol{Y}) = \min_{A \in \mathcal{A}(\boldsymbol{n},\boldsymbol{m})}^{\gamma} \langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

Recursive Computation (Backward)

$$r_{i,j} = \min^{\gamma} (r_{i-1,j-1}, r_{i-1,j}, r_{i,j-1}) + \Delta_{i,j}$$

$$\frac{\partial \operatorname{dtw}_{\gamma}(\boldsymbol{X},\boldsymbol{Y})}{\partial \boldsymbol{X}} = \frac{\partial \operatorname{dtw}_{\gamma}(\boldsymbol{X},\boldsymbol{Y})}{\partial \Delta} \frac{\partial \Delta}{\partial \boldsymbol{X}}$$

$$(\boldsymbol{X},\boldsymbol{Y}) \quad \begin{bmatrix} \partial \boldsymbol{r}_{n} \ m \end{bmatrix} \quad \begin{bmatrix} \partial \boldsymbol{r}_{n} \ m \ \partial \boldsymbol{r}_{i} \ i \end{bmatrix} \quad \begin{bmatrix} \partial \boldsymbol{r}_{n} \ m \end{bmatrix}$$

$$\frac{\partial \operatorname{dtw}_{\gamma}(\boldsymbol{X},\boldsymbol{Y})}{\partial \boldsymbol{\Delta}} = \left[\frac{\partial \boldsymbol{r_{n,m}}}{\partial \boldsymbol{\Delta_{ij}}} \right]_{ij} = \left[\frac{\partial \boldsymbol{r_{n,m}}}{\partial r_{i,j}} \frac{\partial r_{i,j}}{\partial \boldsymbol{\Delta_{ij}}} \right]_{ij} = \left[\frac{\partial \boldsymbol{r_{n,m}}}{\partial r_{i,j}} \right]_{ij}$$

$$e_{i,j} = \frac{\partial \boldsymbol{r_{n,m}}}{\partial r_{i,j}}$$

Computation Graph: Forward

Bellman's recursion has the following computational graph

Backward Pass

with a few simplifications, the backward pass boils down to the following updates.

Backward Recurrence

$$a = e^{\frac{1}{\gamma}(r_{i+1,j} - r_{i,j} - \Delta_{i+1,j})}$$

$$b = e^{\frac{1}{\gamma}(r_{i,j+1} - r_{i,j} - \Delta_{i,j+1})}$$

$$c = e^{\frac{1}{\gamma}(r_{i+1,j+1} - r_{i,j} - \Delta_{i+1,j+1})}$$

$$e_{i,j} = e_{i+1,j} \cdot a + e_{i,j+1} \cdot b + e_{i+1,j+1} \cdot c$$

$$\nabla_X \operatorname{dtw}_{\gamma}(\boldsymbol{X}, \boldsymbol{Y}) = \left(\frac{\partial \Delta(\boldsymbol{X}, \boldsymbol{Y})}{\partial \boldsymbol{X}}\right)^T E$$

Backward Recurrence

$$a = e^{\frac{1}{\gamma}(r_{i+1,j} - r_{i,j} - \Delta_{i+1,j})}$$

$$b = e^{\frac{1}{\gamma}(r_{i,j+1} - r_{i,j} - \Delta_{i,j+1})}$$

$$c = e^{\frac{1}{\gamma}(r_{i+1,j+1} - r_{i,j} - \Delta_{i+1,j+1})}$$

$$e_{i,j} = e_{i+1,j} \cdot a + e_{i,j+1} \cdot b + e_{i+1,j+1} \cdot c$$

$$\nabla_X \operatorname{dtw}_{\gamma}(\boldsymbol{X}, \boldsymbol{Y}) = \left(\frac{\partial \Delta(\boldsymbol{X}, \boldsymbol{Y})}{\partial \boldsymbol{X}}\right)^T E$$

0. The DTW Geometry

1. Soft-DTW

2. Soft-DTW as a Loss Function

Interpolation Between 2 Time Series

$$\min_{\mathbf{X}} \left[\lambda \operatorname{dtw}_{\gamma}(\mathbf{X}, \mathbf{Y}_{1}) + (1 - \lambda) \operatorname{dtw}_{\gamma}(\mathbf{X}, \mathbf{Y}_{2}) \right]$$

sDTW Barycenter

 $\min_{\mathbf{X}} \sum_{j=1}^{\mathbf{\lambda}_j} \frac{\lambda_j}{m_j} \operatorname{dtw}_{\gamma}(\mathbf{X}, \mathbf{Y}_j)$

[DBA] Petitjean et al., A global averaging method for dynamic time warping, with applications to clustering. *Pattern Recognition*, 44 (3):678–693, 2011.

sDTW Barycenter

$$\min_{\mathbf{X}} \sum_{j=1}^{\mathbf{\lambda}_j} \frac{\lambda_j}{m_j} \operatorname{dtw}_{\gamma}(\mathbf{X}, \mathbf{Y}_j)$$

Table 1. Percentage of the datasets on which the proposed soft-DTW barycenter is achieving lower DTW loss (Equation (4) with $\gamma = 0$) than competing methods.

	Random initialization	Euclidean mean initialization			
Comparison with DBA					
$\gamma = 1$	40.51%	3.80%			
$\gamma = 0.1$	93.67%	46.83%			
$\gamma = 0.01$	100%	79.75%			
$\gamma = 0.001$	97.47%	89.87%			
Comparison with subgradient method					
$\gamma = 1$	96.20%	35.44%			
$\gamma = 0.1$	97.47%	72.15%			
$\gamma = 0.01$	97.47%	92.41%			
$\gamma = 0.001$	97.47%	97.47%			

sDTW Barycenter

use $\gamma > 0$ to compute argmin, (Y_j) use $\gamma > 0$ to compute y with dtw_0 , (Y_j) evaluate energy with dtw_0 .					
	Random initialization	Euclidean mean initialization			
Comparison with DBA					
$egin{aligned} &\gamma &= 1 \ &\gamma &= 0.1 \ &\gamma &= 0.01 \ &\gamma &= 0.001 \end{aligned}$	40.51% 93.67% 100% 97.47%	3.80% 46.83% 79.75% 89.87%			
Comparison with subgradient method					
$egin{array}{l} \gamma = 1 \ \gamma = 0.1 \ \gamma = 0.01 \end{array}$	96.20% 97.47%	35.44% 72.15%			

sDTW Clustering

$$\min_{\boldsymbol{X_1},...,\boldsymbol{X_k}} \sum_{j=1}^{N} \min_{\boldsymbol{i}=1,...,k} \operatorname{dtw}_{\gamma}(\boldsymbol{X_i}, \boldsymbol{Y_j})$$

Table 2. Percentage of the datasets on which the proposed soft-DTW based k-means is achieving lower DTW loss (Equation (5) with $\gamma = 0$) than competing methods.

	Random initialization	Euclidean mean initialization			
Comparison with DBA					
$\gamma = 1$	15.78%	29.31%			
$\gamma=0.1$	24.56%	24.13%			
$\gamma = 0.01$	59.64%	55.17%			
$\gamma=0.001$	77.19%	68.97%			
Comparison with subgradient method					
$\gamma = 1$	42.10%	46.44%			
$\gamma = 0.1$	57.89%	50%			
$\gamma = 0.01$	76.43%	65.52%			
$\gamma=0.001$	96.49%	84.48%			

Nearast Controld

sDTW Prediction Loss

sDTW Prediction Loss

Table 3. Averaged rank obtained by a multi-layer perceptron (MLP) under Euclidean and soft-DTW losses. Euclidean initialization means that we initialize the MLP trained with soft-DTW loss by the solution of the MLP trained with Euclidean loss.

Training loss	Random initialization	Euclidean initialization		
When evaluating with DTW loss				
Euclidean	3.46	4.21		
soft-DTW ($\gamma = 1$)	3.55	3.96		
soft-DTW ($\gamma = 0.1$)	3.33	3.42		
soft-DTW ($\gamma = 0.01$)	2.79	2.12		
soft-DTW ($\gamma = 0.001$)	1.87	1.29		
When evaluating with Euclidean loss				
Euclidean	1.05	1.70		
soft-DTW ($\gamma = 1$)	2.41	2.99		
soft-DTW ($\gamma = 0.1$)	3.42	3.38		
soft-DTW ($\gamma = 0.01$)	4.13	3.64		
soft-DTW ($\gamma = 0.001$)	3.99	3.29		

entroid

Poster: #96

- **Dynamic Time Warping** is a natural and flexible discrepancy to compare time series, yet it is **non-differentiable**.
- **Soft-DTW** is a differentiable approximation, with better convexity properties.
- Using **soft-DTW** quantity results in better minima, even when measured with the original DTW
- Code available on

https://github.com/mblondel/soft-dtw