# Soft-DTW: A Differentiable Loss Function for Time Series

**Marco Cuturi**

université
**PARIS-SACLAY**

ENSAE
ParisTech

École nationale
de la statistique
et de l'administration
économique

université
PARIS-SACLAY

**Mathieu Blondel**

NTT

From: Plume App

Follow pollution levels
in real time in your city

**Ground truth (reality)**

From: Plume App

Follow pollution levels in real time in your city

Carrier · 12:00 PM · SAN FRANCISCO

FRESH air

11°C  UV0  21:02  THU 10

How wrong was this prediction?

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

**Ground truth (reality)**

From: Plume App

- 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,
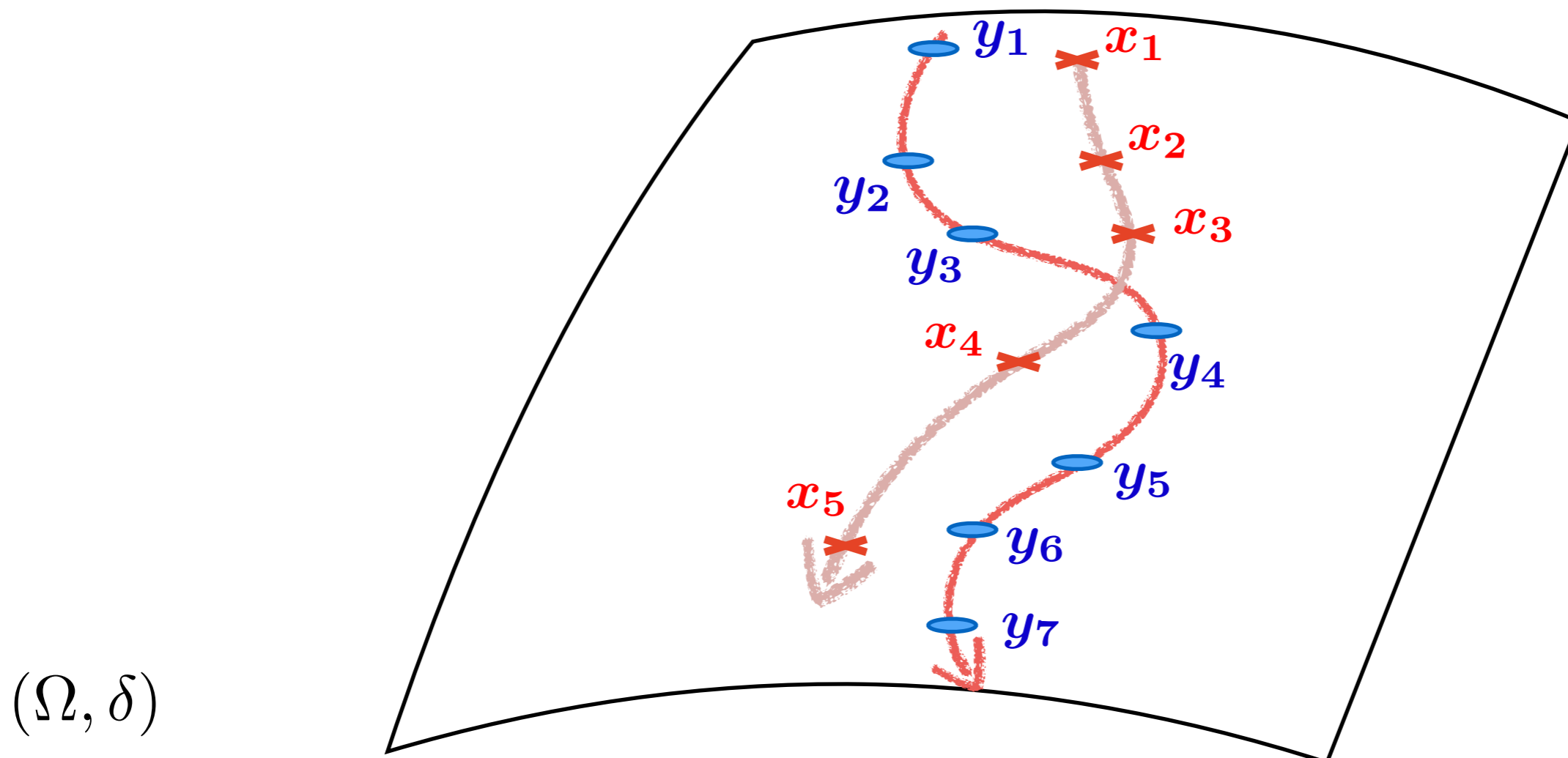  - Learning with structured (time series) output

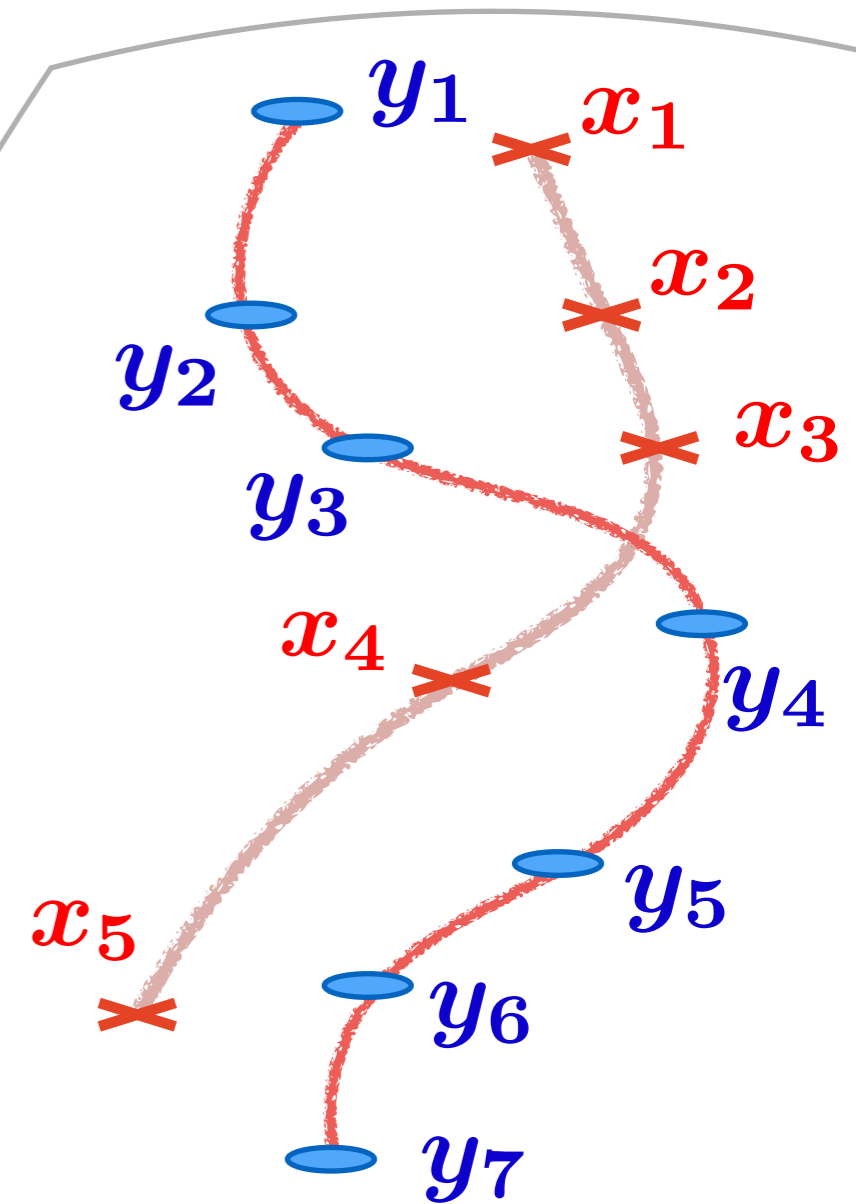# 0. The DTW Geometry

# 1. Soft-DTW

# 2. Soft-DTW as a Loss Function

# Dynamic Time Warping [Sakoe&Chiba'78]

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



$(\Omega, \delta)$

# Pairwise Distance Matrix

|       | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ |
|-------|-------|-------|-------|-------|-------|-------|-------|
| $x_1$ |       |       |       |       |       |       |       |
| $x_2$ |       |       |       |       |       |       |       |
| $x_3$ |       |       |       |       |       |       |       |
| $x_4$ |       |       |       |       |       |       |       |
| $x_5$ |       |       |       |       |       |       |       |

# Pairwise Distance Matrix



$$\Delta_{ij} = \delta(x_i, y_j)$$

# Pairwise Distance Matrix

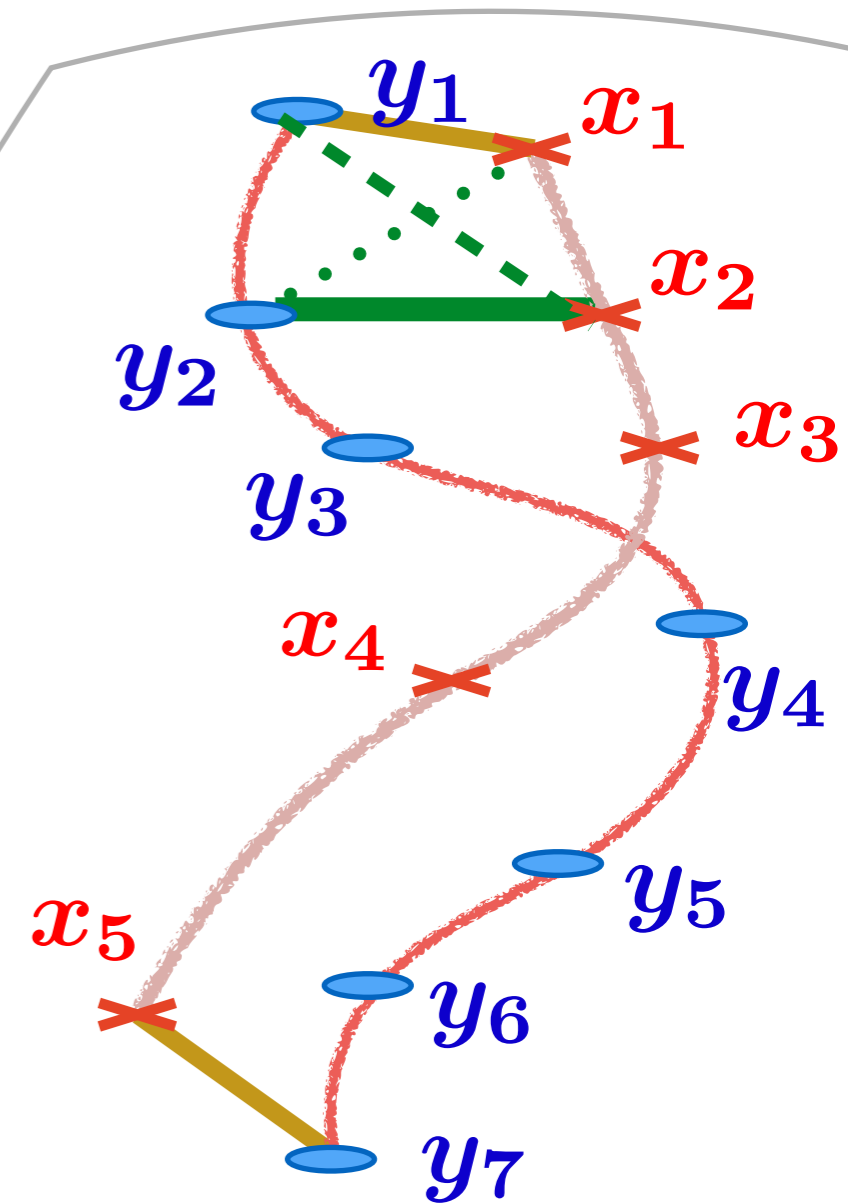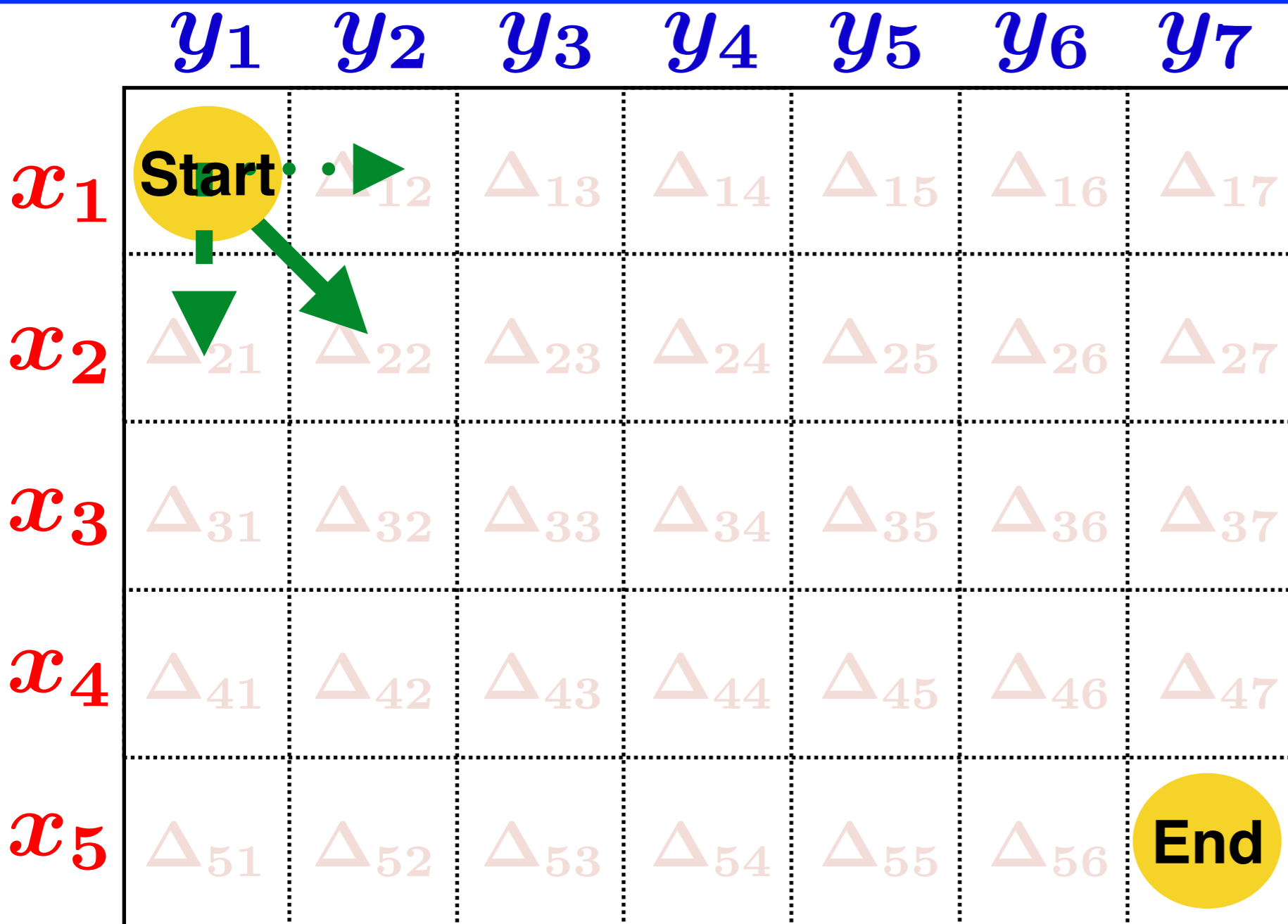|  | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ |
|---|---|---|---|---|---|---|---|
| $x_1$ | $\Delta_{11}$ | $\Delta_{12}$ | $\Delta_{13}$ | $\Delta_{14}$ | $\Delta_{15}$ | $\Delta_{16}$ | $\Delta_{17}$ |
| $x_2$ | $\Delta_{21}$ | $\Delta_{22}$ | $\Delta_{23}$ | $\Delta_{24}$ | $\Delta_{25}$ | $\Delta_{26}$ | $\Delta_{27}$ |
| $x_3$ | $\Delta_{31}$ | $\Delta_{32}$ | $\Delta_{33}$ | $\Delta_{34}$ | $\Delta_{35}$ | $\Delta_{36}$ | $\Delta_{37}$ |
| $x_4$ | $\Delta_{41}$ | $\Delta_{42}$ | $\Delta_{43}$ | $\Delta_{44}$ | $\Delta_{45}$ | $\Delta_{46}$ | $\Delta_{47}$ |
| $x_5$ | $\Delta_{51}$ | $\Delta_{52}$ | $\Delta_{53}$ | $\Delta_{54}$ | $\Delta_{55}$ | $\Delta_{56}$ | $\Delta_{57}$ |

# Alignment Path

# Alignment Path

# Alignment Path

# Alignment Path

# Alignment Path

# Alignment Path

# Alignment Path

# Alignment Path

# Alignment Path

# Alignment Path

# Alignment Path

# Alignment Path

# Path Cost

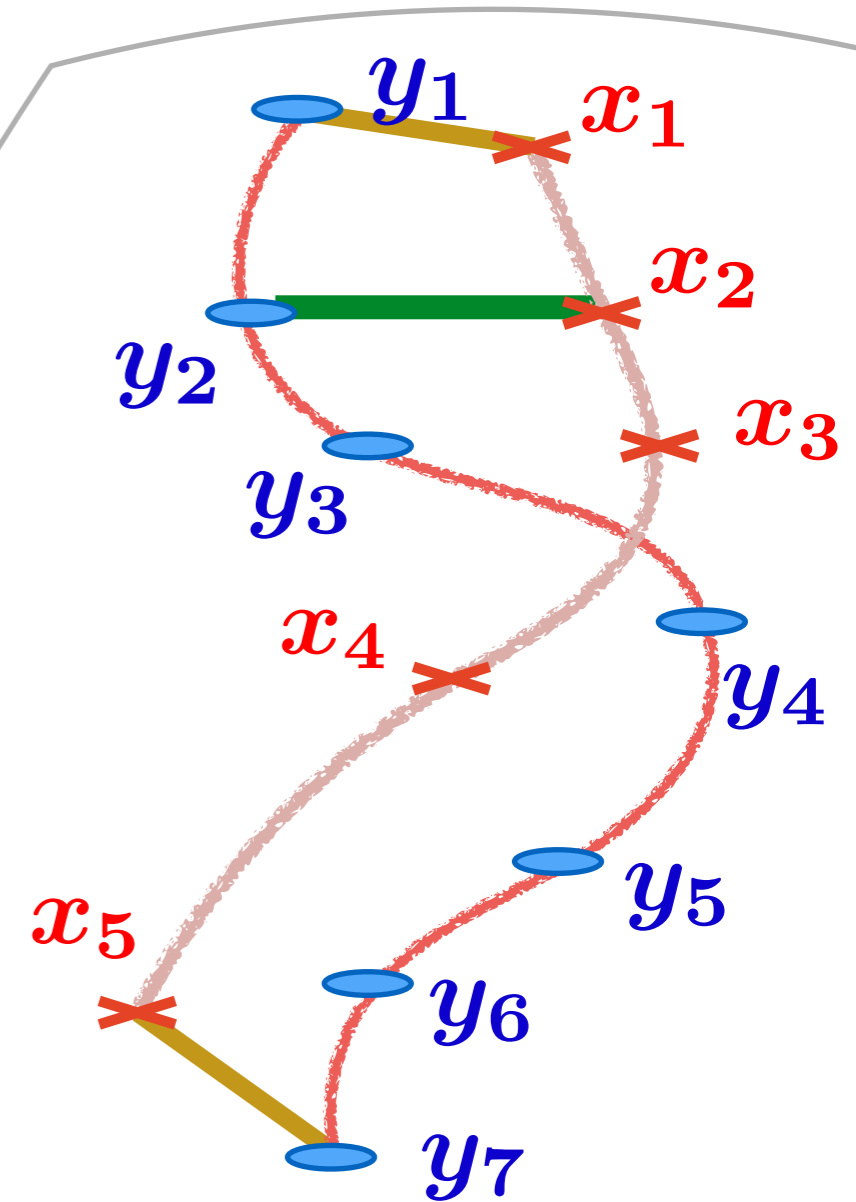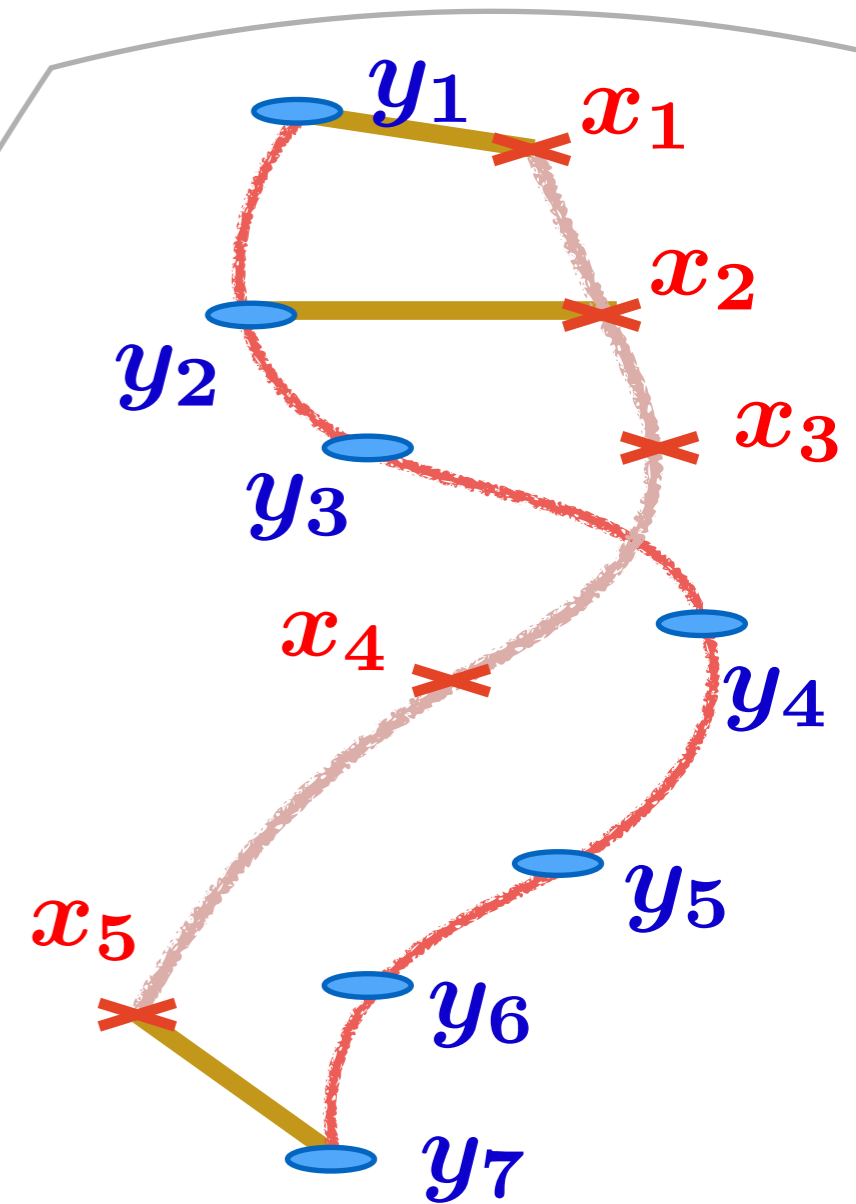|  | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ |
|---|---|---|---|---|---|---|---|
| $x_1$ | $\Delta_{11}$ | $\Delta_{12}$ | $\Delta_{13}$ | $\Delta_{14}$ | $\Delta_{15}$ | $\Delta_{16}$ | $\Delta_{17}$ |
| $x_2$ | $\Delta_{21}$ | $\Delta_{22}$ | $\Delta_{23}$ | $\Delta_{24}$ | $\Delta_{25}$ | $\Delta_{26}$ | $\Delta_{27}$ |
| $x_3$ | $\Delta_{31}$ | $\Delta_{32}$ | $\Delta_{33}$ | $\Delta_{34}$ | $\Delta_{35}$ | $\Delta_{36}$ | $\Delta_{37}$ |
| $x_4$ | $\Delta_{41}$ | $\Delta_{42}$ | $\Delta_{43}$ | $\Delta_{44}$ | $\Delta_{45}$ | $\Delta_{46}$ | $\Delta_{47}$ |
| $x_5$ | $\Delta_{51}$ | $\Delta_{52}$ | $\Delta_{53}$ | $\Delta_{54}$ | $\Delta_{55}$ | $\Delta_{56}$ | $\Delta_{57}$ |

$$\text{Cost} = \Delta_{11} + \Delta_{22} + \Delta_{32} + \Delta_{33} + \Delta_{34} + \Delta_{46} + \Delta_{56} + \Delta_{57}$$

# Path Cost

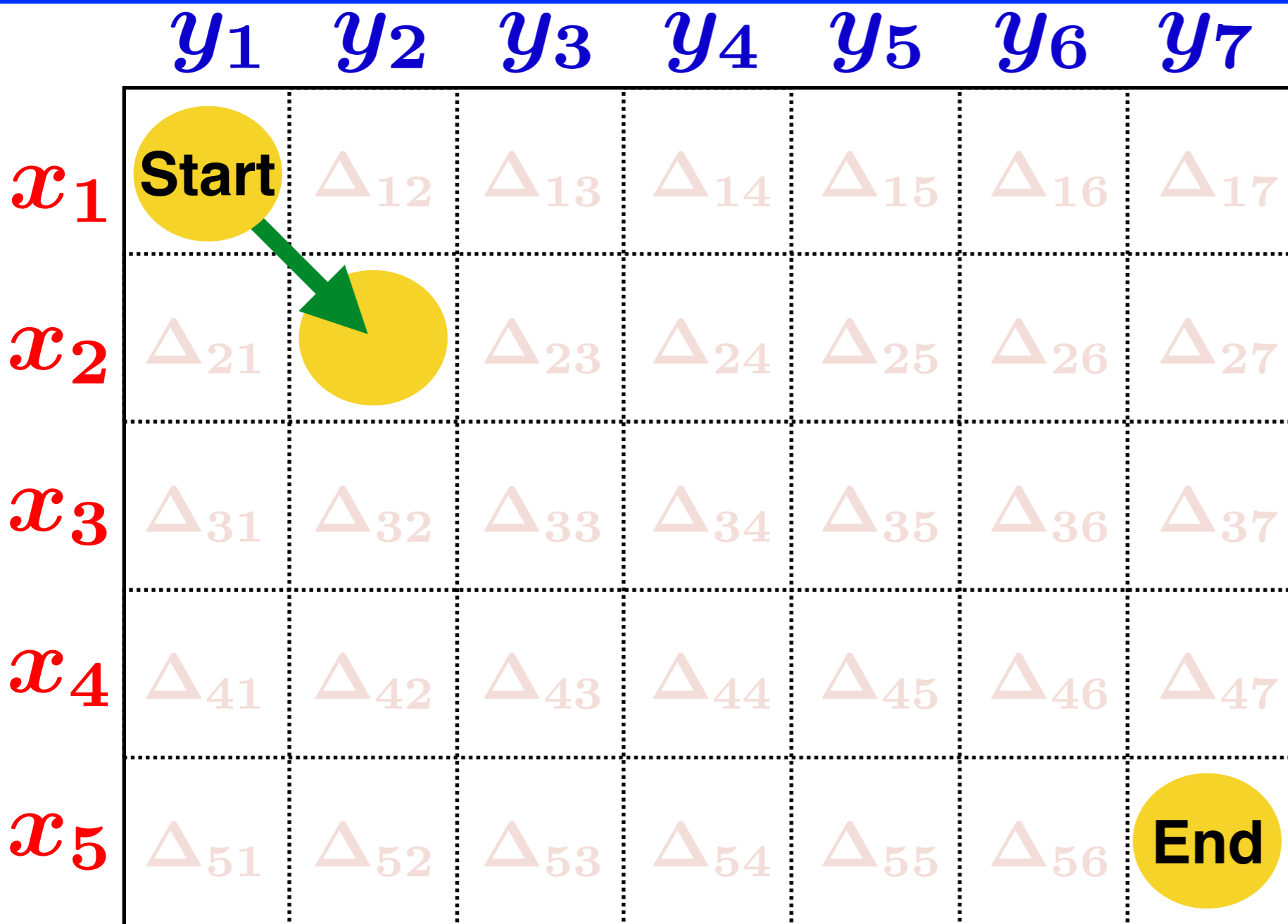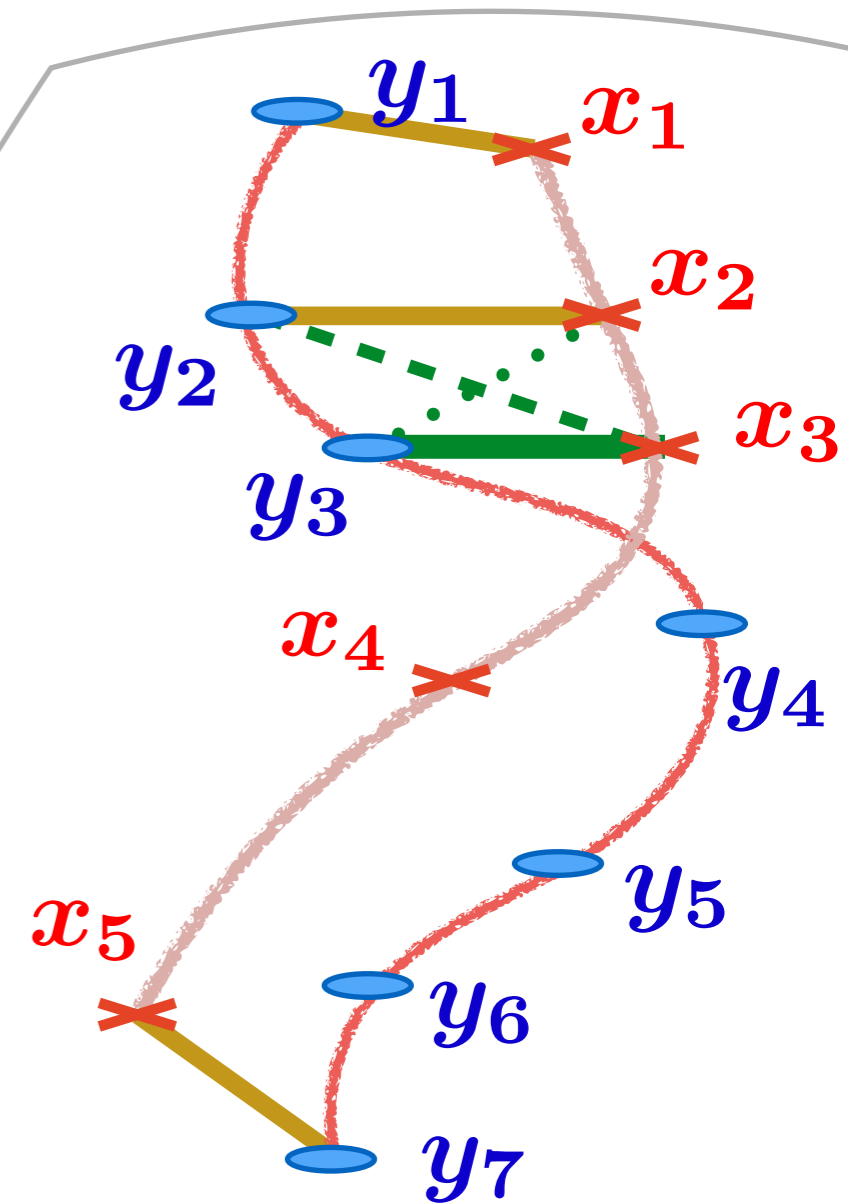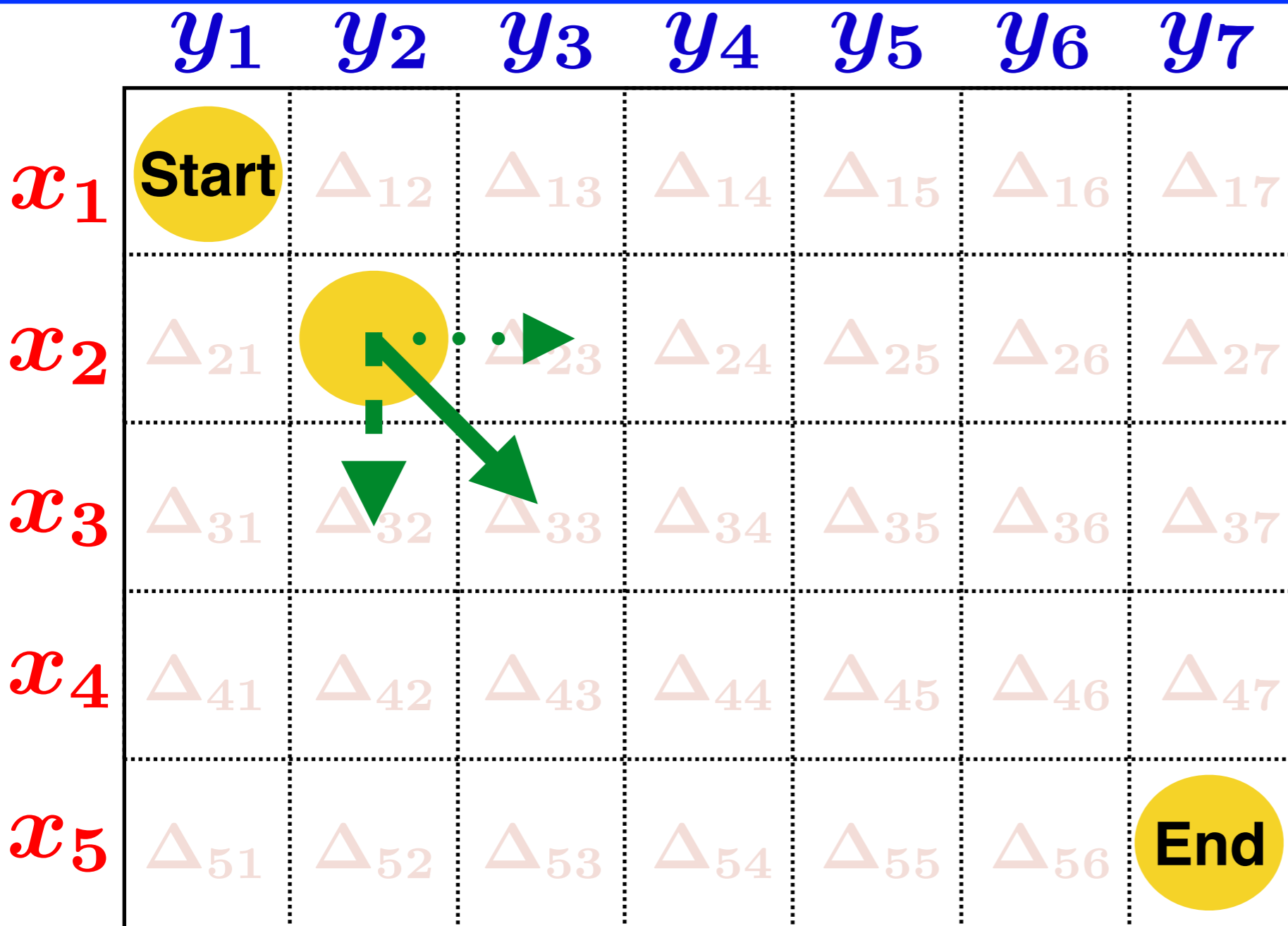|  | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ |
|---|---|---|---|---|---|---|---|
| $x_1$ | $\Delta_{11}$ | $\Delta_{12}$ | $\Delta_{13}$ | $\Delta_{14}$ | $\Delta_{15}$ | $\Delta_{16}$ | $\Delta_{17}$ |
| $x_2$ | $\Delta_{21}$ | $\Delta_{22}$ | $\Delta_{23}$ | $\Delta_{24}$ | $\Delta_{25}$ | $\Delta_{26}$ | $\Delta_{27}$ |
| $x_3$ | $\Delta_{31}$ | $\Delta_{32}$ | $\Delta_{33}$ | $\Delta_{34}$ | $\Delta_{35}$ | $\Delta_{36}$ | $\Delta_{37}$ |
| $x_4$ | $\Delta_{41}$ | $\Delta_{42}$ | $\Delta_{43}$ | $\Delta_{44}$ | $\Delta_{45}$ | $\Delta_{46}$ | $\Delta_{47}$ |
| $x_5$ | $\Delta_{51}$ | $\Delta_{52}$ | $\Delta_{53}$ | $\Delta_{54}$ | $\Delta_{55}$ | $\Delta_{56}$ | $\Delta_{57}$ |

$$\text{Cost} = \Delta_{11} + \Delta_{22} + \Delta_{33} + \Delta_{34} + \Delta_{46} + \Delta_{56} + \Delta_{57}$$

13

# Path Cost

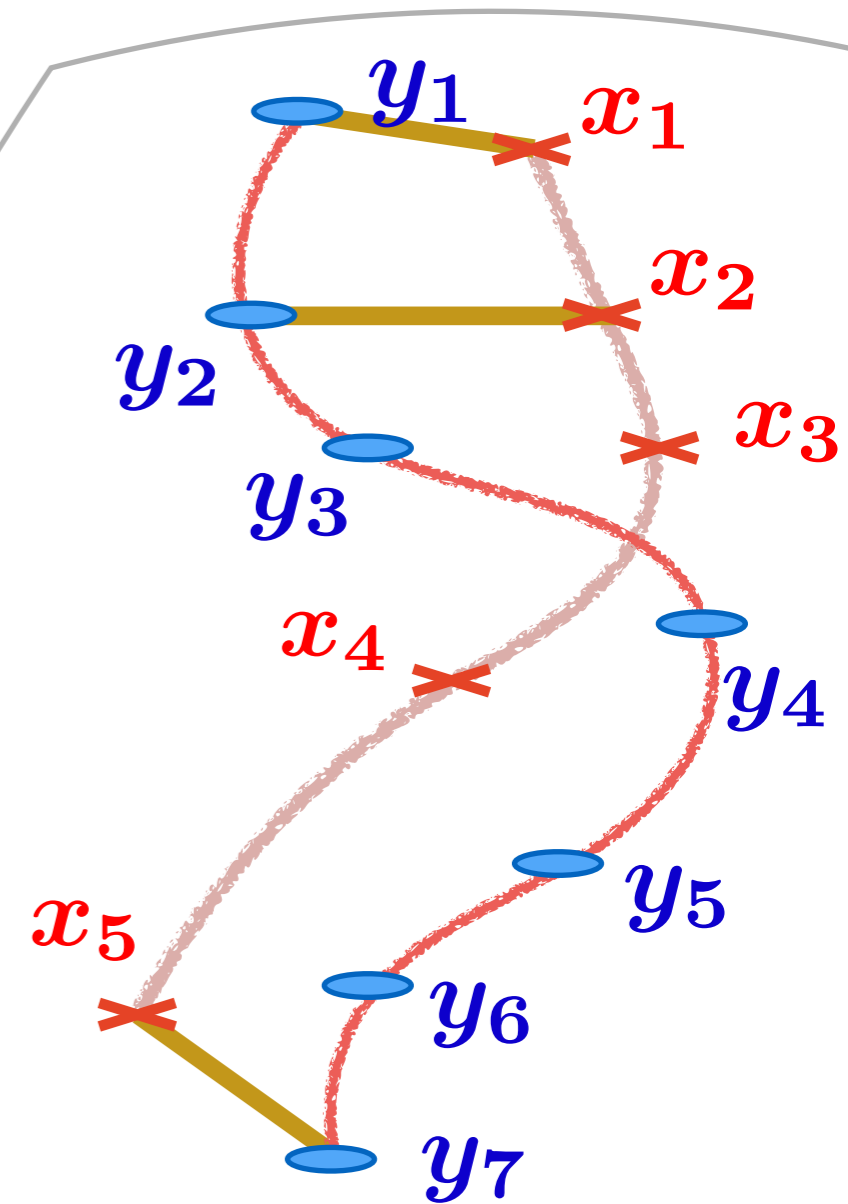| | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ |
|---|---|---|---|---|---|---|---|
| $x_1$ | $\Delta_{11}$ | $\Delta_{12}$ | $\Delta_{13}$ | $\Delta_{14}$ | $\Delta_{15}$ | $\Delta_{16}$ | $\Delta_{17}$ |
| $x_2$ | $\Delta_{21}$ | $\Delta_{22}$ | $\Delta_{23}$ | $\Delta_{24}$ | $\Delta_{25}$ | $\Delta_{26}$ | $\Delta_{27}$ |
| $x_3$ | $\Delta_{31}$ | $\Delta_{32}$ | $\Delta_{33}$ | $\Delta_{34}$ | $\Delta_{35}$ | $\Delta_{36}$ | $\Delta_{37}$ |
| $x_4$ | $\Delta_{41}$ | $\Delta_{42}$ | $\Delta_{43}$ | $\Delta_{44}$ | $\Delta_{45}$ | $\Delta_{46}$ | $\Delta_{47}$ |
| $x_5$ | $\Delta_{51}$ | $\Delta_{52}$ | $\Delta_{53}$ | $\Delta_{54}$ | $\Delta_{55}$ | $\Delta_{56}$ | $\Delta_{57}$ |

$$\text{Cost} = \Delta_{11} + \Delta_{22} + \Delta_{33} + \Delta_{34} + \Delta_{46} + \Delta_{57}$$

13

# Path Cost

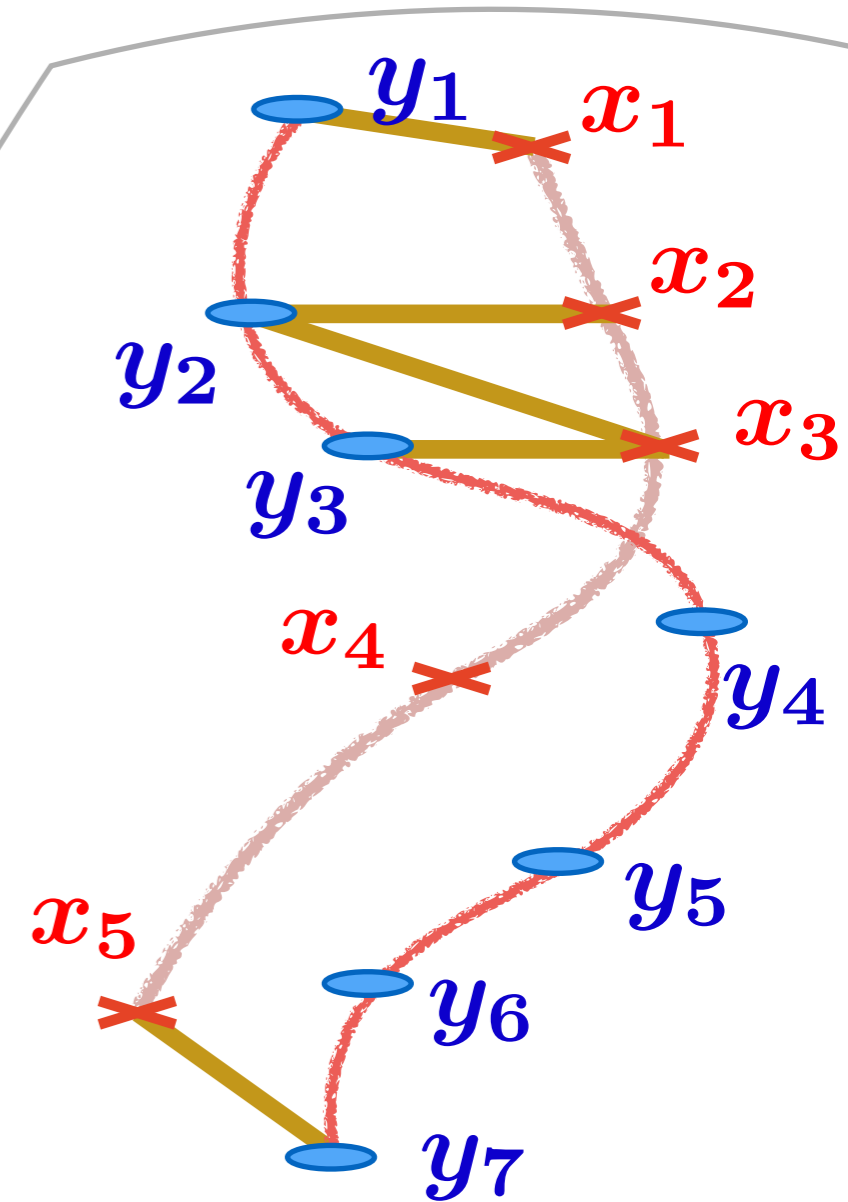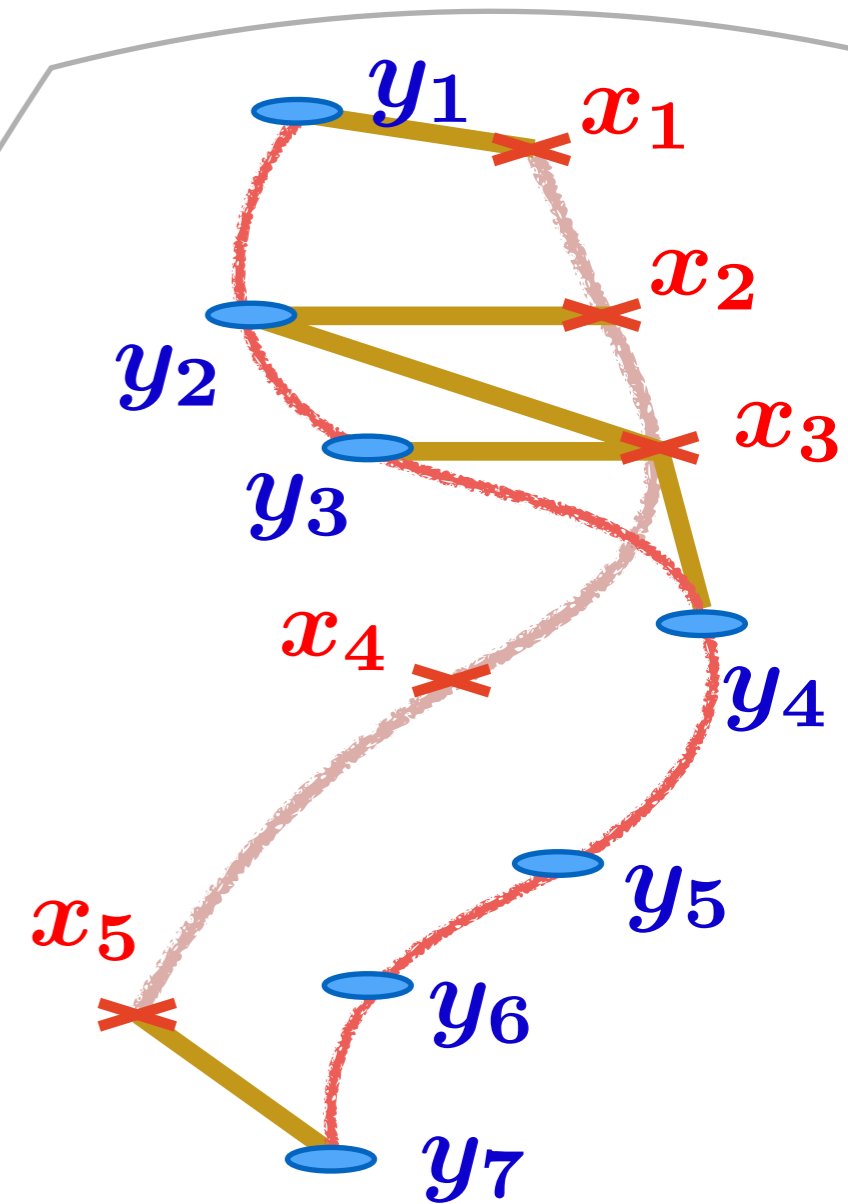|       | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ |
|-------|-------|-------|-------|-------|-------|-------|-------|
| $x_1$ | $\Delta_{11}$ | $\Delta_{12}$ | $\Delta_{13}$ | $\Delta_{14}$ | $\Delta_{15}$ | $\Delta_{16}$ | $\Delta_{17}$ |
| $x_2$ | $\Delta_{21}$ | $\Delta_{22}$ | $\Delta_{23}$ | $\Delta_{24}$ | $\Delta_{25}$ | $\Delta_{26}$ | $\Delta_{27}$ |
| $x_3$ | $\Delta_{31}$ | $\Delta_{32}$ | $\Delta_{33}$ | $\Delta_{34}$ | $\Delta_{35}$ | $\Delta_{36}$ | $\Delta_{37}$ |
| $x_4$ | $\Delta_{41}$ | $\Delta_{42}$ | $\Delta_{43}$ | $\Delta_{44}$ | $\Delta_{45}$ | $\Delta_{46}$ | $\Delta_{47}$ |
| $x_5$ | $\Delta_{51}$ | $\Delta_{52}$ | $\Delta_{53}$ | $\Delta_{54}$ | $\Delta_{55}$ | $\Delta_{56}$ | $\Delta_{57}$ |

# Path Cost

|  | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ |
|---|---|---|---|---|---|---|---|
| $x_1$ | **1** | **0** | **0** | **0** | **0** | **0** | **0** |
| $x_2$ | **0** | **1** | **0** | **0** | **0** | **0** | **0** |
| $x_3$ | **0** | **0** | **1** | **1** | **0** | **0** | **0** |
| $x_4$ | **0** | **0** | **0** | **0** | **1** | **1** | **0** |
| $x_5$ | **0** | **0** | **0** | **0** | **0** | **0** | **1** |

$$= A$$

# Path Cost

|  | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ |
|---|---|---|---|---|---|---|---|
| $x_1$ | **1** | **0** | **0** | **0** | **0** | **0** | **0** |
| $x_2$ | **0** | **1** | **0** | **0** | **0** | **0** | **0** |
| $x_3$ | **0** | **0** | **1** | **1** | **0** | **0** | **0** |
| $x_4$ | **0** | **0** | **0** | **0** | **1** | **1** | **0** |
| $x_5$ | **0** | **0** | **0** | **0** | **0** | **0** | **1** |

$= A$

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

# Minimum Cost Alignment Matrix?

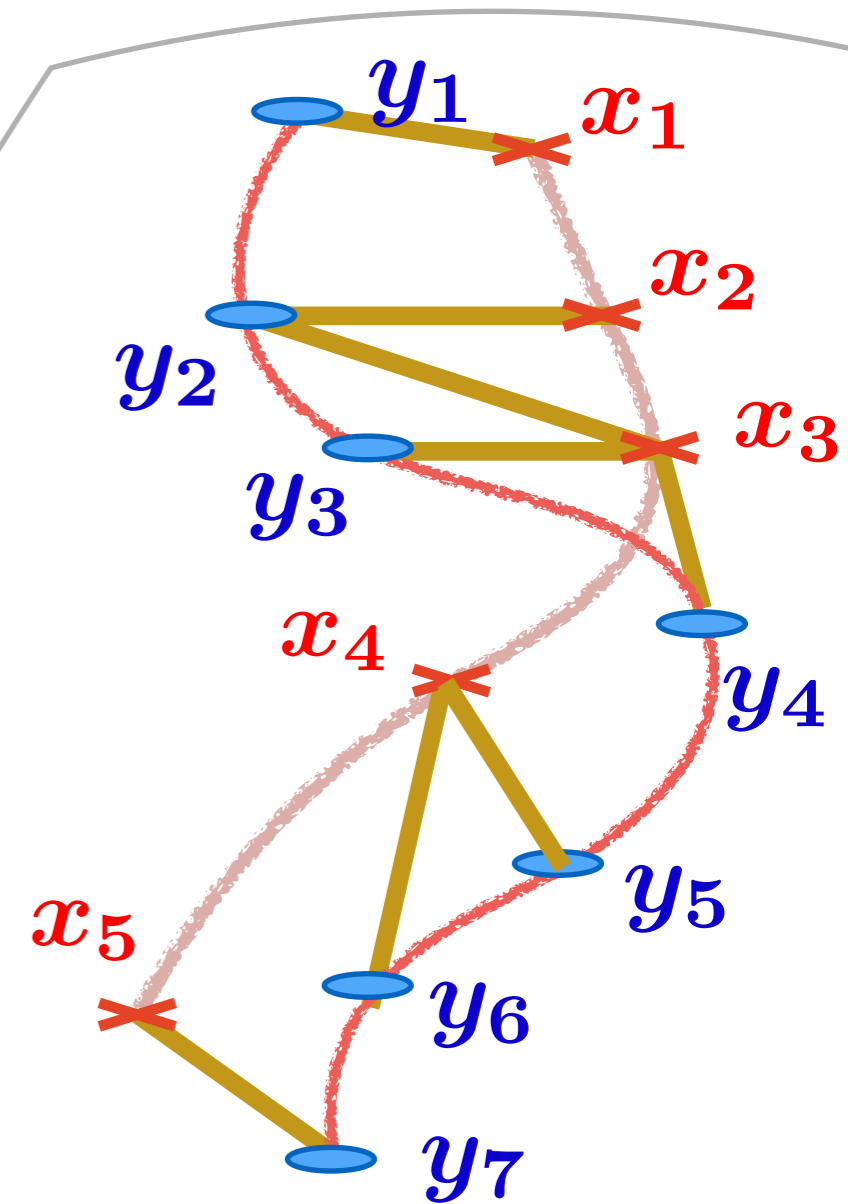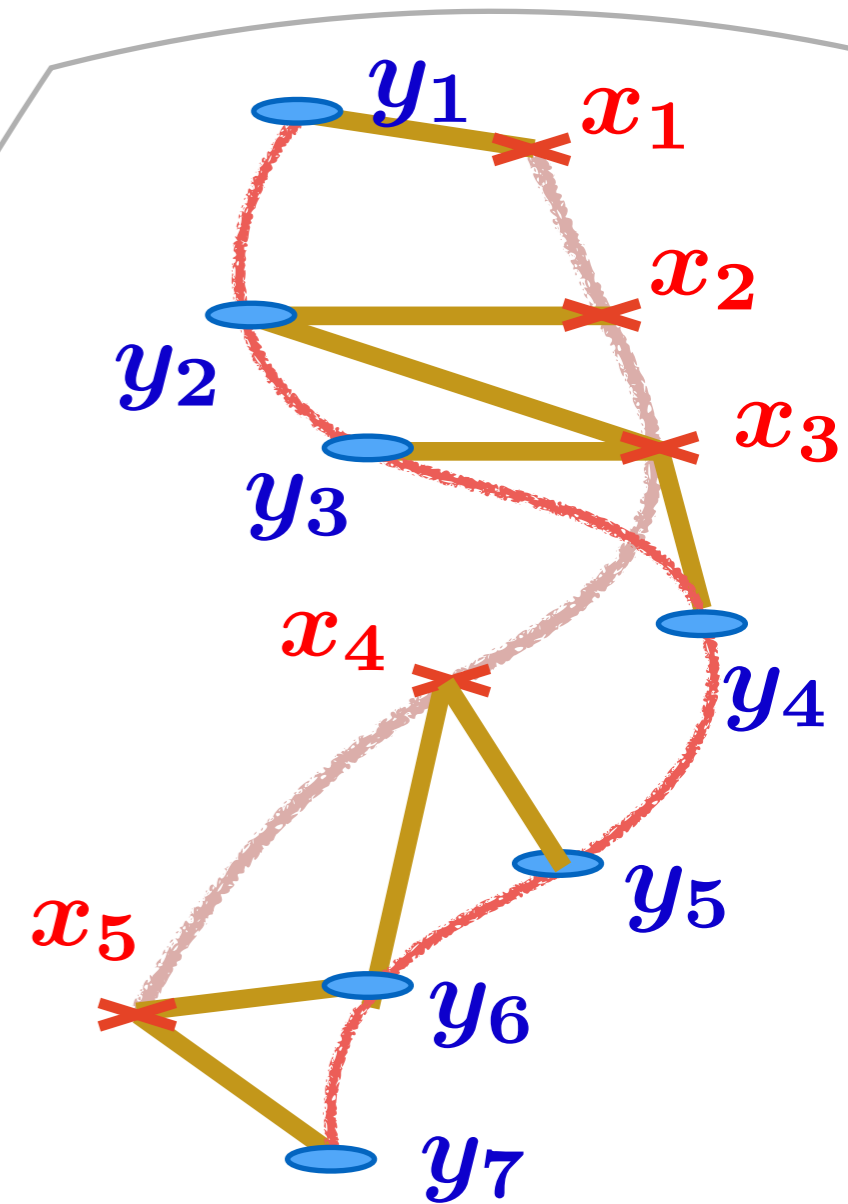|  | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ |
|---|---|---|---|---|---|---|---|
| $x_1$ | **Start** $\Delta_{11}$ | $\Delta_{12}$ | $\Delta_{13}$ | $\Delta_{14}$ | $\Delta_{15}$ | $\Delta_{16}$ | $\Delta_{17}$ |
| $x_2$ | $\Delta_{21}$ | $\Delta_{22}$ | $\Delta_{23}$ | $\Delta_{24}$ | $\Delta_{25}$ | $\Delta_{26}$ | $\Delta_{27}$ |
| $x_3$ | $\Delta_{31}$ | $\Delta_{32}$ | $\Delta_{33}$ | $\Delta_{34}$ | $\Delta_{35}$ | $\Delta_{36}$ | $\Delta_{37}$ |
| $x_4$ | $\Delta_{41}$ | $\Delta_{42}$ | $\Delta_{43}$ | $\Delta_{44}$ | $\Delta_{45}$ | $\Delta_{46}$ | $\Delta_{47}$ |
| $x_5$ | $\Delta_{51}$ | $\Delta_{52}$ | $\Delta_{53}$ | $\Delta_{54}$ | $\Delta_{55}$ | $\Delta_{56}$ | **End** $\Delta_{57}$ |

# Minimum Cost Alignment Matrix?



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

15

# Dynamic Time Warping [Sakoe&Chiba'78]



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

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

16

# Number of valid paths

Size of $\mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})$ is exponential in $\textcolor{red}{n}, \textcolor{blue}{m}$.

$\#\mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m}) = \text{Delannoy}(\textcolor{red}{n} - 1, \textcolor{blue}{m} - 1)$

| | |
|---|---|
| n=m=3 | 13 |
| n=m=5 | 321 |
| n=m=10 | 1462563 |

$\dots$

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

# Best Path: Bellman Recursion

|     | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ |
|-----|-------|-------|-------|-------|-------|-------|-------|
| $x_1$ | $\Delta_{11}$ | $\Delta_{12}$ | $\Delta_{13}$ | $\Delta_{14}$ | $\Delta_{15}$ | $\Delta_{16}$ | $\Delta_{17}$ |
| $x_2$ | $\Delta_{21}$ | $\Delta_{22}$ | $\Delta_{23}$ | $\Delta_{24}$ | $\Delta_{25}$ | $\Delta_{26}$ | $\Delta_{27}$ |
| $x_3$ | $\Delta_{31}$ | $\Delta_{32}$ | $\Delta_{33}$ | $\Delta_{34}$ | $\Delta_{35}$ | $\Delta_{36}$ | $\Delta_{37}$ |
| $x_4$ | $\Delta_{41}$ | $\Delta_{42}$ | $\Delta_{43}$ | $\Delta_{44}$ | $\Delta_{45}$ | $\Delta_{46}$ | $\Delta_{47}$ |
| $x_5$ | $\Delta_{51}$ | $\Delta_{52}$ | $\Delta_{53}$ | $\Delta_{54}$ | $\Delta_{55}$ | $\Delta_{56}$ | $\Delta_{57}$ |

# Best Path: Bellman Recursion



$$r^\star_{3,5} = \min_{A \in \mathcal{A}(3,5)} \langle A, [\boldsymbol{\Delta}_{ij}]_{i \leq 3, j \leq 5} \rangle$$

# Best Path: Bellman Recursion



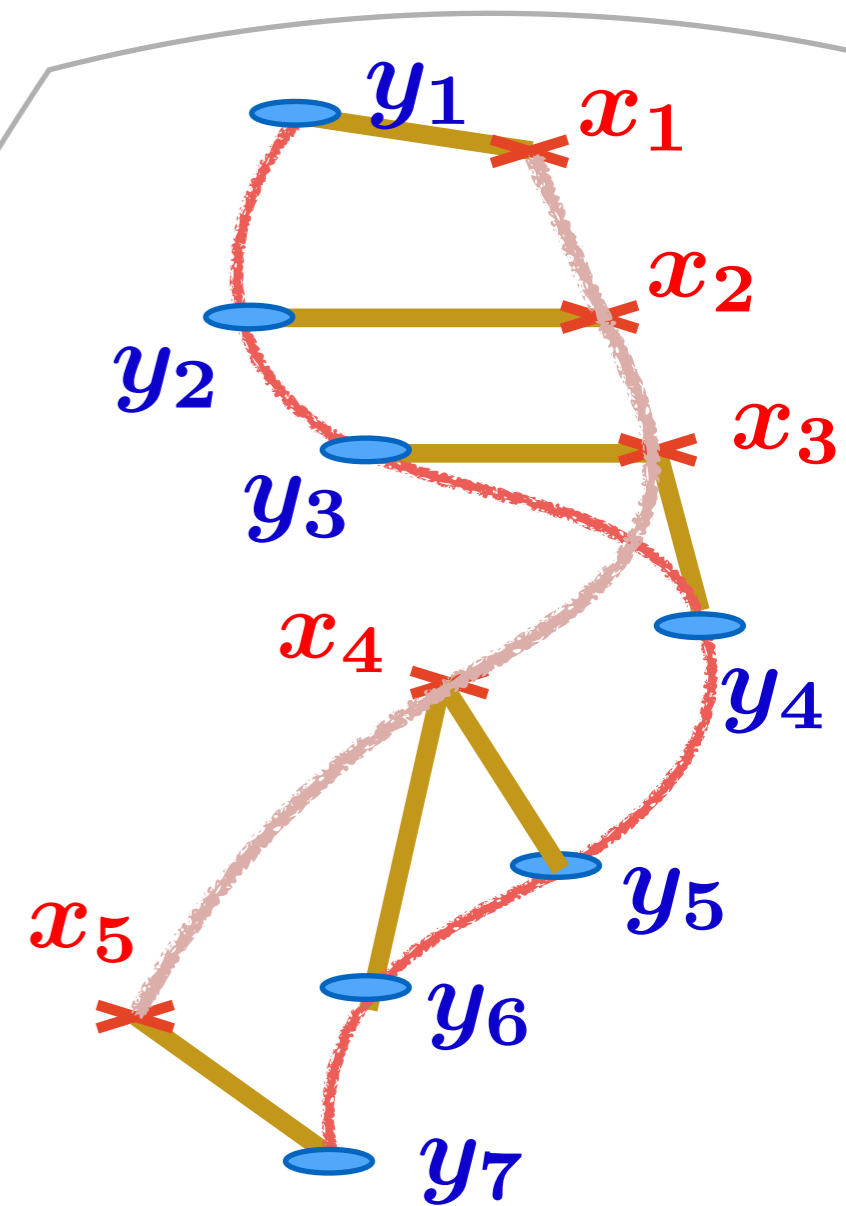|  | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ |
|---|---|---|---|---|---|---|---|
| $x_1$ | $\Delta_{11}$ | $\Delta_{12}$ | $\Delta_{13}$ | $\Delta_{14}$ | $\Delta_{15}$ | $\Delta_{16}$ | $\Delta_{17}$ |
| $x_2$ | $\Delta_{21}$ | $\Delta_{22}$ | $\Delta_{23}$ | $\Delta_{24}$ | $\Delta_{25}$ | $\Delta_{26}$ | $\Delta_{27}$ |
| $x_3$ | $\Delta_{31}$ | $\Delta_{32}$ | $\Delta_{33}$ | $\Delta_{34}$ | $r^{\star}_{3,5}$ | $\Delta_{36}$ | $\Delta_{37}$ |
| $x_4$ | $\Delta_{41}$ | $\Delta_{42}$ | $\Delta_{43}$ | $r^{\star}_{4,4}$ | $\Delta_{45}$ | $\Delta_{46}$ | $\Delta_{47}$ |
| $x_5$ | $\Delta_{51}$ | $\Delta_{52}$ | $\Delta_{53}$ | $\Delta_{54}$ | $\Delta_{55}$ | $\Delta_{56}$ | $\Delta_{57}$ |

# Best Path: Bellman Recursion

# Best Path: Bellman Recursion

|  | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ |
|---|---|---|---|---|---|---|---|
| $x_1$ | $\Delta_{11}$ | $\Delta_{12}$ | $\Delta_{13}$ | $\Delta_{14}$ | $\Delta_{15}$ | $\Delta_{16}$ | $\Delta_{17}$ |
| $x_2$ | $\Delta_{21}$ | $\Delta_{22}$ | $\Delta_{23}$ | $\Delta_{24}$ | $\Delta_{25}$ | $\Delta_{26}$ | $\Delta_{27}$ |
| $x_3$ | $\Delta_{31}$ | $\Delta_{32}$ | $\Delta_{33}$ | $r^{\star}_{3,4}$ $\Delta_{34}$ | $r^{\star}_{3,5}$ $\Delta_{35}$ | $\Delta_{36}$ | $\Delta_{37}$ |
| $x_4$ | $\Delta_{41}$ | $\Delta_{42}$ | $\Delta_{43}$ | $r^{\star}_{4,4}$ $\Delta_{44}$ | $r^{\star}_{4,5}$ $\Delta_{45}$ | $\Delta_{46}$ | $\Delta_{47}$ |
| $x_5$ | $\Delta_{51}$ | $\Delta_{52}$ | $\Delta_{53}$ | $\Delta_{54}$ | $\Delta_{55}$ | $\Delta_{56}$ | $\Delta_{57}$ |

# Best Path: Bellman Recursion

|  | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ |
|---|---|---|---|---|---|---|---|
| $x_1$ | $\Delta_{11}$ | $\Delta_{12}$ | $\Delta_{13}$ | $\Delta_{14}$ | $\Delta_{15}$ | $\Delta_{16}$ | $\Delta_{17}$ |
| $x_2$ | $\Delta_{21}$ | $\Delta_{22}$ | $\Delta_{23}$ | $\Delta_{24}$ | $\Delta_{25}$ | $\Delta_{26}$ | $\Delta_{27}$ |
| $x_3$ | $\Delta_{31}$ | $\Delta_{32}$ | $\Delta_{33}$ | $r^{\star}_{3,4}$ | $r^{\star}_{3,5}$ | $\Delta_{36}$ | $\Delta_{37}$ |
| $x_4$ | $\Delta_{41}$ | $\Delta_{42}$ | $\Delta_{43}$ | $r^{\star}_{4,4}$ | $\Delta_{45}^{r^{\star}_{4,5}}$ | $\Delta_{46}$ | $\Delta_{47}$ |
| $x_5$ | $\Delta_{51}$ | $\Delta_{52}$ | $\Delta_{53}$ | $\Delta_{54}$ | $\Delta_{55}$ | $\Delta_{56}$ | $\Delta_{57}$ |

$$r^{\star}_{4,5} = \min(r^{\star}_{3,5}, r^{\star}_{4,4}, r^{\star}_{3,4}) + \Delta_{4,5}$$

# Best Path: Bellman Recursion

| 0 | $\infty$ | $\infty$ | $\infty$ | $\infty$ | $\infty$ | $\infty$ | $\infty$ |
|---|---|---|---|---|---|---|---|
| $\infty$ | $r_{1,1}$ $\Delta_{11}$ | $\Delta_{12}$ | $\Delta_{13}$ | $\Delta_{14}$ | $\Delta_{15}$ | $\Delta_{16}$ | $\Delta_{17}$ |
| $\infty$ | $\Delta_{21}$ | $\Delta_{22}$ | $\Delta_{23}$ | $\Delta_{24}$ | $\Delta_{25}$ | $\Delta_{26}$ | $\Delta_{27}$ |
| $\infty$ | $\Delta_{31}$ | $\Delta_{32}$ | $\Delta_{33}$ | $\Delta_{34}$ | $\Delta_{35}$ | $\Delta_{36}$ | $\Delta_{37}$ |
| $\infty$ | $\Delta_{41}$ | $\Delta_{42}$ | $\Delta_{43}$ | $\Delta_{44}$ | $\Delta_{45}$ | $\Delta_{46}$ | $\Delta_{47}$ |
| $\infty$ | $\Delta_{51}$ | $\Delta_{52}$ | $\Delta_{53}$ | $\Delta_{54}$ | $\Delta_{55}$ | $\Delta_{56}$ | $\Delta_{57}$ |

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

$$r_{0,0} = 0$$

# Best Path: Bellman Recursion



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

# Best Path: Bellman Recursion



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

# Best Path: Bellman Recursion



$$\mathbf{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \textcolor{red}{r_{n,m}}$$

# Optimal Path

| 0 | ∞ | ∞ | ∞ | ∞ | ∞ | ∞ | ∞ |
|---|---|---|---|---|---|---|---|
| ∞ | **1** | **0** | **0** | **0** | **0** | **0** | **0** |
| ∞ | **0** | **1** | **0** | **0** | **0** | **0** | **0** |
| $A^\star$ ∞ | **0** | **0** | **1** | **1** | **0** | **0** | **0** |
| ∞ | **0** | **0** | **0** | **0** | **1** | **0** | **0** |
| ∞ | **0** | **0** | **0** | **0** | **0** | **1** | **1** |

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

Computational cost: O(nm)

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

# DTW as a Loss: Differentiability?



$$\mathbf{dtw}_0(\boldsymbol{X} + \boldsymbol{dX}, \boldsymbol{Y}) = \boldsymbol{r_{n,m}} + ?$$

# DTW as a Loss: Differentiability?



$$\mathbf{dtw}_0(\boldsymbol{X + dX}, \boldsymbol{Y}) = \langle A^\star, \boldsymbol{\Delta_{X+dXY}} \rangle$$

# DTW as a Loss: Differentiability?



$$\mathbf{dtw}_0(\textcolor{red}{X + dX}, \textcolor{blue}{Y}) \neq \langle A^\star, \Delta_{\textcolor{red}{X+dX}\textcolor{blue}{Y}} \rangle$$

# DTW as a Loss: Differentiability?



$$\mathbf{dtw}_0(\boldsymbol{X} + \boldsymbol{dX}, \boldsymbol{Y}) \neq \langle A^\star, \boldsymbol{\Delta}_{\boldsymbol{X} + \boldsymbol{dX}\boldsymbol{Y}} \rangle$$

# DTW as a Loss: Differentiability?

$$\mathbf{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \textcolor{brown}{\Delta_{\textcolor{red}{X}\textcolor{blue}{Y}}} \rangle$$

- $\mathbf{dtw}_0$ is piecewise **linear** w.r.t $\textcolor{brown}{\Delta}$

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

26

# DTW as a Loss: Differentiability?

$$\mathbf{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \textcolor{brown}{\Delta_{\textcolor{red}{X}\textcolor{blue}{Y}}} \rangle$$

$$\nabla_{\textcolor{red}{X}} \mathbf{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \left( \frac{\partial \textcolor{brown}{\Delta_{\textcolor{red}{X}\textcolor{blue}{Y}}}}{\partial \textcolor{red}{X}} \right)^T \nabla_{\textcolor{brown}{\Delta}} \min_{\mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle \cdot, \textcolor{brown}{\Delta_{\textcolor{red}{X}\textcolor{blue}{Y}}} \rangle$$

27

# DTW as a Loss: Differentiability?

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

$$\nabla_X \mathbf{dtw}_0(X, Y) = \left( \frac{\partial \Delta_{XY}}{\partial X} \right)^T \nabla_\Delta \min_{\mathcal{A}(n,m)} \langle \cdot, \Delta_{XY} \rangle$$

Jacobian matrix of $\Delta$ **w.r.t.** $X$

# DTW as a Loss: Differentiability?

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

$$\nabla_X \mathbf{dtw}_0(X, Y) = \left( \frac{\partial \Delta_{XY}}{\partial X} \right)^T \nabla_\Delta \min_{\mathcal{A}(n,m)} \langle \cdot, \Delta_{XY} \rangle$$

$$= A^\star$$

Jacobian matrix of $\Delta$ **w.r.t.** $X$

iff optimal solution
is unique

# DTW as a Loss: Differentiability?

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

$$\nabla_X \mathbf{dtw}_0(X, Y) = \left( \frac{\partial \Delta_{XY}}{\partial X} \right)^T \quad \nabla_\Delta \min_{\mathcal{A}(n, m)} \langle \cdot, \Delta_{XY} \rangle$$

$$= A^\star$$

Jacobian matrix of $\Delta$ **w.r.t.** $X$

iff optimal solution
is unique

When $A^\star$ is not unique, $\mathbf{dtw}_0$ has a discontinuous gradient!

# Our proposal: smoothing the min

$$\mathbf{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \textcolor{brown}{\Delta_{\textcolor{red}{X}\textcolor{blue}{Y}}} \rangle$$

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

# Our proposal: smoothing the min

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

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

Fix: smoothed min operator

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

# *Example* softmin of quadratic functions

$$f(\boldsymbol{x}) = \min_{i=1,\dots,s}^{\gamma} a_i \boldsymbol{x^2} + b_i \boldsymbol{x} + c_i$$

# *Example* softmin of quadratic functions

$$f(\boldsymbol{x}) = \min_{i=1,\ldots,s}^{\gamma} a_i \boldsymbol{x^2} + b_i \boldsymbol{x} + c_i$$

# Soft-DTW

$$\mathbf{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \textcolor{brown}{\Delta_{\textcolor{red}{X}\textcolor{blue}{Y}}} \rangle$$

# Soft-DTW

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

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

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

30

# Soft-DTW

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

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

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

$$\mathbf{dtw}_{\textcolor{green}{\gamma}}(\textcolor{red}{X}, \textcolor{blue}{Y}) = -\textcolor{green}{\gamma} \log \sum_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} e^{-\frac{\langle A, \textcolor{brown}{\Delta}_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle}{\textcolor{green}{\gamma}}}$$

# Relation to Global Alignment kernels

$$k_{\mathrm{GA}} := \sum_{A \in \mathcal{A}(\textcolor{red}{\boldsymbol{n}}, \textcolor{blue}{\boldsymbol{m}})} e^{-\frac{\langle A, \textcolor{red}{\boldsymbol{\Delta}}_{\textcolor{red}{\boldsymbol{X}}\textcolor{blue}{\boldsymbol{Y}}} \rangle}{\textcolor{green}{\boldsymbol{\gamma}}}}$$

**CVBM'07**

A positive semi-definite **kernel** between time series

# Relation to Global Alignment kernels

$$k_{\mathrm{GA}} := \sum_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} e^{-\frac{\langle A, \textcolor{red}{\Delta}_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle}{\textcolor{green}{\gamma}}}$$

**CVBM'07**

A positive semi-definite **kernel** between time series

$$\mathbf{dtw}_{\textcolor{green}{\gamma}}(\textcolor{red}{X}, \textcolor{blue}{Y}) = -\textcolor{green}{\gamma} \log k_{\mathrm{GA}}$$

Computing soft-DTW is equivalent to computing k$_{\mathrm{GA}}$ in **log domain**

# Recursive Computation

$$\mathbf{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{\boldsymbol{n}}, \textcolor{blue}{\boldsymbol{m}})} \langle A, \textcolor{brown}{\boldsymbol{\Delta_{XY}}} \rangle$$

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

# Recursive Computation

$$\mathbf{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \textcolor{brown}{\Delta_{\textcolor{red}{X}\textcolor{blue}{Y}}} \rangle$$

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

$$\mathbf{dtw}_{\textcolor{green}{\gamma}}(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})}^{\textcolor{green}{\gamma}} \langle A, \textcolor{brown}{\Delta_{\textcolor{red}{X}\textcolor{blue}{Y}}} \rangle$$

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

32

# Recursive Computation

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

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

$$\mathbf{dtw}_\gamma(X, Y) = \min^\gamma_{A \in \mathcal{A}(n,m)} \langle A, \Delta_{XY} \rangle$$

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

Simply replace min operator!

# Recursive Computation

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

$$r_{i,j} = \min(r_{i-1,j-1}, \, r_{i-1,j}, \, r_{i,j-1}) + \boldsymbol{\Delta}_{i,j}$$

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

$$r_{i,j} = \min^{\textcolor{green}{\gamma}}(r_{i-1,j-1}, \, r_{i-1,j}, \, r_{i,j-1}) + \boldsymbol{\Delta}_{i,j}$$

**Simply replace min operator!**

**Stable: recursion in log domain!**

# Differentiation

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

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

# Differentiation

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

$$\nabla_X \mathbf{dtw}_{\gamma}(\textcolor{red}{\boldsymbol{X}}, \textcolor{blue}{\boldsymbol{Y}}) = \left( \frac{\partial \textcolor{brown}{\boldsymbol{\Delta}}(\textcolor{red}{\boldsymbol{X}}, \textcolor{blue}{\boldsymbol{Y}})}{\partial \textcolor{red}{\boldsymbol{X}}} \right)^T \mathbb{E}_{\textcolor{green}{\boldsymbol{\gamma}}}[A]$$

33

# Differentiation

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

$$\nabla_X \mathbf{dtw}_{\gamma}(\boldsymbol{X}, \boldsymbol{Y}) = \left( \frac{\partial \boldsymbol{\Delta}(\boldsymbol{X}, \boldsymbol{Y})}{\partial \boldsymbol{X}} \right)^T \mathbb{E}_{\boldsymbol{\gamma}}[A]$$

$$\mathbb{E}_{\boldsymbol{\gamma}}[A] := \frac{\sum\limits_{A \in \mathcal{A}(\boldsymbol{n}, \boldsymbol{m})} A e^{-\frac{\langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle}{\boldsymbol{\gamma}}}}{\sum\limits_{A \in \mathcal{A}(\boldsymbol{n}, \boldsymbol{m})} e^{-\frac{\langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle}{\boldsymbol{\gamma}}}}$$

Expectation of Path under Gibbs distribution

33

# Differentiation

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

$$\nabla_X \mathbf{dtw}_{\gamma}(X, Y) = \left( \frac{\partial \Delta(X, Y)}{\partial X} \right)^T \underbrace{\mathbb{E}_{\gamma}[A]}_{\nabla_{\Delta} \mathbf{dtw}_{\gamma}(X, Y)}$$

$$\mathbb{E}_{\gamma}[A] := \frac{\sum\limits_{A \in \mathcal{A}(n,m)} A e^{-\frac{\langle A, \Delta_{XY} \rangle}{\gamma}}}{\sum\limits_{A \in \mathcal{A}(n,m)} e^{-\frac{\langle A, \Delta_{XY} \rangle}{\gamma}}}$$

Expectation of Path under Gibbs distribution

33

# Computing the expectation $\mathrm{E}_{\gamma}[A]$

$$\mathbb{E}_{\boldsymbol{\gamma}}[A] := \frac{\displaystyle\sum_{A \in \mathcal{A}(\boldsymbol{n},\boldsymbol{m})} A e^{-\frac{\langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle}{\boldsymbol{\gamma}}}}{\displaystyle\sum_{A \in \mathcal{A}(\boldsymbol{n},\boldsymbol{m})} e^{-\frac{\langle A, \boldsymbol{\Delta}_{\boldsymbol{X}\boldsymbol{Y}} \rangle}{\boldsymbol{\gamma}}}}$$

Naive computation
is intractable

# Computing the expectation $\text{E}_\gamma[A]$

$$\mathbb{E}_{\boldsymbol{\gamma}}[A] := \frac{\displaystyle\sum_{A \in \mathcal{A}(\boldsymbol{n},\boldsymbol{m})} A e^{-\frac{\langle A, \Delta_{\boldsymbol{X}\boldsymbol{Y}} \rangle}{\boldsymbol{\gamma}}}}{\displaystyle\sum_{A \in \mathcal{A}(\boldsymbol{n},\boldsymbol{m})} e^{-\frac{\langle A, \Delta_{\boldsymbol{X}\boldsymbol{Y}} \rangle}{\boldsymbol{\gamma}}}}$$

Naive computation is intractable

$$= \frac{\displaystyle\sum_{A \in \mathcal{A}(\boldsymbol{n},\boldsymbol{m})} A e^{-\frac{\langle A, \Delta_{\boldsymbol{X}\boldsymbol{Y}} \rangle}{\boldsymbol{\gamma}}}}{k_{\text{GA}}}$$

$k_{\text{GA}}$ is the normalization constant (a.k.a. partition function)!

# Computing the expectation $E_\gamma[A]$

$$\mathbb{E}_{\boldsymbol{\gamma}}[A] := \frac{\displaystyle\sum_{A \in \mathcal{A}(\boldsymbol{n},\boldsymbol{m})} A e^{-\frac{\langle A, \Delta_{\boldsymbol{XY}} \rangle}{\boldsymbol{\gamma}}}}{\displaystyle\sum_{A \in \mathcal{A}(\boldsymbol{n},\boldsymbol{m})} e^{-\frac{\langle A, \Delta_{\boldsymbol{XY}} \rangle}{\boldsymbol{\gamma}}}}$$

Naive computation is intractable

$$= \frac{\displaystyle\sum_{A \in \mathcal{A}(\boldsymbol{n},\boldsymbol{m})} A e^{-\frac{\langle A, \Delta_{\boldsymbol{XY}} \rangle}{\boldsymbol{\gamma}}}}{k_{\mathrm{GA}}}$$

$k_{\mathrm{GA}}$ is the normalization constant (a.k.a. partition function)!

$$= \nabla_{\Delta} -\gamma \log k_{\mathrm{GA}}$$

$E_\gamma[A]$ is the gradient of the log partition

# Computing the expectation $E_\gamma[A]$

$$\mathbb{E}_{\boldsymbol{\gamma}}[A] := \frac{\displaystyle\sum_{A\in\mathcal{A}(\boldsymbol{n},\boldsymbol{m})} A e^{-\frac{\langle A, \Delta_{\boldsymbol{X}\boldsymbol{Y}}\rangle}{\boldsymbol{\gamma}}}}{\displaystyle\sum_{A\in\mathcal{A}(\boldsymbol{n},\boldsymbol{m})} e^{-\frac{\langle A, \Delta_{\boldsymbol{X}\boldsymbol{Y}}\rangle}{\boldsymbol{\gamma}}}}$$

Naive computation is intractable

$$= \frac{\displaystyle\sum_{A\in\mathcal{A}(\boldsymbol{n},\boldsymbol{m})} A e^{-\frac{\langle A, \Delta_{\boldsymbol{X}\boldsymbol{Y}}\rangle}{\boldsymbol{\gamma}}}}{k_{\mathrm{GA}}}$$

$k_{\mathrm{GA}}$ is the normalization constant (a.k.a. partition function)!

$$= \nabla_\Delta\text{-}\gamma \log k_{\mathrm{GA}}$$

Classical result of exponential families

~~partition~~

34

# Computing the expectation $E_\gamma[A]$

To summarize, we want to compute:

$$E_\gamma[A] = \nabla_\Delta -\gamma \log k_{GA}$$

# Computing the expectation $E_\gamma[A]$

To summarize, we want to compute:

$$E_\gamma[A] = \nabla_\Delta -\gamma \log k_{GA} = \nabla_\Delta dtw_\gamma$$

# Computing the expectation $E_\gamma[A]$

To summarize, we want to compute:

$$E_\gamma[A] = \nabla_\Delta \; {-\gamma} \log k_{GA} = \nabla_\Delta \; dtw_\gamma$$

$E_\gamma[A]$ can be computed by **backpropagation** in the same $O(nm)$ cost as $dtw_\gamma$

To summarize, we want to compute:

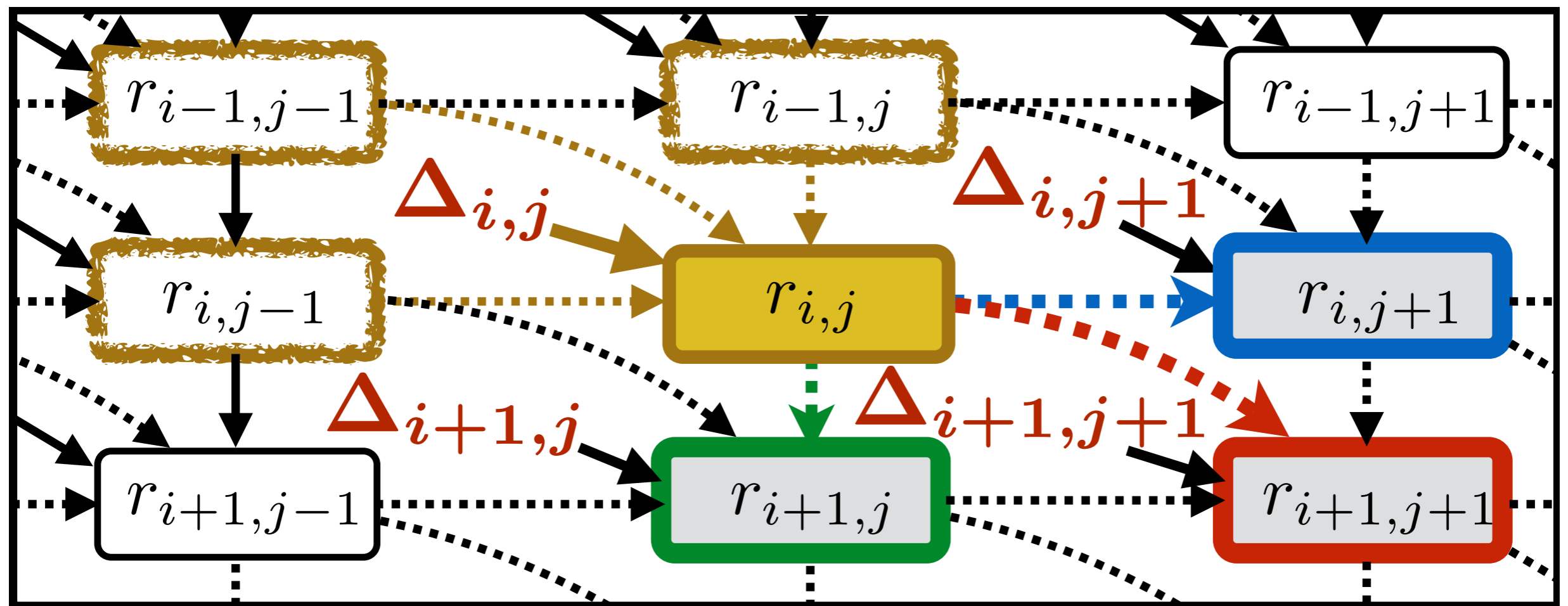$$E_\gamma[A] = \nabla_{\Delta} -\gamma \log k_{GA} = \nabla_{\Delta} dtw_\gamma$$

$E_\gamma[A]$ can be computed by **backpropagation** in the same $O(nm)$ cost as $dtw_\gamma$

We derive a backward recursion **without resorting to autodiff**

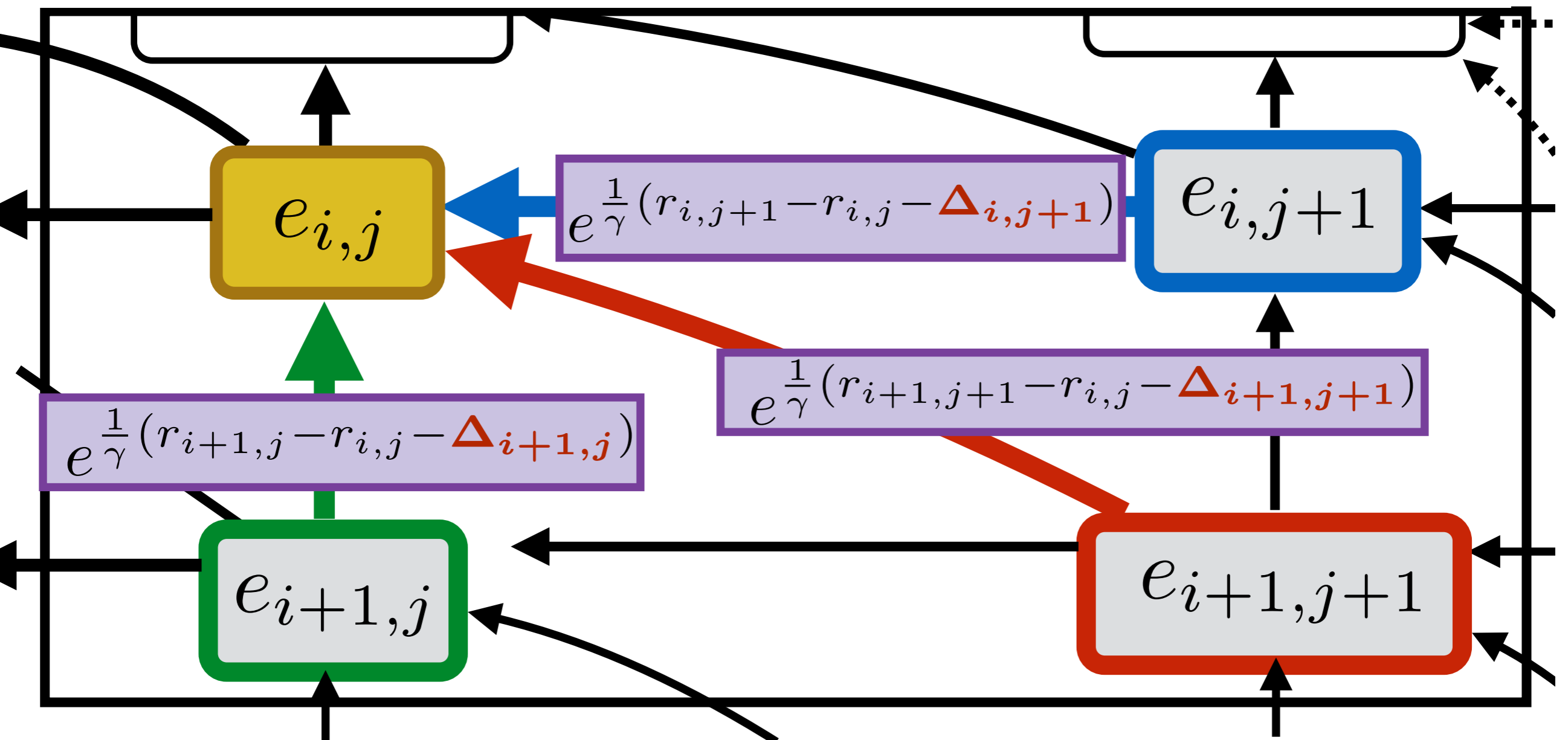Faster and more numerically stable

# Forward Pass

Bellman's recursion has the following computational graph

# Backward Pass

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



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

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

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

$e_{i,j}$

$e_{i,j+1}$

$e_{i+1,j}$

$e_{i+1,j+1}$

# Backward Recursion



$$\begin{array}{ccccccc}
\triangle_{11} & \triangle_{12} & \triangle_{13} & \triangle_{14} & \triangle_{15} & \triangle_{16} & \triangle_{17} \\
\triangle_{21} & \triangle_{22} & \triangle_{23} & \triangle_{24} & \triangle_{25} & \triangle_{26} & \triangle_{27} \\
\triangle_{31} & \triangle_{32} & \triangle_{33} & \triangle_{34} & \triangle_{35} & \triangle_{36} & \triangle_{37} \\
\triangle_{41} & \triangle_{42} & \triangle_{43} & \triangle_{44} & \triangle_{45} & \triangle_{46} & \triangle_{47} \\
\triangle_{51} & \triangle_{52} & \triangle_{53} & \triangle_{54} & \triangle_{55} & \triangle_{56} & e_{5,7}
\end{array}$$

Right column: 0, 0, 0, 0, 0

Bottom row: 0  0  0  0  0  0  0  1

# Backward Recursion

# Backward Recursion



$$E_\gamma[A] = [e]_{ij}$$

# Backward Recursion

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

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

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

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

$$\nabla_X \mathbf{dtw}_\gamma(\mathbf{X}, \mathbf{Y}) = \left( \frac{\partial \mathbf{\Delta}(\mathbf{X}, \mathbf{Y})}{\partial \mathbf{X}} \right)^T E$$
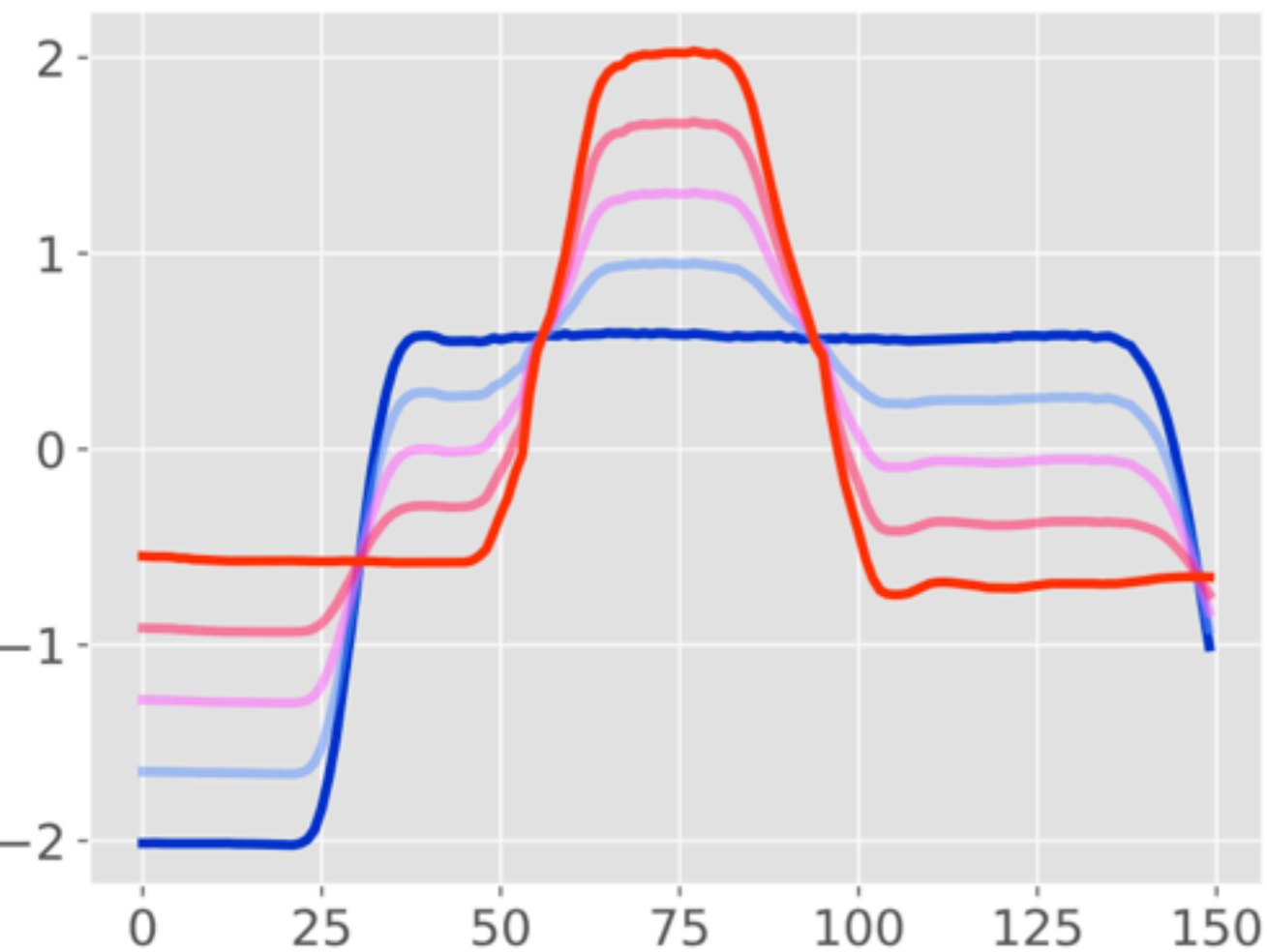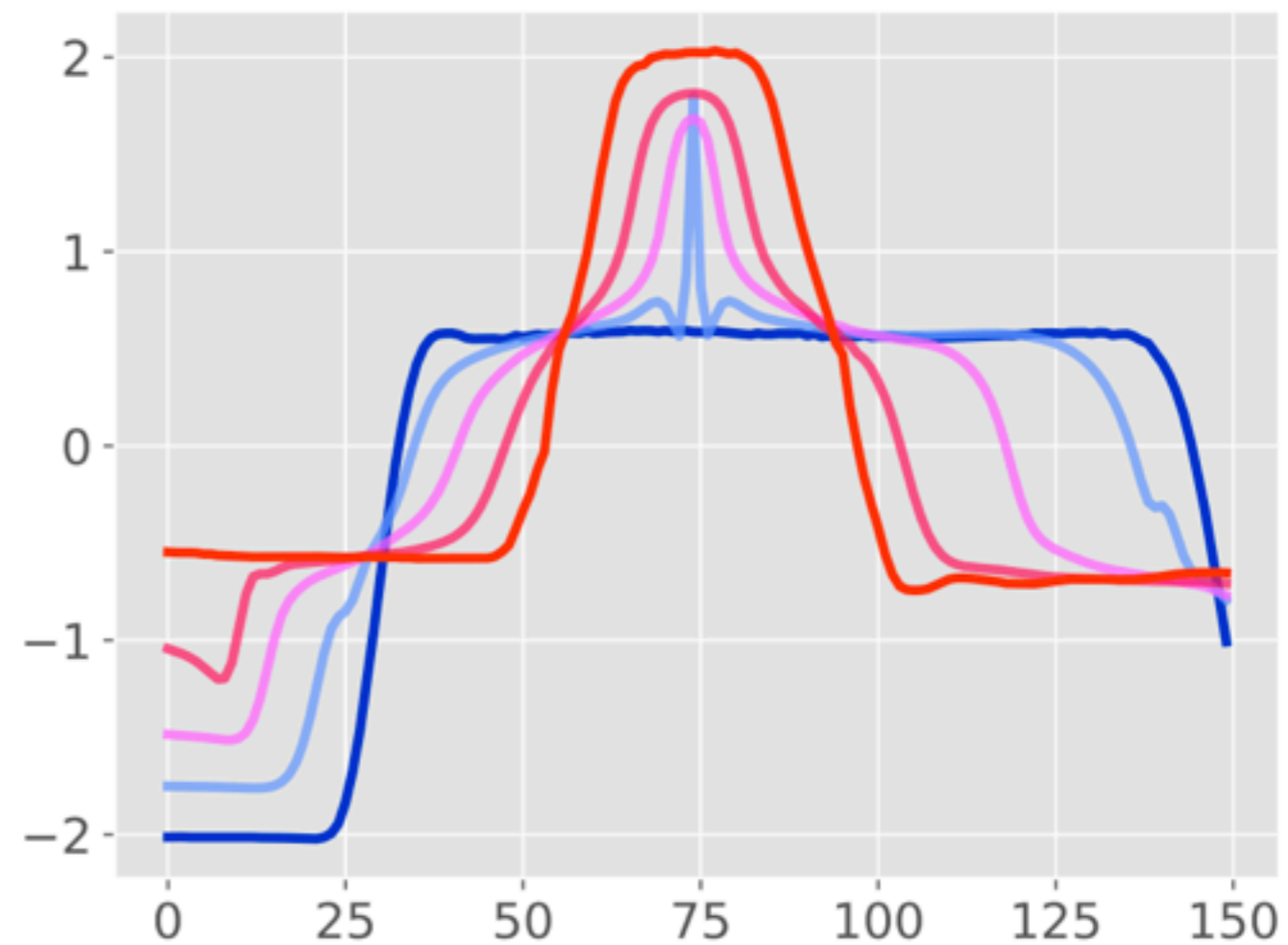
# 0. The DTW Geometry

# 1. Soft-DTW

# 2. Soft-DTW as a Loss Function

# Interpolation Between 2 Time Series

$$\min_{\mathbf{X}} \left[ \lambda \, \mathbf{dtw}_\gamma(\mathbf{X}, {\color{red}Y_1}) + (1 - \lambda) \, \mathbf{dtw}_\gamma(\mathbf{X}, {\color{blue}Y_2}) \right]$$



Euclidean loss          Soft-DTW loss $(\gamma = 1)$

# sDTW Barycenter

$$\min_{\textcolor{red}{X}} \sum_{j=1} \frac{\lambda_j}{m_j} \mathbf{dtw}_\gamma(\textcolor{red}{X}, \textcolor{blue}{Y_j})$$



[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

$$\min_{\textcolor{red}{X}} \sum_{j=1} \frac{\lambda_j}{m_j} \, \mathbf{dtw}_\gamma(\textcolor{red}{X}, \textcolor{blue}{Y_j})$$

*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

$$\min_{\textcolor{red}{X}} \sum_{j=1} \frac{\lambda_j}{m_j} \mathbf{dtw}_\gamma(\textcolor{red}{X}, \textcolor{blue}{Y_j})$$

**Evaluation performed using dtw₀**

*Table 1.* Percentage of the datasets on which the prop... DTW barycenter is achieving lower DTW loss (Equation... $\gamma = 0$) 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% |

% of datasets where soft-dtw is winning

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

Soft-DTW (γ=0.01)

DBA

DBA

Subgradient method

luster 2 (9 points)

Cluster 3 (13 points)

# sDTW Clustering

$$\min_{X_1,\ldots,X_k} \sum_{j=1}^{N} \min_{i=1,\ldots,k} \mathbf{dtw}_\gamma(\textcolor{red}{X_i}, \textcolor{blue}{Y_j})$$

*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% |

47

Nearest Centroid

# sDTW Prediction Loss

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^{N} \frac{1}{m_i} \mathbf{dtw}_{\gamma}(f_{\boldsymbol{\theta}}(x_i), \boldsymbol{Y_i})$$

# sDTW Prediction Loss

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^{N} \frac{1}{m_i} \mathbf{dtw}_{\gamma}(f_{\boldsymbol{\theta}}(x_i), \mathbf{Y_i})$$

*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** ($\mathrm{dtw}_0$) | | |
| 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 |

averaged rank

entroid

# Summary

- **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** typically results in better minima, even when measured with the original DTW

- Python code available on

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