

INTERNATIONAL AUDIO LABORATORIES ERLANGEN
A joint institution of Fraunhofer IIS and Universität Erlangen-Nürnberg



ISMIR Tutorial

Daejeon, Korea, September 21, 2025




Differentiable Alignment Techniques for Music Processing: Techniques and Applications

Part 2: Theoretical Foundations

Meinard Müller, Johannes Zeitler

International Audio Laboratories Erlangen

{meinard.mueller, johannes.zeitler}@audiolabs-erlangen.de



Overview

Part 0: Overview


Part 1: Introduction to Alignment Techniques

Coffee Break

Part 2: Theoretical Foundations & Implementation

© AudiLabs, 2025
Meinard Müller, Johannes Zeitler

ISMIR Tutorial: Differentiable Alignment Techniques for Music Processing
Part 2: Theoretical Foundations. Slide 2



Introduction: Training with Strongly Aligned Targets

- Train DNN-based feature extractor from audio
- Frame-wise annotations (strong targets) are very costly

Input audio

Neural Network


Predicted features, e.g., pitch class

Loss

Frame-wise targets

© AudiLabs, 2025
Meinard Müller, Johannes Zeitler

ISMIR Tutorial: Differentiable Alignment Techniques for Music Processing
Part 2: Theoretical Foundations. Slide 3



Introduction: Training with Weakly Aligned Targets

- Train DNN-based feature extractor from audio
- Frame-wise annotations (strong targets) are very costly
- Only annotate start & end of audio segments
- Retrieve note events from musical score
- Weak targets γ provide information about note event order, but not duration
- Use alignment techniques to train DNN on weakly aligned data

Input audio

Neural Network


Predicted features, e.g., pitch class

? Alignment ?

Weak targets

© AudiLabs, 2025
Meinard Müller, Johannes Zeitler

ISMIR Tutorial: Differentiable Alignment Techniques for Music Processing
Part 2: Theoretical Foundations. Slide 4



Overview

Classical DTW

Differentiability

Convex Regularization

Global formulation

Soft DTW

Recursive formulation

Implementation

Extensions

Relation to CTC

Practical Considerations

© AudiLabs, 2025
Meinard Müller, Johannes Zeitler

ISMIR Tutorial: Differentiable Alignment Techniques for Music Processing
Part 2: Theoretical Foundations. Slide 5



Overview

Classical DTW

Differentiability

Convex Regularization

Global formulation

Soft DTW

Recursive formulation

Implementation

Extensions

Relation to CTC

Practical Considerations

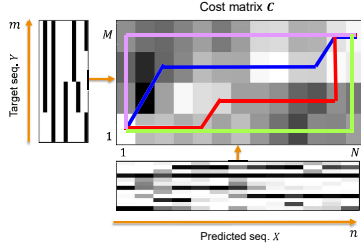
© AudiLabs, 2025
Meinard Müller, Johannes Zeitler

ISMIR Tutorial: Differentiable Alignment Techniques for Music Processing
Part 2: Theoretical Foundations. Slide 6



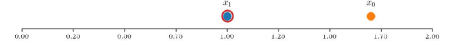
Recap: Dynamic Time Warping

- Compute cost matrix $C \in \mathbb{R}^{N \times M}$
- $C(n, m) = c(x_n, y_m)$ with cost function $c: \mathcal{F}_X \times \mathcal{F}_Y \rightarrow \mathbb{R}$
- Goal: compute minimum cost over the cost matrix, taking valid paths $P \in \mathcal{P} = \{ \text{valid paths} \}$
- $\text{DTW}(C) = \min(\{ \sum_{(n,m) \in P} C(n, m) \mid P \in \mathcal{P} \})$
- Problem: min function does not have a continuous derivative!



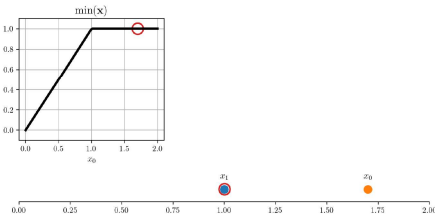
Differentiable Minimum Functions

- Investigate minimum function over $x = [x_0, x_1]$
- Argmin \circ changes when $x_0 > x_1$



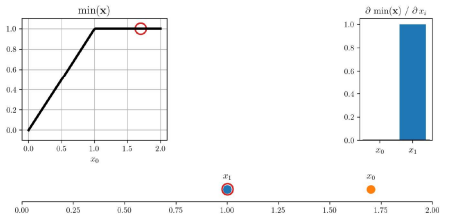
Differentiable Minimum Functions

- Investigate minimum function over $x = [x_0, x_1]$
- Argmin \circ changes when $x_0 > x_1$
- Minimum function: "edge" at $x_0 = 1.0$



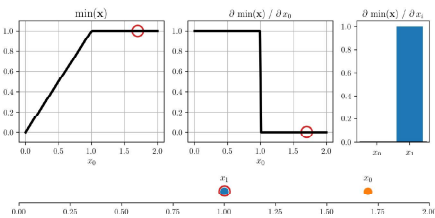
Differentiable Minimum Functions

- Investigate minimum function over $x = [x_0, x_1]$
- Argmin \circ changes when $x_0 > x_1$
- Minimum function: "edge" at $x_0 = 1.0$
- Argmin (derivative): hard decision for x_0 or x_1



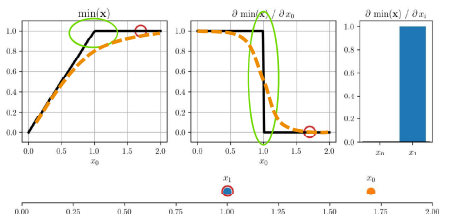
Differentiable Minimum Functions

- Investigate minimum function over $x = [x_0, x_1]$
- Argmin \circ changes when $x_0 > x_1$
- Minimum function: "edge" at $x_0 = 1.0$
- Argmin (derivative): hard decision for x_0 or x_1
- Gradient: discontinuity when argmin changes

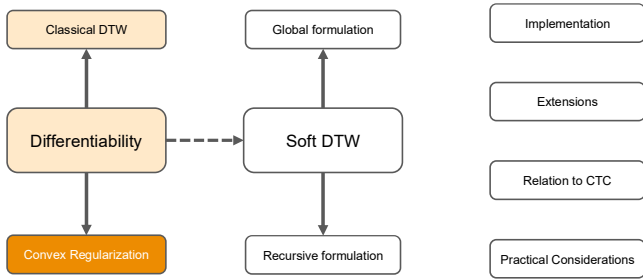


Differentiable Minimum Functions

- Gradient: discontinuity when argmin changes
- Why is the discontinuity problematic?
 - "Winner takes it all"
 - Toy example: hard choice between x_0 and x_1
 - Alignment: hard choice for one path
 - Full gradient flow goes to a single path!
 - What if we are not sure about the best path?
- "Soft Choice" between x_0 and x_1 ?



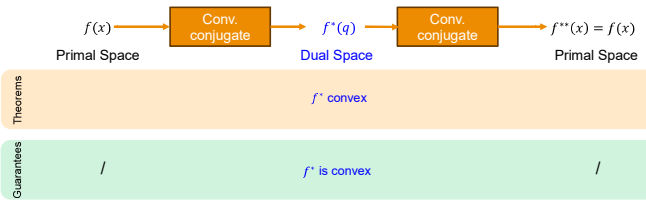
Overview



Smoothing Functions via Convex Regularization

M. Blondel and V. Roulet, "The elements of differentiable programming", arxiv preprint, 2025

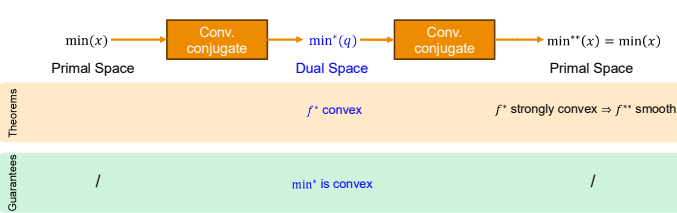
- Represent an optimization problem as a „dual problem“
- Transform: „convex conjugate“



Smoothing Functions via Convex Regularization

M. Blondel and V. Roulet, "The elements of differentiable programming", arxiv preprint, 2025

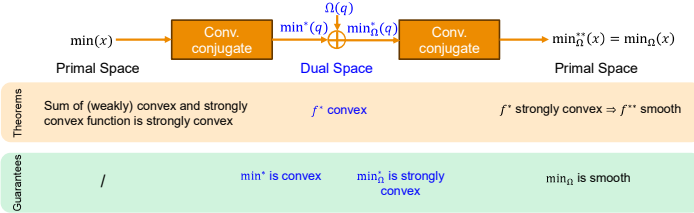
- Calculate convex conjugate for minimum function
- Guarantee: \min^* is convex
- No guarantee for \min^{**}



Smoothing Functions via Convex Regularization

M. Blondel and V. Roulet, "The elements of differentiable programming", arxiv preprint, 2025

- Can we enforce strong convexity in the dual space?
- Add a strongly convex regularizer Ω to \min^*
- \min_Ω is guaranteed to be smooth!



Softmin

- Popular choice for $\Omega(q)$: Entropy function

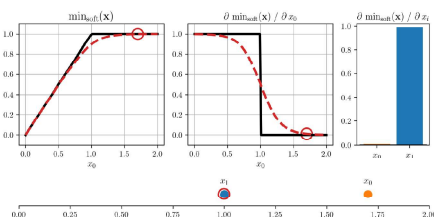
$$\Omega(q) = \sum_{q_i \in q} q_i \log q_i$$

- Solving for optimum yields closed-form "softmin"

$$\min_{\text{soft}}^{\gamma}(x) = -\gamma \log \sum_i \exp\left(-\frac{x_i}{\gamma}\right)$$

$$\dots \text{ with gradient } [\nabla \min_{\text{soft}}^{\gamma}]_i = \frac{\exp(-\frac{x_i}{\gamma})}{\sum_j \exp(-\frac{x_j}{\gamma})}$$

- Temperature parameter γ controls smoothness



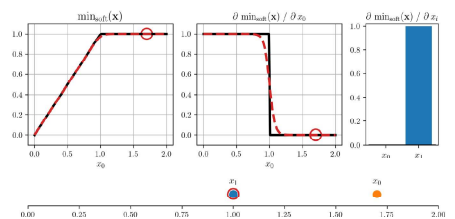
Softmin Temperature

- Softmin:

$$\min_{\text{soft}}^{\gamma}(x) = -\gamma \log \sum_i \exp\left(-\frac{x_i}{\gamma}\right)$$

$$\dots \text{ with gradient } [\nabla \min_{\text{soft}}^{\gamma}]_i = \frac{\exp(-\frac{x_i}{\gamma})}{\sum_j \exp(-\frac{x_j}{\gamma})}$$

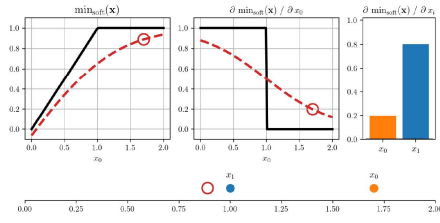
- Small temperature γ : approach hardmin



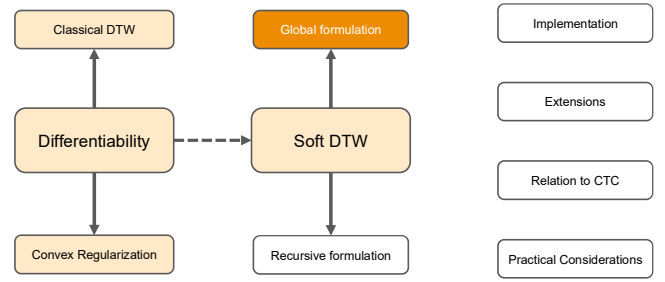
Softmin Temperature

- Softmin:

$$\min_{\text{soft}}^{\gamma}(x) = -\gamma \log \sum_i \exp\left(-\frac{x_i}{\gamma}\right)$$
 ... with gradient $[\nabla \min_{\text{soft}}^{\gamma}]_i = \frac{\exp(-\frac{x_i}{\gamma})}{\sum_j \exp(-\frac{x_j}{\gamma})}$
- Small temperature γ : approach hardmin
- High temperature γ : approach averaging
- We always compute a lower bound for min!



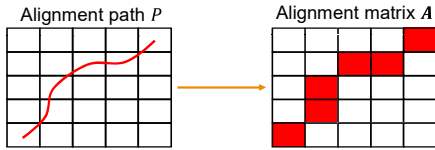
Overview



SoftDTW

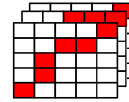
Define alignment paths $P \in \mathcal{P}$ as equivalent alignment matrices $A \in \mathcal{A}$ via a one-hot encoding $A \in \mathbb{R}^{N \times M}$

$$A(n, m) = \begin{cases} 1, & \text{if } (n, m) \in P, \\ 0, & \text{else.} \end{cases}$$



SoftDTW

- Set of valid alignments $\mathcal{A} = \{A_1, \dots, A_I\} =$



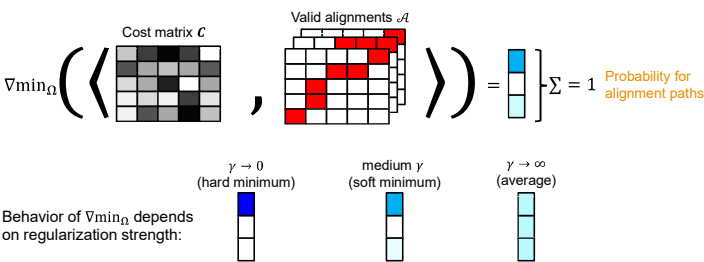
M. Cuturi and M. Blondel, Soft-DTW: a differentiable loss function for time series, ICML 2017

- $\text{SDTW}(C) = \min_{\Omega}(\{C, A\} \mid A \in \mathcal{A})$

$$= \min_{\Omega} \left(\left\langle \begin{matrix} \text{Cost matrix } C \\ \text{Valid alignments } \mathcal{A} \end{matrix}, \begin{matrix} \text{Valid alignments } \mathcal{A} \end{matrix} \right\rangle \mid A \in \mathcal{A} \right)$$

Gradient of SoftDTW

Gradient of minimum function $\nabla \min_{\Omega}$ denotes the influence of individual alignment paths on total cost



Gradient of SoftDTW

- Define gradient $H \in \mathbb{R}^{N \times M}$ as influence of cost cell $C(n, m)$ on total alignment cost $\text{SDTW}(C)$:

$$H(n, m) = \frac{\partial \text{SDTW}(C)}{\partial C(n, m)}$$
- Gradient H is sum of alignment matrices A , weighted with gradient $\nabla \min_{\Omega}$

$$H = \sum \left(\begin{matrix} \text{Valid alignments } \mathcal{A} \\ \nabla \min_{\Omega} \end{matrix} \right) = \begin{matrix} \text{Gradient / expected alignment} \end{matrix}$$

$$= \sum_{i=1}^{|\mathcal{A}|} A^{(i)} \cdot \nabla \min_{\Omega}^{(i)}$$

Gradient for different regularization strengths

$\gamma \rightarrow 0$ (hard minimum) $H = \sum \text{Cost matrix } C \cdot \text{Valid alignments } \mathcal{A} = \text{DTW}$

medium γ (soft minimum) $H = \sum \text{Cost matrix } C \cdot \text{Valid alignments } \mathcal{A} = \text{SoftDTW}$

$\gamma \rightarrow \infty$ (average) $H = \sum \text{Cost matrix } C \cdot \text{Valid alignments } \mathcal{A} = \text{Average over all valid paths}$

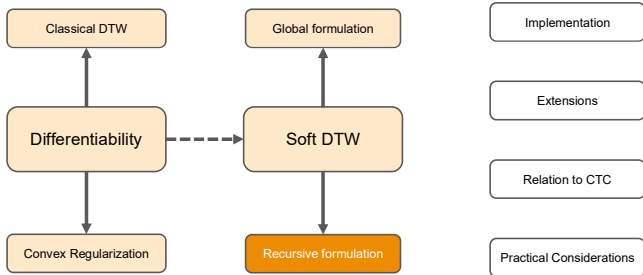
Summary: SDTW in Global Formulation

$\text{SDTW}(C) = \min_{\mathcal{A}} \left(\langle C, \mathcal{A} \rangle \mid \mathcal{A} \in \mathcal{A} \right)$

$H = \sum \text{Cost matrix } C \cdot \text{Valid alignments } \mathcal{A} = \nabla \min_{\mathcal{A}} \langle C, \mathcal{A} \rangle$

Problem: $|\mathcal{A}|$ grows exponentially!

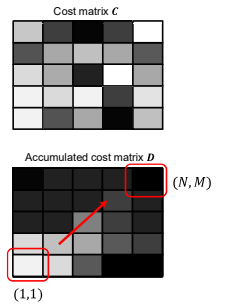
Overview



A Recursive Algorithm for SDTW: Forward

M. Cuturi and M. Blondel, "Soft-DTW: a differentiable loss function for time series", ICML 2017

- Compute SDTW recursively with dynamic programming
- Input: local cost matrix $C \in \mathbb{R}^{N \times M}$
- Output: accumulated cost matrix $D \in \mathbb{R}^{N \times M}$
- $D(n, m)$: minimum cost over all paths leading to (n, m)
- $D(N, M) = \text{SDTW}(C)$
- Requirements:
 - Boundary conditions: start in $(1, 1)$, end in (N, M)
 - Allowed step sizes $\mathcal{S} = \{(1, 0), (0, 1), (1, 1)\}$

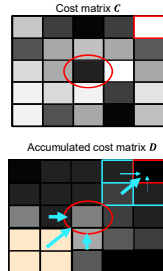


A Recursive Algorithm for SDTW: Forward

M. Cuturi and M. Blondel, "Soft-DTW: a differentiable loss function for time series", ICML 2017

Recursion: $D(n, m) = \min_{(i, j) \in \mathcal{S}} \{ C(n, m) + D(n-i, m-j) \}$

$D(N, M) = \text{SDTW}(C)$



Computational complexity: $\mathcal{O}(NM)$
(linear in sequence lengths)

Relation of Global and Recursive Formulation

A. Mensch and M. Blondel, "Differentiable dynamic programming for structured prediction and attention", ICML 2018

Global formulation: $\text{SDTW}^{\text{gl}}(C) = \min_{\mathcal{A}} \{ \langle C, \mathcal{A} \rangle \mid \mathcal{A} \in \mathcal{A} \}$

Recursive formulation: $D(n, m) = \min_{(i, j) \in \mathcal{S}} \{ C(n, m) + D(n-i, m-j) \}$

$\text{SDTW}^{\text{rec}}(C) = D(N, M)$

$= \min_{\mathcal{A}} \left(\langle C, \mathcal{A} \rangle \right) \leftrightarrow \text{Cost matrix } C \text{ and Acc. cost matrix } D$

Dividing the Global Problem

- 5 possible alignment paths



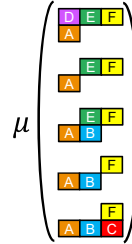
Cost matrix c

D	E	F
A	B	C

Dividing the Global Problem

- 5 possible alignment paths
- Compute minimum cost over these 5 paths

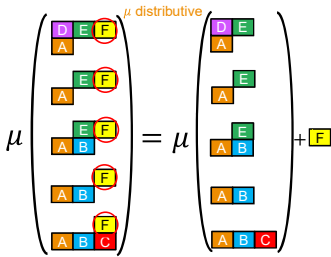
D	E	F
A	B	C



Dividing the Global Problem

- Goal: facilitate the computation
- F is the end of every path
- Move it out of the min function!

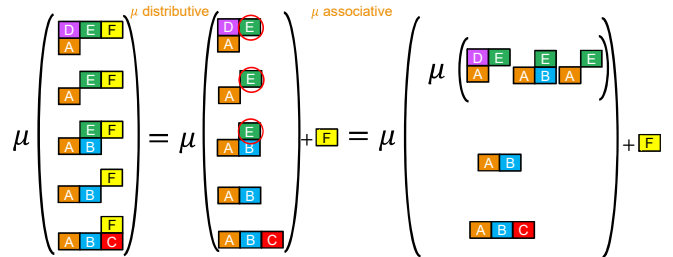
D	E	F
A	B	C



Dividing the Global Problem

- Divide into sub-problems
- For example, all paths ending in E

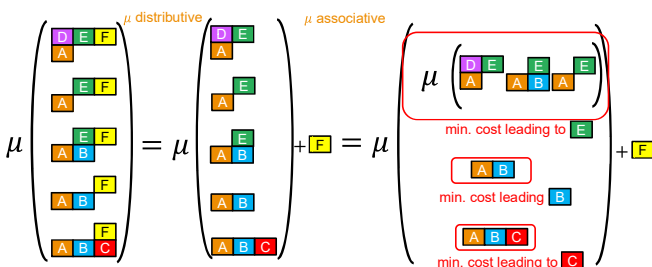
D	E	F
A	B	C



Dividing the Global Problem

- We found solutions for sub-problems!

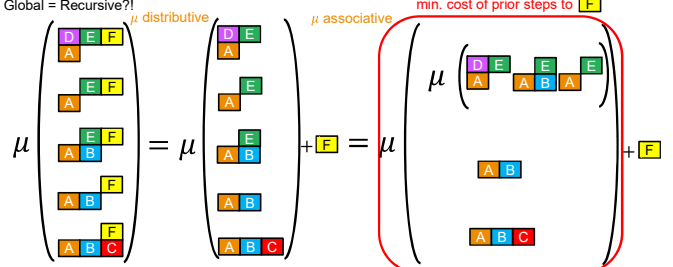
D	E	F
A	B	C



Dividing the Global Problem

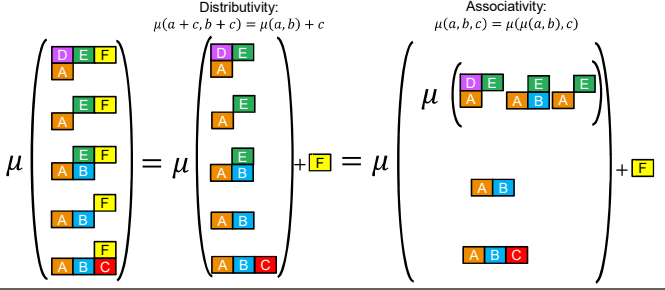
- We found solutions for sub-problems!
- We have established the recursion!
- Global = Recursive?!

D	E	F
A	B	C



Dividing the Global Problem

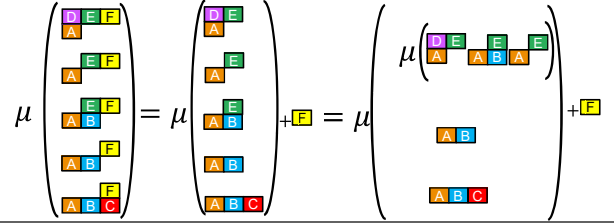
Requirements



Dividing the Global Problem

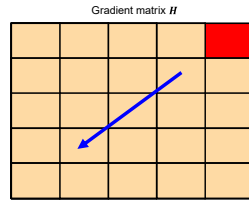
Theoretical guarantees

- Theorem: If μ is a regularized minimum function \min_{Ω} , distributivity and associativity are fulfilled if and only if $\Omega(q) = (q, \log q)$
- Global and recursive solutions are identical for $\mu = \text{softmin}!$



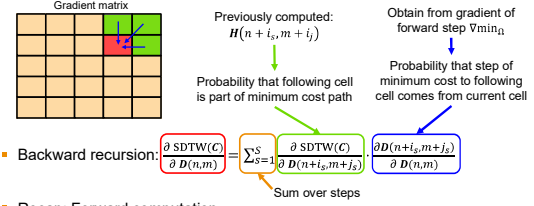
A Recursive Algorithm for SDTW: Backward

- Follow the traditional DTW backtracking algorithm to calculate gradient matrix $H \in \mathbb{R}^{N \times M}$
- Define gradient element $H(n, m)$ as the probability of the minimum cost path going through cell (n, m)
- Initialize the recursion: $H(N, M) = 1$ (all paths end in (N, M))
- Compute cells $H(n, m)$ with a recursion in reverse order



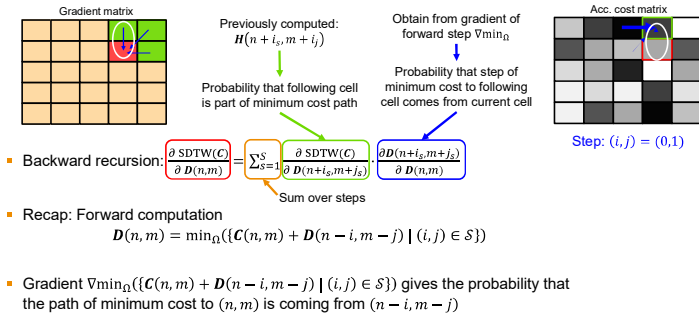
M, Cuturi and M, Børdet, „Soft-DTW: a differentiable loss function for time series, ICML 2017

A Recursive Algorithm for SDTW: Backward

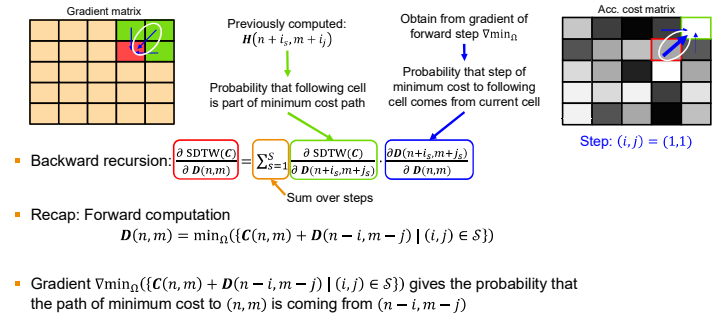


- Recap: Forward computation
 $D(n, m) = \min_{\Omega} \{C(n, m) + D(n - i, m - j) \mid (i, j) \in \mathcal{S}\}$
- Gradient $\nabla \min_{\Omega} \{C(n, m) + D(n - i, m - j) \mid (i, j) \in \mathcal{S}\}$ gives the probability that the path of minimum cost to (n, m) is coming from $(n - i, m - j)$

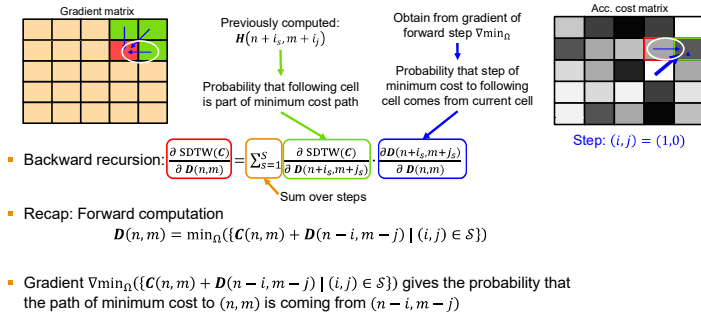
A Recursive Algorithm for SDTW: Backward



A Recursive Algorithm for SDTW: Backward



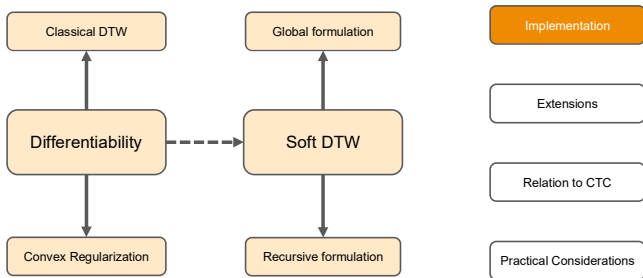
A Recursive Algorithm for SDTW: Backward



Summary: A Recursive Algorithm for SDTW

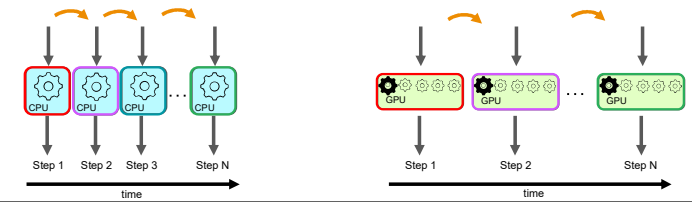
- Recursive forward pass of SDTW = „soft“ version of classical DTW forward pass (differentiable minimum instead of hard minimum)
- Recursive backward pass of SDTW = „soft“ version of classical DTW backtracking (probabilities for paths instead of hard decision)
- Recursion is identical to global formulation if $\min_{\Omega} = \text{softmin}$
- Runtime linear in sequence lengths $\mathcal{O}(NM)$

Overview



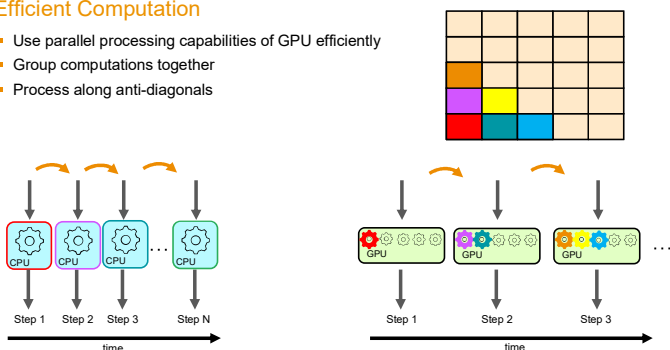
Efficient Computation

- SDTW recursion requires iterative processing
- Well-suited for CPUs
- Not efficient for GPUs



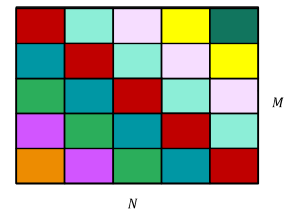
Efficient Computation

- Use parallel processing capabilities of GPU efficiently
- Group computations together
- Process along anti-diagonals



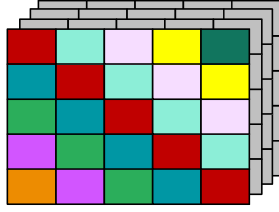
Efficiency & Implementation

- Elements along the anti-diagonals are independent of each other
- Number of „group“ computations: $\#diag = N + M - 1$
- Example: $N = M = 5$
 - Number of individual elements: $N \cdot M = 25$
 - Number of anti-diagonals (groups): $N + M - 1 = 9$
- The same holds for the backward pass



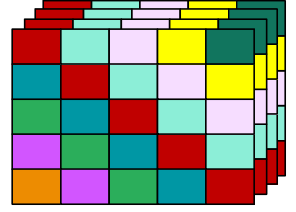
Batch Processing

- Independence along the batch dimension



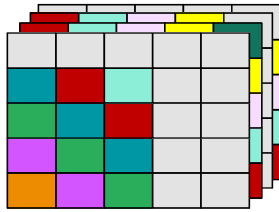
Batch Processing

- Independence along the batch dimension
- Group anti-diagonals together for all batch elements
- Number of groups doesn't change compared to single-matrix processing
- Batch processing over multiple cost matrices comes „free“

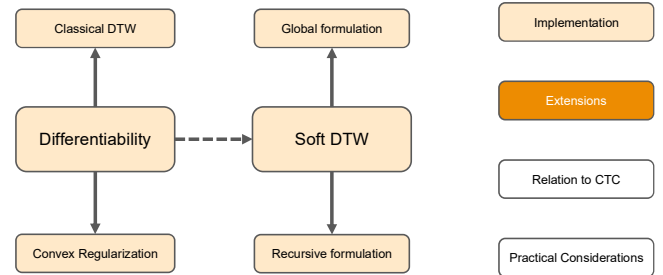


Batch Processing

- How to deal with difference sequence lengths in a batch?
- Pad all cost matrices to same size and concatenate
- Do group processing along anti-diagonals
- Skip computation if outside current sequence length



Overview



SDTW as Generalized Alignment Framework

- Objective: $\text{SDTW}(C) = \min_{\mathcal{A}} \langle C, A \rangle \mid A \in \mathcal{A}$

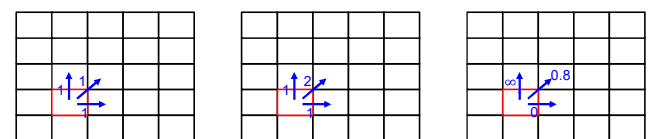
$$\text{SDTW}(C) = \min_{\mathcal{A}} \left(\left\langle \begin{array}{c} \text{Cost matrix } C \\ \begin{array}{|c|c|c|c|c|} \hline \text{[Grid of colored cells]} \\ \hline \end{array} \end{array}, \begin{array}{c} \text{Valid alignment paths } \mathcal{A} \\ \begin{array}{|c|c|c|c|c|} \hline \text{[Grid of red cells]} \\ \hline \end{array} \end{array} \right\rangle \right)$$

- SDTW provides an efficient framework for computing $\min_{\mathcal{A}} \langle C, A \rangle \mid A \in \mathcal{A}$
- We can relax constraints on the alignments \mathcal{A} to make SDTW mor flexible

SDTW with Variable Step Weights

- Choose flexible weights for every step
- Avoid diagonalization for equal sequence lengths
- Control influence of target repetition (horizontal step)
- Include prior knowledge on likelihood of certain steps
- Use step weight ∞ to „block“ certain steps

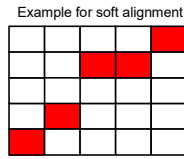
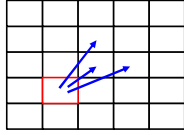
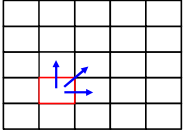
J. Zeile, M. Krause, and M. Müller, „Soft Dynamic Time Warping with Variable Step Weights“, ICASSP 2024



SDTW with Flexible Step Sizes

- Skip certain frames or targets
- 2-1-softDTW

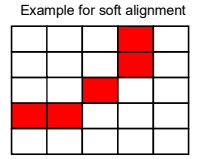
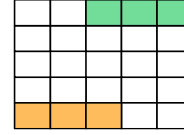
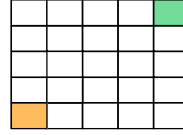
J. Zeitler and M. Müller, „A Unified Perspective on CTC and SDTW using Differentiable DTW“, submitted to IEEE Transactions of Audio, Speech, and Language Processing, 2025



SDTW with Flexible Boundary Conditions

- Subsequence-softDTW
- Prediction and target sequences do not need to align at the boundaries

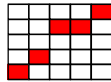
J. Zeitler and M. Müller, „Subsequence SDTW: A Framework for Differentiable Alignment with Flexible Boundary Conditions“, submitted to ICASSP 2026



SDTW as Generalized Alignment Framework

Flexible Step Sizes:

- Skip certain frames or targets
- 2-1-softDTW



J. Zeitler and M. Müller, „A Unified Perspective on CTC and SDTW using Differentiable DTW“, submitted to IEEE Transactions of Audio, Speech, and Language Processing, 2025

Flexible Step Weights:

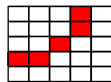
- Choose flexible weights for every step
- Avoid diagonalization for equal sequence lengths
- Control influence of target repetition (horizontal step)
- Include prior knowledge on likelihood of certain steps
- Use step weight ∞ to „block“ certain steps



J. Zeitler, M. Krause, and M. Müller, „Soft Dynamic Time Warping with Variable Step Weights“, ICASSP 2024

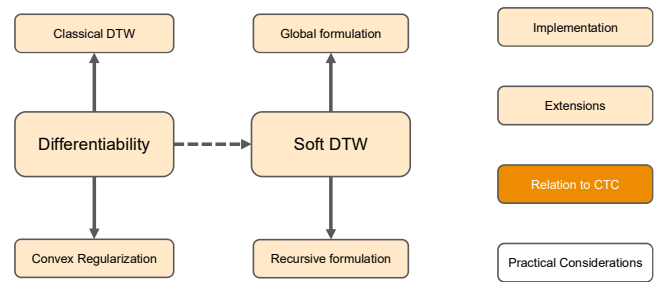
Flexible Boundary Conditions

- Subsequence-softDTW
- Prediction and target sequences do not need to align at the boundaries



J. Zeitler and M. Müller, „Subsequence SDTW: A Framework for Differentiable Alignment with Flexible Boundary Conditions“, submitted to ICASSP 2026

Overview

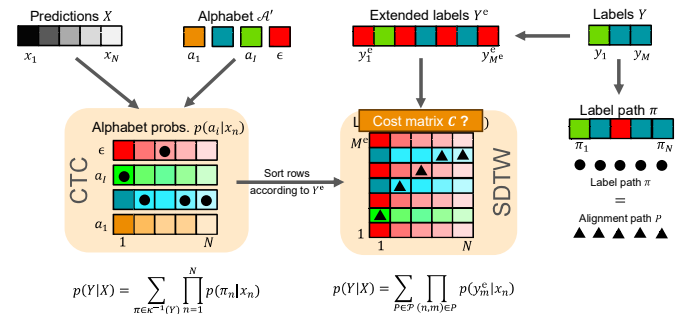


Relation to CTC

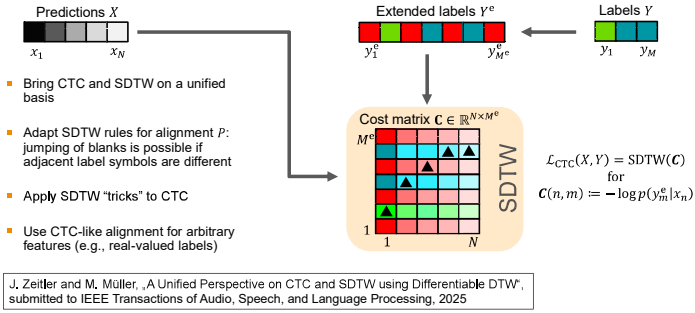
- CTC...
 - has a finite target alphabet
 - is widely used in speech processing
 - has an unintuitive formulation
- SDTW...
 - is based on an arbitrary cost matrix
 - is widely used in signal processing
 - has an intuitive formulation
- Both algorithms align sequences and are fully differentiable
- Can we establish a connection?

CTC Reformulation

A. Graves et al., „Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks“, ICML 2006

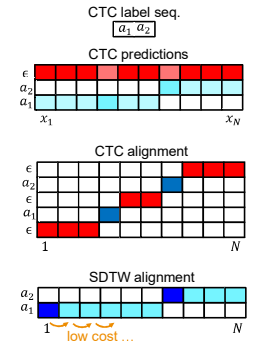


CTC Reformulation

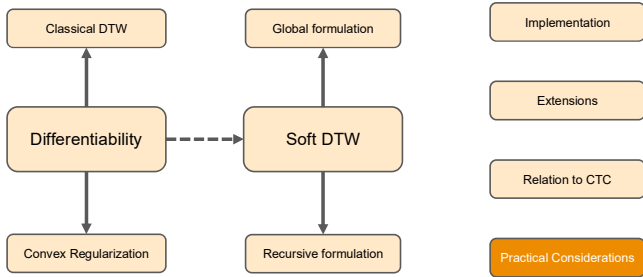


Dominance of Blank Symbol in CTC

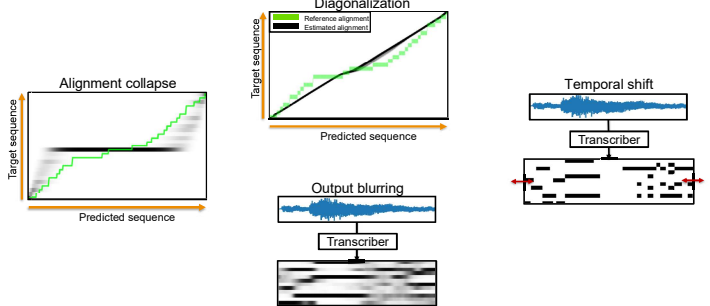
- CTC predictions dominated by blank
- Blank alignment is always "cheap" and leads to stabilization
- Spiky alignment of labels
- Predictions get even more blank-dominated
- Stabilization in SDTW: low cost for horizontal step (label repetition)
- Eliminate need for blank symbol



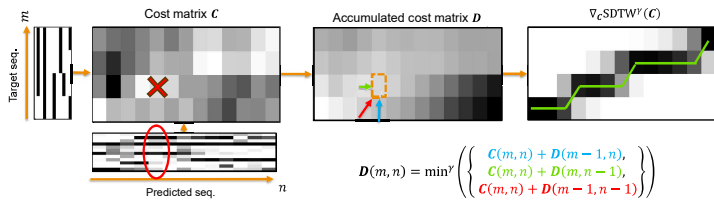
Overview



Common SDTW Problems & Pitfalls

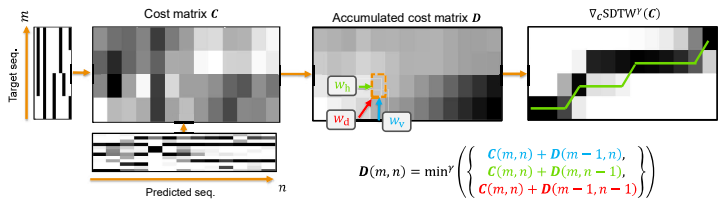


Problem 1: Alignment Collapse



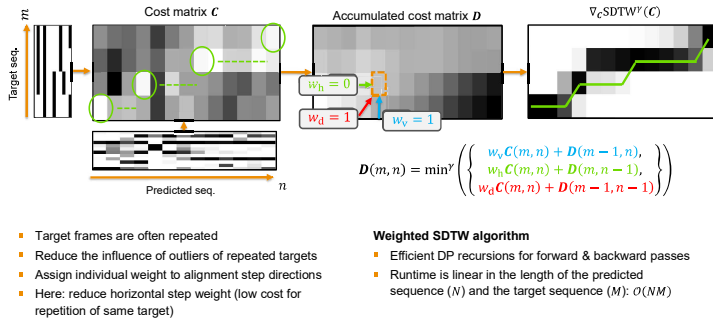
- Single corrupted predictions cause high values in cost matrix
- Alignment collapses to few target frames
- Training diverges

Problem 1: Alignment Collapse



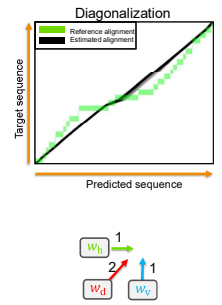
- Target frames are often repeated
- Reduce the influence of outliers of repeated targets
- Assign individual weight to alignment step directions

Problem 1: Alignment Collapse



Problem 2: Diagonalization

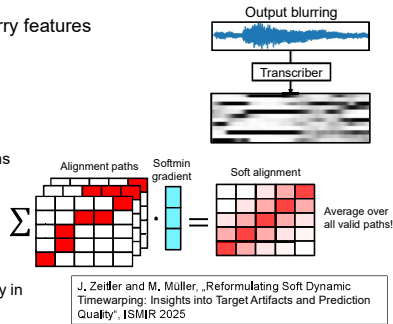
- Problem: computed SDTW alignment focuses only on the main diagonal
- Cause:
 - Equal lengths of prediction and target sequences
 - Sequences of equal length can be aligned using only diagonal steps
 - Taking one diagonal step is cheaper than taking a vertical and horizontal step („around the corner“)
- Solution:
 - Choose SDTW with step weights
 - Set a higher step weight to diagonal step (e.g., 1-1-2)
 - A diagonal step gets the same weight as a horizontal + vertical step



Problem 3: Output Blurring

- Problem: transcriber learns only blurry features

- Cause:
 - Softmin temperature $\gamma \rightarrow \infty$
 - Softmin becomes averaging
 - SDTW gradient is average over all paths
 - Blurry gradient leads to blurry features

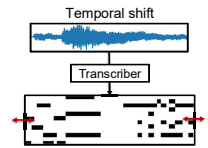


- Solution:
 - Reduce softmin temperature $\gamma \approx 1$
 - If high softmin temperature is necessary in initial training, do gradual reduction

Problem 4: Temporal Shift

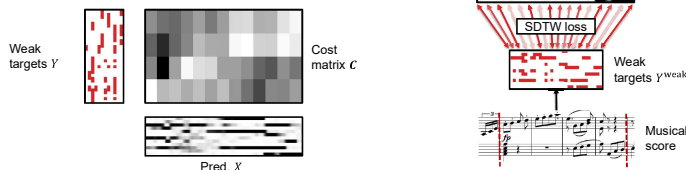
- Problem: small temporal shift between input and predictions

- Cause:
 - SDTW computes flexible alignment between predictions and weak targets
 - Alignment cost is invariant of (small) temporal shift
- Solutions:
 - Identify temporal shift of trained model and compensate during inference
 - Use a DNN with small temporal receptive field (1-1 mapping of input to output frames)
 - Use an auxiliary loss to evaluate smilarity between the predictions and the input



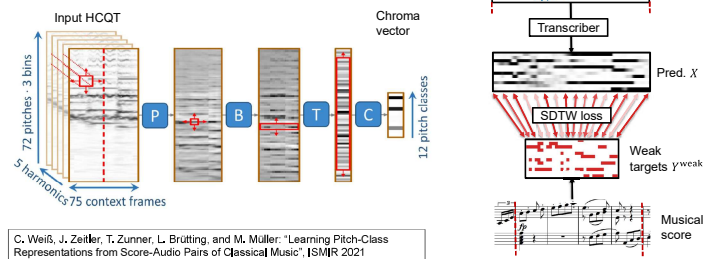
Multi-Pitch & Pitch Class Estimation

- Annotate corresponding segments in the input audio and the musical score (typically 10s – 30s)
- Retrieve weak targets from the musical score
- Weak targets represent sequence of simultaneously active notes, but no information about duration
- Cost function: Binary Cross-Entropy (BCE)
- $C(n, m) = BCE(x_n, y_m)$



Multi-Pitch & Pitch Class Estimation

Parameter-efficient choice for deep learning of pitch (class) activations: musically motivated CNN [Weiss2021]



**AUDIO
LABS**

APPENDIX

Differentiable via Convex Regularization Convex Optimization

- Let $f: \mathbb{R}^D \rightarrow \mathbb{R} \cup \{\infty\}$ denote a function with domain $\text{dom}(f) := \{x | f(x) < \infty\}$
- Definition of convex conjugate: $f^*(y) := \sup_x (\langle x, y \rangle - f(x))$; $\text{dom}(f^*) := \{y | f(y) < \infty\}$
- Define Indicator function $I_C(x) := \begin{cases} 0, & x \in C \\ \infty, & x \notin C \end{cases}$
- Choose $f(x) = \max(x)$
- $f^*(y) = \sup_x (\langle x, y \rangle - \max(x)) = I_{\Delta^D}(y)$

$$= \begin{cases} 0, & \text{if } y \in \Delta^D \\ \infty & \text{else} \end{cases}$$

Differentiable via Convex Regularization Convex Optimization

A. Mensch and M. Blondel, "Differentiable dynamic programming for structured prediction and attention", ICML 2018

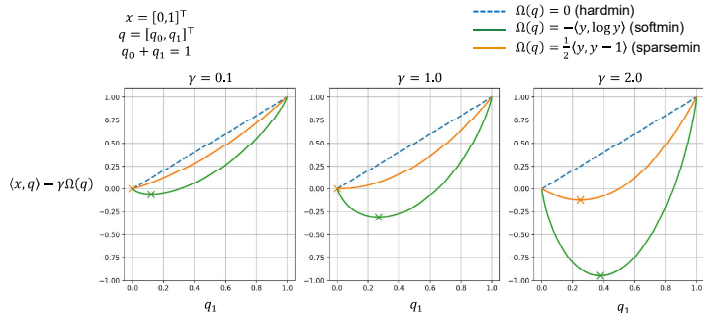
- Theorem: f^* is convex, even if f is non-convex
- Theorem: If f is strongly convex over $\text{dom}(f)$, then f^* is smooth over $\text{dom}(f^*)$
- Add a strongly convex regularizer $\Omega(q)$:
- $f_\Omega^*(q) = I_{\Delta^D} + \Omega(q)$
- Transform to primal space:
- $f_\Omega^*(x) = \sup_q (\langle x, q \rangle - I_{\Delta^D} - \Omega(q)) = \max_{q \in \Delta^D} (\langle x, q \rangle - \Omega(q)) = \max_\Omega(x)$

$$= \begin{cases} -\infty & \text{if } q \notin \Delta^D \\ \langle x, q \rangle - \Omega(q) & \text{if } q \in \Delta^D \end{cases}$$
- $\max_\Omega(x)$ is now smooth, i.e., has a continuous derivative
- As $\max(x) = \max_{q \in \Delta^D} \langle x, q \rangle$, the function $\max_\Omega(x)$ can be seen as the max function plus an additional regularizer
- For minimum functions, we analogously have $\min_\Omega(x) = -\max_\Omega(-x)$
- Add a temperature parameter γ : $\max_\Omega^\gamma(x) = \max_{q \in \Delta^D} (\langle x, q \rangle - \gamma \Omega(q))$

Differentiable via Convex Regularization Common convex regularizers Ω :

- Shannon entropy: $\Omega(y) = -\langle y, \log y \rangle$
 - Solving for optimum yields closed-form "softmax": $\min_{\text{soft}}^\gamma(x) = -\gamma \log \sum_i \exp(-\frac{x_i}{\gamma})$
 - ... with gradient $[\nabla \min_{\text{soft}}^\gamma]_i = \frac{\exp(-\frac{x_i}{\gamma})}{\sum_j \exp(-\frac{x_j}{\gamma})}$
- Gini entropy: $\Omega(y) = \frac{1}{2} \langle y, y - 1 \rangle$
 - Solving for optimum yields "sparsemin": $\min_{\text{sparse}}^\gamma(x) = \langle y^*, x \rangle + \frac{\gamma}{2} \|y^*\|_2^2 - \frac{\gamma}{2}$
 - ... with gradient $\nabla \min_{\text{sparse}}^\gamma(x) = \arg \min_{y \in \Delta^D} \|y + \frac{x}{\gamma}\|_2^2 = y^*$

Minimum functions with convex regularization



Minimum functions with convex regularization

