

1 We thank all the reviewers for their valuable comments. We carefully address all the raised issues accordingly below.

2 **1 Response to Reviewer 1:** We appreciate your positive feedback.

3 **2 Response to Reviewer 2:**

4 Thank you for the comments. We will proofread and improve the readability. Please see details below.

5 **2.1 Q: Definition of the forward pass for the main network $G_{x,w}$. How they go from x to z .**

6 A: $G_{x,w}$ is the network parameterized by the kernel x and linear layers' weights w . Input is y , and output $z = G_{x,w}(y)$.

7 **2.2 Q: In line 4 of Algorithm 1, it refers to $f_t(x, y)$ in Eq. (2), but it is not found there.**

8 A: Sorry for the confusion. $f_t(x, y)$ is the same as $dtw^2(x, y)$ in Eq.2, with the subscript t being the iteration index.

9 **2.3 Q: "Assume that \tilde{x} is close to global minima". This assumption may never hold.**

10 A: We shouldn't abuse the term "assumption". In fact, this "assumption" is not necessary for the analysis. We simply
11 want to emphasize that we are interested in areas having many local minima, which happen to be always close to the
12 global minimum from empirical observations. Note that the local-minima areas of interests can be anywhere, and the
13 analysis still holds. The only issue arises when the left region of w 's stationary point $x^{(w)*}$ is located to the right of
14 region w . In this case, we can combine w and the center region u to form a new quasi-bowl center region $u' = [w, u]$
15 (similar to the analysis at the top of Page 6 in the paper), and the analysis still holds. We will remove this "assumption".

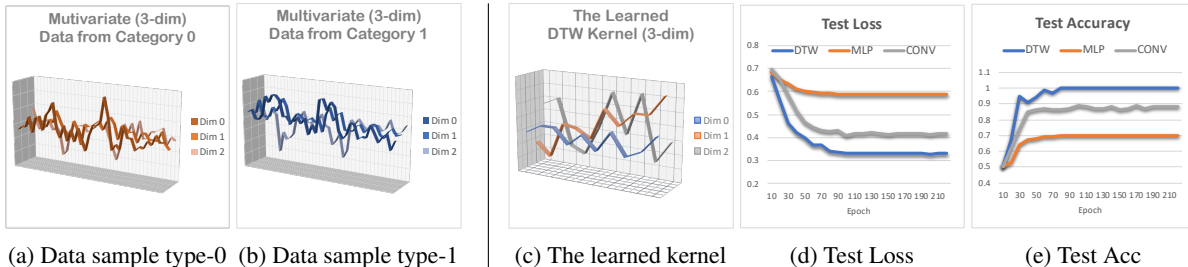
16 **2.4 Q: Proof of convergence: after t iterations, it will be eps-close to the exact DTW.**

17 Since it is highly non-convex and non-smooth, to the best of our knowledge, without strong assumptions it is unlikely to
18 prove global convergence. We also analyze the escaping behavior (from local minima) but not the global convergence.
19 The shape of DTW loss is identified through the alternating algorithm, which is one of the contributions in our paper.

20 **2.5 Q: The idea will only work for the univariate time series.**

21 A: The proposed approach works for both univariate and multivariate time series. Multivariate DTWs are often
22 computed in two forms: MDTW-I and MDTW-D [1]. MDTW-I treats each dimension independently, so it is simply
23 a stack of multiple univariate DTWs, thus directly applies to our method. MDTW-D needs to compute multivariate
24 distance $\|x_i - y_j\|^2$, $x_i \in \mathbf{R}^m$, $y_j \in \mathbf{R}^m$ in the Dynamic Programming step, instead of the scalar version $\|x_i - y_j\|^2$.
25 As long as the norm is well defined, e.g., Euclidean distance, the forward pass and the backpropagation are performed
26 in the same manner. We can even define other distances, as long as their gradients w.r.t. to x can be computed.

27 We run a 3-dim multivariate time series classification task here, using MDTW-D and Euclidean distance in our approach.
28 The experiment settings follow Section 6.1 in the paper. The following figures show: one sample data from each
29 category, the learned kernel, test loss and test acc comparison. Our method (DTW) outperforms others.



(a) Data sample type-0 (b) Data sample type-1

(c) The learned kernel

(d) Test Loss

(e) Test Acc

30 **2.6 Q: The experiments are insufficient. One dataset from UCR repo is not enough.**

31 A: The UCR repo is a collection of a large variety of time series data, being the standard benchmark repo in related
32 publications. We performed comprehensive experiments on 85 datasets from UCR repo (details provided in appendix
33 of the paper), to make a fair comparison with Soft-dtw (they only use the UCR repo in their original paper as well).
34 Additional experiments on synthetic data are also carried out. But we agree that more datasets should be considered in
35 this area and we are exploring them in the final version and subsequent work.

36 **2.7 Q: Some handwaving claims such as interpretability.**

37 A: Some interpretability is revealed by the learned kernel's shape matching true patterns in the data (Fig 4, 5 in paper).

38 **3 Response to Reviewer 3:** Thank you for your comments and please see details below.

39 **3.1 Q: For example, it will be helpful to know how to decide the number of DTW layers.**

40 A: Empirically speaking, 1 or 2 layers of DTW are good enough. We observe no clear improvement with more than 2
41 layers, but the model complexity would be increased dramatically.

42 **3.2 In alg 1's INPUT, kernels are initially set as input. But they are randomly initialised in OUTPUT part.**

43 A: Kernels x and weights w are not inputs but model parameters being randomly initialized. We will clarify it.

44 **3.3 Q: It will be helpful if the proposed method can be tested with more real datasets for the application part.**

45 A: We agree with testing more datasets. Due to page limit, we randomly select the Haptics dataset from UCR repo as an
46 expressive example in the application section, but more datasets will be considered (please also see response 2.6).

47 [1] Mohammad Shokoohi-Yekta, Bing Hu, Hongxia Jin, Jun Wang, and Eamonn Keogh. Generalizing dtw to the
48 multi-dimensional case requires an adaptive approach. *Data mining and knowledge discovery*, 31(1):1–31, 2017.