

1 We thank the reviewers for their constructive comments and suggestions. We respond to each point individually.

2 **R1:** *More description to highlight the unique contribution.*

3 Compared with previous metric learning methods by using low-rank/online/stochastic strategies, that still encounter
4 the scalability problem when handling large data, our paper has two unique contributions as follows. (1) Our method
5 embeds the triplet constraints into a matrix, and further reduces the size of involved matrices by replacing \mathbf{Y} with
6 $\mathbf{B}\mathbf{V}^T$, which ensures the existence of the optimal solution in the reduced matrices. (2) By substituting the closed-form
7 optimal solution of \mathbf{s} , the optimization of positive semidefinite matrices is converted into the optimization on the Stiefel
8 manifold, which can be optimized more efficiently. These two contributions significantly reduce the complexity and the
9 size of involved matrices, which makes our method scalable to both high dimensions and large numbers of samples.

10 **R1:** *Explain the difference with some recent methods.*

11 We will add the missing references in “Related Work” and explain the difference with them in the revision. *Mini-SGD*
12 [Qian et al. ML 2015] uses the mini-batch strategy to optimize the metric matrix in the positive semidefinite cone.
13 Using the online metric learning strategy, *OPML* [Li et al. PR 2018] introduces a closed-form formula for updating the
14 metric matrix. However, the above two methods all require updating a large $D \times D$ matrix (D is the dimensionality of
15 original data) at each iteration, and they can only learn from samples with hundreds of dimensions in their experiments.
16 In contrast, our method models the metric learning problem on the Stiefel manifold with much smaller size $r \times d$, which
17 significantly reduces the complexity and memory usage in optimization. In our experiments, the datasets can be with up
18 to one million dimensions (see Table 1 in our paper).

19 **R1:** *Theoretical results using $\mathbf{B}\mathbf{V}^T$ to replace \mathbf{Y} to ensure performance, and differences with anchor-based strategy.*

20 (1) If the linear equation system $\mathbf{Y} = \mathbf{L}\mathbf{X}$ has a feasible solution of \mathbf{L} , it can always be transformed into a full-rank
21 linear equation system $\mathbf{Y}\mathbf{V} = \mathbf{L}\mathbf{U}\Sigma$. According to Theorem 1 in our paper, since $\mathbf{Y}\mathbf{V} = \mathbf{B}$, all the possible solutions
22 can be covered by $\mathbf{B} \in \mathbb{R}^{d \times r}$, which ensures the performance of accelerated low-rank metric learning.

23 (2) For the anchor-based strategy in *AnchorGraph* [Liu et al. ICML 2010], the anchors are “local samples” with high
24 dimension D , and the algorithm represents each data point as a convex combination of its closest anchors. In contrast,
25 our method optimizes a smaller matrix $\mathbf{B} \in \mathbb{R}^{d \times r}$ to obtain the global optimal solution of a larger matrix \mathbf{L} , which is
26 different with the “local sampling” strategy used in *AnchorGraph*, as will be added in “Related Work” in the revision.

27 **R2:** *The upper bound of r .*

28 As a preprocessing parameter, the upper bound of r is an empirical value. In general, a larger r will increase memory
29 and computational costs, while a smaller r may lose some intrinsic values. In our preliminary experiments under dataset
30 “TDT2”, when r varies in 500, 1000, 2000, 3000 and 4000, the accuracy just slightly changes to 0.948, 0.964, 0.965,
31 0.963 and 0.948, respectively. This shows that the variation of r within a reasonable range merely affects computational
32 time and memory, but has little effect on accuracy. For fair comparison, we use the same r for all compared methods.

33 **R2:** *What about the approximation loss of $\mu(x)$, and how to optimize the problem without using $\mu(x)$?*

34 The authors guess the approximation loss concerned by the reviewer is the difference between $\mu(x) = -\log(\sigma(-x))$
35 and $\max(0, x)$, where its maximum value is $-\log \frac{1}{2} \approx 0.69$ at the point $x = 0$. If not using $\mu(x)$, our model can still
36 be solved. However, due to the discontinuous gradient, the convergence requires more iterations, and the results are
37 less stable. For example, under the dataset “TDT2” in five repeated runs, when using $\mu(x)$, the average number of
38 iterations, average accuracy, and standard deviation are 6.6, 0.962, and 0.0008, respectively. In contrast, when only
39 using $\max(0, x)$, these values become 12.6, 0.955, and 0.012, respectively.

40 **R2:** *Theoretical analysis of stochastic strategy.*

41 The theoretical analysis of the stochastic strategy which updates \mathbf{L} in step sizes by $1/\sqrt{I}$ decay can refer to the reference
42 [14] in our paper. We will add more description for [14] about the stochastic strategy in the revision.

43 **R2:** *Parameter sensitivity examinations.*

- 44 - The examination of different mini-batch sizes has been provided in Fig. 4, where N_t is an indicator of mini-batch size.
- 45 - The examination of different m has been provided in Fig. 3. We will explicitly add the definition of m in the revision.
- 46 - d is the target rank value in low-rank metric learning tasks, which is usually set according to user preference. We use
47 $d = 100$ in all experiments for a fair comparison, which is a moderate value so that most tested methods can achieve
48 good performance on most datasets, as mentioned in Section 4.1.
- 49 - Triplet sizes: 5 triplets are randomly generated for each sample (which is a moderate value for most datasets), and the
50 same triplet sets are used for all tested methods for a fair comparison, as mentioned in Section 4.1.

51 **R3:** *Potential applications and more explicit gaps on the previous literature.*

52 A more detailed highlight of contributions and differences with recent methods can refer to 1st Answer for Reviewer 1.
53 Metric learning has been widely used in various areas, such as dimensionality reduction, feature extraction, and
54 information retrieval. Our method can be applied to the scenario of learning metrics quickly on large numbers of
55 high-dimensional data with limited computing resources. We will introduce some application scenarios in the revision.