

# Supplementary Material: Online Progressive Deep Metric Learning

Paper ID: 1535

## 1 Proof of Theorem 1

**Theorem 1.** Suppose  $M_t$  is positive-definite, then  $M_{t+1}$  given by the MOML update, i.e.,  $M_{t+1} = M_t - \gamma A_{t+1}$  is positive definite by properly setting  $\gamma$ .

*Proof.* As  $A_{t+1} = (x_{t+1} - x_p)(x_{t+1} - x_p)^\top - (x_{t+1} - x_q)(x_{t+1} - x_q)^\top$ , whose rank is 1 or 2, it has at most 2 non-zero eigenvalues. That is to say,  $\text{Tr}(A_{t+1}) = \lambda_1 + \lambda_2$ . Specifically, we can also easily get that,

$$-\|x_{t+1} - x_q\|_2^2 \leq \lambda(A_{t+1}) \leq \|x_{t+1} - x_p\|_2^2, \quad (1)$$

where  $\lambda(A_{t+1})$  means the eigenvalue of  $A_{t+1}$  (i.e.,  $\lambda_1$  or  $\lambda_2$ ). For each sample  $x$  is  $\ell_2$  normalized, the ranges of  $\|x_{t+1} - x_p\|_2^2$  and  $\|x_{t+1} - x_q\|_2^2$  vary from  $[0, 4]$ . Thus,

$$\lambda_{\min}(M_t) - 4\gamma \leq \lambda(M_t - \gamma A_{t+1}) \leq \lambda_{\max}(M_t) + 4\gamma. \quad (2)$$

When  $\gamma \leq \frac{1}{4}\lambda_{\min}(M_t)$ , it is guaranteed that the minimum eigenvalue of  $M_t - \gamma A_{t+1}$  is greater than zero. As the initial matrix  $M_1 = I$  is positive definite (i.e.,  $\lambda_{\min}(M_1) = 1$ ). By properly setting a small  $\gamma$ , the minimum eigenvalue of  $M_t - \gamma A_{t+1}$  is generally large than zero. Thus, the positive definiteness of  $M_{t+1} = M_t - \gamma A_{t+1}$  can be guaranteed. Same theoretical guarantee (i.e., the small perturbations of positive definite matrix) can also be found in the chapter 9.6.12 of [Petersen *et al.*, 2008].  $\square$

## 2 Proof of Theorem 2

**Theorem 2.** Let  $\langle x_1, x_p, x_q \rangle, \dots, \langle x_T, x_p, x_q \rangle$  be a sequence of triplet constraints where each sample  $x_t|_{t=1}^T \in \mathbb{R}^d$  has  $\|x_t\|_2 = 1$  for all  $t$ . Let  $M_t \in \mathbb{R}^{d \times d}$  be the solution of MOML at the  $t$ -th time step, and  $U \in \mathbb{R}^{d \times d}$  denotes an arbitrary parameter matrix. By setting  $\gamma = \frac{1}{R\sqrt{T}}$  (where  $R \in \mathbb{R}^+$ ), the regret bound is

$$\begin{aligned} R(U, T) &= \sum_{t=1}^T \ell(M_t) - \sum_{t=1}^T \ell(U) \\ &\leq \frac{1}{2} \|I - U\|_F^2 + \frac{32}{R^2}. \end{aligned} \quad (3)$$

*Proof.* According to the objective function of MOML, i.e.,

$$\Gamma = \arg \min_{M \succ 0} \frac{1}{2} \|M - M_{t-1}\|_F^2 + \gamma \left[ 1 + \text{Tr}(MA_t) \right]_+, \quad (4)$$

we denote  $\ell_t$  as the instantaneous loss suffered by MOML at each  $t$ -time step with the learnt  $M_t \in \mathbb{R}^{d \times d}$ , and denote by  $\ell_t^*$  the loss suffered by an arbitrary parameter matrix  $U \in \mathbb{R}^{d \times d}$ , which can be formalized as below:

$$\begin{aligned} \ell_t &= \ell(M_t; \langle x_t, x_p, x_q \rangle) = [1 + \text{Tr}(M_t A_t)]_+ \\ \ell_t^* &= \ell(U; \langle x_t, x_p, x_q \rangle) = [1 + \text{Tr}(U A_t)]_+, \end{aligned} \quad (5)$$

where  $A_t = (x_t - x_p)(x_t - x_p)^\top - (x_t - x_q)(x_t - x_q)^\top$ ,  $\text{Tr}$  denotes trace and  $[z]_+ = \max(0, z)$ . As  $\text{Tr}(M_t A_t)$  is a linear function, it is convex w.r.t  $M_t$  by natural. Besides, the hinge loss function  $[z]_+$  is a convex function (but not continuous at  $z = 0$ ) w.r.t  $z$ . Hence, the resulting composite function  $\ell_t(M_t)$  is convex w.r.t  $M_t$ . As  $\ell$  is a convex function, we can introduce the first-order condition as follow:

$$\ell(Y) \geq \ell(X) + \text{VEC}(\nabla \ell(X))^\top \text{VEC}(Y - X), \quad (6)$$

where  $X, Y \in \mathbb{R}^{d \times d}$ ,  $\text{VEC}$  denotes vectorization of a matrix, and  $\nabla \ell(X)$  is the gradient of function  $\ell$  at  $X$ .

Inspired by [Crammer *et al.*, 2006], we define  $\Delta_t$  to be  $\|M_t - U\|_F^2 - \|M_{t+1} - U\|_F^2$ . Then calculating the cumulative sum of  $\Delta_t$  over all  $t \in \{1, 2, \dots, T\}$ , we can easily obtain  $\sum_t \Delta_t$ ,

$$\begin{aligned} \sum_{t=1}^T \Delta_t &= \sum_{t=1}^T (\|M_t - U\|_F^2 - \|M_{t+1} - U\|_F^2) \\ &= \|M_1 - U\|_F^2 - \|M_{T+1} - U\|_F^2 \\ &\leq \|M_1 - U\|_F^2. \end{aligned} \quad (7)$$

For simplicity, we employ stochastic gradient descent (SGD) to update the parameter matrix  $M_t$ . Hence, according to the definition of SGD,  $M_{t+1} = M_t - \eta \nabla \ell(M_t)$ , where  $\eta$  is the learning rate, and  $\nabla \ell(M_t) = \gamma A_{t+1}$ . Then, we can rewrite the  $\Delta_t$  as,

$$\begin{aligned} \Delta_t &= \|M_t - U\|_F^2 - \|M_{t+1} - U\|_F^2 \\ &= \|M_t - U\|_F^2 - \|M_t - \eta \nabla \ell(M_t) - U\|_F^2 \\ &= \|M_t\|_F^2 - 2\langle M_t, U \rangle_F + \|U\|_F^2 - \|M_t - U\|_F^2 \\ &\quad + 2\langle M_t - U, \eta \nabla \ell(M_t) \rangle_F - \eta^2 \|\nabla \ell(M_t)\|_F^2 \\ &= 2\eta \text{VEC}(M_t - U)^\top \text{VEC}(\nabla \ell(M_t)) - \eta^2 \|\nabla \ell(M_t)\|_F^2 \\ &\quad \left( \text{employ the Eq. (6) i.e., } \ell(U) \geq \ell(M_t) + \text{VEC}(\nabla \ell(M_t))^\top \text{VEC}(U - M_t) \right) \\ &\geq 2\eta(\ell_t - \ell_t^*) - \eta^2 \|\nabla \ell(M_t)\|_F^2. \end{aligned} \quad (8)$$

We can easily get that,

$$\sum_{t=1}^T \left[ 2\eta(\ell_t - \ell_t^*) - \eta^2 \|\nabla \ell(\mathbf{M}_t)\|_F^2 \right] \leq \|\mathbf{M}_1 - \mathbf{U}\|_F^2. \quad (9)$$

As all samples are  $\ell_2$  normalized, the 2-norm of each sample is 1, namely  $\|\mathbf{x}_t\|_2 \equiv 1, t \in \{1, 2, \dots, T\}$ . We can easily calculate the Frobenius norm of  $\mathbf{A}_{t+1}$ .

$$\begin{aligned} \|\mathbf{A}_{t+1}\|_F &\leq \|(\mathbf{x}_{t+1} - \mathbf{x}_p)(\mathbf{x}_{t+1} - \mathbf{x}_p)^\top\|_F + \|(\mathbf{x}_{t+1} - \mathbf{x}_q)(\mathbf{x}_{t+1} - \mathbf{x}_q)^\top\|_F \\ &\quad \left( \text{employ } \|\mathbf{ab}^\top\|_F^2 = \left( \sum_{i=1}^d |\mathbf{a}_i|^2 \right) \left( \sum_{j=1}^d |\mathbf{b}_j|^2 \right), \text{ where } \mathbf{a}, \mathbf{b} \in \mathbb{R}^d \right) \\ &= \|\mathbf{x}_{t+1} - \mathbf{x}_p\|_2 \cdot \|\mathbf{x}_{t+1} - \mathbf{x}_p\|_2 + \|\mathbf{x}_{t+1} - \mathbf{x}_q\|_2 \cdot \|\mathbf{x}_{t+1} - \mathbf{x}_q\|_2 \\ &= \|\mathbf{x}_{t+1} - \mathbf{x}_p\|_2^2 + \|\mathbf{x}_{t+1} - \mathbf{x}_q\|_2^2 \\ &\quad \left( \text{for } \|\mathbf{a} - \mathbf{b}\|_2^2 \leq (\|\mathbf{a}\|_2 + \|\mathbf{b}\|_2)^2 \right) \\ &\leq 8. \end{aligned} \quad (10)$$

Thus,

$$\begin{aligned} \sum_{t=1}^T (\ell_t - \ell_t^*) &\leq \frac{1}{2\eta} \|\mathbf{M}_1 - \mathbf{U}\|_F^2 + \frac{\eta}{2} \sum_{t=1}^T \|\nabla \ell(\mathbf{M}_t)\|_F^2 \\ &= \frac{1}{2\eta} \|\mathbf{M}_1 - \mathbf{U}\|_F^2 + \frac{\eta}{2} \sum_{t=1}^T \|\gamma \mathbf{A}_{t+1}\|_F^2 \\ &\leq \frac{1}{2\eta} \|\mathbf{M}_1 - \mathbf{U}\|_F^2 + 32T\eta\gamma^2 \\ &\quad (\mathbf{M}_1 \text{ is initialized to an identity matrix } \mathbf{I}) \\ &= \frac{1}{2\eta} \|\mathbf{I} - \mathbf{U}\|_F^2 + 32T\eta\gamma^2. \end{aligned} \quad (11)$$

In particular, setting  $\eta = \frac{1}{R\sqrt{T}}$  (where  $R > 0$  is a constant) yields the regret bound  $R(\mathbf{U}, T) \leq \left( \frac{R}{2} \|\mathbf{I} - \mathbf{U}\|_F^2 + \frac{32\gamma^2}{R} \right) \sqrt{T}$ . In fact, in this study, as a closed-form solution is employed (i.e.,  $\eta = 1$ ), the regret bound is  $R(\mathbf{U}, T) \leq \frac{1}{2} \|\mathbf{I} - \mathbf{U}\|_F^2 + 32T\gamma^2$ . By setting  $\gamma$  in a decreasing way with the iteration number  $T$ , for example,  $\gamma = \frac{1}{R\sqrt{T}}$ , we can obtain a regret bound  $R(\mathbf{U}, T) \leq \frac{1}{2} \|\mathbf{I} - \mathbf{U}\|_F^2 + \frac{32}{R^2}$ . Hence proved.  $\square$

### 3 Theoretical analysis of Proposition 1

**Proposition 1.** *Let  $\mathbf{M}_1, \dots, \mathbf{M}_n$  be the parameter matrixes learnt by each metric layer of ODML. The subsequent metric layer can learn a feature space that is at least as good as the one learnt by the former metric layer. That is, the composite feature space learnt by both  $\mathbf{M}_1$  and  $\mathbf{M}_2$  is better than the feature space learnt only by  $\mathbf{M}_1$  in most cases (i.e., the feature space is more discriminative for classification).*

*Proof.* For simplicity, we just consider to analyze and prove this proposition of ODML-FP that only uses forward propagation strategy. In fact, as ODML-FP only has forward propagation, each metric layer is a relatively independent MOML algorithm. Thus, Theorem 2 is applicable to each metric layer. In other words, each metric layer (i.e., a MOML algorithm) has its own tight regret bound. As the subsequent metric layer is learnt based on the output of the former metric

layer, the metric space should not be worse according to the theoretical guarantee of regret bound. Moreover, ReLU activation function can introduce nonlinear and sparsity into the feature mapping, which is also beneficial to the exploration of feature space. In some cases, if the latter metric layer is in the wrong direction, backward propagation can be chosen to correct and adjust the direction to some extent.  $\square$

### References

- [Crammer *et al.*, 2006] Koby Crammer, Ofer Dekel, Joseph Keshet, Shai Shalev-Shwartz, and Yoram Singer. Online passive-aggressive algorithms. *Journal of Machine Learning Research*, 7:551–585, 2006.
- [Petersen *et al.*, 2008] Kaare Brandt Petersen, Michael Syskind Pedersen, et al. The matrix cookbook. *Technical University of Denmark*, 7:15, 2008.