

SUPPLEMENTARY MATERIAL FOR "GENERALIZABLE EMBEDDINGS WITH CROSS-BATCH METRIC LEARNING"

Yeti Z. Gürbüz[†] A. Aydın Alatan[‡]

[†]RSiM, Technische Universität Berlin, DE [‡]Center for Image Analysis (OGAM), METU, TR

Appendix

Preliminaries

Definition 1 (Optimal Transport Distance) *The optimal transport (OT) distance between two probability mass distributions (p, X) and (q, Y) is:*

$$\|(p, X) - (q, Y)\|_{OT} = \min_{\substack{\pi \geq 0 \\ \sum_i \pi_{ij} = q_j \\ \sum_j \pi_{ij} = p_i}} \sum_{ij} c_{ij} \pi_{ij} \quad (\text{A.1})$$

where $c_{ij} = \|x_i - y_j\|_2$, and $(p, X) \in \Sigma_n \times \mathbb{R}^{d \times n}$ denotes a probability mass distribution with masses $p \in \Sigma_n$ in the probability simplex (i.e., $\Sigma_n := \{p \in \mathbb{R}_{\geq 0}^n \mid \sum_i p_i = 1\}$), and d -dimensional support $X = [x_i]_{i \in [n]} \in \mathbb{R}^{d \times n}$.

Definition 2 (Maximum Mean Discrepancy) *Maximum mean discrepancy (MMD) between two probability mass distributions (p, X) and (q, Y) is:*

$$\|(p, X) - (q, Y)\|_{MMD} = \max_{f \in \mathcal{C}(X, Y)} \sum_i p_i f(x_i) - \sum_j q_j f(y_j) \quad (\text{A.2})$$

where $\mathcal{C}(X, Y)$ is the set of continuous and bounded functions defined on a set covering the column vectors of X and Y .

Definition 3 (Optimal Transport Distance Dual)

The Lagrangian dual of the optimal transport distance defined in Definition 1 reads:

$$\|(p, X) - (q, Y)\|_{OT} = \max_{f_i + g_j \leq c_{ij}} \sum_i p_i f_i + \sum_j q_j g_j \quad (\text{A.3})$$

with the dual variables $\lambda = \{f, g\}$.

Note that $x_i = y_j$ implies $f_i = -g_j$ and from the fact that $c_{ij} = c_{ji}$, we can express the problem in (A.3) as:

$$\|(p, X) - (q, Y)\|_{OT} = \max_{f \in \mathfrak{L}_1} \sum_i p_i f(x_i) - \sum_j q_j f(x_j) \quad (\text{A.4})$$

where $\mathfrak{L}_1 = \{f \mid \sup_{x, y} \frac{|f(x) - f(y)|}{\|x - y\|_2} \leq 1\}$ is the set of 1-Lipschitz functions.

[†]Affiliated with OGAM-METU during the research.

Proofs

Definition 4 (Histogram Operator) *For n -many d -dimensional features $X = [x_i \in \mathbb{R}^d]_{i=1}^n$ and m -many prototype features $\mathcal{V} = [\nu_i \in \mathbb{R}^d]_{i=1}^m$ of the same dimension, the histogram of X on \mathcal{V} is denoted as z^* which is computed as the minimizer of the following problem:*

$$(z^*, \pi^*) = \arg \max_{z \in \mathcal{S}^m, \pi \geq 0} \sum_{ij} \nu_i^\top x_j \pi_{ij} \text{ s.t. } \sum_i \pi_{ij} = 1/n \quad (\text{A.5})$$

where $\mathcal{S}^m := \{p \in \mathbb{R}_{\geq 0}^m \mid \sum_i p_i = 1\}$.

Claim 1 *The solution of the problem in (A.5) reads:*

$$\pi_{ij}^* = 1/n \mathbf{1}(i = \arg \max_k \{\nu_k^\top x_j\}) \quad (\text{A.6})$$

where $\mathbf{1}(c)$ is 1 whenever c is true and 0 otherwise.

Proof: We prove our claim by contradiction. Denoting $c_{ij} = -\nu_i^\top x_j$, for any j , we express a solution as $\pi_{ij}^* = \epsilon_i$ with $\epsilon_i \geq 0$ and $\sum_i \epsilon_i = 1/n$. Let $i^* = \arg \min_k \{c_{kj}\}$. We can write $\pi_{i^*j}^* = 1/n - \sum_{i \mid i \neq i^*} \epsilon_i$. Our claim states that $\epsilon_i = 0$ for $i \neq i^*$. We assume an optimal solution, π' , with $\epsilon_i > 0$ for some $i \neq i^*$. Since π' is optimal, we must have $\sum_{ij} \pi'_{ij} c_{ij} \leq \sum_{ij} \pi_{ij} c_{ij}$ for any π . For the j^{th} column we have,

$$\begin{aligned} \sum_i \pi'_{ij} c_{ij} &= \left(\frac{1}{n} - \sum_{i' \mid i' \neq i^*} \epsilon_{i'}\right) c_{i^*j} + \sum_{i' \mid i' \neq i^*} \epsilon_{i'} c_{i'j} \\ &= \frac{1}{n} c_{i^*j} + \sum_{i' \mid i' \neq i^*} \epsilon_{i'} (c_{i'j} - c_{i^*j}) \stackrel{(a)}{>} \sum_i \pi_{ij}^* c_{ij} \end{aligned}$$

where in (a) we use the fact that $(c_{i'j} - c_{i^*j}) > 0$ and $\epsilon_{i'} > 0$ for some i' by the assumption. Hence, $\sum_{ij} \pi'_{ij} c_{ij} > \sum_{ij} \pi_{ij}^* c_{ij}$ poses a contradiction. Therefore, $\epsilon_{i'} = 0$ must hold for all $i' \neq i^*$. ■

Lemma 1 *Given n -many convolutional features $X = [x_i \in \mathcal{X}]_{i=1}^n$, and m -many prototype features $\mathcal{V} = [\nu_i]_{i=1}^m$ with $\{\nu_i\}_{i=1}^m$ being δ -cover of \mathcal{X} . If z^* is the histogram of X on \mathcal{V} , defined in (A.5), then we have:*

$$\left\| \sum_{i=1}^m z_i^* \nu_i - \sum_{j=1}^n \frac{1}{n} x_j \right\|_2 \leq \delta$$

Proof: We can express

$$\left\| \sum_{i \in [m]} z_i^* \nu_i - \sum_{j \in [n]} \frac{1}{n} x_j \right\|_2^2 = \sum_{i \in [m]} p_i^* f(\nu_i) - \sum_{j \in [n]} q_j f(x_j)$$

where $f(x) = x^\top (\sum_i z_i^* \nu_i - \sum_j \frac{1}{n} x_j)$, and $[n] = 1, \dots, n$. Note that f is a continuous bounded operator for $\mathcal{X} = \{x \mid \|x\|_2 \leq 1\}$ (We can always map the features inside unit sphere without losing the relative distances). Moreover, the operator norm of f , *i.e.* $\|f\|$, which is $\|\sum_i z_i^* \nu_i - \sum_j \frac{1}{n} x_j\|_2$ is less than or equal to 1. Thus, f lie in the unit sphere of the continuous bounded functions set. Using the definition of MMD distance, we can bound the error as:

$$\sum_{i \in [m]} z_i^* f(\nu_i) - \sum_{j \in [n]} q_j f(x_j) \leq \|(z^*, V) - (q, X)\|_{MMD}$$

where $q_i = 1/n$ for all i . For the continuous and bounded functions of the operator norm less than 1, MMD is lower bound for OT [1]. Namely,

$$\begin{aligned} \sum_{i \in [m]} z_i^* f(\nu_i) - \sum_{j \in [n]} q_j f(x_j) &\leq \|(z^*, V) - (q, X)\|_{MMD} \\ &\leq \|(z^*, V) - (q, X)\|_{OT} \end{aligned}$$

Since columns of V is δ -cover of the set \mathcal{X} , the optimal transport distance between the two distributions are bounded by δ , *i.e.* $\|(z^*, V) - (q, X)\|_{OT} \leq \delta$. Thus, we finally have:

$$\left\| \sum_{i \in [m]} z_i^* \nu_i - \sum_{j \in [n]} \frac{1}{n} x_j \right\|_2 \leq \delta.$$

■

1. REFERENCES

- [1] Bharath K Sriperumbudur, Arthur Gretton, Kenji Fukumizu, Bernhard Schölkopf, and Gert RG Lanckriet, “Hilbert space embeddings and metrics on probability measures,” *The Journal of Machine Learning Research*, vol. 11, pp. 1517–1561, 2010.