Noisy-Correspondence Learning for Text-to-Image Person Re-identification

Yang Qin, Yingke Chen, Dezhong Peng, Xi Peng, Joey Tianyi Zhou, Peng Hu* College of Computer Science, Sichuan University

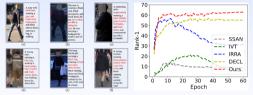


VALSE 2024 重庆

初觉与学习青年学者研讨会 VISION AND LEARNING SEMINAR

Github: https://github.com/OinYang79/RDE (CVPR 2024)

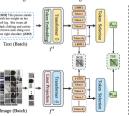
Observation & Motivation



"Noisy correspondences would cause model degradation."

Overall Framework

The overview of our RDE, (a) is the illustration of the cross-modal embedding model used in RDE, which consists of basical global embedding (BGE) and token selection embedding (TSE) modules with different granularity. By integrating them, RDE can capture coarsegrained cross-modal interactions while selecting informative local token features to encode more fine-grained representations for a more accurate similarity. (b) shows the core of RDE to achieve robust similarity learning, which consists of Confident Consensus Division (CCD) and Triplet Alignment Loss (TAL), CCD performs consensus division to obtain confident clean training data, thus avoiding misleading from noisy pairs. Unlike traditional Triplet Ranking Loss (TRL), TAL exploits an upper bound to consider all negative pairs, thus embracing more stable learning.



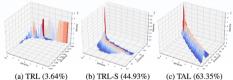
Triplet Alignment Loss (b) Robust Similarity Learning

Confident Consensus Division

BGE: per-sample loss & TSE: per-sample loss $\ell(\mathcal{M}, \mathcal{P}) = \{\ell_i\}_{i=1}^N = \{\mathcal{L}(I_i, T_i)\}_{i=1}^N \Longrightarrow \mathcal{G}_{MM} \Longrightarrow p(k=0|l_i)$

Division and Recalibration

Triplet Alignment Loss



 $\mathcal{L}_{tal}(I_i, T_i) = \left[m - S_{i2t}^+(I_i) + \tau \log(\sum q_{ij} \exp(S(I_i, T_j)/\tau))\right]_+$

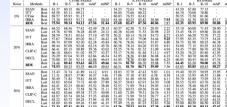
+
$$\left[m - S_{t2i}^{+}(T_i) + \tau \log(\sum_{i=1}^{K} q_{ji} \exp(S(I_j, T_i)/\tau))\right]_{+}$$

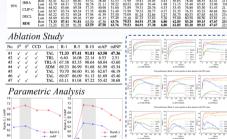
Lemma 1 TAL is the upper bound of TRL, i.e.,

 $\mathcal{L}_{trl}(I_i, T_i) = \left[m - S_{i2t}^+(I_i) + S(I_i, \hat{T}_i)\right]_+ \quad \triangleright More \ Robust$ + $\left[m - S_{t2i}^+(T_i) + S(\hat{I}_i, T_i)\right]_+ \le \mathcal{L}_{tal}(I_i, T_i),$ >No collapse

where $\hat{T}_i \in \{T_j | l_{ij} = 0, \forall j \in \{1, \cdots, K\}\}$ is the hardest negative text for I_i and $\hat{I}_i \in \{I_i | l_{ii} = 0, \forall j \in \{1, \dots, K\}\}$ is the hardest negative image for I:, respectively.

Experiments







(a) Cross-modal Embedding Model