

对比损失

clip

soft max 得到 T22 年

$$\mathcal{L}_{t2i}(i) = -\log \frac{\exp(s(V_i, T_i))}{\sum_{a=1}^B \exp(s(V_a, T_i))}$$

文本对图像的对比损失

即一个确定的文本对所有图像对比损失

对 clip 来说, 找到 (图像-文本) 对 (positive, 匹配的)

得  $\exp(s(V_i, T_i))$

此文本对所有图像求相似度得  $\sum_{a=1}^B \exp(s(V_a, T_i))$

(t2i1b)

$$= \log \sum_{a=1}^B \exp(s(V_a, T_i)) - s(V_i, T_i)$$

越大越好      越小越好

clip-reid stage 1.

$$\mathcal{L}_{t2i}(y_i) = \frac{-1}{|P(y_i)|} \sum_{p \in P(y_i)} \log \frac{\exp(s(V_p, T_{y_i}))}{\sum_{a=1}^B \exp(s(V_a, T_{y_i}))} \quad (4)$$

理解 一个文本可对应多张图片 (多个 positive 对)

则对多个 positive 对求平均 eg  $(V_1, T_1)$

$(V_2, T_1)$

$(V_{10}, T_1)$

图片 1, 2, 10 都跟 T<sub>1</sub> 对应

stage 2.

$$\mathcal{L}_{i2tce}(i) = \sum_{k=1}^N -q_k \log \frac{\exp(s(V_i, T_{y_k}))}{\sum_{y_a=1}^N \exp(s(V_i, T_{y_a}))}$$

图像对文本的交叉熵

图像对每个文本的得分

记为

$$q_1 \log \frac{\exp(s(v, T_{y_1}))}{\exp(s(v, T))} + q_2 \frac{\exp(s(v, T_{y_2}))}{\exp(s(v, T))}$$

↓  
[1 0 0 0 ...]

标准平滑  
第一张图片

↓  
[0 1 0 0 0 ...]

第二张图片

其中  $q_k$  代表的图像与  $T_{y_k}$  代表的文本是配对的

$v_i$

使  $v_i = T_{y_i}$  时最大

$$q_1 [1, 0, 0, 0, \dots]$$

$$v_i = \frac{\exp(s(v_i, T_{y_i}))}{\exp(s(v_i, T))} \rightarrow [0.9, 0.01, 0.01, \dots]$$

直观解释:

对一张图片  $v_i$  把它与每个类别文本特征的相似度做的 softmax 得到预测概率分布, 再用  $q$  (真实平滑标准分布) 作交叉熵。