

DHN-NCE Loss

We are given the loss $\mathcal{L} = \mathcal{L}^{v \rightarrow t} + \mathcal{L}^{t \rightarrow v}$, with hyperparameters γ, β_1, β_2 ,

$$\mathcal{L}^{v \rightarrow t} = - \sum_{i=1}^B \frac{\mathbf{I}_{p_i}^T \mathbf{T}_{p_i}}{\gamma} + \sum_{i=1}^B \log \left(\sum_{j \neq i} \exp \left(\frac{\mathbf{I}_{p_i}^T \mathbf{T}_{p_j}}{\gamma} \right) W_{\mathbf{I}_{p_i} \mathbf{T}_{p_j}}^{v \rightarrow t} \right)$$

$$W_{\mathbf{I}_{p_i} \mathbf{T}_{p_j}}^{v \rightarrow t} = (\beta_1 - 1) \cdot \frac{\exp \left(\frac{\beta_1 \mathbf{I}_{p_i}^T \mathbf{T}_{p_j}}{\gamma} \right)}{\sum_{k \neq i} \exp \left(\frac{\beta_1 \mathbf{I}_{p_i}^T \mathbf{T}_{p_k}}{\gamma} \right)}$$

$$\mathcal{L}^{t \rightarrow v} = - \sum_{i=1}^B \frac{\mathbf{T}_{p_i}^T \mathbf{I}_{p_i}}{\gamma} + \sum_{i=1}^B \log \left(\sum_{j \neq i} \exp \left(\frac{\mathbf{T}_{p_i}^T \mathbf{I}_{p_j}}{\gamma} \right) W_{\mathbf{T}_{p_i} \mathbf{I}_{p_j}}^{t \rightarrow v} \right)$$

$$W_{\mathbf{T}_{p_i} \mathbf{I}_{p_j}}^{t \rightarrow v} = (\beta_2 - 1) \cdot \frac{\exp \left(\frac{\beta_2 \mathbf{T}_{p_i}^T \mathbf{I}_{p_j}}{\gamma} \right)}{\sum_{k \neq i} \exp \left(\frac{\beta_2 \mathbf{T}_{p_i}^T \mathbf{I}_{p_k}}{\gamma} \right)}$$

We aim to vectorise these operations which leverage a GPU.

We start with matrices $\mathbf{I}, \mathbf{T} \in \mathbb{R}^{B \times n}$ which are the matrices of text and image embeddings calculated from the model. These represent 1 batch. We aim to find $f(\mathbf{I}, \mathbf{T}, \gamma, \beta_1, \beta_2)$ which is the loss of matrices \mathbf{I}, \mathbf{V} .

We will leverage the summations that skip the leading diagonal.

First $\mathbf{I}, \mathbf{T} \in \mathbb{R}^{B \times n}$.

Define similarity matrix $S = \mathbf{I}^T \mathbf{T}$.

$$\text{and } A = \exp \left(\frac{1}{\gamma} S \right)$$

$$\text{So } \mathcal{L}^{v \rightarrow t} = - \sum_{i=1}^B \frac{\mathbf{I}_{p_i}^T \mathbf{T}_{p_i}}{\gamma} + \sum_{i=1}^B \log \left(\sum_{j \neq i} \exp \left(\frac{\mathbf{I}_{p_i}^T \mathbf{T}_{p_j}}{\gamma} \right) W_{\mathbf{I}_{p_i} \mathbf{T}_{p_j}}^{v \rightarrow t} \right)$$

$$= \sum_{i=1}^B \log \sum_{j=1}^B \exp \left(\frac{\mathbf{I}_{p_i}^T \mathbf{T}_{p_j}}{\gamma} \right) W_{\mathbf{I}_{p_i} \mathbf{T}_{p_j}}^{v \rightarrow t}$$

As mentioned, the summation ignores the case when $j=i$, which is the leading diagonal of matrix A . So we can add the A_{ii} leading diagonal in, but set its weight to -1 to represent the coefficient.

$$\text{So } \mathcal{L}^{v \rightarrow t} = \sum_{i=1}^B \log \sum_{j=1}^B A_{ij} W_{ij}^{v \rightarrow t}$$

$$\text{Then we have } W_{ij}^{v \rightarrow t} = (B-1) \frac{\exp\left(\frac{B, I_{A_i}^T T_{Pi}}{\tau}\right)}{\sum_{k \neq i} \exp\left(\frac{B, I_{A_k}^T T_{Pi}}{\tau}\right)} \in \mathbb{R}^{B \times B}$$

is equivalent to

The fraction would be a softmax function with the leading diagonal equal to 0. And if $i=j$, we will replace it.

And $W_{ii}^{v \rightarrow t} = -1 \quad \forall i \in \{1, \dots, B\}$ anyway. So calculating $W_{\{p_i\} T_{Pi}}^{v \rightarrow t}$, we set ~~diag(W) = 0~~ the diagonal to 0, apply softmax and then ~~multiply by 1~~ subtract 1 from the diagonal.

$$W_{ij}^{v \rightarrow t} = (B-1) \cdot \text{softmax}\left(A_{ij}^{B_i} - I_B A_{ij}^{B_i}, \text{dim}=-1\right) - I_B \in \mathbb{R}^{B \times B}$$

$A_{ij}^{B_i} - I_B A_{ij}^{B_i}$ applies the B_i scaling element-wise and then replaces the leading diagonal with zeros.

$$\text{So } \mathcal{L}^{v \rightarrow t} = \sum_{i=1}^B \log \sum_{j=1}^B A_{ij} \otimes \left\{ (B-1) \text{softmax}\left(A_{ij}^{B_i} - I_B A_{ij}^{B_i}, \text{dim}=-1\right) - I_B \right\}.$$

We use the Hadamard product since we apply a linear combination to each image, resulting in B total computations. A matrix multiplication will require B^2 computations, which is inefficient.

Since there is a log term between the sums, we cannot turn this into a matrix multiplication as it is non-linear.

As ~~log is one-to-one and strictly increasing~~, removing the log
~~so $\mathcal{L}^{v \rightarrow t} = \text{sum}(\log(\text{sum}(A \otimes (B-1) \text{softmax}(A^{B_i} - I_B A)))$~~

$$L^{v \rightarrow t} = \text{sum}(\log(\text{sum}(A \otimes \{\text{softmax}(A^{\beta_1} - I_B A^{\beta_1}, dm=-1) - I_0\})))$$

By inspection if we say $L^{v \rightarrow t} = g(A)$ (and $A = h(I, T)$),
 $L^{t \rightarrow v} = g(A^T)$. and vice versa.

$$L = \text{sum}(\log(\text{sum}[A \otimes \{\text{softmax}(A^{\beta_1} - I_B A^{\beta_1}, dm=-1) - I_0\}])) \\ + \text{sum}(\text{""} - A^T \text{""} - \{\text{(}(A^{\beta_1})^T \text{)} - I_B A^{\beta_2}, dm=-1\} - I_0\})$$

where $A = \exp\left(\frac{I^T T}{\Sigma}\right)$ and A^{β_1}, A^{β_2} are applied element-wise.

For reference, $A \otimes W \in \mathbb{R}^{B \times B}$

$\text{sum}(A \otimes W) \in \mathbb{R}^{B \times 1}$ as it is applied on 1 axis

$\log(\text{sum}(A \otimes W)) \in \mathbb{R}^{B \times 1}$

$\text{sum}(\log(\text{sum}(A \otimes W))) \in \mathbb{R}^{1 \times 1}$