

---

# Deep Joint-Semantics Reconstructing Hashing for Large-Scale Unsupervised Cross-Modal Retrieval

## Introduction

### Joint-Semantics Reconstructing Hashing

变量定义:

- $m$ : 批的大小;
- $\mathcal{O}$ : 批, 即  $\{o_k = [I_k, T_k]\}_{k=1}^m$ , 每个元素为图片-文本对。图片和文本的特征矩阵分别为  $F_I \in \mathbb{R}^{m \times p_I}$  和  $F_T \in \mathbb{R}^{m \times p_T}$ ;
- $B_I \in \{\pm 1\}^{m \times d}$  和  $B_T \in \{\pm 1\}^{m \times d}$ : 表示  $I_k$  和  $T_k$  分别通过 ImgNet 和 TxtNet 编码后生成的二值向量;
- $\hat{F}_I$  和  $\hat{F}_T$ : 表示归一化后的  $F_I$  和  $F_T$ , 则余弦相似度为

$$S_I = \hat{F}_I \hat{F}_I^\top \in [-1, +1]^{m \times m}$$

$$S_T = \hat{F}_T \hat{F}_T^\top \in [-1, +1]^{m \times m}$$

$B$  可以看作单位超立方体的顶点, 相邻的点具有相似的哈希值。其 Hamming 距离可以用角距离来代替, 即两个向量夹角的  $\cos$  值。

### Constructing Joint-Semantics Matrix

考虑使用上面的  $S_I$  和  $S_T$  指导哈希码的学习。一种传统的方法是

$$\min_B \beta \text{Tr}(B^\top L_I B) + (1 - \beta) \text{Tr}(B^\top L_T B) \quad \text{s.t. } B \in \{\pm 1\}^{m \times d}$$

其中

$$L_I = \text{diag}(S_I \mathbf{1}) - S_I$$

$$L_T = \text{diag}(S_T \mathbf{1}) - S_T$$

本文构造 Joint-semantics Affinity Matrix, 定义组合函数  $\mathcal{C}$ ,

$$S = \mathcal{C}(S_I, S_T) \in [-1, +1]^{m \times m}$$

其中  $S_{ij}$  表示实例  $o_i$  和  $o_j$  的潜在语义相度。令

$$\tilde{S} = \beta S_I + (1 - \beta) S_T$$

$$\begin{aligned} S &= \mathcal{C}(S_I, S_T) \\ &= (1 - \eta) \tilde{S} + \eta \frac{\tilde{S} \tilde{S}^\top}{m} \\ &= (1 - \eta) [\beta S_I + (1 - \beta) S_T] + \frac{\eta}{m} [\beta^2 S_I S_I^\top + \beta(1 - \beta)(S_I S_T^\top + S_T S_I^\top) + (1 - \beta^2) S_T S_T^\top] \end{aligned}$$

---

上述式子可以看成 diffusion process 的一次迭代。

### **Reconstructing with Binary Codes**

目标是最小化

$$\min_{B_I, B_T} \|\mu S - \cos(B_I, B_T)\|_F^2, \quad \text{s.t. } S = \mathcal{C}(S_I, S_T) \in [-1, +1]^{m \times m}$$

这里  $\mu$  是超参数，可以用来调节归约到  $\pm 1$  上的范围程度。相比拉普拉斯约束

$$\text{Tr}(B^\top L B) = \sum_{i,j} S_{ij} \|B_i - B_j\|^2$$

相当于从中间的乘号变成减号，减少了一个 batch 里的过拟合现象，更适合分批学习。进一步加入相同模态内的约束条件，最后的目标函数为

$$\begin{aligned} \min_{B_I, B_T} & \|\mu S - \cos(B_I, B_T)\|_F^2 + \lambda_1 \|\mu S - \cos(B_I, B_I)\|_F^2 + \lambda_2 \|\mu S - \cos(B_T, B_T)\|_F^2, \\ \text{s.t. } & S = \mathcal{C}(S_I, S_T) \in [-1, +1]^{m \times m}, B_I, B_T \in \{-1, +1\}^{m \times d} \end{aligned}$$

### **Optimization**

用  $H \in \mathbb{R}^{m \times d}$  统一表示 ImgNet 和 TxtNet 最后一层的输出（没有激活函数），则

$$B = \text{sgn}(H) \in \{-1, +1\}^{m \times d}$$

由于符号函数的梯度大部分为 0，故改为使用如下函数

$$B = \tanh(\alpha H) \in \{-1, +1\}^{m \times d}, \quad \alpha \in \mathbb{R}^+$$

整个算法的流程图如下所示

---

**Algorithm 1** Deep Joint-Semantics Reconstructing Hashing

---

**Input:**

Training set  $\{\mathbf{o}_k = [\mathbf{I}_k, \mathbf{T}_k]\}_{k=1}^n$  and their corresponding original features  $\mathbf{F}_I$  and  $\mathbf{F}_T$ ; ImgNet  $\mathcal{G}_{\theta_I}$  and TxtNet  $\mathcal{G}_{\theta_T}$  with  $\theta_I$  and  $\theta_T$  denoting the deep network parameters; batch size  $m$ ;

**Output:**

Hashing coding function  $\varphi_I(x) = \text{sgn}(\mathcal{G}_{\theta_I}(x))$  for image input and  $\varphi_T(x) = \text{sgn}(\mathcal{G}_{\theta_T}(x))$  for text input;

- 1: Initialize epoch  $t = 0$ ;
  - 2: **repeat**
  - 3:    $t = t + 1$ ;  $\alpha = \sqrt{t}$ ;
  - 4:   **for**  $\lfloor \frac{n}{m} \rfloor$  iterations **do**
  - 5:     Randomly sample a batch of instances from training set  $\{\mathbf{o}_k = [\mathbf{I}_k, \mathbf{T}_k]\}_{k=1}^m$ ;
  - 6:     Calculate the normalized  $\hat{\mathbf{F}}_I, \hat{\mathbf{F}}_T$  and integrate the cosine matrices  $\mathbf{S}_I = \hat{\mathbf{F}}_I \hat{\mathbf{F}}_I^\top$ ,  $\mathbf{S}_T = \hat{\mathbf{F}}_T \hat{\mathbf{F}}_T^\top$  to the joint-semantics affinity  $\mathbf{S}$  with Equation (3);
  - 7:     Forward propagate  $\mathbf{H}_I = \mathcal{G}_{\theta_I}(\mathbf{I})$ ,  $\mathbf{H}_T = \mathcal{G}_{\theta_T}(\mathbf{T})$ ;
  - 8:     Hash coding with activation function (7)  $\mathbf{B}_I = \tanh(\alpha \mathbf{H}_I)$ ,  $\mathbf{B}_T = \tanh(\alpha \mathbf{H}_T)$ ;
  - 9:     Calculate the objective function (5), back propagate the gradients with the chain rule and update the whole parameters;
  - 10:   **end for**
  - 11: **until** convergence
- 

**Figure 1:** 1.PNG