

特征表示: 一维向量
snippet-level: $x \in \mathbb{R}^D$ 通过 I3D 提取特征
proposal-level: $y \in \mathbb{R}^D$ 通过 BSN 获得 proposal, 对每个 proposal 的 $[start, end]$ 中全部 snippets 进行 max-pooling
video-level: $z \in \mathbb{R}^D$ 将全部的 snippets 进行 max-pooling

$x_j \in \{x_j | j \in S(i)\}$ $S(i)$ i th proposal.

L-Net. 利用余弦相似度, 计算一个 proposal feature 与 ~~proposal~~ ^{snippet} feature 的距离 作为 snippet 的权重 a .

$$a_{ij}^L = \frac{\sigma(s(y_i, x_j))}{\sum_{k \in S(i)} \sigma(s(y_i, x_k))}$$

$s(x, y)$: 余弦相似度
 $\sigma()$: ReLU

根据上下
 文补充 proposal
 细粒度时间
 特征

$$y_i^L = \sigma(W_1^L y_i + \sum_{j \in S(i)} a_{ij}^L W_2^L x_j) \quad \sigma: \text{ReLU}$$

$$W_1^L, W_2^L \in \mathbb{R}^{\frac{D}{2} \times D} \quad y_i^L \in \mathbb{R}^{\frac{D}{2}}$$

G-Net. $a_{ij}^G = \frac{\sigma(s(z, x_j))}{\sum_{k \in S(i)} \sigma(s(z, x_k)) + \sigma(s(z, y_i))}$

$$b_i^G = \frac{\sigma(s(z, y_i))}{\sum_{k \in S(i)} \sigma(s(z, x_k)) + \sigma(s(z, y_i))}$$

a_{ij}^G attention weight between video-level and j th snippet of the i th proposal

b_i^G video-level and i th proposal.

$$z_i^G = \sigma(W_1^G z + W_2^G (\sum_{j \in S(i)} a_{ij}^G x_j + b_i^G y_i))$$

$W_1 \in \mathbb{R}^{\frac{D}{2} \times D}$
 $W_2 \in \mathbb{R}^{\frac{D}{2} \times \frac{D}{2}}$
 $z^G \in \mathbb{R}^{\frac{D}{2}}$

$$y_i^G = y_i^L \oplus z_i^G \quad y_i^G \in \mathbb{R}^D \quad \oplus \text{ concatenation}$$

通过
 融合
 全局
 特征
 来增强
 有 proposal