

$X = (X_t)_{t=1}^T$: RGB or flow features (with a pre-trained model)

[attention-based action localization problem, target is to predict the frame attention. λ .]

MAP: $\max_{\lambda \in [0,1]} \log p(\lambda | X, y)$ $p(\lambda | X, y)$ 给定 X, y unknown probability distribution of λ .
maximum a posteriori.

frame-level labels (ground truth of λ).
 $= \log \frac{P(X, y | \lambda) P(\lambda)}{P(X, y)}$

$$\log p(\lambda | X, y) = \log P(X, y | \lambda) + \log p(\lambda) - \log P(X, y)$$

$$= \log \frac{P(X, y, \lambda)}{P(X, y)} + \log p(\lambda) - \log P(X, y)$$

$$= \log \frac{P(X, y, \lambda)}{P(X, y)}$$

$$= \log \frac{P(X, y, \lambda)}{P(X, y)} \cdot \frac{P(X, \lambda)}{P(X)} + \log p(\lambda) - \log P(X, y)$$

$$= \log P(\lambda | X, y) + \log P(X | \lambda) + \log p(\lambda) - \log P(X, y)$$

$$\propto \log P(y | X, \lambda) + \log P(X | \lambda)$$

prefers λ - 生成模型 从 λ 准确预测帧的特征.

classification.

前景与背景施加不同的关注

$P(X | \lambda)$ 用 generative model 估计 使得模型精准重建特征

discriminative attention and the generative model



扫描全能王 创建

Discriminative Attention Modeling

attention λ as weight:

前导: $x_{fg} = \frac{\sum_{t=1}^T \lambda_t x_t}{\sum_{t=1}^T \lambda_t}$

后导: $x_{bg} = \frac{\sum_{t=1}^T (1-\lambda_t) x_t}{\sum_{t=1}^T (1-\lambda_t)}$ $L_d = L_{fg} + \alpha L_{bg}$

$L_d = L_{fg} + \alpha L_{bg} = -\log P_o(y | x_{fg}) - \alpha \log P_o(o | x_{bg})$

最小化

最大化

Generative Attention Modeling

对用VAE对不同帧特征分布建模. 区分前导与背景

VAE推导: $L = \log P(x)$

$= \sum_z q(z|x) \log P(x)$

$= \sum_z q(z|x) \log \left(\frac{P(z,x)}{P(z|x)} \right)$

$= \sum_z q(z|x) \log \left(\frac{P(z,x)}{q(z|x) P(z|x)} \right)$

$= \sum_z q(z|x) \log \left(\frac{P(z,x)}{q(z|x)} \right) + \sum_z q(z|x) \log \left(\frac{q(z|x)}{P(z|x)} \right)$

$= L^v + D_{KL}(q(z|x) || P(z|x))$

L^v



$L^v = \sum_z q(z|x) \log \frac{P(z,x)}{q(z|x)}$

$= \sum_z q(z|x) \log \frac{P(x|z) P(z)}{q(z|x)}$

$= \sum_z q(z|x) \log P(z) + \sum_z q(z|x) \log P(x|z)$

$= -D_{KL}(q(z|x) || P(z)) + \sum_z q(z|x) \log P(x|z)$



$$P(X|\lambda) = \prod_{t=1}^T P(x_t|\lambda_t)$$

$$P_{\psi}(x_t|\lambda_t) \equiv E_{p_{\psi}(z_t|\lambda_t)} [P_{\psi}(x_t|\lambda_t, z_t)]$$

$$\begin{aligned} L_{CVAE}^{(x)} &= -E_{q_{\phi}(z_t|x_t, \lambda_t)} [\log p_{\psi}(x_t|\lambda_t, z_t)] \\ &\quad + \beta \cdot KL(q_{\phi}(z_t|x_t, \lambda_t) || p_{\psi}(z_t|\lambda_t)) \\ &\approx -\frac{1}{2} \sum_{t=1}^T (\log p_{\psi}(x_t|\lambda_t, z_t)) \end{aligned}$$

$$p_{\psi}(x_t|\lambda_t, z_t) = \mathcal{N}(x_t | f_{\psi}(\lambda_t, z_t), \sigma^2 * I)$$

$$p_{\psi}(z_t|\lambda_t) = \mathcal{N}(z_t | r\lambda_t, I)$$

$$q_{\phi}(z_t|x_t, \lambda_t) = \mathcal{N}(z_t | \mu_{\phi}, \Sigma_{\phi}) \quad \text{outputs of Encoder } f_{\phi}(x_t, \lambda_t)$$

$$L_{re} = -\sum_{t=1}^T \log \{ E_{p_{\psi}(z_t|\lambda_t)} [P_{\psi}(x_t|\lambda_t, z_t)] \} \approx \sum_{t=1}^T \log \left(\frac{1}{2} \sum_{t=1}^T p_{\psi}(x_t|\lambda_t, z_t) \right)$$

$$L=1 \quad L_{re} = -\sum_{t=1}^T \log p_{\psi}(x_t|\lambda_t, z_t) \propto \sum_{t=1}^T \|x_t - f_{\psi}(\lambda_t, z_t)\|^2$$

用之前 Discriminative attention modeling

优化的x. 来更新 CVAE 的分布 $P(X|\lambda)$.

L_{CVAE} 学习参数 λ 更新 λ .

3



扫描全能王 创建

$$\hat{x}_t^H = G(\sigma_s) * \frac{\exp w_g^T x_t}{\sum_{c=0}^C \exp w_c^T x_t} \quad \hat{x}_t^{bg} = G(\sigma_s) * \frac{\sum_{c=1}^C \exp w_c^T x_t}{\sum_{c=0}^C \exp w_c^T x_t}$$

$$L_{guide} = \frac{1}{T} \sum_{t=1}^T |\lambda_t - \hat{x}_t^H| + |\lambda_t - \hat{x}_t^{bg}|$$

update attention and classification modules. with loss

$$L = L_d + \gamma_1 L_{re} + \gamma_2 L_{guide}$$

2-1 where γ_1, γ_2 denote the hyper-parameters

2. update CVAE with loss L_{CVAE}

重建 representation.

↓ $P(X|\lambda)$ 就是具体那帧是动作
更新参数逼近 $P(X|\lambda)$

自底向上方法 ~~维护~~ 从数据产生每一帧注意力 并通过注意力加权。
的验证. 训练分类模型
总得 得尔 识别能力

CVAE 用注意力为条件的来模拟框架特征分布

