

Visual Task Selection Algorithm

August 11, 2021

1 Algorithm

Algorithm 1: Visual Task Selection

Input: $\{w_i\}(i \in ||Q||)$ word embedding of natural language Q_{nl} .
 I an input image.
Require: f : recurrent embedding module.
Require: $Agent$: visual selection system with policy π .
Require: $\mathcal{D} = \{Q_{nl}, I\}$: train set.
Build task pool accordingly;
while *epoch reaches maximum* **do**
 // Perform Agent for all \mathcal{D}
 // Encoder
 $h_i = f_{enc}(w_i, h_{i-1})$; $\triangleright f_{enc}(\cdot)$ is nonlinear function for encoder
 // Decoder with Attention
 $e_{ij} = a(s_{j-1}, h_i)$; $\triangleright e_{ij}$ is the associate energy of probability a_{ij} , s_j is
 an RNN hidden state for time j , $a(\cdot)$ is the alignment model which
 is parametrized as a feedforward neural network
 $a_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{||Q||} \exp(e_{ik})}$; $\triangleright a_{ij}$ is the attention weight
 $c_i = \sum_{j=1}^{||Q||} a_{ij} h_{ij}$; $\triangleright c_i$ is the weighted context vector
 $s_i = f_{dec}(s_{i-1}, a_{i-1}, c_i)$; $\triangleright f_{dec}(\cdot)$ is nonlinear function for decoder
 // Output Action Sequence
 $p(a_i | \{a_1, \dots, a_{i-1}, c_i, Q_{nl}\}) = g(a_{i-1}, s_i, c_i, Q_{nl})$; $\triangleright g(\cdot)$ is a nonlinear
 function
 $A = [a_1, \dots, a_i, \dots, a_t]$;
 // Update Parameters
 Update parameters by Loss $\mathcal{L} = \mathcal{L}_{policy} + \mathcal{L}_{\tau\alpha}$; \triangleright Update parameters
 in $f_{enc}(\cdot), a(\cdot), f_{dec}(\cdot), g(\cdot)$
end
