
A UNROLLING CLASSIFIER

In the following, we use $\text{FC}(v)$ to denote a linear (fully-connected) transformation on vector v . The parameters of different FC's are distinguished by their subscripts.

We took global average pooling followed by a flatten operator on spatial dimension to obtain a 1D vector of y_t (see eq. 17) as \overline{y}_t . The unrolling LSTM is then performed on the \overline{y}_t to unfold the predictions. The hidden and memory states of the unrolling LSTM are initialized by the \overline{y}_t with additional transformations:

$$\xi = \text{ReLU}(\text{FC}_\xi(\overline{y}_t)) \quad (\text{A.1})$$

$$feat = \text{ReLU}(\text{FC}_{feat}(\xi)) \quad (\text{A.2})$$

$$c_0^{url} = \text{FC}_c(feat) \quad (\text{A.3})$$

$$h_0^{url} = \text{FC}_h(feat) \quad (\text{A.4})$$

the states are unrolled T steps, where $T = \tau_a$ (anticipate time) \times fps, with the consistent input \overline{y}_t ,

$$(h_T^{url}, c_T^{url}) = \text{LSTMCell}(\overline{y}_t, (h_{T-1}^{url}, c_{T-1}^{url})) \quad (\text{A.5})$$

We leverage the final unrolling output h_T^{url} for further logits calculation:

$$\text{logit}_a = \text{FC}_a(h_T^{url}) \quad (\text{A.6})$$

$$\text{logit}_v = \text{FC}_v(h_T^{url}) + \text{FC}_{a2v}(\text{logit}_a) \quad (\text{A.7})$$

$$\text{logit}_n = \text{FC}_n(h_T^{url}) + \text{FC}_{a2n}(\text{logit}_a) \quad (\text{A.8})$$

where the logit_a also serves as biases to logit_v and logit_n for explicitly building the relationship of action to the verb and noun.

The cross-entropy loss is used to compare between logits and ground-truth labels. The overall loss function is the summation over the individual loss of each action, verb, and noun, for every anticipation time.