

SEQUENCE METRIC LEARNING AS SYNCHRONIZATION OF RECURRENT NEURAL NETWORKS : SUPPLEMENTARY MATERIAL

Anonymous authors

Paper under double-blind review

1 DIFFERENTIATION OF COUPLED GRU

$\frac{\partial E}{\partial h_t}$ is known, we call dx derivatives of the form $\frac{\partial E}{\partial x}$ with the pair of inputs concatenated (named hereafter left and right inputs). It is sometimes necessary to use separately the left and right sides of the derivative, in which case respectively indicated by \leftarrow and \rightarrow over the indices. Finally, $[x, y, \dots]$ designates a concatenation along the suitable axe. Post-synaptic potentials (before activation function) are represented by the letter a with indices i for input and h for hidden and the associated gate letter. We seek to compute :

- dx_t
- dh_{t-1}
- dW_i where $W_i = [W_{iz}, W_{ir}, W_{in}]$
- dW_h where $W_h = [W_{hz}, W_{hr}, W_{hn}, W_{hc}]$,
- db_i where $b_i = [b_{iz}, b_{ir}, b_{in}]$
- db_h where $b_h = [b_{hz}, b_{hr}, b_{hn}, b_{hc}]$

$$dh_{t-1,o} = dh_t \circ z_t \quad (1)$$

$$dz_t = dh_t \circ (h_{t-1} - \bar{h}_t) \quad (2)$$

$$d\bar{h}_t = dh_t \circ (1 - z_t) \quad (3)$$

$$dc_t = d\bar{h}_{\leftarrow t} \circ (\tilde{h}_{\leftarrow t} - \tilde{h}_{\leftarrow t}) + d\bar{h}_{\leftarrow t} \circ (\tilde{h}_{\leftarrow t} - \tilde{h}_{\leftarrow t}) \quad (4)$$

$$d\tilde{h}_t = d\bar{h}_t + [c_t, c_t] \circ ([d\bar{h}_{\leftarrow t}, d\bar{h}_{\leftarrow t}] - d\bar{h}_t) \quad (5)$$

$$da_z = dz_t \circ ((1 - \sigma(a_z)) \circ \sigma(a_z)) \quad (6)$$

$$da_n = d\tilde{h}_t \circ (1 - \tanh(a_{n,i} + r_t \circ a_{n,h})^2) \quad (7)$$

$$da_{n,h} = da_n \circ r_t \quad (8)$$

$$dr_t = da_{n,h} \circ a_{n,h} \quad (9)$$

$$da_r = dr_t \circ ((1 - \sigma(a_r)) \circ \sigma(a_r)) \quad (10)$$

$$da_c = dc_t \circ ((1 - \sigma(a_{\leftarrow c} + a_{\leftarrow c})) \circ \sigma(a_{\leftarrow c} + a_{\leftarrow c})) \quad (11)$$

$$da_{\leftarrow i} = [da_{\leftarrow z}, da_{\leftarrow r}, da_{\leftarrow n}] \quad (12)$$

$$da_{\leftarrow i} = [da_{\leftarrow z}, da_{\leftarrow r}, da_{\leftarrow n}] \quad (13)$$

$$da_{\leftarrow h} = [da_{\leftarrow z}, da_{\leftarrow r}, da_{\leftarrow n,h}, da_c] \quad (14)$$

$$da_{\leftarrow h} = [da_{\leftarrow z}, da_{\leftarrow r}, da_{\leftarrow n,h}, da_c] \quad (15)$$

$$dW_{iz} = da_{\leftarrow z} x_{\leftarrow t} + da_{\leftarrow z} x_{\leftarrow t} \quad (16)$$

$$dW_{ir} = da_{\leftarrow r} x_{\leftarrow t} + da_{\leftarrow r} x_{\leftarrow t} \quad (17)$$

$$dW_{in} = da_{\leftarrow n} x_{\leftarrow t} + da_{\leftarrow n} x_{\leftarrow t} \quad (18)$$

$$dW_i = [dW_{iz}, dW_{ir}, dW_{in}] \quad \square \quad (19)$$

$$dW_{hz} = da_{\leftarrow z} h_{\leftarrow t} + da_{\leftarrow z} h_{\leftarrow t} \quad (20)$$

$$dW_{hr} = da_{\leftarrow r} h_{\leftarrow t} + da_{\leftarrow r} h_{\leftarrow t} \quad (21)$$

$$dW_{hn} = da_{\leftarrow n,h} h_{\leftarrow t} + da_{\leftarrow n,h} h_{\leftarrow t} \quad (22)$$

$$dW_{hc} = da_c (h_{\leftarrow t} + h_{\leftarrow t}) \quad (23)$$

$$dW_h = [dW_{hz}, dW_{hr}, dW_{hn}, dW_{hc}] \quad \square \quad (24)$$

$$db_i = da_{\leftarrow i} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + da_{\leftarrow i} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \quad \square \quad (25)$$

$$db_h = da_{\leftarrow h} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + da_{\leftarrow h} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \quad \square \quad (26)$$

$$dx_t = [da_{\leftarrow i} W_i, da_{\leftarrow i} W_i] \quad \square \quad (27)$$

$$dh_{t-1} = dh_{t-1,o} + [da_{\leftarrow h} W_h, da_{\leftarrow h} W_h] \quad \square \quad (28)$$

2 STRUCTURAL LOSS

This paper makes use of the loss proposed by Yang et al. Yang et al. (2018) called a *structural loss*. It is a hardness-aware structural loss, composed of two terms a local one and a global one. The local one integrates a system of pair weighting to emphasize the learning on hard-positive samples (see Eq. 31 and 32): all positive distances above a certain class threshold τ_c will contribute more to the learning than the others. The second term is a global loss term used to prevent the similar samples to be disseminated in different places of the feature space and ultimately improve regularization and generalization. To do so, this term acts on the second order statistics of the distances to diminish the variance. The following notations are used: \mathcal{B} , the batch containing the outputs y of the neural network, \mathcal{P} and \mathcal{N} respectively the positive and negative samples of the batch, \mathcal{P}_c the samples of

the class c . The values $\mu_{\mathcal{P}}^t$ and $\mu_{\mathcal{N}}^t$ are the average distances between respectively the positive and negative samples, γ controls the smoothness of the evolution of these values between the previous batch \mathcal{B}_{t-1} and the present one \mathcal{B}_t . Finally, m , m^+ and m^- are margin parameters, η is a scaling parameters and λ controls the magnitude of the global loss term in order to avoid it outweighing the local one.

$$\mathcal{L}_{\text{structural}}(f(\mathcal{B})) = \frac{1}{B} \sum_{(i,j) \in \mathcal{P}} \beta_{ij} \log \left(1 + \sum_{n=1}^{|\mathcal{N}|} \exp(d(y_i, y_j) - d(y_i, y_n) + m) / \eta \right) + \frac{\lambda}{2} ([\sigma_p^2 - m^+]_+ + [\sigma_n^2 - m^-]_+) \quad (29)$$

$$B = \sum_{(i,j) \in \mathcal{P}} \beta_{ij} \quad (30)$$

$$\beta_{ij} = \exp(d(y_i, y_j)) - \tau_c \quad (31)$$

$$\tau_c = \frac{2}{|\mathcal{P}_c|} \sum_{(i,j) \in \mathcal{P}_c} d(y_i, y_j) - \min_{(i,j) \in \mathcal{P}_c} (d(y_i, y_j)) \quad (32)$$

$$\sigma_p^2 = \frac{1}{|\mathcal{P}|} \sum_{(i,j) \in \mathcal{P}} (d(y_i, y_j) - \mu_{\mathcal{P}}^{\mathcal{B}_t})^2 \quad (33)$$

$$\sigma_n^2 = \frac{1}{|\mathcal{N}|} \sum_{(i,n) \in \mathcal{N}} (d(y_i, y_n) - \mu_{\mathcal{N}}^{\mathcal{B}_t})^2 \quad (34)$$

$$\mu_{\mathcal{P}}^{\mathcal{B}_t} = \gamma \mu_{\mathcal{P}}^{\mathcal{B}_{t-1}} + (1 - \gamma) \mu_{\mathcal{P}}^{\mathcal{B}_t} \quad (35)$$

$$\mu_{\mathcal{N}}^{\mathcal{B}_t} = \gamma \mu_{\mathcal{N}}^{\mathcal{B}_{t-1}} + (1 - \gamma) \mu_{\mathcal{N}}^{\mathcal{B}_t}. \quad (36)$$

3 AVERAGE NORMALIZED DISTANCE TO NEAREST NEIGHBORS

We present in Figure 1 the evolution of the average distance to the nearest neighbors depending on the number of nearest neighbors. Each distance has been scaled between 0 and 1 prior to the average. For SGRU, we observe that the average distance grows slowly at first until some point between 150 and 200 (around the average number of samples per class in this validation set, which is about 165) and from there starts to grow faster. This seems to indicate that the training of SGRU has allowed the model to create an important margin between similar and dissimilar samples. The curve for CGRU is very different with a very smooth progression and no apparent margin. This can be explained by the coupling which always combines the New States in some extent even if the input are dissimilar and therefore produces closer outputs. Despite the absence of margin, CGRU gives comparable results as SGRU on UCI HAR dataset which indicates that nevertheless, similar samples are closer than dissimilar samples even if the coupling produces overall lower distances. We attribute this compensation to better performances of CGRU on hard samples.

REFERENCES

Xun Yang, Peicheng Zhou, and Meng Wang. Person reidentification via structural deep metric learning. *IEEE Transactions on Neural Networks and Learning Systems*, 30(10):2987–2998, 2018.

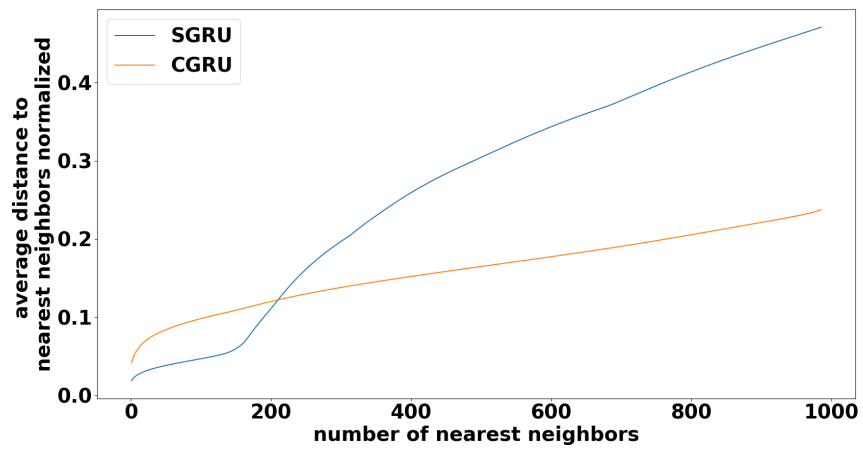


Figure 1: Average normalized distance to nearest neighbors for SGRU and CGRU computed on the validation set (users 1, 3 and 5)