

Supplementary Material

Kyle Min and Jason J. Corso

University of Michigan, Ann Arbor, MI 48109
 {kylemin,jjcorso}@umich.edu

A Center Update Equation

Notations

- D : Dimension of the embedding (feature) space
- $\mathbf{c}_j^r \in \mathbb{R}^D$: center of activity class j
- $\tilde{\mathbf{c}}_j^r \in \mathbb{R}^D$: l_2 -normalized center of activity class j
- $\mathbf{F}_i^r(j) \in \mathbb{R}^D$: video-level feature representation for the activity class j of the i -th training sample
- $\tilde{\mathbf{F}}_i^r(j) \in \mathbb{R}^D$: l_2 -normalized video-level feature representation
- $n_{i,j}^r = \underset{k \neq j}{\operatorname{argmin}} \mathcal{D}(\mathbf{F}_i^r(j), \mathbf{c}_k^r)$: index for the nearest negative center
- $\mathcal{D}(\cdot)$ represents the angular distance:

$$\begin{aligned} \mathcal{D}(\mathbf{F}_i^r(j), \mathbf{c}_j^r) &= \arccos \left(\frac{\mathbf{F}_i^r(j) \cdot \mathbf{c}_j^r}{\|\mathbf{F}_i^r(j)\|_2 \|\mathbf{c}_j^r\|_2} \right) \\ &= \arccos (\tilde{\mathbf{F}}_i^r(j) \cdot \tilde{\mathbf{c}}_j^r) \end{aligned}$$

- \mathbf{f}_i^r : new video-level feature that attends more strongly to the background features than \mathbf{F}_i^r is
- $\tilde{\mathbf{f}}_i^r$: l_2 -normalized new video-level feature
- $\mathbf{y}_i \in \mathbb{R}^{N_c}$: ground-truth annotation for video-level label of the i -th sample where N_c is the number of activity classes. $\mathbf{y}_i(j) = 1$ if the activity class j is present in the sample and $\mathbf{y}_i(j) = 0$ otherwise

Previous works on center loss [3, 1, 2] suggest using an averaged gradient (typically denoted as $\Delta \mathbf{c}_j^r$) to update the centers for better stability. In this section, we show how we update the centers based on this convention. For simplicity, we first look at the RGB stream.

As mentioned in the paper, $\tilde{\mathcal{L}}_{\text{ATCL}_{i,j}}^r$ and $\tilde{\mathcal{L}}_{\text{NT}_{i,j}}^r$ are the loss terms inside the max operation of the i -th sample and of the j -th activity class as follows:

$$\tilde{\mathcal{L}}_{\text{ATCL}_{i,j}}^r = \mathcal{D}(\mathbf{F}_i^r(j), \mathbf{c}_j^r) - \mathcal{D}(\mathbf{F}_i^r(j), \mathbf{c}_{n_{i,j}^r}^r) + m_1 \quad (1)$$

$$\tilde{\mathcal{L}}_{\text{NT}_{i,j}}^r = \mathcal{D}(\mathbf{F}_i^r(j), \mathbf{c}_j^r) - \mathcal{D}(\mathbf{f}_i^r(j), \mathbf{c}_j^r) + m_2 \quad (2)$$

Let $\mathbf{g}_{1,i,j}^r$ and $\mathbf{g}_{2,i,j}^r$ be the derivatives of Eq. 1 with respect to $\tilde{\mathbf{c}}_j^r$ and $\tilde{\mathbf{c}}_{n_{i,j}^r}^r$, respectively; and let $\mathbf{h}_{i,j}^r$ be the derivative of Eq. 2 with respect to $\tilde{\mathbf{c}}_j^r$:

$$\mathbf{g}_{1,i,j}^r = -\frac{\tilde{\mathbf{F}}_i^r(j)}{\sin\left(\mathcal{D}(\mathbf{F}_i^r(j), \mathbf{c}_j^r)\right)} \quad (3)$$

$$\mathbf{g}_{2,i,j}^r = \frac{\tilde{\mathbf{F}}_i^r(j)}{\sin\left(\mathcal{D}(\mathbf{F}_i^r(j), \mathbf{c}_{n_{i,j}^r}^r)\right)} \quad (4)$$

$$\mathbf{h}_{i,j}^r = -\frac{\tilde{\mathbf{F}}_i^r(j)}{\sin\left(\mathcal{D}(\mathbf{F}_i^r(j), \mathbf{c}_j^r)\right)} + \frac{\tilde{\mathbf{f}}_i^r(j)}{\sin\left(\mathcal{D}(\mathbf{f}_i^r(j), \mathbf{c}_j^r)\right)} \quad (5)$$

Then, we can represent the averaged gradient considering the three terms:

$$\Delta\tilde{\mathbf{c}}_j^r = \Delta\mathbf{g}_{1,i,j}^r + \Delta\mathbf{g}_{2,i,j}^r + \Delta\mathbf{h}_{i,j}^r \quad (6)$$

where each term of the right hand side is given by:

$$\Delta\mathbf{g}_{1,i,j}^r = \frac{1}{N} \left(\frac{\sum_{i:\mathbf{y}_i(j)=1} \mathbf{g}_{1,i,j}^r \delta(\tilde{\mathcal{L}}_{\text{ATCL},i,j}^r > 0)}{1 + \sum_{i:\mathbf{y}_i(j)=1} \delta(\tilde{\mathcal{L}}_{\text{ATCL},i,j}^r > 0)} \right) \quad (7)$$

$$\Delta\mathbf{g}_{2,i,j}^r = \frac{1}{N} \left(\frac{\sum_{i:\mathbf{y}_i(n_{i,j}^r)=1} \mathbf{g}_{2,i,j}^r \delta(\tilde{\mathcal{L}}_{\text{ATCL},i,j}^r > 0)}{1 + \sum_{i:\mathbf{y}_i(n_{i,j}^r)=1} \delta(\tilde{\mathcal{L}}_{\text{ATCL},i,j}^r > 0)} \right) \quad (8)$$

$$\Delta\mathbf{h}_{i,j}^r = \frac{1}{N} \left(\frac{\gamma \sum_{i:\mathbf{y}_i(j)=1} \mathbf{h}_{i,j}^r \delta(\tilde{\mathcal{L}}_{\text{NT},i,j}^r > 0)}{1 + \sum_{i:\mathbf{y}_i(j)=1} \delta(\tilde{\mathcal{L}}_{\text{NT},i,j}^r > 0)} \right) \quad (9)$$

Here, $\delta(\text{condition}) = 1$ if the *condition* is true and $\delta(\text{condition}) = 0$ otherwise. Using the chain rule, we can find the general form of the averaged gradient as follows:

$$\Delta\mathbf{c}_j^r = \frac{I_D - \tilde{\mathbf{c}}_j^r \otimes \tilde{\mathbf{c}}_j^r}{\|\mathbf{c}_j^r\|_2} \Delta\tilde{\mathbf{c}}_j^r \quad (10)$$

where I_D is an identity matrix of dimension D and \otimes denotes the outer product. Finally, the centers are updated using $\Delta\mathbf{c}_j^r$ for every iteration of the training process by a gradient descent algorithm. We can update the centers of other streams in a similar manner.

B More Qualitative Results

We provide more qualitative results of A2CL-PT. Please refer to the videos that are included in the subfolder.

References

1. He, X., Zhou, Y., Zhou, Z., Bai, S., Bai, X.: Triplet-center loss for multi-view 3d object retrieval. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1945–1954 (2018)
2. Li, Z., Xu, C., Leng, B.: Angular triplet-center loss for multi-view 3d shape retrieval. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 33, pp. 8682–8689 (2019)
3. Wen, Y., Zhang, K., Li, Z., Qiao, Y.: A discriminative feature learning approach for deep face recognition. In: European conference on computer vision. pp. 499–515. Springer (2016)