

CNN-based Action Recognition and Supervised Domain Adaptation on 3D Body Skeletons via Kernel Feature Maps (Supplementary Material)

Yusuf Tas^{1,2}, Piotr Koniusz^{1,2}
users.cecs.anu.edu.au/~koniusz

¹Data61/CSIRO, Canberra, Australia
²Australian National University

1 Supervised Domain Adaptation [□]

For the full details of the *So-HoT* algorithm, please refer to paper [□]. Below, we review the core part of their algorithm for the reader's convenience. Suppose \mathcal{I}_N and \mathcal{I}_{N^*} are the indexes of N source and N^* target training data points. \mathcal{I}_{N_c} and $\mathcal{I}_{N_c^*}$ are the class-specific indexes for $c \in \mathcal{I}_C$, where C is the number of classes. Furthermore, suppose we have feature vectors from an *FC* layer of the source network stream, one per an action sequence or image, and their associated labels. Such pairs are given by $\Lambda \equiv \{(\phi_n, y_n)\}_{n \in \mathcal{I}_N}$, where $\phi_n \in \mathbb{R}^d$ and $y_n \in \mathcal{I}_C$, $\forall n \in \mathcal{I}_N$. For the target data, by analogy, we define pairs $\Lambda^* \equiv \{(\phi_n^*, y_n^*)\}_{n \in \mathcal{I}_{N^*}}$, where $\phi_n^* \in \mathbb{R}^d$ and $y_n^* \in \mathcal{I}_C$, $\forall n \in \mathcal{I}_{N^*}$. Class-specific sets of feature vectors are given as $\Phi_c \equiv \{\phi_n^c\}_{n \in \mathcal{I}_{N_c}}$ and $\Phi_c^* \equiv \{\phi_n^{*c}\}_{n \in \mathcal{I}_{N_c^*}}$, $\forall c \in \mathcal{I}_C$. Then $\Phi \equiv (\Phi_1, \dots, \Phi_C)$ and $\Phi^* \equiv (\Phi_1^*, \dots, \Phi_C^*)$. The asterisk in superscript (e.g. ϕ^*) denotes variables related to the target network while the source-related variables have no asterisk. The *So-HoT* problem is posed as a trade-off between the classifier and alignment losses ℓ and \hat{h} . Figure 4 (the main submission) shows the setup we use. The loss \hat{h} depends on two sets of variables (Φ_1, \dots, Φ_C) and $(\Phi_1^*, \dots, \Phi_C^*)$ – one set per network stream. Feature vectors $\Phi(\Theta)$ and $\Phi^*(\Theta^*)$ depend on the parameters of the source and target network streams Θ and Θ^* that we optimize over. $\Sigma_c \equiv \Sigma(\Phi_c)$, $\Sigma_c^* \equiv \Sigma(\Phi_c^*)$, $\mu_c(\Phi)$ and $\mu_c^*(\Phi^*)$ denote the covariances and means, respectively, one covariance/mean pair per network stream per class. Specifically, we solve:

$$\begin{aligned} & \arg \min_{W, W^*, \Theta, \Theta^*} \ell(W, \Lambda) + \ell(W^*, \Lambda^*) + \eta \|W - W^*\|_F^2 + & (1) \\ & \text{s. t. } \|\phi_n\|_2^2 \leq \tau, & \frac{\alpha_1}{C} \sum_{c \in \mathcal{I}_C} \|\Sigma_c - \Sigma_c^*\|_F^2 + \frac{\alpha_2}{C} \sum_{c \in \mathcal{I}_C} \|\mu_c - \mu_c^*\|_2^2. \\ & \|\phi_{n'}^*\|_2^2 \leq \tau, & \\ & \forall n \in \mathcal{I}_N, n' \in \mathcal{I}_{N^*} & \underbrace{\hspace{10em}}_{\hat{h}(\Phi, \Phi^*)} \end{aligned}$$

For ℓ , a generic Softmax loss is employed. For the source and target streams, the matrices $W, W^* \in \mathbb{R}^{d \times C}$ contain unnormalized probabilities. In Equation (1), separating the class-specific distributions is addressed by ℓ while attracting the within-class scatters of both network streams is handled by \hat{h} . Variable η controls the proximity between W and W^* which encourages the similarity between decision boundaries of classifiers. Coefficients α_1 , α_2 control the degree of the cov. and mean alignment, τ controls the ℓ_2 -norm of feature vectors.

2 Modifications to the So-HoT Approach

Algorithm 1 details how we perform domain adaptation. We enable the alignment loss \tilde{h} only if the source and target batches correspond to the same class. Otherwise, the alignment loss is disabled and the total loss uses only the classification log-losses ℓ_{src} and ℓ_{trg} . To generate the source and target batches that match w.r.t. the class label, we re-order source and target datasets class-by-class and thus each source/target batch contains only one class label at a time. Once all source and target datapoints with matching class labels are processed, remaining datapoints are processed next. Lastly, we refer readers interested in the details of the So-HoT algorithm and loss \tilde{h} to paper [10].

Algorithm 1 Batch generation + a single epoch of the training procedure on the source and target datasets.

```

1:  $src\_data := sort\_by\_class\_label(src\_data)$ 
2:  $target\_data := sort\_by\_class\_label(target\_data)$ 
3:  $C_s$  ▷ Number of the source classes
4:  $C_t$  ▷ Number of the target classes
5:  $C_{s \cap t}$  ▷ Number of classes in common
6: procedure EPOCH( $src\_data, target\_data, batch\_size$ ) ▷ Training (one epoch)
7:   for  $i \leftarrow 1 : \max(C_s, C_t)$  do
8:     if  $i \leq C_s$  then
9:        $batch_s \leftarrow Choose(src\_data, i, batch\_size)$  ▷ ‘Choose’ pre-fetches data of
       class  $i$ 
10:    else
11:       $batch_s \leftarrow Choose(src\_data, rnd(), batch\_size)$  ▷ ‘Choose’ pre-fetches data
       of random class
12:    if  $i \leq C_t$  then
13:       $batch_t \leftarrow Choose(target\_data, i, batch\_size)$ 
14:    else
15:       $batch_t \leftarrow Choose(target\_data, rnd(), batch\_size)$ 
16:    if  $i \leq C_{s \cap t}$  then
17:       $Loss \leftarrow \ell_{src} + \ell_{trg} + \tilde{h}$ 
18:    else
19:       $Loss \leftarrow \ell_{src} + \ell_{trg}$ 
20:    Forward( $net\_data, batch_s, batch_t$ )
21:    Backward( $net\_data, batch_s, batch_t$ )
22:    Update( $net\_data, batch_s, batch_t$ )

```

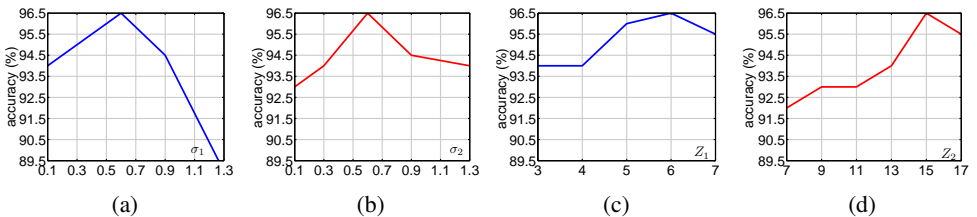


Figure 1: Sensitivity w.r.t. parameters σ_1 and σ_2 on UTK. Figures 1a, 1b, 1c and 1d show the accuracy w.r.t. σ_1 ($\sigma_2 = 0.6$), σ_2 ($\sigma_1 = 0.6$), Z_1 ($Z_2 = 15$) and Z_2 ($Z_1 = 5$), respectively.

References

- [1] Piotr Koniusz, Yusuf Tas, and Fatih Porikli. Domain adaptation by mixture of alignments of second- or higher-order scatter tensors. *CVPR*, 2017.