

**Exercise for MA-INF 2218 Video Analytics SS21**  
**Submission on 06.07.2021**  
**Action Segmentation with TCNs**

In this sheet, you will implement an action segmentation model using temporal convolutional networks (TCNs). Given an input video  $\mathbf{x}_1^T = (x_1, \dots, x_T)$  with  $T$  frames, the goal of action segmentation is to predict the action label for each frame  $\mathbf{c}_1^T = (c_1, \dots, c_T)$ , where  $c_t$  is the action label for frame  $t$ . Your system for this sheet will be based on the model described in [1]. The input  $\mathbf{x}_1^T$  is also the same frame-wise I3D features that are used in [1], where  $x_t \in \mathbb{R}^{2048}$ . Your implementation should be using Python 3 and PyTorch  $\geq 1.1$ .

Download the training data using the provided link in the **README** file and place it in the same directory of your **main.py** file. For details on the features and annotations, please refer to the **README** file.

1. Implement a single-stage TCN as described in Section 3.1 in [1]. However, replace the dilation factors with a dilation factor that increases linearly with the number of layers (*i.e.* 1, 2, 3, 4, 5, ....)
  - Train a network with 10 layers and 64 filters in each layer for 50 epochs with a cross-entropy loss. Use Adam optimizer with a learning rate 0.001 and a batch size of 4.
  - Evaluate the trained model on the test set and report the frame-wise accuracy, edit distance and F1-scores. Use the provided helper functions for evaluation.

(4 Points)

2. Implement a multi-stage TCN as described in Section 3.2 in [1] using the single-stage TCN from Question 1. However, instead of passing only the predicted probabilities to the next stage, concatenate the input features with the predicted probabilities of each stage before passing it to the next one.
  - Train a multi-stage model with 4 stages for 50 epochs with a cross-entropy loss. Use Adam optimizer with a learning rate 0.001 and a batch size of 4.
  - Evaluate the trained model on the test set and report the frame-wise accuracy, edit distance and F1-scores.

(4 points)

3. Repeat Question 2 with an additional video-level loss. The new loss computes the binary cross-entropy between a multi-class video level prediction and a multi-class target that indicates which classes are present in the video. The target is a vector with a dimension equals the number of classes in the dataset. The  $i$ -th element of this vector is 1 if the  $i$ -th class is present in the video, otherwise it should be 0. To get the video level prediction apply a max pooling on the temporal dimension of the predicted frame-wise logits. Note that the final loss function is the sum of the video level loss and the frame-wise cross-entropy loss that is used in Question 2.

(4 points)

4. Implement a multi-scale model using the single-stage TCN with 10 layers from Question 1. The model should contain three parallel branches, where each branch is a

single-stage TCN that receives a downsampled version of the input features and predicts a downsampled version of the frame-wise labels. The first branch operates on the full temporal resolution whereas the other branches operate on a downsampled version with factors 4 and 8 respectively. To get the final prediction, upsample the output of the last dilated layer from each branch (the output before the last  $1 \times 1$  convolution layer, *i.e.* the classification layer) and pass the average of these outputs to another  $1 \times 1$  convolution that predicts the logits for the output classes. Train your model using cross-entropy loss on the final output and on the output of each branch as well. Report the results on the test set.

*(8 points)*

Please exclude the `data` folder from your submission due to its size. However, submit any file you created yourself. Also prepare a table that summerizes the results from all the questions. If you have any questions, feel free to contact me ([abufarha@iai.uni-bonn.de](mailto:abufarha@iai.uni-bonn.de)).

[1] Y. Abu Farha and J. Gall, **MS-TCN: Multi-Stage Temporal Convolutional Network for Action Segmentation**, CVPR 2019.