Activity and Emotion Detection for "Smart Homes"

Inimfon Akpabio

March 2, 2021

**Abstract**

Advances in IoT (Internet of Things) and artificial intelligence have facilitated the deployment of smart home devices that enable hospitals and health workers monitor the state and well-being of patients, particularly elderly people, who often live alone. Such technologies can be crucial in tracking the health of a patient and ultimately, providing emergency care to such individuals if needed. In this project we propose a model to detect hair care related activities like brushing or combing in a smart home setting. Our model starts with a DenseNetV2 network used for feature extraction and then feeds into LSTM + Fully-connected layers to perform the final prediction. The feature extractor uses pre-trained image-net weights. For each video, we sample only 20 frames. The rationale is that an activity like brushing is a perpetual motion, hence, this should be sufficient. A mixture of videos from both the Kinetics700 and HMDB51 datasets were used to train the overall classifier. The model currently has an accuracy of 94% which rivals the state of the art.

## 1. Topic

The selected action for this project is "Combing hair".

## 2. Motivation

Human activity recognition aims to detect common or pre-defined human activities either from sub-sampled video clips or full-length videos. Typically, in the case of full-length videos, the goal is to temporally localize the action. This field of research is particularly attractive due to the potential applications in health care and personalized care for elderly citizens. With the population of elderly citizens projected to rise to 22% by 2050 [1], smart monitoring systems could help such individuals to lead healthy and yet independent lifestyles [2]. There have already been multiple attempts to monitor human activities using signals from wearable devices [3]. While these methods can be useful, they are inevitably limited due to the one-dimensional nature of this approach. Hence, it makes sense to utilize visual data as a means for detection. The benefit of the visual approach is that it affords more features for better detection and can be easily collected using video capture devices. A lot of research work has been done in the field Action/emotion detection in videos in recent years. Advances in computer-vision AI and the advent of competitions worldwide [4] have significantly increased interest in video based approaches for action detection. In this project we will build a model that relies solely on visual data to relate hair care activities such as brushing/combing hair.

## 3. Related Works

While this model design is not directly based on any particular paper, there are multiple research endeavors that have adopted similar designs to tackle video and speech related problems. For instance, [5] uses VGG-16 layers to extract features from a video clip and then feeds this data to Bi-directional LSTM network to generate captions for the video.

## 4. Proposed model

Fig 1 shows the complete model design. The model starts with a DenseNet121 network [6] which has been pre-trained with ImageNet weights, hence, making it ideal for feature extraction. At the end of the DenseNet121 network, we perform a global maxpool which basically takes the maximum value from each filter. This helps to reduce the dimensionality of features and promotes spatial hierarchy. Next using time-distribution, the features are fed into a LSTM network with 64 units. Finally the output the GRU is fed into a fully connected network that performs prediction. The fully connected network consists of two layers with 1024 and two units respectively. The second layer applies Softmax activation on the 2 units to perform binary-classification.

As a pre-processing step, video clip frames must be resized to (224, 224) and scaled to [0, 1]. The input to the model is an array of 20 image frames sub-sampled from a 2 second clip and has a dimension of (20, 224, 224, 3). It outputs an array of size 2 representing the prediction probabilities for "Combing hair" and "not Combing hair".

## 5. Dataset

Feature Extractor:

The feature extractor was pre-trained on ImageNet [8]. Most of the proposed solutions to activity recognition often employ pre-trained models due to their robustness and superior ability to extract features.

Full model:

Finding a video dataset to train the full model proved difficult. The Kinetics700 dataset is a large datasets that contains a myriad of human-performed actions and has been widely explored for video recognition tasks in recent years [9]. Due to its massive size, it could take days to download the entire dataset even with a stable internet connection. Hence, it made sense to download only the "brushing hair" class and combine this with the HMDB51 dataset which similarly contains videos of human-performed actions [10]. For training, the videos were separated into two classes: "brushing/combing hair" and "miscellaneous" which represents the collection of all "not brushing/combing hair" videos. Both classes currently have a size of 538 and 412. As a preprocessing step, the videos were sample-centered to improve convergence during training. It is a common practice to centralize and normalize data before training neural networks.

No data Augmentation was applied, however a future improvement could be to include more hair types (e.g. curly hair) in the dataset.

## 6. Model Training and Performance

Hyper-parameters like batch-size, epochs and learning rate were selected by testing pre-selected values in fixed ranges and choosing the values that yielded the best convergence. Fig. 2 shows a snippet of the learning rate optimization. We use the Adam optimizer with a learning rate of 0.0005 and about 12 epochs to train the model. After this many epochs, the model begins to exhibit overfitting, hence, we employ early stopping to prevent further training. We also curbed the overfitting by adding L2 regularization and dropout of 50% to the fully connected layers. On the validation set, the model achieves an accuracy of about 93% and an overall accuracy of 91% on YouTube videos. The breakdown is about 94.8% on positive (contains desired action) YouTube videos and 87.25% on negative YouTube videos.

Quite noticeably, the model seems to experience dips in performance when dealing with hairstyles that are uncommon in the dataset. For instance, the accuracy on test clips containing wavy, African-American hair is less than 90%. This is obviously due to lack of diversity in our training data. This also alludes to the common issue of lack of representation in machine-learning. If such systems are to be deployed in real-life, they should be robust enough to handle different demographics.

## 7. YouTube

On the curated YouTube test set, the performance on positive videos is a whopping 94.8%. This is probably due to the ability of DenseNet to extract salient features for classification. The FPR and FNR are both 12.75% and 5.2% respectively. Hence there is still room for improvement.

As of right now, the videos are manually segmented into 2 second clips. In the future we will aim to develop models that are capable of temporal localization. Hence, if we feed in an entire video, we expect the model to output time segments that contain the action.

## 8. Improve Accuracy

Next we will explore advanced techniques such as optical flow and Feature Landmarks.
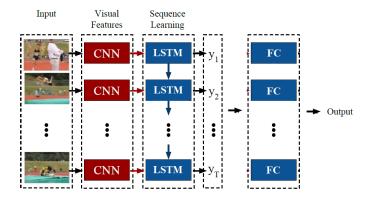
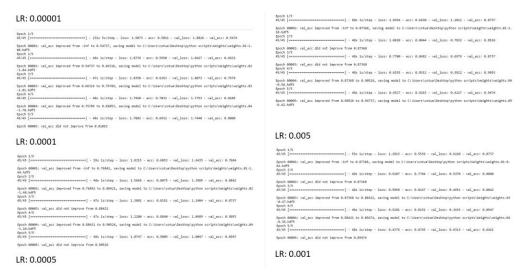Fig 1: Model Diagram

Source: Adapted from [7]



Fig 2: Learning Rate Optimization

Hyper-parameter optimization for Adam optimizer learning rate. The value of 0.0005 was ultimately selected based on validation accuracy.

# References

[1]     ChildStats.gov. "Share of Old Age Population (65 Years and Older) in The Total U.S. Population from 1950 to 2050." Statista, Statista Inc., 23 Sep 2020, https://www.statista.com/statistics/457822/share-of-old-age-population-in-the-total-us-population/.

[2]     K. Davis et al., "Activity recognition based on inertial sensors for Ambient Assisted Living," 2016 19th International Conference on Information Fusion (FUSION), Heidelberg, Germany, 2016, pp. 371-378.

[3]     Huang, Hua. (2016). Toothbrushing Monitoring using Wrist Watch. 202-215. 10.1145/2994551.2994563.

[4]     Dhall, Abhinav & Goecke, Roland & Gedeon, Tom & Sebe, Nicu. (2016). Emotion recognition in the wild. Journal on Multimodal User Interfaces. 10. 10.1007/s12193-016-0213-z.

[5]     Y. Bin, Y. Yang, F. Shen, N. Xie, H. T. Shen and X. Li, "Describing Video With Attention-Based Bidirectional LSTM," in IEEE Transactions on Cybernetics, vol. 49, no. 7, pp. 2631-2641, July 2019, doi: 10.1109/TCYB.2018.2831447.

[6]     G. Huang, Z. Liu, and K. Q. Weinberger. Densely connected convolutional networks. In IEEE Conference on Computer Vision and Pattern Recognition, 2017.

[7]     Donahue, J.; Hendricks, L.; Guadarrama, S.; Rohrbach, M.; Venugopalan, S.; Saenko, K.; Darrell, T. Long-term recurrent convolutional networks for visual recognition and description. IEEE Trans. Pattern Anal. Mach. Intell. 2017, 39, 677–691

[8]     L. Fei-Fei and O. Russakovsky, Analysis of Large-Scale Visual Recognition, Bay Area Vision Meeting, October, 2013

[9]     J. Carreira, E. Noland, C. Hillier, and A. Zisserman. A short note on the kinetics-700 human action dataset. arXiv preprint arXiv:1907.06987, 2019.

[10]    H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre. Hmdb: A large video database for human motion recognition. In 2011 International Conference on Computer Vision, pages 2556–2563, 2011.