
Second Assignment (Deep Learning based on CNNs)

TUTORS: MIRELA POPA

TEACHING ASSISTANT: SUBILAL VATTIMUNDA PURAYIL, AASHUTOSH GANESH

DEADLINE: 29/05/2024 AT 11.59 PM

DESCRIPTION

In this assignment you will get familiarized with deep learning models, namely CNNs and you will develop a solution for a classification problem, namely the emotion recognition task. Given a dataset of images with different facial expressions labeled as one of the basic emotions (e.g. anger, happiness, surprise, sadness, fear, disgust or neutral), you will need to implement a model based on CNNs, which gets trained on a subset of the dataset and is tested on another set. You will analyze the role and effect of different architectures and parameters of your CNN model (e.g. number of layers, filters, pooling, normalization, activation functions, etc.), and report your accuracy and observations in a report. Optionally, you can compare your results to the ones achieved using a pre-trained model (e.g. VGG16), but the scope of this assignment is to understand the effect of the different CNN operations on your results and apply this model on a real-life sample image/video.

This assignment is a group assignment and will be done in pairs.

STEPS

- First, you need to select a deep learning framework on which to perform your experiments. I recommend PyTorch (<https://pytorch.org/>), you will also find on Canvas a short tutorial about it, but you can use a detailed tutorial from the ones available online (e.g. <https://www.tutorialspoint.com/pytorch/index.htm>).
- Then, you need to select the emotion dataset on which to perform your experiments (e.g. you can use the FER2013 dataset consisting of 34k images, which can be easily

downloaded from [Kaggle](#), providing a training and a testing set). You can also use another emotion dataset, but the one I mentioned is small enough to be trained on a CPU (big datasets might require more powerful resources, such as GPU processing power, which goes beyond the scope of this assignment).

- Next, you need to read, parse and load your dataset and for this task you can use PyTorch dataset functionalities such as `torch.utils.data.DataLoader` (https://www.tutorialspoint.com/pytorch/pytorch_loading_data.htm), but also available parsers for csv files (format in which the FER dataset is saved) to [h5py](#). Additionally, transformations can be applied on your images for improving their representation (e.g. cropping, contrast, saturation changes). Consider also augmentation techniques.
- Building your CNN model is one the most important issues and you should try different architectures and configurations, including different number of layers, filters (e.g. check the [Conv2D](#) function from `torch.nn`, which provides many options for configuring CNN parameters). You can use a [Sequential](#) model, for which you can specify the types of desired operations after each convolutional layer, including batch normalization and the activation function (e.g. ReLU). After building your model, you need to select a loss function (e.g. cross entropy), and an [optimization module](#), along with a learning rate. Train and test your CNN model using the usual pipeline (check the PyTorch tutorial available on Canvas for a simple example, but also other [resources](#) online).
- Please try different architectures and parameters and based on the performed experiments select a final model. Please also consider Visual Transformers or Self-Supervised pre-trained models (available on [HuggingFace](#)). For these additional models is important to understand how they work and to explain why you think they would benefit your task. The made decisions and the rationale behind them will be explained in the report (e.g. why did you use specific filter sizes and how did it help at improving the accuracy? In which cases it would be useful to use a self-supervised representation?). You can consider the [Weights & Biases](#) framework for managing your models and hyper-parameters.
- Check the model accuracy, precision, recall on the test data and report it using several measures (quantitative vs. qualitative), also make sure you use a good testing methodology (e.g. n-fold cross validation or at least split the data into training/testing sets).
- Visualize the learnt filters on different layers to check which is the internal representation of your CNN model (e.g. you can check the model parameters for each convolutional layer and transform the filters into an image). Display the relevant activations for different emotions.
- Test your model on two-three real-life videos captured with your web cam (e.g. at least 1 second long) and assess its performance (in the simplest case sample one frame per

second, pre-process it by detecting the face, crop the facial region, feed it to your best model and get the confidence of the detected emotional state). This exercise will help you assess the generalization ability of your best model.

- Formulate your critical reflections about the performed experiment and the lessons learnt (e.g. what is important when constructing a CNN/ViT/Self Supervised model and how to optimize it?).

DELIVERABLES

- A) The code and a well-documented readme file with specific details about installation/running procedure, employed parameters.
- B) A report of maximum 1200-1500 words consisting of your results, plots and findings and an explanation for your final CNN/ViT/Self Supervised architecture.
- C) Upload everything on Canvas as a zip file, having the following format: Surname1_Surname2_Second_Assignment_CV2024.zip. As this is a group assignment, only one of you needs to upload it.

GRADING

A) CODE – 1.5 POINTS

B) REPORT – 1 POINT