**Week - 08**
**Lecture - 27**
**Tutorial: Multimodal Emotion Recognition**

(Refer Slide Time: 00:22)



Hello everyone, welcome to this Affective Computing tutorial on Multi-modal Emotion Recognition. In this tutorial we will try to learn emotion recognition using audio and video information together.

By integrating information from different modality such as audio and video, one can gain a more complete and accurate understanding of person's emotion state and there are multiple different techniques for integrating information from multiple modalities including feature level fusion, decision level fusion and hybrid approaches that combine both.

Each approach has its own strength and weakness and the choice of technique will depend upon specific application and data availability.

(Refer Slide Time: 01:09)



So, in this tutorial we will be using a RAVDESS dataset. The RAVDESS dataset is a popular dataset used for emotion recognition research. It contains recording of actors speaking and performing various emotional states including anger, disgust, fear, happiness, sadness and surprise.

The dataset includes both audio and video recordings which makes it a well-suited for multi-modal emotion recognition research. The audio recordings consist of 24 actors speaking scripted sentences in each of six emotional states. The video recording consist of same actor performing facial expression and body movements to convey the same six emotional states.
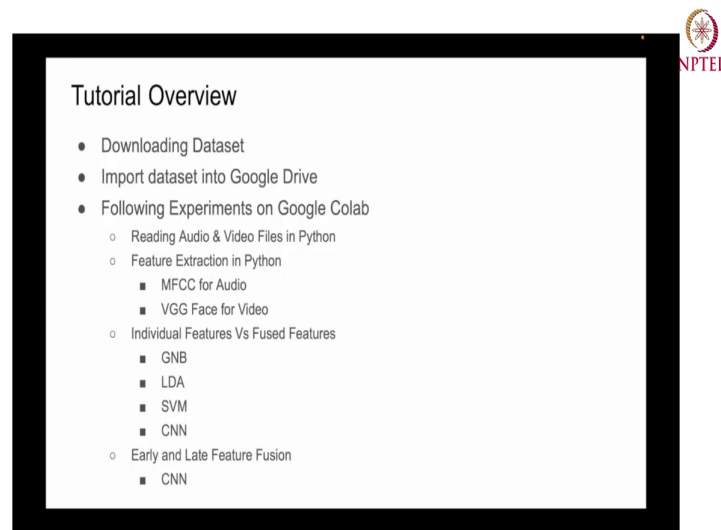
(Refer Slide Time: 01:58)



Filename Identifiers: 03-01-**01**-01-01-01-01.wav

- Modality (01 = full-AV, 02 = video-only, 03 = audio-only).
- Vocal channel (01 = speech, 02 = song).
- Emotion (01 = neutral, 02 = calm, 03 = happy, 04 = sad, 05 = angry, 06 = fearful, 07 = disgust, 08 = surprised).
- Emotional intensity (01 = normal, 02 = strong). NOTE: There is no strong intensity for the 'neutral' emotion.
- Statement (01 = "Kids are talking by the door", 02 = "Dogs are sitting by the door").
- Repetition (01 = 1st repetition, 02 = 2nd repetition).
- Actor (01 to 24. Odd numbered actors are male, even numbered actors are female).

Now, coming to the filename convention in this dataset, each filename consist of seven numerical identifier in which third identifier tells about the emotional class.

So, to get started with the multi-modal emotion recognition using the RAVDESS dataset, we will first need to import audio and video data, an extractor element feature from both modalities. After that we can use variety of machine learning models to classify the emotional states based on the extracted features. It is worth noting that the performance of your multi-modal emotion recognition system will depend on the quality of the feature extracted and fusion techniques used.

As well as the choice of machine learning model so, it is important to carefully evaluate the performance of the system and fine tune the network. So, we will start with downloading and coding the dataset into the Google Drive. After that, we will be using some predefined Python libraries to read audio and video files in Python environment. Later we will do some feature

extraction, particularly this MFCC for audio data. And we will extract video feature using VGG 16 phase model.

After that, we will try to see individual features versus fused features performance. And we will be using Gaussian base linear discriminant analysis, support vector machine and our CNN network for this analysis. And later on we will be doing our early and late fusion and we will be using CNN for this. So, now let us jump directly upon the coding part.

(Refer Slide Time: 03:41)



We will start with installing couple of libraries that might not be pre-installed in the Google colab environment. Once we have installed all these necessary libraries, we can begin with importing them in the code and we will start pre-processing the data and extract given feature for our multi-modal emotional recognition.

And I am also assuming that you have already downloaded the dataset and stored it into your respective Google drives. So, our downloading library code will look something like this. And it might take a couple of seconds to download all these files.

(Refer Slide Time: 04:37)



So, after downloading my relevant packages, I will import all my relevant libraries for this tutorial and the code will look something like this. After downloading and importing the libraries, I will start with defining some path variables. In this tutorial, I will be basically defining three of my path variable. First will be our video path which essentially contains the exact dataset path or video files are stored.

And then I will also define two other paths where I will be storing video only information and audio only information extracted from the original dataset. One more thing, I will be just

using actor 1 data for this analysis and the only purpose of using actor1 data is to reduce the computational demand.

So, as of now, my path's variables are all set. So, I will write first piece of code which will essentially take our video files from original data and extract video and audio information from that video file and save it to two separate folder. And the code will look something like this.
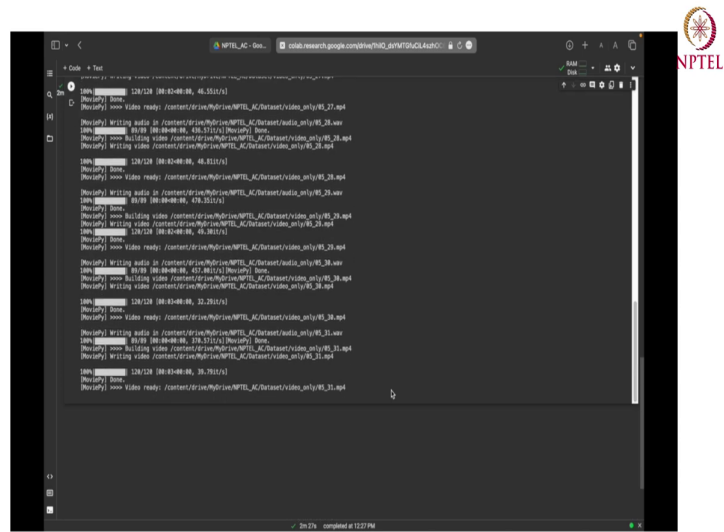
So, let me give you an overall overview of this piece of code but we are essentially doing here is we are iterating over all the file in the dataset folder and we will be only using emotion number 2 and emotion number 5, which essentially corresponds to calm and anger emotions.

The only reason I am using only two emotion classes just reducing down the computational power. As in the general case, Google colab only consist of 12 GB of RAM. For computational purpose only, I will be using just two of the emotion classes. Then after reading the video from their location, I will be resizing them into a smaller size and again, the reason is to reduce the computational requirement.

And after that, I will simply I will be simply considering first four second of the video to reduce the computational requirement and also to make a consistency among all the videos. Like there could be some instance that a particular video might consist of 5 second and other might be consisting of 4 second.

So, I am using a smallest threshold like 4 second and every video so, that all my data will be consisting along the time dimension. Now, I will extract audio information only using this command video clip dot audio. And later I will be storing these separate audio and video content in these two folders. So, now running this code and this code might take good amount of time as we are reading and writing the data.

(Refer Slide Time: 07:35)



Audio and video files are created from the original dataset. Now, I will write another piece of code to bring those audio only and video only files in my Python environment.

(Refer Slide Time: 07:50)



And my code will look something like this.

(Refer Slide Time: 07:55)



So, I will be just declaring to list video label all and video only all these two list will be storing the labels and the exact video only data. And I will be iterating through all the files in my video only path and appending into these two list. Later I will convert these two list into NumPy array for our ease of use.

And I will write another piece of small code which will essentially convert our label into 0 1 binary labels. After execution of this code, I will be using similar sort of code for bringing audio only data in my Python environment and my code will look something like this.

So, now as you can see that I have bought my video only and audio only data in my Python environment and the shape of video only data will look like 32 cross 120 cross 128 cross 227 cross 3.

So, what does these number actually mean is like there are 32 files there are 32 video files and each video files contain 120 frames since the video data was sampled at 30 hertz and we are taking just four second of data. So, there will be 120 frames and each frame will be consisting of 128 cross 227 dimension and there will be three channel red, green and blue.

Similarly, for our audio data we will be considering first four second of audio and all the audio will be sampled at 8000 hertz which essentially means that for one second of data there will be 8000 samples and using this strategy we will be collecting all those 32 files and each file will be consisting of 32000 sample which essentially means 4 second multiplied by 8000

sample rate. You and video only files in my Python environment I will start with the feature extraction part.

In our feature extraction part. So, in our feature extraction part the first feature will be mfcc feature which essentially called as MEL frequency cepstral coefficient and it is a commonly used feature extraction technique for audio analysis including in the text of multimodal emotion recognition. mfcc represent a compact yet efficient way of capturing the spectral characteristics of an audio signal and they have been shown to be effective at capturing emotion information in speed signals.

So, there are some basic steps that involve in I am putting the mfcc from audio signal and these steps consist of the very first step consist of applying a high passed filter to the signal to amplify the higher frequency component and reducing the impact of lower frequency components. Then second step mfcc process is frame segmentation where the signal is divided into short frames usually 20 to 30 milliseconds with some overlap.

Third step is windowing. A window function such as hamming window is applied to each frame to reduce the spectral leakage and to smooth of the signal and then we perform a Fourier transform where that windowed signal is transformed into frequency domain using the discrete Fourier transformation and after that there will be MEL frequency whopping. Transform DFT spectrum will be converted into a MEL frequency scale which essentially approximates the way human perceive the sound.

After that there is a usual cepstral analysis over these whop MEL frequencies. The MEL frequency spectrum is transformed into cepstral domain using inverse of discrete cosine transformation and then we select a subset of resulting cepstral coefficient as the mfcc coefficients.

These mfcc coefficients can be used as a feature for emotion recognition classification either on their own or maybe with the combination of other feature and for our purpose we will be

using a standard library called librosa which will provide the function to easily compute mfcc from audio signals and to do so, our code will look something like this.

In this code we have declared a list where we will be storing our mfcc coefficients and we will be iterating through all the audio only files and we will be using this predefined function librosa dot feature dot mfcc and passing each file in this predefined function and there is a parameter in this function called number of mfcc coefficient and for that we have chosen a value equal to 40.

And after computation of these mfcc coefficients we will convert that list into a array and we will reshape this mfcc list as number of sample followed by number of coefficients. So, after extracting the mfcc features from audio files we will move towards extracting video features from the video only files. To do so, we will take the advantage of VGG16 phase architecture.

So, VGG16 is basically a deep convolution neural network architecture that was originally designed for image classification task. This architecture for image classification consist of 16 layers including 13 convolution layers and 3 fully connected layers. In the standard architecture input to the network is an RTP image of size 224 cross 224 pixels and this input layer will be followed by 13 convolution layers each convolution layer with a 3 cross 3 filter size and a stride of 1 pixel.

In this architecture the number of filter increases as we go deeper into the network basically ranging from 64 in the first layer to 512 in the last layer. Then after every two convolution layer a max pooling operation is applied to reduce the spatial resolution of the feature map through the convolution future maps coming through the convolution layers and after that a fully connected layer is applied which consist of 4096 neuron each and after that in the standard architecture there is a three layered multi-layer perceptron for final prediction and the final prediction is performed using softmax layer.

So, in our architecture we will be using this standard VGG16 phase model it is a pre-trained model and to do so, our code will look something like this and in this code instead of using 224 cross 224 image we will be using 128 cross 227 dimensional image and also I have

downloaded this file this pre-trained weight file in my local drive and I have already given you the exact code for downloading this file.

So, after running this file our model is initiated. So, one essential step that we have to do in this pre-trained model is to freeze down all the convolution layer so, that they can extract relevant feature according to their pre-trained settings. To do so, my code will look something like this where I will iterate through all the layers in my vgg model and I will simply passed an argument layer dot trainable equal to false.

After making our feature learning layers as a non-trainable layers, we can put a another multi-layer perceptron over it on the code will look something like this. So, what does this over the top model is doing and this is taking our vgg model architecture and then we are simply putting a flattened layer which will flat down the extracted features into a single vector and then we will be perform batch normalization and then we put our 21 neuron layer with the activation of relu and we will name this layer as a feature layer.

So, this will be our main layer for extracting the features and then we will simply perform a classification operation on our dataset and extract the feature learned at this layer.

So, you can see that our vgg model has given 4 cross 7 cross 5 12 dimensional feature shape then we have flatten that out and after performing the batch normalization we will be taking this 21 dimensional feature. Now, I can simply compile my model using categorical cross entropy loss and with my Adam optimizer with learning rate equal to 0.0001.

After this I can simply fit my model with the video data for two classes with a batch size of 128 and number of epochs let us say 20. So, fitting this model over the 20 box this thing also might take some time. So, please have some patience.

(Refer Slide Time: 20:17)



So, as you can see we have already trained our model now and we are getting a accuracy somewhat around 86.41 percent which is a good discriminating accuracy saying that we are able to discriminate between the two emotion classes using VGG16 feature. Now, we will simply write a code to extract that feature layer this feature layer for that I will write a simple code there. I will create a feature extractor which is essentially taking model input as our input and output will be the intermediate layer which is our feature layer.

Now, I can pass all my data point to this feature extractor and we will get the corresponding features. For that I can write a simple code and the code will look something like this. Here I have declared a list called video only feature list and I will iterate through all the files and I will pass those file into this feature extractor function over here.

And then I will simply extract the feature convert them into the numpy array and I will append these features to a this video only feature list and I will convert that list into a numpy array for our ease of use. One another thing that you have to look here is after extracting the numpy array I am reshaping it to 120 cross 21.

So, what does this basically means is since our original data consist of 32 files each files containing 120 images of dimension 128 cross 227 cross 3. So, basically for each file I am passing all these 120 images and for one particular image I am getting a 21 dimensional feature from my VGG 16 architecture.

Now, I am simply appending these feature sequentially so, that I can get a long feature for these 120 images and these 120 images basically correspond to our 4 seconds of video. So, in this way I am creating feature for all these 32 images.

So, running this code we got 32 files each having dimension 2520. So, by now we have created our audio features as well as our video features. Now, we will simply first try to see how these audio feature and video features are working with respect to our basic classifier which our Gaussian base, linear discriminant analysis and our support vector machine. To do so, I will import all these classifiers from the standard sklearn library and my classification code will look something like this.

Here first of all I have I will be using my audio features which are mfcc audio only features and I will divide these feature into my train and test plate where my test split size will be 20 percent and I will first call my Gaussian base classifier and fit that data on my training set and then again I will call my linear discriminant analysis classifier then support vector machine with linear kernel and after that support vector machine with my RBF kernel.

Let us see how these feature are actually performing with these basic classifiers. As you can see we are getting a training score of somewhat around 76 percent using Gaussian base linear discriminant analysis and our support vector classifiers and support vector machine with RBF kernel is giving a lower score. So, which can also signify that the discriminating boundary between these two classes are more of a linear boundary.

After this we will perform same thing with our video features and the code will look something like this.

(Refer Slide Time: 25:17)



This is similar code this is similar code as we are doing with the audio feature. Instead of audio feature which I will simply use my video only features and they are corresponding label and the test size will again be 20 percent and we will split it into our train and test splits and

then we will call our Gaussian base linear discriminant analysis support vector machine with linear kernel and support vector machine with our RBF kernel.

So, let us run this code and see how our video features are performing. So, as we can see our video features are really good and they are giving a accuracy of somewhat around 100 percent for each instances including our S codes. But since we are already getting this much of accuracy ah. So, let us try to see how our accuracy changes if we simply combine these two features together.

So, what I am basically trying to say is, I will simply concrete my audio features, audio only features and my video only features together and then I will pass these fused feature in my in these basic machine learning classifier and I will try to see how my accuracy changes with that. So, for my fusion part my code will look something like this. So, I have simply concatenated these two audio features and video features and let us see how the dimension goes ok.

So, we are having these 32 files and each having 5040 dimensional feature and I will also printed these my audio label and video label just to show you that these are same with respect to the indexing. So, we can use them interchangeably. Now, let us try to use our classification code and see how my fused features are working.

So, for this I will simply reuse my code and instead of my video only features here I will use fused features and since our audio labels and video label are same with respect to the index, I can use either of the label again the test set size is 20 percent and let us see how it works.

As we can see here is these are working slightly lesser as compared to our video only feature in our first case which is our GaussianNB case and also for our LdA case the accuracy is dropping, but for the SVM case we are getting a similar accuracies as our video classifier video features are giving.

So, this basically is saying that feature fusion is giving us a better or at least the same amount of performance with our SVM classifier and the reason why these classification accuracy

decreases there could be multiple factors like one potential factor could be the curse of dimensionality since my feature dimension has too much high right now.

So, GaussianNBs and linear discriminant analysis might be suffering from this curse of dimensionality that is why they are giving a poor performance, but in case of a linear SVM I am getting a similar accuracy over here. So, till now we tried with traditional machine learning classifiers and now onwards we will be using 1 D CN architecture to classify the image and audio data encodings into emotion classes.

First, we will separately test the CNN architecture on audio and video data. Later on we will be using the fuse modality as the CNN input and in this case we are hypothesizing that the given audio and video embedding that we generated using MXCC and VGG face respectively are the two representation of the original data and since we have created these embeddings by separately concatenating each time frames we can treat these embedding as a time domain data and use 1D CNN classifier on top of it.

(Refer Slide Time: 30:31)



So, our CNN architecture will look something like this where we will be sequentially where we will be passing the input shape and there will be two convolution layers each convolution layer followed by a max pooling operation then we will flatten the output from the second max pool layer and we will pass through the flattened output through batch normalization layer.

And later we will collect and later we will be using 128 dimensional dense layer followed by a dropout layer for decreasing any chance of over fitting and later with the last dense layer which will be consisting of two neurons using softmax activation we will classify our given data into the emotion class.

Also, we will be using categorical cross entropy as a loss function and the optimizer we will be using is adam optimizer with learning rate of 0.0001 and the metric will be accuracy.

So, let me run this function. So, as our classifier is ready, we just need to reshape our video and mfcc audio features in a way that can be input into this CNN architecture. So, for that I will be using simple reshaping operations and our corresponding data will look something like this for both video and mfcc features.

(Refer Slide Time: 32:14)



Later I will simply train our model and for that I will be first converting our labels into categorical labels as we are using categorical cross entropy function as our loss function and we will divide our data into respective train and test set with 20 percent test size and we will

simply fit our model with batch size equal to 8 and epoch equal to 10. Let us see how it performs.

So, here you can see the summary of our given model and our model ok its already trained.

(Refer Slide Time: 32:50)



So, we can see over here that our CNN architecture is giving a good training accuracy, but the test accuracy is somewhat around 42 percent which is essentially below chance level and this thing is happening with video features only ok. So, let me try the same thing with audio features and let us see how my CNN architecture is working with audio features.

So, in case of audio features we will again divide our audio data into their respective train and test splits with the test size of 20 percent and we will call our model. We will call a fresh instance of our model and then we will fit it on our training data with a batch size of 8 and

number of epoch equal to say 10 and they will try to evaluate it on our test data let us see how it performs ok.

(Refer Slide Time: 33:57)



So, as we can see that using these 1 D CNN are audio features are giving better test accuracy as compared to our video features. Now, we are interested to see how the fusion of these two feature will perform in 1 D CNN architecture. So, for that I will simply use the fused version of the feature and the code will look something like this.

(Refer Slide Time: 34:27)

Here I am using this fuse feature with the set size 20 percent and dividing into my respective filter sets calling a fresh instance of our CNN architecture and fitting our training data with again same batch size 8 and epoch equal to 10 ok.

(Refer Slide Time: 35:03)



And in this case, we can see that our model is able to learn at least 71 percent which is somewhat equivalent to the audio level features. Here we can see like fusing feature I am adding we can get at least the accuracy single test performing feature. And again, here though we have not done any sort of hyper parameter tuning and the data is also very less over here. So, I mean there is still a chance that we can get a better accuracy over here and this is a classic example of our early feature fusion.

So, what is typically happening here is we took two features concatenated it together and passed it to our network. So, after performing early fusion we will now perform late feature fusion where late feature fusion refers to a method where feature extracted from multiple sources such as different sensors or modalities are combined at a later stage in the processing pipeline.

After initial processing of these individual sources, this allows for more flexible and efficient feature representation as it enables the combinations of diverse and complementary information. In late fusion technique the feature are extracted from different sources are first processed independently often using different techniques or architectures and then combined in the later stage of the pipeline such as a fully connected layer.

It can be also be done using various techniques such as concatenation of two features coming from different modalities maybe their element wise additions or element wise multiplication or maybe some sort of a weighted averaging of two features. So, in our case we will simply take our video features and audio features.

And pass them to a different CNN architectures and from the feature representation layer of their respective CNN architecture, we will collect the late features and then we will manually concatenate them to a late feature embedding and using that embedding we will obtain a multi-layer perceptron and see the efficacy of our late fusion technique.

So, to do so, we will first create a model and this will be called as late-CNN model and the agenda of this model will be to extract relevant features from this feature representation layer which essentially consist of 128 neurons and the activation function is relu over here. So, running this code. So, after that we can simply now call our late CNN layer as a audio model and fit it on our mfcc audio only features with a batch size of 8 and the number of epochs will be 10.

After we fit the model, we will again make our feature extractor and extract the features from this feature representation layer and the code will look something like this.

(Refer Slide Time: 39:03)



Now, same thing we will perform with our video features.

And then now since we have our audio and video model train, we will simply extract their feature layer using this code and save them into a separate variables. So, if I show the shape of this variable, it will it will be a tensor of 32 cross 128. So, basically 32 are number of our sample and 128 is our limited representation of that size of that latent representation. Similarly, for the video we have the similar shape.

Now, we can simply concatenate these two features these two late representation and make a late fused latent representation. For that our code will look something like this. Now, we have this late representation we will create an another multi-layer perceptron and we will train that perceptron on our training data which essentially will be coming from these late fused features and test on the respective test data and we will see how this late fusion technique will work.

So, our model will look something like this it is a very simple model where we will simply flatten down the these features and then we perform a batch normalization and there will be two tens layer consisting of 16 neurons and relu activation function and I will be putting a drop-down layer of with 10 probability to regularize this network and later we will classify the emotion classes using softmax classifier.

And I will again training this model with categorical cross entropy loss function with adam optimizer the at the learning rate of 0.0001 percent.

(Refer Slide Time: 41:47)



Now, a model is compiled I can simply divide my late fused data into train and test splits and the code will look something like this and our test size is again 20 percent and now I will simply run my model and evaluate it on test set let us see how does it perform now ok.

(Refer Slide Time: 42:19)



So, yeah as we can see here we are getting better accuracy in this late fusion case as compared to our early fusion. In late fusion we are getting test accuracy of 85 percent whereas, in our early fusion case our test accuracies was 71 percent. So, in line with the literature we were having two different sort of modalities we extracted relevant features of these two modalities using some sort of a machine learning algorithm in our case a convolution neural network.

And later we extracted the latent representation from these machine learning algorithm from the CNN and combined those latent representation and did of late fusion and it is giving a better accuracy than our earlier techniques. So, finally, concluding this tutorial in this tutorial we started with we started with working with audio and video information.

We extracted mfcc features from audio modality and we used a pre-trained VGG16 phase architecture to get features from our video data later using these features we. Firstly,

separately trained our machine learning models and saw how they are performing then we try to fuse the given embeddings and saw how these classified again performed later we did a CNN based approach where we first tested early feature fusion method and then we tested later feature fusion method.

And we saw that in case of data where there are multiple modalities coming these late fusion techniques works better.

Thank you.