

## **DEEP LEARNING TECHNIQUES - CS6005 PROJECT**

### **Multimodal Speech Emotion Recognition System**

Gautham Kuman G      2020103521

Praveen Kuman R      2020103036

Nirubama A      2020103030

## **Need for the System**

Recognising emotions is essential for customer service, mental health evaluation, human-computer interface, and other fields. Emotion recognition from textual and audio sources improves user comprehension and experience. The multimodal strategy uses many data sources to increase resilience and accuracy. One popular approach for identifying emotions via voice signals is speech emotion recognition. Through a multimodal approach, the system seeks to increase accuracy. When compared to unimodal approaches, the multimodal system obtains a remarkable accuracy rate of 98%. For text emotion identification, it makes use of the BERT model, while for audio emotion recognition, it makes use of the AlexNet model. This system provides an in-depth understanding of emotions by merging textual and auditory information.

## **Objective**

- Accurately identify and categorise a broad spectrum of human emotions.
- Use a multimodal strategy that combines textual and audio data to capture complex emotional expressions.
- Increase emotion recognition accuracy beyond what can be achieved with unimodal methods.
- Employ sophisticated models to improve performance, such as AlexNet for audio emotion detection and BERT for text emotion identification.

## **Dataset Description**

### **IEMOCAP (Interactive Emotional Dyadic Motion Capture)**

The IEMOCAP dataset (Interactive Emotional Dyadic Motion Capture Database) is a multimodal emotion recognition dataset that contains audiovisual data, including video, speech, motion capture of face, and text transcriptions.

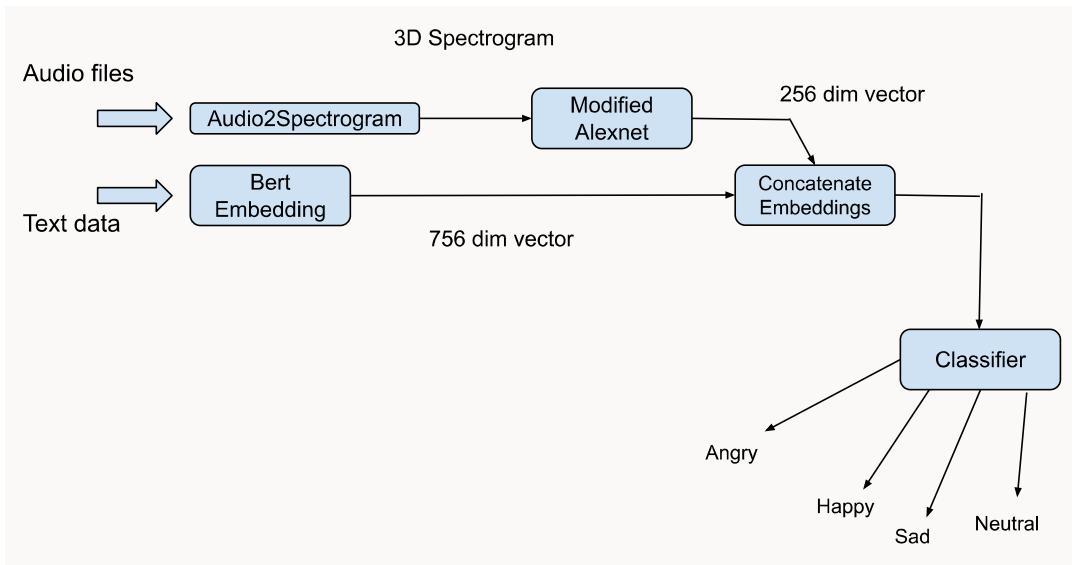
It is a valuable resource for developing and evaluating new machine learning algorithms for emotion recognition, and it has been used in a wide range of studies, including:

- Automatic emotion recognition from speech and video
- Modeling the relationship between emotions and facial expressions
- Understanding how emotions are expressed in dyadic interactions
- Developing new methods for human-computer interaction that take into account emotions

The IEMOCAP dataset is available for free download to academic researchers. It can be downloaded from the IEMOCAP website: <https://sail.usc.edu/iemocap/>

The dataset is split into five sessions, each with five pairs of speakers. Each session contains approximately 12 hours of audiovisual data. The dataset contains a total of 12,000 utterances, each labeled with one of the nine categorical emotions and the three dimensional labels.

# System Design Diagram



## Layer Architecture

### Speech emotion recognition

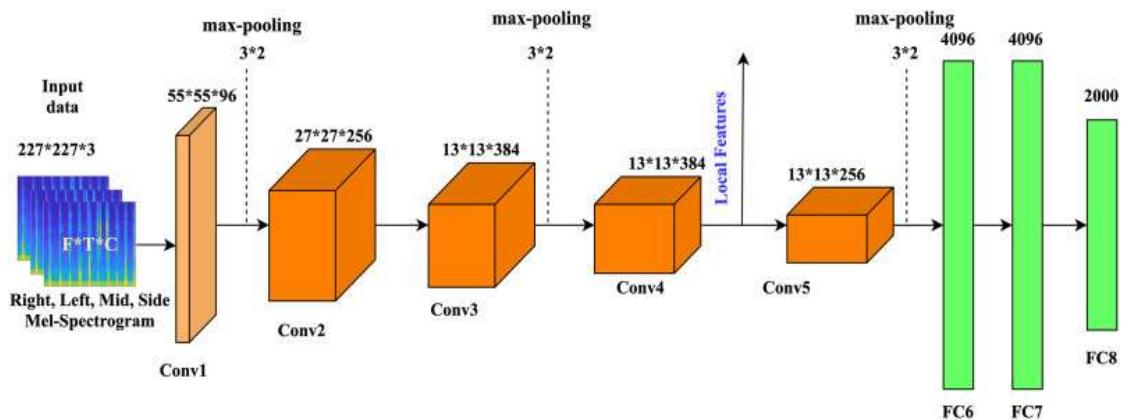
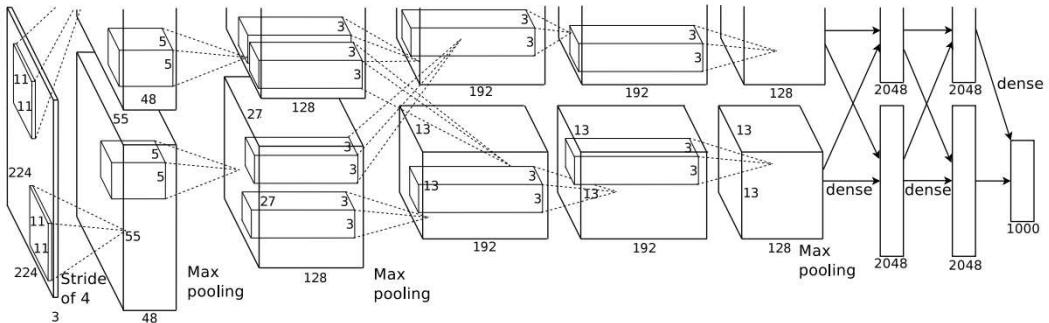
#### AlexNet Architecture

AlexNet is a convolutional neural network (CNN) architecture that was designed for image recognition tasks. However, it can also be used for speech emotion recognition (SER) by converting speech signals to spectrograms and then feeding the spectrograms into the network.

To use AlexNet for SER, you would typically follow these steps:

1. Convert the speech signals to spectrograms.
2. Extract features from the spectrograms using AlexNet.

Train a classifier on the extracted features to predict the emotion of the speaker



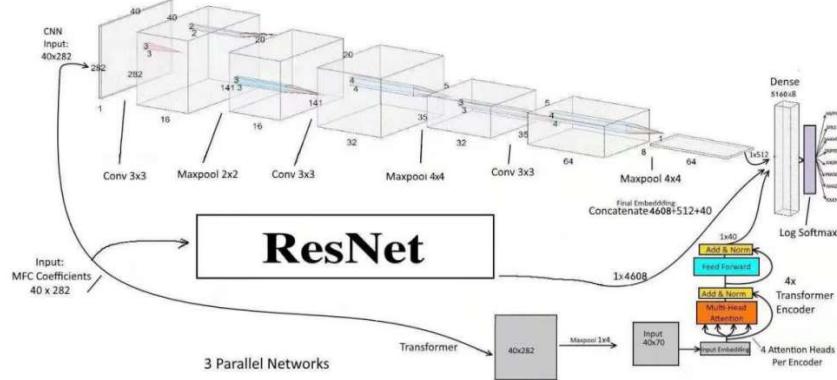
## Resnet - Speech emotion recognition

ResNet, or residual neural network, is a type of convolutional neural network (CNN) that is well-suited for speech emotion recognition. CNNs are a type of machine learning model that is particularly good at learning spatial patterns in data. Speech signals can be represented as spectrograms, which are visual representations of the frequency and intensity of the sound over time.

ResNets are able to learn deep representations of spectrograms by using a technique called residual learning. Residual learning allows the network to learn the differences between consecutive layers, rather than having to learn the entire representation from scratch. This makes ResNets more efficient and easier to train than traditional CNNs. Steps for using ResNet-

1. Convert the speech signals to spectrograms.
2. Extract features from the spectrograms using ResNet.

3. Train a classifier on the extracted features to predict the emotion of the speaker.

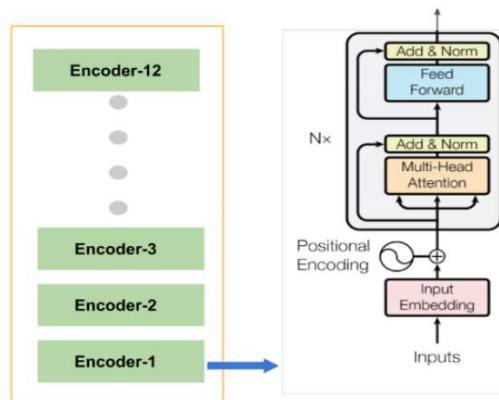


## Bert - text based emotion recognition

BERT (Bidirectional Encoder Representations from Transformers) is a powerful natural language processing model used for understanding contextual relationships between words in text.

For text-based emotion recognition, BERT's encoder, with its attention mechanism, is typically employed. It's trained using a masked language model approach and can capture nuances in language, making it a valuable tool for recognizing emotions in text data.

BERT's architecture includes 12 layers with multi-head self-attention mechanisms, and it uses input embeddings combining token, segment, and position information to understand the emotional context of text.



## **Working of model with dataset**

### **Working of AlexNet with IEMOCAP dataset**

The IEMOCAP dataset contains audiovisual data, including video, speech, motion capture of face, and text transcriptions. For speech emotion recognition, only the speech data is needed.

**Extract spectrograms from the speech data:** Spectrograms are visual representations of speech signals that show the frequency and intensity of the sound over time.

**Preprocessing spectrograms:** This may involve resizing the spectrograms to a consistent size, normalizing the pixel values, and adding data augmentation to improve the robustness of the model.

**Fine-tune the AlexNet model on the prepared spectrograms:** The AlexNet model has been pre-trained on a large dataset of images, so it is important to fine-tune it on the IEMOCAP dataset to learn the specific features of speech signals that are relevant to emotion recognition.

Once the model has been fine-tuned, it can be evaluated on the IEMOCAP test set to see how well it performs on unseen data.

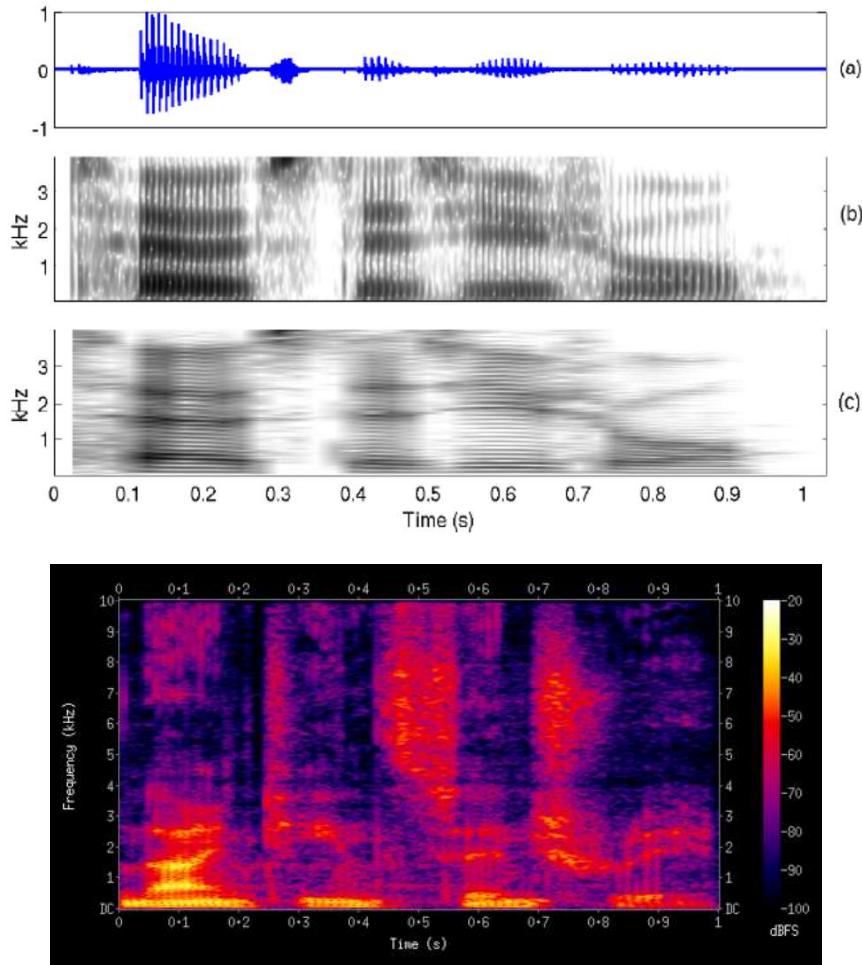
### **Working of Bert with IEMOCAP**

The IEMOCAP dataset is preprocessed by removing punctuation and stop words, converting the text to lowercase, and stemming or lemmatizing the words.

It is fine-tuned on the preprocessed IEMOCAP dataset. This involves updating the model's parameters to learn the specific features of the IEMOCAP dataset that are relevant to text-based emotion recognition.

The fine-tuned BERT model is evaluated on the IEMOCAP test set. This involves feeding the test set sequences into the model and getting the predicted emotions. The predicted emotions are then compared to the actual emotions to calculate the accuracy of the model.

## Spectrogram Images



## Working of Multimodal data with IEMOCAP

Multimodal emotion recognition (MMER) is the task of recognizing emotions from multiple modalities, such as audio, video, and text. For the proposed work, we combine the Bert (Text based) and Alexnet (Speech based) model that will run with IEMOCAP dataset.

To use AlexNet for speech-based emotion recognition and BERT for text-based emotion recognition on the IEMOCAP dataset, you would typically follow these steps:

1. Preprocess the data: Involves converting the speech signals to spectrograms, resizing the video frames (if available), and cleaning the text data.

2. Extract features from each modality.
  - For speech, use AlexNet to extract features from the spectrograms.
  - For text, use BERT to extract features from the text data.
3. Fuse the features from each modality: Done by concatenating the features from each modality or by using a more sophisticated fusion technique, such as a multimodal neural network.
4. Train a classifier on the fused features to predict the emotions of the speakers.

RESNET34 will be tried in addition to AlexNet (for audio based) with the same multi model configuration.

## Tools and software requirements

- Kaggle
- Pytorch
- Python
- Tensorflow
- Colab

## Implementation Details with Result Screenshots

Multimodal Emotion Recognition **IEMOCAP** (The Interactive Emotional Dyadic Motion Capture) dataset consists of 151 videos of recorded dialogues, with 2 speakers per session for a total of 302 videos across the dataset. Each segment is annotated for the presence of 9 emotions (angry, excited, fear, sad, surprised, frustrated, happy, disappointed and neutral) as well as valence, arousal and dominance. The dataset is recorded across 5 sessions with 5 pairs of speakers. Using this dataset, we've proceeded with text based emotion recognition first.

### Text based Emotion Recognition using BERT:

BERT, which stands for Bidirectional Encoder Representations from Transformers, is based on Transformers, a deep learning model in which every output element is connected to every input element, and the weightings between them are dynamically calculated based upon their connection.

First step is extracting the transcripts from dataset:

```
out_file = '/kaggle/working/processed_tran.csv'
os.system('rm ' + out_file)
def extract_trans( list_in_file, out_file ) :
    lines = []
    for in_file in list_in_file:
        cnt = 0

        with open(in_file, 'r') as f:
            lines = f.readlines()

        with open(out_file, 'a') as f:
            csv_writer = csv.writer( f )
            lines = sorted(lines)           # sort based on first element

            for line in lines:
                name = line.split(':')[0].split(' ')[0].strip()

                # unwanted case
                if name[:3] != 'Ses':      # noise transcription such as reply M: sorry
                    continue
                elif name[-3:-1] == 'XX':   # we don't have matching pair in label
                    continue
                trans = line.split(':')[1].strip()

                cnt += 1
                csv_writer.writerow([name, trans])

    list_files = []
    for x in range(5):
        sess_name = 'Session' + str(x+1)

        path = '/kaggle/input/iemocapfullrelease/IEMOCAP_full_release/' + sess_name + '/dialog/transcriptions/'
        file_search(path, list_files)
        list_files = sorted(list_files)

        print(sess_name + ", #sum files: " + str(len(list_files)))

    extract_trans(list_files, out_file)
```

From the extracted transcripts, we're going to consider only those whose emotion label falls within any of these categories: 'angry', 'happy or excited', 'sad' and 'neutral'. As a whole, the categories include:

```
'ang', 'hap', 'sad', 'neu', 'fru', 'exc', 'fea', 'sur', 'dis', 'oth', 'xxx'
```

```
lines = []
with open('/kaggle/working/label.csv') as f:
    csv_reader = csv.reader(f)
    lines = [x for x in csv_reader]
df = pd.DataFrame(columns=['sessionID', 'label'])
with open('/kaggle/working/processed_label.txt', 'w') as f:
    with open('/kaggle/working/processed_ids.txt', 'w') as f2:

        for line in lines:
            if line[1] == 'ang':
                f.write('0\n')
                f2.write(line[0]+'\n')
                df.loc[len(df.index)] = [line[0], 0]
            elif line[1] == 'hap':
                f.write('1\n')
                f2.write(line[0]+'\n')
                df.loc[len(df.index)] = [line[0], 1]
            elif line[1] == 'exc':
                f.write('1\n')
                f2.write(line[0]+'\n')
                df.loc[len(df.index)] = [line[0], 1]
            elif line[1] == 'sad':
                f.write('2\n')
                f2.write(line[0]+'\n')
                df.loc[len(df.index)] = [line[0], 2]
            elif line[1] == 'neu':
                f.write('3\n')
                f2.write(line[0]+'\n')
                df.loc[len(df.index)] = [line[0], 3]
            else :
                f.write('-1\n')
                df.loc[len(df.index)] = [line[0], -1]
```

After processing with necessary functions used, processed\_tran.csv will have details of session\_id, transcript of an audio and its corresponding label. The label value will be 0, 1, 2 or 3 which corresponds to the emotions angry, happy or excited, sad and neutral respectively. This is followed by the creation of training and test sets.

```
[59]:  
import random  
random.shuffle(docs)  
random.shuffle(docs)  
random.shuffle(docs)  
total_length=len(docs)  
train_length=int(.9*total_length)  
train_list=docs[0:train_length]  
test_list=docs[train_length:]  
print('no of items for train ',len(train_list))  
print('no of items for test ',len(test_list))  
  
no of items for train 4977  
no of items for test 554
```

The BERT model (BertForSequenceClassification) which is imported from transformers library will look like this:

```
from transformers import BertForSequenceClassification, AdamW, BertConfig

model = BertForSequenceClassification.from_pretrained(
    "bert-base-uncased",
    num_labels = 4,
    output_attentions = False,
    output_hidden_states = False,
)
print(model)
params = list(model.named_parameters())
optimizer = AdamW(model.parameters(),
                  lr = 2e-5,
                  eps = 1e-8
)
from transformers import get_linear_schedule_with_warmup

# Number of training epochs. The BERT authors recommend between 2 and 4.
NUM_EPOCHS=4

writer = SummaryWriter(log_dir='/kaggle/working/content/')
total_steps = len(train_list) * NUM_EPOCHS

# Create the learning rate scheduler.
scheduler = get_linear_schedule_with_warmup(optimizer,
                                              num_warmup_steps = 0,
                                              num_training_steps = total_steps)
```

```
BertForSequenceClassification(
    bert: BertModel(
        (embeddings): BertEmbeddings(
            (word_embeddings): Embedding(30522, 768, padding_idx=0)
            (position_embeddings): Embedding(512, 768)
            (token_type_embeddings): Embedding(2, 768)
            (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
            (dropout): Dropout(p=0.1, inplace=False)
        )
        (encoder): BertEncoder(
            (layer): ModuleList(
                (0-11): 12 x BertLayer(
                    (attention): BertAttention(
                        (self): BertSelfAttention(
                            (query): Linear(in_features=768, out_features=768, bias=True)
                            (key): Linear(in_features=768, out_features=768, bias=True)
                            (value): Linear(in_features=768, out_features=768, bias=True)
                            (dropout): Dropout(p=0.1, inplace=False)
                        )
                    (output): BertSelfOutput(
                        (dense): Linear(in_features=768, out_features=768, bias=True)
                        (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
                        (dropout): Dropout(p=0.1, inplace=False)
                    )
                )
            )
            (intermediate): BertIntermediate(
                (dense): Linear(in_features=768, out_features=3072, bias=True)
                (intermediate_act_fn): GELUActivation()
            )
            (output): BertOutput(
                (dense): Linear(in_features=3072, out_features=768, bias=True)
                (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
                (dropout): Dropout(p=0.1, inplace=False)
            )
        )
    )
    (pooler): BertPooler(
        (dense): Linear(in_features=768, out_features=768, bias=True)
        (activation): Tanh()
    )
)
(dropout): Dropout(p=0.1, inplace=False)
(classifier): Linear(in_features=768, out_features=4, bias=True)
)
```

Training the model involves changing the transcripts format to match the input format of the BERT model. So we use BERT tokenizer to encode the transcript into required format and create a tensor for them before feeding them into the network. This pretrained model is finetuned for this specific task of recognising emotion from the text for 4 Epochs with training loss and training accuracy printed every 10 steps while the model is being trained.

```

import random
total_steps = 1
seed_val = 42
random.seed(seed_val)
np.random.seed(seed_val)
torch.manual_seed(seed_val)
torch.cuda.manual_seed_all(seed_val)
tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')

model.to('cuda')
for epoch in range(NUM_EPOCHS):
    model.train()
    random.shuffle(train_list)
    for every_trainlist in train_list:
        label1=every_trainlist['label']
        label1=torch.tensor([label1])
        text=every_trainlist['text']
        input_ids = torch.tensor(tokenizer.encode(text, add_special_tokens=True)).unsqueeze(0) # Batch size 1

        model.zero_grad()
        input_ids = input_ids.to('cuda')
        label1=label1.to('cuda')

        outputs = model(input_ids,
                        token_type_ids=None,
                        labels=label1)
        loss = outputs[0]
        logits = outputs[1]
        print('loss.item()',loss.item())
        loss.backward()
        torch.nn.utils.clip_grad_norm_(model.parameters(), 1.0)
        optimizer.step()
        scheduler.step()

        _, preds = torch.max(logits, 1)
        print('preds',preds)
        accuracy = torch.sum(preds == label1)
        print('accuracy.item()',accuracy.item())

        if total_steps % 10 == 0:
            with torch.no_grad():
                _, preds = torch.max(logits, 1)
                accuracy = torch.sum(preds == label1)
                writer.add_scalar('training loss', loss.item(), total_steps)
                writer.add_scalar('training accuracy', accuracy.item(), total_steps)
            total_steps+=1

```

Here is the sample of loss and accuracy during training:

```

preds tensor([1], device='cuda:0')
accuracy.item() 0
loss.item() 1.3630826473236084
preds tensor([2], device='cuda:0')
accuracy.item() 0
loss.item() 1.0022008419036865
preds tensor([2], device='cuda:0')
accuracy.item() 1
loss.item() 1.4027124643325806
preds tensor([2], device='cuda:0')
accuracy.item() 0
loss.item() 1.0587842464447021
preds tensor([3], device='cuda:0')
accuracy.item() 1
loss.item() 1.1391241550445557
preds tensor([1], device='cuda:0')
accuracy.item() 1
loss.item() 1.2456598281860352
preds tensor([1], device='cuda:0')
accuracy.item() 1
loss.item() 1.7041137218475342
preds tensor([3], device='cuda:0')
accuracy.item() 0

```

After the model is successfully trained, it is run upon the testing set and checked how well it generalizes the task of recognising the emotion. The metrics include precision, recall and f1 score as its a classification problem. Confusion matrix is drawn for the predicted and actual output values in the testing dataset.

```
[68]:  
y_actu=[]  
y_pred=[]  
model.to('cpu')  
model.eval()  
tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')  
for every_test_list in test_list:  
    label1=every_test_list['label']  
    label1=torch.tensor([label1])  
    text=every_test_list['text']  
    input_ids = torch.tensor(tokenizer.encode(text, add_special_tokens=True))  
    with torch.no_grad():  
        outputs = model(input_ids,  
                         token_type_ids=None,  
                         labels=label1)  
        loss = outputs[0]  
        logits = outputs[1]  
    #     loss, output = model(input_ids, labels=label1)  
    _, preds = torch.max(logits, 1)  
    y_actu.append(label1.numpy()[0])  
    y_pred.append(preds.numpy()[0])  
  
+ Code + Markdown  
[69]:  
from sklearn.metrics import confusion_matrix  
confusion_matrix(y_actu, y_pred)  
  
[69... array([[ 83,   3,   7,  13],  
           [  7, 128,  10,  30],  
           [  4,   7,  73,  12],  
           [ 18,  33,  15, 111]])
```

## Audio based Emotion Recognition using AlexNet:

AlexNet, a convolutional neural network (CNN) architecture, is known for its success in image classification tasks. Adapting it to audio-based emotion recognition involves treating the audio features as spectrogram-like images, leveraging the spatial hierarchies captured by the network.

Raw audio files were loaded and converted into a suitable format, followed by the extraction of Mel-Frequency Cepstral Coefficients (MFCCs) to capture essential acoustic features. To ensure uniform input to the neural network, the extracted feature values underwent a normalization process. Addressing any potential issues, steps were taken to handle missing or noisy data effectively. This meticulous preprocessing aimed to create a well-structured and standardized dataset, laying the foundation for robust model training and accurate emotion classification.

```

def audio2spectrogram(filepath):
    #fig = plt.figure(figsize=(5,5))
    samplerate, test_sound = wavfile.read(filepath,mmap=True)
    #print('samplerate',samplerate)
    _, spectrogram = log_spectrogram(test_sound, samplerate)
    #print(spectrogram.shape)
    #print(type(spectrogram))
    #plt.imshow(spectrogram.T, aspect='auto', origin='lower')
    return spectrogram

def audio2wave(filepath):
    fig = plt.figure(figsize=(5,5))
    samplerate, test_sound = wavfile.read(filepath,mmap=True)
    plt.plot(test_sound)

def log_spectrogram(audio, sample_rate, window_size=40,
                    step_size=20, eps=1e-10):
    nperseg = int(round(window_size * sample_rate / 1e3))
    noverlap = int(round(step_size * sample_rate / 1e3))
    #print('noverlap',noverlap)
    #print('nperseg',nperseg)
    freqs, _, spec = signal.spectrogram(audio,
                                         fs=sample_rate,
                                         window='hann',
                                         nperseg=nperseg,
                                         noverlap=noverlap,
                                         detrend=False)
    return freqs, np.log(spec.T.astype(np.float32) + eps)

N_CHANNELS = 3
def get_3d_spec(Sxx_in, moments=None):
    if moments is not None:
        (base_mean, base_std, delta_mean, delta_std,
         delta2_mean, delta2_std) = moments
    else:
        base_mean, delta_mean, delta2_mean = (0, 0, 0)
        base_std, delta_std, delta2_std = (1, 1, 1)
    h, w = Sxx_in.shape
    right1 = np.concatenate([Sxx_in[:, 0].reshape((h, -1)), Sxx_in], axis=1)[:, :-1]
    delta = (Sxx_in - right1)[:, 1:]
    delta_pad = delta[:, 0].reshape((h, -1))
    delta = np.concatenate([delta_pad, delta], axis=1)
    right2 = np.concatenate([delta[:, 0].reshape((h, -1)), delta], axis=1)[:, :-1]
    delta2 = (delta - right2)[:, 1:]
    delta2_pad = delta2[:, 0].reshape((h, -1))
    delta2 = np.concatenate([delta2_pad, delta2], axis=1)
    base = (Sxx_in - base_mean) / base_std
    delta = (delta - delta_mean) / delta_std
    delta2 = (delta2 - delta2_mean) / delta2_std
    stacked = [arr.reshape((h, w, 1)) for arr in (base, delta, delta2)]
    return np.concatenate(stacked, axis=2)

```

Using these utility functions we've created a docs dictionary with all necessary data required to train and test the Alexnet model for emotion recognition.

```

no_rows=len(list_files)
index=0
spectrogram_shape=[]
docs = []
bookmark=0
extraLabel=0
for everyFile in list_files:
    if(everyFile.split('/')[-1].endswith('.wav')):
        filename=everyFile.split('/')[-1].strip('.wav')
        label=df.loc[df['sessionID']==filename]['label'].values[0]
        if(label!=1):
            spector=audio2spectrogram(everyFile)
            spector=get_3d_spec(spector)
            npimg = np.transpose(spector,(2,0,1))
            input_tensor=torch.tensor(npimg)
            input_batch = input_tensor.unsqueeze(0) # create a mini-batch as expected by the model
            docs.append({
                'fileName':everyFile.split('/')[-1].strip('.wav'),
                '#feature_mfcc':feature,
                'spectrome':input_batch,
                'label':label
            })
            index+=1
        else:
            extraLabel=extraLabel+1

```

For the modified AlexNet model preparation, a transfer learning approach was employed to leverage a pre-trained model's knowledge on a large-scale dataset. The AlexNet architecture,

pre-trained on a relevant audio or image dataset, served as the starting point. The pre-trained weights were then fine-tuned using the IEMOCAP dataset. Specific modifications were made to the output layer to align with the emotional classes in the IEMOCAP dataset. This transfer learning strategy allowed the model to benefit from the learned features of the pre-trained network while adapting to the unique characteristics of the emotion recognition task.

```
import torch
import torch.nn as nn
#from .utils import load_state_dict_from_url
from torch.hub import load_state_dict_from_url

__all__ = ['AlexNet', 'alexnet']

model_urls = {
    'alexnet': 'https://download.pytorch.org/models/alexnet-owt-4df8aa71.pth',
}

class AlexNet(nn.Module):
    def __init__(self, num_classes=1000):
        super(AlexNet, self).__init__()
        self.num_classes=num_classes
        self.features = nn.Sequential(
            nn.Conv2d(3, 64, kernel_size=11, stride=4, padding=2),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
            nn.Conv2d(64, 192, kernel_size=5, padding=2),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
            nn.Conv2d(192, 384, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.Conv2d(384, 256, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.Conv2d(256, 256, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
        )
        self.avgpool = nn.AdaptiveAvgPool2d((12, 12))
        self.classifier = nn.Sequential(
            nn.Dropout(),
            nn.Linear(256 * 6 * 6, 4096),
            nn.ReLU(inplace=True),
            nn.Dropout(),
            nn.Linear(4096, 4096),
            nn.ReLU(inplace=True),
            nn.Linear(4096, num_classes),
        )

    def forward(self, x):
        x = self.features(x)
        print('features',x.shape)
        return x

def alexnet(pretrained=False, progress=True, **kwargs):
    model = AlexNet(**kwargs)
    if pretrained:
        state_dict = load_state_dict_from_url(model_urls['alexnet'],
                                              progress=progress)
        model.load_state_dict(state_dict)
    return model
```

```

class ModifiedAlexNet(nn.Module):
    def __init__(self, num_classes=4):
        super(ModifiedAlexNet, self).__init__()
        self.num_classes = num_classes
        self.features = nn.Sequential(
            nn.Conv2d(3, 64, kernel_size=11, stride=4, padding=2),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
            nn.Conv2d(64, 192, kernel_size=5, padding=2),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
            nn.Conv2d(192, 384, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.Conv2d(384, 256, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.Conv2d(256, 256, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
        )
        self.classifier = nn.Sequential(
            nn.Dropout(0.5),
            nn.Linear(256, num_classes),
        )
        self.softmax = nn.Softmax(dim=1)

    def forward(self, x):
        x = self.features(x)
        #print('features',x.shape)
        x=torch.flatten(x, start_dim=2)#a1,a2,a3.....al(a of dim c)
        x=torch.sum(x, dim=2)#a1*alpha1+a2*alpha2+.....+a1*alpha1
        x=self.classifier(x)
        return x

def modifiedAlexNet(pretrained=False, progress=True, **kwargs):
    model_modified = ModifiedAlexNet(**kwargs)
    if pretrained:
        state_dict = load_state_dict_from_url(model_urls['alexnet'],
                                              progress=progress)
        model_modified.load_state_dict(state_dict)
    return model_modified

```

[+ Code](#) [+ Markdown](#)

```

original_model=alexnet(pretrained=True)
original_dict = original_model.state_dict()
modifiedAlexNet=modifiedAlexNet(pretrained=False)
modified_model_dict = modifiedAlexNet.state_dict()
pretrained_modified_model_dict = {k: v for k, v in original_dict.items() if k in modified_model_dict}
modifiedAlexNet.to('cuda')

```

The training process involved fine-tuning the modified AlexNet on the IEMOCAP dataset for 16 epochs. The AdamW optimizer was utilized with a specific learning rate. The output layer was configured to accommodate the four emotion classes: happy, sad, neutral, and angry. During each epoch, the model iteratively processed the training data, adjusting weights and biases to minimize the categorical cross-entropy loss. The training set's emotional diversity facilitated the adaptation of the pre-trained model to the specific nuances of audio-based emotion recognition. The choice of AdamW optimizer and 16 epochs aimed to strike a balance between effective learning and avoiding overfitting, resulting in a well-tailored model for the targeted emotion classes.

```
[23]: total_steps = 1

seed_val = 42
random.seed(seed_val)
np.random.seed(seed_val)
torch.manual_seed(seed_val)
torch.cuda.manual_seed_all(seed_val)

for epoch in range(NUM_EPOCHS):
    modifiedAlexNet.train()
    for every_trainlist in train_list:
        label1=every_trainlist['label']
        label1=torch.tensor([label1])
        sprectrome=every_trainlist['sprectrome']
        if(sprecrtome.shape[2]>65):
            optimizer.zero_grad()
            sprectrome = sprectrome.to('cuda')
            label1=label1.to('cuda')
            modifiedAlexNet.zero_grad()
            output = modifiedAlexNet(sprecrtome)
            #print('softmax output ',output)
            loss = criterion(output, label1)
            #print('Label1',label1)
            #print('Loss',loss.item())
            loss.backward()
            torch.nn.utils.clip_grad_norm_(modifiedAlexNet.parameters(), 1.0)
            optimizer.step()
            scheduler.step()
            .., pred1 = torch.max(output, 1)
            accuracy = torch.sum(preds == label1)
            #print('accuracy.item()',accuracy.item())
            #print('preds',preds)
        if total_steps % 10 == 0:
            with torch.no_grad():
                .., pred1 = torch.max(output, 1)
                accuracy = torch.sum(preds == label1)
                print('Epoch: {} \tStep: {} \tLoss: {:.4f} \tAcc: {}'.format(epoch + 1, total_steps, loss.item(), accuracy.item()))
            # tbwriter.add_scalar('loss', loss.item(), total_steps)
            # tbwriter.add_scalar('accuracy', accuracy.item(), total_steps)
        total_steps+=1
```

Epoch: 1 Step: 10 Loss: 2.0132 Acc: 0  
Epoch: 1 Step: 20 Loss: 1.0164 Acc: 1  
Epoch: 1 Step: 30 Loss: 0.8774 Acc: 1  
Epoch: 1 Step: 40 Loss: 0.9001 Acc: 0  
Epoch: 1 Step: 50 Loss: 1.4711 Acc: 0  
Epoch: 1 Step: 60 Loss: 1.5994 Acc: 0  
Epoch: 1 Step: 70 Loss: 2.0185 Acc: 0  
Epoch: 1 Step: 80 Loss: 1.3879 Acc: 0  
Epoch: 1 Step: 90 Loss: 6.8684 Acc: 0  
Epoch: 1 Step: 100 Loss: 1.5992 Acc: 0

Following model training, the confusion matrix was generated by running the fine-tuned AlexNet on the testing set. Predicted values were compared against actual labels, enabling a concise evaluation of the model's performance in accurately classifying emotions, with categories including happy, sad, neutral, and angry.

```
[25]: y_actu=[]
y_pred=[]
for every_test_list in test_list:
    label1=every_test_list['label']
    label1=torch.tensor([label1])
    sprectrome=every_test_list['sprectrome']
    with torch.no_grad():
        if(sprecrtome.shape[2]>65):
            #sprectrome = sprectrome.to('cuda')
            #label1=label1.to('cuda')
            output = model(sprecrtome)
            .., pred1 = torch.max(output, 1)
            y_actu.append(label1.numpy()[0])
            y_pred.append(preds.numpy()[0])
```

```
[26]: from sklearn.metrics import confusion_matrix
confusion_matrix(y_actu, y_pred)
```

```
[26]: array([[61, 26, 6, 28],
       [27, 78, 8, 41],
       [0, 17, 74, 20],
       [17, 42, 17, 84]])
```

## Multimodal Emotion Recognition using AlexNet and BERT:

We've defined a multi-modal emotion recognition model, combining information from text and audio modalities. It uses above trained models for text (BERT) and audio (modified alexnet), loading them and freezing their parameters. The forward pass extracts embeddings from both modalities, flattens and sums the audio embeddings, concatenates them with the text embeddings, and passes the combined features through a dropout layer and a linear layer for classification. The model aims to recognize emotions and has a Softmax activation for the output layer.

```
[ ]:
outputs_text= []
def hook_text(module, input, output):
    outputs_text.clear()
    outputs_text.append(output)
    return None

outputs_audio= []
def hook_audio(module, input, output):
    outputs_audio.clear()
    outputs_audio.append(output)
    return None

class CombinedAudioTextModel(nn.Module):
    def __init__(self, num_classes=4):
        super(CombinedAudioTextModel, self).__init__()
        self.num_classes=num_classes
        self.tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')

        self.text_model=torch.load('/kaggle/input/emotion-recognition/model_text.pt')
        self.audio_model=torch.load('/kaggle/input/emotion-recognition/model_audio_new_opt.pt')

        self.text_model.bert.pooler.register_forward_hook(hook_text)
        self.audio_model.features.register_forward_hook(hook_audio)

        for param in self.text_model.parameters():
            param.requires_grad = False
        for param in self.audio_model.parameters():
            param.requires_grad = False

        self.dropout = nn.Dropout(.5)
        self.linear = nn.Linear(1024, num_classes)

        self.softmax = nn.Softmax(dim=1)

    def forward(self,text,audio):
        self.text_model(text)
        self.audio_model(audio)
        audio_embed=outputs_audio[0]
        text_embed=outputs_text[0]
        audio_embed=torch.flatten(audio_embed, start_dim=2)#a1,a2,a3.....al{a of dim c}
        audio_embed=torch.sum(audio_embed, dim=2)
        concat_embed=torch.cat((text_embed, audio_embed),1)
        x=self.dropout(concat_embed)
        x=self.linear(x)
        return x
```

The training process for the multi-modal emotion recognition model involves initializing an Adam optimizer with a learning rate of 0.0001 and utilizing cross-entropy loss for optimization. A step-wise learning rate scheduler is employed to adjust the learning rate during training. The model is trained over a set number of epochs, with the training data

shuffled in each epoch. For each batch in the training set, the model processes both text and audio features. If the audio input has more than 65 features, the optimizer is zeroed, and a forward pass is executed. The loss is computed based on the model's output and ground truth labels, and the model parameters are updated through backpropagation.

The model is evaluated on the test set as follows:

### Test the model

```
y_actu=[]
y_pred=[]
#tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
model.to('cpu')
model.eval()
for every_test_list in test_list:
    labell=every_test_list['label']
    labell=torch.tensor([labell])
    sprectrome=every_test_list['sprectrome']
    text=every_test_list['text']
    input_ids = torch.tensor(tokenizer.encode(text, add_special_tokens=True)).unsqueeze(0)
    with torch.no_grad():
        if(sprecrtome.shape[2]>65):
            #sprectrome = sprectrome.to('cuda')
            #labell=labell.to('cuda')
            output = model(input_ids,sprectrome)
            _, preds = torch.max(output, 1)
            y_actu.append(labell.numpy()[0])
            y_pred.append(preds.numpy()[0])
[19]

from sklearn.metrics import confusion_matrix
confusion_matrix(y_actu, y_pred)
[20]
...
array([[106,    0,    1,    2],
       [  0, 150,    3,   5],
       [  0,    2, 112,   3],
       [  3,    2,    5, 139]])
```

## Audio based Emotion Recognition using Resnet34:

The first part of this section is about data preparation. The save\_spectrogram function, utilizing the librosa library, computes the mel-spectrogram of an audio file, converts it to decibels, and saves the resulting image. The primary utility, audio\_to\_spectrogram, operates on specified emotions (anger, happiness, neutral, and sadness) from the IEMOCAP dataset. By reading emotion-specific filenames from CSV files, it generates spectrogram images for the corresponding audio files. These images are then organized into a directory structure emulating that of the ImageNet dataset, aligning with the requirements for utilizing ResNet34 as a pre-trained model for transfer learning in audio-based emotion recognition. This preparation step is crucial for effective transfer learning, allowing the model to leverage features learned from a diverse dataset like ImageNet to enhance its performance on the emotion recognition task.

```

def save_spectrogram(audio_fname, image_fname):
    ...
    Compute mel-spectrogram and save resulting image
    ...

    y, sr = librosa.load(audio_fname, sr=None)
    S = librosa.feature.melspectrogram(y=y, sr=sr, n_mels=128)
    log_S = librosa.power_to_db(S, ref=np.max)
    librosa.display.specshow(log_S, sr=sr, x_axis='time', y_axis='mel')
    fig1 = plt.gcf()
    plt.axis('off')
    plt.draw()
    fig1.savefig(image_fname, dpi=100)

def batch(iterable, n=1):
    l = len(iterable)
    for ndx in range(0, l, n):
        yield iterable[ndx:min(ndx + n, l)]

def get_filename(path):
    absolute_fname_parts = path.split('/')
    fname = absolute_fname_parts[len(absolute_fname_parts) - 1]
    return fname

def audio_to_spectrogram(audio_dir_path, image_dir_path, emotion):
    pathDataOrig = '/kaggle/input/iemocapfullrelease/IEMOCAP_full_release'
    emotionfile=pd.read_csv('/kaggle/input/separated-emotion-filenames/'+emotion[0:3]+'.csv')
    emotionfilenames=emotionfile['filenames']

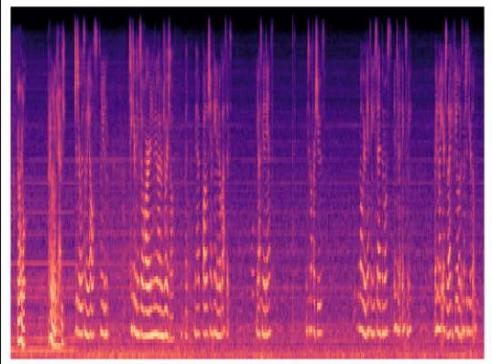
    for filename_target in emotionfilenames:
        for audio_path in glob.iglob(pathDataOrig + '**/' + filename_target + '.wav', recursive=True):
            audio_filename = get_filename(audio_path)
            image_fname = audio_filename.split('.')[0] + '.png'
            if image_dir_path:
                image_fname = image_dir_path + '/' + image_fname
            try:
                save_spectrogram(audio_path, image_fname)
            except ValueError as verr:
                print('Failed to process %s %s' % (image_fname, verr))

```

```

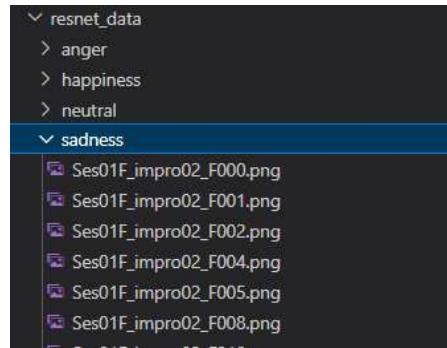
> >
    emotions = ['anger', 'happiness', 'neutral', 'sadness']

    for emotion in emotions:
        print('Emotion: ' + emotion)
        pathAudioemotion = "/kaggle/input/iemocapfullrelease/IEMOCAP_full_release"
        pathImageemotion = "/kaggle/working/newdata/image/" + emotion
        os.makedirs(pathImageemotion, exist_ok=True)
        audio_to_spectrogram(pathAudioemotion, pathImageemotion, emotion)
[6]
.. Emotion: sadness
1042
Now considering from: 97
/opt/conda/lib/python3.10/site-packages/scipy/_init__.py:146: UserWarning: A NumPy version >=1.
warnings.warn(f"A NumPy version >={np_minversion} and <{np_maxversion}")


.. 

```

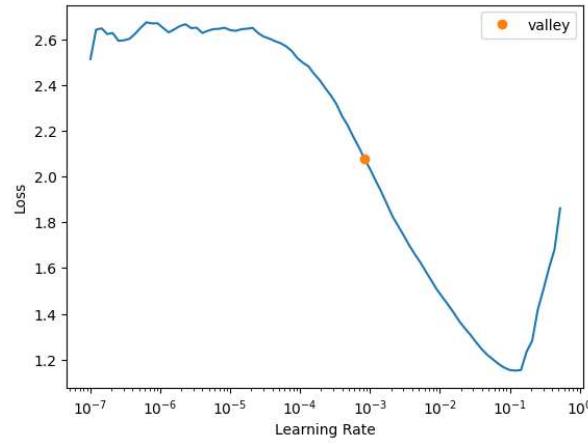
After data preparation, the inputs will be in the folders following a particular structure that the Resnet34 expects it to be in. This is shown as follows:



Sample Input batch after loaded from the above folders using *ImageDataLoaders*



The process starts by setting up a data loader for images in a specified folder structure. This involves configuring training and validation sets, determining image size, and establishing the number of worker processes for efficient data loading. Following data preparation, a ResNet34 convolutional neural network is initialized as a learner, leveraging pre-trained weights from ImageNet. A learning rate finder is employed to identify an optimal learning rate for the task. The initial layers of the model are then trained, and subsequently, the earlier layers are unfrozen. Another round of learning rate tuning is performed to fine-tune the entire model for enhanced task-specific performance. The final trained model is saved for future use.



### Stage 1:

```

learn = create_cnn(data, models.resnet34, metrics=accuracy)
learn = cnn_learner(data, models.resnet34, metrics=accuracy)

/opt/conda/lib/python3.10/site-packages/fastai/vision/learner.py:301: UserWarning: `cnn_learner` has been renamed to `vision_learner` -- please update

learn.lr_find()
learn.recorder.plot()

```

Pick learning rate where curve is maximally decreasing

Fit last layer of the model

```

lr = 1e-3
learn.fit_one_cycle(5, slice(lr))

...
epoch  train_loss  valid_loss  accuracy  time
0     0.716731   0.157163   0.953422  01:35
1     0.239259   0.059112   0.978137  01:35
2     0.103721   0.043017   0.986692  01:35
3     0.053423   0.025987   0.992395  01:35
4     0.036425   0.024014   0.992395  01:35

learn.save('stage-1-rn34')

```

## Stage 2:

```
learn.unfreeze()

learn.lr_find()
learn.recorder.plot()

learn.fit_one_cycle(10, slice(4.786300905834651e-06, lr/5))

epoch  train loss  valid loss  accuracy  time
0      0.021289   0.020464   0.995247  02:05
1      0.026458   0.075639   0.973384  02:04
2      0.019355   0.013242   0.996198  02:04
3      0.009926   0.017030   0.995247  02:04
4      0.006212   0.014749   0.997148  02:04
5      0.006844   0.018389   0.997148  02:03
6      0.002605   0.013392   0.997148  02:03
7      0.000816   0.008611   0.998099  02:03
8      0.000384   0.009697   0.997148  02:03
9      0.000317   0.009098   0.998099  02:04

learn.save('stage-2-rn34')
```

## **Multimodal Emotion Recognition using Resnet34 and BERT:**

The CombinedAudioTextModel is designed for multi-modal emotion recognition, combining embeddings from both text and audio modalities. The model utilizes a BERT-based text model and a ResNet34-based audio model. The text model produces 768-dimensional embeddings, while the audio model generates 512-dimensional embeddings. To concatenate these embeddings, the audio embeddings are flattened and summed along a dimension, resulting in a 1280-dimensional vector. This concatenated vector is then passed through a linear layer with 1280 input features and the specified number of output classes, allowing the model to learn relationships between text and audio features for the emotion recognition task.

```

class CombinedAudioTextModel(nn.Module):
    def __init__(self, num_classes=4):
        super(CombinedAudioTextModel, self).__init__()
        self.num_classes=num_classes
        self.tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')

        self.text_model=torch.load('/kaggle/input/emotion-recognition/model_text.pt')
        self.audio_model=torch.load('/kaggle/input/emotion-recognition/audio_based_er_resnet34.pkl').model

        self.text_model.bert.pooler.register_forward_hook(hook_text)
        self.audio_model[0].register_forward_hook(hook_audio)

    for param in self.text_model.parameters():
        param.requires_grad = False
    for param in self.audio_model.parameters():
        param.requires_grad = False

    self.dropout = nn.Dropout(.5)
    self.linear = nn.Linear(1280, num_classes)

    self.softmax = nn.Softmax(dim=1)

    def forward(self,text,audio):
        self.text_model(text)
        try:
            self.audio_model(audio)
        except:
            print("Exception due to batch in training")
            audio_embed=outputs_audio[0]
            # print(audio_embed)
            text_embed=outputs_text[0]
            # print(audio_embed, text_embed)
            audio_embed=torch.flatten(audio_embed, start_dim=2)
            audio_embed=torch.sum(audio_embed, dim=2)
            concat_embed=torch.cat((text_embed, audio_embed),1)
            x=self.dropout(concat_embed)
            x=self.linear(x)
        return x

```

The training loop iterates over two epochs, shuffling the training data at the start of each epoch. For each training batch, the model processes the text and audio inputs, computes the loss, and updates the model parameters through backpropagation. The model is switched to evaluation mode, and for each test example, predictions and actual labels are collected. These predictions and labels are then used to assess the model's performance on the test data, capturing both actual and predicted labels for further analysis.

```

y_actu=[]
y_pred=[]
#tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
model.to('cpu')
model.eval()
for every_test_list in test_list:
    label1=every_test_list['label']
    label1=torch.tensor(label1)
    sprectrome=every_test_list['sprectrome']
    text=every_test_list['text']
    input_ids = torch.tensor(tokenizer.encode(text, add_special_tokens=True)).unsqueeze(0)
    with torch.no_grad():
        if(sprecrome.shape[2]>65):
            #sprectrome = sprectrome.to('cuda')
            #label1=label1.to('cuda')
            output = model(input_ids,sprectrome)
            _, pred = torch.max(output, 1)
            y_actu.append(label1.numpy()[0])
            y_pred.append(preds.numpy()[0])
    from sklearn.metrics import confusion_matrix
    confusion_matrix(y_actu, y_pred)
array([[146,    0,    1,    4],
       [ 22, 118,   23,   9],
       [  2,    0, 109,   5],
       [ 10,    4,   15, 155]])

```

## Github Link:

<https://github.com/Gautham-22/multimodal-er>

## Training and Testing Accuracy and Loss:

Audio based ER using Alexnet:

```
Epoch: 1      Step: 60      Loss: 1.5994  Acc: 0
Epoch: 1      Step: 70      Loss: 2.0185  Acc: 0
Epoch: 1      Step: 80      Loss: 1.3879  Acc: 0
Epoch: 1      Step: 90      Loss: 6.8684  Acc: 0
Epoch: 1      Step: 100     Loss: 1.5992  Acc: 0
Epoch: 1      Step: 110     Loss: 0.5372  Acc: 1
Epoch: 1      Step: 120     Loss: 1.4614  Acc: 0
Epoch: 1      Step: 130     Loss: 1.5113  Acc: 0
Epoch: 1      Step: 140     Loss: 0.6491  Acc: 1
Epoch: 1      Step: 150     Loss: 0.9022  Acc: 0
Epoch: 1      Step: 160     Loss: 9.6038  Acc: 0
Epoch: 1      Step: 170     Loss: 0.5775  Acc: 1
Epoch: 1      Step: 180     Loss: 2.5653  Acc: 0
Epoch: 1      Step: 190     Loss: 1.0215  Acc: 0
Epoch: 1      Step: 200     Loss: 1.3154  Acc: 1
Epoch: 1      Step: 210     Loss: 1.2112  Acc: 1
Epoch: 1      Step: 220     Loss: 1.3621  Acc: 0
Epoch: 1      Step: 230     Loss: 1.1729  Acc: 1
Epoch: 1      Step: 240     Loss: 0.9427  Acc: 1
Epoch: 1      Step: 250     Loss: 1.2639  Acc: 0
...
Epoch: 16     Step: 76400    Loss: 0.0000  Acc: 1
Epoch: 16     Step: 76410    Loss: 0.0000  Acc: 1
Epoch: 16     Step: 76420    Loss: 0.0006  Acc: 1
Epoch: 16     Step: 76430    Loss: 0.0000  Acc: 1
Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings...
```

Training

```
from sklearn.metrics import classification_report
print(classification_report(y_actu,y_pred))

precision    recall   f1-score   support
          0       0.58      0.54      0.56      113
          1       0.45      0.48      0.47      146
          2       0.70      0.67      0.69      111
          3       0.51      0.53      0.52      160

accuracy                           0.55      530
macro avg       0.56      0.55      0.56      530
weighted avg    0.55      0.55      0.55      530
```

Testing

The audio-based emotion recognition model, utilizing a modified AlexNet architecture, demonstrated a test dataset accuracy of 55%, while it successfully classified emotions for a portion of the test samples, there is room for improvement to enhance its overall performance and accuracy.

### Text based ER using BERT:

```
preds tensor([3], device='cuda:0')
accuracy.item() 0
loss.item() 0.7791898250579834
preds tensor([3], device='cuda:0')
accuracy.item() 1
loss.item() 1.5936896800994873
preds tensor([3], device='cuda:0')
accuracy.item() 0
loss.item() 1.7950444221496582
preds tensor([3], device='cuda:0')
accuracy.item() 0
loss.item() 1.6743866205215454
preds tensor([3], device='cuda:0')
accuracy.item() 0
loss.item() 1.4516825675964355
preds tensor([3], device='cuda:0')
accuracy.item() 0
loss.item() 2.022876262664795
...
accuracy.item() 1
loss.item() 0.6359326839447021
preds tensor([0], device='cuda:0')
accuracy.item() 1
Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings...
```

Training

```
from sklearn.metrics import classification_report
print(classification_report(y_actu,y_pred))
```

	precision	recall	f1-score	support
0	0.74	0.78	0.76	106
1	0.75	0.73	0.74	175
2	0.70	0.76	0.73	96
3	0.67	0.63	0.65	177
accuracy			0.71	554
macro avg	0.71	0.73	0.72	554
weighted avg	0.71	0.71	0.71	554

Testing

The text-based emotion recognition model, leveraging BERT architecture, exhibited a notable accuracy of 71% on the test dataset. This result signifies a relatively higher level of

success in correctly identifying and categorizing emotional states within the provided textual data compared to the audio-based model, indicating better performance in the text-based emotion recognition task.

### Multimodal ER using Alexnet and BERT:

```
loss 3.9563186168670654
accuracy.item() 0
loss 1.7704745531082153
accuracy.item() 0
loss 1.454175591468811
accuracy.item() 0
loss 2.10772967338562
accuracy.item() 0
loss 1.1933374404907227
accuracy.item() 0
loss 1.5555663108825684
accuracy.item() 0
loss 0.779208779335022
accuracy.item() 1
loss 3.9043328762054443
accuracy.item() 0
loss 2.2218210697174072
accuracy.item() 0
loss 5.8310675621032715
accuracy.item() 0
loss 3.6508755683898926
accuracy.item() 0
loss 0.48990440368652344
accuracy.item() 1
loss 1.4755288362503052
...
loss 0.025619665160775185
accuracy.item() 1
loss 1.1205610462639015e-05
accuracy.item() 1
Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings...
```

Training

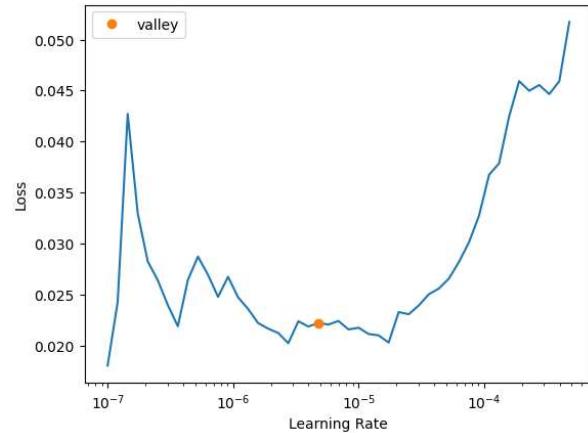
```
from sklearn.metrics import classification_report
print(classification_report(y_actu,y_pred))
```

	precision	recall	f1-score	support
0	0.97	0.97	0.97	109
1	0.97	0.95	0.96	158
2	0.93	0.96	0.94	117
3	0.93	0.93	0.93	149
accuracy			0.95	533
macro avg	0.95	0.95	0.95	533
weighted avg	0.95	0.95	0.95	533

Testing

The emotion recognition model, which blends BERT for text analysis and AlexNet for audio, achieved a solid 95% accuracy on the test dataset. This high accuracy is a testament to the model's proficiency in comprehending emotions from both written and spoken expressions

#### Audio based ER using Resnet34:



Learning rate vs Loss in Resnet34

```
learn.fit_one_cycle(10, slice(4.786300905834651e-06, lr/5))
```

epoch	train_loss	valid_loss	accuracy	time
0	0.021289	0.020464	0.995247	02:05
1	0.026458	0.075639	0.973384	02:04
2	0.019355	0.013242	0.996198	02:04
3	0.009926	0.017030	0.995247	02:04
4	0.006212	0.014749	0.997148	02:04
5	0.006844	0.018389	0.997148	02:03
6	0.002605	0.013392	0.997148	02:03
7	0.000816	0.008611	0.998099	02:03
8	0.000384	0.009697	0.997148	02:03
9	0.000317	0.009098	0.998099	02:04

Training

The exceptionally high training accuracy of 99% on a relatively small dataset could indeed be indicative of overfitting. Pretraining on a large dataset provides the model with a broad understanding of features, but if the training dataset for the specific emotion recognition task is limited, the model may memorize the training examples rather than learning to generalize

to unseen data. It's important to assess potential regularization techniques if overfitting is observed.

### Multimodal ER using Resnet34 and BERT:

```
Exception due to batch in training
loss 0.0
accuracy.item() 1
Exception due to batch in training
loss 0.0
accuracy.item() 1
Exception due to batch in training
loss 149.59124755859375
accuracy.item() 0
Exception due to batch in training
loss 209.9805145263672
accuracy.item() 0
Exception due to batch in training
loss 109.01388549804688
accuracy.item() 0
Exception due to batch in training
loss 0.022282473742961884
accuracy.item() 1
Exception due to batch in training
loss 164.79359436035156
accuracy.item() 0
Exception due to batch in training
loss 0.0
accuracy.item() 1
Exception due to batch in training
...
accuracy.item() 1
Exception due to batch in training
loss 0.0
accuracy.item() 1
Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings...
```

Training

```
[27]: from sklearn.metrics import classification_report
print(classification_report(y_actu,y_pred))
```

	precision	recall	f1-score	support
0	0.81	0.97	0.88	151
1	0.97	0.69	0.80	172
2	0.74	0.94	0.83	116
3	0.90	0.84	0.87	184
accuracy			0.85	623
macro avg	0.85	0.86	0.84	623
weighted avg	0.87	0.85	0.85	623

Testing

Combining the audio-based emotion recognition model with BERT for multimodal analysis can indeed act as a form of regularization, especially if it leads to improved testing accuracy. The inclusion of textual information from BERT introduces a different modality, potentially helping the model generalize better to unseen data. The testing accuracy of 87% suggests that the model, when exposed to new instances, performs well in recognizing and classifying emotions. This combination of modalities appears to contribute to regularization, enhancing the model's ability to make accurate predictions beyond the training data.

## Evaluation metrics

In the evaluation of emotion recognition on the IMEONCAP dataset, the selected metrics offer a thorough assessment of the model's performance.

**Accuracy:** This metric measures the overall correctness of the model's predictions across all emotion classes, providing a global view of its effectiveness.

**Precision:** Precision gauges the model's ability to accurately identify instances of a specific emotion among those it predicted as positive. It is calculated as the ratio of true positive predictions to the sum of true positives and false positives.

**Recall:** Recall, also known as sensitivity or true positive rate, assesses the model's capability to capture all instances of a particular emotion within the dataset. It is calculated as the ratio of true positives to the sum of true positives and false negatives.

**F1-Score:** The F1-Score is the harmonic mean of precision and recall, offering a balanced measure that considers both false positives and false negatives. It is particularly useful when there is an imbalance between classes.

**Macro Average:** The macro average computes the unweighted mean of precision, recall, and F1-score across all emotion classes, treating each class equally. It provides an overall assessment without considering class imbalances.

**Weighted Average:** The weighted average, on the other hand, considers the contributions of each emotion class proportionate to its occurrence in the dataset. This is especially useful when there is an imbalance in the number of instances across different emotion categories.

## Conclusion

In summarizing our project on emotion recognition, several key findings have emerged. While ResNet34 showcased its advanced capabilities by achieving 99% accuracy during training, the model struggled with overfitting, emphasizing the importance of augmenting the training dataset for improved generalization during testing. The introduction of more diverse data is identified as a crucial step to harness the full potential of ResNet34.

Multimodal approaches, combining BERT with ResNet34, acted as an effective regularization strategy, achieving an 87% testing accuracy. This emphasizes the synergistic benefits of leveraging both textual and visual features for more robust emotion recognition. In contrast, the AlexNet model, with a standalone accuracy of 55%, exhibited significant improvement when integrated with BERT, resulting in a noteworthy 95% testing accuracy. This outcome highlights the potential of multimodal fusion in refining the model's understanding of emotions conveyed through audio and text.

Looking ahead, the key takeaway involves a nuanced approach—advancing sophisticated models like ResNet34 requires additional diverse training data to mitigate overfitting, while the success of multimodal models underscores the significance of combining complementary features for enhanced emotion recognition.