

## DL Assignment3

Akshara B

211AI012

### **Dataset:**

I have used an audio dataset for the Multi-class classification task. The dataset consists of wav-format audio files which are of length 3-seconds. They contain the siren sound of Emergency Vehicles - Ambulance and Firetruck. A third category named Traffic also exists where it contains 3-second .wav format audio files of plain traffic sound. Each category contains 200 sound files. Total 600 files.

Link: <https://www.kaggle.com/datasets/vishnu0399/emergency-vehicle-siren-sounds/data>

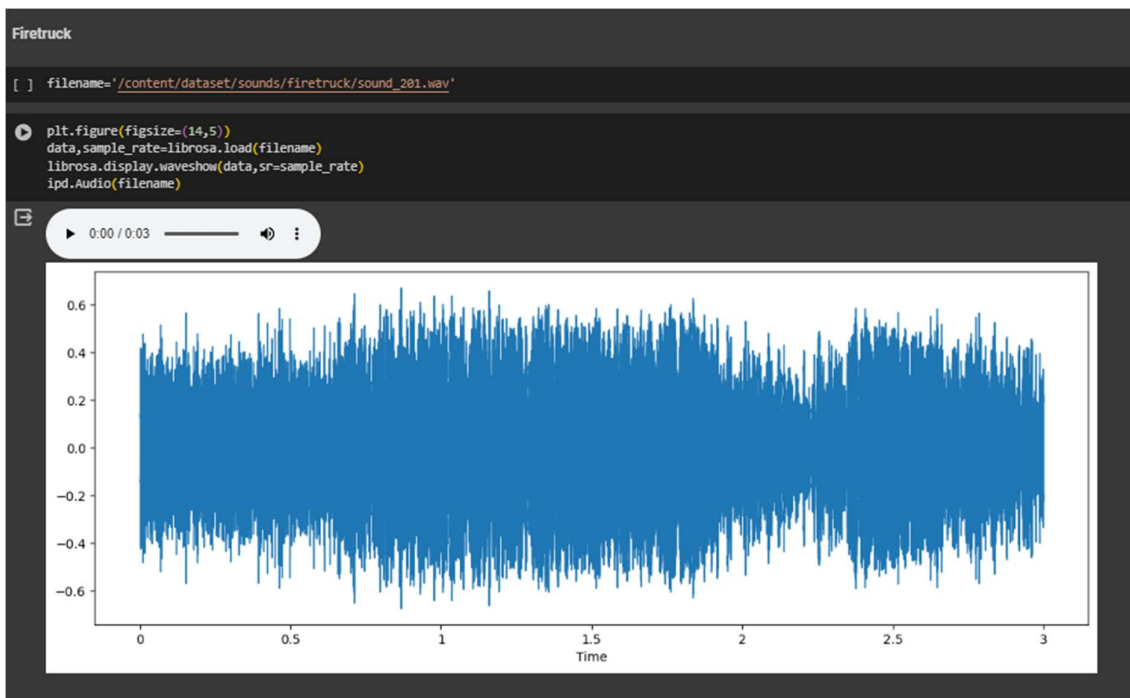
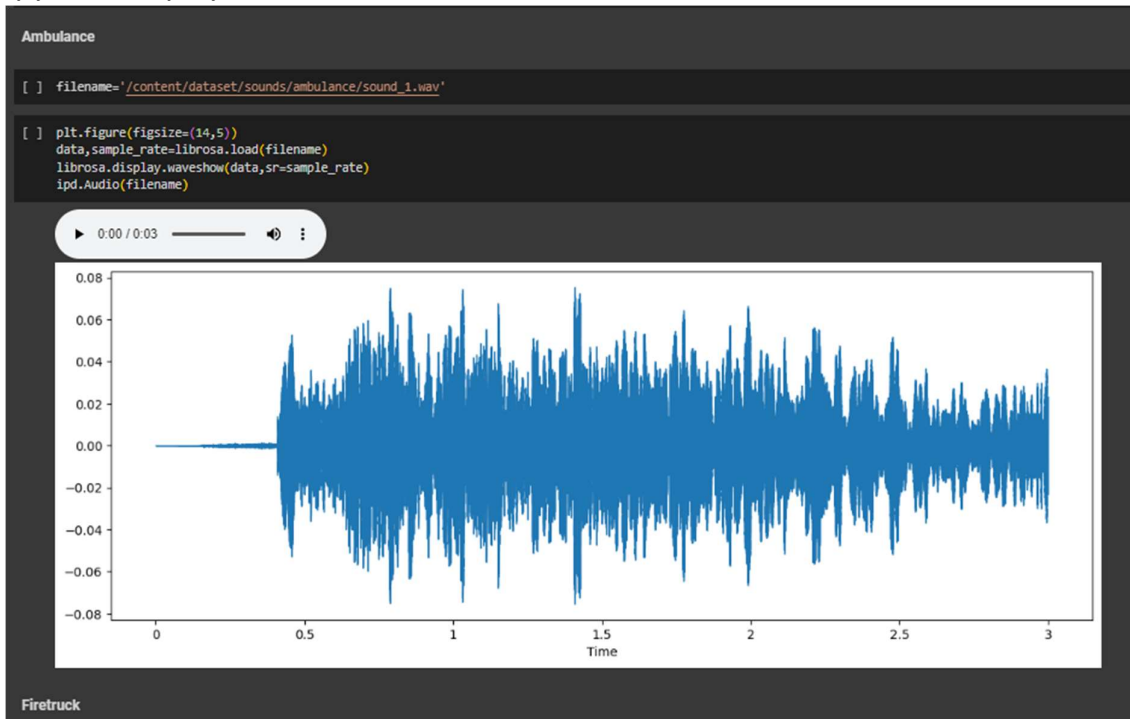
### **Google colab notebook:**

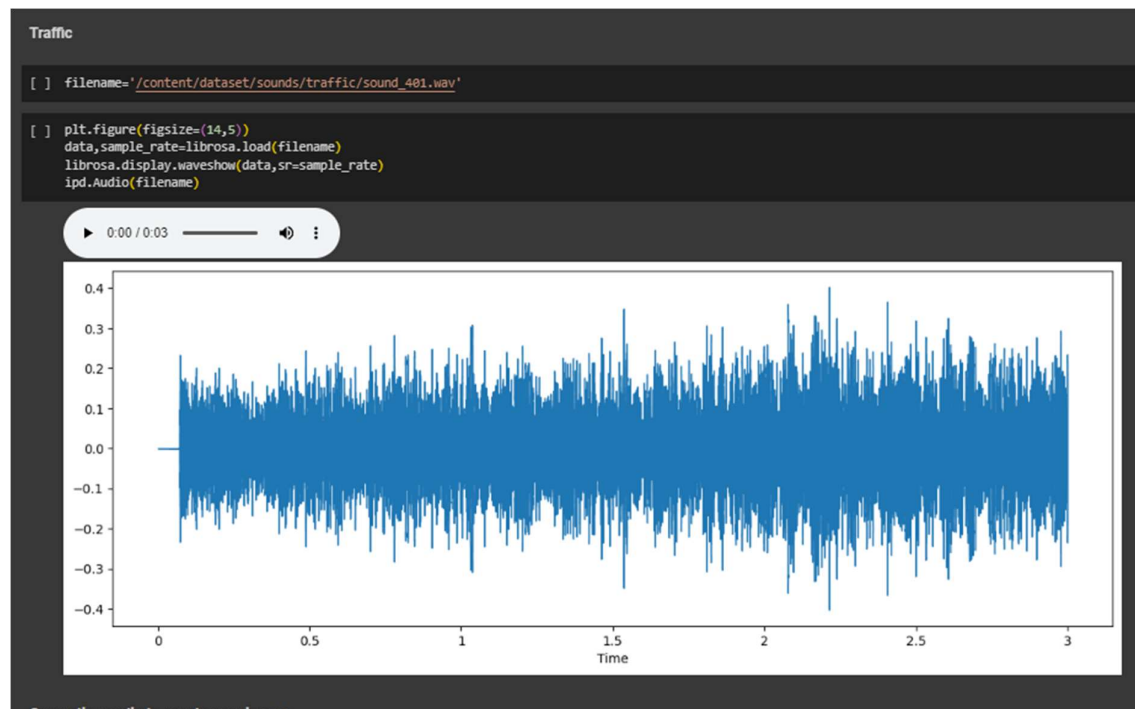
code can be found here –

<https://colab.research.google.com/drive/1cHC1WOxZPR4rUVZffjINT-i54Asv3N7P?usp=sharing>

Task 1 and 3

First we display the audio as a wave using `librosa.display.waveshow` and `lpython.display`





Now we create a function to create spectrograms from audio files.

```
import os
import librosa
import numpy as np
import matplotlib.pyplot as plt

# Define paths
data_folder = '/content/dataset/sounds'
output_folder = '/content/spectrogram_images'

# Function to create spectrogram
def create_spectrogram(file_path, output_path):
    # Load audio file
    audio, sr = librosa.load(file_path, sr=None)

    # Compute spectrogram
    S = librosa.feature.melspectrogram(y=audio, sr=sr)
    S_db = librosa.power_to_db(S, ref=np.max)

    # Plot and save spectrogram
    plt.figure(figsize=(10, 4))
    librosa.display.specshow(S_db, sr=sr, x_axis='time', y_axis='mel')
    plt.colorbar(format='%+2.0f dB')
    plt.title('Mel-frequency spectrogram')
    plt.savefig(output_path)
    plt.close()

# Iterate through all files in the data folder
file_count = 0
total_files = sum([len(files) for r, d, files in os.walk(data_folder)])
for root, dirs, files in os.walk(data_folder):
    for file in files:
        if file.endswith('.wav'):
            file_count += 1
            print(f"Processing file {file_count}/{total_files}...")

            # Extract class from filename
            class_name = os.path.basename(root)

            # Create output folder if not exists
            class_folder = os.path.join(output_folder, class_name)
            os.makedirs(class_folder, exist_ok=True)

            # Generate spectrogram and save
            file_path = os.path.join(root, file)
            output_path = os.path.join(class_folder, file.replace('.wav', '.png'))
            create_spectrogram(file_path, output_path)

print("Spectrogram generation completed.")
```

Using this function we iterate through all 600 audio files and convert them to spectrogram images and make new dataset 'spectrogram\_images'. It utilizes the 'librosa' library for audio processing and 'matplotlib' for visualization.

Here,

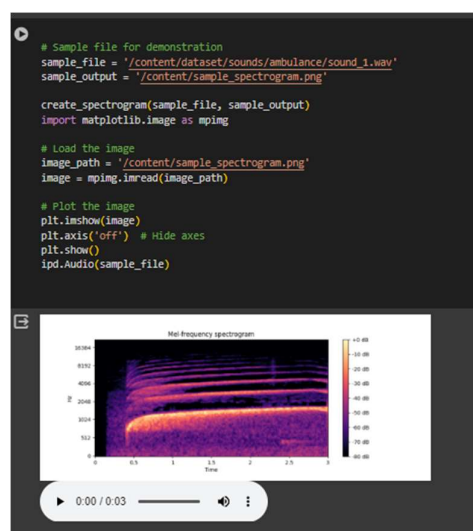
- 'data\_folder': Path to the directory containing the audio files.
- 'output\_folder': Path to the directory where spectrogram images will be saved.

Create Spectrogram Function: This function 'create\_spectrogram' takes a file path as input, loads the audio file using 'librosa', computes the mel-spectrogram, converts it to decibels (dB), plots the spectrogram using 'matplotlib', and saves the plot as an image file (.png) in the specified output path.

We iterate through all files in the 'data\_folder' using 'os.walk()'. For each '.wav' file encountered, it:

- Increments the 'file\_count'.
- Extracts the class (category) name from the directory structure.
- Creates a subfolder in the 'output\_folder' corresponding to the class name if it doesn't exist.
- Generates a spectrogram for the audio file and saves it in the appropriate class subfolder.

The generated spectrogram images can then be used as input data for CNN.



This is a sample audio file and its spectrogram

## Building CNN models

1. Model1 has 4 conv2d layers with filters 32,64,128,256 respectively. Kernel size is 3x3 and stride=1. Maaxpooling size=2x2. This is a very popular and basic CNN architecture. This model gave an accuracy of 98% on train set and 99% on test and validation set. Categorical Crossentropy loss was used and adam was used as optimizer. The model was trained for 15 epochs with batch size=32 and 18 timesteps per epoch. We see that here the model converges at 10 epochs. Accuracy for classes ambulance and traffic are 100% and class firetruck has accuracy 99%. This model takes significant amount of time to run using cpu and gives good results in the end.

Model: "sequential\_1"

Layer (type)	Output Shape	Param #
conv2d_4 (Conv2D)	(None, 124, 124, 32)	2432
max_pooling2d_4 (MaxPooling2D)	(None, 62, 62, 32)	0
conv2d_5 (Conv2D)	(None, 58, 58, 64)	51264
max_pooling2d_5 (MaxPooling2D)	(None, 29, 29, 64)	0
conv2d_6 (Conv2D)	(None, 25, 25, 128)	204928
max_pooling2d_6 (MaxPooling2D)	(None, 12, 12, 128)	0
conv2d_7 (Conv2D)	(None, 8, 8, 256)	819456
max_pooling2d_7 (MaxPooling2D)	(None, 4, 4, 256)	0
flatten_1 (Flatten)	(None, 4096)	0
dense_2 (Dense)	(None, 512)	2097664
dropout_1 (Dropout)	(None, 512)	0
dense_3 (Dense)	(None, 3)	1539

=====  
Total params: 3177283 (12.12 MB)  
Trainable params: 3177283 (12.12 MB)  
Non-trainable params: 0 (0.00 Byte)  
=====  
Epoch 1/15



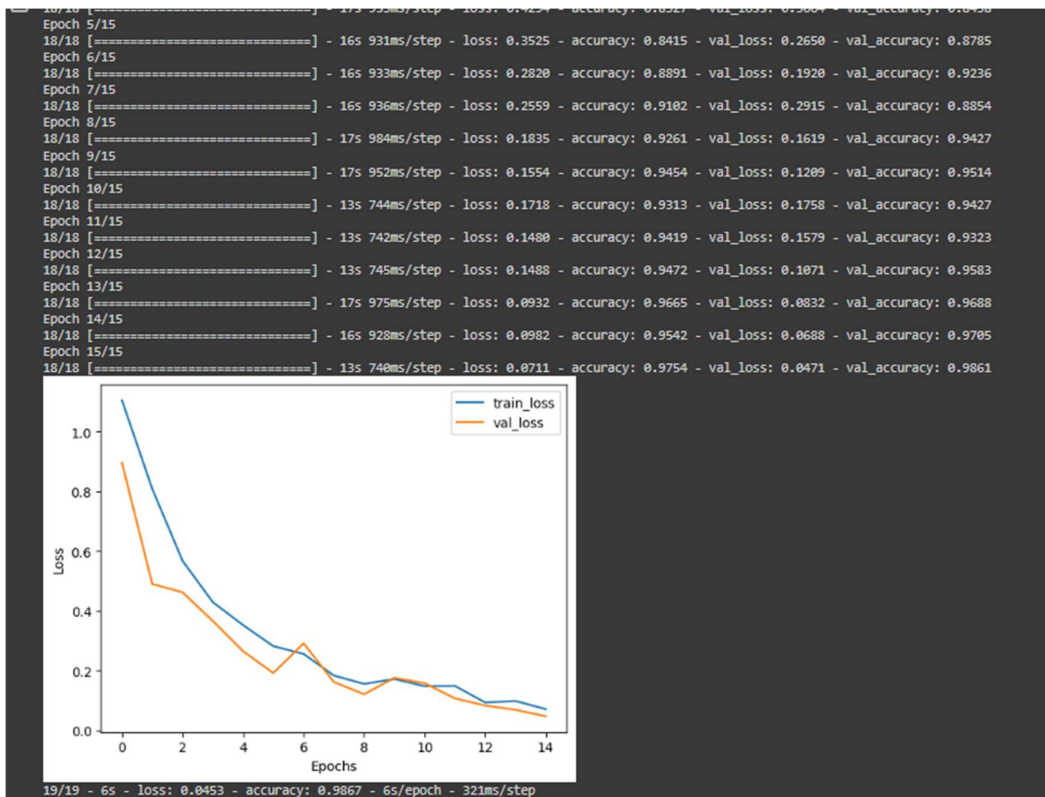
- Model2 has 4 conv2d layers with filters 32,64,128,256 respectively. Kernel size is 5x5 and stride=1. Maxpooling size=2x2. Here we increased the kernel size This model gave an accuracy of 97% on train set and 98% on test and validation set. Categorical Crossentropy loss was used and adam was used as optimizer. The model was trained for 15 epochs with

batch size=32 and 18 timesteps per epoch. We see that here the model doesn't really converge even after 15 epochs it converges properly at 18 epochs. So with increase in kernel size number of epochs to converge increased. This might be due to increase in number of trainable parameters compared to model1. Accuracy for class traffic is 100% and rest of the two classes have 98%. This model takes longest amount of time to run using cpu and gives decent(not too good compared to model 1) results in the end.

Found 600 images belonging to 3 classes.  
Model: "sequential"

Layer (type)	Output Shape	Param #
conv2d (Conv2D)	(None, 126, 126, 32)	896
max_pooling2d (MaxPooling2D)	(None, 63, 63, 32)	0
conv2d_1 (Conv2D)	(None, 61, 61, 64)	18496
max_pooling2d_1 (MaxPooling2D)	(None, 30, 30, 64)	0
conv2d_2 (Conv2D)	(None, 28, 28, 128)	73856
max_pooling2d_2 (MaxPooling2D)	(None, 14, 14, 128)	0
conv2d_3 (Conv2D)	(None, 12, 12, 256)	295168
max_pooling2d_3 (MaxPooling2D)	(None, 6, 6, 256)	0
flatten (Flatten)	(None, 9216)	0
dense (Dense)	(None, 512)	4719104
dropout (Dropout)	(None, 512)	0
dense_1 (Dense)	(None, 3)	1539

=====  
Total params: 5109059 (19.49 MB)  
Trainable params: 5109059 (19.49 MB)  
Non-trainable params: 0 (0.00 Byte)



```

19/19 - 6s - loss: 0.0453 - accuracy: 0.9867 - 6s/epoch - 321ms/step
19/19 [=====] - 7s 359ms/step
Classification Report:
      precision    recall  f1-score   support

     0       0.97       0.99       0.98         200
     1       0.99       0.97       0.98         200
     2       1.00       1.00       1.00         200

 accuracy          0.99
 macro avg         0.99
 weighted avg      0.99

```

Test accuracy: 0.986666793823242

```

] from sklearn.metrics import confusion_matrix
# Calculate confusion matrix
conf_matrix = confusion_matrix(validation_generator.classes, y_pred_classes)

# Print confusion matrix
print("\nConfusion Matrix:")
print(conf_matrix)

Confusion Matrix:
[[198  2  0]
 [  6 194  0]
 [  0  0 200]]

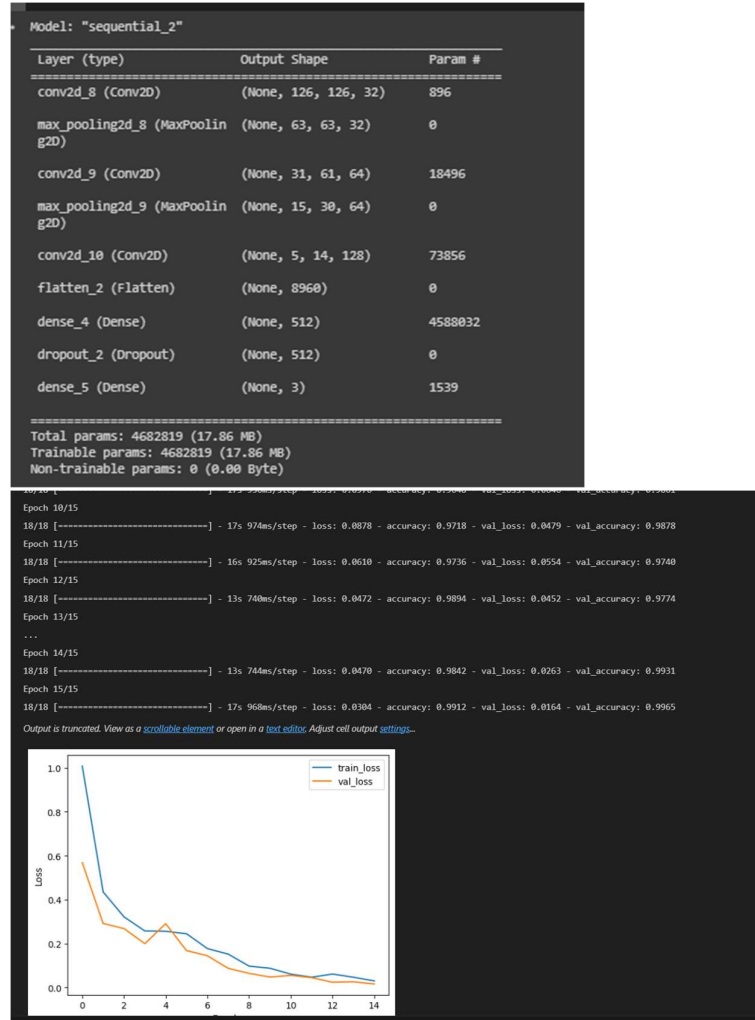
```

1 # CNN model with different stride values

- Model3 has 4 conv2d layers with filters 32,64,128,256 respectively. Kernel size is 3x3. Maaxpooling size=2x2. Here stride value is taken (2,1) for 2<sup>nd</sup> convolution layer and (3,2) for 3<sup>rd</sup> convolution layer. This model performed surprisingly best. gave an accuracy of 99% on train set and valid set and 100% on test set. Categorical Crossentropy loss was used and adam was used as optimizer. The model was trained for 15 epochs



with batch size=32 and 18 timesteps per epoch. We see that here the model converges very quickly at 10 epochs. Accuracy for class ambulance is 99% and rest of the two classes have 100%. This model takes shortest amount of time to run using cpu and gives best results in the end.



```
19/19 - 6s - loss: 0.0158 - accuracy: 0.9967 - 6s/epoch - 322ms/step
19/19 [=====] - 6s 319ms/step
Classification Report:

      precision    recall  f1-score   support

     0       1.00      0.99      0.99       200
     1       0.99      1.00      1.00       200
     2       1.00      1.00      1.00       200

 accuracy          1.00         600
 macro avg          1.00         600
 weighted avg       1.00         600

Test accuracy: 0.996666669845581

from sklearn.metrics import confusion_matrix
# Calculate confusion matrix
conf_matrix = confusion_matrix(validation_generator.classes,

# Print confusion matrix
print("\nConfusion Matrix:")
print(conf_matrix)
```

```
Confusion Matrix:
[[198  2  0]
 [ 0 200  0]
 [ 0  0 200]]
```

4. Model4 has 4 conv2d layers with filters 64,128,256,512 respectively. Kernel size is 3x3. Maaxpooling size=2x2. Here we increased number of filters. This model performed decently. gave an accuracy of 98% on train set ,99% on test set and validation set. Categorical Crossentropy loss was used and adam was used as optimizer. The model was trained for 15 epochs with batch size=32 and 18 timesteps per epoch. We see that here the model converges at 12 epochs. Accuracy for class traffic is 100% and rest of the two classes have 98%. This model takes significant time to run and gives good results.

Model: "sequential\_3"

Layer (type)	Output Shape	Param #
conv2d_11 (Conv2D)	(None, 127, 127, 64)	832
max_pooling2d_10 (MaxPooling2D)	(None, 63, 63, 64)	0
conv2d_12 (Conv2D)	(None, 62, 62, 128)	32896
max_pooling2d_11 (MaxPooling2D)	(None, 31, 31, 128)	0
conv2d_13 (Conv2D)	(None, 30, 30, 256)	131328
max_pooling2d_12 (MaxPooling2D)	(None, 15, 15, 256)	0
conv2d_14 (Conv2D)	(None, 14, 14, 512)	524800
max_pooling2d_13 (MaxPooling2D)	(None, 7, 7, 512)	0
flatten_3 (Flatten)	(None, 25888)	0
dense_6 (Dense)	(None, 512)	12845568
dropout_3 (Dropout)	(None, 512)	0
dense_7 (Dense)	(None, 3)	1539

=====  
Total params: 13536963 (51.64 MB)  
Trainable params: 13536963 (51.64 MB)  
Non-trainable params: 0 (0.00 Byte)

Epoch 11/15

18/18 [=====] - 16s 939ms/step - loss: 0.0770 - accuracy: 0.9683 - val\_loss: 0.1664 - val\_accuracy: 0.9340  
Epoch 12/15

18/18 [=====] - 13s 756ms/step - loss: 0.0809 - accuracy: 0.9736 - val\_loss: 0.0472 - val\_accuracy: 0.9792  
Epoch 13/15

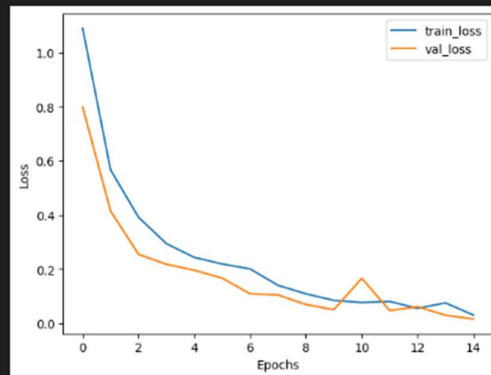
...

Epoch 14/15

18/18 [=====] - 14s 776ms/step - loss: 0.0756 - accuracy: 0.9701 - val\_loss: 0.0301 - val\_accuracy: 0.9931  
Epoch 15/15

18/18 [=====] - 13s 747ms/step - loss: 0.0308 - accuracy: 0.9894 - val\_loss: 0.0161 - val\_accuracy: 0.9948

Output is truncated. View as a [scrollable element](#) or open in a [text editor](#). Adjust cell output [settings](#).



19/19 - 6s - loss: 0.0155 - accuracy: 0.9950 - 6s/epoch - 326ms/step

19/19 [=====] - 6s 324ms/step

Classification Report:

	precision	recall	f1-score	support
0	0.99	0.99	0.99	200
1	0.99	0.99	0.99	200
2	1.00	1.00	1.00	200
accuracy			0.99	600
macro avg	0.99	0.99	0.99	600
weighted avg	0.99	0.99	0.99	600

```
[34]
...
Confusion Matrix:
[[199  1  0]
 [  1 198  1]
 [  0  0 200]]
```

Task2:

We can represent text(sentence) as a matrix of word embeddings of each word. We are using this concept to represent text and we'll use CNN.

Matrix form of text:

```
✓ 0s print("text matrix")
      print(X_train)
      print(X_train.shape)

text matrix
[[ 0  0  0 ... 19 178 32]
 [ 0  0  0 ... 16 145 95]
 [ 0  0  0 ... 7 129 113]
 ...
 [ 0  0  0 ... 4 3586 2]
 [ 0  0  0 ... 12 9 23]
 [ 0  0  0 ... 204 131 9]]
(25000, 450)
```

We are using 1d CNN here. If our matrix version of text is of size  $(m,n)$  where  $n$  are number of words in text and  $m$  is size of embedding then kernel size is  $(n,x)$  here  $n$  is fixed.  $n$  here is 25000 number of imdb text reviews.  $n$  is `max_words = 450`

```
(25000,)  
  
# Padding the data samples to a maximum review length in words  
max_words = 450  
X_train = sequence.pad_sequences(X_train, maxlen=max_words)  
X_test = sequence.pad_sequences(X_test, maxlen=max_words)  
# Building the CNN Model  
model = Sequential() # initializing the Sequential nature for CNN model  
# Adding the embedding layer which will take in maximum of 450 words as input and provide a 32 dimensional output of those words which belong in the top_words dictionary  
model.add(Embedding(top_words, 32, input_length=max_words))  
model.add(Conv1D(32, 3, padding='same', activation='relu'))  
model.add(MaxPooling1D())  
model.add(Flatten())  
model.add(Dense(250, activation='relu'))  
model.add(Dense(1, activation='sigmoid'))  
  
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])  
model.summary()  
  
Model: "sequential_4"  
=====
```

Layer (type)	Output Shape	Param #
embedding (Embedding)	(None, 450, 32)	224000
conv1d (Conv1D)	(None, 450, 32)	3104
max_pooling1d (MaxPooling1D)	(None, 225, 32)	0
flatten_4 (Flatten)	(None, 7200)	0
dense_8 (Dense)	(None, 250)	1800250
dense_9 (Dense)	(None, 1)	251

```
=====
```

Total params: 2027605 (7.73 MB)	
Trainable params: 2027605 (7.73 MB)	
Non-trainable params: 0 (0.00 Byte)	

```
=====
```

```
# Fitting the data onto model  
model.fit(X_train, y_train, validation_data=(X_test, y_test), epochs=2, batch_size=128, verbose=2)  
# Getting score metrics from our model  
scores = model.evaluate(X_test, y_test, verbose=0)  
# Displays the accuracy of correct sentiment prediction over test data  
print("Accuracy: %.2f%%" % (scores[1]*100))  
  
1m 30s completed at 12:01 AM
```

The model first consists of embedding layer in which we will find the embeddings of the top 7000 words into a 32 dimensional embedding and the input we can take in is defined as the maximum length of a review allowed.

Then, we add the convolutional layer and max-pooling layer. Finally, we flatten those matrices into vectors and add dense layers(basically scale,rotating and transform the vector by multiplying Matrix and vector).

The last Dense layer is having one as parameter because we are doing a binary classification and so we need only one output node in our vector.

```
[21] # Fitting the data onto model  
model.fit(X_train, y_train, validation_data=(X_test, y_test), epochs=2, batch_size=128, verbose=2)  
# Getting score metrics from our model  
scores = model.evaluate(X_test, y_test, verbose=0)  
# Displays the accuracy of correct sentiment prediction over test data  
print("Accuracy: %.2f%%" % (scores[1]*100))  
  
Epoch 1/2  
196/196 - 30s - loss: 0.4434 - accuracy: 0.7666 - val_loss: 0.2750 - val_accuracy: 0.8846 - 30s/epoch - 155ms/step  
Epoch 2/2  
196/196 - 25s - loss: 0.2020 - accuracy: 0.9222 - val_loss: 0.2743 - val_accuracy: 0.8864 - 25s/epoch - 125ms/step  
Accuracy: 88.64%
```

We achieved an accuracy of 88% on test set.