# 1 Author

**Student Name**: Dnyanesh Walwadkar Advance Solution Notebook

# 2 Problem formulation

Building Audio Deep Learning Model for forcasting melody name.

Sound Classification is one of the foremost broadly utilized applications in Audio Deep Learning. It includes learning to classify sounds and to anticipate the category of that sound. I am going begin with sound files, convert them into spectrograms, input them into a CNN plus Linear Classifier model, and create forecasts almost the lesson to which the melody has a place.

I started with sound *.wav files, converted them into spectrograms, input them into a CNN plus Linear Classifier model, and produced predictions about the class ( Song ) to which the sound belongs.

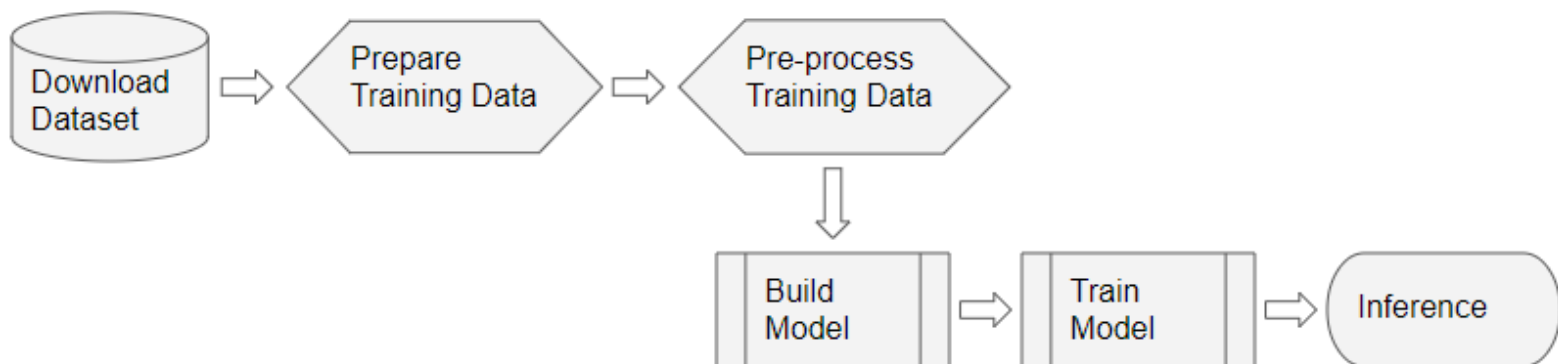I Trained model for 119 Epoch and Got following results:

**Epoch: 200, Loss: 0.14, Traning Accuracy: ( 97% - 98% )**

**Validation Accuracy : ( 52% - 70% )**

( I included range of accuracies beacuse as traning dataset get different number of audio files from different songs, model get affected. If we increse number of audio files per song, this model will get improved. )

We can improve Validation accuracy, for that we have to consider large number of audio files per songClass in our traning dataset.

# 3 Machine Learning pipeline



## Import Dataset from Google Drive

Dataset Contains 8 different types of Songs Audio files. Total 98. We are preprocessing given Dataset to build our required or desired dataset as per follows :

The training data for this problem are classified as:

- The features (X) are the audio file paths.
- The target labels (y) are the class names.

song = "Potter" Class = 1

song = "StarWars" Class = 2

song = "Frozen" Class = 3

song="Panther" Class = 4

song = 'Rain' Class = 5

song = "Showman" Class = 6

song = "Mamma" Class = 7

song = "Hakuna" Class = 8

# 4 Transformation stage

## ▾ Read audio from a file

The first thing we need is to read and load the audio file in ".wav" format as per our project specification. Since we are using Pytorch for this example, the implementation below uses torchaudio for the audio processing, but librosa will work just as well. Librosa use we saw in earlier solution.
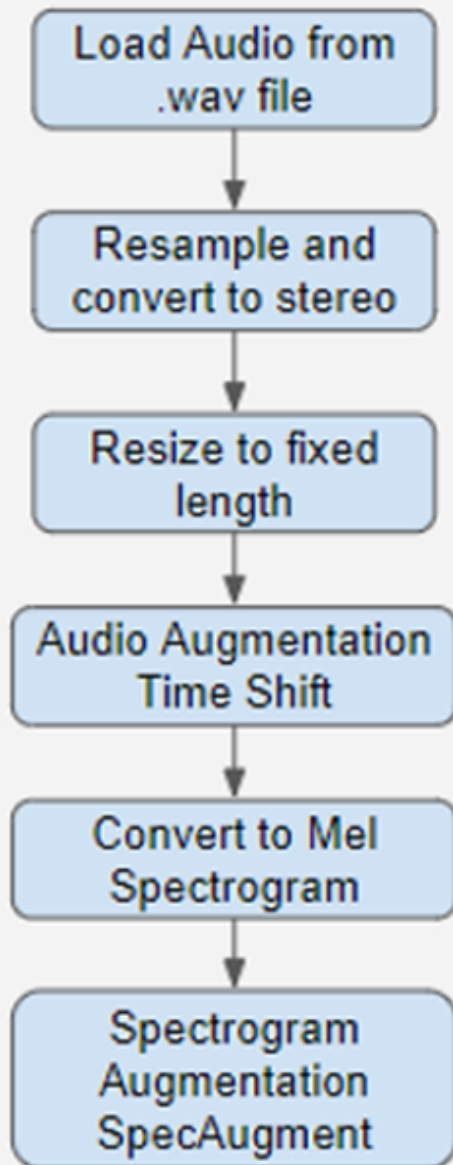
In this code we are taking care of Mono & Stereo Sounds.

The difference between monophonic (mono) and stereophonic (stereo) sound is the number of channels used to record and playback audio. Mono signals are recorded and played back using a single audio channel, while stereo sounds are recorded and played back using two audio channels. As a listener, the most noticeable difference is that stereo sounds are capable of producing the perception of width, whereas mono sounds are not.

We are using following libraries for this task.

- torch
- torchaudio - Torchaudio is a library for audio and signal processing with PyTorch. It provides I/O, signal and data processing functions, datasets, model implementations and application components
- Ipython.display.Audio

## Transcforms

```
┌──────────────────┐
│  Load Audio from │
│    .wav file     │
└──────────────────┘
          │
          ▼
┌──────────────────┐
│  Resample and    │
│ convert to stereo│
└──────────────────┘
          │
          ▼
┌──────────────────┐
│  Resize to fixed │
│     length       │
└──────────────────┘
          │
          ▼
┌──────────────────┐
│Audio Augmentation│
│    Time Shift    │
└──────────────────┘
          │
          ▼
┌──────────────────┐
│  Convert to Mel  │
│   Spectrogram    │
└──────────────────┘
          │
          ▼
┌──────────────────┐
│   Spectrogram    │
│  Augmentation    │
│   SpecAugment    │
└──────────────────┘
```

```python
 1 import math, random
 2 import torch
 3 import torchaudio
 4 from torchaudio import transforms
 5 from IPython.display import Audio
 6
 7 class AudioUtil():
 8   # ----------------------------
 9   # Load an audio file. Return the song audio signal as a tensor and the sample rate
10   # ----------------------------
11   @staticmethod
12   def open(audio_file):
13     sig, sr = torchaudio.load(audio_file)
14     return (sig, sr)
15
16   def rechannel(aud, new_channel):
```

```python
17      sig, sr = aud
18
19      if (sig.shape[0] == new_channel):
20
21        return aud
22
23      if (new_channel == 1):
24        # Convert from stereo to mono by selecting only the first channel
25        resig = sig[:1, :]
26      else:
27        # Convert from mono to stereo by duplicating the first channel
28        resig = torch.cat([sig, sig])
29
30      return ((resig, sr))
31
32    def resample(aud, newsr):
33      sig, sr = aud
34
35      if (sr == newsr):
36        # Nothing to do
37        return aud
38
39      num_channels = sig.shape[0]
40      # Resample first channel
41      resig = torchaudio.transforms.Resample(sr, newsr)(sig[:1,:])
42      if (num_channels > 1):
43        # Resample the second channel and merge both channels
44        retwo = torchaudio.transforms.Resample(sr, newsr)(sig[1:,:])
45        resig = torch.cat([resig, retwo])
46
47      return ((resig, newsr))
48
49    def pad_trunc(aud, max_ms):
50      sig, sr = aud
51      num_rows, sig_len = sig.shape
52      max_len = sr//1000 * max_ms
53
54      if (sig_len > max_len):
55        # Truncate the signal to the given length
56        sig = sig[:,:max_len]
57
58      elif (sig_len < max_len):
59        # Length of padding to add at the beginning and end of the signal
60        pad_begin_len = random.randint(0, max_len - sig_len)
61        pad_end_len = max_len - sig_len - pad_begin_len
62
63        # Pad with 0s
64        pad_begin = torch.zeros((num_rows, pad_begin_len))
65        pad_end = torch.zeros((num_rows, pad_end_len))
66
67        sig = torch.cat((pad_begin, sig, pad_end), 1)
68
69      return (sig, sr)
70
71    def time_shift(aud, shift_limit):
72      sig,sr = aud
```

```
73      _, sig_len = sig.shape
74      shift_amt = int(random.random() * shift_limit * sig_len)
75      return (sig.roll(shift_amt), sr)
76
77    def spectro_gram(aud, n_mels=64, n_fft=1024, hop_len=None):
78      sig,sr = aud
79      top_db = 80
80
81      # spec has shape [channel, n_mels, time], where channel is mono, stereo etc
82      spec = transforms.MelSpectrogram(sr, n_fft=n_fft, hop_length=hop_len, n_mels=n_mels)(sig)
83
84      # Convert to decibels
85      spec = transforms.AmplitudeToDB(top_db=top_db)(spec)
86      return (spec)
87
88    def spectro_augment(spec, max_mask_pct=0.1, n_freq_masks=1, n_time_masks=1):
89      _, n_mels, n_steps = spec.shape
90      mask_value = spec.mean()
91      aug_spec = spec
92
93      freq_mask_param = max_mask_pct * n_mels
94      for _ in range(n_freq_masks):
95        aug_spec = transforms.FrequencyMasking(freq_mask_param)(aug_spec, mask_value)
96
97      time_mask_param = max_mask_pct * n_steps
98      for _ in range(n_time_masks):
99        aug_spec = transforms.TimeMasking(time_mask_param)(aug_spec, mask_value)
100
101     return aug_spec
```

# ▾ Custom Data Loader

Now that we have defined all the pre-processing transform functions we will define a **custom Pytorch Dataset object.** To feed your data to a model with Pytorch, we need two objects:

- A custom Dataset object that uses all the audio transforms to pre-process an audio file and prepares one data item at a time.
- A built-in DataLoader object that uses the Dataset object to fetch individual data items and packages them into a batch of data.

**Demo Steps for Converting Audio Signals to Spectrogram**

```
1 # ----------------------------
2 # Sound Dataset
3 # ----------------------------
4
5 from torch.utils.data import DataLoader, Dataset, random_split
6 import torchaudio
7
8
9 class SoundDS(Dataset):
```

```python
10   def __init__(self, df, song_path):
11     self.df = df
12     self.data_path = str(song_path)
13     self.duration = 4000
14     self.sr =  11000
15     self.channel = 2
16     self.shift_pct = 0.4
17
18   # ----------------------------
19   # Number of items in dataset
20   # ----------------------------
21   def __len__(self):
22     return len(self.df)
23
24   # ----------------------------
25   # Get i'th item in dataset
26   # ----------------------------
27   def __getitem__(self, idx):
28     # Absolute file path of the audio file - concatenate the audio directory with
29     # the relative path
30     song_Audio = self.data_path + str(self.df.loc[idx,'file_id'])
31     # Get the Song ID
32     class_id = int(self.df.loc[idx, 'song'])
33     print(AudioUtil.open(song_Audio))
34
35     song = AudioUtil.open(song_Audio)
36
37     reaud = AudioUtil.resample(song, self.sr)
38     rechan = AudioUtil.rechannel(reaud, self.channel)
39
40     dur_song = AudioUtil.pad_trunc(rechan, self.duration)
41     shift_song = AudioUtil.time_shift(dur_song, self.shift_pct)
42     sgram = AudioUtil.spectro_gram(shift_song, n_mels=64, n_fft=1024, hop_len=None)
43     aug_sgram = AudioUtil.spectro_augment(sgram, max_mask_pct=0.1, n_freq_masks=2, n_time_masks=2)
44
45     return aug_sgram, class_id
```

# Prepared Batches of Data with the Data Loader for traning & Validation

All of the functions we need to input our data to the model have now been defined.

- We use our custom Dataset to load the Features and Labels from our Pandas dataframe and split that data randomly in an 90:10 ratio into training and validation sets.
- We then use them to create our training and validation Data Loaders.

```python
1 sample_path1 = '/content/drive/MyDrive/Data/MLEndHW/sample/MLEndHW_Sample/'
2
```

```python
1 #Run this cell after running Dataset Cell
2 #-------------------------------------------------------------------
3
```

```
 4 from torch.utils.data import random_split
 5 from torch.utils.data import DataLoader
 6
 7 #MLENDHW_df = PandasDataset(MLENDHW_df)
 8 myds = SoundDS(MLENDHW_df, sample_path1)
 9
10 # Random split of 90:10 between training and validation
11 num_items = len(myds)
12 num_train = round(num_items * 0.8)
13 num_val = num_items - num_train
14 train_ds, val_ds = random_split(myds, [num_train, num_val])
15 #train_ds = PandasDataset(train_ds)
16 #val_ds = PandasDataset(val_ds)
17
18 # numTD = round(num_items * 0.6)
19 # numD = num_items - numTD
20 # traD, vaD = random_split(myds, [numTD, numD])
21
22 # Create training and validation data loaders
23 train_dl = DataLoader(train_ds, batch_size=4)
24 val_dl = DataLoader(val_ds, batch_size=4)
```

# 5 Modelling

## ▾ Model

Since our data now consists of Spectrogram images, we build a CNN classification architecture to process them. It has four convolutional blocks which generate the feature maps. That data is then reshaped into the format we need so it can be input into the linear classifier layer, which finally outputs the predictions for the 8 song classes.

**Spectrogram ▶**

A spectrogram is a visual representation of the spectrum of frequencies of a signal as it varies with time. When applied to an audio signal, spectrograms are sometimes called sonographs, voiceprints, or voicegrams.

Spectrograms are used extensively in the fields of music speech processing,seismology, and others. Spectrograms of audio can be used to identify spoken words phonetically, and to analyse the various audio patterns

```python
 1 import torch.nn as nn
 2 import torch.nn.functional as F
 3 from torch.nn import init
 4
 5 # ----------------------------
 6 # Audio Classification Model
 7 # ----------------------------
 8 class AudioClassifier (nn.Module):
 9     # ----------------------------
10     # Build the model architecture
11     # ----------------------------
12     def __init__(self):
13         super().__init__()
14         conv_layers = []
15
16         # First Convolution Block with Relu and Batch Norm. Use Kaiming Initialization
17         self.conv1 = nn.Conv2d(2, 8, kernel_size=(5, 5), stride=(2, 2), padding=(2, 2))
18         self.relu1 = nn.ReLU()
19         self.bn1 = nn.BatchNorm2d(8)
20         init.kaiming_normal_(self.conv1.weight, a=0.1)
21         self.conv1.bias.data.zero_()
22         conv_layers += [self.conv1, self.relu1, self.bn1]
23
24         # Second Convolution Block
25         self.conv2 = nn.Conv2d(8, 16, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1))
26         self.relu2 = nn.ReLU()
27         self.bn2 = nn.BatchNorm2d(16)
28         init.kaiming_normal_(self.conv2.weight, a=0.1)
29         self.conv2.bias.data.zero_()
30         conv_layers += [self.conv2, self.relu2, self.bn2]
31
32         # Second Convolution Block
33         self.conv3 = nn.Conv2d(16, 32, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1))
34         self.relu3 = nn.ReLU()
35         self.bn3 = nn.BatchNorm2d(32)
36         init.kaiming_normal_(self.conv3.weight, a=0.1)
37         self.conv3.bias.data.zero_()
38         conv_layers += [self.conv3, self.relu3, self.bn3]
39
40         # Second Convolution Block
41         self.conv4 = nn.Conv2d(32, 64, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1))
42         self.relu4 = nn.ReLU()
43         self.bn4 = nn.BatchNorm2d(64)
44         init.kaiming_normal_(self.conv4.weight, a=0.1)
45         self.conv4.bias.data.zero_()
46         conv_layers += [self.conv4, self.relu4, self.bn4]
47
48         # Linear Classifier
49         self.ap = nn.AdaptiveAvgPool2d(output_size=1)
50         self.lin = nn.Linear(in_features=64, out_features=10)
51
52         # Wrap the Convolutional Blocks
53         self.conv = nn.Sequential(*conv_layers)
54
55     # ----------------------------
56     # Forward pass computations
```

```
57      # ----------------------------
58      def forward(self, x):
59          # Run the convolutional blocks
60          x = self.conv(x)
61
62          # Adaptive pool and flatten for input to linear layer
63          x = self.ap(x)
64          x = x.view(x.shape[0], -1)
65
66          # Linear layer
67          x = self.lin(x)
68
69          # Final output
70          return x
71
72 # Create the model and put it on the GPU if available
73 myModel = AudioClassifier()
74 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
75 myModel = myModel.to(device)
76 # Check that it is on Cuda
77 next(myModel.parameters()).device
```

```
device(type='cpu')
```

# 6 Methodology

We characterize the capacities for the optimizer, loss, and scheduler to powerfully change our learning rate as preparing advances, which ordinarily permits preparing to focalize in less epochs.

We prepare the demonstrate for a several epochs, handling a clump of information in each emphasis. We keep track of a simple accuracy metric which measures the percentage of correct forcast.



Audio wave → Spectrogram → CNN Architecture → Feature Maps → Linear Classifier → Potter (Class – 01), StarWars (Class – 02), Frozen (Class – 03) ... ... Hakuna (Class – 08)
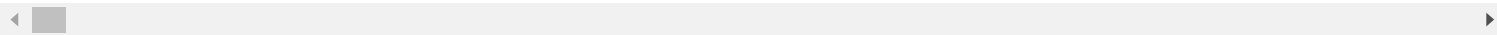
# 7 Dataset

```
1
2
3 sample_path = '/content/drive/MyDrive/Data/MLEndHW/sample/MLEndHW_Sample/*.wav'
4 files = glob.glob(sample_path)
5 print(files)
```

```
1 for _ in range(5):
2   n = np.random.randint(98)
3   display(ipd.Audio(files[n]))
```

0:00 / 0:21

0:00 / 0:15

0:00 / 0:15

0:00 / 0:15

0:00 / 0:14

# ▾ Preprocessing

Replace Song Name Colunm as Distinct Numric Class Value. We will use This Classes as our Target value during Model Building.

This training data with audio file paths cannot be input directly into the model. We have to load the audio data from the file and process it so that it is in a format that the model expects.

```
 1 MLENDHW_table = []
 2 flag=0
 3
 4 for file in files:
 5   file_name = file.split('/')[-1]
 6   flag=0
 7   try:
 8     song = file.split('/')[-1].split('_')[3].split('.')[0]
 9     if(song=="Potter"):
10       song = 1
11     elif(song=="StarWars"):
12       song=2
13     elif(song=="Frozen"):
14       song=3
15     elif(song=="Panther"):
16       song=4
17     elif(song=='Rain'):
18       song=5
19     elif(song=="Showman"):
20       song=6
21     elif(song=="Mamma"):
```

```
22        song=7
23      elif(song=="Hakuna"):
24        song=8
25      else:
26        flag=1
27      if flag==0:
28        MLENDHW_table.append([file_name, song])
29  except:
30    pass
31
32
33
34 MLENDHW_table
```

```
1 MLENDHW_df = pd.DataFrame(MLENDHW_table,columns=['file_id','song'])
2 MLENDHW_df
3 MLENDHW_df.groupby('song').count()
```
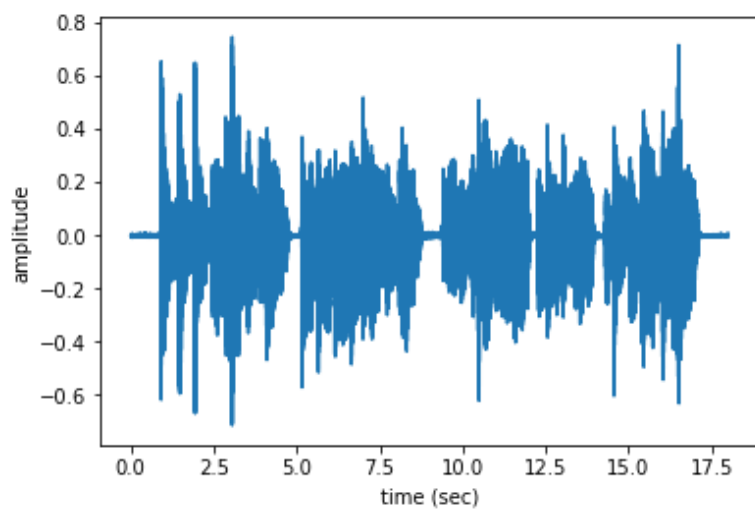
|          | file_id |
|----------|---------|
| **song** |         |
| **1**    | 34      |
| **2**    | 31      |
| **3**    | 32      |
| **4**    | 34      |
| **5**    | 34      |
| **6**    | 33      |
| **7**    | 29      |
| **8**    | 34      |

```
1 n=0
2 fs = None # Sampling frequency. If None, fs would be 22050
3 for i in range(5):
4
5     x, fs = librosa.load(files[i],sr=fs)
6     t = np.arange(len(x))/fs
7     plt.plot(t,x)
8     plt.xlabel('time (sec)')
9     plt.ylabel('amplitude')
10    plt.show()
11    display(ipd.Audio(files[n]))
```
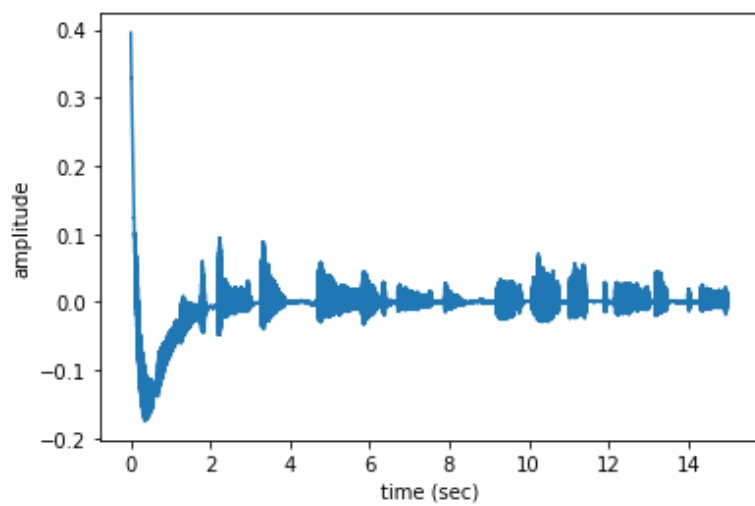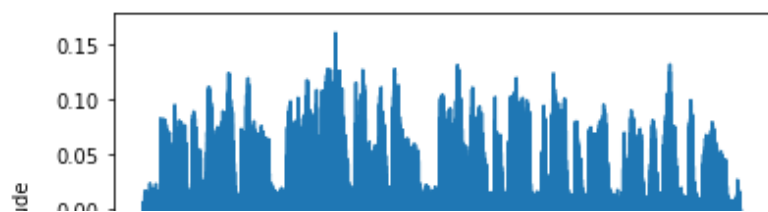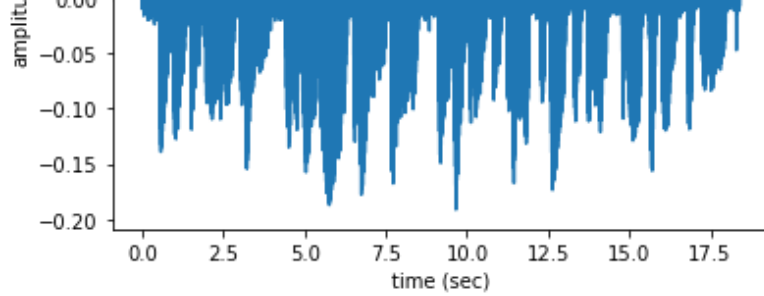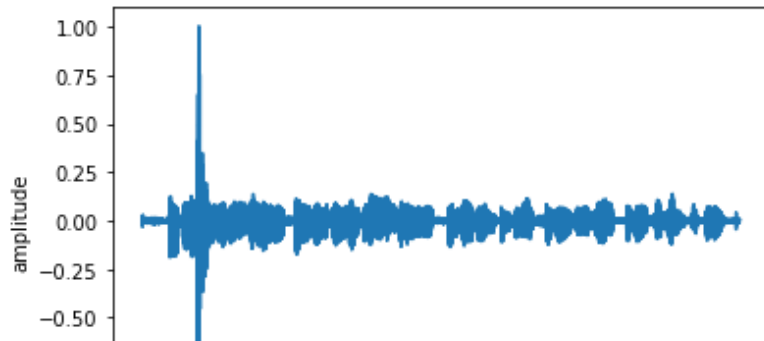
0:00 / 0:17



0:00 / 0:17



0:00 / 0:17

0:00 / 0:17

Traning & Validation

0:00 / 0:17

```
1 def training(model, train_dl, num_epochs):
2    # Loss Function, Optimizer and Scheduler
3    criterion = nn.CrossEntropyLoss()
4    optimizer = torch.optim.Adam(model.parameters(),lr=0.001)
5    scheduler = torch.optim.lr_scheduler.OneCycleLR(optimizer, max_lr=0.001,
6                                                     steps_per_epoch=int(len(train_dl)),
7                                                     epochs=num_epochs,
8                                                     anneal_strategy='linear')
9
10   # Repeat for each epoch
11
12
13   for epoch in range(num_epochs):
14     running_loss = 0.0
15     correct_prediction = 0
16     total_prediction = 0
17     first = 0
18
19
20     # Repeat for each batch in the training set
21     for i, data in enumerate(train_dl):
22        # Get the input features and target labels, and put them on the GPU
23
24        inputs, labels = data[0], data[1]
25        # Normalize the inputs
26        inputs_m, inputs_s = inputs.mean(), inputs.std()
27        inputs = (inputs - inputs_m) / inputs_s
28
29        # Zero the parameter gradients
```

```
30          optimizer.zero_grad()
31
32          # forward + backward + optimize
33          outputs = model(inputs)
34          labels = tuple_of_tensors_to_tensor(labels)
35
36          loss = criterion(outputs, labels)
37          loss.backward()
38          optimizer.step()
39          scheduler.step()
40
41          # Keep stats for Loss and Accuracy
42          running_loss += loss.item()
43
44          # Get the predicted class with the highest score
45          _, prediction = torch.max(outputs,1)
46
47          # Count of predictions that matched the target label
48          correct_prediction += (prediction == labels).sum().item()
49          total_prediction += prediction.shape[0]
50
51          #if i % 10 == 0:    # print every 10 mini-batches
52          #    print('[%d, %5d] loss: %.3f' % (epoch + 1, i + 1, running_loss / 10))
53
54      # Print stats at the end of the epoch
55      num_batches = len(train_dl)
56      avg_loss = running_loss / num_batches
57      acc = correct_prediction/total_prediction
58      print(f'Epoch: {epoch}, Loss: {avg_loss:.2f}, Accuracy: {acc:.2f}')
59
60  print('Finished Training')
61
62 num_epochs=200  # Just for demo, adjust this higher.
63 training(myModel, train_dl, num_epochs)
```

## Validation Part

We run inference loop taking care to disable the gradient updates. The forward pass is performed with the model to get predictions, but we don't need to backpropagate or run the optimizer.

```
1 def inference (model, val_dl):
2   correct_prediction = 0
3   total_prediction = 0
4
5   # Disable gradient updates
6   with torch.no_grad():
7     for data in val_dl:
8       # Get the input features and target labels, and put them on the GPU
9       inputs, labels = data[0], data[1]
10
11      # Normalize the inputs
12      inputs_m, inputs_s = inputs.mean(), inputs.std()
13      inputs = (inputs - inputs_m) / inputs_s
14
15      # Get predictions
```

```
15      # Get predictions
16      outputs = model(inputs)
17      print("Outputs-------------------")
18      print(outputs)
19
20      # Get the predicted class with the highest score
21      _, prediction = torch.max(outputs,1)
22      # Count of predictions that matched the target label
23      correct_prediction += (prediction == labels).sum().item()
24      total_prediction += prediction.shape[0]
25
26   acc = correct_prediction/total_prediction
27   print(f'Accuracy: {acc:.2f}')
28
29 # Run inference on trained model with the validation set
30 inference(myModel, val_dl)
```

# ▾ 9 Conclusions

As Training accuracy is higher than Validation accuracy, there are many opportunities for improvement.

This is a end-to-end case of song classification which is one of the foremost foundational issues in sound Deep learning. Not as it were is this utilized in a wide run of applications, but numerous of the concepts and procedures that we secured here will be pertinent to more complicated sound issues.

I Trained model for 200 Epoch and Got following results:

**Epoch: 119, Loss: 0.14, Traning Accuracy: ( 97% - 98% )**

**Validation Accuracy : ( 52% - 70% )**

( I included range of accuracies beacuse as traning dataset got different number of audio files from different songs, model get affected. If we increse number of audio files per song, this model will get improved.)

Scope for improvement :

- Dataset Size can be increased to cover more features.
- Deep Learning Network can be Improved. Adding more CNN layes can be one option.
- Focus on selecting desirable time durations from audio files can improve feature selection.
- Model building on single type of audio can be one option. ( Selecting only Whistling or Humming audios.