

Music Genre classification

ZZ

November 2023

1 GTZAN dataset

Data set description : Dataset description

The dataset used is GTZAN (the famous GTZAN dataset, the MNIST of sounds)

The GTZAN dataset contains 1000 audio files. Contains a total of 10 genres, each genre contains 100 audio files

1. Blues
2. Classical
3. Country
4. Disco
5. Hip-hop
6. Jazz
7. Metal
8. Pop
9. Reggae
10. Rock

****Genres original****

↳ A compilation of ten genres, each with 100 audio recordings, each lasting 30 seconds (the famous GTZAN dataset, the MNIST of sounds)

****Images original****

↳ Each audio file has a visual representation. Neural networks are one technique to classify data because they usually take in some form of picture representation.

****CSV files****

↳ The audio files' features are contained within. Each song lasts for 30 seconds long has a mean and variance computed across several features taken from an audio file in one file. The songs are separated into 3 second audio files in the other file, which has the same format.

*****Dive deep with me into the world of data! While the dataset has graciously done some of the heavy lifting by extracting certain features for us, we're about to embark on a thrilling journey to uncover the mysteries behind these features. Let's truly understand their essence and the magic of how to obtain them. So, why not take a bold step? Let's delete those pre-existing folders and craft our own path in this adventure!*****

2 properties of sounds

2.1 Wave Form

2.1.1 Definition

Sound is type of energy made by vibrations. When any object vibrates, it causes movement in the air =, these particles bump into the one close to them. This movement call sound waves.

A waveform represents the shape and form of a signal (like a sound) in terms of its amplitude over time

A waveform plot is a visual representation of an audio signal. It displays how the amplitude of the sound signal varies with time

Display Waveform :

```
1 import matplotlib.pyplot as plt
2
3 # Wave form of the audio
4 plt.figure(figsize=(12,6))
5 librosa.display.waveshow(data, color="#2B4F72", alpha
   = 0.5)
6 plt.show()
```

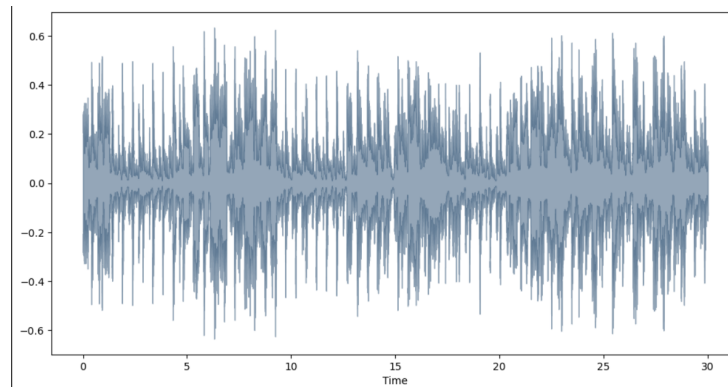


Figure 1: Here is the result

2.2 Sampling Rate

2.2.1 Definition

Sound, in its natural form, is a continuous analog signal. To enable digital processing of these signals, they must be converted into a discrete digital format through a process called 'sampling'. The 'sampling rate', quantified in Hertz

(Hz), denotes the number of samples taken per second from the continuous signal.

2.2.2 Example Code

```
1 data, sr = librosa.load(data_path)
2
3 plt.figure(figsize=(12,6))
4 librosa.display.waveshow(data, color="#2B4F72", alpha
5     = 0.5, sr=22050)
6 plt.show()
```

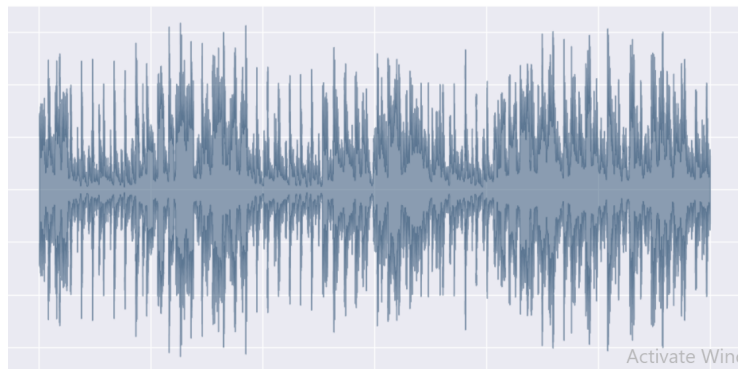


Figure 2: Here is the result

2.3 Root Mean Square Enenrgy

2.3.1 Definition

RMS-E used to measure the magnitude of "Volume" of an audio signal
Why is it important ?

RMS is a crucial metric in audio because our perception of loudness is more closely related to the RMS level of a sound rather than its peak level. Two sounds with the same peak level can have different perceived loudness if one has a higher RMS level. This is why RMS is often used in audio normalization and compression algorithms to ensure consistent perceived loudness across different audio tracks or segments.

```
1 # Let load .wav file with default sampling rate of
2   22,050Hz
3 data, sr = librosa.load(data_path)
4 # Compute RMS
```

```

5 rms = librosa.feature.rms(y=data)
6
7 # Plot RMS
8 plt.figure(figsize=(10, 4))
9 librosa.display.waveshow(data, sr=sr, alpha=0.4, label
10 = 'Audio Waveform')
11 plt.plot(librosa.times_like(rms[0], sr=sr), rms[0],
12 color='r', label='RMS')
13 plt.legend(loc='upper right')
14 plt.title('RMS over Time')
15 plt.tight_layout()
16 plt.show()

```

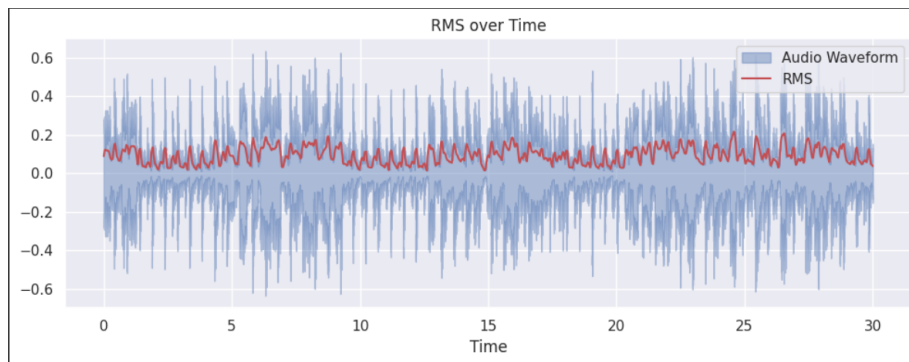


Figure 3: Here is the result

2.4 Zero-Crossing Rate

2.4.1 Definition

The Zero Crossing Rate (ZCR) is a fundamental feature in audio signal processing and analysis. It represents the rate at which a signal changes its sign, i.e., goes from positive to negative or vice versa.

This tell us :

- **Frequency Content:** A higher zero crossing rate often indicates the presence of higher frequencies in the signal, while a lower rate suggests lower frequencies. This makes sense because higher frequency components will cause the signal to oscillate (cross zero) more frequently.
- **Noise Detection:** Noise, especially of the high-frequency variety, can result in increased zero crossings. Thus, an unexpectedly high ZCR can sometimes be an indicator of noise.

- Speech Analysis: In speech processing, the ZCR can help distinguish between voiced and unvoiced speech segments. Voiced segments, like vowels, tend to have a lower ZCR, while unvoiced segments, like consonants such as "s" or "f", tend to have a higher ZCR.

```

1  ## Example of ZCR
2
3  # Let load .wav file with default sampling rate of
4  22,050Hz
5  data, sr = librosa.load(data_path)
6
7  # Calculate ZRC of the first 1000 data point of our
8  song
9  n0 = 0
10 n1 = 1000
11 plt.figure(figsize=(14, 5))
12 plt.plot(data[n0:n1])
13 plt.grid()
14
15 zero_crossings = librosa.zero_crossings(data[n0:n1],
16                                         pad=False)
17 print(f'ZRC = {sum(zero_crossings)}')
```

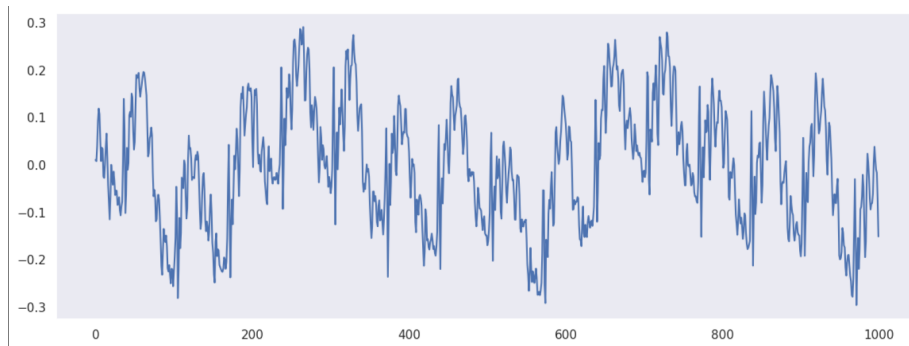


Figure 4: Here is the result

2.5 Special roll-off

2.5.1 Definition

Spectral roll-off is a measure used in audio signal processing and analysis to characterize the spectral shape of an audio signal. It represents a threshold below which a specified percentage of the total spectral energy is contained. The spectral roll-off point is typically defined as the frequency below which a

certain percentage (often 85 to 95 percent) of the total spectral energy lies. In other words, if you were to sum the magnitudes of the spectral components from the lowest frequency up to the roll-off frequency, it would account for, say, 85 percent of the total spectral energy of the signal.

```

1 import sklearn.preprocessing
2
3 # Function to normalize an array
4 def normalize(x, axis=0):
5     return sklearn.preprocessing.minmax_scale(x, axis=
        axis)
6
7 # Load the audio file with default sampling rate of
      22,050Hz
8 data, sr = librosa.load(data_path)
9
10 # Compute and norm the spectral roll-off
11 rolloff_85 = librosa.feature.spectral_rolloff(y=data,
        sr=sr, roll_percent=0.85)[0]
12 rolloff_85_norm = normalize(rolloff_85)
13
14 # Plot the waveform and normalized spectral roll-offs
15 plt.figure(figsize=(12, 6))
16 librosa.display.waveshow(data, sr=sr, alpha=0.4, label=
        'Waveform')
17 times = librosa.times_like(rolloff_85)
18 plt.plot(times, rolloff_85_norm, color='r', label='85%
        roll-off (normalized)')
19
20 plt.title('Waveform with Normalized Spectral Roll-offs
        ')
21 plt.legend(loc='upper right')
22 plt.tight_layout()
23 plt.show()

```

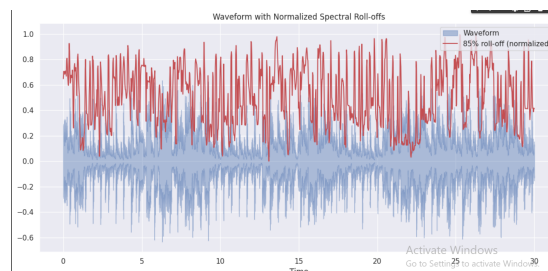


Figure 5: Here is the result

2.6 Spectral Centroid

2.6.1 Definition

The spectral centroid is a measure used in digital signal processing to characterize the "brightness" or "sharpness" of a sound. It provides a center of gravity or a "balance point" for the spectrum of a sound signal, indicating where the "center of mass" of the spectrum is located.

3 Music classification project

3.1 Initialize data set into feature_{3sec}.csv

```
1 from glob import glob
2 import pandas as pd
3
4 num_segment=10
5 num_mfcc=20
6 sample_rate=22050
7 n_fft=2048
8 hop_length=512
9
10 my_csv={"filename":[], "chroma_stft_mean": [], "
11         chroma_stft_var": [], "rms_mean": [], "rms_var":
12         [], "spectral_centroid_mean": [],,
13         "spectral_centroid_var": [],, "
14         spectral_bandwidth_mean": [],, "
15         spectral_bandwidth_var": [],, "rolloff_mean"
16         : [], "rolloff_var": [],,
17         "zero_crossing_rate_mean": [],, "
18         zero_crossing_rate_var": [],, "harmony_mean"
19         : [], "harmony_var": [],, "perceptr_mean":
20         [],,
21         "perceptr_var": [],, "tempo": [],, "mfcc1_mean":
22         [],, "mfcc1_var" : [],, "mfcc2_mean" : [],, "
23         mfcc2_var" : [],,
24         "mfcc3_mean" : [],, "mfcc3_var" : [],, "
25         mfcc4_mean" : [],, "mfcc4_var" : [],, "
26         mfcc5_mean" : [],,
27         "mfcc5_var" : [],, "mfcc6_mean" : [],, "
28         mfcc6_var" : [],, "mfcc7_mean" : [],, "
29         mfcc7_var" : [],,
30         "mfcc8_mean" : [],, "mfcc8_var" : [],, "
31         mfcc9_mean" : [],, "mfcc9_var" : [],, "
32         mfcc10_mean" : [],,
```

```

17         "mfcc10_var" : [], "mfcc11_mean" : [], "
           mfcc11_var" : [], "mfcc12_mean" : [], "
           mfcc12_var" : [],
18         "mfcc13_mean" : [], "mfcc13_var" : [], "
           mfcc14_mean" : [], "mfcc14_var" : [], "
           mfcc15_mean" : [],
19         "mfcc15_var" : [], "mfcc16_mean" : [], "
           mfcc16_var" : [], "mfcc17_mean" : [], "
           mfcc17_var" : [],
20         "mfcc18_mean" : [], "mfcc18_var" : [], "
           mfcc19_mean" : [], "mfcc19_var" : [], "
           mfcc20_mean" : [],
21         "mfcc20_var": [], "label": []}

```

With :

- num segment is number of subdivisions of original data
- num mfcc is the MFCC number system number to be calculated, in the GTZAN set is 20

3.2 find the features and Y

```

1 dataset_path="/content/GTZAN/genres_original"
2 audio_files = glob(dataset_path + "/*/*")
3 genre = glob(dataset_path + "/*")
4 n_genres=len(genre)
5 genre=[genre[x].split('/')[0] for x in range(n_genres
        )]

```

3.3 Calculate mean and var of each features and insert into feature 3 csv file

```

1 samples_per_segment = int(sample_rate*30/num_segment)
2
3 genre=""
4 for f in sorted(audio_files):
5     if genre!=f.split('/')[0]:
6         genre=f.split('/')[0]
7         print("Processing " + genre + "...")
8         fname=f.split('/')[0]
9         try:
10            y, sr = librosa.load(f, sr=sample_rate)
11        except:
12            continue

```



```

13
14     for n in range(num_segment):
15         y_seg = y[samples_per_segment*n:
16                 samples_per_segment*(n+1)]
17         #Chromagram
18         chroma_hop_length = 512
19         chromagram = librosa.feature.chroma_stft(y=
20             y_seg, sr=sample_rate, hop_length=
21             chroma_hop_length)
22         my_csv["chroma_stft_mean"].append(chromagram.
23             mean())
24         my_csv["chroma_stft_var"].append(chromagram.
25             var())
26
27         #Root Mean Square Energy
28         RMSEn= librosa.feature.rms(y=y_seg)
29         my_csv["rms_mean"].append(RMSEn.mean())
30         my_csv["rms_var"].append(RMSEn.var())
31
32         #Spectral Centroid
33         spec_cent=librosa.feature.spectral_centroid(y=
34             y_seg)
35         my_csv["spectral_centroid_mean"].append(
36             spec_cent.mean())
37         my_csv["spectral_centroid_var"].append(
38             spec_cent.var())
39
40         #Spectral Bandwith
41         spec_band=librosa.feature.spectral_bandwidth(y
42             =y_seg, sr=sample_rate)
43         my_csv["spectral_bandwidth_mean"].append(
44             spec_band.mean())
45         my_csv["spectral_bandwidth_var"].append(
46             spec_band.var())
47
48         #Rolloff
49         spec_roll=librosa.feature.spectral_rolloff(y=
50             y_seg, sr=sample_rate)
51         my_csv["rolloff_mean"].append(spec_roll.mean()
52             )
53         my_csv["rolloff_var"].append(spec_roll.var())
54
55         #Zero Crossing Rate
56         zero_crossing=librosa.feature.
57             zero_crossing_rate(y=y_seg)

```

```

44 my_csv["zero_crossing_rate_mean"].append(
    zero_crossing.mean())
45 my_csv["zero_crossing_rate_var"].append(
    zero_crossing.var())
46
47 #Harmonics and Perceptrueal
48 harmony, perceptr = librosa.effects.hpss(y=
    y_seg)
49 my_csv["harmony_mean"].append(harmony.mean())
50 my_csv["harmony_var"].append(harmony.var())
51 my_csv["perceptr_mean"].append(perceptr.mean(
    ))
52 my_csv["perceptr_var"].append(perceptr.var())
53
54 #Tempo
55 tempo, _ = librosa.beat.beat_track(y=y_seg, sr
    =sample_rate)
56 my_csv["tempo"].append(tempo)
57
58 #MFCC
59 mfcc=librosa.feature.mfcc(y=y_seg,sr=
    sample_rate, n_mfcc=num_mfcc, n_fft=n_fft,
    hop_length=hop_length)
60 mfcc=mfcc.T
61
62
63 fseg_name='.'.join(fname.split('.')[2])+f'_{n
    }.wav'
64 my_csv["filename"].append(fseg_name)
65 my_csv["label"].append(genre)
66 for x in range(20):
67     feat1 = "mfcc" + str(x+1) + "_mean"
68     feat2 = "mfcc" + str(x+1) + "_var"
69     my_csv[feat1].append(mfcc[:,x].mean())
70     my_csv[feat2].append(mfcc[:,x].var())
71 print(fname)
72
73 df = pd.DataFrame(my_csv)
74 df.to_csv('/content/GTZAN/features_3_sec.csv', index=
    False)

```

and then after that we get a table :

```
df = pd.read_csv("../content/GT2AN/features_3_sec.csv")
df.head()
```

| | filename | chroma_stft_mean | chroma_stft_var | rms_mean | rms_var | spectral_centroid_mean | spectral_centroid_var | spectral_bandwidth_mean | spect |
|---|-------------------|------------------|-----------------|----------|----------|------------------------|-----------------------|-------------------------|-------|
| 0 | blues.00000.0.wav | 0.335555 | 0.090997 | 0.130189 | 0.003559 | 1773.358004 | 169450.829520 | 1972.334258 | |
| 1 | blues.00000.1.wav | 0.343523 | 0.086782 | 0.112119 | 0.001491 | 1817.244034 | 90766.297254 | 2010.751494 | |
| 2 | blues.00000.2.wav | 0.347746 | 0.092495 | 0.130895 | 0.004552 | 1790.722357 | 110071.206973 | 2088.184750 | |
| 3 | blues.00000.3.wav | 0.363863 | 0.087207 | 0.131349 | 0.002338 | 1660.545231 | 109496.936296 | 1967.920582 | |
| 4 | blues.00000.4.wav | 0.335481 | 0.088482 | 0.142370 | 0.001734 | 1634.465077 | 77425.419232 | 1954.633566 | |

5 rows x 59 columns

Figure 6: Shape = 9990x59

3.4 Encoding y and X

3.4.1 Encoding y

```
1 from sklearn.preprocessing import LabelEncoder
2
3 # Label Encoding - encode the categorical classes with
4   numerical integer values for training
5
6 # Blues - 0
7 # Classical - 1
8 # Country - 2
9 # Disco - 3
10 # Hip-hop - 4
11 # Jazz - 5
12 # Metal - 6
13 # Pop - 7
14 # Reggae - 8
15 # Rock - 9
16
17 encoder=LabelEncoder()
18 y=encoder.fit_transform(y)
19 y
```

3.4.2 Encoding x

```
1 #scaling
2 from sklearn.preprocessing import StandardScaler
3 scaler=StandardScaler()
4 X=scaler.fit_transform(X)
```

it use (Z-score formular) to convert mean = 0, stardard deviation = 1

3.5 MLP model

library :

```

1 import os
2 import numpy as np
3
4 import torch
5 from torch import nn, optim
6 from torch.functional import F
7 from torch.utils.data import DataLoader, TensorDataset

```

Simple MLP model :

```

1 class MLP(nn.Module):
2     def __init__(self, input_size):
3         super(MLP, self).__init__()
4         self.flatten = nn.Flatten()
5         self.fc1 = nn.Linear(input_size, 512)
6         self.fc2 = nn.Linear(512, 256)
7         self.fc3 = nn.Linear(256, 128)
8         self.fc4 = nn.Linear(128, 64)
9         self.fc5 = nn.Linear(64, 32)
10        self.fc6 = nn.Linear(32, 10)
11        self.dropout = nn.Dropout(0.2)
12
13    def forward(self, x):
14        x = self.flatten(x)
15        x = F.relu(self.fc1(x))
16        x = self.dropout(x)
17        x = F.relu(self.fc2(x))
18        x = self.dropout(x)
19        x = F.relu(self.fc3(x))
20        x = self.dropout(x)
21        x = F.relu(self.fc4(x))
22        x = self.dropout(x)
23        x = F.relu(self.fc5(x))
24        x = self.dropout(x)
25        x = self.fc6(x)
26        return x

```

Initialize Hyperparameter :

```

1 input_size = X_train.shape[1]
2 model = MLP(input_size)
3
4 criterion = nn.CrossEntropyLoss()
5 optimizer = optim.Adam(model.parameters(), lr
6                        =0.000146)
7 num_epochs = 300
8 batch_size = 256

```

```

8
9 train_dataset = TensorDataset(torch.tensor(X_train),
    torch.tensor(y_train))
10 train_loader = DataLoader(train_dataset, batch_size=
    batch_size, shuffle=True)
11
12 val_dataset = TensorDataset(torch.tensor(X_test),
    torch.tensor(y_test))
13 val_loader = DataLoader(val_dataset, batch_size=
    batch_size, shuffle=False)
14
15 step = 0

```

Training :

```

1 for epoch in range(num_epochs):
2     model.train()
3     for inputs, labels in train_loader:
4         optimizer.zero_grad()
5         outputs = model(inputs.float())
6         loss = criterion(outputs, labels)
7         loss.backward()
8         optimizer.step()
9
10    if step % 100 == 0:
11        print(f"Step {step}, Train Loss: {loss.item
            ():.4f}")
12        step += 1
13
14    # Validation
15    model.eval()
16    val_loss = 0.0
17    correct = 0
18    total = 0
19    with torch.no_grad():
20        for inputs, labels in val_loader:
21            outputs = model(inputs.float())
22            loss = criterion(outputs, labels)
23            val_loss += loss.item()
24            _, predicted = outputs.max(1)
25            total += labels.size(0)
26            correct += predicted.eq(labels).sum().item
                ()
27
28    val_loss /= len(val_loader)
29    val_accuracy = 100 * correct / total

```

```
30 print(f"Epoch {epoch+1}/{num_epochs}, Validation  
    Loss: {val_loss:.4f}, Validation Accuracy: {  
    val_accuracy:.2f}%")
```