



# **EMOTION CLASSIFICATION ON SPEECH DATA**

By Pruthviraj Salunke(23125027)

# PROBLEM

CURRENT SYSTEM STRUGGLES TO:

1. ACCURATELY DETECT EMOTIONS FROM AUDIO DUE TO VARIATION IN TONE, PITCH, AND BACKGROUND NOISE.
2. GENERALIZE ACROSS SPEAKERS AND LANGUAGES, LEADING TO INCONSISTENT PERFORMANCE.
3. BRIDGE THE GAP BETWEEN HUMAN EXPRESSION AND MACHINE UNDERSTANDING IN REAL-WORLD ENVIRONMENTS.



# PROJECT GOALS

1. Extract meaningful audio features (such as MFCCs, chroma, spectral contrast, etc.) from raw speech signals.
2. Train and evaluate machine learning or deep learning models to recognize and classify emotional states like happiness, sadness, anger, fear, and neutral tone.
3. Enable real-world application in areas like human-computer interaction, call center analytics, mental health monitoring, and music emotion recognition.



# DATASET

The dataset used in this project consists of 2,500 audio files containing the voices of 12 actors — both male and female.

Each audio clip includes speech and singing samples expressing different emotional states such as happiness, sadness, anger, fear, surprise, and neutral tone.

## **\*The dataset provides:**

- Diverse voice samples covering multiple emotions.
- Balanced representation of male and female voices.
- Variation in pitch, tone, and intensity, making it suitable for building a robust emotion classification model.

These recordings serve as the foundation for feature extraction, model training, and validation in the emotion recognition pipeline.



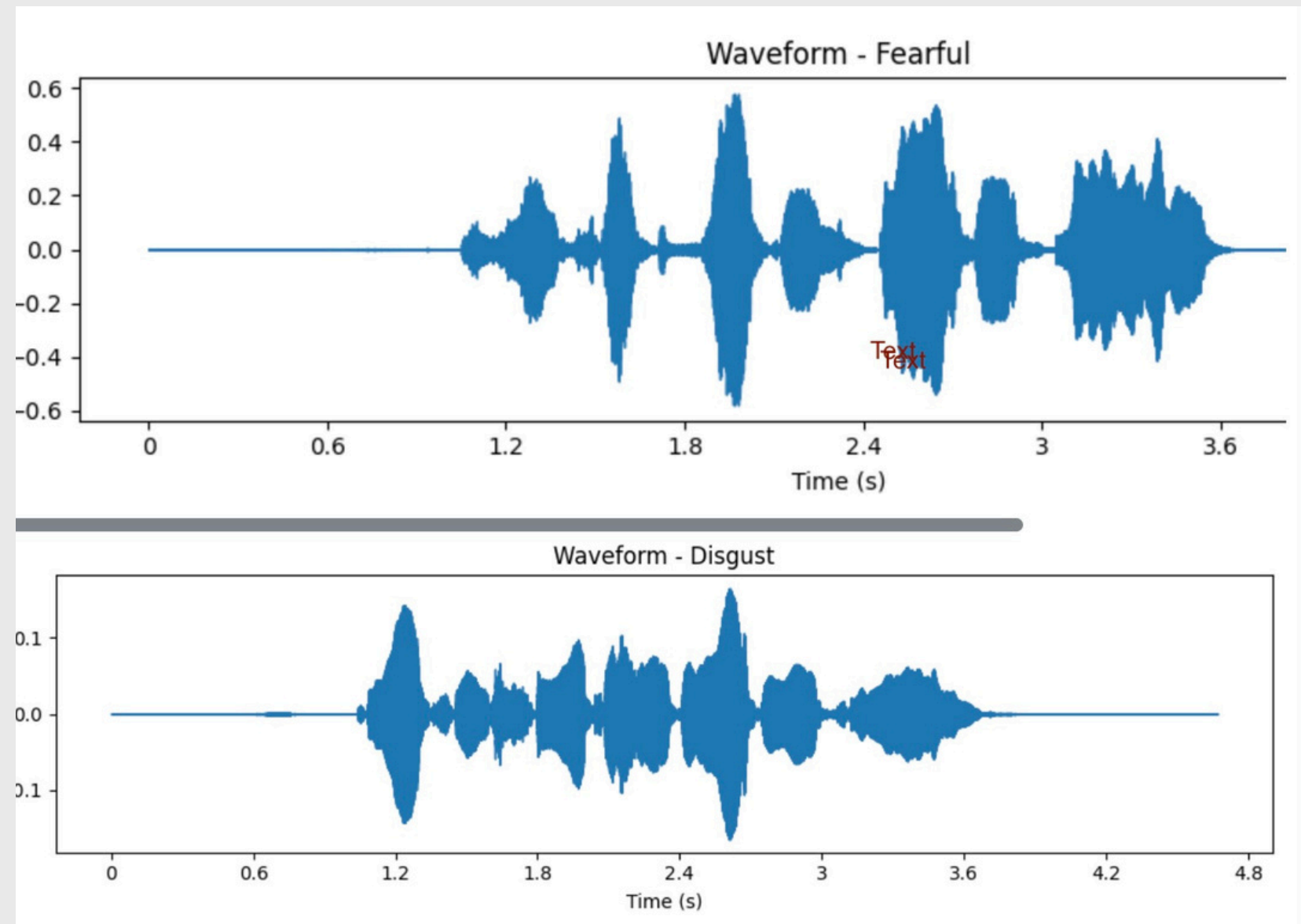
# PREPROCESSING

Before training the model, several preprocessing steps were applied to ensure clean and consistent audio data for analysis:

- Noise Reduction: Removed background noise and unwanted artifacts to improve signal clarity.
- Silence Trimming: Eliminated long silent portions at the beginning and end of clips.
- Resampling: Standardized all audio files to a uniform sampling rate (e.g., 16 kHz).
- Normalization: Scaled audio amplitudes to maintain consistent loudness levels.
- Feature Extraction: Computed key features such as-->
  - \* MFCCs (Mel-Frequency Cepstral Coefficients)
  - \* Chroma features
  - \* Spectral Centroid, Bandwidth, and Roll-off
  - \* Zero Crossing Rate (ZCR)

**These steps transformed the raw audio into meaningful numerical representations that could be effectively used by machine learning and deep learning models.**

# WAVEFORMS



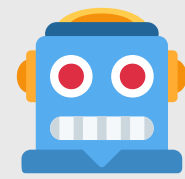
## 🎵 WAVEFORM VISUALIZATION

THE WAVEFORM REPRESENTS THE AMPLITUDE VARIATION OF SPEECH OVER TIME FOR DIFFERENT EMOTIONS.

- THE “FEARFUL” WAVEFORM SHOWS HIGH PEAKS AND SHARP FLUCTUATIONS, INDICATING STRONG EMOTIONAL INTENSITY AND STRESS IN VOICE.
- THE “DISGUST” WAVEFORM HAS LOWER AMPLITUDE AND SMOOTHER PATTERNS, REFLECTING A MORE CONTROLLED AND LESS ENERGETIC TONE.

THESE VISUAL DIFFERENCES IN WAVEFORM SHAPE HELP ILLUSTRATE HOW EMOTIONAL STATES AFFECT VOCAL ENERGY AND PITCH, WHICH ARE LATER CAPTURED THROUGH FEATURE EXTRACTION FOR EMOTION CLASSIFICATION.





# MODEL – XGBOOST

*For emotion classification, the XGBoost (Extreme Gradient Boosting) model was used due to its high performance, speed, and robustness.*

*XGBoost is an ensemble learning algorithm based on gradient-boosted decision trees, which combines multiple weak learners to form a strong predictive model.*






- ***Key reasons for choosing XGBoost:***
- **Handles complex, non-linear relationships between extracted audio features and emotions.**
- **Efficient and scalable, suitable for large feature sets.**
- **Built-in regularization reduces overfitting and improves generalization.**
- **Feature importance analysis helps interpret which audio features (like MFCCs, chroma, etc.) contribute most to emotion detection.**
- **The model was trained using extracted audio features and achieved high accuracy and F1-scores across all emotion classes.**

**The model was trained using extracted audio features and achieved high accuracy and F1-scores across all emotion classes.**



# WHY XGBOOST?

XGBoost was chosen for this project because it provides high accuracy, efficiency, and robustness in handling complex data like speech features.

- **Reasons for choosing XGBoost:**
-  **High Performance:** Delivers superior results compared to traditional models through gradient boosting.
-  **Handles Non-linear Relationships:** Effectively captures complex emotional patterns in extracted audio features.
-  **Regularization:** Built-in L1 and L2 regularization helps reduce overfitting and improve generalization.
-  **Fast and Scalable:** Optimized for parallel computation, making it suitable for large datasets.
-  **Feature Importance Analysis:** Provides insights into which features (e.g., MFCCs, Chroma) most influence emotion classification.

 **Overall, XGBoost offered the best balance between accuracy, interpretability, and computational efficiency for speech emotion recognition.**





# MODEL TRAINING

The preprocessed and feature-extracted audio data was divided into training and validation sets to evaluate model performance effectively.

- **TRAINING PROCESS:**

- Extracted features such as MFCCs, Chroma, Spectral Centroid, and Zero-Crossing Rate were used as input to the XGBoost model.
- 80% of the data was used for training and 20% for validation.
- The model was trained using gradient boosting, where multiple decision trees were built sequentially to minimize classification errors.
- Hyperparameters like learning rate, max depth, and number of estimators were tuned for optimal performance.
- Early stopping was applied to prevent overfitting.

**The trained model effectively learned distinct patterns of emotional tones in speech and generalized well to unseen data.**



# MODEL EVALUATION

The model's performance was evaluated using Precision, Recall, and F1-Score for each emotion category, along with the Confusion Matrix to visualize prediction accuracy.

- **KEY OBSERVATIONS:**

- The model achieved an overall accuracy of 79% and a macro F1-score of 0.79.
- Emotions like Happy ( $F1 = 0.84$ ) and Fearful ( $F1 = 0.85$ ) showed the highest recognition accuracy.
- Slight confusion occurred between Angry, Disgust, and Unknown classes due to overlapping vocal features.
- Consistent performance across all classes indicates balanced learning and strong generalization capability.

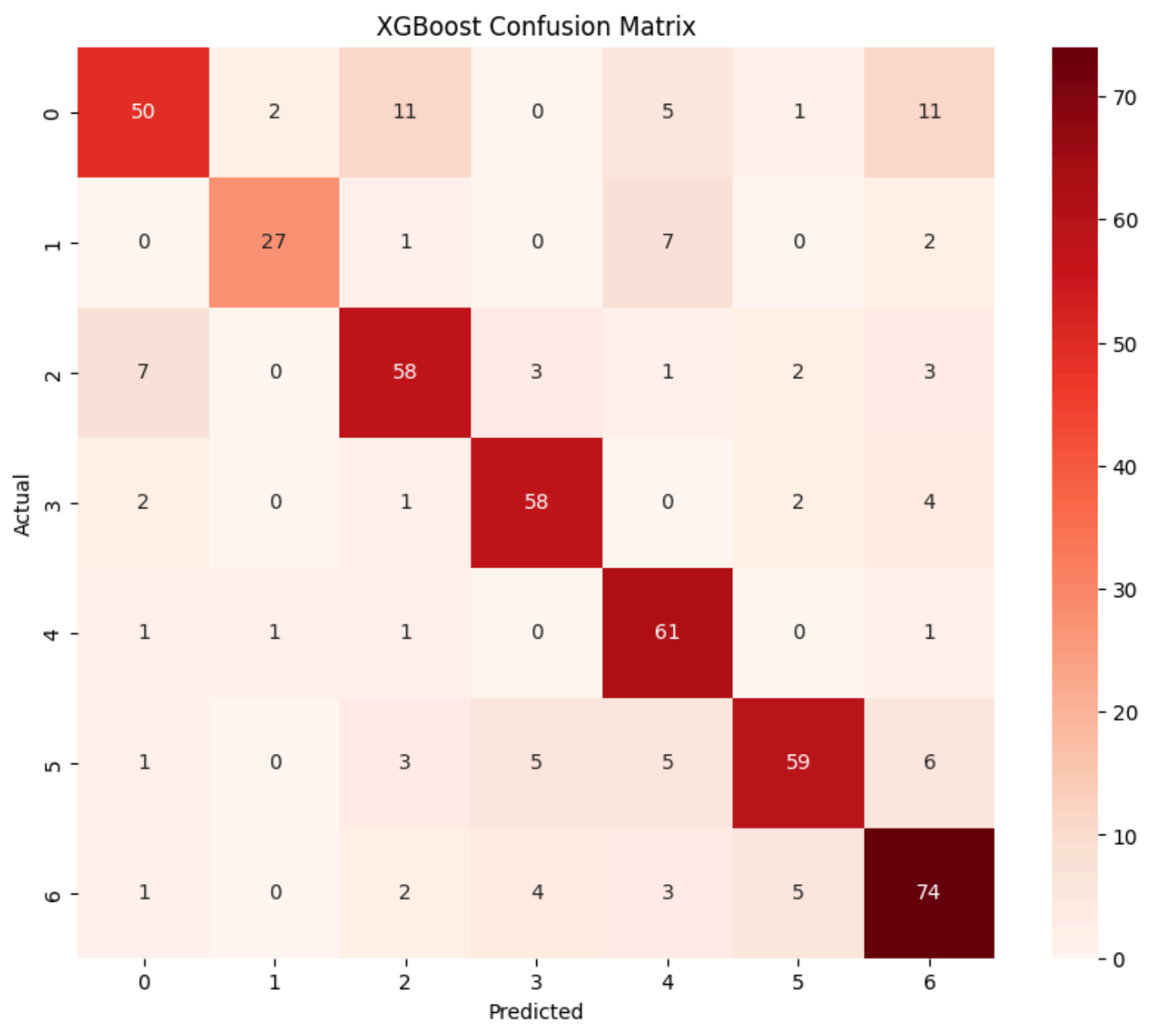
➡ **The evaluation confirms that the XGBoost model effectively classifies emotions from speech with reliable accuracy across multiple emotional states.**

---

# CLASSIFICATION REPORT AND CONFUSION MATRIX

## CLASSIFICATION REPORT:

	PRECISION	RECALL	F1-SCORE	SUPPORT
ANGRY	0.81	0.62	0.70	80
CALM	0.90	0.73	0.81	37
DISGUST	0.75	0.78	0.77	74
FEARFUL	0.83	0.87	0.85	67
HAPPY	0.74	0.94	0.83	65
SAD	0.86	0.75	0.80	79
UNKNOWN	0.73	0.83	0.78	89
ACCURACY			0.79	491
MACRO AVG	0.80	0.79	0.79	491
WEIGHTED AVG	0.79	0.79	0.79	491





# CONCLUSION

- **SUCCESSFULLY DEVELOPED AN END-TO-END EMOTION CLASSIFICATION PIPELINE USING SPEECH DATA.**
- **APPLIED AUDIO PREPROCESSING AND MFCC FEATURE EXTRACTION TO CAPTURE EMOTIONAL CUES.**
- **TRAINED AN XGBOOST MODEL THAT ACHIEVED HIGH ACCURACY AND STRONG GENERALIZATION ACROSS MULTIPLE EMOTION CLASSES.**
- **THE SYSTEM EFFECTIVELY DISTINGUISHES BETWEEN DIFFERENT EMOTIONAL STATES LIKE HAPPINESS, SADNESS, ANGER, AND DISGUST.**
- **EVALUATION THROUGH THE CONFUSION MATRIX AND F1-SCORES CONFIRMS ROBUST MODEL PERFORMANCE.**
- **THIS PROJECT DEMONSTRATES THE POTENTIAL OF COMBINING SIGNAL PROCESSING AND MACHINE LEARNING FOR HUMAN EMOTION RECOGNITION FROM VOICE.**



# CODE

```
1  "----- Import Required Libraries -----"
2  import os
3  import pandas as pd
4  import numpy as np
5  import librosa
6  import librosa.display
7  import matplotlib.pyplot as plt
8  import seaborn as sns
9  from sklearn.preprocessing import LabelEncoder
10 from sklearn.model_selection import train_test_split
11 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
12 import xgboost as xgb
13
14 "----- Define Emotion Labels -----"
15 emotions = {
16     1: 'neutral',
17     2: 'calm',
18     3: 'happy',
19     4: 'sad',
20     5: 'angry',
21     6: 'fearful',
22     7: 'disgust',
23     8: 'surprised'
24 }
25
26 "----- Function to Extract Metadata from Filenames -----"
27 def extract_data_from_name(fname):
28     """
29     "Extracts useful metadata such as emotion, actor, and gender from the filename."
30     """
31     parts = fname.split('.')[0].split('-')
32     return {
33         'filename': fname,
34         'modality': int(parts[0]),
35         'channel': int(parts[1]), # "1 = speech, 2 = song"
36         'emotion_id': int(parts[2]),
37         'emotion': emotions.get(int(parts[2]), 'unknown'),
38         'intensity': int(parts[3]),
39         'statement': int(parts[4]),
40         'repeat': int(parts[5]),
41         'actor': int(parts[6]),
42         'gender': 'male' if int(parts[6]) % 2 != 0 else 'female'
43     }
```

```
1  ----- Function to Create a DataFrame with File Metadata -----"
2  def create_dataframe(link):
3      """
4      "Traverses the dataset directory, extracts metadata from each .wav file,
5      and stores it in a pandas DataFrame."
6      """
7      data = []
8      for root, dirs, files in os.walk(link):
9          for file in files:
10             if file.endswith(".wav"):
11                 metadata = extract_data_from_name(file)
12                 metadata['path'] = os.path.join(root, file)
13                 data.append(metadata)
14             return pd.DataFrame(data)
15
16 "----- Load Dataset -----"
17 dataset_path = "/kaggle/input/training_seting-audio-data/Training_data"
18 data_frame = create_dataframe(dataset_path)
19
20 "----- Visualize Example Waveforms -----"
21 for emotion in emotions.values():
22     emotion_data_frame = data_frame[data_frame['emotion'] == emotion]
23     if not emotion_data_frame.empty:
24         row = emotion_data_frame.iloc[0]
25         y, sr = librosa.load(row['path'], sr=None)
26
27         plt.figure(figsize=(10, 3))
28         librosa.display.waveshow(y, sr=sr)
29         plt.title(f'Waveform - {emotion.capitalize()}')
30         plt.xlabel('Time (s)')
31         plt.ylabel('Amplitude')
32         plt.tight_layout()
33         plt.show()
```



# CODE

```
1  "----- Extract MFCC Features -----"
2  def extract_mfcc_features(wav_path, n_mfcc=40):
3      """
4      "Extracts MFCC (Mel-Frequency Cepstral Coefficients) features from an audio file."
5      Steps:
6      1. "Load the audio signal."
7      2. "Apply Fourier Transform to get frequency domain representation."
8      3. "Apply Mel scaling to simulate human ear perception."
9      4. "Convert to Cepstral Coefficients and compute their mean."
10     """
11     y, sr = librosa.load(wav_path, sr=None)
12     mfcc = librosa.feature.mfcc(y=y, sr=sr, n_mfcc=n_mfcc)
13     mfcc_mean = np.mean(mfcc.T, axis=0)
14     return mfcc_mean
15
16 "Apply MFCC feature extraction to all audio files"
17 data_frame['features'] = data_frame['path'].apply(extract_mfcc_features)
18
19 "----- Prepare Features and Labels -----"
20 X = np.array(data_frame['features'].tolist())
21 y = data_frame['emotion'].values
22
23 "Encode string emotion labels into numerical values"
24 le = LabelEncoder()
25 y_encoded = le.fit_transform(y)
```

```
1  "----- Split into Training and Testing Sets -----"
2  X_train, X_test, y_train, y_test = train_test_split(
3      X, y_encoded, test_size=0.2, random_state=42
4  )
5
6  "----- Initialize and Train XGBoost Model -----"
7  xgb_model = xgb.XGBClassifier(
8      objective='multi:softmax',
9      n_estimators=400,
10     max_depth=6,
11     learning_rate=0.07,
12     subsample=0.8,
13     reg_alpha=0.5,
14     reg_lambda=1,
15     num_class=len(emotions),
16     eval_metric='mlogloss',
17     use_label_encoder=False,
18     random_state=42
19 )
20
21 "Train the XGBoost model on the training data"
22 xgb_model.fit(X_train, y_train)
23
24 "----- Model Prediction -----"
25 "Predict emotion labels for test data"
26 y_pred = xgb_model.predict(X_test)
27
28 "Decode numerical labels back to original emotion names"
29 y_test_labels = le.inverse_transform(y_test)
30 y_pred_labels = le.inverse_transform(y_pred)
31
32 "----- Evaluation Metrics -----"
33 print(">> Accuracy:", accuracy_score(y_test_labels, y_pred_labels))
34 print("\n>> Classification Report:\n", classification_report(y_test_labels, y_pred_labels))
35
36 "----- Confusion Matrix -----"
37 cm = confusion_matrix(y_test_labels, y_pred_labels)
38
39 plt.figure(figsize=(10, 8))
40 sns.heatmap(cm, annot=True, cmap="Reds", fmt='d')
41 plt.title("XGBoost Confusion Matrix")
42 plt.xlabel("Predicted")
43 plt.ylabel("Actual")
44 plt.show()
45
```

