

Speech Emotion Recognition - Week 2

Data Exploration & Feature Engineering

Krishna Kukreja & Garvit Meena

November 2025

Week 2: Resources & Guide

1. Learning Objectives

- Load the entire RAVDESS dataset systematically.
- Extract features from all audio samples.
- Perform Exploratory Data Analysis (EDA).
- Understand feature distributions across emotions.
- Create the complete feature matrix (CSV) for Machine Learning.

2. Reading Materials

- **Feature Engineering:** [Why Mel Spectrograms Perform Better](#)
- **Audio Features:** [Music Information Retrieval - Feature Types](#)
- **EDA Guide:** [Kaggle - Data Visualization Course](#)

3. Video Tutorials

- [Audio Feature Extraction - Complete Walkthrough](#)
- [Python for Data Analysis Playlist](#)

4. Reference Projects

- [Complete SER Implementation on GitHub](#)
- [ProjectPro - SER Tutorial](#)
- [TechVidvan - ML Speech Emotion Recognition](#)

5. Python Code Templates

Batch Feature Extraction

```
import os
import librosa
import numpy as np
import pandas as pd

def extract_features(audio_path):
    y, sr = librosa.load(audio_path, duration=3)

    # MFCC
    mfcc = librosa.feature.mfcc(y=y, sr=sr, n_mfcc=13)
    mfcc_mean = np.mean(mfcc, axis=1)
    mfcc_std = np.std(mfcc, axis=1)

    # Spectral features
    spectral_centroid = np.mean(
        librosa.feature.spectral_centroid(y=y, sr=sr))
    spectral_rolloff = np.mean(
        librosa.feature.spectral_rolloff(y=y, sr=sr))
    zero_crossing_rate = np.mean(
        librosa.feature.zero_crossing_rate(y))

    features = np.hstack([mfcc_mean, mfcc_std,
                          spectral_centroid, spectral_rolloff,
                          zero_crossing_rate])
    return features

# Process all files
all_features = []
all_labels = []

# Ensure your dataset is in data/raw/RAVDESS/
for actor_folder in os.listdir('data/raw/RAVDESS/'):
    actor_path = os.path.join('data/raw/RAVDESS/', actor_folder)
    for filename in os.listdir(actor_path):
        if filename.endswith('.wav'):
            file_path = os.path.join(actor_path, filename)
            # Extract emotion from filename (3rd part of identifier)
            emotion = int(filename.split('-')[2])
            features = extract_features(file_path)
            all_features.append(features)
            all_labels.append(emotion)

df = pd.DataFrame(all_features)
df['emotion'] = all_labels
df.to_csv('features.csv', index=False)
```

EDA Visualization

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

df = pd.read_csv('features.csv')

# Emotion distribution
plt.figure(figsize=(8, 6))
df['emotion'].value_counts().sort_index().plot(kind='bar')
plt.title('Emotion Distribution')
plt.xlabel('Emotion')
plt.ylabel('Count')
plt.show()

# Feature correlation heatmap
plt.figure(figsize=(12, 10))
correlation = df.corr()
sns.heatmap(correlation, cmap='coolwarm', center=0)
plt.title('Feature Correlation Heatmap')
plt.tight_layout()
plt.show()
```

6. Assignments

1. Build complete feature extraction pipeline for all RAVDESS files.
2. Extract MFCC, spectral, and temporal features; save to CSV.
3. Perform EDA: create distribution plots for each emotion.
4. Generate correlation heatmap and identify top 5 most discriminative features.
5. Write 2-3 page EDA report with visualizations.

7. Expected Outputs

- Feature matrix CSV file (total samples \times 100+ features).
- EDA Jupyter notebook with comprehensive analysis.
- 5+ visualization plots (distributions, correlations, box plots).
- Written EDA report (2-3 pages) with insights.

8. Quick Links

- [Pandas Visualization Guide](#)
- [Seaborn Tutorial](#)
- [Librosa Feature Extraction Reference](#)