
Sentiment Analysis from Audio Recordings

Muhammad Qasim¹
(SP – 2022 – BSCS – 048)

Awais Ali¹
(SP – 2022 – BSCS – 080)

Sohaib Ahmad¹
(SP – 2022 – BSCS – 060)

Hammad Ahmad¹
(SP – 2022 – BSCS – 063)

¹Department of Computer Science, Lahore Garrison University, Pakistan

Abstract

The natural language processing field has established sentiment analysis as its subdiscipline which uses machine learning models. The practice of sentiment analysis interprets emotional expressions that appear in text, audio and video data. The project analyses speech emotions through classification methods to study sentiment analysis from audio content. Audio sentiment analysis proves valuable across feedback processing and social media sentiment tracking and audience monitoring in real-time applications. We employ the CREMA-D audio dataset containing 7442 labelled emotional clips to implement a successful approach targeting this issue. The process implements filter-based noise elimination followed by MFCC feature extraction alongside PCA dimension reduction and ends with machine learning classifier deployment. Accuracy, precision, recall, and F1-score are the factors on which comparison of the following classifiers was made: KNN-classifier, Logistic Regression-classifier, SVM-classifier, Naïve Bayes-classifier, and Neural Networks-classifier. The SVM-classifier produced 45.39% accuracy as the highest performing method whereas Logistic Regression-classifier followed with 43.69% accuracy. The moderate performance rates achieved by KNN and Naïve Bayes alongside Neural Networks suffered from data complexity issues alongside emotional intensity variations. The experimental findings both encourage research on MFCC-based features together with machine learning technology in speech emotion analysis and open more applications for this research area.

1. Introduction

This subarea interprets the emotional tones communicated from a few data sources: speech and text. They have largely been directed toward text data where emotions are deduced from indirectly deduced from the words used within the sentence, ignoring the many affective pieces of information spoken within speech. Pitch, intonation, tempo, and rhythm are among the acoustic features that carry the most emotive context into analysis, that would otherwise be lost in text only. This paper seeks to fill this gap in sentiment analysis by proposing a method for identifying and interpreting emotional tones from audio data.

The advantage from an audio context is that sentiment analysis is able to carry more extensive and revealing indicators of emotions, as compared to text-based sentiment analysis. In its own simple existence, audio retains tonal intonations, voice inflections and rhythmic cadences of speech, and as such exposes much more about the emotional state of a speaker than text does, as much of the text is lost in clarity and emotive basis. What makes audio sentiment analysis so valuable across so many applications - from customer service, social media monitoring to the entertainment industry - is the availability of these affective cues. In this paper, we present a step-by-step methodology for audio sentiment analysis through extraction of important acoustic features, and machine learning classification.

The proposed method then pre-processes the audio data such as removing noise from the data and improve the quality of the data. The unwanted background noises as recorded in the motion capture system have been removed using techniques such as low pass filtering. Well anyways we then have the cleaning data that we use to extract

features which have some relationship with audio signals, using techniques like Mel-Frequency Cepstral Coefficients (MFCCs), Fourier Transform, Mel-spectrogram and Chromagram. We extract the features described above and train a machine learning classifier for sentiment analysis.

We test several machine learning models: We used Logistic Regression, Naïve Bayes, Support Vector Machines (SVM), K-Nearest Neighbours (KNN) and Neural Networks to find which one works best at emotion classification.

For this paper, we take advantage of the highly popular CREMA-D dataset for research in emotion recognition. It consists of 7,442 audio files read by 91 actors, male and female of different ethnicities. The six emotions acted are Anger, Disgust, Fear, Happy, Neutral, and Sad, each of which being pronounced at four intensity levels: High, Medium, Low, and Unspecified. The audio data provided by the database is broad and diverse enough to be used for training and testing sentiment analysis models really well. Audio sentiment analysis is a big deal. It is possible to make the models perform better under other than their training environments, by speaker profile and intensity variety. In addition, the analysis identifies many practical applications. They apply sentiment analysis in call monitoring to check the customer satisfaction so that they can monitor whether customers are satisfied or not and the cause for the problem of dissatisfied customers is rectified. For example, in the field of social media analytics, it is used to measure the public's emotional response to a particular topic, brand or event in order for marketers and content providers to gauge that response. In the entertainment sector for example, we also look at what the audience thinks about the films, television shows, etc. so that the producers can fashion their content accordingly as per the thematic interests and emotional needs of the audience. When we have a good measure of the sentiment around an audience, it allows content producers to increase viewer engagement, which leads to better customer satisfaction, and better revenue.

The rest of this paper is divided into the following sections. In Section 2: Mathematical Formulation, we provide the theoretical background and mathematical concepts of the preprocessing, feature extraction, and classification techniques used in this study. Section 3: Identification and Extraction of Features explains in detail the specific audio features that have been extracted from the CREMA-D database, including MFCCs, Mel-spectrograms, and Chromagram, among others. In Section 4: Feature Engineering (Optional), we outline possible techniques for reducing dimensions in order to increase model efficiency and effectiveness. Section 5: Use of Different Classification Algorithms has subheadings for each classifier used in this research—Logistic Regression, Naive Bayes, Support Vector Machines (SVM), K-Nearest Neighbours (KNN), and Neural Networks with implementations, advantages, and disadvantages. Section 6: Performance Evaluation reports the results of our experiments, with performance metrics, plots, and tables charting the effectiveness of each classification model. The qualitative analysis that emerges is rich and allows for a deep understanding of every approach with regard to its strengths and weaknesses. Section 7: Conclusion summarizes key results, discusses implications of findings, and points to possible future directions of work in audio sentiment analysis.

2. Mathematical Formulation

2.1. Preprocessing: Low – Pass Filtering

To remove high-frequency noise from the audio signals, a Butterworth low – pass filter is applied:

Transfer Function:

$$H(s) = \frac{1}{1 + \left(\frac{s}{\omega_c}\right)^{2n}}$$

where:

- ω_c : Cutoff angular frequency ($2\pi f_c$).
- n : Filter Order.
- s : Frequency Variable

Digital Implementation: Using bilinear transformation, the digital filter coefficients b and a are calculated:

$$y[n] = \sum_{i=0}^M b_i x[n-i] - \sum_{j=1}^N a_j y[n-j]$$

where $x[n]$ is the input signal, $y[n]$ is the output signal, and M, N depend on the filter order.

2.2. Feature Extraction

From each filtered audio file, the following features are computed:

a) MFCC (Mel Frequency Cepstral Coefficients):

$$C_m = \sum_{k=1}^K \log(S_k) \cos \left[m \cdot \left(k - \frac{1}{2} \right) \cdot \frac{\pi}{K} \right]$$

where:

- S_k : Power Spectrum of the signal.
- K : Number of frequency bins.
- m : Number of cepstral coefficients.

b) Chrome Features:

$$C_h[k] = \sum_{m=0}^{M-1} P_h[m, k]$$

where P_h is the chromagram matrix, capturing energy distribution across 12 pitch classes.

c) Mel Spectrogram:

$$M(f, t) = \text{MelfilterBank} \cdot \text{STFT}(f, t)$$

where STFT is the Short-Time Fourier Transform.

d) Special Features:

- Centroid: $\mu_c = \frac{\sum_f f \cdot S(f)}{\sum_f S(f)}$
- Bandwidth: $B = \sqrt{\frac{\sum_f (f - \mu_c)^2 \cdot S(f)}{\sum_f S(f)}}$
- Rolloff: $R = \min_f (\sum_f S(f) \geq 0.85 \cdot \text{Total Energy})$

2.3. Dimensionality Reduction

Principal Component Analysis (PCA) is applied to reduce dimensionality while retaining 95% of the variance:

a) Covariance Matrix:

$$\Sigma = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T$$

where x_i is the feature vector, μ is the mean vector, and N is the number of samples.

b) Eigen Decomposition: Solve:

$$\sum v = \lambda v$$

For eigenvalues λ and eigenvectors v .

c) Projection: Retain top k components:

$$X_{\text{reduced}} = X V_k$$

2.4. Classification Algorithms

The following classifiers are implemented:

a) K – Nearest Neighbours (KNN):

Distance between a test point and training points is computed as:

$$d(x, x') = \sqrt{\sum_{i=1}^n (x_i - x'_i)^2}$$

b) Logistic Regression: Model output is given by:

$$P(y|x) = \frac{1}{1 + e^{-(w^T x + b)}}$$

c) Naive Bayes: Assumes conditional independence:

$$P(y|x) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

d) Support Vector Machine (SVM): Optimizes:

$$\min_{w,b} \frac{1}{2} \|w\|^2 \text{ subject to } y_i(w^T x_i + b) \geq 1$$

e) Neural Networks: Uses:

$$z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)}, \quad a^{(l)} = \sigma(z^{(l)})$$

where σ is an activation function, and l is the layer index.

3. Identification and Extraction of Features

In this chapter, we go through the literature and identify characteristics that are suitable for synthetic speech attribution. Features identified and utilized within this project were based on existing studies and practical requirements.

a) Butterworth Low-Pass Filtering

This tends to contain some amount of high frequency artifacts not common in natural speech. The filter used here was Butterworth low pass type with cutoff at 3400 Hz to get rid of unwanted artifacts while saving the important parts with frequency content common in human speech.

- **Implementation:** A Butterworth filter was selected because of the maximally flat frequency response in the passband. It avoids the introduction of distortion.

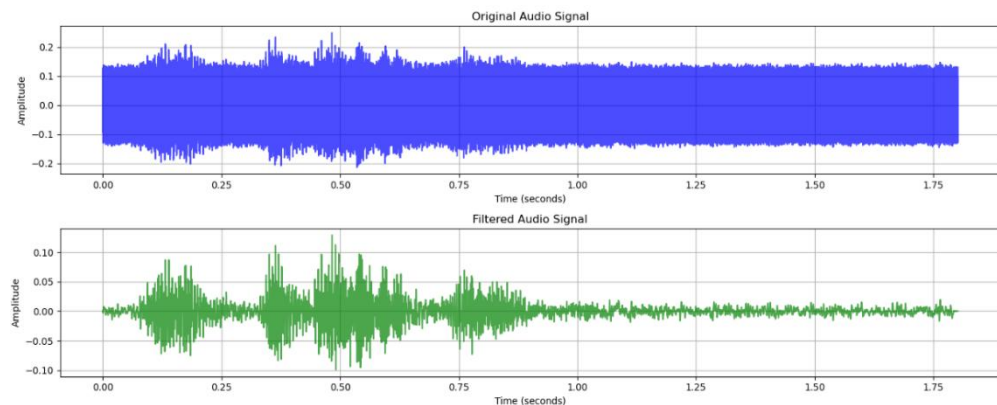


Figure 1: Butterworth Low Pass Filter Visual

b) Mel Frequency Cepstral Coefficients (MFCCs):

MFCCs capture the shape of the spectral envelope, which is critical for distinguishing synthetic speech from natural speech. It also simulates the human auditory system by emphasizing perceptually important features.

- **Implementation:** For each audio sample, 13 MFCCs have been obtained as the number of MFCCs used for standardization is normative.

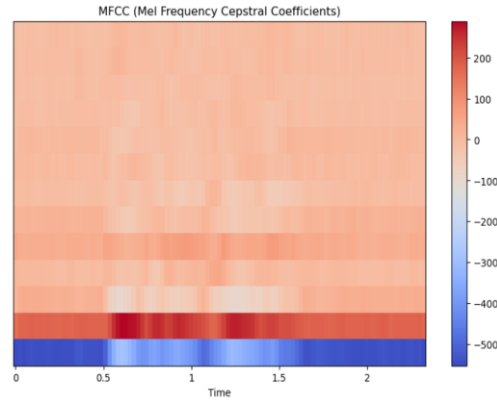


Figure 2: MFCC Visual

c) **Mel Spectrogram:**

The Mel spectrogram gives a time-frequency representation of the audio signal mapped to the Mel scale. This captures both the spectral and temporal properties, which are very useful in speech analysis.

- **Implementation:** Mel spectrogram is computed with the window size and hop length such that temporal resolution and spectral resolution are balanced.

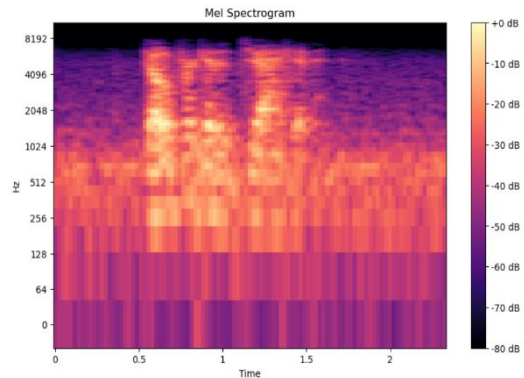


Figure 3: Mel Spectrogram Visual

a) **Chroma Features:**

Chroma features summarize the distribution of energy over the pitch classes (12 semitones of the chromatic scale). These features enable the tonal analysis of speech, which could be different in synthetic and natural speech.

- **Implementation:** Chroma features were obtained through STFT representation.

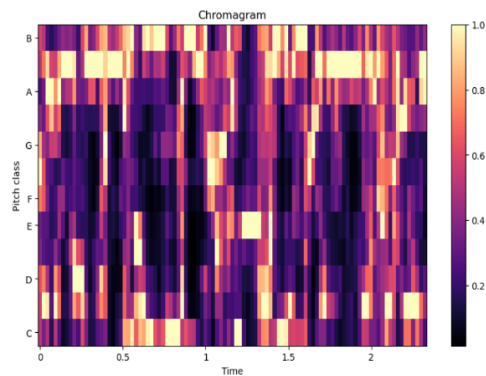


Figure 4: Chromagram Visual

b) Spectral Features:

- **Spectral Centroid:** It indicates the "centre of mass" of the spectrum, which is equivalent to the perceived brightness of the audio.
- **Spectral Bandwidth:** measures the spread of frequencies in a signal.
- **Spectral Contrast:** this feature captures how much the spectral peak and valley variations differ, pointing to the areas of energy change.
- **Reasoning:** These are spectral features providing complementary information concerning the timbre and texture features of speech signals.
- **Implementation:** These features were computed directly from the STFT of the audio signals.

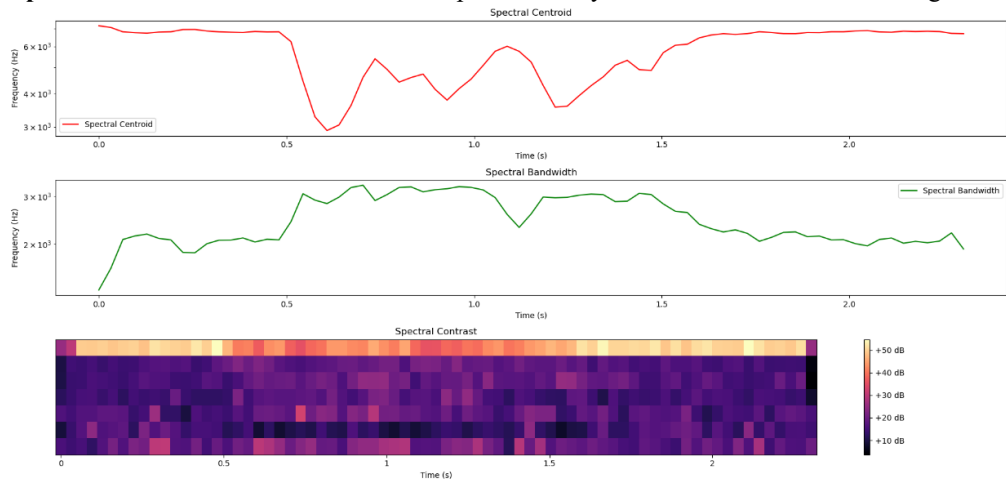


Figure 5: Spectral Features Visual

c) Principal Component Analysis (PCA):

- **Reasoning:** Speech data are typically characterized by very high-dimensional feature representations that increase the intractability of the computation besides contributing to overfitting. PCA was employed with a view of reducing the dimensionality while preserving 95% of the variance, and this actually produced a significantly more compact set of features which is easy to be used for classification.
- **Implementation:** The PCA was first trained on the training set and later used on the validation and test set to prevent bias.

Steps followed in feature extraction using the above techniques are:

a) Preprocessing via Low-Pass Filtering:

- The technique removed the high-frequency noise and artifacts so that the input data is cleaned.
- Audio files pre-processed were saved to maintain consistency and reproducibility.

b) Feature Extraction

- For each filtered audio file, using the librosa library, computed the corresponding MFCCs, Mel spectrogram, Chroma features, and spectral features.
- These features had already been successfully used in other speech-processing tasks, and they had captured all the aspects of the audio signal.

c) PCA Dimensionality Reduction:

- Applying PCA on the flattened transformed high-dimensional feature vectors. It is useful for computation, and it reduces redundant information.

Rationale Against Using Certain Functions:

a) Bicoherence (Magnitude and Phase, with Mean, Variance, Skewness, Kurtosis):

- It is the statistics of order one higher that contain phase relationship between different signals. Again, with lack of processing abilities, the inclusion of the former was impossible, however at the time these were computed.
- b) Spectral Flatness:**
- Although spectral flatness was actually computed during the feature extraction stage, it was left out of the final set of features. Preliminary experiments indicated that it would add very little to the overall classification performance.
- c) Additional Complex Features:**
- It did not use wavelet coefficients or higher-order moments due to this purpose of the research being to focus on interpretable and computationally efficient features.

The selected features have a trade-off between computational cost and discrimination. Low-pass filtering by a Butterworth filter ensures noise-free input to the system. MFCC, Mel spectrogram, and Chroma features are quite effective in presenting the speech features, and then PCA optimizes this feature set to better classify them at later stages.

4. Feature Engineering

Feature engineering is one of the critical steps in the workflows of machine learning, as it transforms raw data into a form appropriate for modelling. In this project, feature engineering includes scaling, dimensionality reduction, and ensuring the process is consistent across the training, validation, and test datasets to avoid data leakage.

- a) Application of Feature Engineering:**
- Feature engineering techniques were first applied on the training dataset.
 - Using learned transformations (like scaling parameters or principal components) the validation and testing datasets were fit so that not to leak out the data in any way
- b) PCA for Reduced Dimensionality:**
- Following feature extraction, the created features were always of high dimensional complexity and were posing problems for computing efficiency and even overfitting issues.
 - PCA has been used so that the problem would be reduced keeping 95 percent of the information of the variance.
 - The PCA model was fitted to the training data and then applied for transforming validation and test data.

4.1. Implementation Details

- a) Feature Flattening:**
- Audio features extracted like MFCCs, chroma, spectral centroid, spectral contrast, etc., were flattened into a single feature vector for every audio sample.
- b) PCA Transformation:**
- Flattened feature vectors are passed through PCA.
 - The number of components automatically selects to retain 95% of the variance. This automatically reduces the dimensionality.

4.2. Results

The feature dimensions were drastically reduced by PCA:

Dataset	Number of Samples	Original Dimensions	Reduced Dimensions
Training Set	5209	156	2
Validation Set	1116	156	2
Test Set	1117	156	2

Dimensionality reduction by PCA had the advantage of reducing computation but preserved the important information, making it ready for further classification steps.

5. Use of diverse algorithms for classification

This section describes the use of different machine learning classifiers. We trained for each algorithm on the training set, and chose optimal hyperparameters on the validation set, using RandomizedSearchCV and then run the final evaluation against the test set. The pipeline is used to ensure that data is adequately used for training, tuning, testing without over fitting. PCA was used as a dimensionality reduction method to keep 95% variance whereas StandardScaler was used to scale the features.

5.1. K-Nearest Neighbour

- a) The KNN algorithm is a nonparametric algorithm that classifies a test sample in the feature space with the majority class of its k nearest neighbours. I preprocessed my data using StandardScaler to normalize the features. RandomizedSearchCV was used to determine optimal number of neighbours k (k), as well as weight method (uniform or distance). We tested the final model on training set and validate it using test set to get the accuracy of final model.
- b) **Advantages:**
 - Easy to apply and easy to understand.
 - No training phase (lazy learner).
- c) **Drawbacks:**
 - Sensitive to the selection of k
 - Computationally intensive at the time of prediction for large data.
- d) **Hyperparameters Tuned:**
 - k : Number of neighbours.
 - Weights: Equal weight (equal) or distance based (a higher weight for closer neighbours)

5.2. Logistic Regression:

- a) A statistical classifier used for binary and multi class classification, Logistic Regression is used. This is another one that gives the probability of a class in a logistic function. It then takes thresholds and labels according to probability.
- b) **Advantages:**
 - Computational efficiency and interpretability.
 - Hard to overfit with small number of data, and is robust to small datasets.
- c) **Problems:**
 - For nonlinear data, the features must be transformed or underperforms.
- d) **Hyperparameters Tuned:**
 - C : High strength of regularization (low C).

Mostly, the tuning of the hyperparameter was focused for the good trade off between the complexity of the model and the generalization.

5.3. Naive Bayes

Allowing feature independence while using normal distribution for each class defines when to apply Gaussian Naive Bayes..

- a) **Pros:**
 - The algorithm operates with exceptional speed for both small and big datasets.
 - Works well in text classification or high-dimensional data.
- b) **Cons:**
 - Real-world datasets exhibit features that inherently depend on each other yet this model ignores this natural relation.
 - Sensitive to class imbalance.

The model required no hyperparameter tuning so it received a default configuration for its direct execution.

5.4. Support Vector Machine (SVM)

The SVM modeling algorithm employs supervised techniques to establish optimal class boundaries through maximal feature space distortion. The model used a linear kernel because of its simple nature..

a) Advantages:

- High-dimensional data can be handled.
- Model performance remains stable when C receives appropriate adjustment.

b) Challenges:

- Computationally expensive for large datasets.
- The selection of kernel functions along with regularization parameters demands great attention because it leads to critical impacts on model performance.

c) Hyperparameters Tuned:

- C: A Regularization parameter determines how maximization of margins interact with risk of misclassification.

Modeled using RandomizedSearchCV methodology which arranged the C value in a manner that delivered improved model generalization capabilities alongside preserved classification performance outcomes.

5.5. Neural Network

The Feedforward Artificial neural network MLP handles complex relationships occurring in datasets as part of its operation.

a) Advantages:

- The model demonstrates advanced functionality through its ability to identify complex non-linear decision boundaries.
- Most tasks can benefit from adjustments to the design of this system

b) Limitations:

- This model demands superior computational power in comparison to prior models since it needs additional data examples
- The response tendency toward overfitting during relationship detection occurs during improper regularization while the model is untuned.

c) Hyperparameters Tuned:

- The number of neurons within this layer decides its maximum computational threshold.
- Learning rate: During gradient descent weight updates the learning rate controls the stepping size.

Through RandomizedSearchCV users obtained the best combination of hyperparameters which optimized model complexity alongside training efficiency.

6. Performance Evaluation

Here, models are compared using performance evaluation parameters such as accuracy, precision, recall, and F1-score. It is created by the use of confusion matrices and bar charts to get visualization for checking the behaviour of a model as well as the efficiency.

6.1. Metrics for Comparison

The metrics that were used in the model evaluation analysis are mentioned below:

- **Accuracy:** Percentage of all the instances of which cases were classified correctly.
- **Precision:** All correct positives divided by the sum of true positives plus false positives, that is true positive predictive value.
- **Recall:** True positives divided by the sum of true positives and false negatives.
- **F1 Score:** Average measure from the harmonic mean of accuracy and sensitivity.

6.2. Visualizations

a) Accuracy Comparison

The bar chart below compares the accuracy of different classifiers.

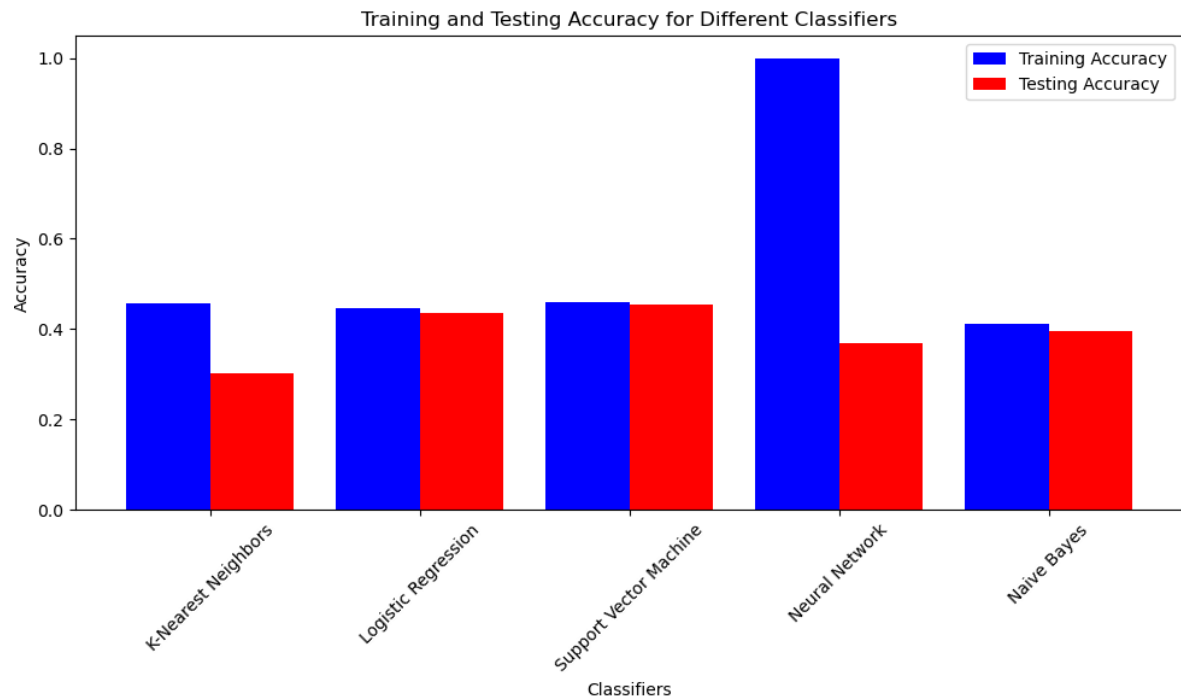


Figure 6: Accuracy Comparison Chart

b) Precision, Recall, and F1-Score Comparison

K-Nearest Neighbors:

Accuracy: 0.3035

Precision: 0.3343

Recall: 0.3035

F1 Score: 0.2987

Logistic Regression:

Accuracy: 0.4369

Precision: 0.4311

Recall: 0.4369

F1 Score: 0.4308

Support Vector Machine:

Accuracy: 0.4539

Precision: 0.4465

Recall: 0.4539

F1 Score: 0.4474

Neural Network:

Accuracy: 0.3688

Precision: 0.3774

Recall: 0.3688

F1 Score: 0.3723

Naive Bayes:

Accuracy: 0.3948

Precision: 0.3980

Recall: 0.3948

F1 Score: 0.3842

c) Confusion Matrices

Confusion matrices for each model illustrate true positive, true negative, false positive, and false negative counts, helping to identify patterns in model performance.

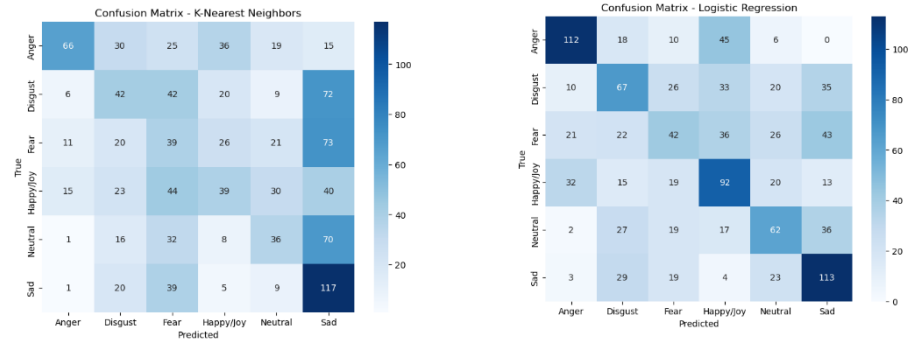


Figure 7.1: Confusion Matrix

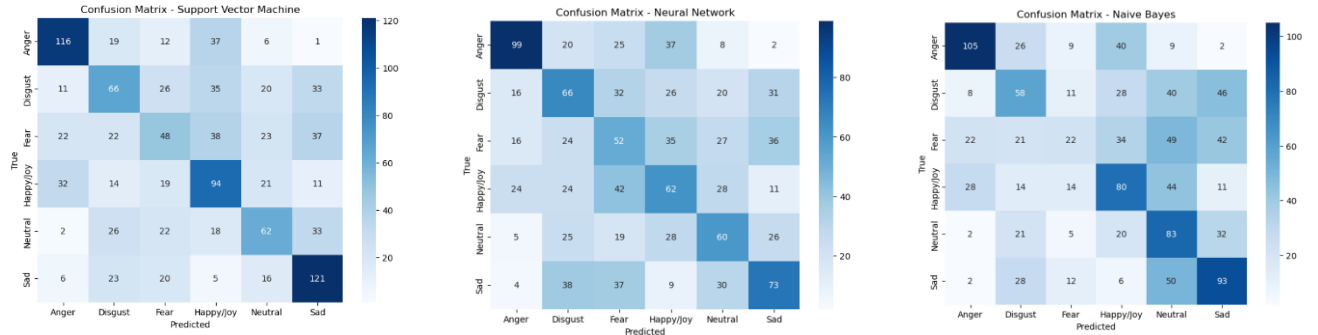


Figure 7.2: Confusion Matrix

6.3. Analysis and Findings

Global Performance

- **Accuracy Maxima:** The accuracy of the presented algorithm was as high as 45.39% by Support Vector Machine, so this verifies its potential to function well in a very high-dimensional feature space.
- **Precision vs. Recall Trade-off:** Both Logistic Regression and SVM exhibited precision as well as recall, therefore both are satisfactory in terms of class sizes almost being equal to each other.

Model Conclusion

- **KNN:** Pretty bad too at 30.35% because of the curse of dimensionality, though sensitivity to feature scaling wasn't predominant even after k tuning.
- **Logistic Regression:** Moderate performance in terms of accuracy at 43.69%. However, struggled with the complex non-linear nature of decision boundaries.
- **Support Vector Machine (SVM):** Has performed better than all other models at 45.39%, as it uses the kernel trick to handle the non-linear boundaries.
- **Neural Network:** Overfitting was noticed with training accuracy (99.87%) much higher than test accuracy (36.88%). This requires regularization or dropout tuning.
- **Naive Bayes:** Fair results were achieved (accuracy: 39.48%), but the assumption of feature independence was a limitation.

Strengths and Weaknesses

- **Strengths:**
 - SVM and Logistic Regression gave consistent precision and recall in all experiments.
 - Neural Network has promise but requires further tuning to generalize.
- **Weaknesses:**
 - KNN's sensitivity to high-dimensional data limited its accuracy.
 - Naive Bayes suffered the problem of feature dependencies in this dataset.

- Neural Network overfitting hindered its usability.

Conclusion:

1. Feature engineering will improve the performance of a model, especially KNN and Naive Bayes.
2. Regularization techniques are the movement to avoid Neural Networks from getting more overfit.
3. The powerful models such as SVM were doing really well, but still, the complexity of an algorithm was a problem while dealing with huge data.

6.4. Recommendation to Future Work

- **Hyperparameter Tuning:** Optimize and further tune hyperparameters using novel optimization algorithms, Bayesian Optimization or Genetic Algorithms.
- **Data Augmentation:** Add diverse audio recordings to the data set for generalization of the model.
- **Feature Selection:** Use advanced techniques such as RFE for selecting the best features.
- **Deep Learning Models:** Apply deep learning architectures, CNNs, or RNNs, to the audio-based sentiment analysis.
- **Cross-validation:** Use more aggressive cross-validation strategies, such as stratified K-fold, to improve the assessment.

7. Conclusion

In summary, this paper focuses on the approach towards sentiment analysis of audio data by using machine learning techniques. Methodology begins with the preprocessing audio recordings from the CREMA-D dataset by including low-pass filtering among other techniques in order to remove noise. Extracted features then go into various models like KNN, Naive Bayes, SVM, Logistic Regression, and Neural Networks, for sentiment classification. Among the tested models, SVM classifier achieved a highest accuracy of 45.39% in comparison to other models. Even though the achieved accuracy is moderate, it also depicts the difficulties inherent in the process of sentiment analysis on audio data due to variability in speech patterns, accents, and expression of emotion.

Even further, the complexity of the task was increased for the classification job since the dataset was ethnically and culturally diversified. Also, Neural Network was overfitting as it had few training samples only; this implies that bigger datasets and regularization techniques are required to make deep learning successful. Thus, the task of the problem was well tackled by the SVM classifier, which implies to have a practical application in sentiment analysis. The future work may focus on the integration of deep learning models, for example, CNNs or RNNs for better feature representation and ensemble learning.

8. Relevant References

8.1. Research Papers

- [1] Fayek, H. M., Lech, M., & Cavedon, L. (2017). *Evaluating deep learning architectures for Speech Emotion Recognition*. Neural Networks, 92, 60-68. ([DOI Link](#))
- [2] Lotfian, R., & Busso, C. (2019). *Curriculum learning for speech emotion recognition from crowdsourced labels*. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 27(4), 815-826. ([DOI Link](#))
- [3] Li, Z., Zhao, Y., & Li, C. (2022). *Speech Emotion Recognition Using Multi-feature Fusion and Ensemble Learning*. Applied Acoustics, 197, 108964. ([DOI Link](#))

8.2. Websites and Articles

- *CREMA-D Dataset Overview*: <https://github.com/CheyneyComputerScience/CREMA-D>
- *Scikit-learn Documentation on SVM*: <https://scikit-learn.org/stable/modules/svm.html>

- *Feature Extraction Techniques for Audio Data:* <https://towardsdatascience.com/audio-feature-extraction-explained-41f50f5551bf>

8.3. Articles on Sentiment Analysis

- *Deep Learning Approaches for Sentiment Analysis in Speech Data:* <https://www.analyticsvidhya.com/blog/2021/06/sentiment-analysis-in-audio-using-deep-learning/>
- *Challenges in Sentiment Analysis from Speech:* <https://www.kdnuggets.com/2021/11/speech-sentiment-analysis-challenges.html>