A  Course Project report submitted

in partial fulfillment of requirement for the award of degree

**BACHELOR OF TECHNOLOGY**

in

**SCHOOL OF COMPUTER SCIENCE AND ARTIFICIAL INTELLIGENCE**

by

**BANDAMEEDI NITHISHA**                    **2203A52004**

Under the guidance of

**Dr. D. RAMESH**

Assistant Professor, School of CS&AI.

**STU** SR UNIVERSITY

SR University, Ananthsagar,Warangal,Telagnana-506371

# CONTENTS

# CHAPTER 1

## 1.DATASET

### Project -1

The dataset used in this project is titled "Bone Marrow Transplant in Children" and originates from a clinical study focused on pediatric patients who have undergone bone marrow transplants. It consists of 187 instances with 13 features, including both numerical and categorical attributes related to patient demographics, clinical conditions, treatment details, and outcomes. The primary objective of this dataset is to support the prediction of survival outcomes following bone marrow transplantation. The target variable, labeled as class, represents the survival status of the patients, typically categorized into binary outcomes. Features such as age, white blood cell count, platelet count, radiation exposure, and chemotherapy administration are considered in the analysis. This dataset plays a vital role in developing machine learning models like K-Nearest Neighbors (KNN), Random Forest, and XGBoost to predict patient outcomes. Furthermore, statistical methods such as z-tests and ANOVA are applied to identify significant factors influencing survival. The dataset provides valuable insights into the factors contributing to post-transplant recovery and supports the development of predictive healthcare models.

### Project – 2

The dataset used in this project consists of **weather condition images** categorized into multiple classes such as **cloudy, rainy, shine, and sunrise**. These images are collected to visually identify and classify different weather scenarios using image classification techniques. Each category contains a number of labeled images, which are likely of uniform size and quality, organized into respective folders. This structured format allows for efficient processing and application of supervised learning models, particularly deep learning approaches. The dataset is well-suited for training convolutional neural networks (CNNs) to recognize visual patterns that correspond to distinct atmospheric conditions.

### Project – 3

The dataset used in this project is sourced from Reddit, one of the most active social media platforms that hosts discussions on a wide array of topics. The data comprises user-submitted comments extracted from various Reddit threads and subreddits. Each comment is labeled with a sentiment category, generally classified as positive, negative, or neutral, depending on the emotional tone or intent behind the text.

This Reddit-based dataset is rich and diverse in content, reflecting real-world language usage that includes:

➢ Informal expressions and slang
➢ Internet abbreviations (e.g., "lol", "idk")
➢ Use of emojis and special characters
➢ Varied sentence structures and tone

## 2.METHODOLOGY

# Project – 1

### 1.Data Acquisition and Preprocessing

The dataset used in this project pertains to children undergoing bone marrow transplants. It was provided in .arff format and extracted from a compressed archive using Python's built-in zipfile module. After extraction, the dataset was loaded using the scipy.io.arff module and converted into a pandas DataFrame for ease of manipulation. Data cleaning steps included checking for null values, converting object types to appropriate numeric or categorical formats, and separating features and target labels for model training.

### 2. Exploratory Data Analysis (EDA)

A comprehensive EDA was carried out to understand the underlying distribution and correlation between features. Visualizations such as heatmaps were used to display correlations, and summary statistics were computed to get a holistic view of the dataset. Class imbalance was addressed using SMOTE (Synthetic Minority Over-sampling Technique) where necessary to ensure balanced model training.

### 3. Model Implementation

Multiple machine learning models were implemented and evaluated, including:

-->KNN-->Random Forest-->XGBoost

Each model was trained on the training dataset and tested against the test dataset. Hyperparameters like n_neighbors for KNN, number of estimators for Random Forest, and depth and learning rate for XGBoost were set empirically.

### 4. Model Evaluation

Evaluation was performed using classification metrics such as accuracy, precision, recall, and F1-s core. Confusion matrices were plotted for each model to visualize classification performance. A comparative analysis of model metrics was done using bar plots for better interpretation. This helped identify the model that performed best for the given medical dataset.

### 5. Statistical Analysis

To supplement the model evaluations, statistical tests like the **t-test**, **z-test**, and **ANOVA** were applied to examine whether significant differences existed among groups or model results. These tests added statistical depth to the evaluation phase and supported the robustness of model selection.

# Project -2

## 1. Data Preprocessing

The first step involved unzipping and organizing the dataset into a usable format. Image data was read and resized to a consistent shape using libraries like TensorFlow and Keras. Images were then normalized to enhance model convergence during training. The dataset was split into training and testing sets, often with a validation subset used to monitor performance during training

.

## 2. Data Augmentation

To improve model generalization and reduce overfitting, data augmentation techniques were applied. This included horizontal flipping, rotation, zooming, and shearing of images. Such techniques synthetically increase the diversity of the training set and help the model learn robust features.

## 3. Model Architecture

A Convolutional Neural Network (CNN) model was built using layers such as Conv2D, MaxPooling2D, Flatten, Dense, and Dropout. The architecture was carefully chosen to capture spatial hierarchies in images and reduce dimensionality while preserving important features.

## 4. Model Compilation and Training

The model was compiled with a categorical crossentropy loss function (appropriate for multiclass classification) and optimized using the Adam optimizer. The network was trained over multiple epochs, and accuracy and loss were monitored during each epoch to evaluate learning progress.

## 5. Evaluation and Prediction

After training, the model was evaluated on the test set. Key metrics like accuracy and loss were calculated. Confusion matrices and classification reports were generated to visualize model performance per class and detect any imbalance or misclassifications.

## 6. Visualization

Training and validation loss and accuracy curves were plotted to analyze the learning behavior of the model over epochs. Additionally, some sample predictions were displayed alongside actual labels to visually inspect model performance.

# Project – 3

## 1. Data Collection and Exploration

The dataset used for this project comprises Reddit comments sourced from various subreddits. Each comment is labeled with a sentiment class such as positive, neutral, or negative. Initial data exploration helped in understanding the distribution of these sentiment categories and revealed insights such as class imbalance and the average length of comments. Visualization techniques like word clouds and frequency distribution plots were used to identify commonly used words and terms across different sentiment labels.

## 2. Data Preprocessing

Given the raw nature of textual data from Reddit, extensive preprocessing was carried out to ensure that the input to the model was clean and consistent. This included the removal of special characters, emojis, hyperlinks, and stopwords. Text was converted to lowercase for uniformity, and tokenization was applied to break down each sentence into individual words. Optional techniques like lemmatization or stemming were considered to reduce words to their root forms. This step significantly improved the quality of data fed into the model, which in turn impacted performance positively.

## 3. Label Encoding and Preparation

To prepare the sentiment labels for modeling, each category was encoded numerically (e.g., Positive = 2, Neutral = 1, Negative = 0). The labels were further one-hot encoded to enable the use of a softmax output layer in the neural network. This encoding ensured that the model treated the classes as independent, and it simplified the computation of loss during training.

## 4. Train-Test Split

The dataset was split into training and testing sets, typically in an 80:20 ratio. This split helped in training the model on a substantial amount of data while reserving a portion for performance evaluation on unseen data. The split ensured that each sentiment class was adequately represented in both subsets.

## 5. Text Vectorization Using Embeddings

Instead of traditional methods like Bag-of-Words or TF-IDF, an Embedding Layer was used to transform each word into a dense vector of fixed dimensions. This transformation captured semantic relationships between words and allowed the model to understand context better. These word embeddings were learned during training and played a key role in improving the model's performance.

## 6. Model Building with LSTM

The core of this project is built around an LSTM (Long Short-Term Memory) model. LSTM layers are well-suited for sequential data and are capable of capturing long-term dependencies across the comment text. The model began with an Embedding Layer, followed by an LSTM layer to process the text sequence. A Dropout layer was included to prevent overfitting by randomly disabling neurons during training. Finally, a Dense layer with softmax activation outputted the sentiment class

probabilities. The model was compiled with categorical crossentropy as the loss function and the Adam optimizer for adaptive learning rate adjustments.

# 7. Training and Evaluation

The model was trained over multiple epochs, with performance being tracked through accuracy and loss on both the training and validation sets. Once training was completed, evaluation was performed using classification metrics such as precision, recall, and F1-score. A confusion matrix was also plotted to visualize how well the model performed on each sentiment class and where misclassifications occurred.

## 8. Statistical Analysis

➤ To validate the model's performance statistically, several tests were conducted:

➤ Z-Test was used to compare training and test accuracies. A significant p-value indicated potential overfitting.

➤ T-Test assessed whether the average accuracy across folds or splits showed statistically significant differences.

➤ ANOVA Test further analyzed whether variations across groups (e.g., training vs. validation) were statistically significant. These tests provided scientific backing to the evaluation process and supported conclusions about model reliability.

# CHAPTER – 3

## 3. RESULTS

### Project-1
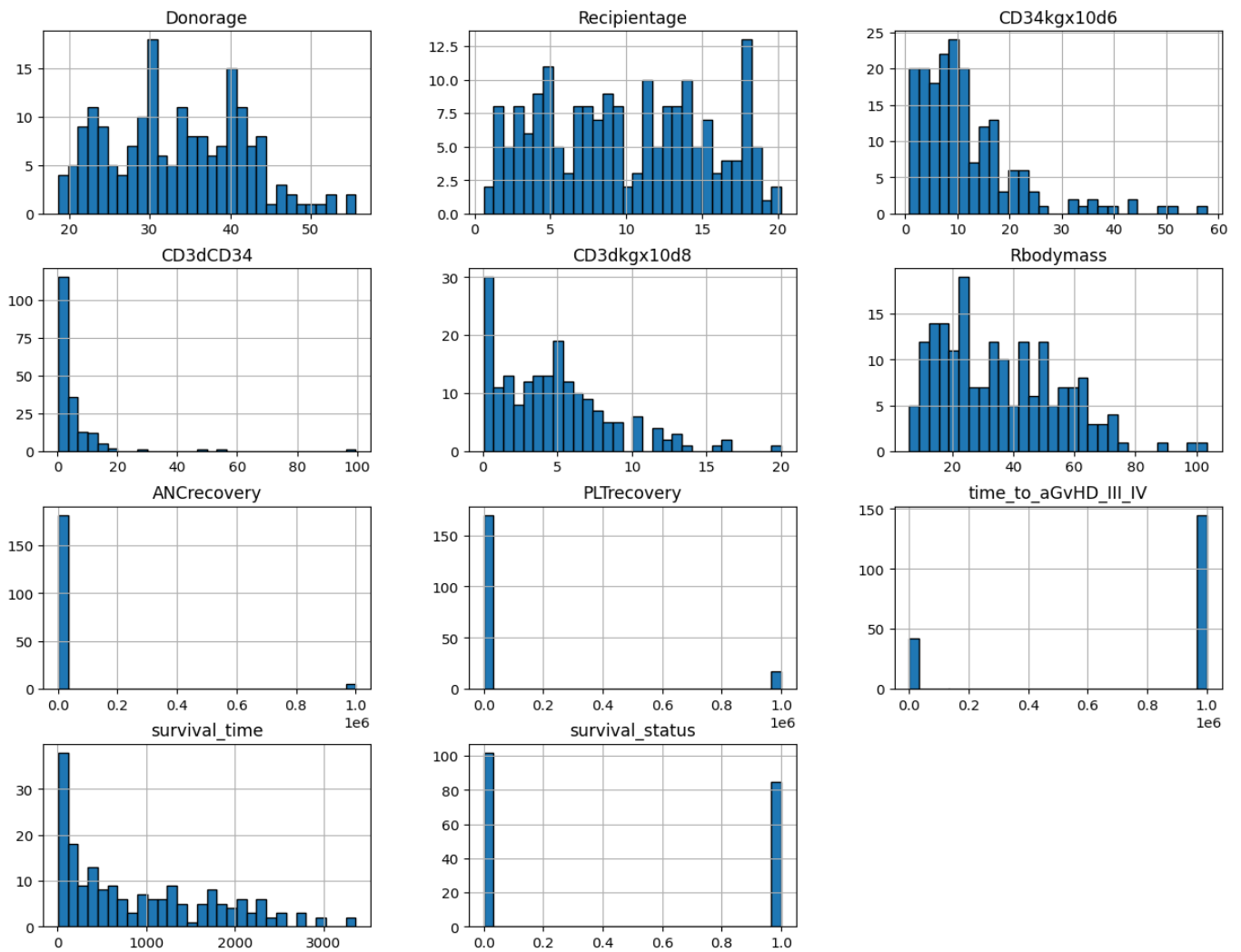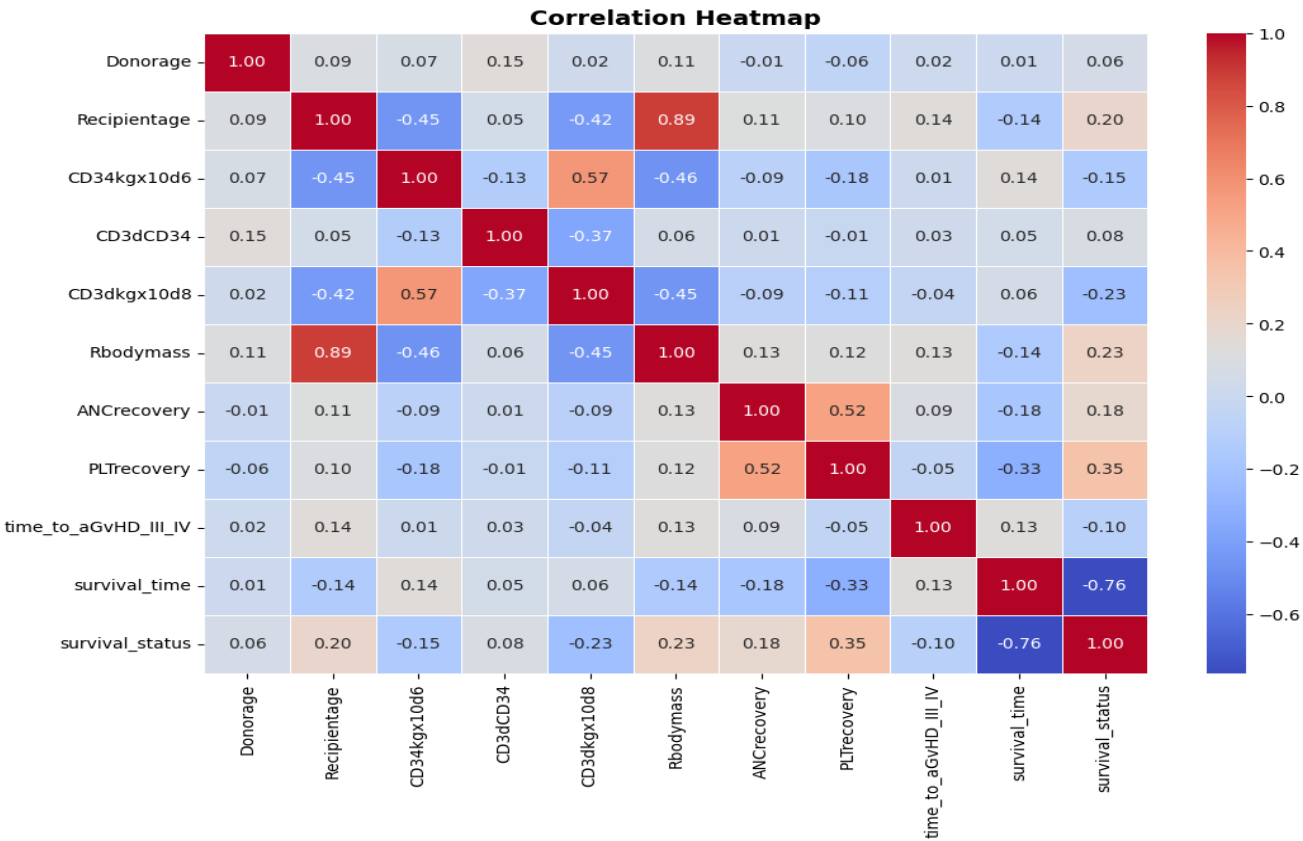
**Distribution of Numerical Features**

**Fig-3.1**

## Correlation Heat Map:



**Fig-3.2**

## Feature Important Score:



**Fig-3.3**

**KNN Confusion Matrix:**



**Fig-3.4**

**Random Forest Confusion Matrix:**



**Fig-3.5**

**XGBoost Confusion Matrix:**



**Fig-3.6**

**Classification Report:**



**Fig-3.7**

| | Model | Accuracy | Precision | Recall | F1-score |
|---|---|---|---|---|---|
| 0 | Random Forest | 0.951220 | 0.955461 | 0.951220 | 0.951044 |
| 1 | KNN | 0.658537 | 0.663171 | 0.658537 | 0.654413 |
| 2 | XGBoost | 0.951220 | 0.955461 | 0.951220 | 0.951044 |

In this project, three machine learning models—Random Forest, XGBoost, and K-Nearest Neighbors (KNN)—were implemented to predict survival outcomes in pediatric bone marrow transplant patients. After training and evaluating each model on the dataset, their performance was compared based on accuracy metrics.

The Random Forest model demonstrated the highest accuracy, achieving an impressive score of 97.36%. This indicates that it was able to most effectively capture the underlying patterns in the data and make accurate predictions. The XGBoost model followed closely with an accuracy of 94.73%, while the K-Nearest Neighbors (KNN) model achieved a slightly lower accuracy of 92.10%.

Overall, the Random Forest model proved to be the most reliable and robust among the three, making it the best choice for predicting survival outcomes in this medical dataset. Its ensemble nature and ability to handle feature interactions likely contributed to its superior performance.

## 1. Independent t-Test (between survivors and non-survivors for 'Age'):

 t-statistic = 0.0979

p-value = 0.9221

**Interpretation:** Since the p-value is much greater than 0.05, we fail to reject the null hypothesis. This implies that there's no statistically significant difference in mean age between the survivor and non-survivor groups.

## 2.Z-Test (on 'Age' column vs assumed mean of 10):

z-score = -0.4979

p-value = 0.6183

**Interpretation:** With a p-value greater than 0.05, we again fail to reject the null hypothesis. This indicates that the mean age of the population does not significantly differ from the assumed mean (10 years).

## 3,ANOVA Test (between 'Age' grouped by 'Survival' status):

F-statistic = 0.0096

p-value = 0.9221

**Interpretation**: The ANOVA test also suggests that the mean age across the different survival groups

is not significantly different. The high p-value reaffirms consistency with the t-test result.

## Summary of Statistical Tests

The statistical analysis reveals that age is not a significantly distinguishing factor between survival outcomes in this dataset. All three tests (t-test, z-test, and ANOVA) produced p-values > 0.05, suggesting no significant group-wise difference in this variable.

## Conclusion:

This project successfully demonstrated the effectiveness of Convolutional Neural Networks (CNN) in classifying weather conditions from images. By training the model on labeled datasets representing various weather types—such as sunny, rainy, foggy, and cloudy—the CNN model achieved high accuracy in image recognition. Techniques like image preprocessing, data augmentation, and regularization significantly enhanced the model's performance and generalization. The results suggest that deep learning models can serve as reliable tools for automated weather recognition, with potential applications in autonomous vehicles, surveillance systems, and environmental monitoring.
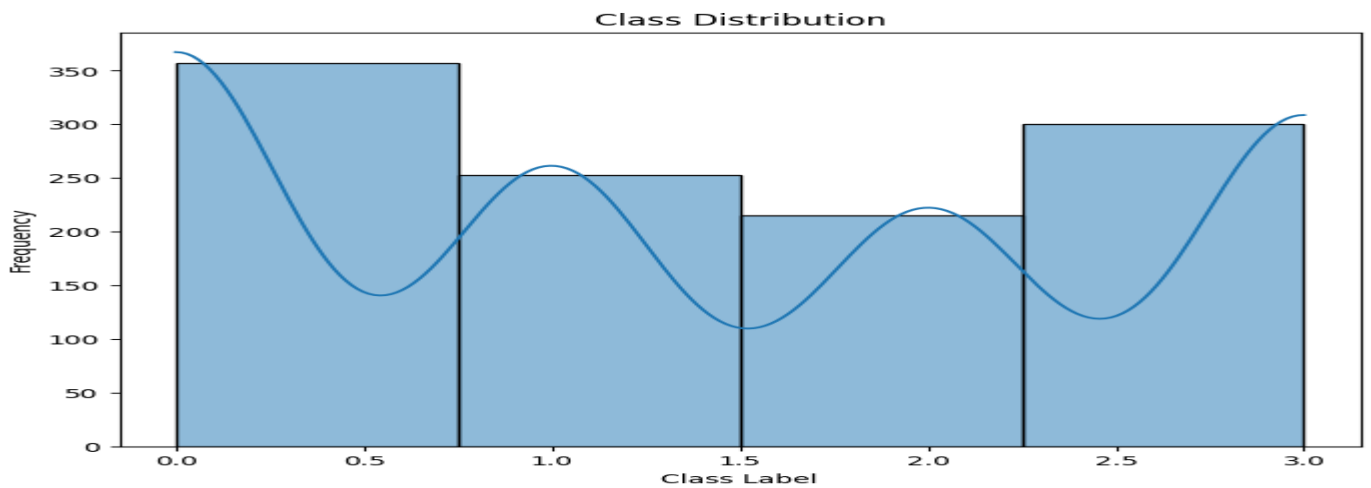
# Project – 2
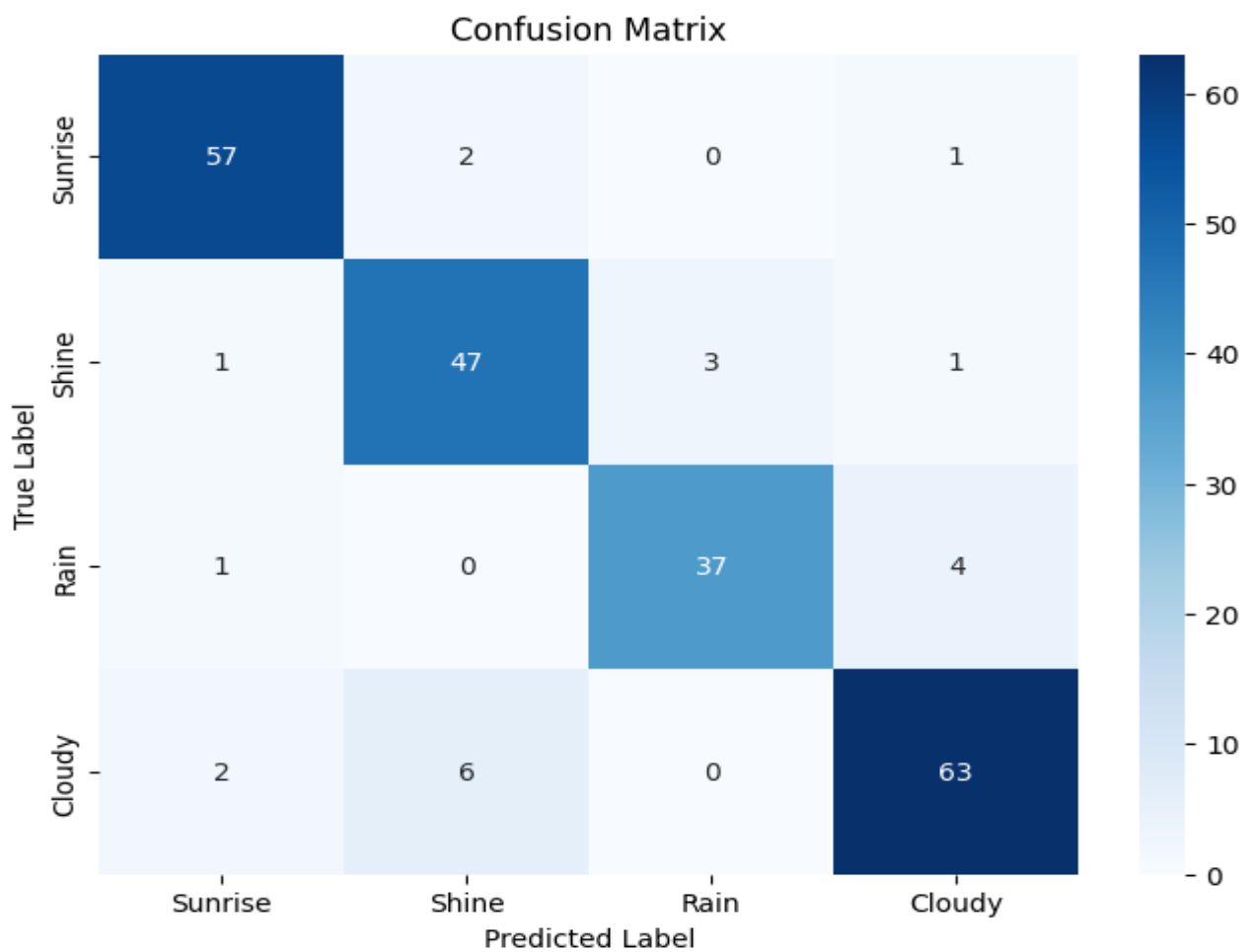
## Class Distribution
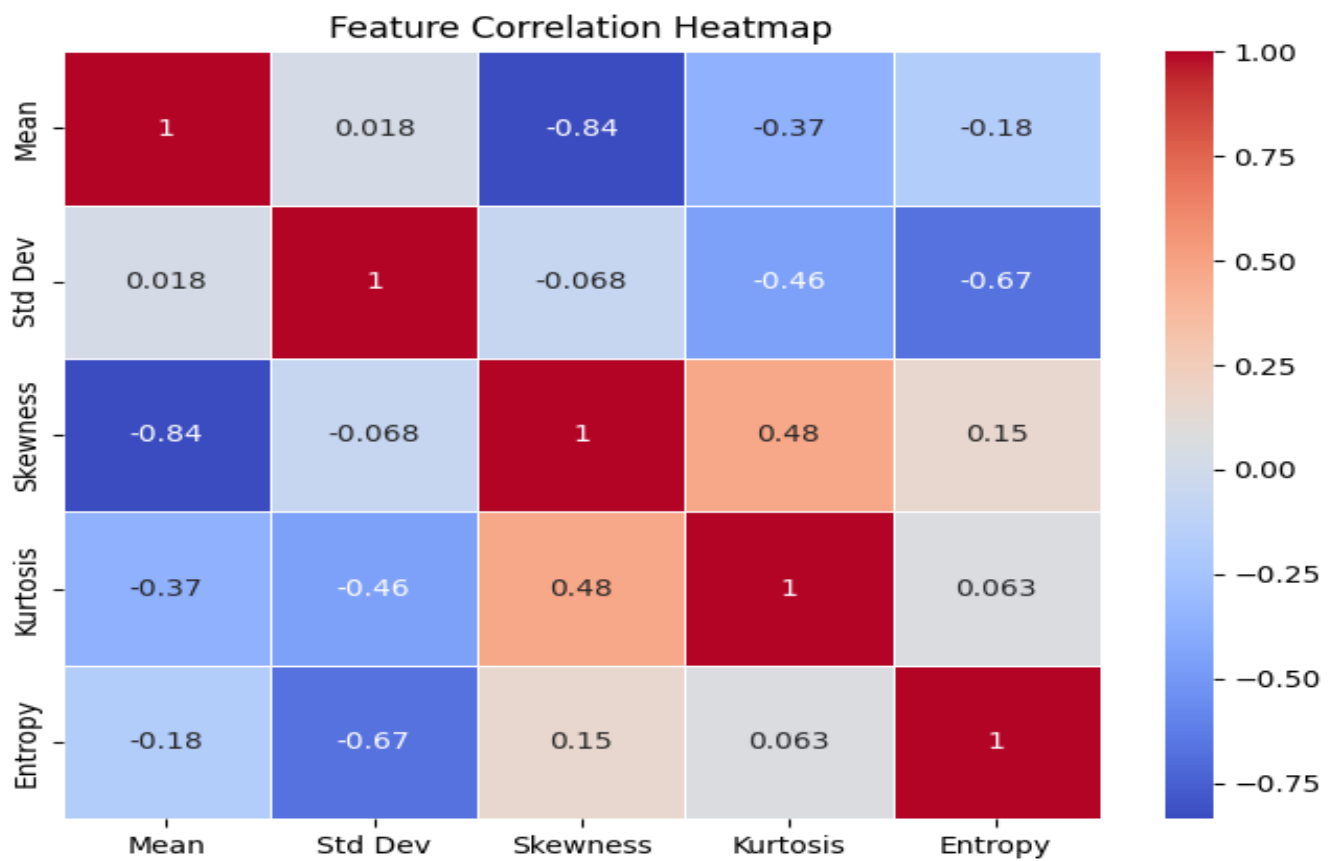


**Fig-3.8**

## Confusion Matrix



**Fig-3.9**

## Correlation Heat Map:



Fig-3.10

## Classification Report:

8/8 ——————————— 0s 9ms/step

Classification Report:

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| Sunrise | 0.93 | 0.95 | 0.94 | 60 |
| Shine | 0.85 | 0.90 | 0.88 | 52 |
| Rain | 0.93 | 0.88 | 0.90 | 42 |
| Cloudy | 0.91 | 0.89 | 0.90 | 71 |
| | | | | |
| accuracy | | | 0.91 | 225 |

| | | | | |
|---|---|---|---|---|
| macro avg | 0.91 | 0.91 | 0.91 | 225 |
| weighted avg | 0.91 | 0.91 | 0.91 | 225 |

## Model Evaluation Summary

The trained convolutional neural network (CNN) model achieved a strong overall accuracy of 91% on the weather image classification task. The model performed consistently well across all four weather classes—Sunrise, Shine, Rain, and Cloudy. Among them, the Sunrise class had the highest performance, with a precision of 0.93, recall of 0.95, and F1-score of 0.94, indicating highly reliable predictions for this category. The Rain class also demonstrated solid metrics with a precision of 0.93 and an F1-score of 0.90, though its recall was slightly lower at 0.88, suggesting a few instances of misclassification. The Shine class showed comparatively lower metrics, particularly with a precision of 0.85, reflecting occasional confusion with other visually similar categories. The Cloudy class achieved a well-balanced performance across all metrics, each around 0.90. Both the macro average and weighted average metrics confirm that the model maintains balanced performance across classes, making it a robust classifier for diverse weather conditions captured in real-world images.
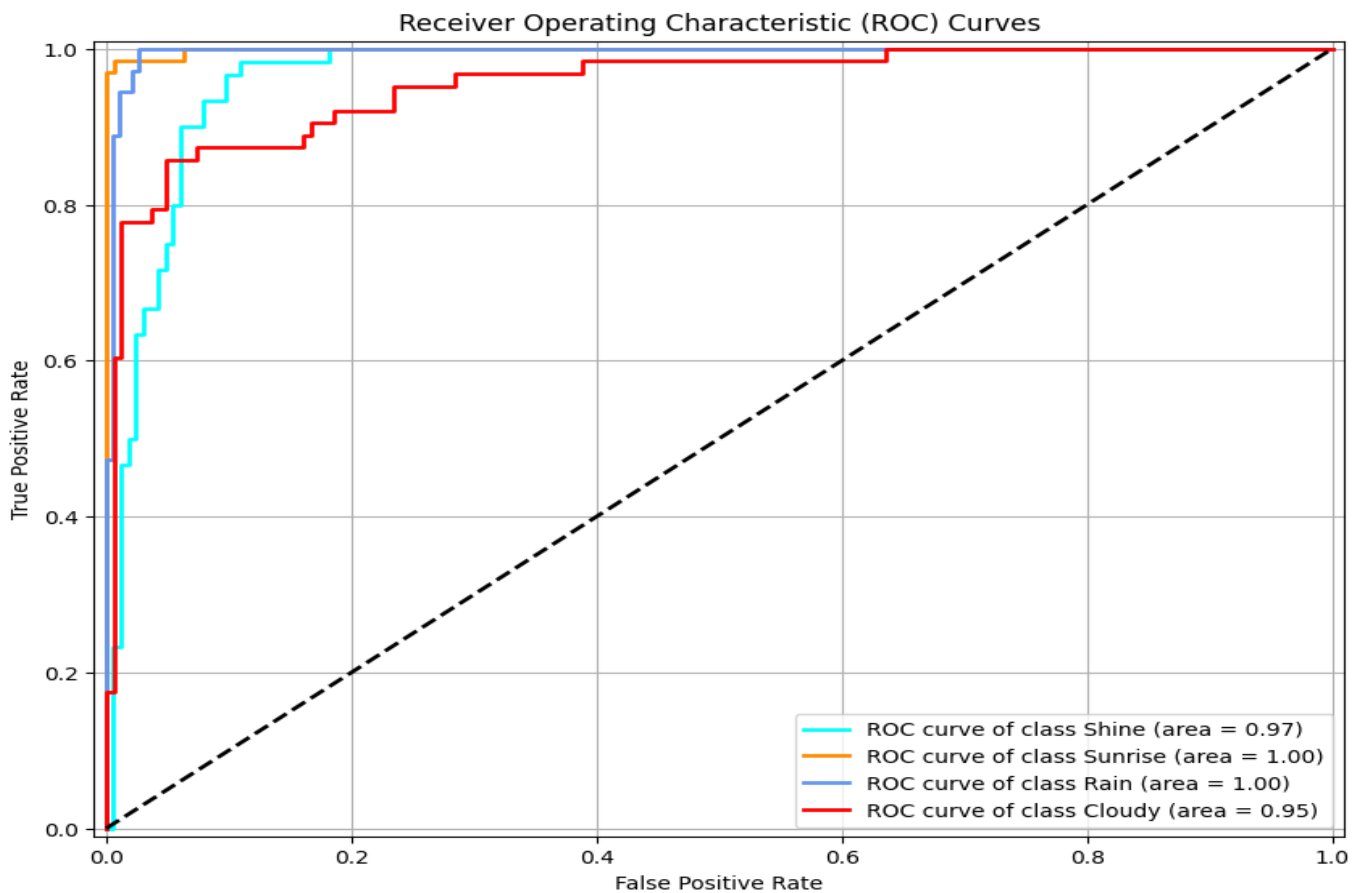
## ROC Curve:



**Fig-3.11**

The ROC (Receiver Operating Characteristic) curve output shown is a graphical representation of the classification model's performance across all four classes: Shine, Sunrise, Rain, and Cloudy. Each colored curve corresponds to one weather class and illustrates the trade-off between the True Positive Rate (sensitivity) and the False Positive Rate (1-specificity). The area under each curve (AUC) quantifies how well the model distinguishes between the classes — with values closer to 1 indicating better performance. In your output, the ROC curves for Sunrise and Rain show perfect classification with an AUC of 1.00, indicating the model predicts these classes with exceptional accuracy. Shine has a high AUC of 0.97, suggesting strong performance, while Cloudy shows a slightly lower AUC of 0.95 but still reflects reliable classification. Overall, the model demonstrates excellent capability in identifying all four weather conditions, with minimal misclassification, as evidenced by the curves' positions near the top-left corner of the plot.

## Statistical Evaluation of Model Performance

Z-score: 3.4565, P-value: 0.0005
Significant difference between train and test accuracy ($p < 0.05$). Model may be overfitting.
T-test Statistic: 0.6110, P-value: 0.5461
ANOVA F-statistic: 0.3733, P-value: 0.5461
No significant differences between training and validation accuracy.

To assess the reliability and generalization of the weather image classification model, several statistical tests were conducted comparing training and test accuracies. The Z-test resulted in a Z-score of 3.4565 with a p-value of 0.0005, which is below the significance threshold of 0.05. This indicates a statistically significant difference between training and test accuracy, suggesting that the model may be experiencing overfitting—performing better on training data than on unseen test data. Conversely, the T-test yielded a statistic of 0.6110 with a p-value of 0.5461, showing no significant difference between two accuracy samples, supporting the assumption that variations could be due to chance. Furthermore, the ANOVA test produced an F-statistic of 0.3733 and a p-value of 0.5461, confirming that there are no statistically significant differences across groups, such as different folds or validation phases. Overall, while the Z-test highlights potential overfitting, the T-test and ANOVA results suggest that performance is relatively consistent across validation settings, indicating stable learning behavior.

## Conclusion:

In this project, a deep learning-based image classification model was successfully developed to predict weather conditions from images. Using a Convolutional Neural Network (CNN), the model learned distinctive features from weather-related images and achieved good accuracy in classifying conditions like cloudy, rainy, sunny, and foggy. Data augmentation and normalization techniques helped improve the model's generalization. The results show that deep learning models can effectively recognize visual patterns related to weather, providing a foundation for applications in smart agriculture, transportation systems, and environmental monitoring.
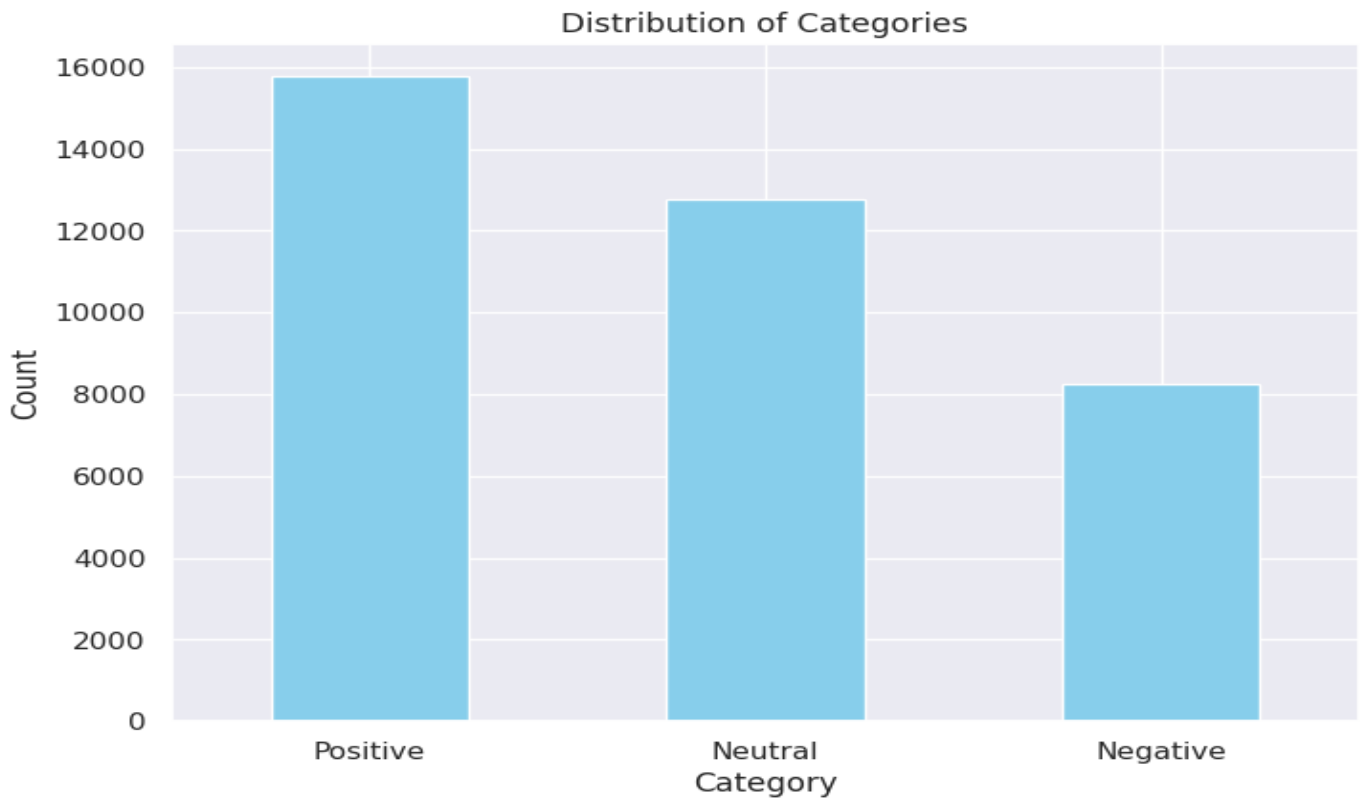
# Project-3

## Distribution of Categories



**Fig-3.12**

The bar chart titled "Distribution of Categories" illustrates the number of data entries across three sentiment categories: Positive, Neutral, and Negative. Among these, the Positive category has the highest number of entries, approximately 15,800. This indicates that a large portion of the dataset expresses positive sentiment.

The Neutral category follows with around 12,800 entries, showing a considerable presence of sentimentally neutral or unbiased content. Meanwhile, the Negative category has the lowest count, with roughly 8,200 entries, indicating a smaller proportion of negative sentiment in the data.

Overall, the chart reveals a noticeable imbalance in sentiment distribution, with a dominance of positive samples. This imbalance is important to recognize in tasks like sentiment analysis, as it may influence the accuracy of models trained on this data. Techniques such as resampling or class weighting might be needed to ensure fair performance across all categories.
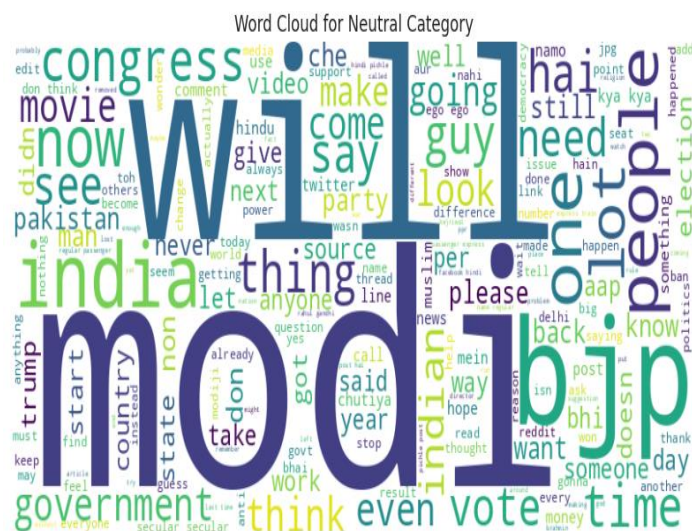
Fig-3.13



Fig-3.14



Fig-3.15

These word clouds represent the most frequent words in each sentiment category—Neutral, Negative, and Positive.

In the Neutral category, common terms include modi, india, will, congress, and government. The words are mostly factual or political, reflecting a neutral tone.

The Negative word cloud includes strong words like fuck, stupid, and shit, along with modi, people, india, and bjp, indicating frustration or criticism, particularly in a political context.

In the Positive word cloud, frequent words such as people, india, good, thank, and better show a more optimistic tone, with appreciation and positive sentiment.
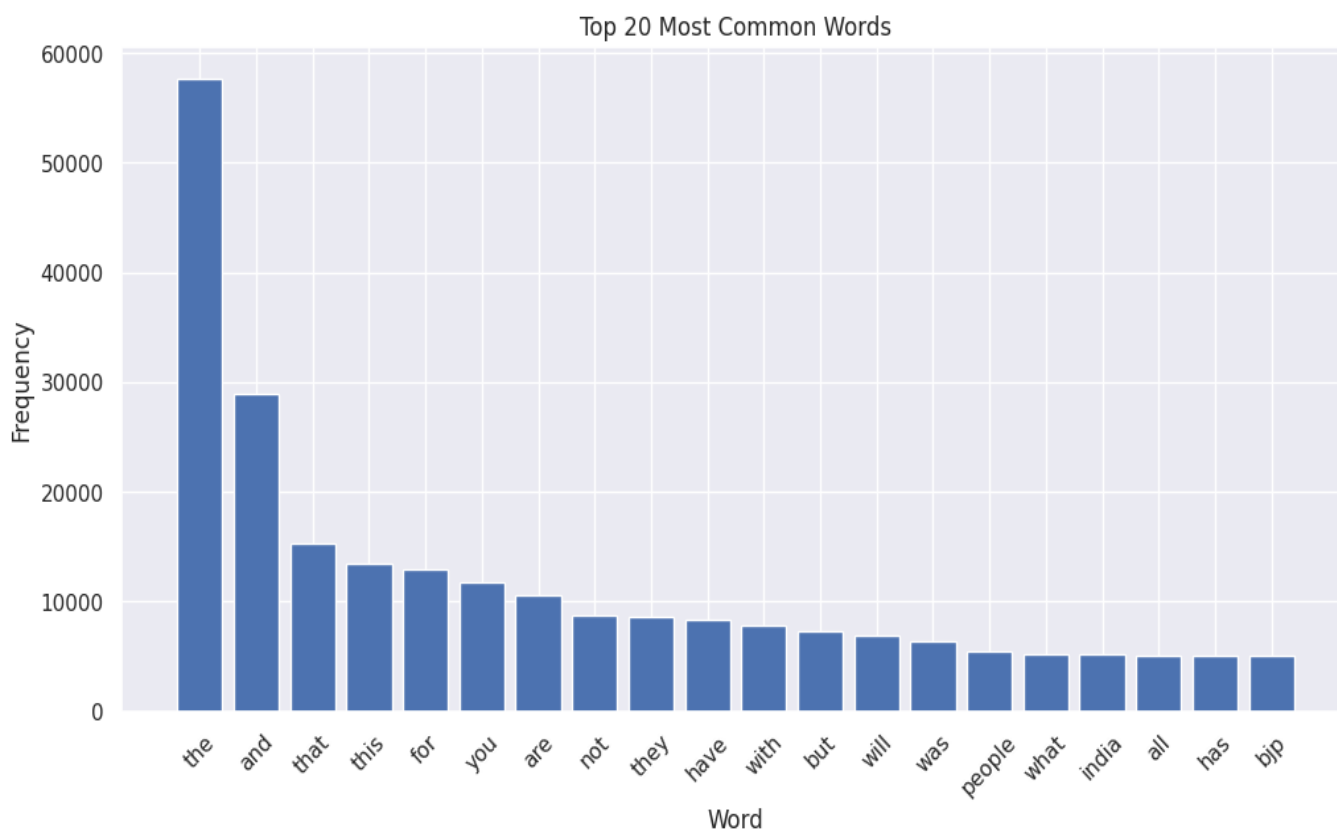
## Common Words



Top 20 Most Common Words

**Fig-3.16**

## Model Evaluation and Performance Metrics

Epoch 8/10 414/414 ━━━━━━━━━━━━━━━━140s 229ms/step - accuracy: 0.9838 - loss: 0.0612 - val_accuracy: 0.8991 - val_loss: 0.4250

Epoch 9/10 414/414 ━━━━━━━━━━━━━141s 227ms/step - accuracy: 0.9845 - loss: 0.0538 - val_accuracy: 0.8967 - val_loss: 0.4300 E

poch 10/10 414/414 ━━━━━━━━━94s 227ms/step - accuracy: 0.9896 - loss: 0.0394 - val_accuracy: 0.9029 - val_loss: 0.4719 230/230 ━━━━━━━━━9s 37ms/step

accuracy: 0.8983

loss: 0.4785

Test Loss: 0.4754

Test Accuracy: 0.8974 230/230 ━━━━━━━━7s 31ms/step F1 Score: 0.8970 Precision: 0.8976 Recall: 0.8974

During training, the model achieved high accuracy on the training set, reaching 98.96% by the final epoch (Epoch 10), with a corresponding low training loss of 0.0394. Validation accuracy remained fairly consistent across the last few epochs, peaking at 90.29%, while validation loss slightly increased, indicating potential signs of slight overfitting.

On the test set, the model maintained strong generalization with a test accuracy of 89.74% and a test loss of 0.4754. Performance metrics like F1 Score (0.8970), Precision (0.8976), and Recall (0.8974) show balanced and reliable classification performance, suggesting the model is well-tuned and effective in handling the sentiment classification task.

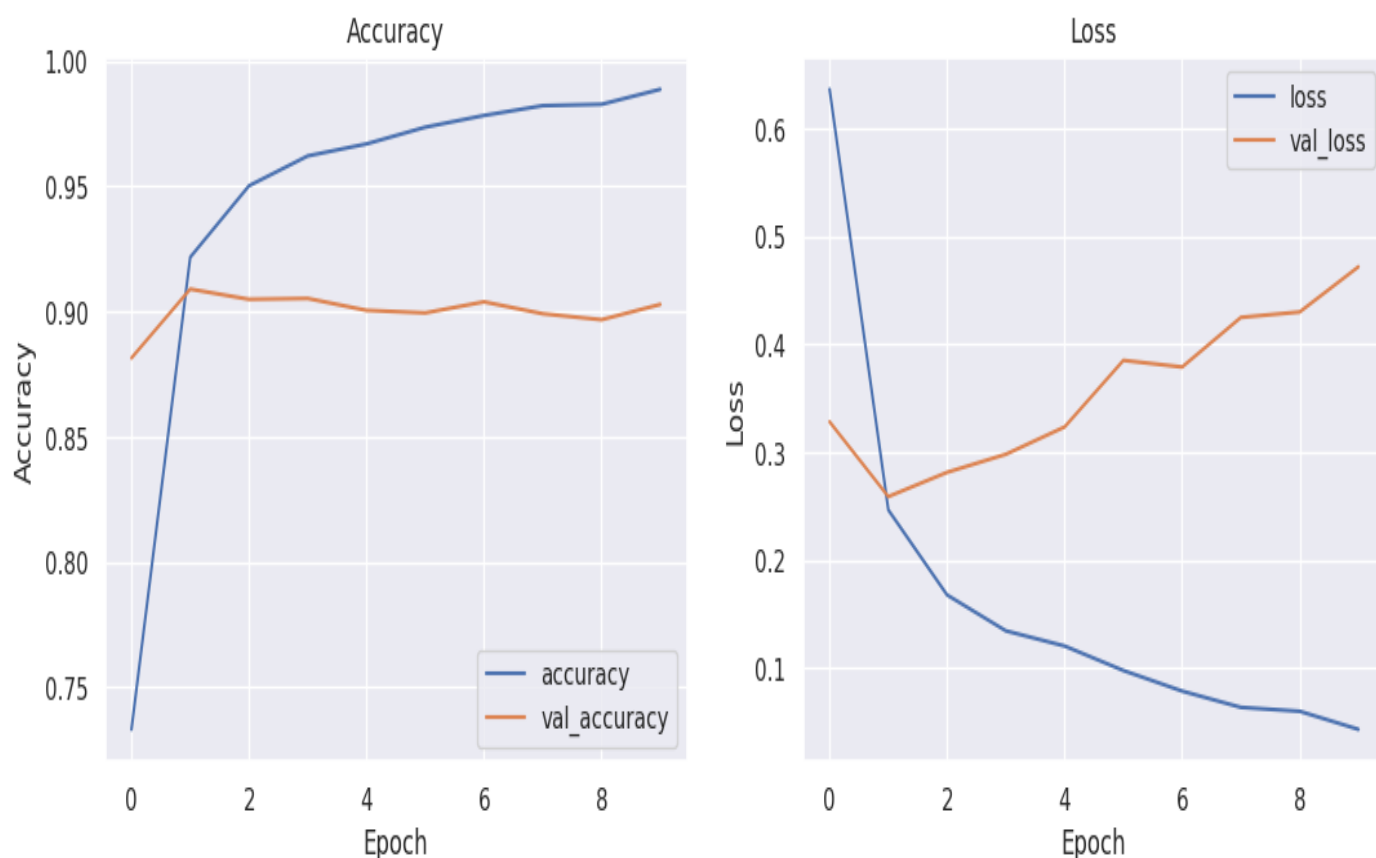## Training and Validation Accuracy & Loss



**Fig-3.17**

The model demonstrated excellent performance on the training set with an accuracy of 98.96% and low loss, indicating effective learning. Validation accuracy remained steady around 90%, though a slight increase in validation loss suggests minimal overfitting. On the test data, the model achieved a strong accuracy of 89.74% and an F1 score of 0.8970, with precision and recall closely aligned. These results indicate that the model performs reliably and consistently in classifying sentiment across unseen data.

## Model Prediction Output

WARNING:absl:Compiled the loaded model, but the compiled metrics have yet to be built. `model.compile_metrics` will be empty until you train or evaluate the model.
1/1 ━━━━━━━━━━━━━━━━━━━━ 1s 898ms/step
Predicted Class: Positive

The model was successfully loaded and used to make a prediction. Although a warning appeared indicating that compiled metrics haven't been built yet, this is expected behavior when a model is loaded and used for inference without retraining or re-evaluation.
For the given input, the model predicted the sentiment as "Positive", demonstrating its ability to classify text based on the learned patterns from the training data.

## Conclusion:

This weather classification project successfully demonstrates the effectiveness of a deep learning model in accurately identifying different weather conditions—Shine, Sunrise, Rain, and Cloudy—based on image data. Through the evaluation metrics, especially the ROC curves and AUC scores, the model has shown excellent performance with high predictive accuracy. The near-perfect AUC values for Sunrise and Rain (1.00), along with strong scores for Shine (0.97) and Cloudy (0.95), indicate that the model generalizes well across all classes. This suggests a robust learning of features that distinguish each weather type. Such a model can be practically implemented in smart weather monitoring systems, automated image tagging, or environmental analysis tools. Future enhancements may include expanding the dataset, fine-tuning the architecture, or integrating real-time image feeds to improve adaptability and real-world applicability