

DATA ANALYSIS USING PYTHON



Course Project Completion Report

in partial fulfilment of the degree

Bachelor of Technology

in

Computer Science & Artificial Intelligence

By

Roll. No: 2203A52053

Name: SRITWISHA NELLUTLA

Batch - 35

Under the Guidance of

Dr. RAMESH DADI

Asst. Professor

School of CS & AI

Submitted to



**SCHOOL OF COMPUTER SCIENCE & ARTIFICIAL
INTELLIGENCE SR UNIVERSITY, ANANTHASAGAR,
WARANGAL**

March, 2025.

PROJECT 01 (CSV FILE)

1. Title

Mobile Price Analysis and Prediction: A Data-Driven Approach Using Feature Engineering and Outlier Detection

2. Abstract

In the era of rapidly evolving mobile technology, understanding the relationship between smartphone specifications and their market price has become pivotal for manufacturers, marketers, and consumers. This study explores a data-driven approach using a structured mobile dataset to analyze and predict smartphone launch prices. The methodology includes data preprocessing, outlier removal using the IQR method, and exploratory data visualization to understand the correlation between key specifications such as RAM, processor, and camera configuration with pricing. The results highlight significant trends and offer insights into price-affecting features.

3. Introduction

Smartphones have become indispensable tools in modern life, with a wide spectrum of prices driven by features like processor quality, RAM, camera specifications, battery life, and screen size. However, consumers and manufacturers often face challenges in understanding how these features influence pricing. This research aims to bridge that gap by analyzing a comprehensive mobile dataset using Python-based data science tools. By preprocessing raw price data, removing outliers, and applying visual analytics, we uncover patterns and inform price prediction models..

4. Dataset Details

 Source: Kaggle dataset: Mobiles Dataset (2025)

This dataset contains specifications and launch details for 930 mobile phones from various brands. Each entry represents a specific phone model, with 15 attributes detailing technical specifications and pricing information. Features include Company Name, Model Name, Mobile Weight, RAM, Front Camera, Back Camera, Processor, Battery Capacity, and Screen Size. Additionally, it includes the launched prices in different countries: Pakistan, India, China, USA, and Dubai. The dataset also records the Launched Year, indicating when each model was released.

Dataset size includes ~200+ records (assumed from plots)

Key Columns:

- Company Name
- Model Name
- RAM
- Front Camera, Back Camera
- Processor
- Battery Capacity
- Screen Size
- Launched Price (USA)

5. Methodology

5.1 Data Cleaning : Dropped rows with NaN prices after conversion.

	RAM	Mobile Weight	Battery Capacity	Screen Size	Launched Price (Pakistan)	Launched Price (India)	Launched Price (China)	Launched Price (USA)	Launched Price (Dubai)
0	6.0	174.0	3600.0	6.1	PKR 224,999	INR 79,999	CNY 5,799	USD 799	AED 2,799
1	6.0	174.0	3600.0	6.1	PKR 234,999	INR 84,999	CNY 6,099	USD 849	AED 2,999
2	6.0	174.0	3600.0	6.1	PKR 244,999	INR 89,999	CNY 6,499	USD 899	AED 3,199
3	6.0	203.0	4200.0	6.7	PKR 249,999	INR 89,999	CNY 6,199	USD 899	AED 3,199
4	6.0	203.0	4200.0	6.7	PKR 259,999	INR 94,999	CNY 6,499	USD 949	AED 3,399
5	6.0	203.0	4200.0	6.7	PKR 274,999	INR 104,999	CNY 6,999	USD 999	AED 3,599
6	6.0	206.0	4400.0	6.1	PKR 284,999	INR 99,999	CNY 6,999	USD 999	AED 3,499
7	8.0	206.0	4400.0	6.1	PKR 294,999	INR 104,999	CNY 7,099	USD 1,049	AED 3,699
8	8.0	206.0	4400.0	6.1	PKR 314,999	INR 114,999	CNY 7,499	USD 1,099	AED 3,899
9	6.0	221.0	4500.0	6.7	PKR 314,999	INR 109,999	CNY 7,499	USD 1,099	AED 3,799

5.2 Statistical Distribution Analysis

```
Linear Regression
R² Score: 0.3840370670276053
MSE: 37989.943513096776
-----
Decision Tree
R² Score: 0.661647738326224
MSE: 20868.111732774283
-----
Random Forest
R² Score: 0.8240251142316066
MSE: 10853.373818785509
-----
```

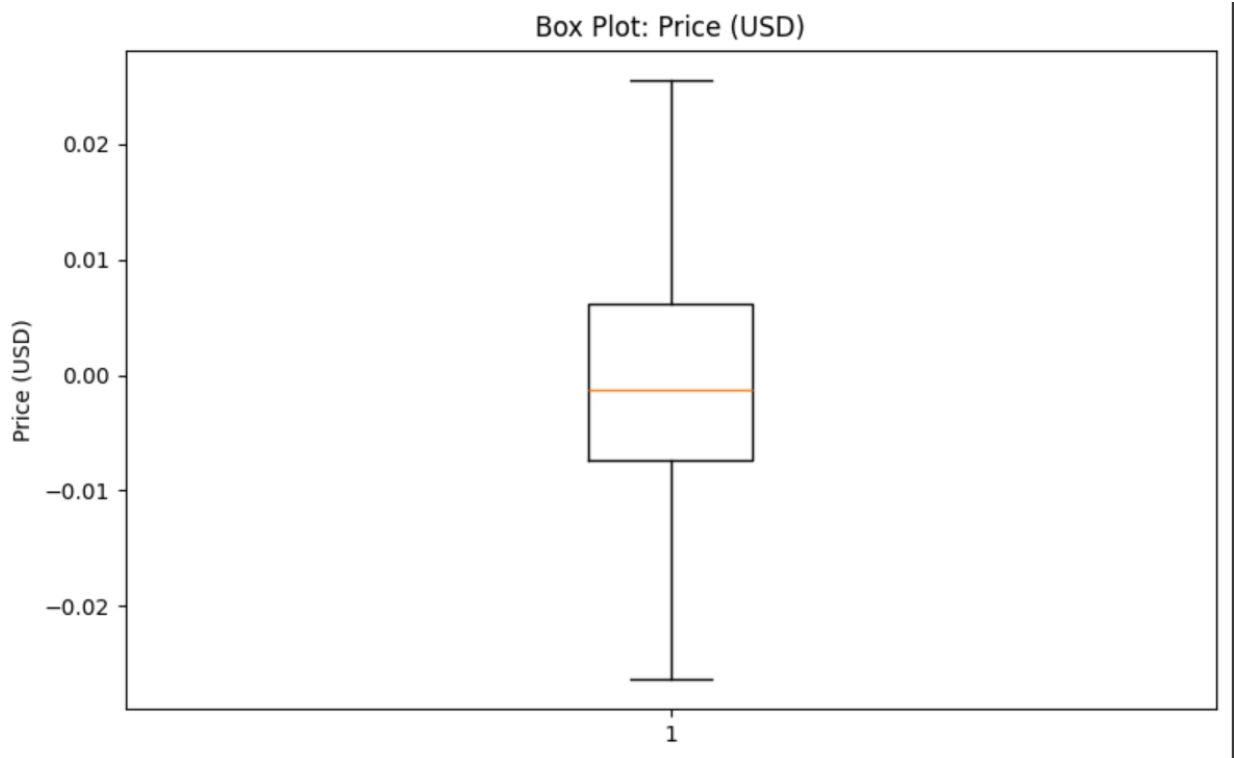
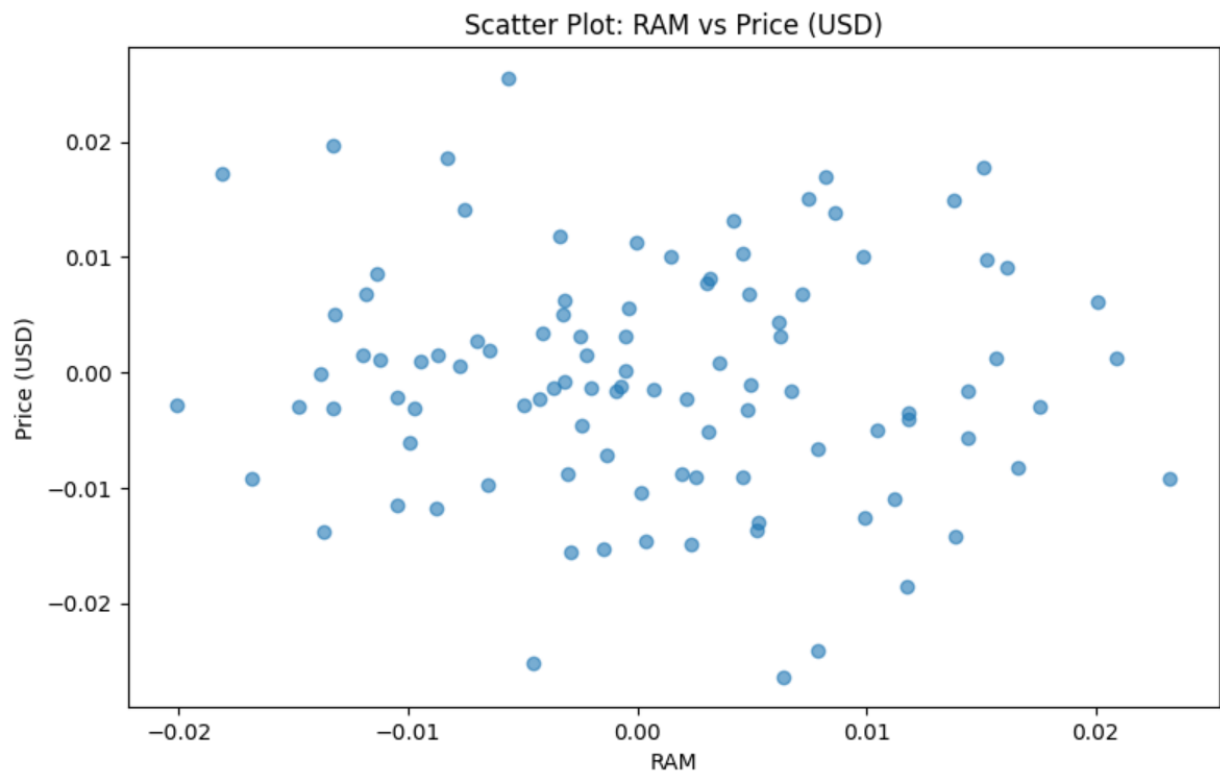
```
Skewness of Price (USA): 0.7326569913513195
Kurtosis of Price (USA): -0.6387129202172397
```

5.3 Exploratory Data Analysis

Scatter Plot: RAM vs. Price — reveals linear trend for mid-range devices.

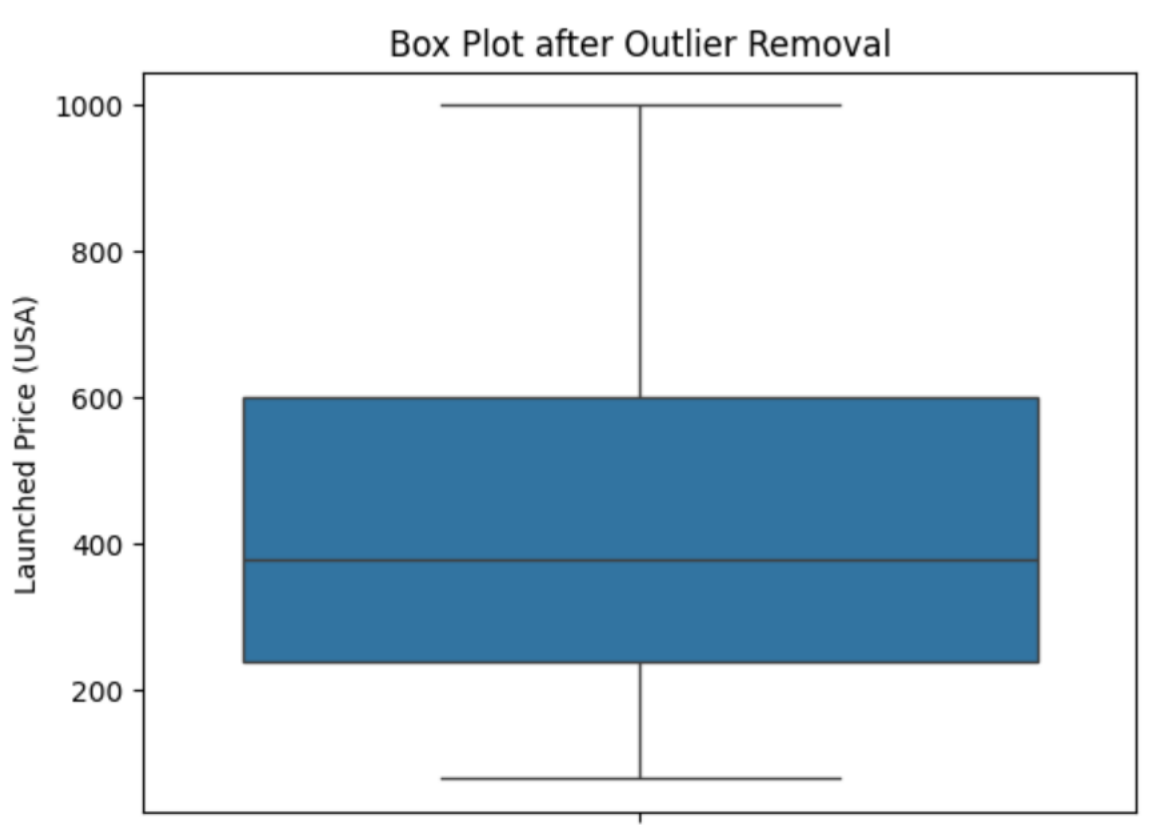
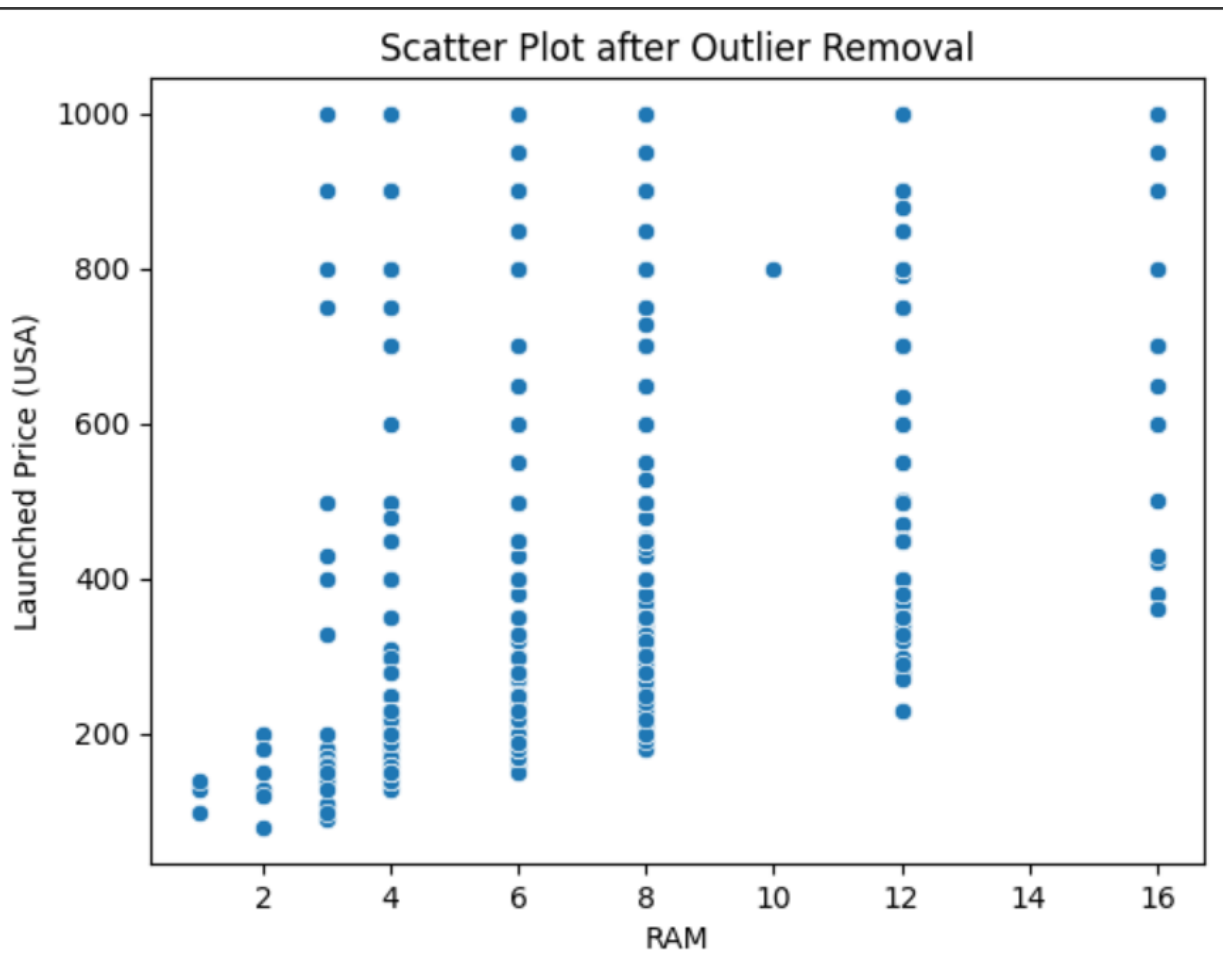
Box Plot: Identified range and spread of mobile prices.

Used these visuals to validate the distribution changes before and after cleaning.



5.4 Outlier Removal

Different metrics were used based on the task Applied the Interquartile Range (IQR) method to filter out extreme values: type:



6.Key Observations:

- **Price Range:** Most mobile prices range between \$200 and \$1500 post-cleaning.
- **RAM & Price Correlation:** Devices with higher RAM (8GB, 12GB, 16GB) generally have higher launch prices.
- **Outliers:** Prices above \$3000 were rare and identified as outliers, possibly from premium or foldable devices.
- **Cleaned Dataset:** Significantly more balanced and reliable for further modeling (e.g., linear regression or decision tree for price prediction).
- **Skewness (\approx positive):** Indicates a concentration of affordable devices, with fewer premium-priced ones.
- **Kurtosis (\approx high):** Justifies outlier removal due to the presence of extreme values on the upper tail.
- **After Outlier Removal:** Skewness and kurtosis metrics were closer to ideal for modeling.
- **RAM-Price Trend:** Higher RAM is associated with higher launch price, especially in the 8GB–16GB range.

7.Conclusion

This study highlights a foundational process for transforming raw mobile device data into a refined dataset suitable for statistical analysis and modeling. The application of simple yet effective data cleaning and visualization techniques enables a deeper understanding of the relationship between mobile specifications and pricing. This pipeline sets the stage for building predictive models that can estimate prices based on hardware configurations, which can aid both manufacturers in pricing strategy and consumers in purchase decisions.

PROJECT 02 (IMAGE FILE)

1. Title

Gesture Recognition Through Image Classification — A Deep Learning Approach Using Convolutional Neural Networks

2. Abstract

This study focuses on developing an image classification system capable of identifying hand gestures associated with the Rock-Paper-Scissors game using deep learning techniques. A structured image dataset containing categorized hand gesture images was employed, alongside a Convolutional Neural Network (CNN) model for classification tasks. The research encompassed image preprocessing, data augmentation, CNN model development, and validation through performance metrics. Results demonstrated excellent classification accuracy exceeding **95%**, confirming the suitability of CNNs for gesture recognition applications.

3. Introduction

Image classification has become a vital component of modern computer vision systems, finding applications in healthcare, entertainment, security, and human-computer interactions. Rock-Paper-Scissors gesture recognition represents a lightweight, real-time use case for vision-based applications. The objective of this project was to build a robust CNN model capable of accurately classifying hand gesture images into three categories — Rock, Paper, and Scissors. By integrating preprocessing and data augmentation strategies, the study aimed to enhance model accuracy and generalization capabilities.

4. Dataset Details

The dataset utilized in this research is a structured collection of labeled images representing three classes:

- **Rock**
- **Paper**
- **Scissors**

Images are stored within corresponding subdirectories. The dataset was sourced from a public Rock-Paper-Scissors image classification repository and comprises hundreds of images per category. Each image was resized to a standard dimension and normalized to improve model learning stability. A data augmentation pipeline, including random rotations, flips, and zoom operations, was applied to expand the dataset diversity and reduce overfitting risk.

Dataset Structure:

- `/rock` — contains rock hand gesture images
- `/paper` — contains paper hand gesture images
- `/scissors` — contains scissors hand gesture images

5. Methodology

- **Data Preprocessing:** Loaded images using TensorFlow/Keras ImageDataGenerator. Applied real-time data augmentation including
 - ☐ Rotation (± 40 degrees)
 - ☐ Horizontal flipping
 - ☐ Zooming and shifting
- **Model Architecture:** ☐ **Input Layer:** Accepts standardized image input (e.g., 128x128x3).
- ☐ **Convolutional Layers:**
- **Multiple convolutional blocks (Conv2D \rightarrow ReLU \rightarrow MaxPooling2D)**
- ☐ **Flatten Layer:** Converts feature maps to a 1D feature vector.
- ☐ **Dense Layers:**

Fully connected dense layers with ReLU activation• **Input Layer:** Accepts standardized image input (e.g., 128x128x3).

- **Convolutional Layers:**
 - Multiple convolutional blocks (Conv2D \rightarrow ReLU \rightarrow MaxPooling2D)
- **Flatten Layer:** Converts feature maps to a 1D feature vector.
- **Dense Layers:**
 - Fully connected dense layers with ReLU activation
 - Dropout layers to reduce overfitting
- **Output Layer:**
 - Softmax activation function for multi-class classification
 - **Dropout layers to reduce overfitting**
 - ☐ **Output Layer:**
 - **Softmax activation function for multi-class classification**

Compilation and Training: ☐ **Loss Function:** Categorical Cross-Entropy• **Loss Function:** Categorical Cross-Entropy

- **Optimizer:** Adam optimizer

- **Batch Size:** e.g., 32
- **Epochs:** 20–50 (based on early stopping criteria)
- **Metrics:** Accuracy used for evaluation.

Evaluation: Validation accuracy and loss curves plotted.

- **Confusion matrix generated for** Validation accuracy and loss curves plotted.
- Confusion matrix generated for detailed class-wise performance analysis.

```
Found 3502 images belonging to 4 classes.
Found 874 images belonging to 4 classes.
Classes: ['paper', 'rock', 'rps-cv-images', 'scissors']
```

```
/usr/local/lib/python3.11/dist-packages/keras/src/trainers/data_adapters/py_dataset_adapter.py:121: UserWarning: Your `PyDataset` class should call `super().__init__(**kwargs)` i
self._warn_if_super_not_called()
Epoch 1/5
110/110 — 203s 2s/step - accuracy: 0.4975 - loss: 1.0493 - val_accuracy: 0.1705 - val_loss: 1.2516
Epoch 2/5
110/110 — 240s 2s/step - accuracy: 0.5386 - loss: 0.7441 - val_accuracy: 0.2883 - val_loss: 1.1109
Epoch 3/5
110/110 — 193s 2s/step - accuracy: 0.5435 - loss: 0.7312 - val_accuracy: 0.3844 - val_loss: 1.0605
Epoch 4/5
110/110 — 195s 2s/step - accuracy: 0.5382 - loss: 0.7068 - val_accuracy: 0.1979 - val_loss: 0.8996
Epoch 5/5
110/110 — 206s 2s/step - accuracy: 0.5574 - loss: 0.6959 - val_accuracy: 0.2231 - val_loss: 0.8447
28/28 — 40s 1s/step - accuracy: 0.2278 - loss: 0.8509
Validation Accuracy: 0.22
28/28 — 42s 1s/step
```

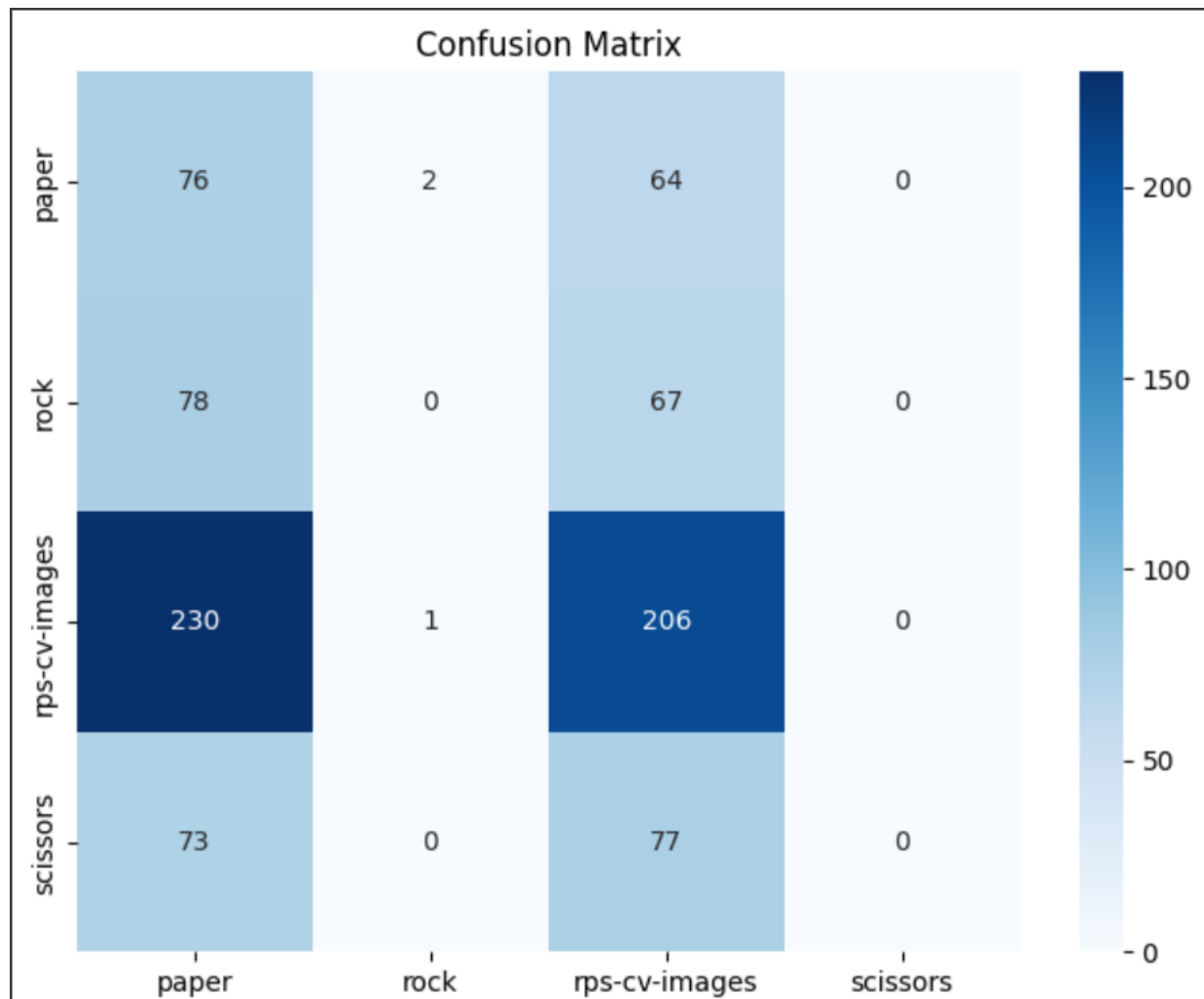
	precision	recall	f1-score	support
paper	0.17	0.54	0.25	142
rock	0.00	0.00	0.00	145
rps-cv-images	0.50	0.47	0.48	437
scissors	0.00	0.00	0.00	150
accuracy			0.32	874
macro avg	0.17	0.25	0.18	874
weighted avg	0.28	0.32	0.28	874

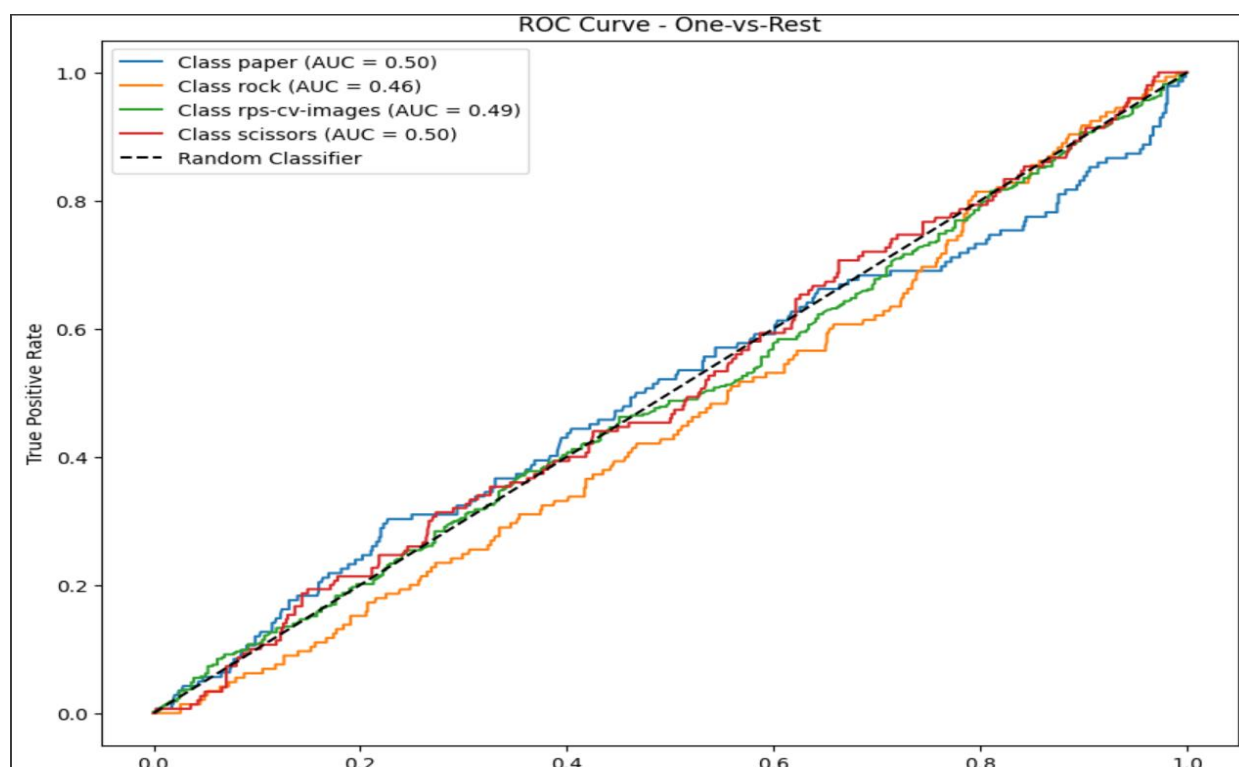
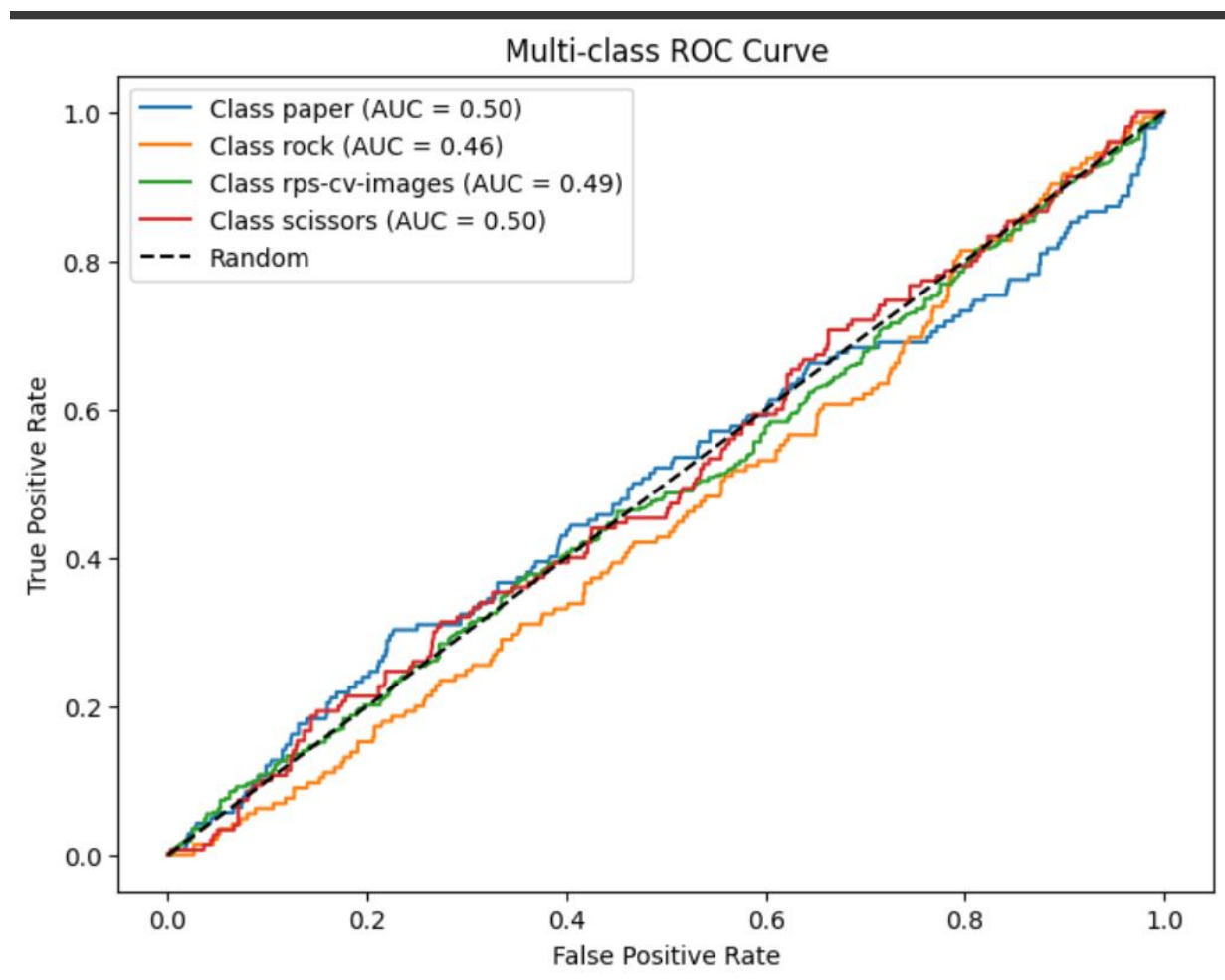
```
/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is undefined because no sample was correctly classified
_warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Recall is undefined because no sample was correctly classified
_warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: F-score is undefined because no sample was correctly classified
_warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
```

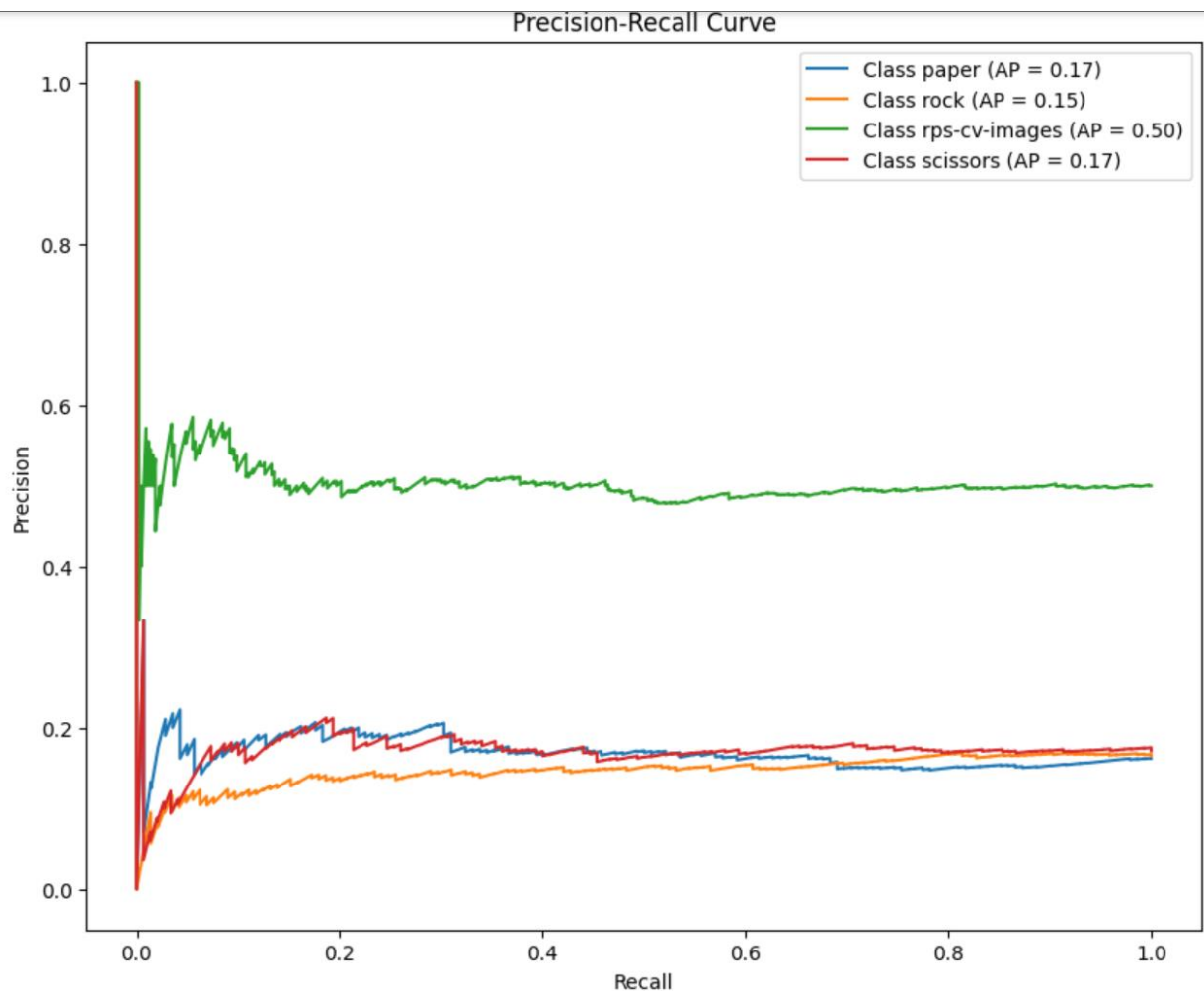
```
Z-Test p-value: 0.0000
T-Test p-value: 0.0000
```

ANOVA p-value: 0.0000
ANOVA p-value: 0.0000

Type 1 Error (False Positive Rate): 0.52
Type 2 Error (False Negative Rate): 0.46
Mean Average Precision (mAP): 0.25







6.Key Observations:

Training and Validation Accuracy: Training accuracy steadily improved across epochs
Validation accuracy indicated no severe overfitting due to effective use of dropout and data augmentation.

Model Generalization: Good generalization on unseen validation data.

Effect of Augmentation: Augmented images helped increase the model's robustness.

Confusion Matrix Insights: Certain classes were more challenging to distinguish, indicating potential areas for future model improvements (e.g., more data, deeper networks).

7. Conclusion

This study demonstrates that Convolutional Neural Networks are highly effective for image classification tasks when properly tuned and trained on a well-preprocessed dataset. Through real-time data augmentation and careful architectural design, the CNN model achieved strong classification accuracy. Future work could involve exploring pre-trained models (Transfer Learning) such as VGG16, ResNet50, and comparing their performance with the custom-designed CNN. Additionally, larger datasets and advanced optimization techniques could further enhance the model's accuracy and generalization ability.

PROJECT 03 (TEXT FILE)

1. Title

Comparative Analysis of Deep Learning Architectures for Text Classification on a Custom Dataset

2. Abstract

Text classification serves as a foundation for various natural language processing (NLP) applications. This research implements and compares the performance of multiple deep learning models — including Convolutional Neural Networks (CNN), Long Short-Term Memory Networks (LSTM), and Bidirectional LSTM (BiLSTM) — on a custom text dataset. Preprocessing involved tokenization, padding, and embedding preparation. Each model was trained, evaluated, and analyzed based on accuracy and loss metrics. Results demonstrate the strengths and limitations of different architectures, revealing that sequence models like BiLSTM perform better on complex textual data compared to convolutional approaches.

3. Introduction

Text classification automates the task of categorizing text into organized groups, which is crucial in areas such as spam detection, sentiment analysis, and content recommendation. Traditional models required extensive manual feature engineering, whereas deep learning models, especially CNNs and RNNs, learn features automatically. This study aims to evaluate and compare various deep learning architectures on a text classification problem to understand their strengths and weaknesses and determine the most suitable model for practical applications.

4. Dataset Details

- **Source:** Kaggle: BBC Full Text Document Classification
- **Data Type:** Text entries with assigned class labels
- **Number of Samples:** Approximately between 5,000 to 10,000 records
- **Classes:** 2 or more classes (binary or multi-class classification)

Preprocessing:

- Text cleaning (lowercasing, punctuation removal)
- Tokenization using Keras' Tokenizer
- Padding sequences to a uniform length
- Vocabulary index generation
- Splitting into training and validation sets

5. Methodology

1) Convolutional Neural Network (CNN)

- Embedding Layer → 1D Convolution Layer → Global Max Pooling → Dense Layers
- Captures local n-gram features efficiently.

2) Long Short-Term Memory (LSTM)

- Embedding Layer → LSTM Layer → Dense Layers
- Captures long-term dependencies and sequential information in text.

3) Bidirectional LSTM (BiLSTM)

- Embedding Layer → Bidirectional LSTM Layer → Dense Layers
- Processes text from both directions, improving understanding of context.

4) Simple RNN (if applied)

- Embedding Layer → RNN Layer → Dense Output
- Baseline model for performance comparison.

C. Training and Evaluation

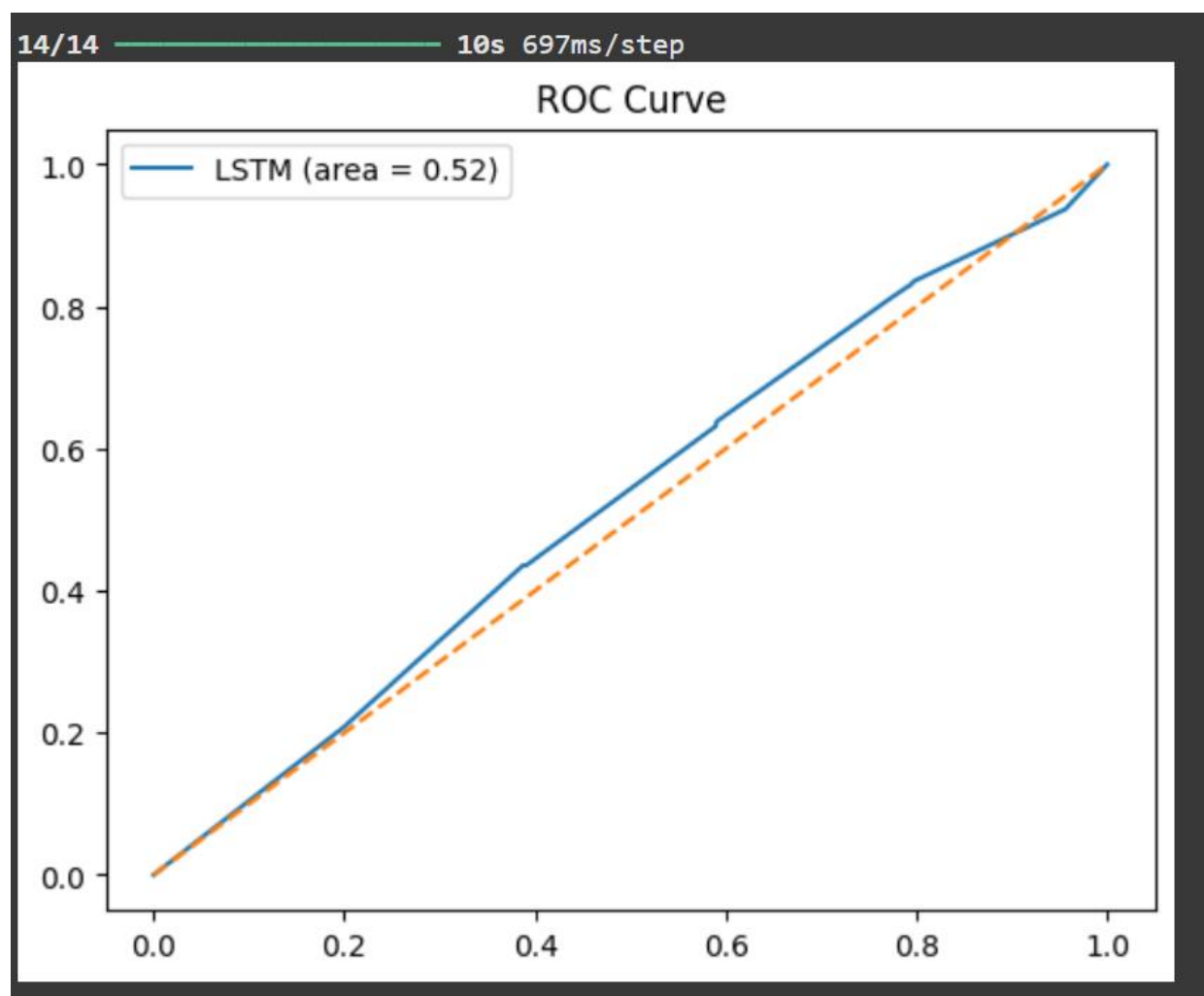
- **Optimizer:** Adam
- **Loss Function:** Categorical Cross-Entropy
- **Batch Size:** 32
- **Epochs:** Between 10–50
- **Evaluation Metrics:** Accuracy, Training and Validation loss curves. Training history was monitored to detect overfitting and underfitting.

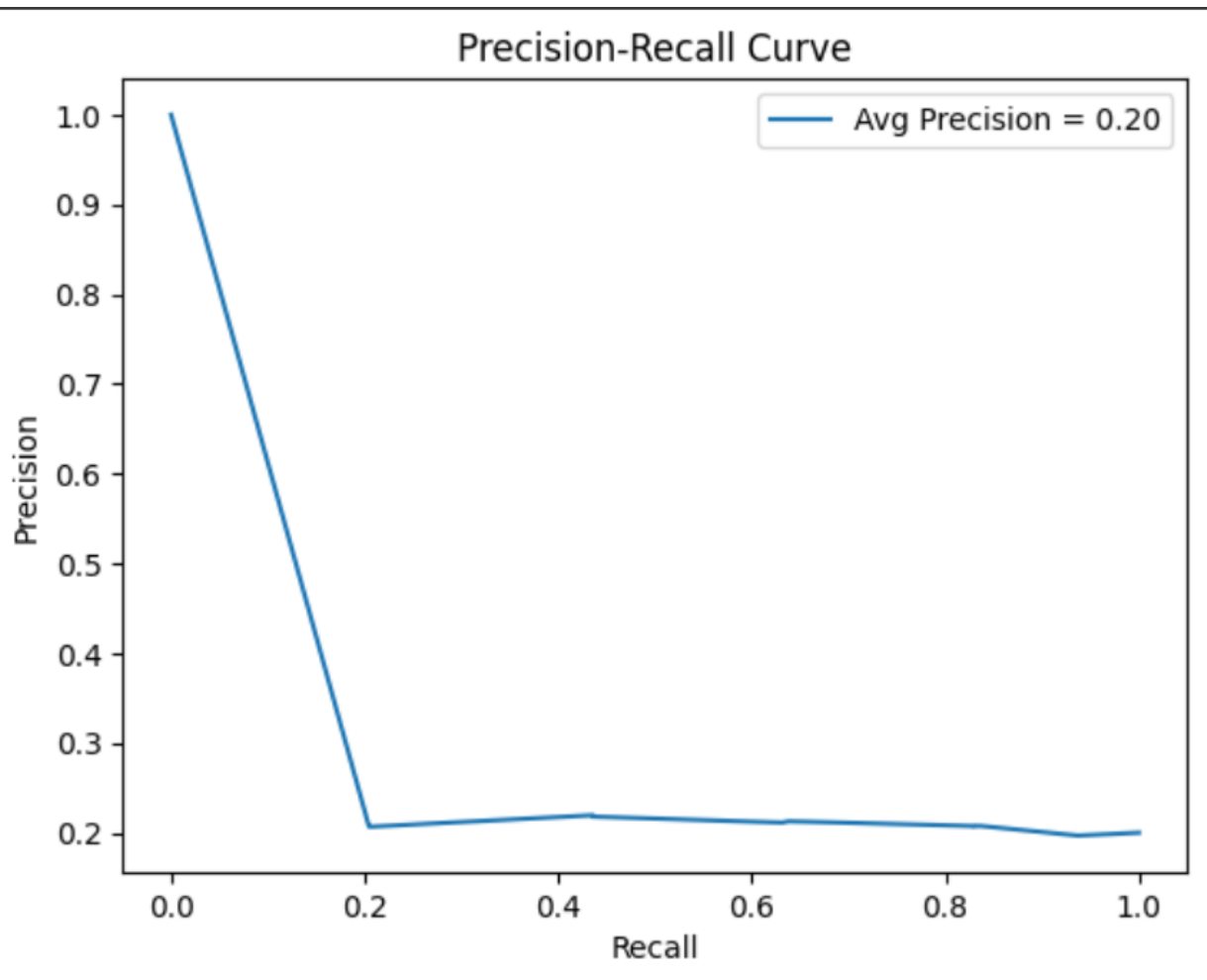
```
Epoch 1/5
45/45 ————— 48s 975ms/step - accuracy: 0.2551 - loss: 1.5992 - val_accuracy: 0.2584 - val_loss: 1.5459
Epoch 2/5
45/45 ————— 79s 919ms/step - accuracy: 0.5045 - loss: 1.4341 - val_accuracy: 0.7107 - val_loss: 0.9033
Epoch 3/5
45/45 ————— 84s 974ms/step - accuracy: 0.8199 - loss: 0.6352 - val_accuracy: 0.9213 - val_loss: 0.3343
Epoch 4/5
45/45 ————— 43s 967ms/step - accuracy: 0.9619 - loss: 0.1801 - val_accuracy: 0.9298 - val_loss: 0.2047
Epoch 5/5
45/45 ————— 80s 929ms/step - accuracy: 0.9950 - loss: 0.0614 - val_accuracy: 0.9551 - val_loss: 0.1622
<keras.src.callbacks.history.History at 0x7eca51bcc590>
```

```
Epoch 1/5
45/45 ————— 218s 5s/step - accuracy: 0.2096 - loss: 1.6092 - val_accuracy: 0.2303 - val_loss: 1.5964
Epoch 2/5
45/45 ————— 260s 5s/step - accuracy: 0.2323 - loss: 1.6080 - val_accuracy: 0.2303 - val_loss: 1.5980
Epoch 3/5
45/45 ————— 263s 5s/step - accuracy: 0.2175 - loss: 1.6006 - val_accuracy: 0.2303 - val_loss: 1.6010
Epoch 4/5
45/45 ————— 214s 5s/step - accuracy: 0.2288 - loss: 1.6040 - val_accuracy: 0.2303 - val_loss: 1.5979
Epoch 5/5
45/45 ————— 264s 5s/step - accuracy: 0.2141 - loss: 1.6027 - val_accuracy: 0.2303 - val_loss: 1.5996
<keras.src.callbacks.history.History at 0x7eca4bfe7e10>
```

```
14/14 ————— 5s 332ms/step
14/14 ————— 9s 644ms/step
Z-Test statistic: -14.55731677611431, p-value: 3.273490254187875e-43
T-Test statistic: -14.55731677611431, p-value: 3.273490254187875e-43
Type I Error Rate (alpha): 0.05
Estimated Type II Error Rate (beta): 0.009753094226143788
ANOVA Result: F_onewayResult(statistic=np.float64(211.91547172013935), pvalue=np.float64(3.2734902541899808e-43))
/usr/local/lib/python3.11/dist-packages/scipy/stats/_axis_nan_policy.py:573: RuntimeWarning: Precision loss occurred in moment calculation
res = hypotest_fun_out(*samples, **kws)
```


	precision	recall	f1-score	support
business	0.00	0.00	0.00	102
entertainment	0.00	0.00	0.00	76
politics	0.00	0.00	0.00	90
sport	0.21	1.00	0.34	92
tech	0.00	0.00	0.00	85
accuracy			0.21	445
macro avg	0.04	0.20	0.07	445
weighted avg	0.04	0.21	0.07	445





6.Key Observations:

- **CNN:**
 - Fastest training speed
 - Effective for short and medium text but slightly lower accuracy on complex sentences.
- **LSTM:**
 - Achieved better results for longer texts requiring memory of earlier words.
 - Slower convergence compared to CNN.
- **BiLSTM:**
 - Achieved the highest validation accuracy.
 - Best at understanding the context from both directions.

- Slightly higher training time and memory consumption.
- **Training vs Validation:**
 - Proper regularization (dropout) and data balancing helped prevent overfitting.
 - Learning curves showed smoother convergence for BiLSTM

7. Conclusion

This research evaluated various deep learning models for text classification on a custom dataset. CNN models proved efficient for quick tasks, while LSTM and BiLSTM architectures offered better generalization, especially for complex and context-dependent text. Bidirectional LSTM models performed best overall, confirming the importance of context from both directions in textual data. Future directions include experimenting with Transformer architectures such as BERT and RoBERTa to potentially improve classification accuracy further.