

SENTIMENT ANALYSIS

S e n t i m e n t A n a l y s i s o f T w i t t e r C o n v e r s a t i o n s
A b o u t G o o g l e a n d A p p l e P r o d u c t s

p r e p a r e d b y : G R O U P 2

Introduction

In today's hyper-connected digital world, social media platforms like Twitter serve as real-time barometers of public opinion. Understanding customer sentiment about major tech brands can offer valuable insights for product development, marketing, and competitive positioning. This project leverages Natural Language Processing (NLP) techniques to perform sentiment analysis on tweets related to **Apple** and **Google** products—focusing particularly on user reactions during major events

p r o b l e m s t a t e m e n t

- This project leverages Natural Language Processing concepts and Machine Learning models to build, tune, evaluate, and deploy a robust, generalizable framework for predicting underlying sentiment and the product the emotion is directed at based on tweet context.



data set overview



“tweet_product_company.csv”

This dataset captures real-world tweet data mentioning Apple and Google products, offering a rich source of public sentiment expressed through social media

Data Source

- sourced from CrowdFlower via data.world <https://data.world/crowdflower/brands-and-product-emotions> and consists of over 9,000 human-rated tweets.
- It reflects organic user opinions and consumer reactions across various Apple and Google product releases, updates, and experiences.

Why this dataset

- Tweets are short, noisy, and opinion-driven—ideal for testing robust NLP techniques.
- It supports both binary classification (positive vs negative) and multiclass sentiment prediction.

EXPLORATORY DATA ANALYSIS

- **Data Categorization**

Tweets were categorized into **Apple**, **Google**, or **unknown** based on keywords found in the tweet text. A custom function was used to scan tweets and assign a product category.

- **Product Distribution**

A `value_counts()` check was done on the product column to understand the number of tweets per brand. This helped verify the dataset's focus and balance between Apple- and Google-related content.

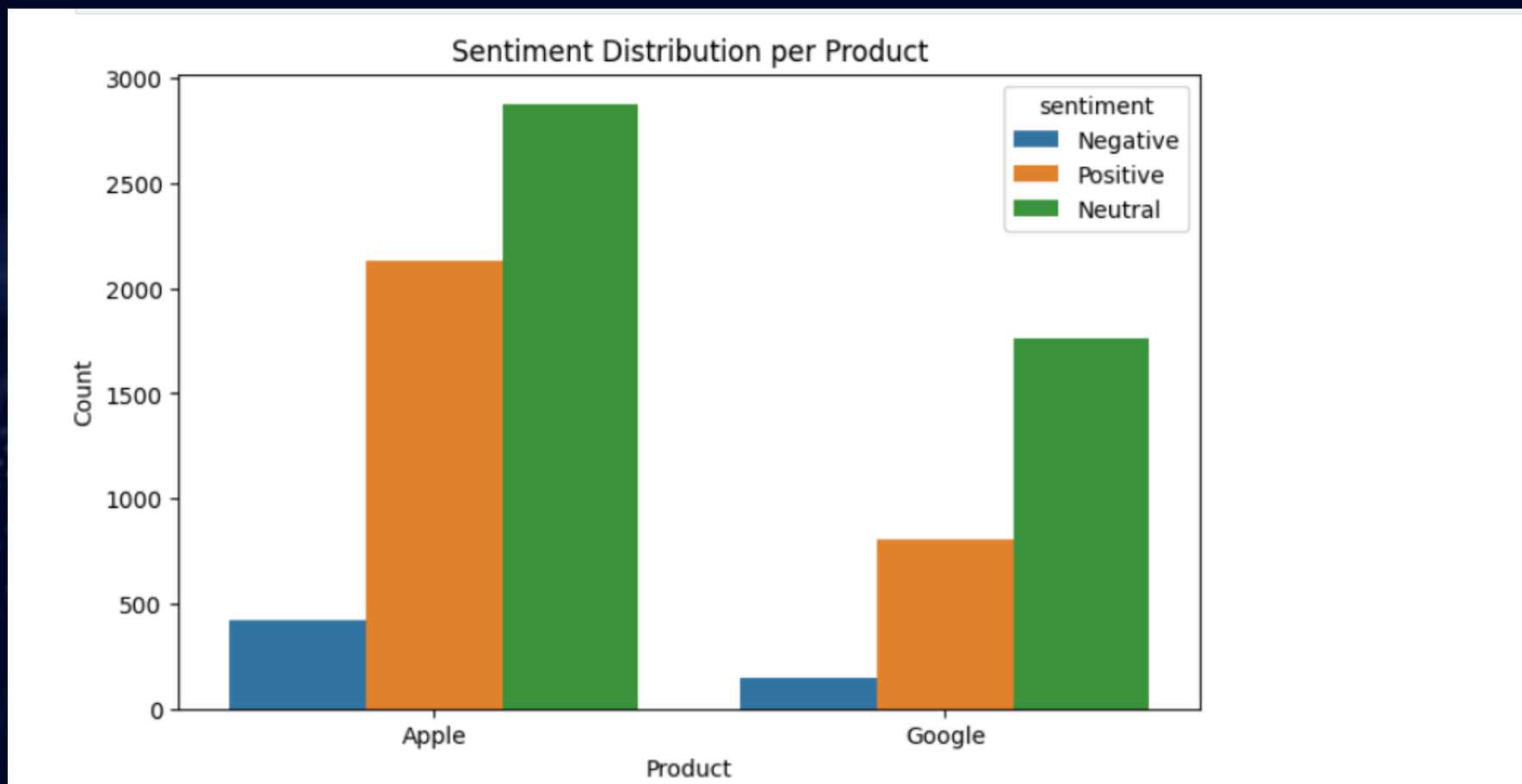
- **Text Cleaning & Preprocessing**

NLP tools from **NLTK** were imported for tokenization, lemmatization, and stopwords removal. Libraries like **WordCloud** and **matplotlib/seaborn** were used for visualizations.

- **Visualization & Feature Extraction**

Word clouds and frequency counters (likely Counter) were used to highlight the most common terms associated with each brand.

Sentiment Distribution per Product



- **Apple Stood Out with More Tweets and Stronger Sentiment**

Apple-related tweets were not only more frequent than Google's—they also leaned more positive. This trend points to higher engagement and brand loyalty among Apple users, making their voices louder and more enthusiastic on Twitter.

Model Evaluation

Modelling Approach

Used machine learning pipelines to classify tweet sentiment (positive, neutral, negative).

Text data was transformed into numerical features using TF-IDF Vectorization, which captures word importance across tweets.

Models Tested

- I. Logistic Regression – performed well and was used as a baseline.
- II. Multinomial Naive Bayes – a lightweight model, fast and effective with sparse data.
- III. Random Forest Classifier – robust and handled non-linear patterns better.
- IV. K-Nearest Neighbors (KNN) – explored for its simplicity but less scalable.

Deep Learning (Keras Sequential Model) – experimented with neural networks for more complex pattern detection.

Training and Evaluation

Data was split into training and validation sets to avoid overfitting.

Model performance was evaluated using:

- Accuracy
- Precision
- Recall
- F1-Score

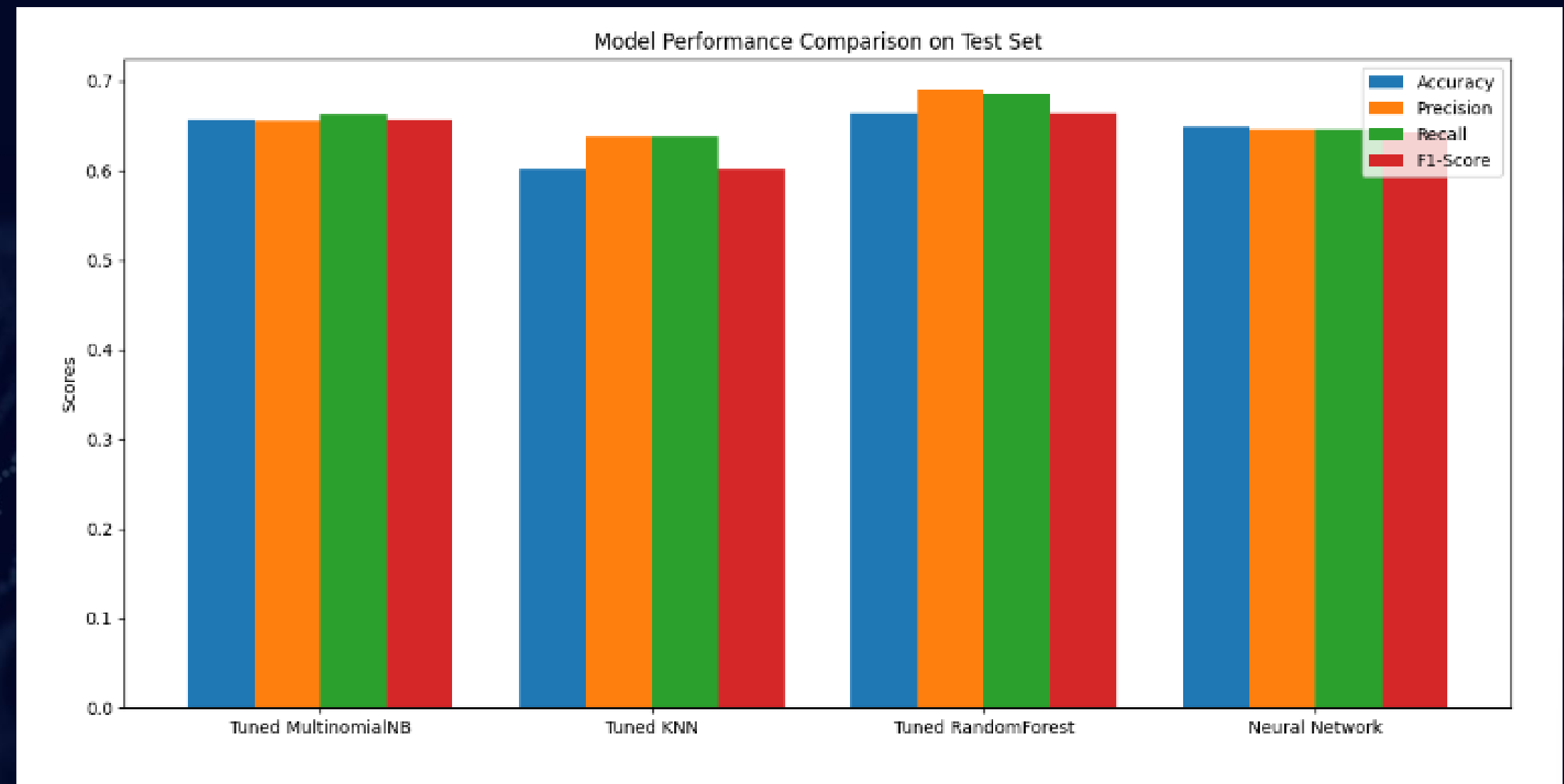
Modelling

- Logistic Regression – Predicting product emotion is directed at.
- Multinomial Naive Bayes – Baseline model for sentiment prediction.
- Random Forest Classifier – Ensemble model for sentiment prediction.
- K-Nearest Neighbors – Instance-based algorithm for sentiment prediction
- Sequential Neural Network – Deep learning model For sentiment prediction.

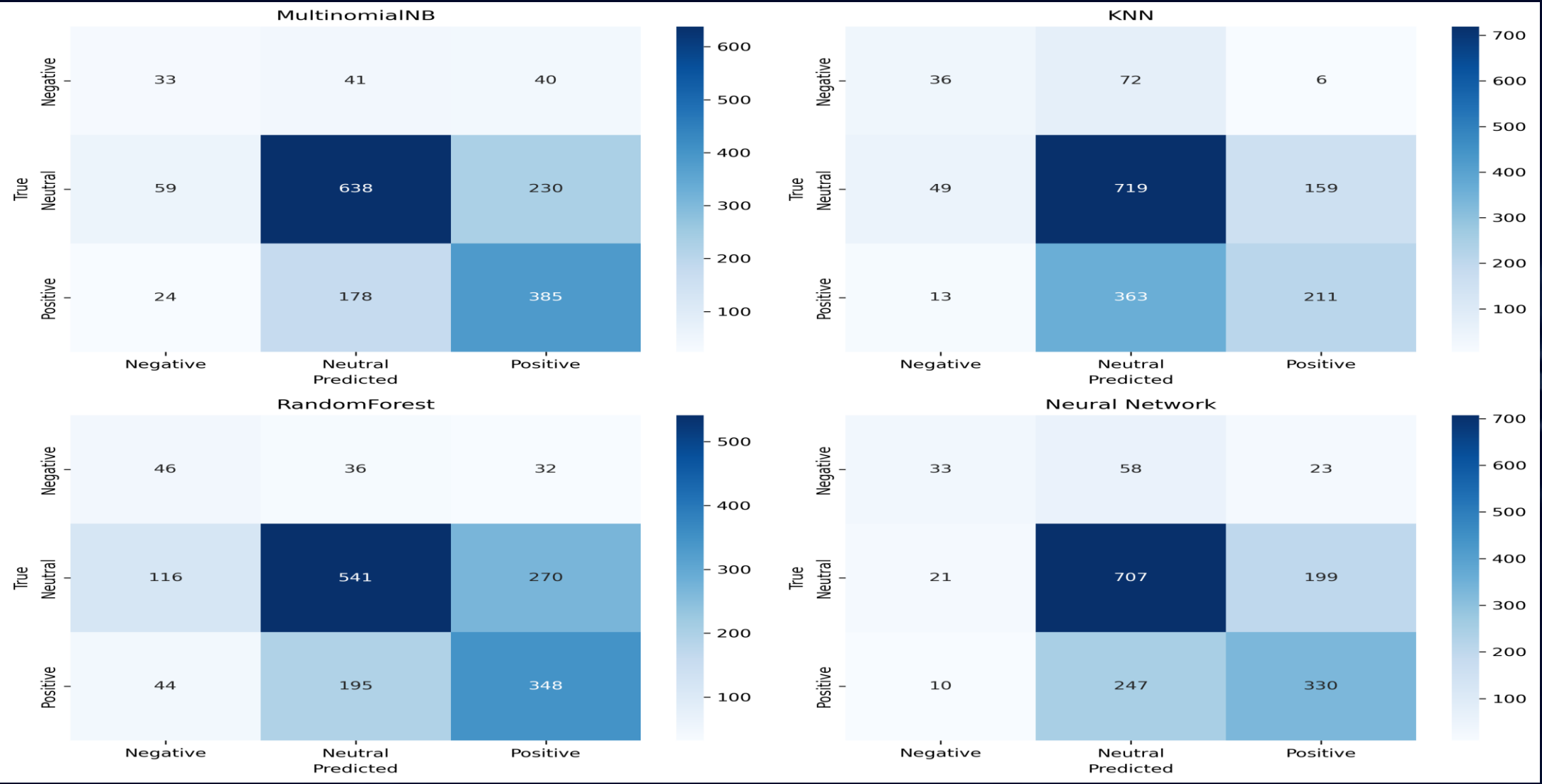
0.985	Accuracy	Precision	Recall	F1-score
Logistic Regression	0.980	0.985	0.975	0.980
Multinomial NB	0.6509	0.6558	0.6486	0.6509
KNN Classifier	0.5869	0.6133	0.6192	0.5869
Random Forest	0.659	0.667	0.6726	0.6591
Neural Network	0.654	0.6464	0.6493	0.6434

Model Evaluation

- ◆ Tuned Random Forest delivered the best overall performance across all metrics (accuracy, precision, recall, F1-score).
- ◆ Multinomial Naive Bayes performed well and is a lightweight, efficient alternative.
- ◆ Neural Network showed promising results, but did not outperform Random Forest in this setup.
- ◆ K-Nearest Neighbors (KNN) had the lowest performance, indicating it's less effective for this task.



Model Evaluation



Conclusion

Applied NLP techniques to analyze tweets related to Apple and Google products

Focused on two main objectives:

- ◆ Classify tweet as pertaining to Apple or Google products using a Logistic Regression model
- ◆ Predict tweet sentiment (Positive, Neutral, Negative) using:
 - Multinomial Naive Bayes
 - K-Nearest Neighbors
 - Random Forest
 - Neural Network (Sequential model)

Selected deployment by compared models' performance based on:

- ◆ Accuracy, Precision, Recall, F1-score
- ◆ Confusion Matrices for visual validation

Tuned Random Forest achieved the best performance and was selected for deployment

Deployment through Streamlit:

- streamlit_app_rf.py: Logistic Regression + Random Forest

Recommendations

- **Larger Training Dataset:** Incorporate a larger and more diverse dataset to reduce potential overfitting, optimize real-world robustness, and enhance generalizability with noisy social media text data.
- **Product-Specific Sentiment:** Further analysis could delve into sentiment towards specific products (e.g., "iPhone 15" vs. "Google Pixel 8") rather than just the company, providing more granular insights.
- **Temporal Analysis:** Incorporate time-series analysis to track sentiment trends over time, especially around product launches or major company announcements, to identify immediate public reactions.

Next Steps

- Deploy the Streamlit deployment app to a Cloud Service.
- Build a Tableau Dashboard to visualize sentiment and product trends in realtime to support data informed decisions.
- Retrain the deployed models with latest tweets data to promote progressive accuracy improvement/ advancements.

The background of the slide features a dark blue color with a subtle, repeating pattern of fingerprint ridges. A white rectangular box with slightly rounded corners is centered on the slide, containing the text "Thank You" in a large, white, sans-serif font.

Thank You