NLP Project for

Disaster Tweet Classification

In this project, I built a **Tweet Classification Model** to detect whether a tweet is related to a disaster or not. I started by cleaning the text—removing URLs, special characters, digits, and stopwords, and applied lemmatization to standardize the words.

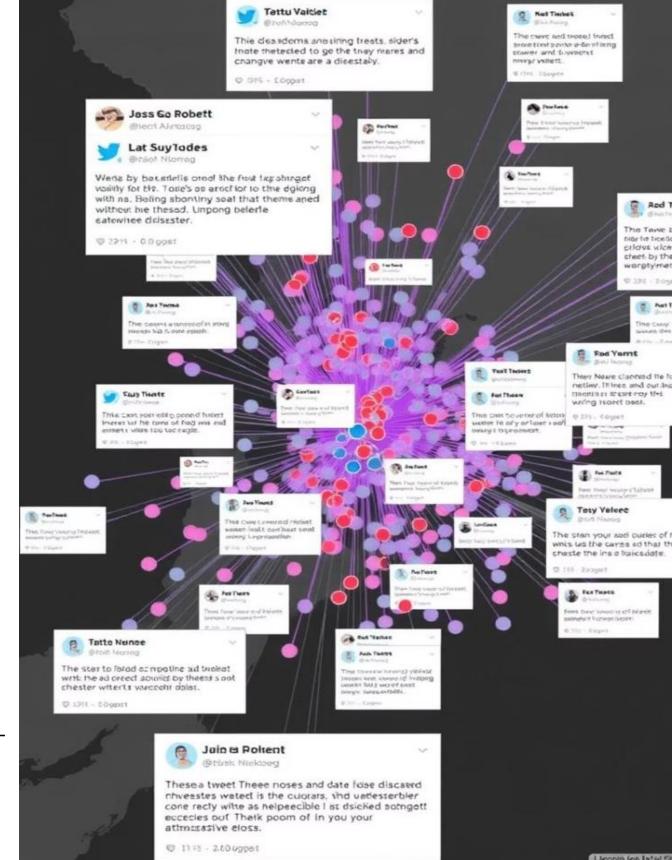
I engineered **additional features** like sentiment polarity (using **TextBlob**), tweet length, hashtag count, and mention presence. These were combined with **TF-IDF vectors** to form the final input for the model.

I experimented with several machine learning algorithms and finalized **Logistic Regression** based on its accuracy and interpretability. After training, I created a **real-time prediction system** that takes user input, processes it, and classifies the tweet with a **confidence score**.

To make the model more accessible, I deployed it using **Streamlit**, building an interactive web app where users can test the model, explore example tweets, and download results as a CSV.

This project brought together NLP, feature engineering, machine learning, real-time prediction, and deployment to deliver a complete end-to-end tweet classification solution.





Contributions of Key Libraries & Modules to the Project



Pandas & Numpy

I used these for loading, cleaning, and manipulating structured data and handling numerical computations.





re, NLTK & TextBlob

I used these for cleaning tweet text, removing stopwords, lemmatizing, and calculating sentiment polarity.





Matplotlib, Seaborn & WordCloud

I used these for visualizing EDA insights like class distribution, word frequencies, and generating word clouds.





Scikit-learn (sklearn)

I used this for TF-IDF vectorization, building ML models, evaluating performance using accuracy, F1-score, & confusion matrix.





Scipy.sparse

I used this for combining sparse TF-IDF features with engineered features like sentiment score, tweet length, and hashtag count.





Joblib & Pickle

I used these for saving and loading trained models and vectorizers for real-time prediction.





Tqdm & Collections

I used these for monitoring loop progress and analyzing frequency of tokens and hashtags.





OS & IPython.display

I used these for handling files/directories & improving output readability in the notebook.





Streamlit

I used this to build and deploy an interactive web app for real-time tweet classification.





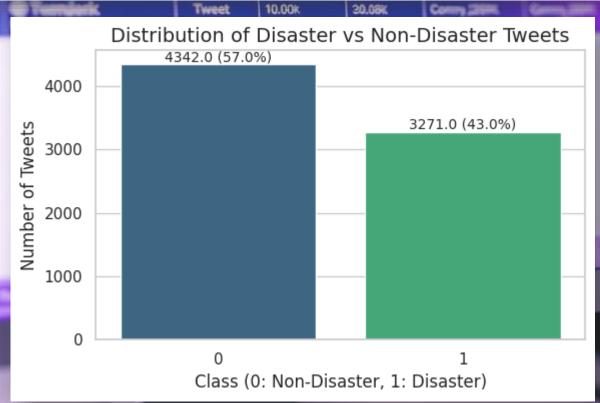
Transformers (Hugging Face)

I used this optionally to compare traditional models with state-ofthe-art deep learning models.



Identifying Missing Values & Handling

```
1 # Handling Missing Values
      print("Missing Values Before Handling:")
      print(df.isnull().sum()[df.isnull().sum() > 0])
      df.fillna("", inplace=True) # Filling missing values with an empty string
      print("\nMissing Values After Handling:")
     print(df.isnull().sum())
Missing Values Before Handling:
keyword
              61
location
           2533
dtype: int64
Missing Values After Handling:
keyword
location
text
target
dtype: int64
```



Dataset Overview

	Da	taset Overview
		Loaded dataset with 7,613 rows and 5 columns: 'id', 'keyword'
		`location`, `text`, `target`
		Used `df.head()` to preview top 10 rows
	Da	ta Types & Memory Usage
		`id` and `target`: Integer
		`keyword`, `location`, `text`: Object
		Missing values found in `keyword` and `location`
<u>^</u>	На	ndling Missing Values
		`keyword`: 61 missing values
		`location`: 2,533 missing values
		Filled missing entries with empty string ""
		After handling: No missing values remained
1	Dis	saster vs. Non-Disaster Distribution
		Created bar plot using **Seaborn + Matplotlib**
		57% Non-Disaster Tweets (4,342)
		43% Disaster Tweets (3,271)
		Slight imbalance, manageable without resampling

Data Cleaning & Frequency Analysis

Data Preparation & Cleaning

- ✓ Removed URLs, special characters, & numbers
- ✓ Removed boring stopwords
- **Outcome:** Clean, sharp text with only the *good stuff* left!

Sentiment Analysis

- Used **TextBlob** to score each tweet's *emotional tone*
- Polarity score from -1 (negative) to +1 (positive)
- Added a new sentiment column
- Helps the model *understand feelings*, not just words!

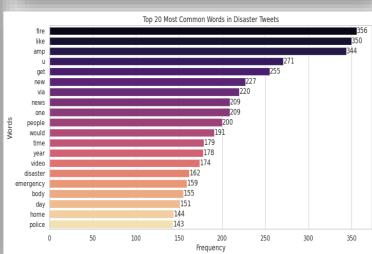
Keyword Frequency Analysis

- Tokenized + Lemmatized + Counted via `Counter()`
- Top Words:
- fire" (356), "like" (350), "amp" (333), "get" (255), "new" (227)
 - bisaster words like "emergency," "disaster," "fire" dominate!

Word Cloud Visualization

- Built a Word Cloud from tokenized words
- Bigger words = More mentions
- Spotlight on: "fire," "emergency," "disaster," "police," "news"
- **P** Confirms the theme: **Disaster-Focused Tweets**

```
1 # Initialize NLP tools
  ps = PorterStemmer()
   lemmatizer = WordNetLemmatizer(
   stop words = set(stopwords.words('english'))
       phrase = re.sub(r"'s\b", " is", phrase) # Only replaces 's if it's a word ending
       nbrase = re.sub(r"\b[ull]\b", "us", nbrase) # Fixes 'u' to 'us
       phrase = re.sub(r"\bu\b", "us", phrase, flags=re.IGNORECASE) # Fixes 'u' to 'us
                                   phrase) # Removes digits like phone numbers
       sentence = decontracted(sentence) # Apply phrase normalization
       tokens = [lemmatizer.lemmatize(word.lower()) for word in tokens if word.lower() not in stop_words]
       # Rejoining processed tokens into a sentence
```



Word Cloud Visualization of Word Frequencies in Disaster Tweets







Tweet Location Analysis



"Unknown" Rules the Map

- Most tweets have no location data
- ◆ 1460 Normal | 1077 Disaster tweets = "Unknown"
- Shows a **huge gap in geotagging** on Twitter



us **USA: Top Disaster Tweeter (Excl. Unknown)**

- **70 Disaster Tweets** vs. **39 Normal** from the USA
- Suggests stronger disaster awareness, reporting, or tagging in the U.S.



New York Talks, But Not About Disasters

- □ 56 Normal Tweets vs. only 17 Disaster Tweets
- High general activity, but fewer tweets tagged as disasters



Disaster-Only Locations Appear

- ¶ Nigeria, India, UK, Washington DC, Mumbai
- These regions only pop up during disasters, not in normal tweeting



Global South: Loud in Crisis

- Countries like India & Nigeria show up in disasters
- May reflect higher alertness or public concern during emergencies



Urban West = Normal Tweet Hotspots

- **March London, NYC, LA, Canada** lead normal tweeting
- Suggests **routine Twitter use is higher** in urban, developed regions



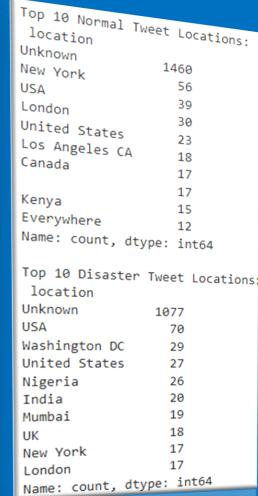
Washington DC: All About Disaster

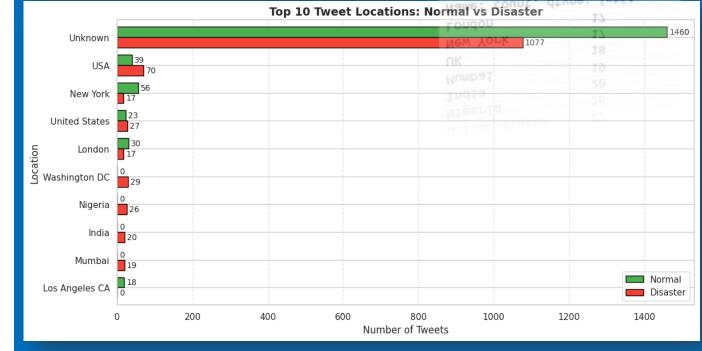
- 29 Disaster Tweets | 0 Normal Tweets
- Possibly due to news media, alerts, or official comms

Tweet Location Hotspot Analysis # Ensure 'Unknown' is properly assigned to missing or blank locations tweet_data['location'] = tweet_data['location'].fillna('Unknown') # Get the Top 10 locations for each category normal_tweets = tweet_data[tweet_data["target"] == 0]["location"].value_counts().nlargest(10) disaster_tweets = tweet_data[tweet_data["target"] == 1]["location"].value_counts().nlargest(10) # Combine results into a DataFrame & fill missing values locations = pd.DataFrame({"Normal": normal_tweets, "Disaster": disaster_tweets}).fillna(0) # Keep only the top 10 locations based on total tweet count locations["Total"] = locations["Normal"] + locations["Disaster"] locations = locations.nlargest(10, "Total").drop(columns=["Total"]) fig, ax = plt.subplots(figsize=(12, 6)) locations.plot(kind='barh', ax=ax, width=0.75, color=['#4CAF50', '#F44336'], edgecolor='black' ax.set title("Top 10 Tweet Locations: Normal vs Disaster", fontsize=14, fontweight='bold' ax.set xlabel("Number of Tweets", fontsize=12) ax.set_ylabel("Location", fontsize=12) ax.invert vaxis() # Ensure the highest value is at the top ax.grid(axis='x', linestyle='--', alpha=0.5) # Add data labels for container in ax.containers: ax.bar_label(container, fmt='%.0f', label_type='edge', fontsize=10, padding=3) plt.show(# Displaying the top 10 locations print("\nTop 10 Normal Tweet Locations:\n", normal_tweets) print("\nTop 10 Disaster Tweet Locations:\n", disaster_tweets) 38 # Clean location data in the main DataFrame df['location'] = df['location'].astype(str).str.strip() 40 df.loc[df['location'] == '', 'location'] = 'Unknown' # Replace blanks with 'Unknown df['location'] = df['location'].str.replace(r'[^\w\s]', '', regex=True) 43 # Create a copy of the relevant columns to avoid modification warnings tweet data = df[['location', 'target']].copy()

46 # Visualize hotspots

47 visualize hotspots(tweet data)







Splitting the Dataset into Training & Testing Sets for Model **Development & Evaluation**

```
💢 Splitting the Dataset into Training & Testing Sets for Model Development & Evaluation.
```

```
# Defining features and target
# clean text: used for TF-IDF to capture semantic meaning
# sentiment: polarity score from TextBlob, useful for understanding tone
X = df[['clean_text', 'sentiment']]
v = df['target']
# Printing the variable names
pd.set option('display.max columns', 5)
print("X contains:\n")
X.head()
```

→ x contains:

	clean_text	sentiment
0	deed reason earthquake may allah forgive u	0.0
1	forest fire near la ronge sask canada	0.1
2	resident asked shelter place notified officer	-0.1
3	people receive wildfire evacuation order calif	0.0
4	got sent photo ruby alaska smoke wildfire pour	0.0

```
# Importing library for splitting dataset
      from sklearn.model_selection import train_test_split
      # Splitting dataset
      X train, X test, y train, y test = train test split(X, y, test size=0.2, random state=42)
      # Printing the shapes of the resulting datasets
      print("Training set size:", X train.shape[0])
      print("Testing set size:", X test.shape[0])
Training set size: 6090
Testing set size: 1523
```

Feature Setup Features (X) include: Cleaned Tweets ✓ Sentiment Scores from TextBlob Together, they form a rich dataset of 7613 entries ♣ Sentiment adds emotional context to raw text! **Target Variable** y holds the binary labels: \square 1 \rightarrow Disaster Tweet \triangle \square 0 \rightarrow Normal Tweet A clean and simple setup for supervised learning! Train-Test Split Summary Used an 80/20 ratio for training and testing: Training set: 6090 entries Testing set: 1523 entries Ensures the model learns well and is tested fairly **Bottom Line** ✓ The dataset is now model-ready, complete with smart Features & clear labels

accurate evaluation 2

✓ Split cleanly to enable robust model development &

TF-IDF & Sentiment-Based Feature Engineering & Class Balancing





Cleaned tweets transformed into numerical vectors using TfidfVectorizer

Feature cap set to **5000** to maintain balance Matrix sizes:

 \square Training \rightarrow (6090, 5000)

 \square Testing \rightarrow (1523, 5000)

Printed sample features like 'aa', 'ab', 'aba' for transparency

Vectorizer saved in saved models/ for reuse



Sentiment Added as a Feature

Used **TextBlob polarity** to score tweet tone

Converted sentiment column into a sparse matrix Merged with TF-IDF vectors

Gives tweets an emotional dimension — not just text!





- ☐ A powerful combination of semantic (TF-IDF) and emotional (sentiment)
- features enhances the model's understanding of tweet content.

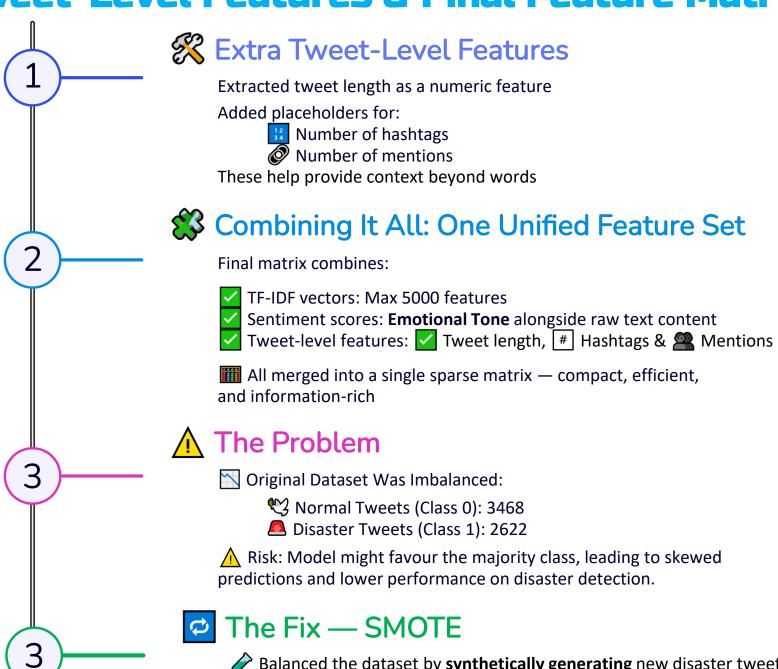


```
1 # Import necessary libraries
       2 from sklearn.feature_extraction.text import TfidfVectorizer
         from scipy.sparse import csr matrix, hstack
       5 # Convert text data into numerical features using TF-IDF
         vectorizer = TfidfVectorizer(max features=5000)
      7  X train tfidf = vectorizer.fit_transform(X train['clean text'])
       8  X test tfidf = vectorizer.transform(X test['clean text'])
Combining TF-IDF Features with Sentiment Analysis for Enhanced Text Representation
      1 # Import necessary libraries
       2 from scipy.sparse import csr_matrix, hstack
       4 # Convert sentiment column to sparse matrix format
       5 # Cast sentiment values to float before reshaping
         extra train = csr matrix(X train['sentiment'].astype(float).values.reshape(-1, 1))
         extra_test = csr_matrix(X_test['sentiment'].astype(float).values.reshape(-1, 1))
         # Combine TF-IDF features with sentiment features
      10 X train combined = hstack([X train tfidf, extra train])
     11 X test combined = hstack([X test tfidf, extra test])
     13 # Print the shape of the combined feature matrices
     14 print("Shape of X_train_combined:", X_train_combined.shape)
     15 print("Shape of X test_combined:", X test_combined.shape)
     17 # Printing a few TF-IDF feature names for transparency
      18 print("\nTF-IDF Feature Names:", vectorizer.get feature names out()[:7]) # Displaying first 7 feature names
→ Shape of X train combined: (6090, 5001)
    Shape of X_test_combined: (1523, 5001)
              import pickle
              import os
              # Create directory if it doesn't exist
              os.makedirs("saved_models", exist_ok=True)
              # Save the vectorizer and combined features
              with open("saved_models/vectorizer.pkl", "wb") as f:
                   pickle.dump(vectorizer, f)
              # Print confirmation
             print("Vectorizer saved successfully in 'saved models/vectorizer.pkl'")
```

Vectorizer saved successfully in 'saved models/vectorizer.pkl'

```
Experimenting with additional features like tweet length, presence of hashtags, or user mentions.
              # Function to extract additional tweet-level features
              def extract_features(df_subset):
                   # Initialize a DataFrame to hold the features
                   # Calculate the length of each tweet's cleaned text
                   features['tweet_length'] = df_subset['clean_text'].apply(len)
                    # Placeholder for number of hashtags
           10
                    features['num_hashtags'] = 0
           11
           12
                    # Placeholder for mentions
           13
                    features['has_mention'] = 0
           14
           15
                    # Return features as a sparse matrix
           16
                    return csr_matrix(features.values)
           17
                # Extract features for training and testing datasets
                 tweet_feats_train = extract_features(X_train)
                 tweet_feats_test = extract_features(X_test)
            21
                 # Combine various feature sets into a single sparse matrix for training and testing datasets
                X_train_combined = hstack([X_train_tfidf, extra_train, tweet_feats_train])
                 X_test_combined = hstack([X_test_tfidf, extra_test, tweet_feats_test])
            25
                 # Print the shape of the combined feature matrices
                 print("Combined training features shape:", X_train_combined.shape)
                 print("Combined testing features shape:", X_test_combined.shape)
         Combined training features shape: (6090, 5004)
          Combined testing features shape: (1523, 5004)
         Combined testing features shape: (1523, 5004)
    Combined training features shape: (6898, 5884)
```

Tweet-Level Features & Final Feature Matrix



Balanced the dataset by **synthetically generating** new disaster tweet samples to match the majority class.

Post-SMOTE Class Counts:

- Class 0 (Normal): 3468
- Class 1 (Disaster): 3468

This ensures **fairer learning** and reduces model bias toward the majority class.

Selecting the Right Model

01 **STEP**

Models Explored for Training:

To identify the best fit, I tested these three powerful classifiers:

✓ Logistic Regression

Random Forest

Support Vector Machine (SVM)

All models were initialized with basic parameters & verified for correct setup.

STEP

Q Logistic Regression:

- ✓ Set max iter =1000 to ensure smooth convergence
- Extremely efficient with TF-IDF and added features
- ✓ Quick, accurate, and reliable

03 STEP

🛕 Random Forest

- ✓ Great at handling complex feature sets
- Should perform decently but is a bit resource-intensive & slower
- Best used when interpretability isn't a top priority

04 STEP

✓ SVM

- Known for strong performance in text classification
- Required more time and hyperparameter tuning
- Less practical here due to slower training time

Task: Model Selection and Training

```
# Import necessary libraries for model training and evaluation
from sklearn.linear_model import LogisticRegression
 from sklearn.ensemble import RandomForestClassifier
  from sklearn.svm import SVC
  # Define models to evaluate
       "Logistic Regression": LogisticRegression(max_iter=1000),
   models = {
       "Random Forest": RandomForestClassifier(),
       "SVM": SVC()
10
11
    # Print the models to evaluate
12
    print("Models to evaluate:")
    for model_name, model in models.items():
         print(f"- {model_name}: {model}")
```

- Logistic Regression: LogisticRegression(max_iter=1000) Models to evaluate: - Random Forest: RandomForestClassifier()
- SVM: SVC()

Model Performance – Metrics at a Glance

Model Performance & Visual Analysis

```
1 # Set up the figure size for the plots
   plt.figure(figsize=(18, 12))
    # Initialize a list to store metrics data for each model
    metrics_data = []
   # Variables to track the best F1 score and the corresponding model
   best model = None
   # Iterate through each model defined in the models dictionary
    for idx, (name, model) in enumerate(models.items()):
       # Fit the model on the combined training data
       model.fit(X_train_combined, y_train)
       # Make predictions on the combined test data
       y_pred = model.predict(X_test_combined)
       # Determine predicted probabilities based on model capabilities
       if hasattr(model, "predict_proba"):
           # Use predicted probabilities for the positive class if available
           y_proba = model.predict_proba(X_test_combined)[:, 1]
        elif hasattr(model, "decision_function")
           # Use decision function output if available
           y_proba = model.decision_function(X_test_combined)
           # Default to predicted class labels if neither is available
           v proba = v pred
        # Calculate evaluation metrics for the model
        from sklearn.metrics import (accuracy_score, precision_score, recall_score,
                                fl score, classification report, confusion matrix,
                                roc_curve, auc, precision_recall_curve)
        acc = accuracy score(y test, y pred) # Accuracy
        prec = precision_score(y_test, y_pred) # Precision
        rec = recall_score(y_test, y_pred) # Recall
        f1 = f1_score(y_test, y_pred) # F1 Score
       # Store metrics in the list for later analysis
       metrics_data.append({"Model": name, "Accuracy": acc, "Precision": prec, "Recall": rec, "F1 Score": f1})
        # Print classification report for the current model
       print(f"\n{name} Classification Report:")
```

```
Logistic Regression Classification Report:
            precision recall f1-score support
                                            874
                                            649
                0.77
                         0.76
                                  0.76
   accuracy
                                   0.80
                                            1523
                 0.79
                         0.79
                                   0.79
                                            1523
  macro avg
weighted avg
                0.80
                         0.80
                                            1523
Random Forest Classification Report:
            precision recall f1-score support
                 0.77
                         0.85
                                   0.81
                                             874
                 0.77
                         0.66
                                  0.71
                                            649
   accuracy
                                   0.77
                                            1523
                 0.77
                         0.76
                                   0.76
                                            1523
  macro avg
                0.77
                         0.77
                                   0.77
                                            1523
weighted avg
SVM Classification Report:
            precision recall f1-score support
                0.71
                0.49
                         0.77
                                   0.60
                                            649
                                   0.57
                                            1523
                 0.60
                         0.59
                                   0.56
                                            1523
  macro avg
weighted avg
```

```
# Update best F1 score and model if the current F1 score is higher
        if f1 > best_f1:
            best f1 = f1
            best_model = model
        # --- Confusion Matrix ---
        # Create a subplot for the confusion matrix
        plt.subplot(3, len(models), idx + 1)
        cm = confusion_matrix(y_test, y_pred) # Compute confusion matrix
        sns.heatmap(cm, annot=True, fmt='d', cmap='Blues') # Plot confusion matrix as a heatmap
        plt.title(f'{name}\nConfusion Matrix') # Title for the subplot
        plt.xlabel('Predicted') # X-axis label
        plt.ylabel('Actual') # Y-axis label
        # --- ROC Curve ---
61
        # Create a subplot for the ROC curve
        plt.subplot(3, len(models), idx + 1 + len(models))
        fpr, tpr, _ = roc_curve(y_test, y_proba) # Compute ROC curve
        plt.plot(fpr, tpr, label=f'AUC = {auc(fpr, tpr):.2f}') # Plot ROC curve with AUC
        plt.plot([0, 1], [0, 1], 'k--', alpha=0.7) # Diagonal line for reference
        plt.title(f'{name}\nROC Curve') # Title for the subplot
        plt.xlabel('False Positive Rate') # X-axis label
        plt.ylabel('True Positive Rate') # Y-axis label
        plt.legend() # Show legend
        # --- Precision-Recall Curve ---
        # Create a subplot for the Precision-Recall curve
        plt.subplot(3, len(models), idx + 1 + 2 * len(models))
        precision, recall, _ = precision_recall_curve(y_test, y_proba) # Compute precision-recall curve
        plt.plot(recall, precision) # Plot Precision-Recall curve
        plt.title(f'{name}\nPrecision-Recall') # Title for the subplot
        plt.xlabel('Recall') # X-axis label
        plt.vlabel('Precision') # Y-axis label
81 # Adjust layout for better spacing and display the plots
82 plt.tight layout()
83 plt.show()
```

01

1 Logistic Regression Leads with Consistency

✓ Accuracy: ~80%
✓ F1 Score: ~0.80
✓ ROC AUC: ~0.86

✓ Precision-Recall: High & Stable

02

2 Random Forest Delivers Competitive Results

✓ Accuracy: ~77%
✓ F1 Score: ~0.77
✓ ROC AUC: ~0.84

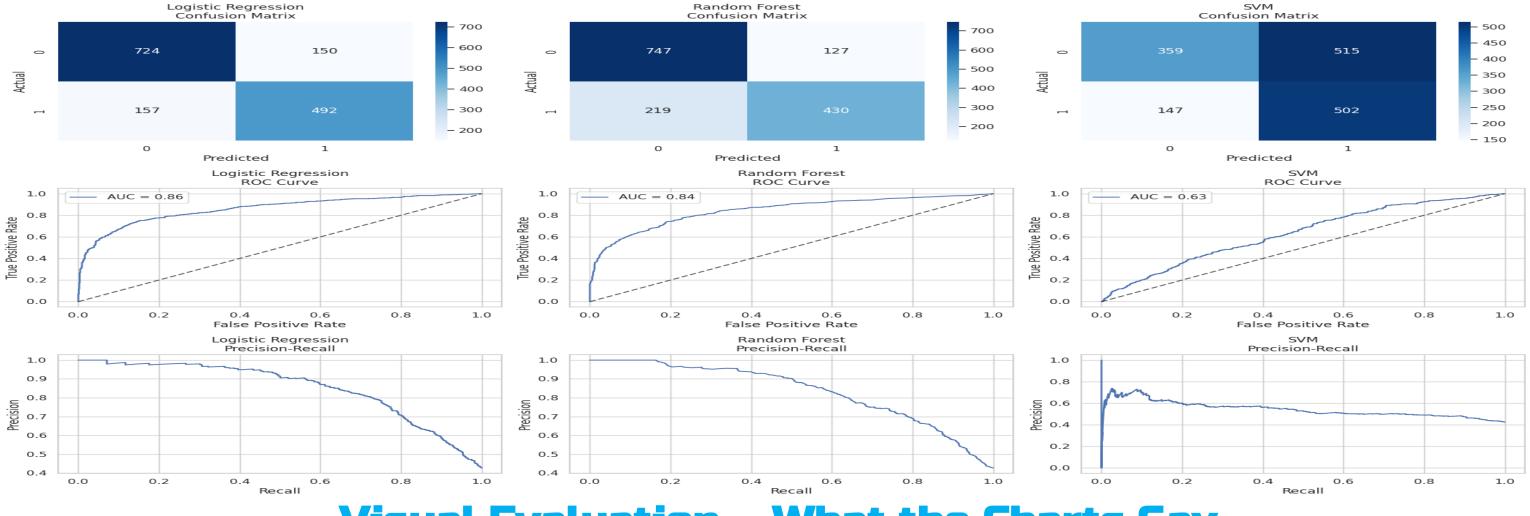
↑ Slightly more misclassifications for disaster tweets

03

3 SVM Falls Short for This Problem

Accuracy: ~57%F1 Score: ~0.56ROC AUC: ~0.63

X Precision-Recall curve drops sharply



Visual Evaluation – What the Charts Say

- **✓** Confusion Matrices:
- ☐ Logistic Regression:
 Strong diagonals, better
 class distinction
- ☐ SVM: Lighter matrix, more errors & confusion

- **ROC Curves:**
- ☐ Logistic Regression:

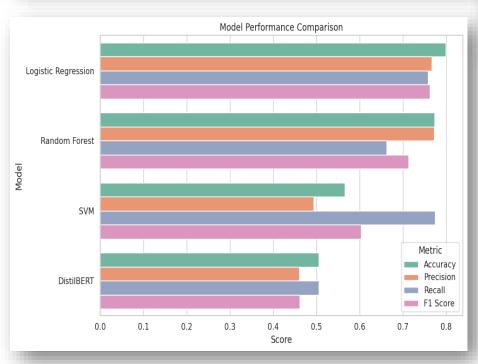
 Tops with AUC ~0.86
- ☐ Random Forest:

 Close second with AUC ~0.84
- □ SVM:
 Lowest, AUC ~0.63

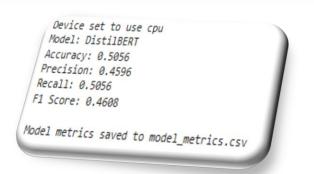
- Precision-Recall Curves:
- ☐ Logistic Regression:
 High precision across recalls
- ☐ Random Forest:
 Steady, but dips slightly
- ☐ SVM:
 Sharp drop, lower overall precision

Model Comparison & Best Model Determination & Auto-Saving

```
🧠 Evaluating Pretrained DistilBERT Model from Hugging Face
      1 # Import necessary libraries for Hugging Face transformer model
          from huggingface_hub import login
          from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
         # Authenticate with your Hugging Face token
          login(token="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")
         # Print message indicating model evaluation start
        print("Evaluating DistilBERT model from Hugging Face..."
         classifier = pipeline("text-classification", model="distilbert-base-uncased-finetuned-sst-2-english")
         # Use the classifier to make predictions on the test set
         predictions = classifier(X test['clean text'].tolist())
         # Extract predicted labels and convert them to a list
         predicted labels = [pred['label'] for pred in predictions]
         predicted labels = [1 if label == "POSITIVE" else 0 for label in predicted labels] # FIX
          acc = accuracy_score(y_test, predicted_labels)
         prec = precision score(v test, predicted labels, average='weighted')
          rec = recall score(y test, predicted labels, average='weighted')
          f1 = f1_score(y_test, predicted_labels, average='weighted')
         # Print the evaluation metrics
         print(f"Model: DistilBERT")
          print(f"Accuracy: {acc:.4f}"
         print(f"Precision: {prec:.4f}"
          print(f"Recall: {rec:.4f}")
         print(f"F1 Score: {f1:.4f}\n")
          metrics_data.append({"Model": "DistilBERT", "Accuracy": acc, "Precision": prec, "Recall": rec, "F1 Score": f1})
         # Export model evaluation results to CSV
         metrics df = pd.DataFrame(metrics data)
        # Export the DataFrame to a CSV file
        metrics_df.to_csv("model_metrics.csv", index=False)
     44 # Print confirmation message
     45 print("Model metrics saved to model_metrics.csv"
```







■ Logistic Regression model saved to: saved_models/Logistic_Regression.pkl





I compared my tailored models (Logistic Regression, Random Forest, SVM) against the pretrained DistilBERT from Hugging Face.

Custom Model F1 Score: ~0.80

DistilBERT F1 Score: ~0.46 (no fine-tuning)

Visuals (Confusion Matrix, ROC, Precision-Recall) clearly favoured the custom setup \bigcirc

02

Y Best Model Determination & Saving **□**

ightharpoonup The Custom Model emerged as the top performer ightharpoonup

■ Logged performance metrics in model_metrics.csv

Saved the best model in the saved_models directory using pickle

03

Bottom Line

A domain-tuned model outperformed a generic transformer

The best model is stored & good to go for deployment



Real-Time Tweet Classification Performance

01

Live Testing with User-Entered Tweets

Ran the model on fresh, unseen tweets to check real-world performance.

Tweet: "There was an earthquake in Myanmar" \rightarrow Predicted as Disaster Tweet with 84.05% confidence.

Tweet: "it's an honour to be invited in President Trump's Birthday" \rightarrow Predicted as Non-Disaster Tweet with 73.86% confidence.

02

Boosted Accuracy with Rich Features

I enriched the classification by combining TF-IDF vectors with smart features like:

Tweet length

Sentiment score

Hashtag count

Mentions (@) presence

03

Results

Fast, consistent, and reliable predictions

Acted as a mini deployment simulation, showing strong performance outside the test set

The model confidently handles real-world inputs — it's not just trained, it's ready to **serve** insights on live data!

```
Part 4: Real-Time Tweet Classification
       1 print("\n=== Real-Time Tweet Classification ===")
           with open("saved models/vectorizer.pkl", "rb") as f:
              loaded_vectorizer = pickle.load(f)
           loaded_model = joblib.load("/content/saved_models/Logistic_Regression.pkl")
           def preprocess tweet(tweet):
              # Remove URLs
              tweet = re.sub(r"http\S+", "", tweet)
              # Expand contractions (custom function assumed defined elsewhere)
              tweet = decontracted(tweet)
              # Remove words with digits
              tweet = re.sub(r"\S*\d\S*", "", tweet).strip()
              # Remove non-alphabetic characters
              tweet = re.sub(r'[^A-Za-z]+', ' ', tweet)
              # Clean special characters from hashtags and underscores
              tweet = tweet.replace("#", "").replace("_", " ")
              tokens = word tokenize(tweet)
               # Lemmatize and remove stop words
               tokens = [lemmatizer.lemmatize(word.lower()) for word in tokens if word.lower() not in stop words]
              clean = ' '.join(tokens).strip()
              return clean
               user_input = input("\nEnter a tweet to classify (or 'exit' to stop): ")
              if user_input.lower() == "exit":
              clean = preprocess_tweet(user_input)
              # Calculate sentiment polarity using TextBlob
               sentiment = TextBlob(clean).sentiment.polarity
              # Vectorize cleaned tweet text
               tfidf_input = loaded_vectorizer.transform([clean])
               tweet_len = len(clean)
                                                      # Length of cleaned tweet
               num hashtags = user input.count("#") # Count hashtags in original tweet
               has_mention = int("@" in user_input) # Check if tweet contains mention
              # Create sparse matrix for extra features
               extra feat = csr_matrix([[sentiment, tweet_len, num_hashtags, has_mention]])
               # Combine TF-IDF features with extra features
               final_input = hstack([tfidf_input, extra_feat])
               pred = loaded_model.predict(final_input)[0]
               # Get prediction probabilities and extract confidence for predicted class
               proba = loaded model.predict proba(final input)[0]
               confidence = proba[pred]
               label = " ■ Disaster Tweet 6 " if pred == 1 else " V Non-Disaster Tweet 8 "
               # Print prediction with confidence score as percentage
               print(f"Prediction: {label} with confidence score: {confidence * 100:.2f}%")
```

```
=== Real-Time Tweet Classification ===

Enter a tweet to classify (or 'exit' to stop): There was an earthquake in Myanmar.

Prediction: ▲ Disaster Tweet ♠ with confidence score: 88.31%

Enter a tweet to classify (or 'exit' to stop): it's an honour to be invited in President Trump's Birthday.

Prediction: ✔ Non-Disaster Tweet ఄ with confidence score: 68.61%

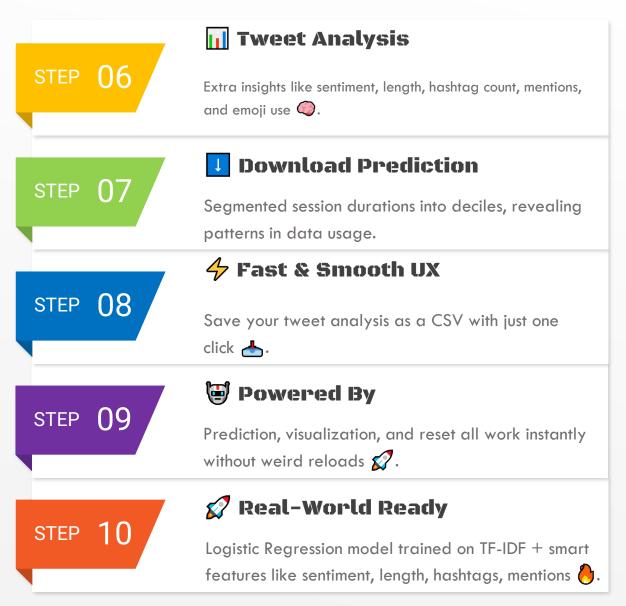
Enter a tweet to classify (or 'exit' to stop): exit
```



Disaster Tweet Detector — Streamlit App Overview

Live App: <u>Disaster Tweet Detector</u> All About My Streamlit App: <u>Click to Watch (CTRL + Click)</u>

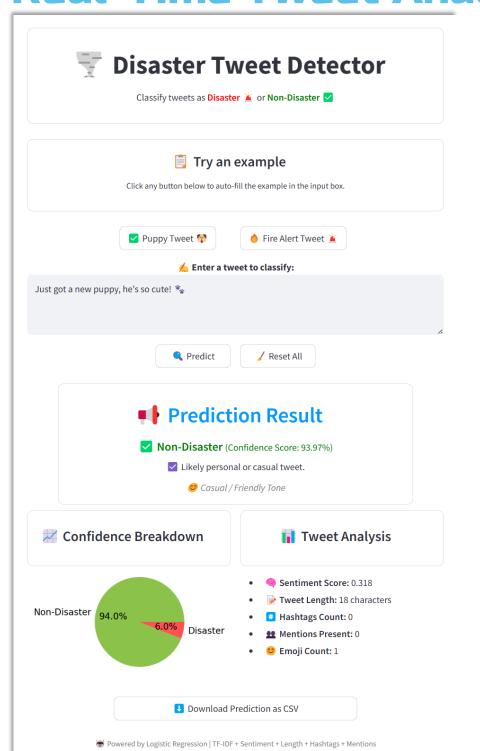






Sleek & User-Friendly Streamlit Interface for Real-Time Tweet Analysis





O1 Centralized Layout

Every button, example, & output is thoughtfully center-aligned for a clean, professional look

Instant Example Loading

Example tweets load immediately into the input box without refreshing the page.

OConsistent Design

Matching boxed sections for prediction, analysis & pie charts keep the UI smooth & intuitive.

04

Data Insertion

Used a MERGE statement to insert or update user scores based on the existing records in the table.



I appreciate your time and attention & hope this presentation was informative and helpful.

For any questions or further assistance, please feel free to reach out.