

SEC 162 GENERATIVE AI

Submitted by

V.Dinesh AP23110011594

Y. Dharani AP23110010030

K.VanajaAP23110011661

G.Geethika AP23110011104

V.Siva Nikshya AP23110010285

SEM: V

SECTION: Y

BACHELOR OF TECHNOLOGY

in

COMPUTER SCIENCE AND ENGINEERING

of

SCHOOL OF ENGINEERING AND SCIENCES



SRM University-AP, Neerukonda, Andhra Pradesh 522240

December 2025

ClassiNews Pro Solution

Project Description:

The **ClassiNews Pro: Advanced News Categorization System** is an advanced machine learning project designed to efficiently **classify, rank, and match** news articles against professional domains (**World, Sports, Business, Science/Tech**). It utilizes Natural Language Processing (NLP) techniques and supervised learning algorithms to automatically screen and categorize large volumes of news, significantly reducing manual organization time and improving information management efficiency.

This intelligent system leverages **TF-IDF vectorization** with trigrams, and an **Ensemble Model** combining **Multinomial Naive Bayes** and **Logistic Regression** for robust and accurate prediction. By analyzing article content, extracting key features, and comparing them against known domain categories, the model assists analysts, editors, and researchers in identifying relevant information quickly. The system provides category classification and **Confidence Scores** (ranking) based on prediction probability.

Project Scenarios

The ClassiNews Pro application, accessible via the interactive Gradio web interface, directly addresses the described scenarios:

Scenario 1: High-Volume Article Processing (Batch Classification)

A large media monitoring firm needs to rapidly process a feed of thousands of articles to direct them to the correct editorial departments (e.g., World news goes to the Global Desk, financial news goes to the Business Desk).

ClassiNews Pro Solution:

The system's **Batch Processing** capability allows users to input multiple articles simultaneously. The model processes all submissions within minutes, automatically classifying and returning detailed results, including predicted category and confidence scores, reducing manual sorting time significantly.

Article Example	Text Snippet	Predicted Category
Article 4	Stock market reaches all-time high amid economic recovery...	Business
Article 5	United Nations discusses climate change policies at global summit...	World
Article 8	New AI technology promises to revolutionize healthcare industry...	Business (The model correctly identifies the commercial/economic focus of the AI application)

Scenario 2: Research Agency Operations (Matching and Comparison)

A market research agency manages multiple reports simultaneously, needing to match and compare two different articles to ensure they cover the same domain or to contrast their content.

ClassiNews Pro Solution:

The system's **Compare Articles** feature handles this matching. It automatically classifies both pieces of text and determines if they belong to the same category, offering clear differentiation between domains.

Input	Text Snippet	Predicted Category	Matching Result
Article 2	The football team won the championship after a thrilling match...	Sports	Articles belong to different categories: Sports vs Business
Article 7	Microsoft acquires gaming company in multi-billion dollar deal...	Business	

Scenario 3: Continuous Content Pipeline Management (History Tracking)

A data journalism team requires continuous tracking and logging of all articles classified over time to ensure compliance and maintain an auditable database of prediction outcomes.

Prerequisites:

To complete this project, you must require the following software, concepts, and packages

Software Requirements

The software requirements for the **classinews** are primarily related to the development environment, python libraries, and application deployment tools.

Based on the typical architecture of such a project (as detailed in the provided document snippets), here are the key software and library requirements:

Software Requirements

1. Development Environment

1. Python: Version **3.8 or higher** is required to ensure compatibility with all modern libraries.

2. Integrated Development Environment (IDE):

3. Google Colab or Jupyter Notebook for interactive development, data exploration, and model training.

4. Anaconda Navigator (Recommended) for managing Python environments and packages.

- Refer to installation guide: <https://youtu.be/1ra4zH2G4o0>

Python Packages

Install the following packages using pip in Anaconda Prompt (run as administrator):

- `pip install numpy` - Numerical computing library
- `pip install pandas` - Data manipulation and analysis
- `pip install scikit-learn` - Machine learning algorithms and utilities
- `pip install nltk` - Natural Language Toolkit for text processing
- `pip install matplotlib` - Data visualization library
- `pip install seaborn` - Statistical data visualization
- `pip install flask` - Web framework for deployment
- `pip install kagglehub` - Kaggle dataset downloading utility
- `pip install gradio` - Interface for machine learning models
- `pip install joblib` - Model persistence and serialization

Prior Knowledge Requirements

To successfully complete this project, you should have understanding of:

Machine Learning Concepts

- **Supervised Learning:** <https://www.javatpoint.com/supervised-machine-learning>
- **Unsupervised Learning:** <https://www.javatpoint.com/unsupervised-machine-learning>
- **Logistic Regression:** <https://www.javatpoint.com/logistic-regression-in-machine-learning>
- **Decision Tree:** <https://www.geeksforgeeks.org/python-decision-tree-regression-using-sklearn/>
- **Random Forest:** <https://www.javatpoint.com/machine-learning-random-forest-algorithm>
- **Evaluation Metrics:** <https://www.analyticsvidhya.com/blog/2019/08/11-important-model-evaluation-error-metrics/>
- **Naive Bayes Classification:** <https://www.javatpoint.com/machine-learning-naive-bayes-classifier>
- **TF-IDF Vectorization:** Text feature extraction technique for machine learning
- **Cosine Similarity:** Similarity measurement between document vectors

Natural Language Processing

- Text preprocessing and cleaning techniques
- Tokenization and stopword removal
- Feature extraction from unstructured text data

Web Framework Knowledge

- **Flask Basics:** https://www.youtube.com/watch?v=lj4I_CvBnto
- **Gradio Interface:** Interactive model deployment

Project Flow

1. **Resume Dataset Collection** - Gather resume data from Kaggle or other sources
2. **Data Preprocessing** - Clean and normalize resume text
 - Remove special characters and formatting
 - Tokenize text into words
 - Remove stopwords
 - Convert to lowercase
3. **Feature Extraction** - Convert text to numerical features
 - Apply TF-IDF vectorization to extract key terms and their importance
 - Generate feature vectors for each resume
4. **Model Training** - Train classification and ranking models
 - Split data into training and testing sets
 - Train Logistic Regression classifier for category prediction
 - Store TF-IDF vectorizer for consistent feature transformation
5. **Resume Ranking** - Score resumes against job descriptions
 - Calculate cosine similarity between resume and job description vectors
 - Rank resumes by compatibility score
 - Classify resumes by professional category
6. **Results Display** - Present ranked candidates to recruitment team
 - Display top-matching resumes
 - Show confidence scores and relevance percentages
 - Provide actionable insights for hiring decisions

Project Activities

Activity 1: Data Collection & Preparation

- **Gather raw data** (e.g., from Kaggle).
- **Clean and normalize** all text data.
- **Address data quality** (missing/duplicates).
- **Balance data** across professional categori.

Activity 2: Exploratory Data Analysis

- **Calculate text statistics** (length, word count).
- **Analyze category distribution** and balance.
- **Identify skill frequency** (key terms).
- **Generate visualizations** (e.g., Word Clouds).

Activity 3: Data Preprocessing

- **Normalize text** (lower-case, remove noise).
- **Perform Tokenization** (using NLTK).
- **Remove stopwords** for better signal.
- **Execute feature engineering** for text data.

Activity 4: Model Building

- Train **TF-IDF Vectorizer** for numerical conversion.
- Train the **Logistic Regression classifier**.
- Develop the **Cosine Similarity** ranking function.
- Create the final resume **Scoring Mechanism**.

Activity 5: Model Evaluation

- Check **Accuracy, Precision, and F1-Score**.
- Test ranking performance with sample job descriptions.
- Validate fit scores for matching reliability.
- Compare model performance configurations.

Activity 6: Model Deployment

- Save trained models (Vectorizers/Classifiers).
- Develop the **Gradio/Flask web interface**.
- Launch application for recruitment team us

Project Structure

```
resume-screening/
├── data/
│   └── resumes/ # Raw resume dataset
├── notebooks/
│   └── resume_screening.ipynb # Jupyter notebook with full pipeline
├── models/
│   ├── log_reg_model.pkl # Trained Logistic Regression model
│   └── tfidf_vectorizer.pkl # Fitted TF-IDF vectorizer
├── app.py # Flask web application
└── requirements.txt # Python dependencies
    └── templates/
        ├── index.html # Landing page
        └── results.html # Results display
    └── static/
        └── style.css # Application styling
└── README.md # Project documentation
```

Project Structure Explanation

Step 1: Data Acquisition and Preparation

1. The process starts by gathering and cleaning the input data.

2. Collect Dataset: Acquire the raw resume dataset, which is categorized by professional roles, typically from sources like **Kaggle**.

3. Initial Cleaning: Load the data (using **pandas**), perform preliminary data cleaning, and **normalize** text (e.g., combining title and description fields).

Ensure Data Quality: Handle **missing values** and remove **duplicate entries** to guarantee a clean training pool.

Balance Categories: Verify that the data is well-distributed across all desired professional categories to prevent model bias toward a single role.

Step 2: Exploratory Data Analysis (EDA)

Analyze the dataset to understand its characteristics.

1. Analyze Text Stats: Calculate and review key statistics like **text length** and **word count** for all resumes.

2. Examine Distribution: Analyze the distribution of resumes across the target **categories**.

3. Identify Key Terms: Perform **skill frequency analysis** and generate visualizations (like **Word Clouds**) to highlight important words and phrases.

Step 3: Data Preprocessing and Feature Engineering

Transform the raw text into a numerical format the model can understand.

1. Text Normalization: Apply the preprocessing function (`preprocess_text`) to make text **lower-case** and remove **punctuation** and noise (e.g., using **regular expressions**).

2. Tokenization & Filtering: Break down the cleaned text into individual words or tokens, and remove non-informative **stopwords** using **NLTK**.

3. Vectorization: Implement and train the **TF-IDF Vectorizer** (Term Frequency-Inverse Document Frequency). This converts the text corpus into a numerical matrix, effectively creating **features** that represent the importance of each word.

Step 4: Model Training and Scoring Logic

This stage builds the intelligent core of the system.

1.Train Classifier: Train the primary classification model, which is a **Logistic Regression classifier** (often combined into an **Ensemble Model** for better performance, as seen in the code log).

2.Develop Ranking Logic: Develop a **Cosine Similarity function** to compute the degree of match between the job description vector and the resume vector.

3.Create Scoring Mechanism: Integrate the model's **classification confidence score** (the probability of the resume belonging to the target category) with the **cosine similarity score** to produce a final, comprehensive resume fit score.

Step 5: Model Evaluation and Validation

Measure the model's performance and ensure reliability.

1. Assess Performance: Use the test dataset to calculate key metrics: **Accuracy, Precision, Recall, and F1-Score** (e.g., achieving **89.70% Accuracy**).

2. Test Ranking: Run the system against sample job descriptions to test the validity of the generated ranking/fit scores.

3. Validate Results: Review the performance dashboard and classification reports to confirm the model is robust and reliable across all categories.

Milestone 1: Data Collection & Preparation

This phase focuses on acquiring, cleaning, and structuring the raw resume data.

1.1 Importing Required Libraries: Imports core libraries such as **Pandas** for data handling, **NLTK** for natural language processing, **Scikit-learn** for machine learning components, and **KaggleHub** for data access.

1.2 Download NLTK Resources: Executes commands like `nltk.download('stopwords')` to prepare necessary lexical resources for text cleaning.

1.3 Read and Load the Dataset: Uses **KaggleHub** to download and `pd.read_csv` to load the raw resume data into a DataFrame (df).

1.4 Data Preparation Overview: Outlines the quality plan: handling missing values, removing duplicates, and text normalization.

1.5 Handling Missing Values: Checks for nulls using `df.isnull().sum()` and fills critical missing resume content with an **empty string** (`df['resume'].fillna("")`).

1.6 Handling Duplicates: Identifies and removes redundant resume entries using `df.drop_duplicates(subset=['resume'], keep='first')`.

1.7 Filtering Valid Data: Enforces a **minimum length** constraint (`df['resume'].str.len() > min_length`) to discard incomplete or meaningless resume entries.

Milestone 2: Exploratory Data Analysis (EDA)

This phase analyzes the characteristics and distribution of the cleaned data.

2.1 Descriptive Statistics: Calculates total resumes, unique professional categories, and the value_counts() for category distribution.

2.2 Visual Analysis - Category Distribution: Calculates resume length (df['resume_length']) and plots a **Bar Plot** to check for potential class imbalance.

2.3 Resume Length Distribution: Generates a **Histogram** and a **Box Plot** to visualize the typical range and spread of resume lengths across different categories.

2.4 Skill Frequency Analysis: Uses Python's Counter on all resume words to identify and visualize the **Top 20 Most Common Words (Skills)**.

Milestone 3: Text Preprocessing & Feature Extraction

This critical phase transforms text into numerical features.

3.1 Text Preprocessing Function: Defines preprocess_text to handle: **Lowercase** conversion, **Tokenization**, **Stopword Removal**, and filtering.

3.2 Apply Preprocessing: Executes the function using df["resume"].apply(preprocess_text) to create the cleaned text column (clean_resume).

3.3 TF-IDF Vectorization: Initializes the TfidfVectorizer (e.g., max_features=1500, ngram_range=(1, 2)). The vectorizer is fitted to the cleaned text, creating the sparse numerical feature matrix **X**.

Milestone 4: Model Building & Training

This phase trains and saves the core intelligence.

4.1 Train-Test Split: Divides the feature matrix **X** and labels **y** into **80% training** and **20% testing** subsets using train_test_split, employing stratify=y to maintain category balance.

4.2 Logistic Regression Model: Initializes and trains the primary classification model (LogisticRegression).

4.3 Model Evaluation: Calculates **Accuracy** and prints the **Classification Report** and **Confusion Matrix** using the test set to assess model performance.

4.4 Testing Additional Algorithms: Performs a brief comparison against other classifiers (e.g., Random Forest, SVM) to confirm the model choice.

4.5 Save Trained Models: Serializes the final trained model (log_reg_model) and the tfidf_vectorizer using joblib.dump for persistence and deployment.

Milestone 5: Resume Ranking & Matching Functions

This phase integrates the model with real-world matching logic.

5.1 Fit Score Computation: Defines compute_fit_score, which calculates the **Cosine Similarity** (sklearn.metrics.pairwise.cosine_similarity) between a resume's vector and a job description's vector, returning a 0-100% fit score.

5.2 Resume Ranking Function: Defines rank_resumes, which scores a batch of resumes using the fit score, predicts the professional category, and **sorts the entire list** for prioritization.

5.3 Live Testing: Executes a test with sample job descriptions and resumes to demonstrate the **real-time ranking capability**.

Milestone 6: Model Deployment

The final phase makes the system operational.

6.1 Flask Web Application (API): Sets up the backend using **Flask**, loads the saved models at startup, and defines a **RESTful API endpoint** (/rank) for processing requests.

6.2 Gradio Interface (Frontend): Creates an alternative, simpler web interface using **Gradio** (gr.Interface) for rapid prototyping and demonstration.

6.3 System Architecture: Defines the system structure as a **three-tier architecture**: Frontend (Gradio/HTML), Backend (Flask/Model Inference), and Data Layer (Model Persistence).

6.4 Deployment Workflow: Outlines the request flow: Load Models \rightarrow Preprocess Input \rightarrow TF-IDF Transformation \rightarrow Calculate Similarity \rightarrow Rank Results \rightarrow Display.

II. Comparative Analysis: Resume Screening vs. ClassiNews Pro

The provided log details the **ClassiNews Pro** implementation (AG News dataset, 89.70% accuracy). This comparison highlights the different technical decisions made for each project, despite following similar ML milestones.

Feature	Automated Resume Screening System (Plan)	ClassiNews Pro (Implementation Log)
Primary Goal	Ranking & Matching resumes to a specific job description.	Classification of articles into 4 distinct categories.
Dataset	Resume data (imbalance is common).	AG News data (perfectly balanced: 30,000 per category).
Model	Logistic Regression (Primary Classifier).	Voting Ensemble of MultinomialNB + Logistic Regression .

Key Takeaways for Departmental Code Review:

1. Ensemble Strength: The ClassiNews Pro project demonstrates a more robust modeling approach using an **Ensemble Model** (MultinomialNB + Logistic Regression) for classification. This technique should be considered for the Resume Screening System to potentially boost classification accuracy above the single Logistic Regression model proposed.

2. Deployment Flexibility: The Gradio interface design in ClassiNews, with its dedicated tabs for **Batch Processing** and **Comparison**, provides an excellent template for the Resume Screening System's user interface, enhancing usability for HR staff.

Activity 1.2: Importing Required Libraries

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from nltk.corpus import stopwords
from sklearn.feature_extraction.
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from sklearn.metrics.pairwise import cosine_similarity
```

Download NLTK resource

```
nltk.download('punkt_tab')
nltk.download('stopwords')
```

Library Purposes:

- **NumPy** - Numerical operations and array manipulation
- **Pandas** - Data loading, manipulation, and analysis
- **Matplotlib & Seaborn** - Data visualization and exploration
- **NLTK** - Natural language processing and text tokenization

Activity 1.3: Read and Load the Dataset

```
import kagglehub
import os
```

Download dataset

```
print("🚀 Downloading dataset...")
path = kagglehub.dataset_download("palaksood97/resume-dataset")
print("✅ Dataset root path:", path)
```

Load resume data

```
import glob
csv_files = glob.glob(os.path.join(path, "*.csv"), recursive=True)
if csv_files:
    df = pd.read_csv(csv_files[0])
    print(f" Loaded dataset: {csv_files[0]}")
else:
    print("Dataset loading requires proper Kaggle authentication")
```

Display basic information

```
print(df.head())
print(df.shape)
print(df.columns)
```

The `head()` function displays the first 5 rows of the dataset, providing immediate understanding of the data structure and content format. The dataset typically contains columns for resume text and professional category.

Activity 1.4: Data Preparation

Before using the data to train machine learning models, thorough data preparation is essential. This includes:

- **Handling Missing Values** - Identify and manage null or incomplete entries
- **Removing Duplicates** - Eliminate duplicate resume entries
- **Checking Types** - Verify appropriate mistake any line

Code Action	Code Reference
Data Acquisition	<code>kagglehub.dataset_download(file_path)</code>
DataFrame Initialization	<code>df = pd.read_csv('resume_data.csv')</code>

1.5 Handling Missing Values (Data Imputation)

Missing values (`NaN`) in the text column are a major source of runtime errors in NLP pipelines. They must be handled explicitly to ensure pipeline continuity.

Python

```
# Check for nulls: df.isnull().sum()# Impute missing values:
df['resume'].fillna("", inplace=True)
```

Professional Rationale:

Substituting null values with an empty string ("") is the safest form of imputation for text data. It ensures that the subsequent string-based functions (e.g., tokenization, lowercasing) execute without failure.

Resumes with no content are effectively treated as zero-length documents, which the TF-IDF vectorizer can process without crashing, thereby preventing pipeline failure.

1.6 Handling Duplicates (Data Integrity)

Duplicate resume submissions must be eliminated to prevent **model bias** and **data leakage**, which artificially inflate performance metrics.

Python

```
# Check for duplicates: df.duplicated(subset=['resume']).sum()# Remove duplicates:
df.drop_duplicates(subset=['resume'], keep='first', inplace=True)
```

Professional Rationale:

The use of `subset=['resume']` is crucial as it enforces the uniqueness check based solely on the content of the resume, ignoring potentially non-unique row IDs or submission timestamps. This action ensures that the model is trained on a set of unique candidates, promoting realistic performance evaluation and ensuring that the learned patterns generalize well to unseen data.

1.7 Filtering Valid Data (Quality Assurance)

This step enforces a minimum quality threshold, removing truncated or non-meaningful records that would otherwise introduce noise into the model's learning process.

Python

```
# Define minimum character length
min_length = 50
# Apply filter to retain only high-signal resumes
df = df[df['resume'].str.len() > min_length]
```

Professional Rationale:

By calculating the string length using `df['resume'].str.len()` and applying a minimum filter, we perform Data Quality Assurance. This ensures that the model learns from high-signal, meaningful text relevant to skills and experience. Excluding text below the threshold minimizes noise and concentrates the model's learning capacity on valid input, which is key for accurate TF-IDF vectorization

The Machine Learning Process



1. Application Architecture & Key Components

The ClassiNews Pro system follows a classic **Machine Learning Pipeline Architecture** for Natural Language Processing (NLP), deployed via a user-friendly **Gradio Web Interface**.

Component	Technology	Role
Data Source	AG News Dataset (Kaggle)	Provides 120,000 training and 7,600 test samples across 4 balanced categories.
Data Processing	Python, Pandas, NLTK, Regex	Handles text cleaning, tokenization, stop word removal (implicitly via TF-IDF), and category mapping.

Component	Technology	Role
Feature Extraction	TF-IDF Vectorizer (with <code>\$\text{ngram_range}=(1, 2)\$</code>)	Transforms text into a numerical feature matrix, capturing the importance of unigrams and bigrams (a key enhancement).
Classification Model	Multinomial Naive Bayes (MNB)	A probabilistic classifier highly effective for text classification with sparse features like TF-IDF vectors.
Deployment/UI	Gradio	Creates the interactive, multi-tab web application for real-time and batch predictions.
Persistence	joblib, pickle	Saves the trained model (<code>news_classifier_model.pkl</code>) and the vectorizer (<code>tfidf_vectorizer.pkl</code>) for future use and deployment.

Key Architectural Enhancement: The use of **bigrams** (via `$ngram_range=(1, 2)$` in `TfidfVectorizer`) significantly enriches the feature set, allowing the model to capture two-word phrases and improve context understanding, which is crucial for distinguishing similar categories (e.g., "market crash" in Business vs. "market analysis" in Sci/Tech).

2. Deployment Workflow

The application is designed for both local execution and sharing via Gradio's built-in sharing feature.

1.Preparation: Required libraries are installed and imported. NLTK stopwords and tokenizers are downloaded.

Model Training (Offline):

1.Dataset is loaded from Kaggle.

2.Text is preprocessed and cleaned.

3.The `TfidfVectorizer` is fitted on the training data, capturing the vocabulary and feature weights.

4.The **Multinomial Naive Bayes** model is trained on the TF-IDF features and category labels.

5.Model Persistence: The trained **model** and the essential **vectorizer** are saved to disk (.pkl files). This is vital for production, allowing the system to load the model without retraining every time.

Application Startup (Online):

1.The Gradio interface (gr.Blocks) is defined with multiple tabs for different user tasks (Quick Classify, Batch, Compare)

Launch: The `demo.launch(share=True)` command makes the application accessible via a publicly shareable link (if successful) or a local URL, completing the deployment cycle.

3. Classification Accuracy and Evaluation

The model's performance is evaluated on the dedicated, unseen **Test Set** (7,600 articles) to ensure generalization.

A. Performance Summary

Metric	Value (Test Set)	Interpretation
Accuracy	\$\text{92.49}\%\$	$\frac{\text{Correct Predictions}}{\text{Total Predictions}}$. The model correctly categorizes news articles over 92% of the time.
Precision	\$\text{92.51}\%\$	$\frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$. Out of all articles predicted to be a category, 92.51% were actually that category.
Recall	\$\text{92.49}\%\$	$\frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$. Out of all articles that <i>should</i> have been a category, 92.49% were correctly identified.
F1-Score	\$\text{92.49}\%\$	Harmonic mean of Precision and Recall. An excellent single score showing robust performance.

B. Confusion Matrix Analysis

The Confusion Matrix (Figure 7: Confusion Matrix Heatmap) would visually confirm that:

High True Positives: The diagonal elements (correct classifications) are significantly larger than off-diagonal elements (errors).

Low Confusion: Misclassification rates between categories are low. For example, confusion might be slightly higher between **World** and **Business** compared to **Sports** and **Science/Tech**, but overall separation is strong due to the balanced, distinct nature of the categories.

4. Ranking Effectiveness & Scalability Improvement

A. Ranking Effectiveness (Feature Importance)

The system includes a **Feature Importance Analysis** (Step 8), a critical step for interpretability and model debugging.

1. Methodology: The log-probabilities of features ($\log P(\text{feature} | \text{category})$) in the Multinomial Naive Bayes model are used to rank the most important words/bigrams for each category.

2. Effectiveness: This ranking validates the model's learning process:

3. Sports: Features like 'game', 'win', 'team', 'nba finals' would rank highly.

4. Business: Features like 'stock', 'market', 'shares', 'firm', 'billion dollar' would rank highly.

This confirms the model is making classifications based on **contextually relevant keywords**, ensuring the predictions are not arbitrary.

B. Scalability Improvements

Computational Efficiency: Multinomial Naive Bayes is known for its **fast training and prediction** times, even on large datasets. Its low computational overhead makes it highly scalable for classifying massive streams of news articles.

Vectorization Optimization: The `TfidfVectorizer` is limited to a controlled vocabulary of **10,000 features** and filters out very rare features (`min_df=2`) and very common ones (`max_df=0.95`). This prevents the feature space from exploding, keeping both memory usage and prediction latency low.

1. Model Persistence: Saving the model and vectorizer eliminates the need to reprocess \$120,000\$ training articles every time the application starts, dramatically improving deployment speed and operational scalability.

2. Batch Processing: The Gradio interface directly supports **Batch Processing**, allowing the user to classify hundreds of articles simultaneously with a single vectorization and prediction call,

5. Conclusion :

The **ClassiNews Pro** project successfully delivers an intelligent, high-performing, and feature-rich news article categorization system.

Advanced Features: The integration of advanced visualizations (Word Clouds, Confusion Matrix), interpretability (Feature Importance), and user-centric features (Batch Processing, History, Comparison) via Gradio elevates the system from a simple script to a **production-ready prototype**.

Robust Architecture: The use of saved models (.pkl files) ensures rapid deployment and scalability, making it suitable for integrating into larger news aggregation or content management platforms.