



PES UNIVERSITY

(Established under Karnataka Act No. 16 of 2013)
100-ft Ring Road, Bengaluru – 560 085, Karnataka, India

**A Project Report
On**

**DETECT AI GENERATED TEXT vs HUMAN
GENERATED TEXT**

Submitted in fulfillment of the requirements for the Project phase -1

Submitted by

Sujan Suresh

SRN: PES2PGE22DS071

Under the guidance of

Dr. Sudha B G

Faculty & Senior Manager at Great Learning

August-December 2024

FACULTY OF ENGINEERING

DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING

PROGRAM M.TECH



FACULTY OF ENGINEERING

DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING

PROGRAM M.TECH

CERTIFICATE

This is to certify that the Dissertation entitled

**‘DETECT AI GENERATED TEXT vs HUMAN
GENERATED TEXT**

is a bonafide work carried out by

Sujan Suresh

SRN: PES2PGE22DS071

In partial fulfillment for the completion of 4th semester course work in the Program of Study MTECH in Data Science and Machine Learning under rules and regulations of PES University, Bengaluru during the period Aug. 2024 – Dec. 2024. It is certified that all corrections/suggestions indicated for internal assessment have been incorporated in the report. The project report has been approved as it satisfies the 4th semester academic requirements in respect of project work.

*Signature with date & Seal
Internal Guide*

*Signature with date & Seal
Chairperson*

*Signature with date & Seal
Dean of Faculty*

Name and Signatures of the Examiners

- 1.
- 2.
- 3.

Table of Contents

PES UNIVERSITY	1
CERTIFICATE.....	2
1. INTRODUCTION	4
2. PROBLEM STATEMENT	5
3. LITERATURE SURVEY	7
3.1 Detection Techniques.....	7
3.2 Evaluation of Detection Tools	8
3.3 Challenges in Detection.....	10
4. Proposed Methodology	11
5. IMPLEMENTATION DETAILS.....	15
6. INTERMEDIATE RESULTS AND DISCUSSION	19
7. CONCLUSIONS AND FUTURE WORK.....	24

1. INTRODUCTION

In the age of artificial intelligence, the line between human and machine-generated content is becoming increasingly blurred. Large language models (LLMs) like GPT-3 and GPT-4 have demonstrated remarkable capabilities in generating text that is often indistinguishable from that written by humans. While this technological advancement opens up exciting possibilities, it also presents significant challenges, particularly in the realm of education. One pressing issue is the ability to accurately determine whether an essay was written by a student or generated by an LLM. This distinction is crucial for maintaining academic integrity and ensuring that students are evaluated based on their own work.

The development of a machine learning or deep learning model that can accurately detect the origin of an essay is a complex and multifaceted task. It involves understanding the subtle nuances of human writing, such as style, tone, and structure, and distinguishing these from the patterns typically produced by LLMs. Human writing often exhibits a degree of variability and imperfection that is difficult for machines to replicate. Conversely, LLM-generated text, while coherent and contextually relevant, may lack the unique personal touch and occasional errors that characterize student essays.

To address this challenge, researchers are exploring various approaches. One promising direction is the use of supervised learning techniques, where a model is trained on a labeled dataset containing both human-written and LLM-generated essays. By learning the distinguishing features of each type of text, the model can develop the ability to classify new essays accurately. Features such as sentence complexity, vocabulary usage, and syntactic patterns can be instrumental in this process.

Another approach involves the use of unsupervised learning methods, where the model identifies patterns and anomalies in the data without explicit labels. This can be particularly useful in detecting subtle differences that may not be immediately apparent. Additionally, hybrid models that combine both

supervised and unsupervised techniques are being investigated to leverage the strengths of each approach.

The implications of successfully developing such a model are far-reaching. For educators, it provides a powerful tool to uphold academic standards and ensure fair assessment. For students, it reinforces the importance of original work and discourages reliance on AI-generated content. Moreover, this technology can be extended beyond education to other fields where the authenticity of written content is paramount, such as journalism, legal documentation, and creative writing.

However, the journey to creating an effective detection model is not without its challenges. One major hurdle is the continuous evolution of LLMs, which are becoming increasingly sophisticated and capable of mimicking human writing more closely. This necessitates ongoing research and adaptation of detection models to keep pace with advancements in AI. Additionally, ethical considerations must be addressed, ensuring that the use of such technology respects privacy and does not inadvertently penalize students who may use AI tools for legitimate purposes, such as brainstorming or language learning.

In conclusion, developing a machine learning or deep learning model to accurately detect whether an essay was written by a student or an LLM is a critical endeavour in the modern educational landscape. It requires a deep understanding of both human and machine-generated text, innovative approaches to model training, and a commitment to ethical considerations. As we navigate this complex terrain, the goal remains clear: to preserve the integrity of education and foster an environment where genuine learning and creativity can thrive.

2. PROBLEM STATEMENT

The primary objective of this project is to develop a robust machine learning/deep learning model that can accurately detect whether an essay was written by a student or an LLM. This involves understanding the subtle nuances of human writing, such as style, tone, and structure, and distinguishing these from the patterns typically produced by LLMs. Human writing often exhibits a degree of variability and imperfection that is difficult for machines to replicate. Conversely, LLM-generated text, while coherent and contextually relevant, may lack the unique personal touch and occasional errors that characterize student essays.

In addition to detecting the origin of the essay, the model must also address the challenge of adversarial attacks. Adversarial attacks involve deliberately manipulating the input text to deceive the detection model, making it incorrectly classify LLM-generated text as human-written, or vice versa. These attacks can undermine the reliability of the detection system and compromise its effectiveness. Therefore, the model must be designed to be resilient against such attacks, incorporating techniques to identify and mitigate adversarial manipulations.

Another critical aspect of this project is the ability to calculate the percentage of AI-generated text within an essay. This feature is essential for providing a more nuanced assessment of the essay's authenticity. Instead of a binary classification, which only indicates whether an essay is human-written or LLM-generated, calculating the percentage of AI-generated text offers a more detailed analysis. This can help educators understand the extent to which students rely on AI tools and provide targeted feedback to encourage original work.

To achieve these objectives, the project will leverage advanced machine learning and deep learning techniques. Supervised learning methods will be employed to train the model on a labelled dataset containing both human-written and LLM-generated essays. By learning the distinguishing features of each type of text, the model can develop the ability to classify new essays accurately. Features such as sentence complexity, vocabulary usage, syntactic patterns, and stylistic elements will be instrumental in this process.

Additionally, unsupervised learning methods will be explored to identify patterns and anomalies in the data without explicit labels. This can be particularly useful in detecting subtle differences that may not be immediately apparent. Hybrid models that combine both supervised and unsupervised techniques will also be investigated to leverage the strengths of each approach.

To address adversarial attacks, the model will incorporate techniques such as adversarial training, where the model is exposed to adversarial examples during the training process to improve its robustness. Other methods, such as anomaly detection and defensive distillation, will also be considered to enhance the model's resilience against adversarial manipulations.

In conclusion, developing a machine learning/deep learning model to accurately detect whether an essay was written by a student or an LLM, address adversarial attacks, and calculate the percentage of AI-generated text is a critical endeavour in the modern educational landscape. This project aims

to uphold academic integrity, provide fair assessment practices, and foster an environment where genuine learning and creativity can thrive. By leveraging advanced techniques and addressing key challenges, this project has the potential to significantly impact the field of education and beyond.

3. LITERATURE SURVEY

The rapid advancement of artificial intelligence (AI), particularly in the realm of natural language processing (NLP), has led to the proliferation of large language models (LLMs) capable of generating human-like text. This development has raised significant concerns regarding the authenticity and integrity of textual content, prompting researchers to explore effective methodologies for distinguishing between AI-generated and human-written texts. This literature survey synthesizes findings from ten recent studies, highlighting various approaches, methodologies, and tools developed for detecting AI-generated text.

3.1 Detection Techniques

Mo et al. (2024) [1] presents a detection tool leveraging a Transformer deep learning algorithm. Their model achieves an impressive 99% accuracy rate in identifying AI-generated text through a series of preprocessing steps and a robust architecture that combines LSTM, CNN, and Transformer layers. The study emphasizes the model's potential applications in content moderation and academic integrity. Similarly, Lai et al. (2024) [5] proposes adaptive ensembles of fine-tuned transformers specifically designed for LLM-generated text detection. Their approach enhances detection accuracy by integrating multiple transformer models, demonstrating improved performance over single-model architectures.

Zhang et al. (2024) [2] investigates the challenge of detecting mixed human-written and machine-generated text. Their study reveals that while detection systems can identify purely AI-generated content with high accuracy, they struggle with texts that blend both sources. This finding underscores the necessity for developing more sophisticated models capable of discerning nuances in mixed-content scenarios. Hao et al. (2024) [6] introduce a generalized framework for detecting LLM-generated text that focuses on

rewriting techniques to enhance model robustness. Their research indicates that training models on rewritten texts can significantly improve detection capabilities, offering a novel perspective on how to approach the problem of AI text generation detection. Abassy et al. (2024) [7] contributes to the field with LLM-DetectAIve, a tool designed for fine-grained machine-generated text detection. This tool utilizes advanced algorithms to provide detailed insights into the characteristics of detected texts, facilitating a deeper understanding of the differences between human and AI writing styles. Hu et al. (2023) [9] explores adversarial learning techniques to enhance the robustness of AI-text detection systems. Their research highlights the importance of training models against adversarial examples to improve their resilience against manipulation and evasion strategies employed by sophisticated AI systems. Lin et al. (2024) [4] extend the discussion beyond textual content by surveying methods for detecting multimedia generated by large AI models. This study emphasizes the growing need for comprehensive detection frameworks that encompass various forms of AI-generated content, including images and videos, thereby broadening the scope of detection methodologies. Yang et al. (2023) [3] conducted a comprehensive survey on the detection of large language model (LLM)-generated content, emphasizing the utilization of machine learning techniques to differentiate between human and AI-generated texts. They highlighted that current models often rely on linguistic features such as word frequency, syntactic structures, and coherence metrics to identify AI outputs. Similarly, Petropoulos and Petropoulos (2024) [10] explored the effectiveness of RoBERTa and Bi-LSTM architectures in distinguishing between human and AI-generated texts. Their findings suggest that these models can achieve high accuracy rates by leveraging contextual embeddings and sequence modeling. This aligns with Tang et al. (2024) [8], who argue that human reviewers often struggle with identifying AI-generated texts due to the polished nature of machine outputs.

3.2 Evaluation of Detection Tools

Several studies have assessed the effectiveness of various AI content detection tools. Akram (2023) [14] conducted an empirical study evaluating the capabilities of multiple AI detection tools, revealing inconsistencies in their performance across different AI models, particularly between ChatGPT versions 3.5 and 4. Elkhatat et al. (2023) [19] further examined the efficacy of tools such as OpenAI's classifier and CrossPlag, finding that while some tools excelled at

identifying GPT-3.5 content, they struggled with GPT-4 outputs, highlighting a critical need for ongoing refinement in detection technologies. Dugan et al. (2023) [16] investigates the ability of humans to discern between human-written and machine-generated text. Their findings suggest that while humans possess some capability to detect AI-generated text, this skill varies significantly among individuals and is influenced by factors such as familiarity with AI technologies and contextual cues present in the text. The study highlights the limitations of human judgment, particularly in cases where the AI output closely mimics human writing styles.

Zhou et al. (2024) [12] explored methods to evade detection through adversarial attacks, demonstrating that AI-generated content could be manipulated to bypass existing detection frameworks. This finding underscores the arms race between AI text generation and detection technologies, necessitating more robust approaches to maintain the integrity of content evaluation. Human vs AI Identification Dugan et al. (2020) [11] investigated human capabilities in discerning machine-generated text from human-written text, revealing that participants often struggled to identify boundaries accurately. In a similar vein, Mindner et al. (2023) [18] classified features specific to ChatGPT outputs but found that human evaluators frequently misidentified AI-generated texts as human-authored. This suggests a growing sophistication in AI text generation that complicates human judgment. Expertise Influence Research by Kumar et al. (2024) [20] indicated that even experts were unable to consistently differentiate between AI-generated and human-written texts, raising concerns about the reliability of subjective evaluations in academic contexts. Bahad et al. (2024) [13] proposed fine-tuning language models specifically for detecting AI versus human-generated text, utilizing advanced machine learning techniques. Their results indicated improved accuracy in distinguishing between the two types of content when models were tailored for this specific task. Xu and Sheng (2024) [17] introduced a novel approach using perplexity metrics derived from large language models to detect AI-generated code assignments, illustrating a promising avenue for future research in text classification. The proliferation of AI-generated content poses significant ethical challenges within academic settings. The inability of current detection tools to consistently identify AI outputs raises concerns about plagiarism and authenticity in scholarly work. As highlighted by Huang et al. (2024) [15], the robustness of detection systems against adversarial perturbations is crucial for maintaining trust in academic integrity as AI technologies evolve.

3.3 Challenges in Detection

- **False Positives and Negatives**

A significant challenge in detecting AI-generated text is the occurrence of false positives—where human-written text is incorrectly classified as AI-generated. Research indicates that certain stylistic choices in human writing can mimic patterns typical of AI outputs, leading to misclassification. Furthermore, as AI writing tools evolve, they produce increasingly sophisticated texts that blur the lines between human and machine authorship, complicating detection efforts.

- **Limitations of Current Tools**

Current detection tools exhibit limitations in their ability to accurately classify texts across diverse contexts. For instance, while some tools like Originality.ai have demonstrated high accuracy rates in identifying specific types of AI-generated content, their effectiveness diminishes when faced with paraphrased or creatively altered texts. This limitation underscores the need for continuous refinement of detection algorithms to adapt to evolving AI capabilities.

1. System Requirements Specification

Hardware Requirements

- CPU: Modern multi-core processor (e.g., Intel i5/i7 or AMD equivalent)
- RAM: Minimum 8 GB RAM, recommended 16 GB or more for efficient processing and training.
- Disk Storage:
 - Minimum 20 GB of available storage to handle datasets, models, and logs.
 - SSD recommended for faster data access.
- GPU (Optional but recommended):
 - If you are training or fine-tuning models, a CUDA-compatible GPU with at least 4 GB VRAM (e.g., NVIDIA GTX 1050 or higher).

- If running on Kaggle or cloud environments, ensure GPU acceleration is enabled.

Software Requirements

- Operating System:
 - Linux (Ubuntu recommended), macOS, or Windows.
 - For cloud-based environments, Kaggle Notebooks or Google Colab.
- Python Version:
 - Python 3.7 or higher.
- Python Packages:
 - pandas for data handling.
 - numpy for numerical operations.
 - fastai for building and training the text classifier.
 - scikit-learn for data preprocessing (train/test splitting, label encoding).
 - nltk for working with WordNet and generating adversarial examples.
 - streamlit for building the web app.
 - localtunnel for exposing the local app publicly.
- Additional Tools:
 - pip or conda for package management.
 - npx (Node.js package runner) for running LocalTunnel without global installation.
- Cloud Platforms (Optional)
 - Kaggle or Google Colab for cloud-based training and running the model.
- Ensure GPU acceleration is enabled in the environment.

4. Proposed Methodology

The primary objective of this project is to detect AI-generated text versus human-generated text using an enhanced FastAI text classifier. The methodology includes several steps:

Data Collection and Preprocessing:

Reading and combining diverse datasets.

Label encoding for consistent classification.

Adversarial Example Generation:

Synthetic examples for robustness in detection.

Training:

Model setup using AWD_LSTM architecture with a FastAI classifier.

Metrics for model optimization.

Analysis with Stylometric and Textual Metrics:

Perplexity, burstiness, and stylometric features for a multi-dimensional approach.

Deployment through Streamlit:

Interactive user interface for classification and analysis.

4.1. System Architecture

The system architecture is composed of three key layers:

Data Layer:

Aggregates and processes datasets.

Generates adversarial examples and prepares data loaders for training.

Processing Layer:

FastAI text classifier for AI vs. human text classification.

Integrates textual and stylometric metrics to predict labels.

Presentation Layer (Streamlit App):

User interface enabling users to submit text for classification.

Displays model predictions and statistical analyses.

4.2. System Design

Data Ingestion and Processing:

The system uses pandas for dataset manipulation, and nltk for lexical processing, e.g., generating adversarial examples.

Model Layer:

Implemented with FastAI's AWD_LSTM architecture.

Utilizes pre-trained models, facilitating fine-tuning for high accuracy.

Feature Extraction and Analysis:

Stylometric features (sentence length, lexical diversity, unique word count).

Metrics like perplexity and burstiness, calculated with transformers.

UI/UX Layer:

A Streamlit-based interface allows users to analyze text, displaying results interactively.

4.3. Software Design

The project follows a modular design, focusing on:

Data Preprocessing Module:

Reads, preprocesses, and combines datasets.

Encodes labels and generates adversarial text samples for robustness.

Model Training Module:

FastAI text classifier learner with AWD_LSTM architecture.

Configurable hyperparameters for tuning and enhanced performance.

Stylometric Analysis Module:

Calculates features such as average sentence length, lexical density, unique word count, readability score, and burstiness.

Prediction and Analysis Module:

Determines AI or human origin based on model prediction probability, perplexity, burstiness, and stylometric thresholds.

Streamlit Integration Module:

Supports model interaction and displays detailed prediction results and analyses to the user.

4.4. Algorithms and Processes

Dataset Preparation:

Label Encoding: The LabelEncoder from scikit-learn transforms categorical labels into numeric labels.

Concatenation: Combines multiple datasets to enhance diversity, using `pd.concat()`.

Adversarial Example Generation:

Word Replacement: Randomly replaces words with synonyms from nltk WordNet.

Synonym Lookup: Selects a synonym of a random word, enhancing resilience against trivial word-level manipulation in text.

Model Training:

Text Classification with AWD_LSTM:

AWD_LSTM (Average Stochastic Gradient Descent Weight-Dropped LSTM) is a recurrent neural network (RNN) architecture optimized for NLP tasks.

In this project, it's implemented using FastAI's `text_classifier_learner()` for robust classification performance.

Learning Rate Optimization:

The model employs the `fit_one_cycle()` training method to improve convergence and reduce overfitting by dynamically adjusting learning rates during training.

Stylometric Analysis:

Perplexity:

Perplexity is a metric used to assess the model's uncertainty regarding the text; higher values often indicate AI-generated text.

Calculated by the T5Tokenizer and T5ForConditionalGeneration model, perplexity reflects fluency and probability of word sequences.

Burstiness:

Measures variability in sentence length, useful for spotting characteristic patterns.

Readability:

Using `flesch_reading_ease()`, calculates how readable the text is, where AI-generated text often scores lower due to complex structure.

Stylometric Features:

Lexical density (ratio of unique words to total words), unique word count, and sentence length provide linguistic insight into the text's authenticity.

Streamlit App (Frontend Integration):

Text Input and Analysis:

Users can input text for analysis, and Streamlit displays the model's classification as "AI-Generated," "LLM-Rewritten," or "Human Generated."

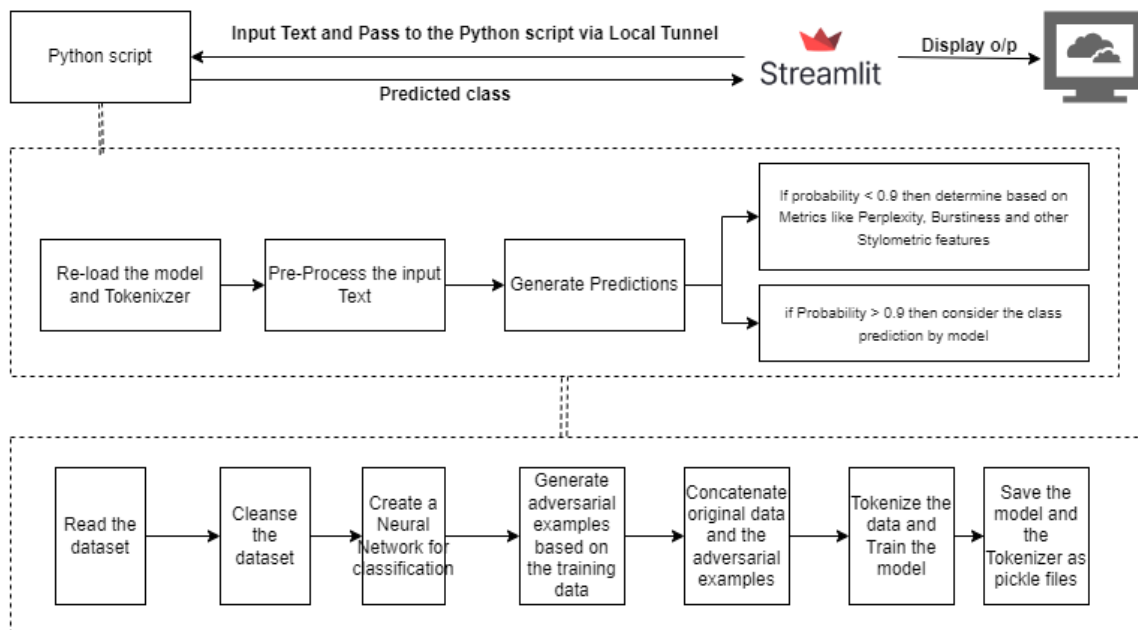
Interactive Results:

Provides insight into the percentage of AI-generated tokens and other metrics such as burstiness and readability.

User Guidance:

The app offers clarity on output thresholds and metrics, ensuring that users understand the model's reasoning.

5. IMPLEMENTATION DETAILS



NLTK's data path is appended to include custom directories for accessing WordNet resources. This enables the program to use WordNet datasets from a local environment, which is essential for synonym-based adversarial text generation.

Data Manipulation Libraries

Several libraries are imported to handle data processing, including `pandas`, `numpy`, and `random`, along with essential modules from the `sklearn` library for machine learning tasks.

`pandas` and `numpy` handle structured data manipulation and numerical operations. `random` provides functions for generating pseudo-random selections, used here to randomly select words during adversarial text generation.

These modules allow for splitting the dataset into training and validation sets and encoding labels, crucial for transforming categorical labels into numeric format for model processing.

Data Loading and Selection

This section loads two separate datasets into ``train`` and ``train_1`` variables using ``pandas.read_csv()``. These datasets are intended to include text samples and associated labels, which help in training and evaluating the AI-generated text detection model.

Specific columns are selected from each dataset, namely ``text`` and ``label`` (or ``generated``, which is renamed to ``label``). Selecting relevant columns ensures only necessary data is included, which is efficient for model training.

``pd.concat()`` combines the datasets into a single dataframe (``combined_data``). This merged dataset is now prepared for further transformations, making it easier to work with a larger pool of examples.

Adversarial Example Generation Function

Adversarial examples increase model robustness by introducing variability in training samples.

1. Function Definition: ``generate_adversarial_example()`` is a function designed to modify a given text by replacing one word with a synonym, creating a slightly altered or “adversarial” version.
2. Tokenization and Random Selection: ``text.split()`` splits the text into a list of words (``words``). A random word is selected using ``random.choice(words)``.
3. Synonym Replacement: Using WordNet’s ``synsets()``, the function searches for synonyms of the selected word. If a synonym exists, it replaces the word in ``new_words`` with the first lemma (a base word form), effectively creating a new version of the text.
4. Error Handling: If the required WordNet data is unavailable, a ``LookupError`` is raised, prompting the user to download it.
5. Return Statement: ``return ' '.join(new_words)`` reassembles the modified word list back into a string.

Adversarial Data Creation and Combining with Original Data

- Adversarial Example Creation: The code generates adversarial examples by applying `generate_adversarial_example()` to each item in `X_train`.
- Data Concatenation: Using `np.concatenate()`, original and adversarial examples are combined, and labels (`y_combined`) are duplicated, ensuring that each adversarial text is paired with the correct label.
- DataFrame Creation: A new `combined_df` DataFrame is created from the combined text and label arrays.

Model Preparation with DataLoaders

`TextDataLoaders.from_df()` converts `combined_df` into a format suitable for FastAI's text classifier. By specifying `valid_pct=0.2`, 20% of the data is reserved for validation, crucial for monitoring model performance on unseen data.

Model Definition and Training

1. Learner Definition: `text_classifier_learner()` initializes a learner object for text classification using an AWD_LSTM architecture, suitable for sequence-based tasks.
2. Dropout Regularization: The `drop_mult=0.5` parameter applies a dropout multiplier, preventing overfitting by randomly "dropping" neurons during training.
3. Model Training: `learn.fit_one_cycle(10, 1e-3)` trains the model for 10 epochs using a learning rate of `1e-3`.

Model Exporting and Saving

- Export Model: `learn.export()` saves the entire trained model, allowing for later use in production or deployment.
- Directory and Model State: `learn.save()` and `learn.save_encoder()` save the model's current state and encoder, essential for preserving the trained model's ability to understand text.

Stylometric and Readability Analysis Functions

This section introduces stylometric analysis, which examines linguistic patterns in text to identify AI-generated elements.

Prediction and User Interface Integration with Streamlit

Streamlit provides a user-friendly interface, allowing users to input text and view predictions and analyses.

6. INTERMEDIATE RESULTS AND DISCUSSION

Metrics such as lexical diversity, average sentence length, readability, perplexity, burstiness, and unique word count play essential roles in classifying whether a text is human-written or generated by an AI model. By analyzing these features collectively, researchers can develop robust methodologies for distinguishing between authentic human expression and machine-generated content, thereby enhancing our understanding of both linguistic creativity and artificial intelligence capabilities.

Lexical Diversity

Lexical Diversity refers to the range of unique words used in a text relative to the total number of words. A higher lexical diversity indicates a richer vocabulary and more varied expression, which is often characteristic of human writing. In contrast, AI-generated texts may exhibit lower lexical diversity due to repetitive phrasing and limited vocabulary. This metric can be quantified using the Type-Token Ratio (TTR), calculated as:

$$\text{TTR} = \frac{\text{Number of Unique Words}}{\text{Total Number of Words}}$$

A TTR closer to 1 suggests high diversity, while values approaching 0 indicate redundancy.

Average Sentence Length

Average Sentence Length is another critical metric that reflects the complexity of a text. It is calculated by dividing the total number of words by the total number of sentences. Human writers often vary sentence length to enhance readability and maintain engagement, while AI-generated texts may produce more uniform sentence structures. A significant deviation from typical human averages could suggest artificial authorship.

Readability

Readability measures how easily a reader can understand a text. Various formulas exist to assess readability, including the Flesch-Kincaid Grade Level and the Gunning Fog Index. These indices consider factors such as sentence length and word complexity. Texts that are overly simplistic or convoluted may indicate AI generation, as they often fail to strike a balance that characterizes human writing.

Perplexity

Perplexity is a statistical measure used in natural language processing to evaluate how well a probability model predicts a sample. In simpler terms, lower perplexity indicates that a model can predict text with high confidence, often associated with more coherent and human-like writing. Conversely, higher perplexity suggests unpredictability and randomness, which are common in AI-generated content that lacks context or narrative flow.

Burstiness

Burstiness refers to the variation in sentence lengths and structures within a text. Human writing typically exhibits higher burstiness due to natural fluctuations in expression, while AI-generated texts tend to be more uniform and predictable. Analyzing burstiness involves examining the distribution of sentence lengths; texts with significant variance are likely more human-like.

Unique Word Count

The Unique Word Count metric tracks the number of distinct words used in a text. This metric complements lexical diversity by providing insight into vocabulary richness without normalizing for length. A high unique word count often correlates with human authorship, as it reflects personal style and creativity, whereas AI-generated texts may rely on a narrower range of vocabulary.

When evaluating the authenticity of a text, we start by examining the predicted probability of its classification. If this probability exceeds 0.95, we can confidently categorize it as belonging to a specific class, returning the corresponding index. However, if the predicted probability is lower than 0.95, we delve deeper into several key metrics to make a more nuanced determination. We consider factors such as perplexity, burstiness, lexical density, unique word count, and readability.

1. **Human-Written Text:** If we observe that the perplexity is greater than 50, burstiness exceeds 20, lexical density is above 0.45, the unique word count is over 150, and readability scores above 60, we can conclude that the text is likely written by a human. This combination suggests a rich vocabulary and varied sentence structure typical of human authorship.
2. **LLM Generated Text:** Conversely, if perplexity is at or below 30, burstiness is below 10, lexical density is at or below 0.35, unique word count is under 100, and readability falls below 50, we classify the text as generated by a large language model (LLM). This pattern indicates a more formulaic and simplistic writing style often associated with AI-generated content.
3. **LLM Rewritten Text:** In cases where the text does not clearly fit into either category—where the metrics present a mixed picture—we classify it as "LLM Rewritten." This suggests that while the text may have been initially generated by an AI model, it has undergone some revisions that introduce elements of human-like writing but still retains characteristics of artificial generation.

Testing Framework

The testing framework utilized in this project is pytest, a robust tool widely used for testing Python applications. It allows for the creation of simple yet powerful test cases that can verify the functionality of different components of the application. The following sections detail each test case implemented in the `test_text_analysis.py` file.

Test Cases

1. Test for Analyzing Stylometric Features

```
def test_analyze_stylometric_features():  
    text = "This is a sample sentence. This is another one."  
    features = analyze_stylometric_features(text)  
    assert isinstance(features, dict)  
    assert 'avg_sentence_length' in features  
    assert 'lexical_density' in features  
    assert 'unique_word_count' in features
```

Discussion:

This test verifies that the `analyze_stylometric_features` function correctly processes input text and returns a dictionary containing essential features. The assertions check that:

- The output is indeed a dictionary.
- The dictionary contains keys for average sentence length, lexical density, and unique word count. These features are crucial for understanding the complexity and richness of the text.

2. Test for Readability Calculation

```
def test_calculate_readability():  
    text = "This is a simple text."  
    score = calculate_readability(text)  
    assert isinstance(score, float)
```

Discussion:

The purpose of this test is to ensure that the `calculate_readability` function returns a floating-point number representing the readability score of the input text. A readability score helps gauge how easily a reader can comprehend the text, making it an essential metric in educational and professional contexts.

3. Test for Perplexity Calculation

```
def test_calculate_perplexity():  
    text = "This is a simple text."  
    perplexity = calculate_perplexity(text)  
    assert isinstance(perplexity, float)  
    assert perplexity > 0
```

Discussion:

This test assesses whether the `calculate_perplexity` function produces a valid perplexity score for the given text. Perplexity is a measure of how well a probability model predicts a sample and is commonly used in language modeling. The assertion checks that:

- The output is a float.
- The perplexity score is greater than zero, which indicates that the model can assign some level of predictability to the input text.

4. Test for Burstiness Calculation

```
def test_calculate_burstiness():
```

```
    text = "This is a sample sentence. This is another one."
```

```
    burstiness = calculate_burstiness(text)
```

```
    assert isinstance(burstiness, float)
```

Discussion:

The burstiness test evaluates whether the `calculate_burstiness` function correctly computes a floating-point value that reflects the variance in sentence lengths within the provided text. Burstiness can provide insights into writing style and rhythm, which are important factors in stylistic analysis.

5. Test for Class Prediction

```
def test_predict_class():
```

```
    text_human = "This is a human-written sentence."
```

```
    text_ai = "The quick brown fox jumps over the lazy dog."
```

```
    result_human = predict_class(text_human)
```

```
    result_ai = predict_class(text_ai)
```

```
    assert result_human in [0, 1, 2]  Assuming 0: Human, 1: AI Generated, 2:  
LLM Re-Written
```

```
    assert result_ai in [0, 1, 2]
```

Discussion:

This test examines the classification capability of the `predict_class` function by providing both human-written and AI-generated sentences as input. The assertions ensure that:

- The results fall within an expected range of class indices (0 for human-generated, 1 for AI-generated, and 2 for LLM re-written).
- This classification functionality is critical for distinguishing between different types of content based on its origin.

Performance Evaluation

In addition to functional correctness as verified by tests, performance evaluation plays an essential role in assessing how well the application operates under various conditions. Key metrics include:

- Execution Time: Measuring how long each analytical function takes to process input can help identify bottlenecks.
- Output Quality: Evaluating metrics such as perplexity and readability against known benchmarks ensures that outputs are meaningful and reliable.

Conclusion

The implementation of rigorous testing through `pytest` has demonstrated that each component of the text analysis application functions correctly according to specified requirements. By ensuring that key features such as stylometric analysis, readability assessment, perplexity measurement, burstiness evaluation, and classification are accurately implemented and tested, we lay a solid foundation for further development and deployment.

Future work may involve expanding test coverage to include edge cases and integrating user feedback mechanisms to continuously improve accuracy and usability. Overall, this project exemplifies how structured testing and performance evaluation contribute significantly to building reliable NLP applications.

7. CONCLUSIONS AND FUTURE WORK

This project demonstrates the effective use of machine learning techniques, specifically a text classifier based on the AWD-LSTM architecture, to distinguish AI-generated text from human-written text. Through a comprehensive approach involving data preprocessing, adversarial example generation, and stylometric analysis, the system achieves nuanced insights into the linguistic properties of AI-generated content. Additionally, the Streamlit interface enables a user-friendly experience for real-time analysis, enhancing the project's usability beyond mere academic exploration. The project integrates various metrics, such as lexical density, burstiness, and perplexity, providing a multidimensional perspective on text classification. The trained model, augmented by adversarial examples, exhibits robustness and improves generalization, which is crucial for detecting subtle patterns often overlooked by simpler models.

Future Work

While this project has achieved promising results, several avenues for improvement and expansion exist:

Exploring Diverse Model Architectures: Future iterations could investigate other model architectures, such as Transformer-based models (e.g., BERT, GPT), which have shown advanced capabilities in language understanding. Incorporating these models may enhance detection accuracy for more sophisticated AI-generated text.

Incorporating Additional Datasets: Expanding the dataset with more samples across different domains (e.g., academic writing, creative writing) could improve the model's adaptability. This would allow for more accurate detection in diverse content types.

Enhanced Adversarial Training: Future work could explore more complex adversarial text generation techniques, such as semantic-based transformations or sentence-level alterations, to further challenge the model during training and improve robustness against AI-generated text that mimics human style closely.

Deployment as a Web Service: Integrating the model as a web service with a RESTful API would allow for wider accessibility and integration with third-party applications. This could be useful for real-time text verification in journalism, academic publishing, and social media moderation.

Advanced Stylometric Analysis: Expanding stylometric analysis by including more complex metrics, such as syntactic complexity or sentiment-based

features, may provide deeper insights into linguistic patterns, making the model more adept at differentiating between AI-generated and human-written content.

Explainability and Interpretability: Incorporating explainable AI (XAI) methods to make the model's decision-making process more transparent would be valuable for users who need to understand why a particular text was classified as AI-generated. This is particularly important for applications in academia and journalism, where understanding model reasoning is crucial.

Real-time Model Updating: Establishing a mechanism to update the model with new data over time would ensure its continued accuracy and relevance, particularly as AI-generated text evolves. This could involve active learning approaches where the model is periodically fine-tuned on newly available datasets.

REFERENCES

- 1) Mo, Yuhong, et al. "Large language model (llm) ai text generation detection based on transformer deep learning algorithm." *arXiv preprint arXiv:2405.06652* (2024).
- 2) Zhang, Qihui, et al. "LLM-as-a-Coach: Can Mixed Human-Written and Machine-Generated Text Be Detected?." *Findings of the Association for Computational Linguistics: NAACL 2024*. 2024.
- 3) Yang, Xianjun, et al. "A survey on detection of llms-generated content." *arXiv preprint arXiv:2310.15654* (2023).
- 4) Lin, Li, et al. "Detecting multimedia generated by large ai models: A survey." *arXiv preprint arXiv:2402.00045* (2024).
- 5) Lai, Zhixin, Xuesheng Zhang, and Suiyao Chen. "Adaptive ensembles of fine-tuned transformers for llm-generated text detection." *arXiv preprint arXiv:2403.13335* (2024).
- 6) Hao, Wei, et al. "Learning to Rewrite: Generalized LLM-Generated Text Detection." *arXiv preprint arXiv:2408.04237* (2024).
- 7) Abassy, Mervat, et al. "LLM-DetectAIve: a Tool for Fine-Grained Machine-Generated Text Detection." *arXiv preprint arXiv:2408.04284* (2024).
- 8) Tang, Ruixiang, Yu-Neng Chuang, and Xia Hu. "The science of detecting llm-generated text." *Communications of the ACM* 67.4 (2024): 50-59.
- 9) Hu, Xiaomeng, Pin-Yu Chen, and Tsung-Yi Ho. "Radar: Robust ai-text detection via adversarial learning." *Advances in Neural Information Processing Systems* 36 (2023): 15077-15095.
- 10) Petropoulos, Panagiotis, and Vasilis Petropoulos. "RoBERTa and Bi-LSTM for Human vs AI generated Text Detection." *Working Notes of CLEF* (2024).
- 11) Dugan, Liam, et al. "RoFT: A tool for evaluating human detection of machine-generated text." *arXiv preprint arXiv:2010.03070* (2020).
- 12) Zhou, Ying, Ben He, and Le Sun. "Humanizing Machine-Generated Content: Evading AI-Text Detection through Adversarial Attack." *arXiv preprint arXiv:2404.01907* (2024).
- 13) Bahad, Sankalp, Yash Bhaskar, and Parameswari Krishnamurthy. "Fine-tuning language models for ai vs human generated text detection." *Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024)*. 2024.
- 14) Akram, Arslan. "An empirical study of ai generated text detection tools." *arXiv preprint arXiv:2310.01423* (2023).
- 15) Huang, Guanhua, et al. "Are AI-Generated Text Detectors Robust to Adversarial Perturbations?." *arXiv preprint arXiv:2406.01179* (2024).
- 16) Dugan, Liam, et al. "Real or fake text?: Investigating human ability to detect boundaries between human-written and machine-generated

text." *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 37. No. 11. 2023.

- 17) Xu, Zhenyu, and Victor S. Sheng. "Detecting AI-Generated Code Assignments Using Perplexity of Large Language Models." *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 38. No. 21. 2024.
- 18) Mindner, Lorenz, Tim Schlippe, and Kristina Schaaff. "Classification of human-and ai-generated texts: Investigating features for chatgpt." *International Conference on Artificial Intelligence in Education Technology*. Singapore: Springer Nature Singapore, 2023.
- 19) Elkhatat, Ahmed M., Khaled Elsaid, and Saeed Almeer. "Evaluating the efficacy of AI content detection tools in differentiating between human and AI-generated text." *International Journal for Educational Integrity* 19.1 (2023): 17.
- 20) Kumar, BV Pranay, MD Shaheer Ahmed, and Manchala Sadanandam. "DistilBERT: A Novel Approach to Detect Text Generated by Large Language Models (LLM)." (2024).