
Classifying AI LLM vs Human-Written Texts

Bryce Hiraoka
Boston University, CS 506
bhiraoka@bu.edu

1 Introduction

For my 506 project, I chose to tackle the challenge of detecting human-written text versus text generated by large language models (LLMs). I worked with six datasets: five sourced from Kaggle (LLM-Detect-AI-Generated-Text, PaLm-Generated-Essays, Combined-Set, AI-vs-Human-Text, and Human-vs-LLM-Corpus-Bloom-7B-and-GPT) and one from Hugging Face (AI-Text-Detection-Pile). To build the initial model, I combined multiple datasets to cover a diverse range of topics and LLMs, ensuring comprehensive detection capabilities. However, for subsequent models, I focused on a single dataset to achieve more refined and specific results.

Github Link:

1.1 Method & Relevant Results

I chose logistic regression for this project due to its effectiveness in binary classification and its widespread use in similar research problems. Using a single dataset containing 64,000 texts and five distinct features, one of my models achieved an overall accuracy of 80

1.2 Why is AI detection important?

Distinguishing between human-written and LLM-generated text has become increasingly important as AI continues to integrate into various aspects of our daily lives. AI detection is relevant to everyone, from children in schools to the elderly consuming news. One key application of AI detection is protecting readers from misinformation. Since free access to large language models like ChatGPT became available, the number of LLM-generated fake articles has risen by over 1000 percent [1]. With the growing prevalence of AI-generated fake news, the number of people impacted by such misinformation has likely increased as well. Implementing an AI detection model to alert readers about potentially inaccurate LLM-generated news could help protect the public and ensure they are better informed.

2 Related Work (Model and Features)

2.1 Decision Trees (Model)

The key distinction between logistic regression and decision trees lies in how they fit to data. Decision trees partition the data into progressively smaller regions, while logistic regression fits a single line to separate the space into two categories. I chose logistic regression over decision trees because it is more efficient to train, particularly for binary classification problems like mine. Decision trees excel at handling complex, non-linear relationships, but logistic regression is better suited for predicting binary outcomes. Ultimately, logistic regression aligned more closely with both my project's objectives and the training dataset I used.

33 **2.2 Support Vector Machines (Model)**

34 The main difference between support vector machines (SVMs) and logistic regression is that SVMs
35 use support vectors to identify the optimal dividing line between data points. While the two approaches
36 are not vastly different, I chose logistic regression because it aligns better with the typical NLP
37 features outlined in existing research.

38 **2.3 Neural Networks (Model)**

39 The main difference between neural networks and logistic regression is that logistic regression forms
40 the foundation of a neural network, representing a single layer with one decision boundary, whereas
41 neural networks can have multiple layers and decision boundaries. I opted not to use neural networks
42 due to their complexity and the lack of necessary hardware to train the model efficiently. Additionally,
43 logistic regression could theoretically achieve comparable results for this problem without the added
44 complexity.

45 **2.4 Uppercase Word Count (Feature)**

46 Some popular AI detectors use uppercase word count to predict whether or not a piece of text is
47 human or LLM generated. I decided not to include this feature in my model based on the previous
48 experiments in "How to Detect AI-Generated Texts?" [3]

49 **2.5 Number of Parts of Speech (Feature)**

50 Some popular AI detectors use the number of different parts of speech (nouns, verbs, adjectives,
51 etc.) to identify LLM generated text. I decided not to include this feature in my model based on the
52 previous experiments in "How to Detect AI-Generated Text?" [3]

53 **2.6 Readability (Feature)**

54 I do have a readability score as one of my features however, I only use the Coleman Liau score
55 to determine its readability. There are many other readability scores including but not limited to
56 Flesch, Gunning Fog, Dale Chall, etc. The reason I decided to use Coleman Liau is it was the most
57 accurate in deciphering whether a text was written by an LLM or a human. This was also based on
58 the experiments outlined in "How to Detect AI-Generated Texts?" [3]

59 **3 Resources**

60 **3.1 Hardware**

61 Laptop: Apple Macbook M1 Pro 2021 CPU: Apple M1 (10 Cores) Memory: 32 GB unified memory
62 (shared between GPU and CPU)

63 All programming was done in VSCode in Python, code being run on .ipynb files, pandas dataframe
64 manipulation and sklearn logistic regression ran on CPU.

65 **3.2 Software**

66 Python Libraries: NumPy, Matplotlib, Pandas, Seaborn, textwrap, NLTK, collections, textstat, sci-kit
67 learn, language-tool-python, datasets

68 LLM Tools Used: GitHub Copilot for writing repetitive code blocks i.e. generating graphs/statistics,
69 Chat-GPT for generating custom test cases

70 **3.3 Datasets**

71 AI-Text-Detection-Pile (Hugging Face) - 1.418m [990k:340k]

72 [https://huggingface.co/datasets/artem9k/ai-text-detection-pile/viewer/](https://huggingface.co/datasets/artem9k/ai-text-detection-pile/viewer/default/train?p=5)
73 [default/train?p=5](https://huggingface.co/datasets/artem9k/ai-text-detection-pile/viewer/default/train?p=5)

74 LLM-Detect-AI-Generated-Text (Kaggle) - 27k [17k:11k]
 75 <https://www.kaggle.com/datasets/sunilthite/llm-detect-ai-generated-text-dataset>
 76 PaLm-Generated-Essays (Kaggle) - 1.3k [0:1.3k]
 77 <https://www.kaggle.com/datasets/kingki19/llm-generated-essay-using-palm-from-google-gen-ai>
 78 Combined-Set (Kaggle) - 87k [55k:32k]
 79 <https://www.kaggle.com/datasets/jdragonxherrera/augmented-data-for-llm-detect-ai-generated-text>
 80 AI-vs-Human-Text (Kaggle) - 500k [305k:195k]
 81 <https://www.kaggle.com/datasets/shanegerami/ai-vs-human-text>
 82 Human-vs-LLM-Corpus-Bloom-7B-and-GPT (Kaggle) - 800k [360k:440k]
 83 <https://www.kaggle.com/datasets/starblasters8/human-vs-llm-text-corpus>
 84 Total Dataset Distribution [1.73m:1m] [Human:AI]

85 **4 Methods**

86 **4.1 Featurization**

87 I used 5 standard NLP features that have the largest influence on a logistic regression binary classifier
 88 for LLM vs Human written text. [3] [2]

89 The following describes how each feature was calculated in python. Everything was done within one
 90 function, and applied to the dataset's dataframe using df.apply().

91 The text is also tokenized, where special characters, linking words, and stop words are removed. This
 92 is what specifies token vs word in the following.

93 **4.1.1 Coleman Liau Index (Readability)**

94 Using the textstat library, calculate the readability of the given text.

95 **4.1.2 Word Density**

96 Divide the number of characters by the number of words in the given text.

97 **4.1.3 Matches (Grammatical Errors)**

98 Using the language-tool-python library, count the number of grammatical errors in the given text.

99 **4.1.4 Title Word Count**

100 Count the number of tokens that start a sentence (title words) in the given text.

101 **4.1.5 Text Words (Text Length)**

102 The number of words in the given text.

103 **4.2 Logistic Regression**

104 Logistic regression is a type of regression that is often used to predict classes (typically binary like in
 105 this case). It does so by predicting the probability of a certain outcome (or class) based on predictor
 106 variables. The reason it's called "logistic" regression is due to its use of the logistic or "sigmoid"
 107 function:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

108 where z is the linear combination of the predictor variables and their coefficients:

$$z = w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n$$

109 $\sigma(z)$ represents the probability that the dependent variable y belongs to the class 1 (the positive class,
110 which in this case is LLM generated).

111 Given these two equations, this is the function that is minimized in logistic regression (logistic loss)
112 [4]:

$$f(w) = -\frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i(w^T x_i))) + \frac{\lambda}{2} \|w\|^2$$

113 where: w represents the weights, $\frac{\lambda}{2} \|w\|^2$ is the l2 regularization term, and $\log(1 + \exp(-y_i(w^T x_i)))$
114 is the logistic loss for each point (x_i, y_i) where x_i is the feature vector and y_i is the actual class. $\frac{1}{n}$ is
115 done to get the mean loss across the entire dataset.

116 During training, the model manipulates the coefficients w to minimize the equation.

117 The model then predicts the new sample's class based on this sigmoid function using the calculated
118 weights:

$$\hat{y} = \sigma(w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n)$$

119 where \hat{y} is the predicted probability and x_1, x_2, \dots, x_n are the individual features (not the feature
120 vector) of the new sample.

121 5 Experimental Results

122 5.0.1 Dataset Selection Methodology

123 I trained three separate models using three separate datasets:

- 124 • (1) A random sample of 10000 texts from a concatenated set of all datasets listed previously
125 (LLM-Detect-AI-Generated-Text, PaLm-Generated-Essays, Combined-Set, AI-vs-Human-
126 Text, Human-vs-LLM-Corpus-Bloom-7B-and-GPT, AI-Text-Detection-Pile). Around 6500
127 essays were human written, 3500 were LLM written. This ratio is a result of the ratio of
128 human to LLM written texts found in the total concatenated data set of around 3.1m texts.
- 129 • (2) 20000 texts, evenly split between human and LLM-written text, randomly sampled from
130 the AI-Text-Detection-Pile dataset.
- 131 • (3) 64000 texts, evenly split between human and LLM-written text, randomly sampled from
132 the Combined-Set dataset.

133 The models were trained in the given order, based off of intuition gained from the previous one. The
134 first model was trained with a portion of the total concatenated dataset (due to computing power
135 constraints). The second model was trained on the largest individual dataset, and the third model was
136 trained on an already curated dataset.

137 5.0.2 Parameters

138 Each model was evaluated with a train-test split of 3:2 using the scikit-learn python library's built-
139 in Logistic Regression function initialized with default parameters, with random-state set to 42
140 (arbitrary) to maintain consistency between runs.

141 Briefly, the default parameters are the use of the L2 penalty term, a stopping tolerance of $1e^{-4}$, a
142 regularization value $C = 1.0$, 100 max iterations, a class weight of 1 for both classes (Human vs
143 LLM), and the Limited-memory BFGS (LBFGS) solver to optimize.

144 From testing, changing the parameters made little tangible difference, so everything was left as
145 default.

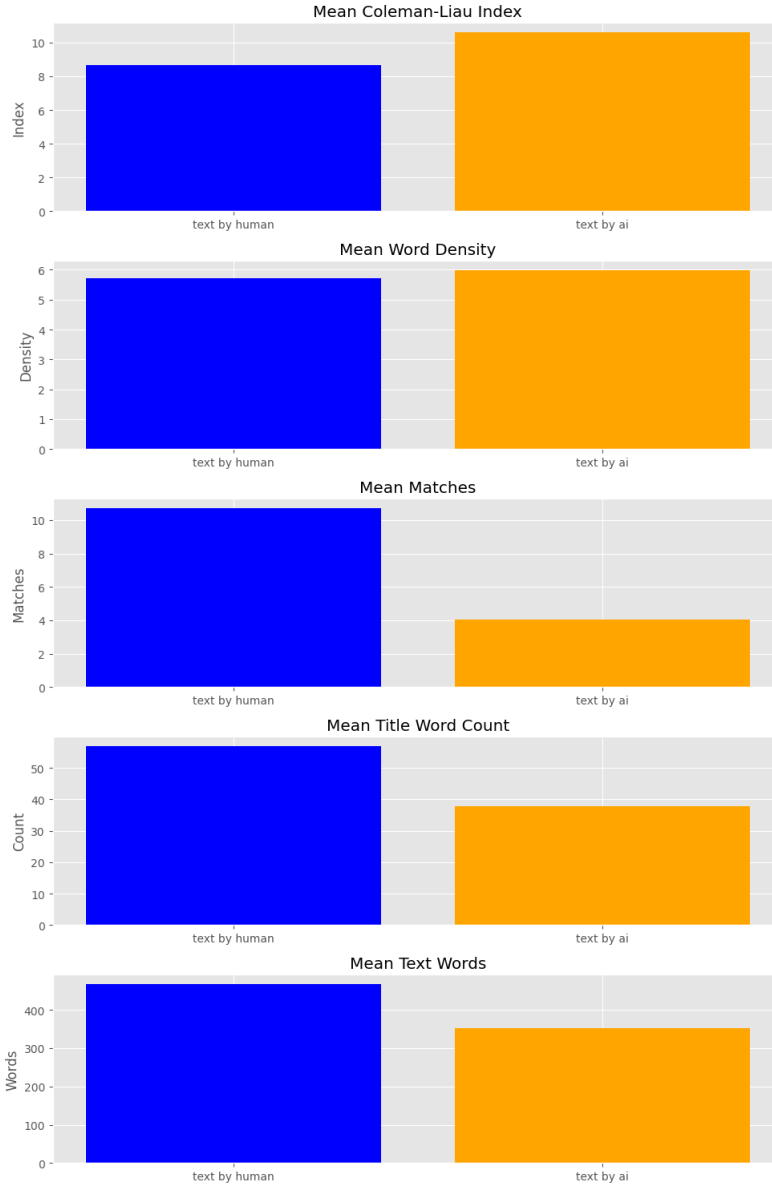


Figure 1: Mean Feature Values for all Datasets

Table 1: Logistic Regression Model w/ All Datasets Metrics

Accuracy: 0.731625

Train Loss: 0.5807722830077537

Test Loss: 0.5770358406982319

Name	Precision	Recall	f1-Score	Support
0 (Human)	0.736433	0.894161	0.807668	25208
1 (LLM)	0.715959	0.454638	0.556130	14792
Macro Avg	0.726196	0.674399	0.681899	40000
Weight Avg	0.728862	0.731625	0.714649	40000

Table 2: Logistic Regression Model w/ AI-Text-Detection-Pile Dataset Metrics

Accuracy: 0.678875				
Train Loss: 0.6397425870016508				
Test Loss: 0.6450721579241364				
Name	Precision	Recall	f1-Score	Support
0 (Human)	0.707965	0.615996	0.658786	4026
1 (LLM)	0.656215	0.742577	0.696730	3974
Macro Avg	0.682090	0.679286	0.677758	8000
Weight Avg	0.682258	0.678875	0.677635	8000

147 5.2 AI-Text-Detection-Pile 10000:10000

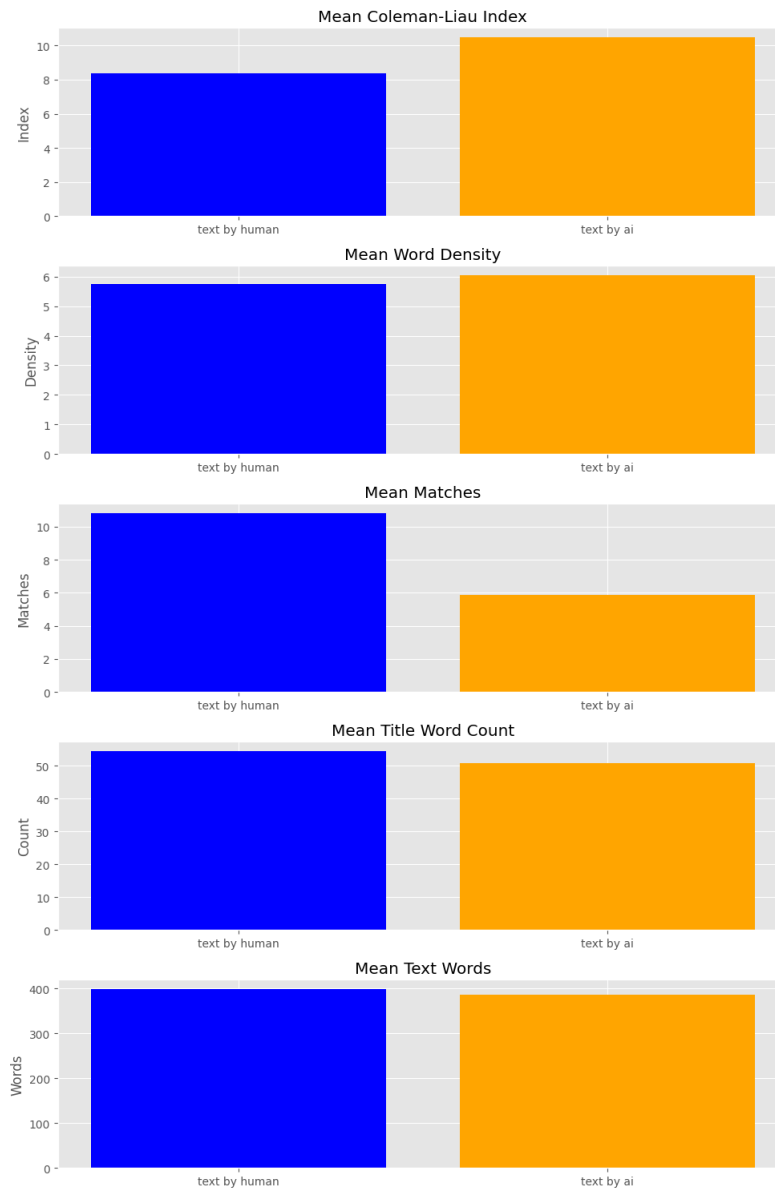


Figure 2: Mean Feature Values for AI-Text-Detection-Pile dataset

Table 3: Logistic Regression Model w/ Combined-Set Dataset Metrics

Accuracy: 0.8012890625				
Train Loss: 0.42126754359148816				
Test Loss: 0.42277091258149474				
Name	Precision	Recall	f1-Score	Support
0 (Human)	0.799938	0.803127	0.801529	12790
1 (LLM)	0.802649	0.799454	0.801048	12810
Macro Avg	0.801293	0.801290	0.801289	25600
Weight Avg	0.801294	0.801289	0.801289	25600

148 5.3 Combined-Set 32000:32000

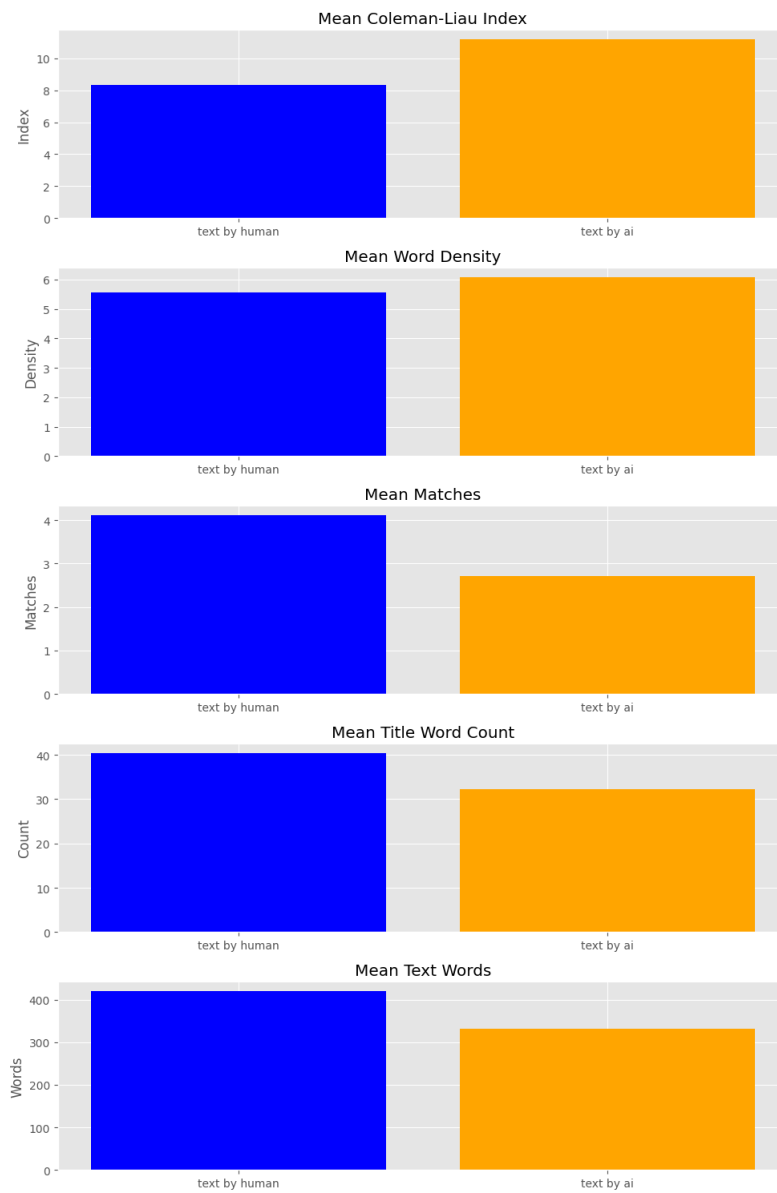


Figure 3: Mean Feature Values for Combined-Set

149 5.4 Results Analysis

150 5.4.1 Dataset Influence

151 Based purely on the mean, all 5 features seem to have a noticeable influence on whether the text is
152 AI-generated or Human-generated.

153 The worst performing model in terms of accuracy, the model trained on the AI-Text-Detection-Pile,
154 had the least differentiation between the mean for all 5 features.

155 The best-performing model, on the other hand, had the greatest differentiation between the mean for
156 all 5 features.

157 Throughout all 3 models, however, the difference was the same: AI text had a higher readability
158 score, greater word density, fewer grammatical errors, fewer title words (and therefore sentences),
159 and shorter length. The text length is only present ultimately to determine the ratio between itself and
160 the number of errors and title words, so just observing the mean in this way doesn't indicate anything.
161 However, given the clear differences they indicate, the model will pick up on its influence on if the
162 text is LLM or human-written.

163 Undersampling was used for the second and third datasets to attempt to rectify the first model's poor
164 accuracy.

165 5.4.2 Accuracy and Loss

166 The most accurate model with the least loss was the last model trained on the Combined-Set dataset,
167 potentially due to 3 reasons:

168 (1) The given dataset was curated already for a text detection competition, and may have already been
169 optimized for such a task

170 (2) The number of texts used was the largest, 64000 compared to 20000 and 10000. The first two
171 models were most likely underfitting the data.

172 (3) The use of undersampling compared to the first dataset to prevent drastic differences in accuracy
173 in classifying both texts.

174 Due to time and computational power constraints, no further testing could be conducted to see how
175 each of these reasons directly influenced the accuracy, but a combination of the 3 in training future
176 models would likely yield better results.

177 6 Conclusion

178 One approach that can be taken in the future is to use datasets with text generated exclusively by one
179 LLM, as this would prevent the model from needing to compare human-written text to text written by
180 several LLMs, which likely all have their own differences in features. This obviously reduces the use
181 case of the individual model, but when used in conjunction with multiple models, one might yield
182 better results. In addition to training on a larger dataset, this is what I believe to be the most effective
183 technique based on my results.

184 References

185 [1] The washington post, Dec 2023.

186 [2] Lorenz Mindner, Tim Schlippe, and Kristina Schaaff. Classification of human- and ai-generated
187 texts: Investigating features for chatgpt. In Tim Schlippe, Eric C. K. Cheng, and Tianchong Wang,
188 editors, *Artificial Intelligence in Education Technologies: New Development and Innovative
189 Practices*, pages 152–170, Singapore, 2023. Springer Nature Singapore.

190 [3] Trung Nguyen, Amartya Hatua, and Andrew Sung. How to detect ai-generated texts? pages
191 0464–0471, 10 2023.

192 [4] Mark Schmidt. Tutorial 8 logistic regression and stochastic gradient descent.