
Automated Stress Detection in Reddit Posts Using Traditional and Deep Learning Approaches

Ximing Wan

University of Melbourne, Text Analytics for Health (COMP90090-2025-SM1)

ABSTRACT

Objective: Using the Dreddit data set to find out how well deep learning and traditional models of machine learning perform in distinguishing between stressful and non-stressful Reddit posts.

Materials and Methods: The Dreddit dataset was analyzed, comprising 3,000 Reddit posts manually annotated as stressful or not-stressful. Cleaning, tokenization, and stopword removal were all part of the preprocessing. Feature representations were generated using TF-IDF, GloVe (100d embeddings from Stanford NLP), DistilBERT embeddings, and the domain-specific “mental-roberta-base” model from Hugging Face. Models tested included SVM, Random Forest, BiLSTM, and fine-tuned transformer models. Accuracy, precision, recall, and F1-score were among the evaluation metrics.

Results: TF-IDF + SVM achieved an accuracy of 72%, while Random Forest + TF-IDF performed comparably (72%). BiLSTM and BERT embeddings + SVM attained around 70% accuracy. Fine-tuning the mental-roberta-base model outperformed others, achieving 77% accuracy and superior macro-averaged F1-scores.

Discussion: The Transformer-based mental-roberta-base model consistently outperforms other models, demonstrating the advantages of domain-specific pre-training. However, while deep learning models can capture richer contextual features, simpler traditional models such as SVM are significantly faster to train in small sample environments, perform comparably well, and offer greater interpretability. For medical institutions without specialized equipment, this method may be more suitable for large-scale deployment.

Conclusion: Although transformer-based models optimised for domain data show promise in stress detection tasks, computational cost is still a major barrier to practical implementation.

Key words: TF-IDF, GloVe, SVM, Random Forest, BiLSTM, Transformer

INTRODUCTION

With the rise of social media platforms, many people openly “vent” their emotions online, making these platforms a valuable source of data for monitoring stress levels in the population. Reddit in particular, a popular forum where users can anonymously discuss personal challenges, often contains posts expressing stress. As for stress, it is a widespread public health issue. If stress is excessive or long-term, it can lead to many negative physical and mental health problems. Some authors even describe stress as a “silent killer”, emphasizing the importance of early stress management¹. Therefore, it is positive to detect psychological stress from personal network communication in a timely manner through machine learning, which can be used for early intervention before it leads to serious diseases such as anxiety, depression or cardiovascular disease.

Research on automatic stress detection emerges at the nexus of mental health and natural language process-

ing (NLP). Simple classifiers and dictionary-based techniques were the mainstays of early approaches. For instance, TensiStrength, a rule-based system created by Thelwall (2017), infers the degree of stress and relaxation in brief texts with a stress dictionary². Additional study has examined the linguistic characteristics linked to stress. Turcan and McKeown (2019) examined lexical patterns in stressed and unstressed texts using the Linguistic Inquiry Dictionary (LIWC)³. Subsequently, traditional machine learning techniques trained using bag-of-words models or TF-IDF features, such as logistic regression, naive Bayes, support vector machines (SVM), and random forests, have also been applied to stress classification in social media data. These methods have had some degree of success; for example, a recent study of academic Reddit communities by Oryngozha et al. used a logistic regression classifier on bag-of-words features with an accuracy of approximately 77% to 78%⁴.

More recently, deep learning approaches have advanced the state of the art in text classification, in-

cluding mental health domains. Bidirectional LSTMs (BiLSTMs) and transformer-based models can capture complex language patterns. Transformers like BERT (Bidirectional Encoder Representations from Transformers) have achieved state-of-the-art performance on many NLP tasks by learning rich contextual representations of words in context. In the realm of stress and mental health detection, researchers have begun fine-tuning such models on domain-specific data. For example, models initially pre-trained on general language (e.g. BERT) or on mental health-related corpora (e.g. MentalBERT/MentalRoBERTa) can be adapted to detect stress in social media posts. Ji et al. (2022) introduced MentalBERT and MentalRoBERTa, transformer models pre-trained on Reddit posts from mental health support communities, which improved performance on downstream mental health detection tasks⁵.

In this study, we address the task of binary stress detection (stress vs. no stress) in Reddit posts using the Dreaddit dataset. Dreaddit provides a collection of posts from five Reddit domains (such as r/relationships, r/college, etc.), with ground truth labels indicating whether each post contains a stressor or the author is expressing stress. In total, around 3,000 posts are human-labeled (training and test data) for supervised learning. We build and compare multiple models: two traditional classifiers (SVM and Random Forest) using TF-IDF features, and two deep learning models (a BiLSTM and a fine-tuned transformer). By evaluating their performance, we aim to quantify the benefits and drawbacks of each approach. We also generate visualizations (word clouds, learning curves, confusion matrices) to interpret model behavior and important features. Ultimately, we discuss how such models could be deployed in real-world clinical or public health scenarios—for example, as tools for monitoring stress levels in specific communities or for flagging high-stress posts to moderators or health professionals—while considering the practicality, resource requirements, and ethical implications.

METHODS

Data and Preprocessing

We used the Dreaddit dataset introduced by Turcan and McKeown which were split into a training set (dreaddit-train.csv) and test set (dreaddit-test.csv)³. The dataset consists of roughly 3,000 Reddit posts labeled by annotators for stress or no-stress, and were drawn from five different subreddit categories (relationships, academic, etc.), ensuring a mix of contexts in which users might express stress. Basic characteristics of the dataset included an average post length of a few sentences and an approximate balance between the stress and no-stress classes (Figure 1). For text preprocessing, we applied standard NLP cleaning steps. All text was

lowercased, tokenized and removed English stop words which carry little semantic content. We also removed or normalized punctuation, URLs, user mentions, and other non-alphanumeric characters, since these were not directly useful for detecting stress. We then performed POS tagging on each word and lemmatized it, thereby restoring the word to its basic dictionary form to improve semantic consistency and modeling effect. After preprocessing, each post was represented in multiple ways for input into different models, as described below.

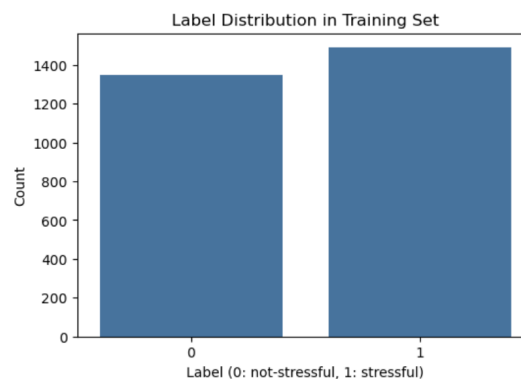


Figure 1. Label Distribution in Training Set.

Feature Extraction

For traditional machine learning models (Support Vector Machine (SVM) and Random Forest), we used Term Frequency-Inverse Document Frequency (TF-IDF). We built a vocabulary based on the preprocessed training set and converted each post into a high-dimensional sparse vector of TF-IDF weights for each word. This approach captures the importance of words (words that appear frequently in a post but rarely in other posts receive higher weights). We limited the size of the vocabulary to the top 5,000 words by frequency to reduce the dimensionality and filter out very rare words. Figure 2 shows the word cloud of the words with the highest weights extracted based on the TF-IDF feature.

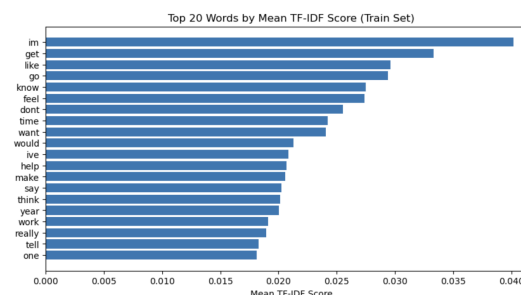


Figure 2. Top 20 Words by Mean TF-IDF Score (Train Set)

We used GloVe vectors to represent each word as a dense vector with 100 dimensions for the BiLSTM model. GloVe (global vectors) encode semantic similarity and are pre-trained on a large corpus (Pennington, Socher, & Manning, 2014)⁶. During LSTM training, we updated (fine-tuned) pre-trained GloVe 6B corpus vectors to fit our domain. Every post was transformed into a series of word vectors that were either padded or truncated to a predetermined length of 100 tokens, and BiLSTM used these sequences as input.

For Transformer-based models, we fine-tune mental-roberta-base. This model belongs to the MentalBERT family (Domain-Specific Language Model for Mental Health), initialized by RoBERTa architecture (12-layer Transformer), and trained on a large corpus of online posts on mental health issues (Ji et al., 2021)⁷. We loaded the mental-roberta-base model and tokenizer via Hugging Face Transformers. During fine-tuning, the model generates context vectors for the entire article, as well as a special CLS token vector that represents the aggregate meaning of the article and feeds it into the classifier layer. We also tried using BERT-base (pre-trained on general English) as a baseline to compare the effect of domain-specific pre-training. Preliminary results show that mental-roberta-base outperforms the general BERT-base, which is consistent with the finding that domain-specific pre-training can improve performance on mental health NLP tasks. Therefore, we report results for mental-roberta-based models as the primary Transformer approach.

Classifiers and Model Training

Four classification models were trained to identify signs of stress in Reddit posts. First, a SVM with a linear kernel and the default value of the regularization parameter C (1.0) was trained on TF-IDF features. This model is robust to high-dimensional sparse text data and is able to quickly and effectively classify “stressed” and “non-stressed” categories.

Second, a random forest model was trained on TF-IDF features with 100 decision trees and default settings (e.g., Gini impurity and maximum number of features as the square root of the total number of features). This model captures nonlinear interactions between words and has some ability to explain feature importance, although this is limited in the case of very high dimensional text.

The third model is a deep learning-based bidirectional LSTM (BiLSTM) neural network. We initialized the embedding layer with pre-trained GloVe word vectors and constructed an LSTM layer with 128 units in the forward and backward passes. The bidirectional structure can capture the contextual information of the sentence at the same time, such as expressions that require contextual judgment, such as “not coping well”. The final concatenated hidden state is input into a fully connected

layer to complete the binary classification. During the training process, the Adam optimizer (learning rate of 0.001) and dropout (0.3) are used to prevent overfitting, and the training rounds are controlled by the early stopping mechanism.

Finally, we fine-tuned the RoBERTa model (mental-roberta-base) based on the Transformer architecture. We added a linear classification layer to the output of its CLS token to determine whether the post expresses stress. Fine-tuning was performed using Hugging Face’s Trainer, set to 3 epochs, batch size of 16, learning rate of 2e-5, and a linear learning rate decay strategy. Despite the small amount of training data, the model still showed good generalization ability because it was originally pre-trained on mental health data. For comparison, we also fine-tuned the general bert-base-uncased model with the same settings as RoBERTa (e.g., learning rate 2e-5, batch size 16, training for 3 epochs). Although this model was not specifically pre-trained on mental health corpus, it still performed well in this task, verifying the wide applicability of pre-trained Transformer models in stress detection.

Ethics Statement

This research exclusively uses publicly available Reddit posts from the Dreddit dataset. The data was originally collected and released in prior work with annotations, and we comply with the terms of use of the dataset. No additional private or identifiable user information was collected. The Reddit platform policy allows the research use of public comments; however, we acknowledge that content can be sensitive as it pertains to mental health. We took care to anonymize the data (removing any usernames or personal references if present) in our processing. Because the data are public and anonymized, this study was exempt from institutional review by the IRB. We do not intend to redistribute the data or any content beyond derived analyses. We discuss only aggregate results and patterns (e.g., common stress indicators) to avoid any potential harm or stigma to individuals. All experiments were carried out according to relevant ethical guidelines for research on human communication. Finally, we note that automated stress detection is intended as a supporting tool for mental health outreach; any deployed system should include human oversight and robust privacy protections.

RESULTS

For the classification performance of each model is summarised in Figure 3-7. First, with precision, recall, and F1 scores balanced at 0.72 across both categories, the SVM model with TF-IDF features achieves 72% accuracy. likewise, the random forest model with TF-IDF

features attains the same accuracy (72%), yet it has a slightly higher macro-average accuracy (0.74) and a higher recall for the stress category (0.87). Furthermore, the BiLSTM model with max pooling and a deeper classifier achieves 72% accuracy with a balance of precision and recall (both 0.72).

As can be seen, traditional machine learning models are usually fast and have fairly high accuracy. However, the enhanced domain-specific Transformer model (mental-roberta-base) outperformed traditional methods in this task (achieving an accuracy of 77%). The model has high macro-average precision (0.78) and recall (0.77), indicating that it is able to accurately identify stressful posts. In contrast, the performance of BERT embeddings alone (without task-specific fine-tuning) dropped to 69% overall accuracy and did not outperform traditional methods. This shows that fine-tuning and domain adaptation are critical to fully exploit the power of Transformer models in stress detection.

```

=== Test Set Performance (SVM + TF-IDF) ===
      precision    recall  f1-score   support

     0       0.72       0.68       0.70       346
     1       0.72       0.76       0.74       369

 accuracy          0.72          0.72          0.72       715
 macro avg       0.72       0.72       0.72       715
 weighted avg    0.72       0.72       0.72       715

```

Figure 3. Test Set Performance (SVM + TF-IDF)

```

=== Test Set Performance (Random Forest + TF-IDF) ===
      precision    recall  f1-score   support

     0       0.80       0.56       0.66       346
     1       0.68       0.87       0.76       369

 accuracy          0.72          0.72          0.72       715
 macro avg       0.74       0.72       0.71       715
 weighted avg    0.74       0.72       0.71       715

```

Figure 4. Test Set Performance (Random Forest + TF-IDF)

```

=== Test Set Performance (BiLSTM + Max Pooling + Deeper Classifier) ===
      precision    recall  f1-score   support

     0       0.69       0.74       0.72       346
     1       0.74       0.69       0.72       369

 accuracy          0.72          0.72          0.72       715
 macro avg       0.72       0.72       0.72       715
 weighted avg    0.72       0.72       0.72       715

```

Figure 5. Test Set Performance (BiLSTM + Max Pooling + Deeper Classifier)

```

=== Test Set Performance (BERT Embeddings + SVM) ===
      precision    recall  f1-score   support

     0       0.69       0.66       0.68       346
     1       0.69       0.72       0.70       369

 accuracy          0.69          0.69          0.69       715
 macro avg       0.69       0.69       0.69       715
 weighted avg    0.69       0.69       0.69       715

```

Figure 6. Test Set Performance (BERT Embeddings + SVM)

```

=== Classification Report on Test Set (mental-roberta-base) ===
      precision    recall  f1-score   support

     0       0.79       0.73       0.76       346
     1       0.76       0.82       0.79       369

 accuracy          0.77          0.77          0.77       715
 macro avg       0.78       0.77       0.77       715
 weighted avg    0.78       0.77       0.77       715

```

Figure 7. Classification Report on Test Set (mental-roberta-base)

DISCUSSION

The performance differences observed between different models in this experiment may stem from their ability to capture complex stress-related language features. To explain, the excellent performance of mental-roberta-base suggests that pre-training on mental health-related data can enhance contextual understanding. In contrast, the poor performance of individual BiLSTM and BERT embeddings may be due to their limited domain adaptation due to lack of fine-tuning, which may make them misled by the context; for example, "using stress to make cookies" is understood as psychological stress. However, traditional models such as SVM perform surprisingly well on small datasets, which may benefit from simpler decision boundaries that are able to fit the data without overfitting. As the result, these patterns highlight the importance of domain specific pre-training and the trade-off between model complexity and training data size in stress detection.

CONCLUSION

In summary, this study finds that the fine-tuned mental-roberta-base outperforms traditional methods as well as unspecifically optimized deep learning for stress detection on the Dreddit dataset; this result highlights the advantages of domain-specific Transformers in capturing complex stressful linguistic cues. However, classical models sacrifice some accuracy for faster training and higher interpretability. Future research should build on these findings by training on larger and more diverse datasets to improve generalization as well as exploring domain adaptation to maintain performance across different social platforms. These measures will contribute

to the responsible deployment of automated stress detection systems in real-world environments.

REFERENCES

1. Balwan WK, Kour S. A Systematic Review of Hypertension and Stress – The Silent Killers. *Scholars Academic Journal of Biosciences*. 2021;9(6):150-4. Available from: <https://saspublishers.com/media/articles/SAJB96154158.pdf>.
2. Thelwall M. TensiStrength: Stress and relaxation magnitude detection for social media texts. *Information Processing & Management*. 2017;53(1):106-21. Available from: <https://www.sciencedirect.com/science/article/pii/S0306457316302321>.
3. Turcan E, McKeown K. Dreddit: A Reddit Dataset for Stress Analysis in Social Media. In: *Proceedings of the Tenth International Workshop on Health Text Mining and Information Analysis (LOUHI 2019)*. Hong Kong: Association for Computational Linguistics; 2019. p. 97-107. Available from: <https://aclanthology.org/D19-6213/>.
4. Oryngozha N, Shamo P, Igali A. Detection and Analysis of Stress-Related Posts in Reddit's Academic Communities. *IEEE Access*. 2024;12:14932-48. Available from: <https://ieeexplore.ieee.org/document/10412045>.
5. Ji S, Zhang T, Ansari L, Fu J, Tiwari P, Cambria E. MentalBERT: Publicly Available Pretrained Language Models for Mental Healthcare. In: *Proceedings of the Thirteenth Language Resources and Evaluation Conference*. European Language Resources Association; 2022. p. 4346-52.
6. Pennington J, Socher R, Manning CD. GloVe: Global Vectors for Word Representation. In: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics; 2014. p. 1532-43. Available from: <https://aclanthology.org/D14-1162/>.
7. Ji S, Zhang T, Ansari L, Fu J, Tiwari P, Cambria E. MentalBERT: Publicly Available Pretrained Language Models for Mental Healthcare. *arXiv preprint arXiv:2110.15621*. 2021. Available from: <https://huggingface.co/papers/2110.15621>.

APPENDIX: KEY JUPYTER NOTEBOOK CODE

GitHub:<https://github.com/Nooob233/nlp4health-a3.git>

```

1 # ----- Import necessary libraries -----
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import re
7 from sklearn.model_selection import train_test_split
8 import nltk
9 from nltk.corpus import stopwords, wordnet
10 from nltk.stem import WordNetLemmatizer
11 from nltk.tokenize import word_tokenize
12 from nltk import pos_tag
13
14 # == Read the dataset ==
15 # Load Dreddit training and test sets
16 train_path = "CORPORA/REDDIT/DREADDIT/dreddit-train.csv"
17 test_path = "CORPORA/REDDIT/DREADDIT/dreddit-test.csv"
18
19 train_raw = pd.read_csv(train_path)
20 test_raw = pd.read_csv(test_path)
21
22 # Keep only the 'text' and 'label' columns
23 train_df = train_raw[['text', 'label']].copy()
24 test_df = test_raw[['text', 'label']].copy()
25
26 print("Number of training samples:", train_df.shape[0])
27 print("Number of test samples:", test_df.shape[0])
28 print("\nFirst 3 rows of training data:\n", train_df.head(3))
29
30 # ----- Exploratory Data Analysis -----
31 # Visualize the label distribution in the training set
32 plt.figure(figsize=(6,4))
33 sns.countplot(data=train_df, x='label')
34 plt.title("Label Distribution in Training Set")
35 plt.xlabel("Label (0: not-stressful, 1: stressful)")
36 plt.ylabel("Count")
37 plt.show()
38
39 # ----- Text Cleaning and Preprocessing -----
40 # Download stopwords list and WordNet resources
41 nltk.download('stopwords')
42 nltk.download('wordnet')
43 nltk.download('punkt')
44 nltk.download('averaged_perceptron_tagger')
45
46 # Function to map Penn Treebank tags to WordNet POS tags
47 def get_wordnet_pos(treebank_tag):
48     """
49     Map Penn Treebank POS tags to WordNet POS tags
50     """
51     if treebank_tag.startswith('J'):
52         return wordnet.ADJ
53     elif treebank_tag.startswith('V'):
54

```

```

58         return wordnet.VERB
59     elif treebank_tag.startswith('N'):
60         return wordnet.NOUN
61     elif treebank_tag.startswith('R'):
62         return wordnet.ADV
63     else:
64         return wordnet.NOUN # fallback to
65         noun
66 # Load stop words and initialize
67     lemmatizer
68 stop_words = set(stopwords.words('english'))
69 lemmatizer = WordNetLemmatizer()
70 # Define a function to clean the text data
71 def clean_text(text):
72     """
73     Full cleaning process: Remove HTML tags
74     and URLs, Remove non-alphabetic
75     characters, Lowercase text,
76     Tokenize and remove stop words, POS
77     tagging and lemmatization
78     """
79     text = re.sub(r'<.*?>', '', str(text))
80     # Remove HTML tags
81     text = re.sub(r'http\S+|www\S+|https\S+', '', text) # Remove URLs
82     text = re.sub(r'[^a-zA-Z\s]', '', text) # Remove non-alphabetic
83     characters
84     text = text.lower() # Lowercase
85     tokens = word_tokenize(text) #
86     Tokenize
87     tokens = [word for word in tokens if
88     word not in stop_words] # Remove
89     stopwords
90     pos_tags = pos_tag(tokens) # POS
91     tagging
92     # Lemmatization with accurate POS
93     lemmatized_tokens = [
94         lemmatizer.lemmatize(word,
95         get_wordnet_pos(tag))
96         for word, tag in pos_tags
97     ]
98     return ' '.join(lemmatized_tokens)
99 # Apply the cleaning function to training
100 and test data
101 train_df['clean_text'] = train_df['text'].
102 apply(clean_text)
103 test_df['clean_text'] = test_df['text'].
104 apply(clean_text)
105 # Display a few cleaned examples for
106 sanity check
107 print("\nExample of cleaned text:\n",
108 train_df[['text', 'clean_text']].head

```

```

(3))
# ----- Split the training data
into train and validation sets (80/20)
-----
# Stratify ensures the label distribution
is preserved
train_final, valid_final =
train_test_split(
    train_df[['clean_text', 'label']],
    test_size=0.2,
    random_state=42,
    stratify=train_df['label']
)
print("\nSize of training set after split:
", train_final.shape)
print("Size of validation set after split:
", valid_final.shape)
# == Save cleaned datasets for future use
==
# Save cleaned training, validation, and
test sets to CSV files
train_final.to_csv("clean_train.csv",
    index=False)
valid_final.to_csv("clean_valid.csv",
    index=False)
test_df[['clean_text', 'label']].to_csv("
clean_test.csv", index=False)
print("\nCleaned datasets saved
successfully.")
# ----- Visualize the distribution
of text lengths before and after
cleaning in the same plot, with
annotations -----
# Compute the number of words before and
after cleaning
train_df['text_length'] = train_df['text']
.apply(lambda x: len(str(x).split()))
train_df['clean_text_length'] = train_df['
clean_text'].apply(lambda x: len(str(x)
.split()))
plt.figure(figsize=(8,5))
# Plot original text length distribution
sns.histplot(train_df['text_length'], bins
=30, kde=True, color='blue', label='
Before Cleaning', alpha=0.5)
# Plot cleaned text length distribution
sns.histplot(train_df['clean_text_length']
], bins=30, kde=True, color='green',
label='After Cleaning', alpha=0.5)
# Calculate most frequent text lengths (

```

```

mode)
134 mode_before = train_df['text_length'].mode()
135 mode_after = train_df['clean_text_length'].mode()
136 # Calculate maximum text lengths
137 max_before = train_df['text_length'].max()
138 max_after = train_df['clean_text_length'].max()
139 # Annotate mode (most frequent)
140 plt.axvline(mode_before, color='blue',
141             linestyle='--', label=f"Most Frequent (Before): {mode_before}")
142 plt.axvline(mode_after, color='green',
143             linestyle='--', label=f"Most Frequent (After): {mode_after}")
144 # Annotate maximum length
145 plt.scatter(max_before, 0, color='blue',
146             marker='o', s=100, label=f"Max (Before): {max_before}")
147 plt.scatter(max_after, 0, color='green',
148             marker='o', s=100, label=f"Max (After): {max_after}")
149 plt.title("Comparison of Text Length
150           Distribution (Before vs After Cleaning)")
151 plt.xlabel("Number of Words")
152 plt.ylabel("Number of Samples")
153 plt.legend()
154 plt.show()
155 # ----- Generate word cloud
156 # visualization for high-frequency words -----
157 from wordcloud import WordCloud
158 # Concatenate all cleaned text into a
159 # single string
160 all_text = ' '.join(train_df['clean_text']
161                      ).dropna().tolist()
162 # Create and configure the word cloud
163 # object
164 wordcloud = WordCloud(width=800, height
165                       =400,
166                       background_color='white',
167                       max_words=100, # show top 100
168                       words
169                       contour_width=1,
170                       contour_color='steelblue').
171                       generate(all_text)
172 # Plot the word cloud
173 plt.figure(figsize=(10, 5))
174 plt.imshow(wordcloud, interpolation='bilinear')
175 plt.axis('off')
176 plt.title("Word Cloud of Cleaned Training
177           Text")
178 plt.show()
179 # -----Re-import pandas and
180 # read the cleaned datasets -----
181 import pandas as pd
182 # Read cleaned training, validation, and
183 # test sets
184 clean_train = pd.read_csv("clean_train.csv")
185 clean_valid = pd.read_csv("clean_valid.csv")
186 clean_test = pd.read_csv("clean_test.csv")
187 # Display the first few rows of each to
188 # confirm successful loading
189 print("First 3 rows of clean_train:")
190 print(clean_train.head(3))
191 print("\nFirst 3 rows of clean_valid:")
192 print(clean_valid.head(3))
193 print("\nFirst 3 rows of clean_test:")
194 print(clean_test.head(3))
195 # *----- TF-IDF Feature
196 # Extraction -----
197 from sklearn.feature_extraction.text
198 import TfidfVectorizer
199 # Initialize TF-IDF vectorizer for
200 # unigrams
201 tfidf_vectorizer = TfidfVectorizer(
202     max_features=5000)
203 # Fit on the training set and transform
204 # train/valid/test sets
205 X_train_tfidf = tfidf_vectorizer.
206 fit_transform(clean_train['clean_text'])
207 X_valid_tfidf = tfidf_vectorizer.transform(
208     clean_valid['clean_text'])
209 X_test_tfidf = tfidf_vectorizer.transform(
210     clean_test['clean_text'])
211 print("\nTF-IDF features generated
212       successfully.")
213 print("Shape of training TF-IDF matrix:",
214       X_train_tfidf.shape)
215 print("Shape of validation TF-IDF matrix:")

```



```

, X_valid_tfidf.shape)
210 print("Shape of test TF-IDF matrix:",
      X_test_tfidf.shape)
211
212
213 # ----- Visualize top 20 high-
      frequency words in training data
214 import numpy as np
215 import matplotlib.pyplot as plt
216
217 feature_names = tfidf_vectorizer.
      get_feature_names_out()
218 mean_tfidf = np.array(X_train_tfidf.mean(
      axis=0)).flatten()
219 top_indices = mean_tfidf.argsort()
      [::-1][:20]
220
221 top_words = feature_names[top_indices]
222 top_scores = mean_tfidf[top_indices]
223
224 plt.figure(figsize=(10, 5))
225 plt.barh(top_words[::-1], top_scores
      [::-1])
226 plt.xlabel("Mean TF-IDF Score")
227 plt.title("Top 20 Words by Mean TF-IDF
      Score (Train Set)")
228 plt.show()
229
230
231 # ----- SVM model with TFIDF
      training & validation
232 from sklearn.svm import LinearSVC
233 from sklearn.metrics import
      classification_report, confusion_matrix
234 , ConfusionMatrixDisplay
235
236 # Prepare target labels
237 y_train = clean_train['label']
238 y_valid = clean_valid['label']
239 y_test = clean_test['label']
240
241 # Initialize and train SVM
242 svm_tfidf = LinearSVC(random_state=42)
243 svm_tfidf.fit(X_train_tfidf, y_train)
244
245 # Predict on validation set
246 y_valid_pred = svm_tfidf.predict(
      X_valid_tfidf)
247
248 # Evaluate validation set
249 print("\n=== Validation Set Performance (
      SVM + TF-IDF) ===")
250 print(classification_report(y_valid,
      y_valid_pred))
251
252 # Visualize confusion matrix
253 cm = confusion_matrix(y_valid,
      y_valid_pred)
254 disp = ConfusionMatrixDisplay(
      confusion_matrix=cm)
255 disp.plot()
256 plt.title("Confusion Matrix: SVM + TF-IDF
      on Validation Set")
257 plt.show()
258
259 # ----- Test set evaluation
260 y_test_pred = svm_tfidf.predict(
      X_test_tfidf)
261 print("\n=== Test Set Performance (SVM +
      TF-IDF) ===")
262 print(classification_report(y_test,
      y_test_pred))
263
264 # Load cleaned data for TF-IDF features
265 import pandas as pd
266 from sklearn.feature_extraction.text
      import TfidfVectorizer
267 from sklearn.ensemble import
      RandomForestClassifier
268 from sklearn.metrics import
      classification_report
269
270 clean_train = pd.read_csv("clean_train.csv")
271 clean_valid = pd.read_csv("clean_valid.csv")
272 clean_test = pd.read_csv("clean_test.csv")
273
274 # TF-IDF Feature Extraction
275 tfidf_vectorizer = TfidfVectorizer(
      max_features=5000)
276 X_train_tfidf = tfidf_vectorizer.
      fit_transform(clean_train['clean_text']
      ])
277 X_valid_tfidf = tfidf_vectorizer.transform(
      clean_valid['clean_text'])
278 X_test_tfidf = tfidf_vectorizer.transform(
      clean_test['clean_text'])
279
280 # Train Random Forest on TF-IDF
281 y_train = clean_train['label']
282 y_valid = clean_valid['label']
283 y_test = clean_test['label']
284
285 rf_tfidf = RandomForestClassifier(
      n_estimators=100, random_state=42)
286 rf_tfidf.fit(X_train_tfidf, y_train)
287
288 # Evaluate on test set
289 y_test_pred_tfidf = rf_tfidf.predict(
      X_test_tfidf)
290 print("\n=== Test Set Performance (Random
      Forest + TF-IDF) ===")
291 print(classification_report(y_test,
      y_test_pred_tfidf, digits=2))
292
293
294
295 # ----- Load Cleaned Data & Build

```



```

Vocab-----
296 import pandas as pd
297 import numpy as np
298 import torch
299 import torch.nn as nn
300 import torch.optim as optim
301 from torch.utils.data import Dataset,
    DataLoader
302 from torch.nn.utils.rnn import
    pad_sequence
303 from collections import Counter
304
305 # Load cleaned datasets
306 clean_train = pd.read_csv("clean_train.csv"
307                             ")
308 clean_valid = pd.read_csv("clean_valid.csv"
309                             ")
310 clean_test = pd.read_csv("clean_test.csv"
311                             ")
312
313 # Build vocab manually
314 word_counter = Counter()
315 for text in clean_train['clean_text']:
316     word_counter.update(text.split())
317
318 vocab = {"<pad>": 0, "<unk>": 1}
319 for idx, (word, _) in enumerate(
320     word_counter.most_common(), start=2):
321     vocab[word] = idx
322
323 def tokenize(text):
324     return text.split()
325
326 def text_to_ids(text, vocab, max_len=100):
327     tokens = tokenize(text)
328     ids = [vocab.get(token, vocab["<unk>"]
329         ) for token in tokens]
330     return ids[:max_len]
331
332 # ----- Load GloVe and Build Embedding
333 Matrix-----
334 def load_glove_vectors(glove_file_path,
335     embedding_dim=100):
336     word_vectors = {}
337     with open(glove_file_path, 'r',
338         encoding='utf-8') as f:
339         for line in f:
340             parts = line.strip().split()
341             word = parts[0]
342             vector = np.asarray(parts[1:],
343                 dtype='float32')
344             word_vectors[word] = vector
345     return word_vectors
346
347 glove_vectors = load_glove_vectors("glove
348     .6B.100d.txt", embedding_dim=100)
349 print("Loaded GloVe vectors. Example:",
350     list(glove_vectors.items())[:1])
351
352 vocab_size = len(vocab)
353 embedding_dim = 100
354
355 embedding_matrix = np.random.normal(scale
356     =0.6, size=(vocab_size, embedding_dim))
357 for word, idx in vocab.items():
358     vector = glove_vectors.get(word)
359     if vector is not None:
360         embedding_matrix[idx] = vector
361 embedding_matrix = torch.tensor(
362     embedding_matrix, dtype=torch.float32)
363
364 # ----- Dataset & DataLoader -----
365 class RedditDataset(Dataset):
366     def __init__(self, texts, labels,
367         vocab, max_len=100):
368         self.texts = texts
369         self.labels = labels
370         self.vocab = vocab
371         self.max_len = max_len
372
373     def __len__(self):
374         return len(self.texts)
375
376     def __getitem__(self, idx):
377         token_ids = text_to_ids(self.texts
378             [idx], self.vocab, self.max_len
379             )
380         return torch.tensor(token_ids),
381             torch.tensor(self.labels[idx])
382
383 def collate_fn(batch):
384     texts, labels = zip(*batch)
385     texts_padded = pad_sequence(texts,
386         batch_first=True, padding_value=
387         vocab["<pad>"])
388     return texts_padded, torch.tensor(
389         labels)
390
391 train_dataset = RedditDataset(clean_train[
392     'clean_text'].tolist(), clean_train['
393     label'].tolist(), vocab)
394 valid_dataset = RedditDataset(clean_valid[
395     'clean_text'].tolist(), clean_valid['
396     label'].tolist(), vocab)
397 test_dataset = RedditDataset(clean_test[
398     'clean_text'].tolist(), clean_test['
399     label'].tolist(), vocab)
400
401 train_loader = DataLoader(train_dataset,
402     batch_size=32, shuffle=True, collate_fn
403     =collate_fn)
404 valid_loader = DataLoader(valid_dataset,
405     batch_size=32, collate_fn=collate_fn)
406 test_loader = DataLoader(test_dataset,
407     batch_size=32, collate_fn=collate_fn)
408
409 #----- BiLSTM Model with Max Pooling
410 and Deeper Classifier-----
411 class BiLSTMMaxPoolClassifier(nn.Module):
412     def __init__(self, vocab_size,
413         embedding_dim, hidden_dim,

```

```

        output_dim, pad_idx,
        embedding_matrix):
382     super(BiLSTMMaxPoolClassifier,
        self).__init__()
383     self.embedding = nn.Embedding(
        vocab_size, embedding_dim,
        padding_idx=pad_idx)
384     self.embedding.weight.data.copy_(
        embedding_matrix)
385     self.embedding.weight.
        requires_grad = True
386
387     self.lstm = nn.LSTM(embedding_dim,
        hidden_dim, bidirectional=True
        , batch_first=True)
388
389     self.classifier = nn.Sequential(
        nn.Linear(hidden_dim * 2,
        hidden_dim),
391         nn.ReLU(),
392         nn.Dropout(0.3),
393         nn.Linear(hidden_dim,
        output_dim)
394     )
395
396     def forward(self, text):
397         embedded = self.embedding(text)
398         output, _ = self.lstm(embedded)
399         pooled, _ = torch.max(output, dim
        =1)
400         return self.classifier(pooled)
401
402 hidden_dim = 128
403 output_dim = 2
404 pad_idx = vocab["<pad>"]
405 model = BiLSTMMaxPoolClassifier(len(vocab)
        , embedding_dim, hidden_dim, output_dim,
        pad_idx, embedding_matrix)
406
407 # ----- Training & Evaluation -----
408 device = torch.device('cuda' if torch.cuda
        .is_available() else 'cpu')
409 model = model.to(device)
410 criterion = nn.CrossEntropyLoss()
411 optimizer = optim.Adam(model.parameters(),
        lr=1e-3)
412
413 def train(model, iterator):
414     model.train()
415     total_loss = 0
416     for texts, labels in iterator:
417         texts, labels = texts.to(device),
        labels.to(device)
418         optimizer.zero_grad()
419         predictions = model(texts)
420         loss = criterion(predictions,
        labels)
421         loss.backward()
422         optimizer.step()
423         total_loss += loss.item()
424
425     return total_loss / len(iterator)
426
427 def evaluate(model, iterator):
428     model.eval()
429     total_loss = 0
430     correct = 0
431     total = 0
432     with torch.no_grad():
433         for texts, labels in iterator:
434             texts, labels = texts.to(
        device), labels.to(device)
435             predictions = model(texts)
436             loss = criterion(predictions,
        labels)
437             total_loss += loss.item()
438             preds = predictions.argmax(dim
        =1)
439             correct += (preds == labels).
        sum().item()
440             total += labels.size(0)
441     return total_loss / len(iterator),
        correct / total
442
443 for epoch in range(5):
444     train_loss = train(model, train_loader
        )
445     valid_loss, valid_acc = evaluate(model
        , valid_loader)
446     print(f"Epoch {epoch+1}: Train Loss={
        train_loss:.4f}, Valid Loss={
        valid_loss:.4f}, Valid Acc={
        valid_acc:.4f}")
447
448 # ----- Test Set Performance (
        Classification Report) -----
449 from sklearn.metrics import
        classification_report
450
451 def get_preds_and_labels(model, iterator):
452     model.eval()
453     all_preds = []
454     all_labels = []
455     with torch.no_grad():
456         for texts, labels in iterator:
457             texts = texts.to(device)
458             predictions = model(texts)
459             preds = predictions.argmax(dim
        =1).cpu().numpy()
460             all_preds.extend(preds)
461             all_labels.extend(labels.numpy
        ())
462     return all_labels, all_preds
463
464 y_test_true, y_test_pred =
        get_preds_and_labels(model, test_loader
        )
465 print("\n=== Test Set Performance (BiLSTM
        + Max Pooling + Deeper Classifier) ==="
        )

```

```

466 print(classification_report(y_test_true, 502
    y_test_pred, digits=2)) 503
467 504
468 # Load Cleaned Data 505
469 import pandas as pd
470
471 clean_train = pd.read_csv("clean_train.csv" 506
    ")
472 clean_valid = pd.read_csv("clean_valid.csv" 507
    ")
473 clean_test = pd.read_csv("clean_test.csv")
474 # Load pretrained BERT tokenizer & model
475 from transformers import AutoTokenizer, 508
    AutoModel 509
476 import torch 510
477
478 device = torch.device('cuda' if torch.cuda 511
    .is_available() else 'cpu') 512
479
480 # Using DistilBERT (lighter, faster) 513
481 tokenizer = AutoTokenizer.from_pretrained(" 514
    distilbert-base-uncased") 515
482 bert_model = AutoModel.from_pretrained(" 516
    distilbert-base-uncased").to(device) 517
483 518
484 # Function to get sentence embeddings from 519
    BERT 520
485 def get_bert_embeddings(texts, tokenizer, 521
    model, batch_size=32, max_length=128): 522
486     """ 523
487     Given a list of texts, compute the [ 524
        CLS] token embeddings from BERT. 525
488     """ 526
489     all_embeddings = [] 527
490
491     for i in range(0, len(texts), 528
        batch_size): 529
492         batch_texts = texts[i:i+batch_size] 530
493         encoded = tokenizer(batch_texts, 531
            padding=True, truncation=True, 532
            max_length=max_length, 533
            return_tensors="pt") 534
494         input_ids = encoded['input_ids']. 535
            to(device) 536
495         attention_mask = encoded[' 537
            attention_mask'].to(device) 538
496
497         with torch.no_grad(): 539
498             outputs = model(input_ids, 540
                attention_mask=
                attention_mask)
499             last_hidden_state = outputs.
                last_hidden_state
500             cls_embeddings =
                last_hidden_state[:, 0, :]
                cpu().numpy() # Take [CLS]
                token embedding 541
501             all_embeddings.append( 542
                cls_embeddings) 543

```

```

    return np.vstack(all_embeddings)
# Generate sentence embeddings
X_train_bert = get_bert_embeddings(
    clean_train['clean_text'].tolist(),
    tokenizer, bert_model)
X_valid_bert = get_bert_embeddings(
    clean_valid['clean_text'].tolist(),
    tokenizer, bert_model)
X_test_bert = get_bert_embeddings(
    clean_test['clean_text'].tolist(),
    tokenizer, bert_model)

y_train = clean_train['label']
y_valid = clean_valid['label']
y_test = clean_test['label']

from sklearn.svm import SVC

# Initialize a SVM classifier
clf_svm = SVC(kernel='linear', C=1,
    random_state=42)
clf_svm.fit(X_train_bert, y_train)

# Predict labels for the test set using
    the trained SVM
y_test_pred_svm = clf_svm.predict(
    X_test_bert)

# Print the classification report (
    precision, recall, f1-score) for each
    class
from sklearn.metrics import
    classification_report
print("\n=== Test Set Performance (BERT
    Embeddings + SVM) ===")
print(classification_report(y_test,
    y_test_pred_svm, digits=2))

# Loading data & preprocessing
import pandas as pd
from datasets import Dataset

# Loading the cleaned Reddit dataset
train_df = pd.read_csv("clean_train.csv")
valid_df = pd.read_csv("clean_valid.csv")
test_df = pd.read_csv("clean_test.csv")

# Convert to Hugging Face Dataset format
train_dataset = Dataset.from_pandas(
    train_df)
valid_dataset = Dataset.from_pandas(
    valid_df)
test_dataset = Dataset.from_pandas(test_df
    )

# Load the pre-trained tokenizer for the
    mental-roberta-base model

```

```

544 from transformers import AutoTokenizer
545 # Initialize the tokenizer
546 tokenizer = AutoTokenizer.from_pretrained(
547     "mental/mental-roberta-base")
548 # Define a tokenization function for the dataset
549 def tokenize(batch):
550     """
551     Tokenize each example in the dataset
552     batch.
553     - batch["clean_text"]: list of text
554       samples to be tokenized
555     - truncation: cuts text to max_length
556       if needed
557     - padding: ensures all sequences in
558       the batch have the same length (
559       max_length)
560     - max_length: 128 tokens (fixed for
561       all samples)
562     """
563     return tokenizer(batch["clean_text"],
564                       truncation=True, padding="
565                       max_length", max_length=128)
566 # Apply tokenization to the training,
567 validation, and test sets
568 train_dataset = train_dataset.map(tokenize,
569                                   batched=True)
570 valid_dataset = valid_dataset.map(tokenize,
571                                   batched=True)
572 test_dataset = test_dataset.map(tokenize,
573                                 batched=True)
574 # Remove the original text column as it's
575 no longer needed after tokenization
576 train_dataset = train_dataset.
577 remove_columns(["clean_text"])
578 valid_dataset = valid_dataset.
579 remove_columns(["clean_text"])
580 test_dataset = test_dataset.remove_columns
581 (["clean_text"])
582 # Rename the label column to 'labels' that
583 required by Hugging Face Trainer API
584 train_dataset = train_dataset.
585 rename_column("label", "labels")
586 valid_dataset = valid_dataset.
587 rename_column("label", "labels")
588 test_dataset = test_dataset.rename_column(
589     "label", "labels")
590 # Set dataset format to PyTorch tensors
591 for compatibility with Trainer
592 train_dataset.set_format("torch")
593 valid_dataset.set_format("torch")
594 test_dataset.set_format("torch")
595 # Load the pre-trained mental-roberta-base
596 model with a classification head
597
598 from transformers import
599     AutoModelForSequenceClassification
600 # Load the pre-trained RoBERTa-based model
601 .
602 # Since doing binary classification set
603 num_labels=2.
604 model = AutoModelForSequenceClassification
605 .from_pretrained("mental/mental-roberta
606 -base", num_labels=2)
607 # Set up training arguments
608 from transformers import TrainingArguments
609 , Trainer
610 # Initialize TrainingArguments with basic
611 training configuration.
612 training_args = TrainingArguments(
613     output_dir="./results_mental_roberta",
614     # Directory to save model
615     checkpoints and logs
616     do_train=True,
617     # Enable
618     training
619     do_eval=True,
620     # Enable
621     evaluation during training
622     learning_rate=2e-5,
623     # Common
624     learning rate for fine-tuning
625     transformers
626     per_device_train_batch_size=16,
627     # Batch size for training
628     per_device_eval_batch_size=16,
629     # Batch size for
630     evaluation
631     num_train_epochs=3,
632     # Number of
633     training epochs
634     weight_decay=0.01,
635     # Weight decay
636     for regularization
637     logging_dir="./logs"
638     # Directory to
639     store logs
640 )
641 # Define the evaluation metrics
642 import numpy as np
643 from sklearn.metrics import accuracy_score
644 , precision_recall_fscore_support
645
646 def compute_metrics(eval_pred):
647     """
648     Compute standard classification
649     metrics (accuracy, precision,
650     recall, f1).
651
652     Args:
653     - eval_pred: a tuple (logits, labels)

```

```

        from the Trainer
612
613 Returns:
614 - A dictionary containing the metrics
        for logging and evaluation.
615 """
616 logits, labels = eval_pred
617 # Convert logits to predicted labels
618 predictions = np.argmax(logits, axis
        =-1)
619 # Compute precision, recall, f1-score
        with macro averaging (equal weight
        for each class)
620 precision, recall, f1, _ =
        precision_recall_fscore_support(
        labels, predictions, average="macro
        ")
621 # Compute overall accuracy
622 acc = accuracy_score(labels,
        predictions)
623 # Return metrics in a dictionary
624 return {"accuracy": acc, "precision":
        precision, "recall": recall, "f1":
        f1}
625 # Create a Trainer instance and start fine
        -tuning
626 trainer = Trainer(
627     model=model,
        #
        The model to fine-tune
628     args=training_args,
        #
        Training parameters
629     train_dataset=train_dataset,
        #
        Dataset for training
630     eval_dataset=valid_dataset,
        #
        Dataset for validation
631     tokenizer=tokenizer,
        #
        Tokenizer (needed for saving &
        loading pipeline)
632     compute_metrics=compute_metrics
        #
        Function to compute evaluation
        metrics
633 )
634
635 # Start the training process
636 trainer.train()
637
638 # Evaluate the model on the test dataset
639 print("\n== Test Set Evaluation (mental-
        roberta-base) ==")
640 # Evaluate the model on the held-out test
        dataset
641 metrics = trainer.evaluate(test_dataset)
642 print(metrics)
643
644
645 # Generate a classification report on the
        test dataset
646 from sklearn.metrics import
        classification_report
647
        # Predict the labels for the test dataset
648 preds_output = trainer.predict(
649     test_dataset)
        y_true = preds_output.label_ids
        y_pred = np.argmax(preds_output.
        predictions, axis=1)
650
651 # Print detailed classification report
652 print("\n== Classification Report on Test
        Set (mental-roberta-base) ==")
653 print(classification_report(y_true, y_pred
        , digits=2))

```