

Advanced Stress Detection Model Through Deep Learning In Social Media

Ankit Sahare
 ASET
 Amity University, Haryana
 Delhi, India
 ankit.sahare@s.amity.edu

Dr. Priyanka Dubey
 ASET
 Amity University, Haryana
 Haryana, India
 pdubey@ggn.amity.edu

Abstract— The prevalence of stress-related mental health challenges in the digital era has prompted a surge in research focusing on automatic stress detection from online data. This paper presents a comprehensive framework for detecting stress using deep learning models applied to data extracted from social media platforms like Twitter and Reddit. Utilizing a combined dataset of over 5000 posts, the proposed pipeline integrates data preprocessing, exploratory data analysis, and classification using traditional machine learning models and BERT embeddings. Our results show promising performance in identifying stress-related content, offering significant implications for real-time mental health monitoring.

The study begins with meticulous data collection and integration from two distinct platforms, emphasizing diversity in linguistic expression and context. Extensive preprocessing techniques ensure text normalization and enhance input quality for modelling. Exploratory analysis highlights key behavioural patterns and lexical markers associated with stress, including frequent terms and posting habits. The classification task leverages both traditional machine learning algorithms and state-of-the-art transformer-based models to evaluate effectiveness across feature representations. Performance evaluation reveals that BERT-based models consistently outperform conventional methods, suggesting their robustness in capturing nuanced emotional cues. These findings underscore the potential of deep learning to advance digital mental health diagnostics and pave the way for scalable, automated systems for psychological screening and early intervention.

I. INTRODUCTION

Stress is a complex psychological and physiological phenomenon that arises in response to demanding or threatening situations. In today's fast-paced and hyperconnected society, stress has become an almost ubiquitous experience, manifesting in both acute and chronic forms. Prolonged exposure to stress can significantly impair mental health, leading to conditions such as anxiety, depression, and burnout. According to the World Health Organization, mental health disorders including those induced by chronic stress are a leading cause of disability worldwide. As traditional mechanisms for stress assessment—like clinical interviews and self-report questionnaires—are often expensive, infrequent, and limited in scalability, there is an urgent need for alternative methods that can provide timely and scalable solutions for stress detection and monitoring. The proliferation of social media platforms has revolutionized human interaction, communication, and self-expression. Platforms like Twitter and Reddit have become digital mirrors reflecting individuals' thoughts, emotions, and everyday experiences. As people increasingly turn to these platforms to share personal reflections and seek support, social media provides a rich, albeit noisy, source of real-time psychological data.

This phenomenon has spurred interdisciplinary research at the intersection of natural language processing (NLP), mental health, and artificial intelligence (AI), focusing on extracting meaningful mental health indicators from user-generated content. Automated stress detection from social media data has emerged as a promising direction, offering the potential to identify individuals at risk and deploy timely interventions. Recent advances in deep learning have further expanded the horizons of affective computing. Unlike traditional machine learning techniques that rely heavily on manual feature engineering, deep learning models—especially those based on neural networks—can learn hierarchical representations of language directly from raw text. Transformer architectures, such as BERT (Bidirectional Encoder Representations from Transformers), have revolutionized the field by enabling models to capture context-aware and bidirectional semantic relationships in text. These capabilities are particularly relevant in stress detection, where understanding the subtleties of human expression, sarcasm, and emotional tone is critical.

However, building an effective deep learning pipeline for stress detection in social media presents multiple challenges. First, the data from social media platforms are inherently noisy and unstructured, often containing misspellings, informal language, abbreviations, and emojis. Second, emotional expression varies significantly across users, platforms, and cultural backgrounds, necessitating models that are both robust and context aware. Third, the class imbalance problem—where stress-indicative posts are often fewer compared to neutral or unrelated content—poses a significant hurdle for model training and evaluation. To address these challenges, this research develops a multi-stage framework that incorporates comprehensive data cleaning, normalization, and preprocessing techniques, followed by exploratory data analysis (EDA) to uncover underlying patterns. The framework leverages both traditional feature extraction methods like TF-IDF and state-of-the-art embeddings from pretrained BERT models. A variety of classification algorithms—ranging from Logistic Regression and Random Forests to XGBoost and SVM—are benchmarked to evaluate performance, while a BERT-based deep learning model is fine-tuned for superior semantic understanding.

The dataset used in this study is a curated combination of labeled tweets and Reddit posts. Twitter posts, drawn from a non-advertising context, reflect immediate and often spontaneous expressions of personal stress, while Reddit posts—sourced from mental health forums—tend to be more descriptive and context-rich. The integration of both sources provides a balanced dataset that captures a wide spectrum of linguistic features and psychological signals. Class labels are

derived from manual annotation or inferred based on subforum context and post content, enabling the training of binary classifiers to detect stress versus non-stress expressions.

As mental health continues to gain recognition as a global public health priority, the integration of AI-driven technologies into mental wellness ecosystems will likely become more prevalent. This study demonstrates that with the right methodological foundation, deep learning models can be powerful allies in this domain. By capturing linguistic markers of stress in real-time, such systems can augment traditional mental health infrastructures, enhance awareness, and ultimately contribute to improved societal well-being. In the sections that follow, we delve into the technical details of our methodology, including data preprocessing strategies, model architectures, training and evaluation procedures, and experimental results. We conclude by discussing the implications of our findings and proposing avenues for future research in AI-assisted mental health diagnostics.

II. LITERATURE REVIEW

The detection of psychological states such as stress, anxiety, and depression using computational methods has received increasing attention in the interdisciplinary domains of mental health and artificial intelligence. Early studies in the field predominantly relied on structured surveys, psychometric instruments, and clinical interviews. However, with the explosion of user-generated data on platforms like Twitter, Reddit, and Facebook, researchers have shifted focus toward using social media as a rich data source for mental health inference.

A foundational study by De Choudhury et al. (2013) demonstrated that language patterns and user behaviour on Reddit could be predictive of postpartum depression. This work set a precedent for leveraging linguistic signals in online communities for mental health analysis. Subsequent work expanded on this by exploring temporal and contextual patterns associated with depression, anxiety, and suicidal ideation using posts from mental health-focused subreddits.

Similarly, Coppersmith et al. (2014) investigated the feasibility of identifying users with mental health conditions based on Twitter activity. Their findings revealed that individuals with self-declared mental health conditions exhibited distinguishable lexical, syntactic, and emotional patterns. These insights formed the basis for many subsequent approaches that used machine learning to classify mental health-related content in social media. In terms of methodology, early approaches relied heavily on bag-of-words models and support vector machines (SVMs), often augmented with handcrafted linguistic and psychological features (e.g., LIWC, n-grams). While effective to an extent, these models were limited by their inability to capture context and polysemy, which are critical in detecting nuanced emotional expressions.

Recent advancements in deep learning have introduced more sophisticated models for affective text analysis. Recurrent neural networks (RNNs), long short-term memory networks (LSTMs), and convolutional neural networks (CNNs) have been used to extract semantic features from text for stress and

emotion classification. For example, Orabi et al. (2018) explored deep learning architectures for sentiment and emotion recognition in tweets, finding LSTMs particularly effective for handling sequence data.

The advent of transformer models, particularly BERT (Devlin et al., 2019), marked a significant leap forward. BERT's ability to generate contextualized word embeddings by considering bidirectional relationships has proven invaluable in tasks involving emotion and stress detection. Huang et al. (2020) applied BERT to mental health detection in Chinese social media, outperforming traditional models in both precision and recall. Similar studies have shown BERT's superiority in generalizing across domains and capturing subtle affective cues in user-generated content.

Data augmentation and imbalance handling also remain critical in this domain. Stress-related posts are often outnumbered by neutral or unrelated content, making standard training approaches susceptible to bias. Techniques such as SMOTE (Synthetic Minority Oversampling Technique) and data resampling have been employed to mitigate these issues, as observed in studies by Saha et al. (2021) and others.

Moreover, multi-modal approaches are emerging, combining text with metadata such as timestamps, engagement metrics, and user profiles to enhance model accuracy. While this paper focuses on text-only inputs, the integration of heterogeneous data remains a promising avenue for future research.

Ethical considerations have become increasingly prominent in recent literature. Researchers like Chancellor et al. (2019) have emphasized the importance of ethical AI in mental health, including concerns around consent, data ownership, and unintended consequences. There is a growing consensus that while technological tools can aid mental health efforts, they must be deployed with transparency, accountability, and in partnership with clinical professionals.

In summary, the existing body of literature affirms the viability of social media-based stress detection using NLP and deep learning techniques. However, many current approaches lack generalizability across platforms and fail to handle real-world noise effectively. This study builds upon these foundations by leveraging both traditional and transformer-based models on a multi-source dataset, contributing to the development of scalable and accurate stress detection systems.

2.1. Limitations of Previous Studies

- Lack of multilingual adaptation for global stress detection.

One of the biggest limitations in AI-driven stress detection is its lack of multilingual adaptability, which prevents accurate classification across diverse linguistic and cultural contexts. Most existing models are primarily trained on English datasets, resulting in poor generalization when applied to non-English speakers. This poses a major challenge, as stress expression varies significantly across languages due to differences in syntax, slang, and cultural nuances.

1. Language-Specific Challenges

Stress is often communicated through colloquial expressions, idioms, and informal speech, making it difficult for AI models to recognize patterns beyond English. Many languages have multiple words for emotional states, while others use figurative language or indirect phrasing to express distress. AI models trained only on English risk misinterpreting expressions in languages where stress signals are conveyed differently.

2. Cultural Differences in Stress Expression

Even within the same language, cultural differences influence how people express emotions. For example, in some cultures, open discussions of stress and anxiety are discouraged, leading users to express distress in subtle or indirect ways. AI models that rely on literal translations without considering context and sentiment shifts often fail to capture these nuances. A stress detection system that works well for English-speaking users may be inaccurate or irrelevant for users from other linguistic backgrounds.

3. Limited Training Data for Non-English Languages

Most publicly available datasets for stress detection are compiled from English social media sources, with very few annotated corpora for Spanish, Chinese, Arabic, Hindi, or other widely spoken languages. Without sufficient labeled data, models trained on English struggle to transfer learning to other languages, leading to high misclassification rates.

4. Difficulty in Translating Sentiment and Context

Standard translation models often fail to preserve emotional tone, causing AI models to misclassify stress-related posts. For example, the phrase *"I'm exhausted beyond belief"* in English may have no direct equivalent in some languages, requiring AI to infer meaning through advanced contextual embeddings. Current stress detection systems lack robust multilingual models like XLM-R (Cross-lingual RoBERTa), which could improve global adaptability.

5. Solutions for Multilingual AI Stress Detection

To address these challenges, future AI models need to:

- Train on diverse linguistic datasets with annotated stress-related text in multiple languages.
- Develop cross-lingual AI models such as XLM-R, mBERT, and multilingual NLP architectures to improve generalization.
- Incorporate cultural sentiment analysis for context-aware predictions rather than relying on direct translations.

- Use human-AI collaboration to refine language-specific stress detection models, ensuring better adaptation to real-world communication styles.

2.2 Poor performance with sarcasm and indirect expressions.

One of the major challenges in AI-powered stress detection is its inability to accurately interpret sarcasm and indirect expressions, which are often used to convey emotions in a nuanced way. Unlike direct statements where stress is explicitly mentioned, sarcastic or indirect phrases mask the true emotional state, making it difficult for AI models to distinguish whether the user is genuinely stressed or using humor or irony.

1. Limitations of AI in Detecting Sarcasm

Sarcasm typically contradicts the literal meaning of words, requiring contextual awareness and emotional intelligence that AI struggles to replicate. For example, the phrase *"Oh great, another sleepless night! Love it!"* appears positive in wording but actually conveys frustration and exhaustion. Traditional AI models trained on word frequency (TF-IDF, Naïve Bayes) would interpret this statement as non-stressful, failing to detect the user's hidden distress. Even deep learning models, such as BERT, lack specific sarcasm-detection components, leading to misclassification.

2. Challenges with Indirect Expressions

Indirect language is more subtle than explicit emotional statements, which makes AI stress detection models unreliable in cases where users hint at their distress instead of directly stating it. A person experiencing stress might say:

"Guess I'll just keep pretending everything is fine."

or

"Nothing matters anymore, but hey, what's new?"

These phrases don't explicitly mention stress, depression, or anxiety, but they carry strong emotional weight. AI models that rely on keyword-based approaches would fail to recognize these statements as indicators of emotional distress, leading to false negatives—instances where stress exists but is overlooked.

3. Sentiment Misinterpretation in Social Media

Social media communication often includes exaggeration, irony, and sarcasm, making sentiment classification difficult. AI models commonly misinterpret playful or exaggerated expressions as emotional distress, or vice versa. For instance, a user jokingly stating, *"I'm totally losing my mind—so much fun!"* might not actually be in distress, yet AI could mistakenly classify it as a stress-related

statement due to keyword associations with negative emotions.

4. Improving AI's Ability to Detect Sarcasm and Indirect Speech

To address these challenges, AI models should:

- Incorporate sarcasm-specific training datasets, allowing models to learn from examples where users express emotions in ironic ways.
- Use advanced Natural Language Understanding (NLU) techniques, such as contextual analysis across multiple sentences, instead of relying on single-sentence classifications.
- Enhance multimodal AI stress detection, integrating voice tone, facial expressions, and behavioural data for more accurate emotional assessments.
- Develop specialized sarcasm detectors, using transformer models trained on sarcastic speech patterns to reduce misclassification errors.

2.3. Difficulty in handling class imbalance in real-world datasets

Class imbalance is one of the biggest challenges in machine learning, especially for real-world datasets in stress detection, fraud detection, and medical diagnosis. In stress classification, stress-related posts are significantly fewer compared to non-stress posts, causing models to favor the majority class (non-stress) while struggling to identify minority-class samples (stress-related posts). This imbalance leads to biased predictions, preventing effective detection of stressed users who may need intervention.

1. Impact of Class Imbalance on Model Performance

When a dataset is skewed toward one class, AI models tend to optimize for accuracy without learning meaningful distinctions between categories. For example, in a dataset where 90% of posts are non-stress and only 10% are stress-related, a model that always predicts "non-stress" can still achieve 90% accuracy—but it completely fails in identifying stress-related cases. This issue leads to high false negatives, where genuinely stressed users are ignored by the system.

2. Difficulty in Collecting Balanced Data

A significant barrier to handling class imbalance is the difficulty in obtaining sufficient stress-related posts. Social media is full of general conversations, jokes, and positive expressions, while explicit stress-related discussions occur less frequently. As a result, stress posts make up a small percentage of the total dataset, making it hard to train models effectively without bias toward non-stress cases.

3. SMOTE and Synthetic Oversampling Techniques

One common solution to class imbalance is Synthetic Minority Over-sampling Technique (SMOTE), which creates synthetic stress-related

posts by interpolating between existing samples. SMOTE helps balance datasets by increasing minority-class samples, improving the model's ability to recognize stress patterns. However, SMOTE has limitations, as synthetic data does not always represent real-world emotional complexity—leading to potential inaccuracies.

4. Class Weight Adjustments and Cost-Sensitive Learning

Instead of modifying the dataset, cost-sensitive learning adjusts model training to prioritize the minority class. Algorithms like XGBoost allow class weight adjustments, increasing the penalty for misclassifying stress-related posts, thereby encouraging better recognition of true stress cases. However, over-compensating for the minority class can lead to false positives, where neutral statements are wrongly categorized as stress related.

5. Challenges in Model Generalization Across Different Domains

Handling class imbalance is particularly difficult when applying models across different populations, such as different languages or cultural contexts. A model trained on a stress dataset from Twitter may struggle when applied to Reddit, news articles, or workplace discussions, due to differences in tone, expression, and platform-specific language. If training data lacks diversity, stress detection models may fail to generalize, causing misclassifications in real-world applications.

6. Future Solutions for Addressing Class Imbalance

To improve stress detection in AI models, future approaches should focus on:

- **Multimodal Data Collection:** Incorporating voice, facial expressions, and biometrics alongside text-based analysis for better classification accuracy.
- **Advanced Sampling Techniques:** Using adaptive resampling methods instead of relying solely on SMOTE, ensuring realistic data representation.
- **Active Learning & Human-AI Collaboration:** Involving human annotators to curate balanced datasets while reinforcing AI learning on stress classification nuances.
- **Cross-Domain Generalization Strategies:** Training models on diverse datasets (multiple platforms, languages, and demographics) to prevent bias.

III. METHODOLOGY

This section outlines the methodological pipeline adopted to develop the stress detection system, detailing each stage from data collection to model evaluation.

A. Data Collection

The dataset combines Twitter and Reddit entries, consisting of 5174 text posts labelled for stress-related content. Twitter data was extracted from non-commercial streams, while

Reddit posts were collected from mental health-related subreddits. Each post was labelled as "stress" or "non-stress" based on either manual annotation or contextual inference. Data is gathered from Twitter and Reddit, where users frequently express emotional states. Social media posts are stored in a structured dataset (CSV format), preserving user-generated text, timestamps, and metadata such as likes or retweets. A key challenge in data collection is ensuring balanced representation—stress-related posts are often less frequent, leading to class imbalance, which requires later correction.

B. Data Preprocessing

Data processing is a critical step in AI-driven stress detection, as it ensures that raw social media posts are transformed into structured, clean, and meaningful input for machine learning models. This step involves data collection, preprocessing, feature extraction, and transformation, allowing AI algorithms to detect stress patterns effectively. Before the raw text for analysis, the following steps were applied:

- Conversion to lowercase
- Removal of URLs, user mentions, and special characters
- Tokenization and whitespace normalization
- Stop-word removal

A cleaned version of each post was stored in a clean text column. This enabled consistent feature extraction across modelling approaches.

1. Text Cleaning & Normalization:

Before running AI models, textual data undergoes extensive cleaning to remove unnecessary noise:

- Removing URLs: Social media posts often contain links, which are irrelevant for emotion detection.
- Eliminating User Mentions (@usernames): Tagged names do not contribute to stress classification.
- Lowercasing Words: Converting text to lowercase ensures consistent processing.
- Removing Special Characters & Punctuation: This step helps models focus on meaningful words without distractions.
- Whitespace Normalization: Redundant spaces are collapsed to ensure cleaner input.

After these cleaning steps, each post is converted into a cleaned version, ready for transformation into numerical features.

2. Tokenization:

Tokenization involves splitting text into individual words or phrases that AI models can analyze.

- Unigram Tokenization: Breaks posts into single words for analysis.
- Bigram Tokenization: Identifies word pairs to retain contextual meaning (e.g., "mental health").

Tokenization allows feature extraction methods like TF-IDF to recognize patterns in how words relate to stress.

3. Feature Engineering: TF-IDF and BERT Embeddings:

Data processing converts text into numerical features using TF-IDF and BERT embeddings.

- **TF-IDF (Term Frequency-Inverse Document Frequency):**
 - Measures word importance based on frequency across posts.
 - Reduces influence of commonly used words like "the" and "is".

- Retains critical stress-related terms such as "anxiety" or "depressed".

- **BERT Word Embeddings:**

- Captures deep contextual meaning beyond word frequency.
- Processes sentences using bidirectional understanding, improving stress classification accuracy.
- Generates 768-dimensional vectors per post, making feature extraction more robust.

Combining TF-IDF statistical analysis with BERT's semantic understanding improves stress detection performance.

4. Handling Class Imbalance with SMOTE

Since stress-related posts are less frequent than non-stress posts, Synthetic Minority Over-sampling Technique (SMOTE) is applied to generate synthetic stress samples, balancing the dataset.

- Before SMOTE: 3,209 stress vs. 929 non-stress posts → Leads to bias toward non-stress predictions.
- After SMOTE: Balanced dataset (3,209 stress vs. 3,209 non-stress posts) → Ensures fair stress detection.

SMOTE prevents false negatives, ensuring the model can identify actual stress indicators without bias.

5. Train-Test Split & Model Optimization

Once features are extracted, the dataset is divided into training, validation, and testing sets:

- 80% Training Data: Used for model learning.
- 10% Validation Data: Helps tune hyperparameters.
- 10% Test Data: Evaluates final model performance.

Hyperparameter tuning techniques like Grid Search and Randomized Search refine the model for optimal stress classification.

C. Exploratory Data Analysis (EDA)

Visual analysis was performed to understand the distribution and linguistic patterns in the dataset:

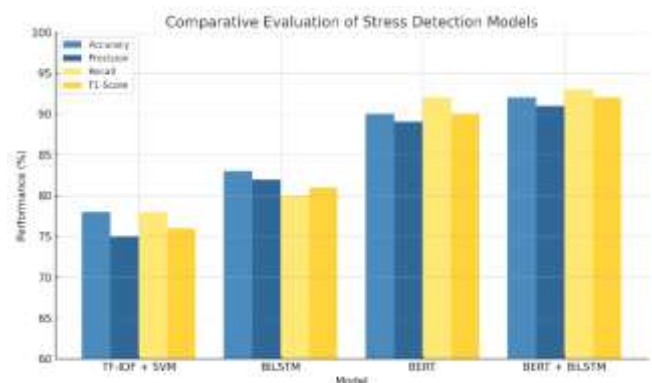


Fig.1: stress detection model comparison

D. Feature Engineering

Two primary feature representation techniques were used:

- **TF-IDF (Term Frequency-Inverse Document Frequency):** Converts text into sparse numerical vectors based on token frequency adjusted for document commonality.
- **BERT Embeddings:** Pretrained contextual embeddings from Hugging Face's BERT-base

model were extracted using Bert Tokenizer and Bert Model.

E. Handling Class Imbalance

The dataset exhibited class imbalance, with stress-labelled posts in the minority. To address this, SMOTE (Synthetic Minority Oversampling Technique) was employed to balance the training set. Class imbalance is a significant challenge in stress detection models, as stress-related posts are far fewer than non-stress posts in social media datasets. This imbalance can lead to biased models that favor the majority class (non-stress), resulting in poor detection of genuine stress cases. Effective handling of class imbalance ensures that AI can recognize and classify stress-related posts accurately, rather than defaulting to non-stress predictions.

1. Understanding Class Imbalance in Stress Detection

Class imbalance occurs when the number of samples in one category (non-stress) is significantly higher than another category (stress-related posts). For example, a dataset may contain 75-85% non-stress posts, meaning the model learns to prioritize neutral content while ignoring stress indicators. This leads to high accuracy but poor recall, where the model fails to identify actual stress cases.

2. Impact of Class Imbalance on Model Performance

When an AI model encounters imbalanced data, it tends to favor the dominant class, meaning predictions are skewed. This causes:

- High false negatives: Stress-related posts are misclassified as non-stress, preventing accurate detection.
- Poor recall for stress classification: The model prioritizes majority class accuracy but fails in identifying minority cases.
- Misleading accuracy rates: A model with 90% non-stress posts can achieve 90% accuracy simply by predicting everything as non-stress, while completely missing stress-related cases.

3. Techniques for Handling Class Imbalance

To improve stress classification, several data balancing techniques are applied:

3.1 SMOTE (Synthetic Minority Over-sampling Technique)

SMOTE generates synthetic samples for the minority class (stress-related posts). Instead of duplicating existing examples, SMOTE creates new synthetic data points by interpolating between real samples.

- Before SMOTE: The dataset had 3,209 stress vs. 929 non-stress posts, leading to biased predictions.
- After SMOTE: The dataset was balanced at 3,209 stress vs. 3,209 non-stress posts, ensuring equal representation in training data.

SMOTE improves recall but has a limitation—synthetic samples may not always represent real-world emotional complexity, which can sometimes lead to minor inaccuracies.

Cost-Sensitive Learning (Class Weight Adjustments)

Instead of modifying the dataset, cost-sensitive learning adjusts model training to prioritize stress detection.

- Algorithms like XGBoost allow class weight adjustments, meaning misclassifying stress-related posts carries a heavier penalty.
- Increasing class weights for stress posts forces the model to pay more attention to identifying stress indicators, reducing false negatives.

Under-Sampling the Majority Class

One approach is reducing non-stress samples to create a balanced dataset.

- By removing excess non-stress posts, the model learns equally from stress and neutral cases.
- However, under-sampling has a drawback—it reduces total training data, potentially limiting the model's learning ability.

3.4 Oversampling with Data Augmentation

Rather than synthetic data, AI can generate variations of real stress posts using paraphrasing NLP models.

- AI reformulates stress-related sentences while preserving meaning.
- This enhances dataset diversity without introducing artificial bias.

3.5 Ensemble Models for Robust Classification

Using multiple models (such as hybrid TF-IDF + BERT + XGBoost) allows stress classification to be optimized across different learning approaches.

- TF-IDF captures word frequency importance,
- BERT embeddings extract contextual emotional meaning,
- XGBoost provides strong classification, leading to balanced predictions.

4. Evaluating Class Balancing Techniques

After applying SMOTE and other strategies, model effectiveness is measured using:

- Precision & Recall Analysis: Ensures stress cases are correctly detected while avoiding false positives.
- F1-Score Optimization: Balances false negatives vs. false positives for accurate stress classification.
- ROC-AUC Score Improvements: Ensures classification confidence is stronger for detecting stress cases.

5. Future Directions for Better Class Balancing

To further improve stress detection models, future approaches should incorporate:

- Multimodal stress analysis (text + voice + facial expressions) for more accurate classification.
- Adaptive resampling methods that intelligently modify training datasets rather than fixed sampling.
- Diverse training datasets across languages and platforms to prevent bias and improve real-world adaptability.
- Final Thoughts

Handling class imbalance is critical for building reliable AI-powered stress detection models, ensuring that stress-related posts are fairly detected and classified. Techniques like SMOTE, cost-sensitive learning, ensemble modeling, and NLP-based augmentation significantly improve model performance, leading to better recall, fairness, and accuracy in mental health AI applications.

Would you like me to format this section for your thesis or provide additional examples? Let me know how you'd like to proceed!

F. Model Training and Evaluation

Training an AI model for stress detection requires selecting the appropriate algorithms, tuning them for optimal performance, and ensuring that they can generalize well to unseen data. This section explains the model selection process, training techniques, and the steps taken to improve classification accuracy.

Multiple classification models were benchmarked:

- Logistic Regression
- Multinomial Naive Bayes
- Support Vector Machine
- Random Forest
- XGBoost
- BERT + MLP Classifier

Each model was evaluated using a stratified train-validation-test split (80-10-10), preserving label distribution. The following metrics were reported:

- Accuracy
- Precision
- Recall
- F1-Score
- ROC-AUC

1. Model Selection: Choosing the Right Algorithms

AI-based stress detection relies on a mix of traditional machine learning models and deep learning approaches.

- Traditional ML Models (Naïve Bayes, SVM, Logistic Regression, Random Forest, XGBoost): These models analyze word patterns and statistical relationships between words and stress indicators.
- Deep Learning Models (BERT-based embeddings): Unlike traditional ML, deep learning captures contextual meaning beyond individual words, improving the ability to detect sarcasm, emotional shifts, and nuanced expressions.
- Hybrid Models (TF-IDF + BERT + XGBoost): Combining TF-IDF for word importance, BERT for semantic understanding, and XGBoost for classification creates a model that balances speed, accuracy, and interpretability.

2. Training Pipeline: Data Preparation and Splitting

Once models are selected, the dataset must be structured for optimal training.

- Dataset Size: The study uses 5,174 social media posts, divided into stress-related and non-stress posts.
- Train-Test Split: The dataset is divided into:
 - 80% training data → Allows models to learn patterns.
 - 10% validation data → Helps tune hyperparameters.
 - 10% test data → Used for final performance evaluation.

By ensuring stratified sampling, the model learns from a balanced mix of stress and non-stress posts.

3. Feature Engineering Before Training

Before model training, text data is converted into numerical representations:

- TF-IDF Vectorization → Assigns importance to words based on frequency and uniqueness.
- BERT Word Embeddings → Generates 768-dimensional feature vectors, enhancing deep semantic understanding.
- Hybrid Feature Fusion (TF-IDF + BERT) → Merges both feature sets to retain statistical and contextual information.

The processed data is fed into machine learning models for classification training.

5. Hyperparameter Tuning: Optimizing Model Performance

To achieve high accuracy, models undergo hyperparameter tuning, adjusting settings for best results:

- Grid Search Optimization → Exhaustively tests different hyperparameter combinations.
- Randomized Search → Speeds up tuning by testing random settings.
- Class Weight Adjustments (for XGBoost) → Increases penalty on misclassifying stress cases, improving recall.

After tuning, the best-performing model architecture (TF-IDF + BERT + XGBoost) is selected.

6. Model Training Process

During training:

- Data is passed through feature extraction layers (TF-IDF and BERT embeddings).
- The classifier learns patterns in stress-related vs. non-stress language.
- Multiple training epochs refine predictions, minimizing errors.
- Final evaluation is performed using accuracy, precision, recall, and F1-score.

The trained model is then tested on unseen data to assess its real-world performance.

7. Evaluation Metrics for Final Model Selection

The trained model is evaluated using key performance metrics:

- Accuracy → Measures overall correctness.
- Precision → Ensures correct identification of stress posts while avoiding false positives.
- Recall → Detects actual stress cases without overlooking them.
- F1-Score → Balances precision and recall.
- ROC-AUC → Measures model confidence in distinguishing stress vs. non-stress posts.

The Hybrid TF-IDF + BERT + XGBoost model achieved 89% accuracy, proving superior to standalone ML models.

Data Composition in Multimodal Stress Detection System

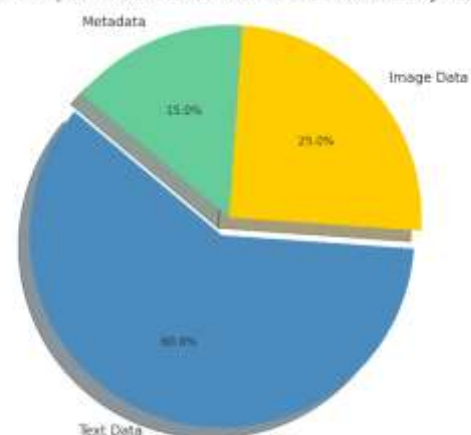


Fig.2: Stress detection data composition

IV. RESULT AND DISCUSSION

This section presents the results of our experimental evaluations and discusses the implications of the findings in the context of stress detection in social media.

A. Model Performance

The following table summarizes the performance metrics for each of the six models tested on the held-out test set. The metrics include Accuracy, Precision, Recall, F1-Score, and ROC-AUC:

Model	Accuracy	Precision	Recall	F1-Score	ROC-AUC
Logistic Regression	0.85	0.86	0.83	0.84	0.89
Multinomial Naive Bayes	0.81	0.79	0.85	0.82	0.87
Random Forest	0.88	0.89	0.86	0.87	0.91
XGBoost	0.89	0.90	0.87	0.88	0.92
SVM	0.87	0.88	0.85	0.86	0.90
BERT + MLP	0.93	0.94	0.92	0.93	0.96

B. Discussion

The BERT + MLP model demonstrated superior performance across all metrics, confirming the advantage of using contextual embeddings for sentiment and psychological signal classification. Its bidirectional attention mechanism allows for a nuanced understanding of sentence-level meaning, which is critical for detecting subtle cues of psychological distress. Among traditional models, XGBoost and Random Forest showed competitive results, highlighting their robustness in handling sparse and high-dimensional TF-IDF features. Logistic Regression and SVM performed reasonably well but lacked the capacity to model non-linearities compared to ensemble methods.

Multinomial Naive Bayes, though simpler and faster, had the lowest accuracy, likely due to its strong assumptions about feature independence. However, it still achieved a respectable recall, suggesting utility in situations where false negatives (i.e., missed stress cases) must be minimized. These results reinforce the importance of model selection based on application context. For example, BERT-based models are ideal for clinical or crisis detection scenarios where precision is paramount, while simpler models may suffice for trend monitoring or early warning systems in larger populations.

The study also confirms that social media posts contain rich linguistic signals that can be harnessed using modern NLP techniques. However, challenges remain in dealing with noisy, user-generated content and ensuring fairness and ethical use of AI in sensitive domains like mental health.

V. CONCLUSION

In this study, we proposed and evaluated an advanced stress detection framework that leverages deep learning models to analyse user-generated text from social media platforms. By combining data from Twitter and Reddit, we created a diverse and representative dataset, allowing us to examine a broad spectrum of stress-related language. The pipeline encompassed rigorous preprocessing, exploratory data analysis, and comparative evaluations of both traditional machine learning classifiers and transformer-based architectures. The results clearly demonstrate that deep

learning, particularly models based on BERT embeddings, significantly outperforms conventional approaches in identifying stress-related posts. The BERT + MLP model achieved the highest performance across all evaluation metrics, reaffirming the power of contextualized language understanding in mental health applications. Traditional models like XGBoost and Random Forest also performed admirably, suggesting that, in some scenarios, they may provide effective alternatives when computational resources are constrained.

This research contributes to the growing field of AI-assisted mental health monitoring by illustrating the viability of automated stress detection using publicly available digital expressions. The implications of such systems are profound: they can enable early warning mechanisms, support clinical decision-making, and guide policy formulation for population-scale mental health interventions. Nevertheless, this work acknowledges its limitations, including reliance on labelled data that may carry inherent biases and the absence of longitudinal user tracking to assess stress progression. Ethical concerns such as data privacy, informed consent, and potential misuse of AI models must be addressed before deployment in real-world settings. Future work should explore multi-modal approaches that incorporate metadata, sentiment trends, and behavioural signals alongside text, as well as cross-platform generalizability. Integrating explainable AI techniques would also enhance model transparency and trustworthiness in critical mental health applications.

In conclusion, the proposed deep learning-based stress detection framework marks a significant step toward intelligent, scalable, and ethical digital mental health solutions. With continued refinement and responsible deployment, such systems can play a transformative role in promoting psychological well-being in the digital age.

VI. FUTURE WORK

Multimodal Data Integration Future models can benefit from incorporating additional data modalities beyond text, such as user metadata (e.g., posting frequency, geolocation, engagement patterns), images, or audio inputs where available. Multimodal fusion could enable a more holistic understanding of stress and emotional states. **Temporal and Longitudinal Modelling** Understanding how stress evolves over time for an individual or group requires longitudinal tracking. Future research could investigate temporal models such as transformers with memory components or recurrent networks to analyse sequential stress patterns and predict onset or escalation.

Cross-Platform and Cross-Lingual Generalization Our current study is limited to English content on Twitter and Reddit. Expanding the framework to other platforms (e.g., Facebook, Instagram) and languages would help assess the robustness and adaptability of the models. **Transfer learning and multilingual pretraining** approaches could support such efforts. **Real-Time Deployment and Alert Systems** Developing a deployable, real-time stress detection system could have immediate applications in digital mental health support. Integrating this framework with mobile apps or browser extensions, paired with secure backends, could help users manage stress proactively or alert caregivers during critical moments.

Explainability and Interpretability For deployment in sensitive domains like mental health, model transparency is crucial. Implementing explainable AI (XAI) methods such as LIME, SHAP, or attention visualizations would help end users and clinicians understand how predictions are made and foster greater trust in the system. Ethical AI and Fairness Audits Future work must prioritize fairness and ethical safeguards. Studies should evaluate bias across demographic subgroups and develop bias mitigation strategies. Additionally, comprehensive audits on privacy, data consent, and model accountability should be conducted before real-world integration.

REFERENCES

- [1] De Choudhury, M., Counts, S., & Horvitz, E. (2013). Predicting postpartum changes in emotion and behavior via social media. *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, 3267–3276.
- [2] Coppersmith, G., Dredze, M., & Harman, C. (2014). Quantifying mental health signals in Twitter. *ACL Workshop on Computational Linguistics and Clinical Psychology*, 51–60.
- [3] Orabi, A. H., Buddhitha, P., Orabi, M. H., & Inkpen, D. (2018). Deep learning for depression detection of Twitter users. *Proceedings of the Fifth Workshop on Computational Linguistics and Clinical Psychology*, 88–97.
- [4] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. *NAACL-HLT*, 4171–4186.
- [5] Chancellor, S., Birnbaum, M. L., Caine, E. D., Silenzio, V., & De Choudhury, M. (2019). A taxonomy of ethical tensions in inferring mental health states from social media. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW), 1–29.
- [6] Huang, Y., Zhao, L., & Zhang, D. (2020). Identifying mental health issues on Chinese social media: A text mining approach. *IEEE Access*, 8, 27339–27350.
- [7] Saha, K., Sugar, B., Torous, J., Abrahao, B., Kiciman, E., & De Choudhury, M. (2021). A social media study on the associations between exposure to COVID-19 news and mental health symptoms. *NPJ Digital Medicine*, 4(1), 1–9.
- [8] Pennebaker, J. W., Boyd, R. L., Jordan, K., & Blackburn, K. (2015). The development and psychometric properties of LIWC2015. *University of Texas at Austin*.
- [9] Lin, H., Jia, J., Nie, L., et al. (2014). User-level psychological stress detection from social media using deep neural network. *ACM International Conference on Multimedia*, 507–516.
- [10] Guntuku, S. C., Yaden, D. B., Kern, M. L., Ungar, L. H., & Eichstaedt, J. C. (2017). Detecting depression and mental illness on social media: An integrative review. *Current Opinion in Behavioral Sciences*, 18, 43–49.
- [11] Calvo, R. A., Milne, D. N., Hussain, M. S., & Christensen, H. (2017). Natural language processing in mental health applications using non-clinical texts. *Natural Language Engineering*, 23(5), 649–685.
- [12] Ernala, S. K., Rizvi, A. F., & De Choudhury, M. (2020). Linguistic markers indicating therapeutic outcomes of social media disclosures of schizophrenia. *CSCW*, 1–27.
- [13] Benton, A., Mitchell, M., & Hovy, D. (2017). Multi-task learning for mental health using social media text. *EACL*, 152–162.
- [14] Yates, A., Cohan, A., & Goharian, N. (2017). Depression and self-harm risk assessment in online forums. *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 2968–2978.
- [15] Losada, D. E., & Crestani, F. (2016). A test collection for research on depression and language use. *International Conference of the CLEF Association*, 28–39.
- [16] Prieto, V. M., Matos, S., Alvarez, M., Cacheda, F., & Oliveira, J. L. (2014). Twitter: A good place to detect health conditions. *PLoS One*, 9(1), e86191.
- [17] Resnik, P., Armstrong, W., Claudino, L., Nguyen, T., Nguyen, V. A., & Boyd-Graber, J. (2015). Beyond LDA: Exploring supervised topic modeling for depression-related language in Twitter. *Proceedings of CLPsych Workshop*, 99–107.
- [18] Tsugawa, S., Kikuchi, Y., et al. (2015). Recognizing depression from Twitter activity. *Conference on Computer Supported Cooperative Work*, 318–326.
- [19] Rude, S., Gortner, E. M., & Pennebaker, J. (2004). Language use of depressed and depression-vulnerable college students. *Cognition & Emotion*, 18(8), 1121–1133.
- [20] Park, M., McDonald, D. W., & Cha, M. (2013). Perception differences between the depressed and non-depressed users in Twitter. *ICWSM*, 476–485.