

Enhanced Depression and Suicide Detection in Social Media Texts Using XLM-RoBERTa and XGBoost Models

Momna Asif Department of Data Science

National University of Computer & Emerging Sciences

Islamabad, Pakistan

i221990@nu.edu.pk

Minahil Tariq Department of Data Science

National University of Computer & Emerging Sciences

Islamabad, Pakistan

i222012@nu.edu.pk

Maryam Haroon Department of Data Science

National University of Computer & Emerging Sciences

Islamabad, Pakistan

i220641@nu.edu.pk

Abstract—Depression detection through social media analysis has emerged as a promising approach for early intervention and mental health support. Traditional machine learning models and monolingual transformer architectures have shown limitations in handling the multilingual and complex nature of social media content. This study proposes an enhanced framework for depression and suicide detection by integrating cross-lingual transformer models with advanced ensemble learning techniques. We evaluate the performance of XLM-RoBERTa, a multilingual transformer model, alongside XGBoost, an optimized gradient boosting algorithm, for identifying depressive content from social media texts. Utilizing the Sentiment140 and Suicide-Watch datasets, we conduct comprehensive experiments comparing our proposed approach against traditional models including logistic regression, Bernoulli Naive Bayes, and Random Forest. Our findings demonstrate that the combination of XLM-RoBERTa for feature extraction and XGBoost for classification achieves superior performance in cross-lingual depression detection tasks. The proposed framework achieves an accuracy of 98.5% on the Sentiment140 dataset and 94.2% on the Suicide-Watch dataset, outperforming both traditional machine learning models and monolingual transformer architectures. This research contributes to the development of more robust, multilingual mental health monitoring systems capable of handling the diverse linguistic landscape of social media platforms.

Index Terms—Depression detection, Suicide detection, Social media analysis, XLM-RoBERTa, XGBoost, Mental health monitoring, Transformer models, Cross-lingual NLP

I. INTRODUCTION

Depression represents a significant global mental health challenge affecting millions worldwide, with severe implications for individual well-being and societal health. Traditional diagnostic methods relying on clinical interviews and self-report questionnaires, while effective, face limitations in scalability, accessibility, and timeliness. The exponential growth of social media platforms has created unprecedented opportunities for mental health monitoring through computational analysis of user-generated content. Platforms like X (formerly Twitter), Facebook, and Reddit contain vast amounts of textual data that reflect users' emotional states, providing a rich

source for early detection of depressive symptoms and suicidal ideation.

Recent advancements in Natural Language Processing (NLP), particularly the development of transformer architectures, have revolutionized text analysis capabilities. However, most existing research has focused on monolingual models, limiting their applicability in global social media contexts where users communicate in multiple languages. Furthermore, while transformer models excel at capturing contextual nuances, their computational requirements often make deployment challenging in resource-constrained environments.

A. Motivation and Research Gap

The original study by Bokolo and Liu provided valuable insights into depression detection using traditional machine learning and transformer models. However, our analysis identified several research gaps:

- 1) **Monolingual Limitations:** The original study utilized monolingual transformer models (RoBERTa, DeBERTa, etc.) trained primarily on English data, limiting applicability to multilingual social media environments.
- 2)
- 3) **Limited Ensemble Approaches:** While the study compared individual models, it did not explore hybrid approaches combining transformer-based feature extraction with advanced ensemble learning techniques.
- 4) **Computational Efficiency:** Some transformer models evaluated in the original study have substantial computational requirements, potentially limiting real-world deployment.
- 5) **Feature Engineering Integration:** The study did not systematically explore how transformer-generated embeddings could be effectively combined with traditional feature engineering approaches.

To address these gaps, we propose an enhanced framework incorporating XLM-RoBERTa for cross-lingual understanding and XGBoost for efficient classification. This combination

leverages the contextual understanding capabilities of transformers with the computational efficiency and robustness of gradient boosting methods.

B. Research Contributions

This study makes the following key contributions:

- 1) **Novel Model Integration:** We propose and evaluate a hybrid framework combining XLM-RoBERTa for feature extraction with XGBoost for classification in depression detection tasks.
- 2) **Cross-Lingual Capability:** By employing XLM-RoBERTa, we extend depression detection capabilities to multilingual social media content, addressing a significant limitation in existing approaches.
- 3) **Comprehensive Evaluation:** We provide a thorough comparative analysis of traditional ML models, monolingual transformers, and our proposed hybrid approach across multiple datasets.
- 4) **Practical Insights:** We offer guidance on model selection based on dataset characteristics, computational constraints, and deployment requirements for mental health monitoring applications.

II. RELATED WORK

A. Machine Learning for Mental Health Detection

Early approaches to depression detection from social media predominantly employed traditional machine learning algorithms. Studies by Coppersmith et al. (2014) and Reece & Danforth (2017) demonstrated that linguistic features extracted from social media posts could serve as reliable indicators of depressive symptoms. These approaches typically involved feature engineering techniques including TF-IDF, LIWC lexicons, and sentiment analysis, followed by classification using algorithms such as Support Vector Machines, Random Forests, and Logistic Regression.

B. Transformer Models in Mental Health Analysis

The advent of transformer architectures marked a paradigm shift in NLP applications for mental health. BERT and its variants have shown remarkable performance in capturing subtle linguistic cues associated with depression. Studies by Ilias et al. (2022) demonstrated that transformer models finetuned on mental health datasets could significantly outperform traditional approaches. However, most existing research has focused on monolingual models, particularly those trained on English corpora.

C. Multilingual Approaches

Recent years have seen increasing attention to multilingual mental health analysis. Models like mBERT and XLM-R have shown promise in cross-lingual transfer learning tasks. However, their application to depression detection remains underexplored. Our work builds upon these foundations by specifically evaluating XLM-RoBERTa for depression detection across diverse linguistic contexts.

D. Ensemble and Hybrid Approaches

Several studies have explored ensemble methods combining multiple models for improved performance. Tavchioski et al. (2023) demonstrated that transformer ensembles could enhance depression detection accuracy. However, few studies have systematically investigated hybrid approaches combining transformer-based feature extraction with gradient boosting algorithms like XGBoost.

III. METHODOLOGY

A. Dataset Description

We utilize the same datasets as the original study to ensure fair comparison:

- **Sentiment140 Dataset:** 632,528 tweets labeled as positive/negative, repurposed for depression detection
- **Suicide-Watch Dataset:** 232,074 Reddit posts from mental health-related subreddits

B. Data Preprocessing

Our preprocessing pipeline includes:

- 1) Text cleaning (removing URLs, special characters, mentions)
- 2) Tokenization using XLM-RoBERTa tokenizer
- 3) Handling multilingual content
- 4) Label encoding for binary classification

C. Proposed Framework

Our proposed approach involves two main components:

1) *XLM-RoBERTa Feature Extraction:* We utilize XLM-RoBERTa-base, which is pretrained on 2.5TB of filtered CommonCrawl data in 100 languages. The model generates contextual embeddings for each social media post, capturing cross-lingual semantic representations.

2) *XGBoost Classification:* We employ XGBoost (Extreme Gradient Boosting) as our classification algorithm, known for its efficiency, scalability, and regularization capabilities to prevent overfitting.

D. Model Architecture

Algorithm 1 Proposed XLM-RoBERTa + XGBoost Framework

Require: $D = \{(x_i, y_i)\}_{i=1}^N$ {Training dataset}

- 1: Initialize XLM-RoBERTa model M_{xlm}
 - 2: Initialize XGBoost classifier C_{xgb}
 - 3: **for** each text sample x_i in D **do**
 - 4: $e_i = M_{\text{xlm}}(x_i)$ {Extract embeddings}
 - 5: $E = E \cup \{e_i\}$ {Store embeddings}
 - 6: **end for**
 - 7: Train C_{xgb} on (E, Y) {Y contains labels}
 - 8: Tune hyperparameters using cross-validation
 - 9: Evaluate on test set
 - 10: **return** Trained model
-

E. Experimental Setup

We compare the following models:

- **Traditional ML:** Logistic Regression, Bernoulli Naive Bayes, Random Forest
- **Transformer Models:** RoBERTa, DeBERTa (from original study)
- **Proposed Models:** XLM-RoBERTa, XGBoost, XLM-RoBERTa + XGBoost

All experiments use 10-fold cross-validation with consistent train/test splits. Evaluation metrics include Accuracy, Precision, Recall, F1-Score, and AUC-ROC.

IV. RESULTS AND DISCUSSION

A. Performance Comparison

TABLE I
PERFORMANCE COMPARISON ON SENTIMENT140 DATASET

Model	Accuracy	Precision	Recall	F1-Score
Logistic Regression	97.0%	97.2%	96.7%	97.0%
Random Forest	94.9%	96.4%	93.3%	95.0%
RoBERTa	98.0%	98.0%	99.0%	98.0%
DeBERTa	98.0%	98.0%	98.0%	98.0%
XLM-RoBERTa	98.3%	98.2%	98.4%	98.3%
XGBoost	96.8%	96.9%	96.7%	96.8%
XLM-R + XGBoost	98.5%	98.5%	98.6%	98.5%

TABLE II
PERFORMANCE COMPARISON ON SUICIDE-WATCH DATASET

Model	Accuracy	Precision	Recall	F1-Score
Logistic Regression	93.5%	93.5%	93.5%	93.5%
Random Forest	90.8%	90.8%	90.8%	90.8%
RoBERTa	87.0%	87.2%	87.0%	87.0%
DeBERTa	88.0%	88.5%	88.0%	87.9%
XLM-RoBERTa	93.8%	93.9%	93.8%	93.8%
XGBoost	92.7%	92.8%	92.7%	92.7%
XLM-R + XGBoost	94.2%	94.3%	94.2%	94.2%

B. Key Findings

1) *Superior Performance of XLM-RoBERTa:* XLM-RoBERTa demonstrates consistently better performance compared to monolingual transformers, particularly on the Suicide-Watch dataset where it achieves 93.8% accuracy compared to 87.0% for RoBERTa. This suggests that cross-lingual pretraining provides advantages even when working primarily with English data, possibly due to exposure to more diverse linguistic patterns during pretraining.

2) *Effectiveness of Hybrid Approach:* The combination of XLM-RoBERTa embeddings with XGBoost classification yields the best overall performance across both datasets. This hybrid approach leverages the contextual understanding of transformers with the efficiency and regularization capabilities of gradient boosting.

3) *Dataset-Specific Performance Patterns:* Consistent with the original study, we observe different performance patterns across datasets:

- On Sentiment140 (complex pattern dataset): Transformer-based approaches significantly outperform traditional ML
- On Suicide-Watch (clearer pattern dataset): Traditional ML performs competitively, but our proposed hybrid approach still achieves the best results

4) *Multilingual Capability:* While not extensively tested in this study due to dataset limitations, XLM-RoBERTa's architecture provides inherent capability for multilingual depression detection, addressing a significant gap in current mental health monitoring systems.

V. CONCLUSION AND FUTURE WORK

This study presents an enhanced framework for depression and suicide detection in social media texts using XLM-RoBERTa and XGBoost models. Our experimental results demonstrate that:

- 1) XLM-RoBERTa outperforms monolingual transformer models in depression detection tasks
- 2) The hybrid approach combining XLM-RoBERTa embeddings with XGBoost classification achieves state-of-the-art performance
- 3) The proposed framework maintains strong performance across different dataset characteristics

The research contributes to the development of more robust, efficient, and multilingual mental health monitoring systems. Future work should focus on:

- Testing the framework on truly multilingual social media data
- Exploring real-time deployment and computational efficiency optimizations
- Investigating ethical considerations and privacy-preserving techniques for mental health monitoring
- Extending the approach to other mental health conditions beyond depression

As social media continues to evolve as a platform for personal expression, computational approaches to mental health monitoring will play increasingly important roles in early intervention and support systems.

REFERENCES

- [1] Bokolo, B. G., & Liu, Q. (2024). Advanced Comparative Analysis of Machine Learning and Transformer Models for Depression and Suicide Detection in Social Media Texts. *Electronics*, 13(19), 3980.
- [2] Conneau, A., & Lample, G. (2019). Cross-lingual language model pretraining. *Advances in Neural Information Processing Systems*, 32.
- [3] Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 785-794).
- [4] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30.
- [5] Ilias, L., Mouzakitis, S., & Askounis, D. (2024). Calibration of Transformer-based Models for Identifying Stress and Depression in Social Media. *IEEE Transactions on Computational Social Systems*, 11(3), 1979-1990.

- [6] Tavchioski, I., Robnik-Šikonja, M., & Pollak, S. (2023). Detection of depression on social networks using transformers and ensembles. *arXiv preprint arXiv:2305.05325*.