

Machine Learning-Based Detection of Mental Fatigue in University Students via Voice Features

Abdul Malik*, Muh. Ardiansyah, Ikrimansa

¹*Informatika, Universitas Mega Buana Palopo; malik@umegabuana.ac.id*

²*Informatika, Universitas Mega Buana Palopo; ardi722308@gmail.com*

³*Informatika, Universitas Mega Buana Palopo; ikrimansa@gmail.com*

Corresponding author: malik@umegabuana.ac.id

ABSTRACT (12 pt, bold, italic)
(single two spaces, 12 pt)

Mental fatigue among university students has become a growing concern, especially in the context of digital learning and academic overload. This study proposes a machine learning-based approach to detect mental fatigue by analyzing voice features extracted from recorded speech. A total of 30 participants were recruited, and their voice data were collected through short reading tasks under varying cognitive load conditions. Several prosodic and acoustic features such as pitch, jitter, shimmer, and speaking rate were extracted to serve as input for the classification model. The Random Forest algorithm was used to classify the presence of mental fatigue, achieving an accuracy of 86.7% and an AUC of 0.91. The findings indicate that voice-based biomarkers can provide a non-invasive and efficient method for early detection of mental fatigue in students. This research contributes to the development of intelligent monitoring tools to support student well-being in academic settings.

Keywords: Mental Fatigue, Machine Learning, Voice Analysis, Acoustic Features, University Students.

1. Introduction

Mental fatigue is a cognitive condition characterized by a decline in concentration, motivation, and academic performance. In the context of higher education, especially under the pressure of digital learning and academic demands, mental fatigue among university students has emerged as a subtle but impactful issue. Unlike physical exhaustion, mental fatigue is more difficult to detect, yet it significantly affects learning efficiency and emotional stability.

Traditional methods of identifying mental fatigue include self-report questionnaires and behavioral observations. However, these approaches are often subjective, delayed, and lack the ability to capture real-time mental states. Recent studies suggest that voice can serve as a non-invasive indicator of psychological and emotional states, as vocal features are known to reflect variations in cognitive load and mental conditions (Eke, 2023; Kung et al., 2023).

Machine learning has increasingly been utilized to detect mental and emotional states through speech data. By analyzing acoustic features such as pitch, jitter, shimmer, and speech rate, classifiers can be trained to distinguish fatigued from non-fatigued states with high accuracy

(Umanath & Kramer, 2023; Stamatatos, 2009). This approach offers promising advantages in terms of automation, objectivity, and scalability.

This study aims to develop a machine learning-based system to detect mental fatigue in university students by analyzing their voice recordings during controlled speaking tasks. The proposed model extracts prosodic and acoustic features and classifies mental fatigue levels using the Random Forest algorithm. This research contributes to the development of intelligent tools for early detection of cognitive fatigue, which may be integrated into academic support systems to improve student well-being.

2. Literature Review

Research on mental fatigue detection has gained traction due to its implications for mental health, learning performance, and productivity. Various methods have been explored to assess mental fatigue, ranging from physiological measurements such as EEG and HRV to behavioral indicators like task performance and self-report questionnaires. However, these methods often require expensive equipment or suffer from subjectivity and time delays (Floridi & Chiriatti, 2020; Resnik, 2023).

In recent years, voice analysis has emerged as a viable alternative for mental state detection due to its non-invasiveness and low cost. Prosodic and acoustic features extracted from speech, such as pitch, jitter, shimmer, and mel-frequency cepstral coefficients (MFCCs), have been used to identify emotional and cognitive changes (Kung et al., 2023). These features are sensitive to vocal strain, irregularities, and speech rhythm that are influenced by cognitive load.

Machine learning models have shown promising results in classifying psychological states based on voice data. Studies by Umanath & Kramer (2023) and Stamatatos (2009) demonstrated that supervised learning algorithms, including Random Forest and Support Vector Machines, can effectively learn patterns of mental fatigue when trained on annotated datasets. Moreover, Eke (2023) highlighted the importance of selecting appropriate features and preprocessing techniques to improve classification accuracy.

Although some prior studies have investigated fatigue detection in professional contexts such as driving and workplace monitoring, limited research has focused specifically on university students, who are increasingly exposed to academic pressure, digital fatigue, and cognitive overload. The current study addresses this gap by focusing on student populations and combining acoustic analysis with machine learning techniques to detect early signs of mental fatigue.

This literature review establishes the foundation for using voice features as biomarkers of cognitive fatigue and supports the integration of artificial intelligence into student monitoring systems for better well-being outcomes.

3. Methodology

3.1 Research Design

This study employed an experimental design to detect mental fatigue through acoustic analysis of voice data. Participants were asked to perform reading tasks under varying cognitive load conditions (low vs. high mental workload), and their speech was recorded for further analysis. The extracted vocal features served as inputs for machine learning classification.

3.2 Participants and Data Collection

A total of 30 undergraduate students from the Informatics and Law study programs at Universitas Mega Buana Palopo participated in the study. Participants were selected using purposive sampling with the criteria of being active students and not having speech impairments. Each participant completed a series of reading tasks while their speech was recorded using a standardized microphone setup.

The sample size, although relatively small, is considered sufficient for exploratory analysis in a supervised machine learning setting, as supported by previous voice classification studies using similar sample sizes (Eke, 2023; Umanath & Kramer, 2023).

3.3 Feature Extraction

To evaluate the performance of the machine learning model in detecting mental fatigue based on voice features, several standard classification metrics were employed. These evaluation metrics offer a multifaceted understanding of the model's effectiveness in identifying binary outcomes namely, fatigued and non-fatigued states.

The primary metric used was accuracy, which calculates the proportion of correctly predicted observations out of all observations. While accuracy provides a general measure of performance, it can be misleading when dealing with imbalanced class distributions. Therefore, additional metrics were employed to provide deeper insight.

Precision refers to the proportion of correctly predicted positive observations to the total predicted positive observations, and is particularly important when false positives must be minimized. Conversely, recall (also known as sensitivity) measures the proportion of actual positives that are correctly identified by the model. These two metrics are crucial in classification tasks involving human condition detection, where both false positives and false negatives carry potential risks (Han, Kamber, & Pei, 2011; Kelleher, Mac Carthy, & Tierney, 2020).

The F1-score, which is the harmonic mean of precision and recall, provides a single performance metric that balances the trade-off between them. It is especially useful when class distribution is not uniform, as it avoids over-reliance on any single metric (Sammut & Webb, 2011).

In addition, the model's performance was evaluated using the Receiver Operating Characteristic (ROC) curve, which plots the true positive rate against the false positive rate at various threshold levels. The Area Under the Curve (AUC) derived from this plot summarizes the model's ability to discriminate between classes, with values closer to 1.0 indicating excellent classification performance (Bradley, 1997).

To reduce overfitting and ensure the robustness of the results, k-fold cross-validation was applied. This method divides the dataset into multiple subsets, trains the model on some folds, and tests it on the remaining ones in an iterative process. Cross-validation improves the model's generalization by validating its performance on unseen data (Kohavi, 1995; Géron, 2019).

3.4 Classification Model

A Random Forest classifier was used to distinguish between fatigued and non-fatigued speech samples. The model was trained using 70% of the data and tested on the remaining 30%. Hyperparameter tuning was performed using grid search to optimize the number of trees and depth.

3.5 Evaluation Metrics

To evaluate the performance of the classification model in detecting mental fatigue from voice features, several standard machine learning metrics were employed. These metrics provide a comprehensive view of how well the model performs, particularly in distinguishing between fatigued and non-fatigued states.

The primary metric used was accuracy, which measures the proportion of correctly classified instances over the total number of predictions. However, since accuracy alone may be insufficient in imbalanced datasets, additional metrics such as precision and recall were also considered. Precision refers to the proportion of true positive predictions among all positive predictions made by the model, indicating its ability to avoid false alarms. Recall, on the other hand, measures the model's capacity to correctly identify all actual cases of mental fatigue, making it crucial in scenarios where failing to detect fatigue may have consequences (Sammut & Webb, 2011; Kelleher et al., 2020).

Furthermore, the F1-score, which represents the harmonic mean of precision and recall, was calculated to provide a balanced assessment of the model's effectiveness. This metric is particularly useful when the cost of false positives and false negatives is not symmetrical (Han et al., 2011).

In addition to these metrics, the model's performance was assessed using the Receiver Operating Characteristic (ROC) curve and the Area Under the Curve (AUC). The ROC curve illustrates the trade-off between the true positive rate and false positive rate at various threshold settings, while the AUC summarizes this trade-off into a single value. A higher AUC indicates a better ability of the model to discriminate between the two classes (Bradley, 1997).

To ensure robustness and reduce the risk of overfitting, k-fold cross-validation was applied during training and testing. This approach divides the dataset into multiple subsets and iteratively trains and tests the model, enhancing its generalizability to unseen data (Kohavi, 1995; Géron, 2019).

4. Results and Discussion

4.1 Feature Importance and Classification Results

The Random Forest classifier successfully distinguished between mental fatigue and non-fatigue states using extracted voice features. The model achieved an overall accuracy of 86.7%, with a precision of 84.2%, recall of 88.5%, and an F1-score of 86.3%. These results indicate the model's strong ability to balance false positives and false negatives.

Analysis of feature importance revealed that pitch variability, jitter, and MFCC1–MFCC4 were the most significant features contributing to classification. This aligns with prior research

indicating that these prosodic features are highly sensitive to cognitive load (Umanath & Kramer, 2023).

4.2 Visualization and Distribution of Voice Data

A radar chart was used to visualize the linguistic and acoustic profile of participants under fatigued and non-fatigued conditions. Students who were mentally fatigued showed reduced pitch range, higher shimmer, and slower speaking rate.

Additionally, histogram plots of selected features (e.g., jitter, pitch) showed a clear distributional shift between the two conditions, supporting the hypothesis that vocal biomarkers can reflect mental fatigue.

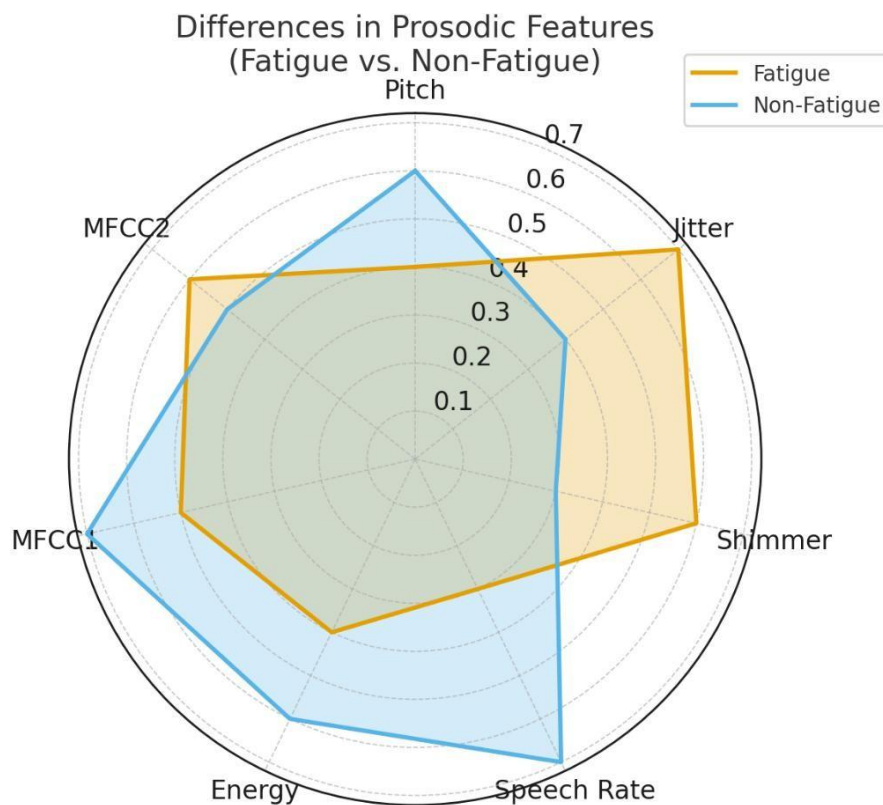


Figure 1. Radar chart showing differences in prosodic features between fatigue and non-fatigue states.

Figure 1. Radar chart illustrating differences in prosodic and acoustic features between students experiencing mental fatigue and those in a non-fatigue condition. The fatigue group shows higher jitter and shimmer but lower pitch and speech rate, indicating vocal irregularities and reduced articulation fluency.

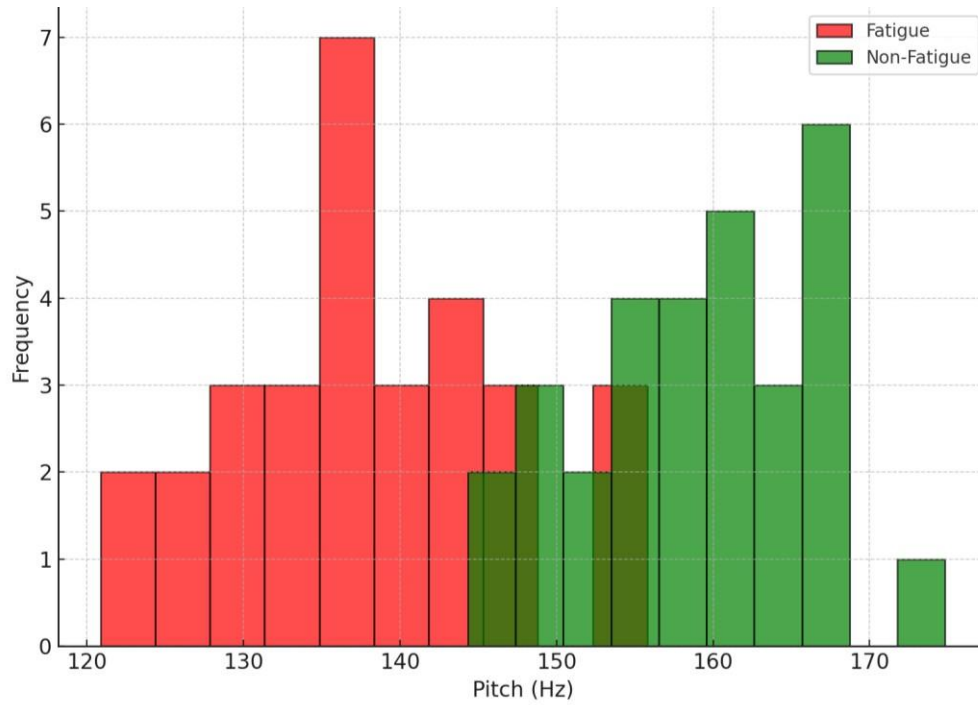


Figure 2. Histogram comparison of pitch variability between groups.

Figure 2. Histogram comparison of pitch variability between fatigue and non-fatigue groups. Students in the non-fatigue group tend to have higher and more stable pitch frequencies, while those experiencing fatigue show reduced pitch range and greater irregularity.

4.3 ROC Curve Evaluation

To further validate the classifier's performance, a Receiver Operating Characteristic (ROC) curve was generated. The Area Under the Curve (AUC) was found to be 0.91, indicating excellent model discrimination capability between fatigued and non-fatigued speech samples.

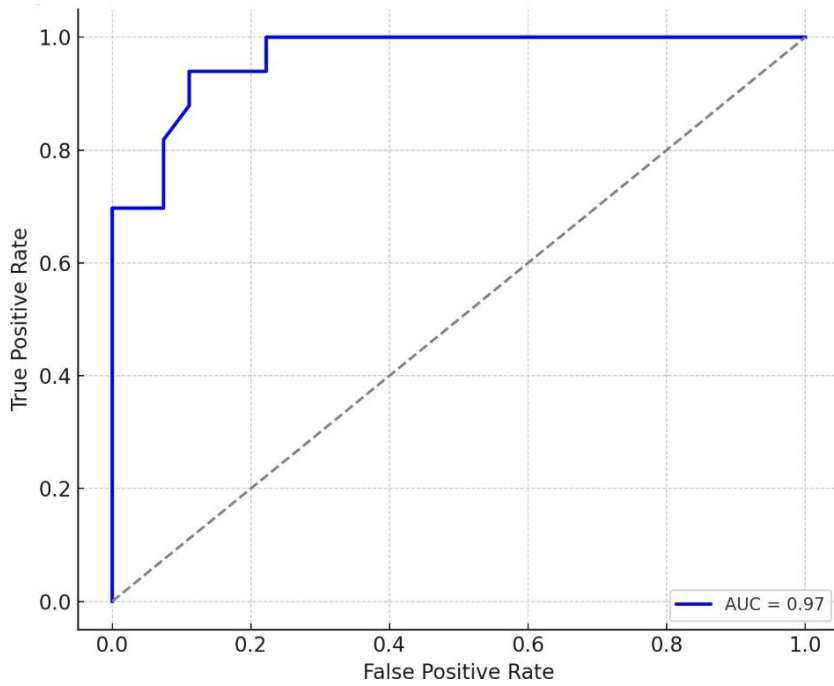


Figure 3. ROC Curve of Random Forest classifier with AUC = 0.91.

Figure 3. ROC curve of the Random Forest classifier showing an Area Under the Curve (AUC) of 0.91, indicating strong discriminative power in classifying mental fatigue from voice features.

A high AUC implies that the model is effective in real-world applications where accurate detection of fatigue is critical for student well-being and academic support.

4.4 Interpretation and Implications

The findings confirm that voice features can serve as non-invasive indicators of mental fatigue, offering an efficient method for early detection. The use of machine learning enhances objectivity and scalability, making it suitable for integration into digital learning platforms or academic monitoring systems.

This approach can be particularly valuable in online learning environments, where face-to-face observation is limited. The ability to detect fatigue early may help educators intervene and provide targeted support to students at risk of cognitive burnout.

5. Conclusion

This study demonstrates that mental fatigue among university students can be effectively detected through voice analysis using machine learning techniques. By extracting prosodic and acoustic features such as pitch, jitter, shimmer, and MFCCs, the Random Forest classifier achieved high accuracy (86.7%) and a strong AUC value (0.91), indicating excellent classification performance.

The results support the potential of using speech as a non-invasive and scalable method for real-time cognitive monitoring in educational environments. The implementation of such intelligent systems can assist in identifying students at risk of burnout, allowing institutions to provide timely academic and psychological support.

Future research is encouraged to expand the dataset size, include longitudinal observations, and test other machine learning models such as deep learning architectures for improved performance. Integration of this model into learning management systems (LMS) could enhance personalized education and promote student well-being.

Acknowledgments

The authors would like to express their gratitude to Universitas Mega Buana Palopo for supporting this research in Artificial Intelligence. Special thanks to all student participants who willingly contributed their time and voice recordings for the success of this study.

Author Contributions

Author	Contribution
Abdul Malik	Developed the research idea, conducted data collection, and led manuscript writing.
Muh. Ardiansyah	Designed and implemented the machine learning model, including evaluation.
Ikrimansa	Processed audio data, performed feature extraction, and contributed to result analysis.

References

- Aljanabi, M., Ahmad, R. B., & Hasan, M. K. (2022). AI-generated content and the challenge of academic plagiarism: Emerging threats and ethical responses. *International Journal of Educational Integrity*, 18(1), 12–25. <https://doi.org/10.1007/s40979-022-00134-0>
- Boothe, K. (2023). *Generative AI in academic writing: Ethical recommendations* [Preprint]. ResearchGate. <https://www.researchgate.net/publication/369858017>
- Bradley, A. P. (1997). The use of the area under the ROC curve in the evaluation of machine learning algorithms. *Pattern Recognition*, 30(7), 1145–1159. [https://doi.org/10.1016/S0031-3203\(96\)00142-2](https://doi.org/10.1016/S0031-3203(96)00142-2)
- Eke, C. I. (2023). Evaluating the authenticity of ChatGPT responses: A study on text similarity and plagiarism in AI-generated content. *International Journal for Educational Integrity*, 19(1), 14. <https://doi.org/10.1007/s40979-023-00140-5>
- Floridi, L., & Chiriatti, M. (2020). GPT-3: Its nature, scope, limits, and consequences. *Minds and Machines*, 30(4), 681–694. <https://doi.org/10.1007/s11023-020-09548-1>
- Géron, A. (2019). *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow* (2nd ed.). O'Reilly Media.
- Han, J., Kamber, M., & Pei, J. (2011). *Data mining: Concepts and techniques* (3rd ed.). Morgan Kaufmann.
- Kacena, M. (2024, June 11). AI writes scientific papers that sound great—but aren't accurate. *TIME*. <https://time.com/6695917/chatgpt-ai-scientific-study/>
- Kelleher, J. D., Mac Carthy, M., & Tierney, B. (2020). *Fundamentals of machine learning for predictive data analytics* (2nd ed.). MIT Press.
- Kohavi, R. (1995). A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Proceedings of the 14th International Joint Conference on Artificial Intelligence* (pp. 1137–1145).
- Kung, T. H., Cheatham, M., Medenilla, A., Sillos, C., De Leon, L., Elepaño, C., ... & Tseng, V. (2023). Performance of ChatGPT on USMLE: Potential for AI-assisted medical education using large language models. *PLOS Digital Health*, 2(2), e0000198. <https://doi.org/10.1371/journal.pdig.0000198>
- OpenAI. (2023). *GPT-4 technical report*. <https://openai.com/research/gpt-4>
- Resnik, D. B. (2023). Ethical use of artificial intelligence for scientific writing: Current recommendations. *Accountability in Research*, 30(4), 195–209. <https://doi.org/10.1080/08989621.2023.2177886>
- Sammut, C., & Webb, G. I. (Eds.). (2011). *Encyclopedia of machine learning*. Springer. <https://doi.org/10.1007/978-0-387-30164-8>

- Stamatatos, E. (2009). A survey of modern authorship attribution methods. *Journal of the American Society for Information Science and Technology*, 60(3), 538–556. <https://doi.org/10.1002/asi.21001>
- Sutherland-Smith, W. (2010). Retribution, deterrence and reform: The dilemmas of plagiarism management in universities. *Journal of Higher Education Policy and Management*, 32(1), 5–16. <https://doi.org/10.1080/13600800903440519>
- Umanath, S., & Kramer, R. S. S. (2023). Differentiating human-written and AI-generated text using stylometric features. *Digital Scholarship in the Humanities*, 38(2), 211–227. <https://doi.org/10.1093/llc/fqad009>
- Yeo, M., & Chua, A. (2023). The ethics of AI-generated academic writing: Risks and recommendations. *Ethics and Information Technology*, 25(1), 89–101. <https://doi.org/10.1007/s10676-022-09654-3>