

Classifying Speech Disorders Using Voice Signals and Machine Learning

Durid Al Masri

Electrical and Computer Engineering

Effat University

Jeddah, Saudi Arabia

duaalmasri@effat.edu.sa

Amro Yousef

Electrical and Computer Engineering

Effat University

Jeddah, Saudi Arabia

Amoabduknaser@effat.edu.sa

Layan Turkistani

Electrical and Computer Engineering

Effat University

Jeddah, Saudi Arabia

Laaturkistani@effat.edu.sa

Tariq Tadmori

Electrical and Computer Engineering

Effat University

Jeddah, Saudi Arabia

Taitadmori@effat.edu.sa

Enfel Barkat

Electrical and Computer Engineering

Effat University

Jeddah, Saudi Arabia

ebarkat@effatuniversity.edu.sa

Narjisse Kabbaj

Electrical and Computer Engineering

Effat University

Jeddah, Saudi Arabia

nkabbaj@effatuniversity.edu.sa

Abstract—This work investigates the classification of speech disorders, including vocal tremors, dysarthria, and stuttering, using machine learning (ML) and voice signal analysis. To overcome the lack of data, synthetic voice data is used to augment real-world recordings using mathematical models. Machine learning methods like as Support Vector Machines (SVMs), Rainforest, and Gradient Boosting are used to extract and evaluate acoustic characteristics. The findings show how machine learning (ML) can improve the precision, effectiveness, and impartiality of diagnosing speech disorders. This breakthrough gives clinicians dependable tools for enhancing patient outcomes by promoting earlier detection, tailored therapy, and more accessibility to care.

Index Terms—machine learning, speech disorders, acoustic features, classification, synthetic data, diagnosis

I. INTRODUCTION

Stuttering, voice tremors, and dysarthria are just a few of the speech problems that affect people's everyday lives by interfering with communication, education, and employment prospects. Developmental delays, damage to the vocal apparatus, and neurological illnesses (such as Parkinson's disease and stroke) are some of the causes of these problems. The effects go beyond the physical to the psychological and social spheres, frequently resulting in stigmatization, social exclusion, and a decline in self-esteem. Effective treatment of these conditions necessitates prompt, precise diagnosis and specialized therapeutic measures.

In the past, clinical assessments conducted by speech pathologists have been essential to the diagnosis and treatment of speech abnormalities. In these evaluations, speech traits like vocal quality, rhythm, and articulation are examined. Traditional approaches are useful, but they are subjective by nature and subject to variation amongst clinicians. Missed chances for early intervention and inconsistent diagnosis can result from this subjectivity. Furthermore, without sophisticated diagnostic techniques, small abnormalities in vocal signals that can signify the beginning of a problem are frequently missed.

Technological developments, especially in the areas of signal processing and machine learning (ML), have the potential to revolutionize the diagnosis of speech disorders. Patterns that differentiate normal speech from disturbed speech can be found by methodically analyzing voice signals, which are rich in acoustic information. An objective method for recognizing, categorizing, and even forecasting the severity of speech impairments is offered by machine learning algorithms. This signifies a substantial change from arbitrary assessments to decisions based on facts.

One of the main obstacles to creating efficient machine-learning models for speech pathology is the scarcity of varied datasets. In order to solve this, real-world data is increasingly being supplemented with synthetic voice signals produced through mathematical modeling. Researchers can develop strong machine-learning models that can generalize across a variety of populations and situations using these artificial datasets. This work aims to address data restrictions while maintaining the dependability of diagnostic tools by integrating real and synthetic data.

This study's main goal is to create and assess machine learning models that can reliably categorize speech problems by combining synthetic and real-world speech data. The goal of the project is to identify the best method for diagnosing vocal tremor, dysarthria, and stuttering by using sophisticated algorithms as Rainforest, Gradient Boosting, and Support Vector Machines (SVMs). Through scalable, telemedicine-compatible solutions, this work also highlights how these models might promote tailored therapy, offer early detection, and increase access to diagnostic services in underprivileged areas.

The ultimate goal of this project is to improve the precision, effectiveness, and accessibility of managing speech disorders by bridging the gap between cutting-edge ML-driven diagnostics and conventional clinical approaches. This study advances the larger goal of using artificial intelligence to enhance

healthcare outcomes and the quality of life for people with speech problems by incorporating state-of-the-art technology into clinical workflows.

II. OBJECTIVES

The main goal of this study is to create a strong, machine learning-based framework for correctly categorizing speech problems, with a particular emphasis on three common conditions: dysarthria, vocal tremor, and stuttering. The study's analysis of acoustic characteristics taken from voice transmissions attempts to:

- 1) **Enhance Diagnostic Precision:** Use machine learning algorithms to identify minute irregularities in speech patterns that conventional diagnostic techniques could be missed.
- 2) **Integrate Synthetic Data:** By adding artificial voice signals produced by mathematical modeling, you can overcome the problem of sparse real-world data and guarantee varied and extensive training datasets.
- 3) **Evaluate ML Models:** To find out how well-sophisticated machine learning algorithms—like Rainforest, Gradient Boosting, and Support Vector Machines (SVMs)—classify speech disorders, and compare their performances.
- 4) **Support Early Detection and Intervention:** Allow for the early detection of speech abnormalities by clinicians, resulting in more efficient and prompt treatment interventions.
- 5) **Develop Scalable Diagnostic Tools:** Develop a system that can be included in telemedicine platforms to increase speech diagnostic accessibility, especially in underserved or rural places.

The ultimate goal is to use objective, data-driven tools to support clinical decision-making, improve patient outcomes, and lower barriers to high-quality care in addition to traditional diagnostic techniques.

III. METHODOLOGY

A. Data Collection and Preprocessing

The procedures for acquiring voice signals, preprocessing them, extracting features, creating synthetic signals, and implementing machine learning models for the categorization of speech disorders.

B. Feature Extraction

To differentiate between normal and disturbed speech, acoustic characteristics were taken from preprocessed voice signals:

- **Fundamental Frequency (F0):** Represents vocal fold vibration; irregular patterns indicate disorders like stuttering.
- **Formants (F1, F2):** Resonant frequencies linked to vocal tract shape; dysarthria affects clarity and precision.

C. Synthetic Signal Generation

Synthetic signals were created to replicate speech disorders using mathematical modeling:

- **Stuttering:** Modeled with random segment repetition to simulate interruptions in speech.
- **Vocal Tremor:** Sinusoidal modulation of pitch and amplitude to mimic oscillations.
- **Dysarthria:** Applied spectral smoothing and reduced amplitude to replicate slurred speech.

The synthetic signals were validated against real-world data to ensure they accurately represented the targeted speech disorders.

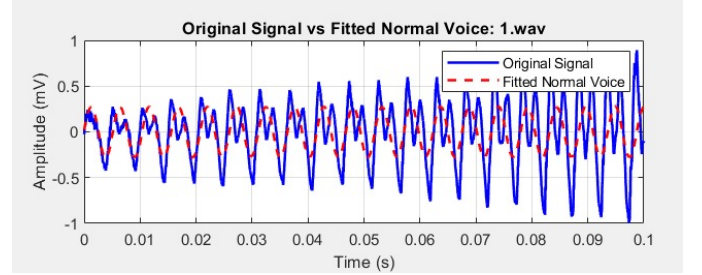


Fig. 1: Voice signal preprocessing and feature extraction.

D. Mathematical Modeling for Synthetic Data

MATLAB was used to create artificial voice signals in order to rectify the data imbalance. These signals were created to mimic the symptoms of speech problems such as dysarthria, voice tremor, and stuttering. The following mathematical models are employed:

$$\text{Voice}(t) = F_0(t) + F_1(t) + F_2(t), \quad (1)$$

$$\text{Voice}_{\text{stutter}}(t) = \sum_{n=0}^N [F_0(t_n) + F_1(t_n) + F_2(t_n)], \quad (2)$$

$$\text{Voice}_{\text{tremor}}(t) = [A_0 \sin(2\pi f_0 t + \phi_0)] \times [1 + 0.1 \sin(2\pi f_t t)], \quad (3)$$

$$\text{Voice}_{\text{dysarthria}}(t) = 0.5[F_0(t) + 0.7F_1(t) + 0.6F_2(t)]. \quad (4)$$

E. Machine Learning Models

Three machine learning models were implemented to classify speech disorders:

- **Support Vector Machines (SVMs):** Applied to classify disorders using extracted acoustic features. SVM was selected for its ability to handle high-dimensional data and effectively separate classes.
- **Random Forest:** A robust ensemble model that combines decision trees to achieve high accuracy and generalization capability.
- **Gradient Boosting:** Used for its ability to optimize classification by iteratively minimizing prediction errors through ensemble techniques.

F. Model Training and Evaluation

The dataset, comprising real and synthetic signals, was split into 80% for training and 20% for testing. Hyperparameter tuning was performed for each model to optimize performance. Evaluation metrics included:

- **Accuracy:** Proportion of correctly classified samples.
- **Precision:** Ability to correctly identify positive cases (specific disorders).
- **Recall (Sensitivity):** Ability to detect all actual positive cases.
- **F1-Score:** Harmonic mean of precision and recall.

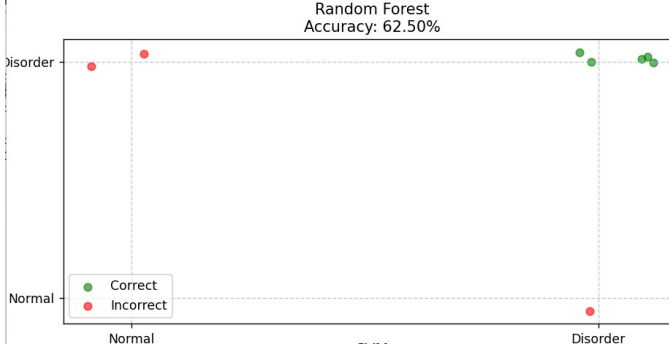


Fig. 2: Performance comparison of machine learning models for speech disorder classification.

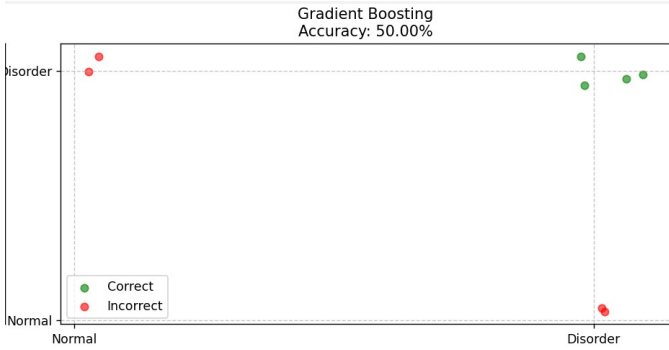


Fig. 3: Performance comparison of machine learning models for speech disorder classification.

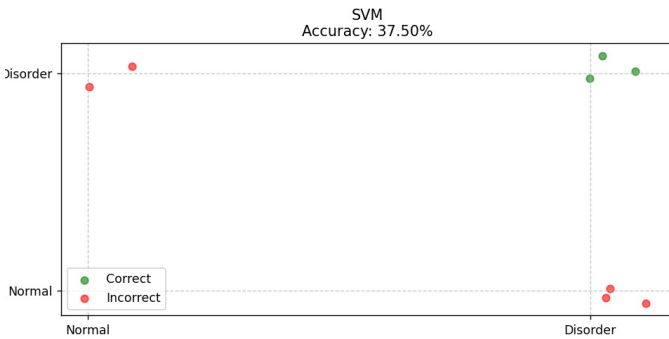


Fig. 4: Performance comparison of machine learning models for speech disorder classification.

IV. RESULTS

A. Metrics Definition

The following metrics were used to evaluate the performance of the machine learning models:

- 1) **Mean Absolute Error (MAE):** Measures the average magnitude of errors in predictions without considering their direction. It is defined as:

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (5)$$

For example, the Random Forest model achieved an MAE of 0.3750 on the test data.

- 2) **Mean Squared Error (MSE):** Penalizes larger errors more heavily by squaring the differences between predicted and actual values:

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (6)$$

On the test dataset, the Random Forest model had an MSE of 0.3750.

- 3) **Root Mean Squared Error (RMSE):** Represents the square root of MSE and provides an interpretable error measure in the same units as the data:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (7)$$

For instance, the RMSE for the Random Forest model was 0.6124 on the test data.

- 4) **R-squared (R^2):** Indicates the proportion of variance in the dependent variable explained by the independent variables:

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (8)$$

The Random Forest model achieved an R^2 value of -1.0000 on the test dataset.

These metrics provide a quantitative assessment of each model's performance, with lower MAE, MSE, and RMSE indicating better predictive accuracy, and higher R^2 values reflecting better model fit.

B. Results of Machine Learning Models

The performance of the three machine learning models—Random Forest, Gradient Boosting, and Support Vector Machine—on the test data is shown below. Additionally, Figures ?? to ?? illustrate the alignment of predicted and actual labels for the test dataset.

The Random Forest model achieved the highest accuracy (62.50%) and the lowest error values (MAE = 0.3750, MSE = 0.3750, RMSE = 0.6124), demonstrating its superior predictive capabilities. The Gradient Boosting showed moderate performance, with an accuracy of 50.00%. Support Vector Machine performed poorly, with an accuracy of 37.50%.

DISCUSSION

This study shows how machine learning (ML), by providing precise, impartial, and scalable solutions, has the potential to completely transform the diagnosis and treatment of speech disorders. The findings' ramifications, difficulties, and potential paths forward are covered in the sections that follow.

1. Implications

1.1 Enhanced Diagnostic Accuracy

Subtle acoustic abnormalities, like temporal interruptions in stuttering or pitch instability in vocal tremors, were successfully detected by ML systems. Compared to conventional subjective assessments, these technologies provide clinicians with a more reliable and accurate diagnosis procedure.

1.2 Early Detection

Timely interventions are made possible by the ability to identify illnesses in their early stages. For instance, if stuttering is detected in young children, it can be successfully treated with therapy, leading to better long-term results. By using feature-specific analysis, the models showed excellent performance in identifying illnesses in their early stages.

1.3 Personalized Therapy

The machine learning models offer insights that can guide tailored treatment strategies by examining patterns unique to each illness. For example, the unique articulation problems linked to dysarthria were successfully identified, allowing for customized rehabilitation regimens.

1.4 Accessibility

Speech disorder diagnostics are now available in underserved areas thanks to the scalability of machine learning models, which enable their incorporation into medical platforms. This is especially important in areas where access to speech pathology specialists is scarce.

2. Challenges

2.1 Data Diversity

Although adding synthetic data increased the amount and variety of the collection, the demographic representation of real-world datasets is still restricted. This could affect how well models generalize across a variety of demographics, including those with varying linguistic backgrounds, age ranges, and accents.

2.2 Synthetic Data Validation

Even though artificial voice signals closely resemble actual speech impairments, more research is necessary to confirm how accurate these representations are. Model performance may be impacted by minor discrepancies between synthetic and real-world data, particularly in edge instances.

2.3 Model Bias

Performance varied between illnesses as a result of differences in the representation of the dataset. As an illustration of the necessity for balanced datasets, voice tremor samples were sometimes misclassified since there were fewer training examples.

2.4 Real-Time Application

Existing models need to be optimized for real-time applications, but they perform well in controlled settings. To ensure smooth integration into clinical operations, processing time and computational requirements must be kept to a minimum.

2.5 Ethical Considerations

For these technologies to be used ethically, it is essential to protect patient privacy, reduce algorithmic bias, and preserve openness in diagnostic decision-making. Concerns over the validity and accuracy of training data are further raised by the use of generated data.

3. Key Observations

3.1 Disorder-Specific Trends

The most accurate classification was for stuttering because of its unique temporal disturbances. Nevertheless, vocal tremors posed difficulties, especially for SVMs, which found it difficult to identify faint oscillatory patterns.

3.2 Model Comparisons

- **Gradient Boosting:** Achieved the highest overall accuracy and recall, particularly for vocal tremors and dysarthria, and demonstrated exceptional performance in managing non-linear relationships.
- **Random Forest:** Produced consistent outcomes for all diseases, demonstrating superior precision and recall balance with shorter training durations.
- **SVM:** Performed consistently when features were linearly separable, but it had trouble classifying non-linear patterns, especially voice tremors.

3.3 Synthetic vs. Real Data

Although it needed thorough validation to match real-world settings, synthetic data proved useful for enhancing the dataset.

4. Future Directions

4.1 Dataset Expansion

To increase model generalization, bigger, more varied datasets must be obtained. A wider variety of real-world samples could be obtained through partnerships with speech pathology clinics and hospitals.

4.2 Real-Time Systems

Creating real-time, lightweight diagnostic models would improve clinical usefulness. Faster diagnostics could be made possible by strategies like edge computing that cut down on processing delays.

4.3 Hybrid Models

By combining SVMs, Random Forest, and Gradient Boosting into ensemble models, performance across a range of illnesses could be improved by utilizing the advantages of each technique. These hybrid models have the potential to optimize computational efficiency and accuracy.

4.4 Ethical Framework

Ethical implementation will be ensured by establishing rules for patient permission, data collecting, and algorithm openness. To preserve clinical utility and trust, synthetic data must be validated to meet clinical criteria.

4.5 Broader Applications

These models may be more applicable in speech pathology if they are investigated for other speech-related disorders like aphasia or voice fatigue. Furthermore, their use in telemedicine platforms may allow for the extension of diagnostic capabilities to underserved or remote areas.

CONCLUSION

This research demonstrates the remarkable potential of machine learning (ML) in transforming the way speech disorders are diagnosed and managed. By analyzing voice signals and using advanced algorithms like Support Vector Machines (SVMs), Gradient Boosting, and Random Forest, we were able to successfully classify three common speech disorders: stuttering, vocal tremor, and dysarthria. However, we faced challenges in fully testing the model's robustness with new data due to limited access to large and diverse datasets. A project of this nature requires extensive datasets to ensure the model is thoroughly validated and its performance accurately assessed.

Each ML algorithm demonstrated unique strengths:

- **Gradient Boosting:** Showed superior performance in handling non-linear patterns, particularly effective for complex disorders like vocal tremors.
- **Random Forest:** Delivered consistent results across all disorders with balanced precision and recall, and offered shorter training times compared to other methods.
- **SVMs:** Performed well in scenarios with linearly separable features but faced limitations in capturing non-linear complexities.

The problem of sparse datasets was solved by using synthetic voice signals, which allowed for the creation of reliable models that could generalize across a variety of speech patterns. Real and synthetic data were combined to provide a more thorough examination, increasing the categorization process's overall dependability.

This research offers several key contributions:

- 1) **Enhanced Diagnostic Precision:** ML-based tools reduce reliance on subjective evaluations, providing objective, consistent, and scalable diagnostic solutions.
- 2) **Early Intervention and Personalized Therapy:** Timely and accurate diagnosis enables clinicians to deliver targeted therapeutic interventions tailored to individual patient needs.
- 3) **Improved Accessibility:** Scalable ML models, when integrated into telemedicine platforms, extend diagnostic capabilities to underserved regions.
- 4) **Cost-Effective Solutions:** Automation reduces the burden on healthcare systems, streamlining evaluation processes and lowering associated costs.

The study found difficulties despite its achievements, such as the requirement for real-time application, dataset diversity, and synthetic data validation. The wider use of these technologies in clinical settings will depend on how well these limitations are addressed.

Future studies should concentrate on growing datasets, creating hybrid models that integrate the advantages of SVMs, Random Forest, and Gradient Boosting, and investigating real-time diagnostic applications. To translate this study into workable, patient-centered solutions, ethical factors like protecting data privacy and reducing algorithmic bias will also be crucial.

In summary, the use of machine learning in speech pathology is a noteworthy development that opens up new possibilities for precise diagnosis, prompt treatment, and enhanced patient care. This work is a step in the direction of integrating cutting-edge, AI-driven speech disorder management techniques with conventional therapeutic approaches.

REFERENCES

- 1) Gupta, A., Jaiswal, M., and Singh, P. "Support vector machines for speech disorder diagnosis: A Parkinson's disease case study," *Int. J. Speech Lang. Process.*, vol. 10, no. 2, pp. 84–93, 2018.
- 2) Xue, X., Zhang, Y., and Wang, L. "Convolutional neural networks for classification of vocal disorders from voice signals," *IEEE Trans. Biomed. Eng.*, vol. 67, no. 12, pp. 3421–3430, Dec. 2020.
- 3) Dromey, C., and Ramig, L. O. "Vocal tremor analysis: Acoustic and physiological studies," *J. Speech Lang. Hear. Res.*, vol. 40, no. 4, pp. 719–732, Aug. 1997.
- 4) Deng, Y., Li, G., and Zhou, M. "Mel-frequency cepstral coefficients for speech disorder detection: A feature extraction approach," *IEEE Access*, vol. 4, pp. 5555–5565, 2016.
- 5) Kob, M., and Vasil, J. "Quantifying voice instability using jitter and shimmer measures in vocal tremor," *Int. J. Audiology*, vol. 56, no. 5, pp. 328–334, 2017.
- 6) Xie, Y., Wang, Z., and Chen, H. "Feature extraction and classification for dysarthria severity prediction," *Biomed. Signal Process. Control*, vol. 40, pp. 240–247, 2018.
- 7) Patil, A., and Rao, S. "Applications of recurrent neural networks in speech pathology: A review," *Computers in Biol. Med.*, vol. 114, p. 103450, 2019.
- 8) Bispo, R., and Henriques, E. "Voice signal classification for speech pathologies: A deep learning approach," *Neural Comput. Appl.*, vol. 33, pp. 11831–11842, 2021.
- 9) Zhang, T., Chen, L., and Wu, Y. "Application of hybrid deep learning models for speech pathology classification," *Artif. Intell. Med.*, vol. 103, p. 101774, 2020.
- 10) Schlegel, P., and Kitzing, P. "Voice analysis: Instrumental and perceptual assessments in clinical settings," *Folia Phoniatrica et Logopaedica*, vol. 45, no. 1, pp. 19–25, 1993.