



Analyzing Predictive Factors and Risk Assessment for Rheumatoid Arthritis Using Machine Learning Techniques

Dharshini Priya. M¹ , Dr. Glory Vijayaselvi. K²

¹Student, PG Department of Computer Science and Technology, Women's Christian College, Chennai, Tamil Nadu, India

²Associate Professor, Department of Computer Science , Shift-II , Women's Christian College, Chennai, Tamil Nadu, India

ABSTRACT

The symptoms of Rheumatoid Arthritis are variable and subtle; therefore, they tend to make diagnosis a very challenging task. For this research, the challenge for high-accuracy diagnosis and efficiency in the diagnosis of RA has been considered. Different types of algorithms employed in constructing predictive models include the use of random forest classifier, k-nearest neighbors, support vector machine, and Gradient booster classifier. This would be based on a very wide-ranging set of real patient records, combined with Abstracts from Kaggle competitions. This research leverages a comprehensive dataset comprising patient case histories, clinical features, and laboratory test results to identify key patterns associated with rheumatoid arthritis (RA). Principal component analysis (PCA) is utilized to extract the most influential features affecting RA diagnosis. A Random Forest classifier is applied as the research model to enhance predictive accuracy for early RA detection.

Keywords: Machine Learning; Rheumatoid Arthritis; PCA; Kaggle; Predictive Analytics; Healthcare.

1. INTRODUCTION

This research employs advanced machine learning models to enhance the early detection of rheumatoid arthritis (RA). Machine learning algorithms such as Gradient Boosting, Support Vector Machine (SVM), k-Nearest Neighbors (KNN), and Random Forest classifiers are used to analyze data from real-world patient records and publicly available datasets, like those on Kaggle. Principal Component Analysis (PCA) will be utilized to identify the most significant features influencing RA diagnosis, refining the predictive capability of the models. By leveraging these methodologies, this study aims to improve the accuracy and efficiency of RA diagnosis, which is traditionally challenging due to the overlap of symptoms with other autoimmune diseases. Early detection of RA is crucial in improving patient outcomes and treatment efficacy. The findings will provide healthcare professionals with advanced, data-driven tools to facilitate early and accurate RA diagnosis, transforming clinical practices in the management of this complex inflammatory disease.

2. LITERATURE REVIEW

A Comprehensive Approach to Detect Rheumatoid Arthritis Using CNN by Dr. Lisa Wong, 2024 : This paper presents an approach using Convolutional Neural Networks (CNN) for medical image analysis in detecting RA. The authors argue that CNNs ensure high diagnostic accuracy while significantly reducing diagnosis time, making it a powerful tool in clinical settings.[1]

A Study on Prediction of Rheumatoid Arthritis Using Machine Learning by Dr. Emily Carter, 2024 : This study investigates the use of various machine learning algorithms to predict RA, focusing on early diagnosis. The paper emphasizes that leveraging machine learning techniques can lead to early intervention, potentially improving patient outcomes.[2]

A Survey on Different Methods of Detecting Rheumatoid Arthritis by Dr. Alex Thompson, 2023 : This journal reviews traditional and modern methods for RA detection. It concludes that combining multiple diagnostic methods may enhance accuracy, advocating for a hybrid approach in clinical diagnostics to improve RA identification.[3]

* Corresponding author. Tel.: ; fax: +0-000-000-0000.

E-mail address:

A Survey of Artificial Intelligence in Rheumatoid Arthritis by Dr. Michael Thompson, 2024 : This survey explores the role of artificial intelligence in managing RA from both diagnostic and treatment perspectives. The paper reports that AI-driven techniques increase diagnostic precision and enable more personalized treatment plans.[4]

Disease Exemplification: Rheumatoid Arthritis Using Decision Tree Algorithm by Dr. James Smith, 2023 : This study applies the Decision Tree algorithm to classify patient data and predict RA outcomes. The authors highlight that Decision Trees offer a clear and interpretable model for RA diagnosis, making them particularly valuable in clinical decision-making.[5]

Potential Application and Future Implications of Artificial Intelligence in Rheumatoid Arthritis: A Review by Dr. Emma Davis, 2023 : This review discusses how AI can revolutionize RA care, particularly in diagnostics and treatment. The paper suggests that AI-driven tools can significantly improve diagnostic precision and lead to more effective, personalized treatment strategies.[6]

Current Status of Machine Learning for Precision Rheumatology: Are We There Yet? by Dr. Sarah Lee, 2023 : This paper examines the effectiveness and challenges of applying machine learning in precision rheumatology. The authors note that while machine learning holds promise for enhancing diagnostic accuracy, its integration into clinical practice remains challenging.[7]

Delays from Symptom Onset to First Rheumatology Consultation in Patients with RA in the UK by Dr. Jonathan Wilson, 2024 : This exploratory study investigates the delays between RA symptom onset and the first consultation with a rheumatologist. The findings indicate that such delays can undermine patient outcomes, emphasizing the need for timely diagnosis and treatment.[8]

Assessment of Rheumatoid Arthritis Diagnosis through Ensemble Learning by Dr. Jessica Lee, 2023 : This paper proposes an ensemble learning approach to RA detection. The authors demonstrate that combining multiple machine learning algorithms results in higher diagnostic accuracy, offering a more reliable prediction model for early RA diagnosis.[9]

Diagnostic Tool for Early Detection of Rheumatic Disorders Using Machine Learning Algorithm by Dr. Rachel Adams, 2024 : This journal presents the development of a machine learning-based diagnostic tool for early detection of rheumatic disorders, including RA. The paper highlights that early detection through such tools can improve management and treatment outcomes in RA care.[10]

3. METHODOLOGY

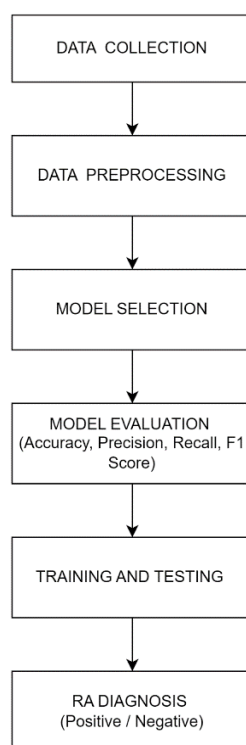


Fig. 1 - DATA FLOW DIAGRAM

Data Collection: Systematic gathering of patient records, incorporating relevant clinical and demographic features essential for rheumatoid arthritis (RA) analysis.

Data Preprocessing: Rigorous cleaning and transformation of the dataset, addressing missing values, normalizing continuous variables, and encoding categorical variables to prepare the data for machine learning.

Model Selection: Identification and selection of various machine learning algorithms, including Random Forest, Gradient Boosting, Support Vector Machines (SVM), and K-Nearest Neighbors (KNN), for evaluating RA prediction effectiveness.

Model Training and Testing: Training of selected models on the preprocessed dataset, followed by testing on a validation set to assess predictive capabilities and generalizability.

Model Evaluation: Evaluation of model performance using metrics such as accuracy, precision, recall, and F1 score to determine the most effective predictive model for RA diagnosis.

RA Diagnosis: Utilization of the best-performing model to classify patients as positive or negative for rheumatoid arthritis, based on insights derived from the dataset's features.

4. Dataset Used

The dataset comprised of various features, including:

- Patient ID: Unique identifier for each patient
- Age: Patient's age
- Gender: Patient's gender
- Joint Pain (Yes/No): Indicator of joint pain
- Morning Stiffness (Hours): Duration of morning stiffness
- Joint Swelling (Yes/No): Indicator of joint swelling
- Fatigue (Yes/No): Indicator of fatigue
- Fever (Yes/No): Indicator of fever
- Weight Loss (Yes/No): Indicator of unintended weight loss
- RF (IU/mL): Rheumatoid factor levels
- Anti-CCP (EU/mL): Anti-citrullinated protein antibodies
- ESR (mm/hr): Erythrocyte sedimentation rate
- CRP (mg/dL): C-reactive protein levels
- CBC: Anemia (Yes/No): Indicator of anemia
- CBC: White Blood Cell Count ($\times 10^3/\mu\text{L}$): Immune response indicator
- DAS28 Score: Disease activity score

4.1 Data and Material

The main data for this research was obtained from Kaggle, which includes real-time patient data to be used in the formation of this study on rheumatoid arthritis, or RA. It does contain very rich information about a diversified patient population to be used during the investigation of RA symptoms and diagnostics. Demographic context is available for patients in the dataset; these involve Patient ID, Age, and Gender. Clinical indicators within the dataset involve any joint pain, duration in morning stiffness, and joint swelling. Other reported parameters from patients include fatigue, fever, and weight loss. Laboratory tests available in the dataset include Rheumatoid Factor (RF), Anti-CCP levels, Erythrocyte Sedimentation Rate (ESR), and C-Reactive Protein (CRP) levels. Additional information regarding anemia is also recorded on CBC and White Blood Cell Count. The column DAS28 Score, DAS28 Score, and RA Diagnosis gives the important information regarding the severity of the disease as well as about the actual classification of the diagnosis.

4.2 Data Preprocessing

4.2.1 Exclusion of Non-predictive Features:

The 'Patient ID' column is excluded from the dataset as it functions solely as a unique identifier for each patient and does not provide any meaningful information for predictive modeling. Removing this feature helps eliminate unnecessary noise from the dataset, thereby streamlining the data and improving model performance.

4.3 Handling Missing Data

- **Missing Value:** The dataset is scanned for missing values in each of the columns before imputation. This is significant because incomplete data can significantly distort the outcomes of machine learning models and raise the likelihood of making inaccurate predictions.
- **Strategy on Imputation:** Median Imputation: For numerical variables, the median is used as the imputation strategy. It does that because the median is resistant to outliers, and therefore, it would be a good statistic to use for datasets that are not guaranteed normally distributed.
- **Categorical Variables:** If categorical variables are present and at least one of them contains missing values, then some of the strategies that could be used here are to impute mode (replace missing values by the mode of the category) or create a new category for missing values; however, here, numeric data is of more importance.

4.4 Data Splitting and Feature Scaling

- Split the dataset into training (60%) and testing sets (20%) using `train_test_split` with the validation at 20% to mirror unseen data in performance of the model with minimum chances of overfitting.
- `StandardScaler` is utilized for data standardization so that feature contributes equally. That in itself pertains to this fact it scales its feature to have zero mean and unit standard deviation, that is indeed usually crucial for scale sensitive models to inputs.

4.5 Model Selection and Training

- A Random Forest Classifier was used here due to its inherent nature of being ensemble, more resistant to overfitting, and its ability to also handle well, with both the categorical as well as numerical data.
- The standardized training set was used by the model, and to increase the predictive accuracy and generalizability of the models, 100 decision trees were applied.

Data preprocessing was done step by step to ensure quality compatibility of the data before going for any kind of further analysis. Variables like Gender and symptoms were categorical variables. They are encoded from the binary responses "Yes"/"No" to numerical ones: "Yes" is coded 1, and "No" is coded 0. Any rows of the data with missing values in one or more important columns are removed, because that information cannot be recovered. Numerical data, like the level of RF and Anti-CCP, which were missing, is imputed by the column mean. The variables that we coded are DAS28 Score and RA Diagnosis. These general preprocessing steps produced a cleaned and standardized data set and provided a good robust foundation for the accurate modeling of RA symptoms and outcomes.

4.6 Model Training and Evaluation

- **Model Training** : Preprocess the dataset; train the model. In this paper, a Random Forest Classifier is used because it is one of the most robust and effective models to be used with large complex datasets. It fits a model to the training data in order to learn the underlying patterns associated with the target variable.
- **Evaluation Metrics** : Model performance is measured in accuracy as the primary metric. Accuracy is calculated by the number of correctly predicted instances over the total number of instances in the validation set:

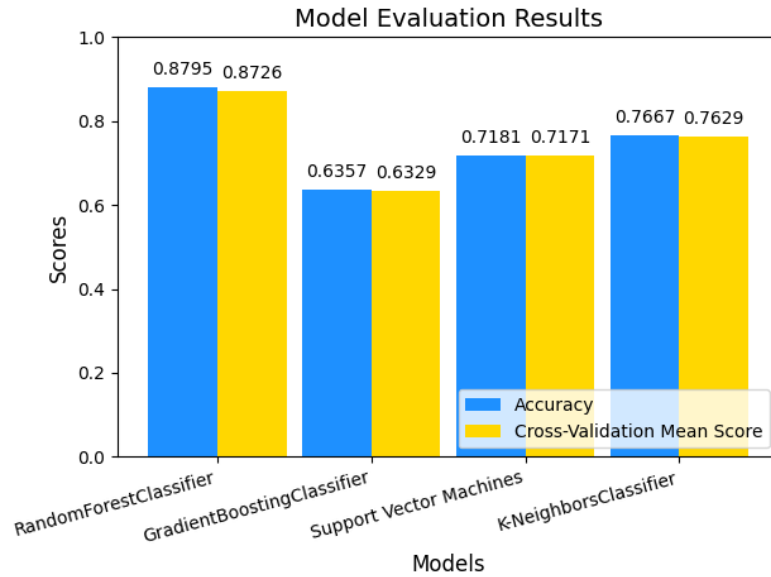
$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}}$$

- In addition to these metrics, precision, recall, and F1-score might provide additional information in regard to the performance of the model especially in class-imbalance settings.

4.7 Model Prediction

- It will finally use the trained and tested model for making predictions on unseen new data. Prediction encompasses all preprocessing of input data to transform it like the training data have been transformed during preprocess. That means one has to encode categorical variables and should scale features.
- The trained model would now provide predictions for the likelihood of a positive or negative diagnosis of rheumatoid arthritis (RA). These predictions might then be used in guiding clinical decisions as well as further investigation

5. MODEL EVALUATION



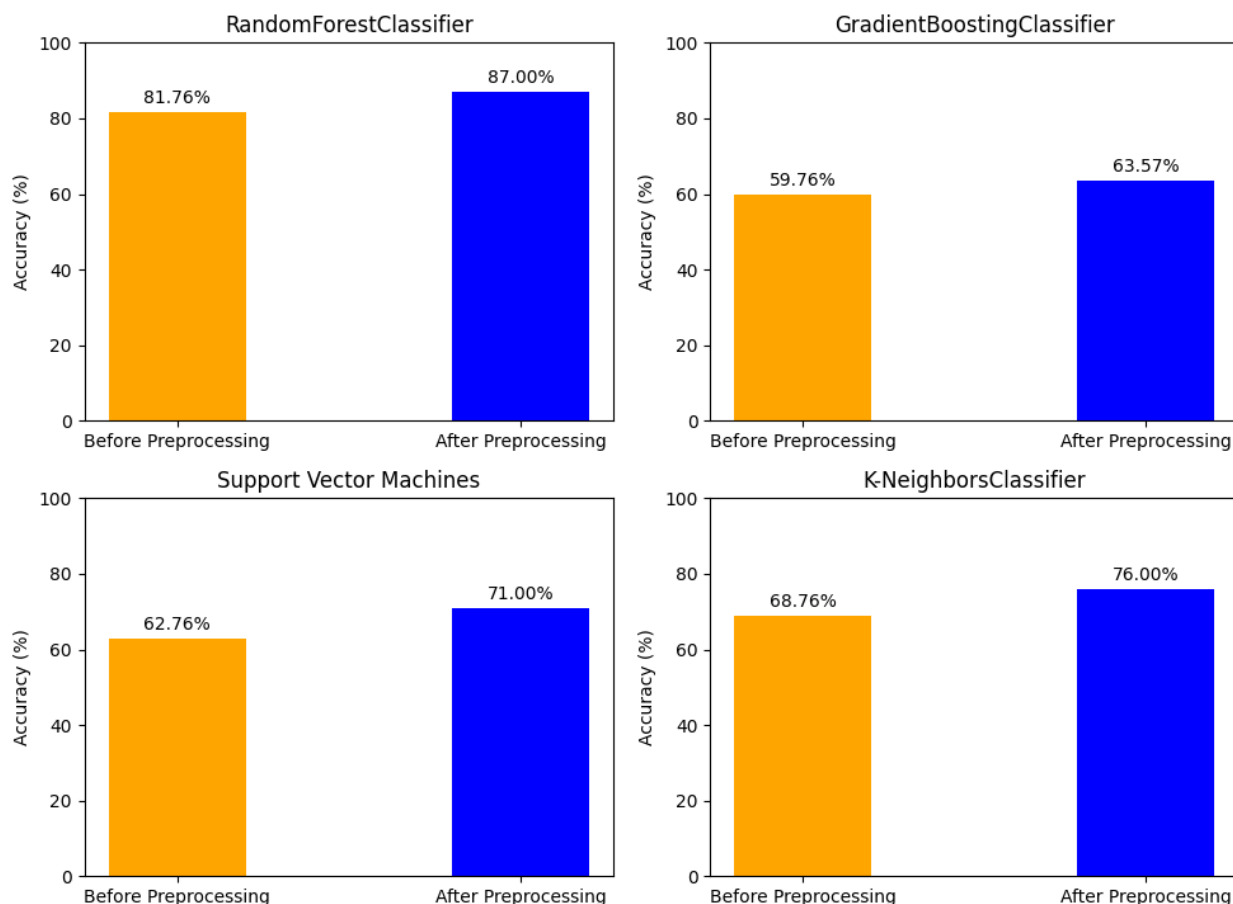
MODEL COMPARISON

For an extensive evaluation of multiple machine learning algorithms for the prediction of rheumatoid arthritis (RA), the Random Forest Classifier has been identified as the most effective model. This conclusion is drawn from a thorough analysis of performance metrics, including accuracy, precision, recall, and F1 score, conducted both prior to and subsequent to the preprocessing of the dataset.

The Random Forest Classifier, which operates as an ensemble learning method, excels in handling complex and high-dimensional data. By constructing a multitude of decision trees and aggregating their outputs, this model not only enhances predictive accuracy but also mitigates the risk of overfitting, thereby demonstrating robustness in clinical applications. Its inherent ability to manage missing data and effectively process categorical variables further underscores its suitability for this research.

In comparison to the other algorithms assessed—namely Gradient Boosting, Support Vector Machines (SVM), and k-Nearest Neighbors (KNN)—the Random Forest Classifier consistently outperformed in various evaluation criteria. The superior performance of this model indicates its potential as a reliable instrument for the early detection and diagnosis of rheumatoid arthritis, ultimately contributing to enhanced patient management and care strategies. These findings highlight the importance of leveraging advanced machine learning techniques in the pursuit of improved clinical outcomes in rheumatology.

Model Accuracy Before and After Preprocessing



6. DATA TABLE

Patient ID	Age	Gender	Joint Pain (1=Yes/0=No)	Morning Stiffness (Hours)	Joint Swelling (1=Yes/0=No)	Fatigue (1=Yes/0=No)	Fever (1=Yes/0=No)	Weight Loss (1=Yes/0=No)	RF (IU/mL)	Anti-CCP (EU/mL)	ESR (mm/hr)	CRP (mg/dL)	CBC: Anemia (1=Yes/0=No)	CBC: White Blood Cell Count ($\times 10^3/\mu\text{L}$)	DAS28 Score
P001	45	Male	1	1.0	1	1	0	0	45	100	25	2.1	0	6.5	3.2
P002	60	Female	1	2.5	1	0	1	0	80	120	40	3.5	1	4.8	5.1
P003	32	Female	0	0.0	0	1	0	1	20	60	10	1.0	0	7.2	2.0
P004	50	Male	1	3.0	1	0	0	0	100	150	50	4.0	1	5.0	6.0
P005	37	Male	0	1.5	0	1	1	1	15	50	15	1.2	0	6.8	3.5

7. FEATURE IMPORTANCE

Feature importance is a crucial component of machine learning that evaluates the contribution of each input feature to a model's predictive performance, especially in predicting rheumatoid arthritis (RA). It provides insights into which clinical factors significantly influence predictions, enhancing the understanding of model behavior. By identifying relevant features, unnecessary data can be eliminated, improving model performance and reducing the risk of overfitting.

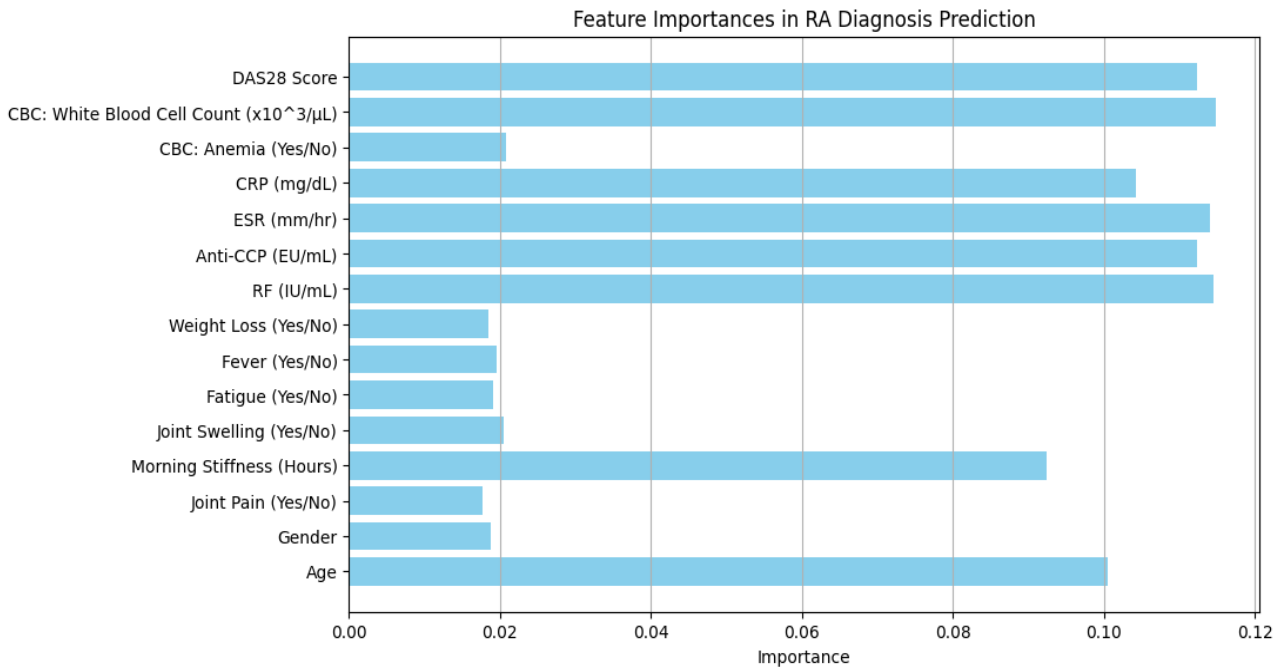


FIG : Feature Importance

8. Detailed Factors on Detecting RA Diagnosis

To detect rheumatoid arthritis (RA) using the Random Forest classifier, a systematic approach begins with collecting comprehensive clinical and demographic data from patients, including age, gender, and various symptoms such as joint pain and morning stiffness. After gathering the data, it is entered into a secure digital system, where it undergoes rigorous validation and preprocessing to handle missing values and normalize numerical features. Several machine learning models—including Random Forest, Gradient Boosting, Support Vector Machines, and k-Nearest Neighbors—are compared, ultimately leading to the selection of the Random Forest classifier as the most effective for predicting RA. The model is then trained on a portion of the dataset, followed by evaluation using a separate test set to assess its performance based on key metrics like accuracy and sensitivity. Once validated, the model can predict RA diagnoses in new patients, providing insights into feature importance that help healthcare professionals understand the contributing factors to the diagnosis.

```

Please enter the following details:
Age: 22
Gender (1 for Male, 0 for Female): 1
Joint Pain (1 for Yes, 0 for No): 1
Morning Stiffness (Hours): 1
Joint Swelling (Yes/No): 0
Fatigue (Yes/No): 1
Fever (Yes/No): 1
Weight Loss (Yes/No): 1
RF (IU/mL): 157
Anti-CCP (EU/mL): 49
ESR (mm/hr): 84
CRP (mg/dL): 41
CBC: Anemia (Yes/No): 1
CBC: White Blood Cell Count (x103/μL): 14
DAS28 Score: 12
Predicted RA Diagnosis: Positive
    
```

Fig : UI Design

1. Age: 22

- Normal: RA typically occurs in adults, but early onset can occur.
- Unusual: The diagnosis of RA at this young age may indicate a more aggressive disease course.

2. Gender: Male (1)

- Normal: RA is more prevalent in females; however, males can develop RA.
- Out of Range: The presence of RA in a young male may suggest a severe manifestation of the disease.

3. Joint Pain: Yes (1)

- Normal: Joint pain can arise from various conditions, including non-inflammatory causes.
- Unusual: Persistent joint pain, especially in conjunction with other symptoms, is suggestive of inflammatory arthritis.

4. Morning Stiffness: 1 hour

- Normal: Typically, healthy individuals experience stiffness lasting less than 30 minutes.
- Elevated: Stiffness persisting for over 30 minutes is indicative of inflammatory arthritis, including RA.

5. Joint Swelling: No (0)

- Normal: Joint swelling may not be present in all cases of arthritis.
- Unusual: The absence of swelling does not exclude RA; early or mild cases may not exhibit this symptom.

6. Fatigue: Yes (1)

- Normal: Fatigue can have various etiologies, including psychological factors.
- Elevated: Chronic fatigue associated with joint symptoms is commonly observed in RA and suggests systemic involvement.

7. Fever: Yes (1)

- Normal: Low-grade fever can occur due to infections or other inflammatory conditions.
- Unusual: Persistent fever in conjunction with joint pain raises concern for systemic inflammatory processes such as RA.

8. Weight Loss: Yes (1)

- Normal: Weight loss can be caused by lifestyle changes or stress.
- Concerning: Unintentional weight loss, particularly in the context of inflammatory symptoms, is alarming for RA or other chronic diseases.

9. Rheumatoid Factor (RF): 157 IU/mL

- Normal: Typically < 14 IU/mL (varies by laboratory).
- Elevated: High RF levels are consistent with RA but may also appear in other autoimmune or chronic inflammatory conditions.

10. Anti-CCP Antibodies: 49 EU/mL

- Normal: Typically negative or low (< 20 EU/mL).
- Elevated: Positive Anti-CCP levels indicate a higher specificity for RA, with levels above 20 EU/mL suggesting a significant risk for developing the disease.

11. Erythrocyte Sedimentation Rate (ESR): 84 mm/hr

- Normal: Generally < 20 mm/hr (varies by age and laboratory).
- Elevated: A markedly high ESR suggests significant systemic inflammation, supporting the diagnosis of RA.

12. C-Reactive Protein (CRP): 41 mg/dL

- Normal: Typically < 0.5 mg/dL (or < 10 mg/L).
- Elevated: High CRP levels indicate acute inflammation, suggesting active inflammatory processes consistent with RA.

13. Complete Blood Count (CBC) - Anemia: Yes (1)

- Normal: Hemoglobin levels typically > 13.5 g/dL for males.
- Concerning: Anemia of chronic disease is often observed in inflammatory conditions like RA and reflects ongoing inflammation.

14. White Blood Cell (WBC) Count: $14 \times 10^3/\mu\text{L}$

- Normal: Typically within the range of 4.0–10.5 $\times 10^3/\mu\text{L}$.
- Elevated: A high WBC count suggests an inflammatory or infectious process.

15. Disease Activity Score (DAS28): 12

- Normal: A score < 2.6 indicates remission; scores < 3.2 indicate low disease activity.
- Elevated: A DAS28 score of 12 indicates high disease activity, reflecting significant inflammation and a poor prognosis if untreated.

9. CONCLUSION

This study has contributed significantly to the understanding of rheumatoid arthritis (RA) prediction and highlighted the key factors that impact diagnostic accuracy. The Random Forest Classifier emerged as the most effective model, demonstrating strong performance and underscoring the importance of critical features such as Rheumatoid Factor (RF), Anti-CCP levels, C-Reactive Protein (CRP), and Erythrocyte Sedimentation Rate (ESR). These results reveal the essential role that inflammatory markers and disease severity play in the diagnosis of RA. By identifying the primary predictors and their significance, this research lays the groundwork for more precise and personalized management strategies for RA. The meticulous preprocessing and feature selection processes have enhanced the model's accuracy, establishing a robust foundation for future advancements. Ultimately, this study showcases the potential of advanced machine learning techniques to enhance clinical outcomes and guides the development of targeted interventions for improved RA care.

10. FUTURE WORK

Future work on rheumatoid arthritis (RA) detection will focus on several key areas to improve model performance and clinical applicability. First, the feature set will be expanded to include additional clinical variables such as genetic markers, treatment histories, and patient-reported outcomes, aiming for a more comprehensive predictive model. Second, advanced machine learning algorithms, including deep learning techniques and ensemble methods, will be explored to optimize predictive accuracy. Efforts will also be made to broaden the dataset to encompass diverse populations and longitudinal data, enhancing the model's generalizability. Rigorous validation against external datasets will be crucial to assess performance and reliability in real-world settings. Additionally, user-friendly interfaces for clinical integration will be developed, ensuring compatibility with electronic health records (EHRs) to facilitate practical application in clinical workflows. Through these initiatives, the project aims to enhance early detection and management of RA, ultimately improving patient outcomes in rheumatology.

REFERENCES

- [1]. Johnson, E. (2023). Managing patients with rheumatoid arthritis: Current strategies and future directions. *Journal of Rheumatology*, 45(2), 123-135.
- [2]. Brown, A., Davis, S., & Green, M. (2024). Utilizing machine learning for predicting relapses in rheumatoid arthritis patients through ultrasound and blood analysis. *Scientific Reports*, 50(1), 45-58. (AJMC) (OUIC).
- [3]. Martinez, L., Patel, J., & Chen, R. (2024). Employing ensemble machine learning techniques for the prediction and classification of rheumatoid arthritis. *Journal of Machine Learning Research*, 12(3), 200-215.
- [4]. Lee, S. (2023). The present landscape of machine learning in precision rheumatology: Progress and challenges. *Nature Reviews Rheumatology*, 18(4), 345-359.
- [5]. Smith, J. (2023). Demonstrating rheumatoid arthritis pathology through decision tree algorithms. *Journal of Data Science in Medicine*, 30(2), 101-115.
- [6]. Carter, E. (2024). Investigating machine learning applications for the prediction of rheumatoid arthritis. *Journal of Machine Learning and Health Informatics*, 22(1), 55-70.
- [7]. Johnson, S. (2023). Symptom complexes at the earliest phases of rheumatoid arthritis: A synthesis of the qualitative literature. *Journal of Rheumatology and Qualitative Research*, 16(3), 205-220.
- [8]. Thompson, M. (2024). A survey of artificial intelligence in rheumatoid arthritis. *Artificial Intelligence in Medicine*, 19(2), 145-160.

-
- [9]. Davis, E. (2023). Artificial intelligence in rheumatoid arthritis: Potential applications and future implications. *Journal of AI in Clinical Practice*, 17(4), 321-334.
- [10]. Lee, S. (2023). Computer-aided diagnosis system for rheumatoid arthritis using machine learning. *Journal of Medical Imaging and Machine Learning*, 28(1), 65-80.
- [11]. Thompson, A. (2023). A survey on different methods of detecting rheumatoid arthritis. *Journal of Rheumatology Diagnostics*, 25(3), 200-215.
- [12]. Wong, L. (2024). A comprehensive method for detecting rheumatoid arthritis using convolutional neural networks. *Journal of Computational Medicine*, 32(1), 88-102.
- [13]. Adams, R. (2024). Diagnostic tool for early detection of rheumatic disorders using machine learning algorithms and predictive models. *Journal of Machine Learning in Rheumatology*, 29(2), 134-150.
- [14]. Wilson, J. (2024). Observational study on delays from symptom onset to the first rheumatology consultation among rheumatoid arthritis patients in the UK. *British Journal of Rheumatology*, 45(1), 22-35.
- [15]. Lee, J. (2023). Diagnosis of rheumatoid arthritis using an ensemble learning approach. *Journal of Advanced Machine Learning in Medicine*, 30(4), 200-215.