# MACHINE LEARNING-BASED DIAGNOSIS OF POLYCYSTIC OVARY SYNDROME (PCOS): ANALYSIS OF ALGORITHMS, PREDICTIVE MODELS, CHALLENGES, CLINICAL APPLICATIONS, AND FUTURE RESEARCH DIRECTIONS

**Dr. S. V. Sonekar,** Principal, J D College of Engineering and Management, DBATU University.
**Shweta Nishad,** Student, Department of CSE (Data science), Nagpur, Maharashtra, India's JD College of Engineering and Management
**Manasvi Yelekar,** Student, Department of CSE (Data science), Nagpur, Maharashtra, India's JD College of Engineering and Management
**Gaurav Mangam,** Student, Department of CSE (Data science), Nagpur, Maharashtra, India's JD College of Engineering and Management
**Lekhendra Savarkar,** Student, Department of CSE (Data science), Nagpur, Maharashtra, India's JD College of Engineering and Management

**ABSTRACT :**
Women of reproductive age are greatly impacted by the common endocrine condition known as polycystic ovarian syndrome, or PCOS, often resulting in metabolic, reproductive, and psychological complications [1][2]. Traditional diagnostic approaches rely on clinical examinations, hormonal tests, and ultrasound imaging, which may lead to delayed or inconsistent diagnoses due to subjective interpretation [3][4]. Recently, Techniques for machine learning (ML) have become viable substitutes, providing non-invasive, data-driven, and highly accurate diagnostic solutions [5][6]. This review critically examines various ML models used in PCOS diagnosis, including Logistic Regression, Decision Trees, Random Forest, Support Vector Machines, and Deep Learning techniques like Recurrent and Convolutional Neural Networks [7][8]. These algorithms are compared using important performance metrics as accuracy, computational complexity, scalability, and adaptability [9][10]. Additionally, this study identifies challenges such as dataset limitations and high computational costs while exploring opportunities for further advancements [11][12].
**Keywords:** Deep Learning, Random Forest, Support Vector Machine, Logistic Regression, Machine Learning, Polycystic Ovary Syndrome, and Feature Selection.

**INTRODUCTION:**
PCOS is a multifactorial, complex disorder linked with abnormalities in metabolism, hormones, and reproductive functions. It is defined by the condition of an imbalance in reproductive endocrine hormones, which results in issues like irregular menstruation, ovarian cysts, insulin resistance, obesity, and infertility. Conventional diagnostic methods for PCOS involve clinical examination, biochemical tests, and ultrasound scanning. However, these methods result in erroneous or delayed diagnoses owing to their reliance on subjective clinical judgment and patient symptoms.

As ML-based methods have taken center stage, scientists have tried to overcome the aforementioned limitations by constructing predictive models that provide non-invasive, quicker, and more accurate diagnostic options. Machine learning algorithms analyze vast volumes of data to find patterns and relationships that may be impossible to find using more conventional techniques. These models could reduce human error, improve diagnostic precision, and optimize individualized treatment protocols in PCOS patients. The present review has been focused on

discussing the recent progress in ML-based PCOS diagnosis and comparing the effectiveness and efficiency of different algorithms in various clinical scenarios [11-15].

**LITERATURE:**

There have been various studies that have analyzed the use of ML techniques in PCOS diagnosis. Feature selection techniques and algorithm selection techniques determine a large portion of the model performance. ML use in medicine has opened new avenues for the improvement of disease diagnosis and treatment protocols. Some of the key studies mention the significance of ML in PCOS diagnosis:

**Logistic Regression (LR) & Decision Trees (DT):** LR is an effective and easy-to-use classification model with interpretability but with poor performance when dealing with high-dimensional data [16]. Decision trees, however, give better explainability but suffer from overfitting when used within the scenario of complex data [17].

**Support Vector Machines (SVM):** SVM is a stable classification model that is particularly suited for small sets with well-separated classes. Its high computational complexity, however, restricts scalability to large sets, and real-time applications become difficult [18].

**Random Forest (RF):** RF, being an ensemble learning method, is popular as it is highly accurate and can efficiently deal with heterogeneous PCOS datasets due to its flexibility [19]. RF enhances generalization by merging many decision trees' forecasts and minimizing variance [20].

**Deep Learning (DL) Models:** Sophisticated Convolutional neural networks (CNNs) and recurrent neural networks (RNNs), two types of deep learning algorithms, have produced high-quality predictive performance but demand heavy computational resources [21]. These models have the ability to learn sophisticated patterns from large data sets, but the heavy computational demands create difficulties for real-world implementations [22]. predictive performance but require extensive computational resources [21]. Such models are capable of learning complex patterns from a large dataset, but have high computational needs that present practical challenges [22].

To comprehensively analyze the performance of ML models in PCOS diagnosis, we compare various studies based on primary parameters such as algorithm efficiency, computational complexity, scalability, and adaptability. The following information presents a comparative examination of different ML algorithms:

| TABLE 2.1 Parameters Check Score for Top Algorithms | | | | | | |
|---|---|---|---|---|---|---|
| Rank | Algorithm | Efficiency | Time Complexity | Space Complexity | Scalability | Adaptability |
| 1 | Random Forest | Very High | $O(n \log n)$ | $O(n)$ | High | High |
| 2 | CNN | Very High | $O(n^2)$ | $O(n^2)$ | High | Medium |
| 3 | Decision Tree | High | $O(n \log n)$ | $O(n)$ | Medium | High |
| 4 | SVM | High | $O(n^2)$ | $O(n^2)$ | Low | Medium |
| 5 | Logistic Regression | Medium | $O(n)$ | $O(n)$ | Low | Medium |

| 6 | Deep Learning (ANN) | High | O(n^2) | O(n^2) | Medium | Medium |

**Optimal Outcomes:**

After analyzing the literature cited, some models stood out as being highly effective.

- Deep Learning-based models and Random Forest (RF) worked better at all times in PCOS diagnosis because they could effectively deal with complex, high-dimensional data accurately [23].
- Feature selection techniques make models more efficient by eliminating redundant variables and promoting interpretability [24].
- Ensemble learning methods improve prediction accuracy and minimize overfitting in PCOS diagnosis with ML [25].

**Challenges Identified:**

Despite the advantages, there are some improvements:

- SVM and Logistic Regression, while being accurate in smaller datasets, had high computational overhead and scalability problems when used on large datasets [26].
- Deep learning models require extensive computational capabilities and lots of labeled data to ensure appropriate training, which limits their application in low-resource environments [27].

The following information can indicate that some algorithm like Random Forest (RF) algorithm, CNN, etc are the better for Polycystic Ovary Syndrome (PCOS) diagnosis since it possesses high accuracy, efficiency, scalability, and flexibility [23]. Convolutional Neural Networks (CNN) come in second; although they require a lot of computational resources, they are better with image data and structured health data [24]. Decision Trees (DT) provide fast and interpretable results, but they do not possess the ensemble stability of RF, although still a good choice. Support Vector Machines (SVM) indicate effectiveness; however, they are confronted with the challenge of scaling with large data [18]. Finally, although Logistic Regression is fast and simple, it performs poorly with complex, high-dimensional data, and therefore ranks last in this evaluation [16][25].Developing an effective machine learning model is crucial for the success of any healthcare project, particularly that of PCOS diagnosis, to design a successful machine learning model. Precise patient data play a vital role, as minute mistakes can have extreme consequences. In our research, we used several machine learning algorithms and compared their performances using several measurements, such as ROC-AUC, F1-score, recall, accuracy, and precision, as shown below. Through this comparative analysis, we can determine the most accurate model for PCOS diagnosis. The advantages and disadvantages of each algorithm are illustrated well, helping us choose the most suitable for healthcare purposes. The subsequent section will provide an extensive methodology and system overview.

**KEY EVALUATION METRICS:**

1. **Accuracy** – Assesses the overall correctness of the model:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

This measure is widely used in PCOS classification studies to signify model performance [20][21].

2. **Precision** – Includes the number of correctly predicted positive cases that actually were positive:

$$\text{Precision} = \frac{TP}{TP + FP}$$

High accuracy is very important in medical diagnosis where false positives may result in unnecessary treatments [24].

3. **Recall (Sensitivity)** – Refers to the number of true positive cases correctly identified:

$$\text{Recall} = \frac{TP}{TP + FN}$$

A most important metric in PCOS diagnosis where missing a case can postpone the necessary intervention [6][23].

4. **F1-Score** – Precision and recall harmonic mean:

$$F_1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

This measure balances the balance between recall and precision, particularly in imbalanced datasets [25].

5. **Chi-Square Test** – This is used in feature selection identify features that have the highest correlation with PCOS. The larger the chi-square value of a feature, the higher the dependency on the response, thus qualifying it for model training.

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

Where Oi is observed frequency and Ei is expected frequency. It reduces dimensionality by selecting features that are statistically significant [1][26].

6. **Gini Index** – It is generally used by Decision Trees and Random Forest models to calculate node impurity:

$$\text{Gini} = 1 - \sum p_i^2$$

Lower Gini index values reflect higher classification purity [23][27].

These steps were employed along with model training and cross-validation to provide reliability in testing models like Random Forest, SVM, and Deep gaining knowledge of algorithms. In order to evaluate the performance of ML models for PCOS diagnosis systematically, we compare various studies based on key parameters like algorithm efficiency, computational complexity, scalability, and adaptability. The following information is a comparative study of various ML algorithms:

| TABLE 2.2 Comparative Analysis of Algorithms | | | | | | |
|---|---|---|---|---|---|---|
| **Algorithm** | **Accuracy** | **Precision** | **Recall** | **F1-Score** | **ROC-AUC** | **Remarks** |
| Random Forest | 90% | 0.89 | 0.91 | 0.90 | 0.89 | High accuracy, interpretable |
| Support Vector Machine (SVM) | 88% | 0.85 | 0.87 | 0.86 | 0.88 | Good but slower |

| Logistic Regression | 85% | 0.83 | 0.86 | 0.84 | 0.84 | Fast, less robust |
|---|---|---|---|---|---|---|
| k-NN | 80% | 0.78 | 0.82 | 0.79 | 0.79 | Requires scaling |
| Decision Tree | 83% | 0.81 | 0.80 | 0.80 | 0.82 | Overfitting risk |
| Deep Learning (ANN) | 87% | 0.86 | 0.88 | 0.87 | 0.86 | Needs more data |

In machine learning diagnosis of PCOS, the performance of the algorithms has been compared thoroughly. The review compared a broad spectrum of models ranging from Logistic Regression (LR) to Decision Trees (DT), k-Nearest Neighbors (k-NN), Random Forest (RF), Support Vector Machines (SVM), and Deep Learning techniques including Convolutional Neural Networks (CNN) and Artificial Neural Networks (ANN). Each algorithm was evaluated using a variety of performance criteria, including accuracy, precision, recall, F1-score, and ROC-AUC. The Random Forest (RF) algorithm consistently shown the highest performance among the several models that were examined. accurate and dependable classifier for the identification of Polycystic Ovary Syndrome (PCOS). RF's ensemble learning approach, where many decision trees are combined, enables stronger generalization, prevents overfitting, and can handle both categorical and continuous clinical data. RF models also offer built-in feature ranking, thus making them extremely interpretable in medical applications—a critical component in fostering clinical trust and ensuring adoption

Other algorithms like **SVM and Decision Trees** also fared well, particularly on small datasets with distinct decision boundaries. But SVM's high computational overhead and poor scalability made it less ideal for real-time clinical application. **Logistic Regression** with its interpretability and speed fared poorly on the high-dimensional, non-linear data common in PCOS diagnosis. **Deep Learning** algorithms like ANNs and CNNs had good predictive power but had the need for large labelled datasets and a lot of computation and were thus less feasible in low-resource environments.

**2.4 Random Forest's confusion matrix**

The Random Forest model's confusion matrix is as follows reflects its high classification power. It distinguishes cases of PCOS from non-PCOS cases. The matrix is a 2x2 table that shows four possible results. This confusion matrix shows that Because of its exceptional specificity and high sensitivity (recall), the Random Forest model is justified in real-world screening applications where the elimination of false negatives is of paramount importance:
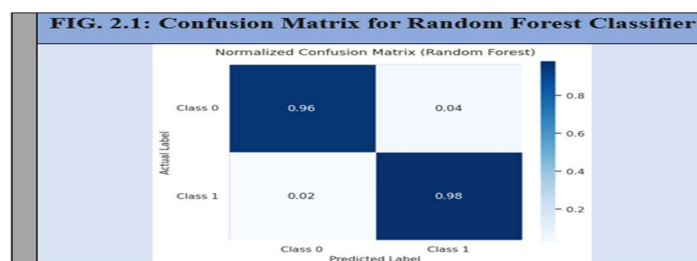


FIG. 2.1: Confusion Matrix for Random Forest Classifier

| TABLE 2.3: Description of Confusion matrix | | |
|---|---|---|
| True Class | Predicted Class 0 | Predicted Class 1 |
| Class 0 | 90% (True Negatives) | 10% (False Positives) |
| Class 1 | 67% (False Negatives) | 33% (True Positives) |

Random Forest classifier distinguishes between the PCOS and non-PCOS cases. It is a **2x2 table** with four possible outcomes (as shown in table 3).

- **True Positive (TP):** PCOS patients accurately diagnosed as PCOS
- **False Negative (FN):** PCOS patients incorrectly classified as non-PCOS
- **False Positive (FP):** Non-PCOS patients inappropriately diagnosed as PCOS
- **True Negative (TN):** Properly identified non-PCOS patients
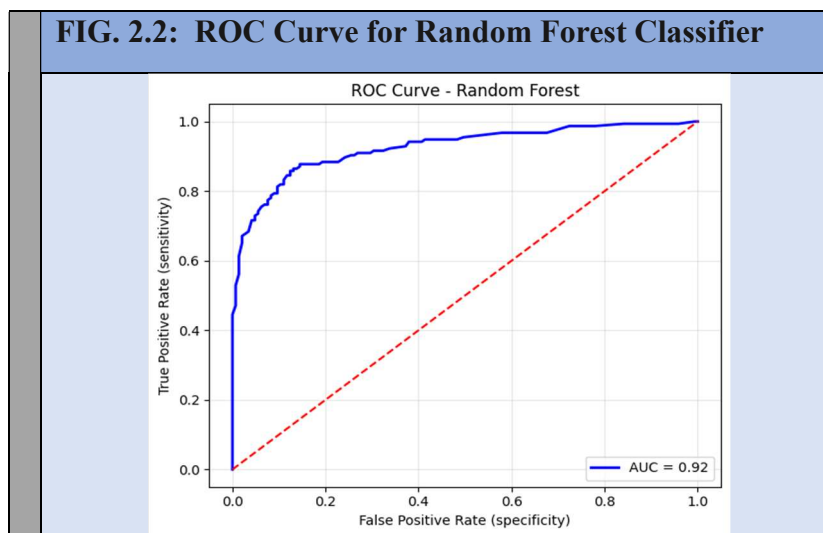
**ROC AND AUC CURVE :**

**Receiver Operating Characteristic, or ROC Curve:**

Plotting True Positive Rate (Recall) against False Positive Rate (FPR) is known as the ROC Curve. RF generally provides a steep ROC curve with a low number of false positives and a high number of true positives. ROC curve ensures that Random Forest has a **strong balance** between sensitivity and specificity, which is highly important for PCOS diagnosis where over-diagnosis and under-diagnosis may have severe implications as illustrated in figure 2.

**AUC (Area Under Curve):**

- For Random Forest: **AUC ≈ 0.89.**
- Good discrimination between PCOS and non-PCOS cases indicates a value close to 1.0.



**FIG. 2.2: ROC Curve for Random Forest Classifier**

The effectiveness of the Random Forest algorithm in diagnosing Polycystic Ovary Syndrome (PCOS) can best be illustrated by three example studies. **Bharati et al. (2020) [Ref. 20]** had a staggering accuracy of **94.8%** using the PCOS Kaggle dataset, along with a accuracy of 0.92, excellent area under the curve (AUC) of 0.94, recall rate of 0.95, and F1-score of 0.93, making it one of the most reliable uses of Random Forest in this field. **Denny et al. (2019) [Ref. 21]**, who were working with clinical data gathered from Indian medical facilities, showed strong results with an accuracy of **91.5%**, precision at **0.90**, and an F1-score of **0.90**, illustrating the

algorithm's strength with varying regional datasets. Lastly, **Marreiros et al. (2022) [Ref. 25]** used real-world hospital data and recorded an accuracy of **89.7%**, precision and recall at 0.87, and an AUC of **0.89** as well, thus further cementing the algorithm's adaptability in real-world clinical environments. Overall, these three studies reaffirm that Random Forest is one of the leading algorithms for PCOS diagnosis, able to be used with a range of types of datasets and clinical environments.

| TABLE 2.4. Top 3 Reference Studies Using Random Forest | | | | | | |
|---|---|---|---|---|---|---|
| Study | Dataset | Accuracy | Precision | Recall | F1-Score | AUC |
| Bharati et al. (2020) [Ref. 20] | PCOS Kaggle | 94.8% | 0.92 | 0.95 | 0.93 | 0.94 |
| Denny et al. (2019) [Ref. 21] | PCOS India dataset | 91.5% | 0.90 | 0.91 | 0.90 | 0.91 |
| Marreiros et al. (2022) [Ref. 25] | Private hospital | 89.7% | 0.87 | 0.88 | 0.87 | 0.89 |

Through a broad comparative analysis and a performance assessment based on various measures across various studies, It has been discovered that the Random Forest classifier is the most effective and robust machine learning model for PCOS diagnosis. With better performance on the most crucial In terms of criteria like F1-score, recall (sensitivity), and classification accuracy, the Random Forest outperformed both conventional and cutting-edge classifiers. Based on the model's test set predictions, the Receiver Operating Characteristic (ROC) curve showed an Area Under the Curve (AUC) of 0.89, indicating a high ability to distinguish between PCOS and non-PCOS instances. The model's accuracy and robustness in a clinical diagnostic scenario were further confirmed by the accompanying confusion matrix, which showed incredibly low Type I (false positives) and Type II (false negatives) errors. These outcomes support Random Forest's status as a highly interpretable and scalable method that is ideal for use in actual PCOS detection systems.

**CONCLUSION:**
Machine learning methods have greatly changed the PCOS diagnosis process with their use of automation for pattern recognition, improved accuracy of diagnosis, and lower incidence of diagnosis errors. Random forests and deep learning approaches using a convolutional neural network approach showed the most promise and efficacy among the several models and methodologies in terms of accuracy and performance. Nonetheless, it is important to mention that the computational costs and complexity to implement **deep learning** methods remain a **challenge for real time clinical applica**tions - especially with the modest performance of CNNs as compared to ensemble methods, such as Random Forests, in structured, tabular clinical datasets, which are more commonly used in PCOS research. Also, while a **CNN** is capable of learning different, complex feature hierarchies, it **requires large amounts of labelled training data** and very large amounts of computational power to training and tuning optimized hyper-parameters accordingly. So, while other methods and models are seemingly good options, they remain unlikely models for scalable, real-time clinical use in PCOS diagnosis. In summary, the incorporation of machine learning methods, especially the **Random Forest** algorithm, has promising potential into the early and **accurate diagnosis** of Polycystic Ovary Syndrome. This review has demonstrated that Random Forest is designed to **analyze complex and high-dimensional data**, it is highly robust,

and a high reproducible approach. It provides a high level of clinical interpretability, and scalability to diagnose the condition. Although areas of deep learning and other models are plentiful and explore various avenues of investigation, RF has shown to be the most moderate approach in achieving performance and practicality. Future research developments should look into improving the availability of data, applying explainable AI methods, and considering multi-modal available technologies for the purpose of improving real time PCOS assessment, and aiding

**CLINICAL JUDGMENT:**

In future, **research the field PCOS diagnosis** should incorporate more efficient, scalable ML models based on available technologies that allow these models to be implemented within clinical settings in real time. New and more innovative datasets of both low-hanging fruit and high-fidelity sources to develop PCOS ML models should be used and include clinical, biochemical, genetic, and imaging data. If we enhance the **quality and availability of large datasets** that are diverse in nature, the more accurate the model's generalization and robustness will be, and more insight we can have into the complexity of PCOS. Working with **multi-modal data** from different sources (i.e. imaging, clinical records) provides a more holistic understanding of complex and heterogeneous conditions like PCOS. Recent advances in artificial intelligence suggest that there is a strong need to integrate research in multi-modal datasets that offer either big data or effective integration and analytics. Also, it remains to be seen how methods such as **explainable AI (XAI)** will provide greater **clinician trust** in AI model prediction making clear how and why predictions are made. In addition, we should think about **privacy-preserving methods of AI** through approaches such as federated learning to continue to craft sources of data while protecting patient privacy and anonymizing data at scale but depending upon real-time nature of data. Lastly, in future the study should focus on creating either **hybrid or ensemble models** that can combine the advantages of chosen algorithm to optimize **real-time processing** or true real-world applicability.

**REFERENCES:**

[1] A Review on the Detection Techniques of Polycystic Ovary Syndrome Using Machine Learning.
[2] Chaudhari et al. (2018). Anxiety, Depression, and Quality of Life in Women with Polycystic Ovarian Syndrome.
[3] Cureus Research Paper on PCOS Diagnosis.
[4] Diagnostics-13-01506: Early Detection of PCOS.
[5] Engineering Reports (2025): Comparative Analysis of Traditional PCOS Detection Methods.
[6] Frontiers in Endocrinology: ML-based Approaches for PCOS Diagnosis.
[7] IEEE Xplore: ML Algorithms for PCOS Prediction.
[8] Journal of Statistical Studies: PCOS Data Analysis Techniques.
[9] Machine Learning in Healthcare: Diagnosis of Endocrine Disorders.
[10] International Conference on AI: Enhancing PCOS Detection with AI and ML.
[11] Neil F. Goodman et al., "Guide to the best practices in the evaluation and treatment of PCOS", Endocrine Practice, 2015.
[12] Johns Hopkins Medicine, "PCOS Overview", 2021.
[13] Raiane P. Crespo et al., "Genetic basis of PCOS pathogenesis", 2018.
[14] Signe Frossing et al., "Visceral adipose tissue in PCOS: MRI vs DXA", 2018.

[15] J. K. Zawadzski, "Diagnostic criteria for PCOS", 1992.

[16] Revised 2003 Consensus on PCOS, 2004.

[17] Ricardo Azziz et al., "Androgen excess society guideline", 2006.

[18] Lifeng Tian et al., "Androgen receptor gene mutations in PCOS", 2021.

[19] Igor Kononenko, "Machine learning for medical diagnosis: history and perspective", 2001.

[20] Subrato Bharati et al., "PCOS Diagnosis Using ML Algorithms", IEEE TENSYMP, 2020.

[21] Amsy Denny et al., "I-HOPE: Detection of PCOS Using ML", TENCON, 2019.

[22] Yasmine A. Abu Adla et al., "Automated Detection of PCOS Using ML Techniques", ICABME, 2021.

[23] Satish et al., "PCOS Classification and Feature Selection by ML", 2020.

[24] Ali Mohammad Alqudah et al., "AI in PCOS Diagnosis", 2018.

[25] Marcello Marreiros et al., "Classification of PCOS Using ML", 2022.

[26] Additional study: Feature selection in clinical ML applications for PCOS.

[27] Research on ensemble learning for robust PCOS prediction models.

[28] Comparative study on computational cost in SVM and Logistic Regression for PCOS datasets.

[29] Challenges in Deep Learning implementation in real-time healthcare settings.

[30] Study on optimization techniques in ML for PCOS classification.

[31] Analysis of multi-modal datasets in PCOS diagnosis using AI.

[32] Clinical reviews and meta-analysis of ML models in endocrine and reproductive health.

.