

Dr. S. Venkatalakshmi<sup>1</sup>,  
Dr. M. Regina<sup>2</sup>

## Leveraging Hybrid Machine Learning Models for Early Diagnosis and Prediction of Polycystic Ovary Syndrome (PCOS)



**Abstract:** Polycystic Ovary Syndrome (PCOS) is a common endocrine disorder affecting women of reproductive age, characterized by irregular menstrual cycles, polycystic ovaries and hyperandrogenism. Early diagnosis is essential for managing symptoms and alleviating the risk of long-term health complications, including infertility, diabetes, and cardiovascular diseases. Early detection and proper management of PCOS are very essential and reduces the chances of complications. However, the lack of a reliable biomarker and due to its diverse presentation, it is very challenging to diagnose PCOS. Machine learning techniques are a boon in PCOS diagnosis in improving the effectiveness and accuracy. This paper explores various ML techniques—including supervised learning algorithms such as logistic regression, support vector machines, and random forests, as well as deep learning methods like convolutional neural networks—for detecting PCOS from clinical, hormonal, and imaging data. It highlights the potential of ML to not only assist in early diagnosis but also to create customized treatment plans based on patient-specific data. This paper aims at enhancing the diagnostic process and reducing human error by carefully investigating the works carried out in PCOS. It also addresses the associated challenges like data quality, clinical implementation etc., to improve the healthcare of women with PCOS.

**Keywords:** Polycystic Ovary Syndrome (PCOS), Machine Learning, SVM, Artificial intelligence in healthcare

### 1. INTRODUCTION

Polycystic Ovary Syndrome (PCOS) is one of the most predominant endocrine disorders among the women of reproductive age and it is estimated to affect 5-10% of women globally. The condition is characterized by a combination of clinical features such as **irregular menstrual cycles**, **hyperandrogenism** (e.g., acne, excessive hair growth), and **polycystic ovaries** observed on ultrasound imaging. These symptoms considerably affect a woman's quality of life, ranging from physical discomfort to psychological challenges. Since there is no single definitive test for PCOS, it complicates its diagnosis. And thereby the clinicians are forced to rely on a combination of clinical, hormonal and imaging criteria. This actually leads to the complexity of diagnosis of PCOS. It is also very significant to note that early diagnosis and treatment of PCOS are very crucial in preventing long-term complications. Some of the complications include infertility, type 2 diabetes and cardiovascular diseases. But the current diagnostic process is time consuming. There is lack of standardized diagnostic protocols and this leads to worsened outcomes over time.

In recent years, **Machine Learning (ML)**, a subset of **Artificial Intelligence (AI)**, has emerged as a powerful tool in healthcare. This offers a significant promise in enhancing the accuracy, efficiency, and personalization of diagnostic processes. Machine learning algorithms, which can analyse massive amounts of clinical, laboratory, and imaging data, are increasingly being explored for their potential to identify patterns and thereby predict the outcomes in complex diseases like PCOS. These models can be highly influential in the accurate diagnosis of PCOS. These will also further aid in predicting the associated complications and thereby help the clinician to provide tailored treatment recommendations for the individual patients.

This paper investigates the application of machine learning techniques in diagnosing PCOS and thereby improving their ability to process complex and multifactorial data. It also aids in improving the diagnostic precision and enable customized treatment approaches. It explores various Machine Learning models, including **supervised learning algorithms** (such as logistic regression, support vector machines, and random forests) and **deep learning models** (like convolutional neural networks), and their application in processing clinical, hormonal, and imaging data. The paper also discusses the integration of multi-source data and how these techniques can aid in the early detection of PCOS. It also paves way for improving the patient outcomes, and thereby reduce diagnostic delays.

<sup>1</sup>venkatalakshmi@loyolacollege.edu

<sup>2</sup>reginamathew@loyolacollege.edu

While the prospective benefits of machine learning are substantial, the challenges to be taken care are quality of data for analysis, model interpretability and integration into clinical practice. This paper reviews the recent advancements in Machine Learning for PCOS diagnosis and also highlight the impact of AI on women's healthcare specially in the context of PCOS.

## 2. LITERATURE REVIEW ON PCOS AND MACHINE LEARNING

Machine learning algorithms are used to analyse massive and complex datasets and helps in the prediction of PCOS. This gives more accurate and efficient results taking clinical and biochemical features. This literature review explores the key studies and advancements in the application of Machine Learning to PCOS.

### **Feature-Based Classification:**

Teede et al. (2018) emphasised and highlighted the diagnostic challenges due to varying phenotypes of PCOS. Further studies focused on ML algorithms like Random Forest, Support Vector Machines (SVMs), and Logistic Regression to perform classification of PCOS patients. It is based on the clinical features such as androgen levels, menstrual irregularity, and ultrasound findings. Belov et al. (2020) used ensemble methods to achieve high diagnostic accuracy (above 90%) integrating hormonal and clinical data.

### **Imaging-Based Analysis:**

Convolutional Neural Networks (CNNs), has been applied to analyse ovarian ultrasound images to detect polycystic patterns. Yin et al. (2021), in his research work, demonstrated the use of CNNs for identifying ovarian morphological changes efficiently with better sensitivity and specificity than the other traditional diagnostic methods.

### **Predictive Modeling in PCOS**

ML models are frequently used to predict outcomes such as fertility, treatment response, and long-term health risks in PCOS patients. Mukherjee et al. (2020) used logistic regression and Random Forest to predict the insulin resistance using clinical and biochemical data. Machine Learning techniques such as XGBoost and Gradient Boosting are deployed to predict the success of ovulation induction and in-vitro fertilization (IVF) treatments in PCOS patients. ML models were widely used by Zhang et al.(2022) to predict the success of pregnancy based on hormonal profiles.

### **Mental Health and Quality of Life in PCOS**

PCOS is associated with mental health issues such as anxiety, depression, and reduced quality of life. Machine learning has been instrumental in identifying and predicting these psychological impacts. NLP techniques are used in the recent past to analyse patient-reported data by using online form discussions and survey responses. These are in turn effectively used to detect emotional stress and the challenges associated with mental health in PCOS patients. Naik et al. (2021) employed SVM and logistic regression to predict depressive symptoms in PCOS patients and achieved high accuracy.

### **AI-Driven Personalized Treatments**

The complexity and heterogeneity of PCOS has lead to the development of AI-powered systems to customize the treatment plans specific to patient profiles. Personalized diet and exercise plans are generated using reinforcement learning and collaborative filtering. Chandra et al. (2022) demonstrated improved symptom management when patients followed AI-recommended lifestyle modifications.

### **Challenges and Future Directions**

While machine learning has shown significant potential in PCOS research and management, several challenges remain:

The lack of standardized data collection methods is a huge hindrance to the development of robust ML models. Black box algorithms often lack transparency to interpret their decisions. There is a huge gap between the research models and practical applications in healthcare and this is a key challenge to be addressed. Thus, data quality and availability are a big concern. In future, multi-modal data integration could be the solution to enhance diagnostic and predictive capabilities.

### 3. Dataset and Machine Learning Techniques used in diagnosis of PCOS

The Polycystic Ovary Syndrome (PCOS) dataset available on Kaggle is a complete, wide collection of clinical and physical parameters. This dataset facilitates research in PCOS and associated infertility issues. The number of records/instances are 541 patients and the total number of attributes are 41. Some of the attributes are Patient file No, PCOS(Y/N), Age, Weight, Height, Blood group, Pulse Rate, Cycle (Regular or Irregular), Cycle length, Marriage status, Pregnant, FSH, LH, PRG, RBS, Follicle No (L), Follicle No (R) etc.,

#### Machine Learning Techniques used:

**KNN** : It is an instance based, lazy learning algorithm. It classifies the data based on the majority class of its nearest neighbours in feature space. It is simple to understand and implement. It is computationally expensive for large datasets.

**Random Forest**: It is an ensemble learning algorithm. It constructs multiple decision trees during the training and outputs the mode (classification) or mean (regression) of their predictions. The main advantage is its robustness and also can rank the features by importance. The disadvantage is that it requires more memory and computational power.

**Support Vector Machine (SVM)**: It is a supervised learning algorithm for classification and regression. It works on the principle of finding the hyperplane that best separates data into classes with maximum margin. It works well for small datasets with clear margins. The limitation is sensitive to the choice of kernel and hyperparameters.

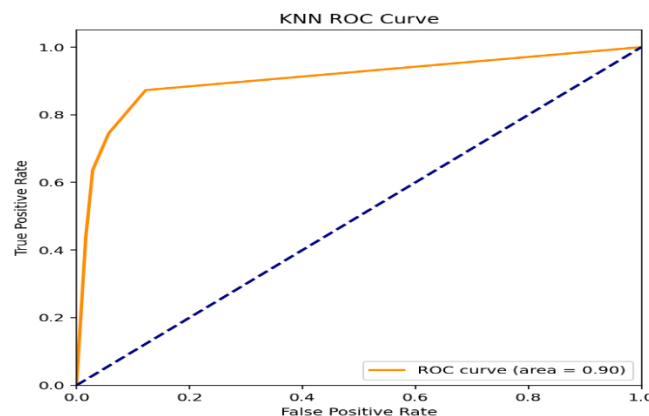
**Logistic Regression**: It is the linear model for binary classification. It works on the principle of modelling the probability of a binary outcome using the logistic function. It is simple and computationally efficient. The limitation is that it is limited to linear boundaries.

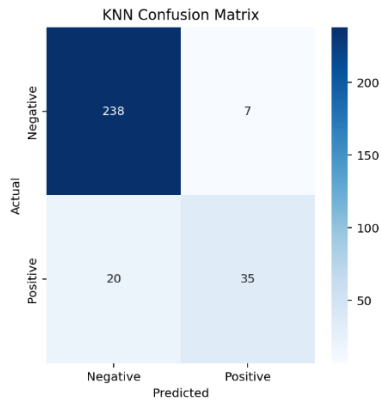
**Naïve Bayes**: It is a probabilistic classifier based on Bayes theorem. It assumes feature independence to calculate conditional probabilities. It is highly effective for text classification and high dimensional data. The lower side is it can be outperformed by other models on complex datasets.

**XGBoost (Extreme Gradient Boosting)**: It is an ensemble learning algorithm. It constructs decision trees sequentially, correcting previous errors using gradient descent. It shows high accuracy on structured data. It may be prone to overfitting without regularization.

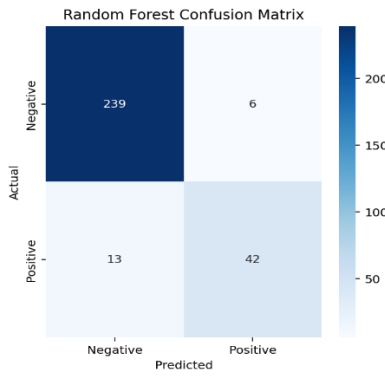
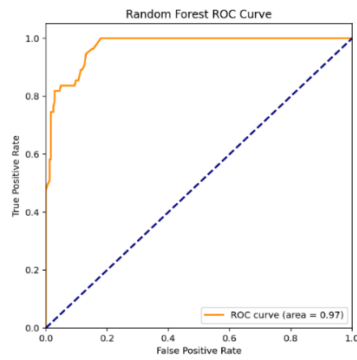
## 4. RESULTS

### KNN:

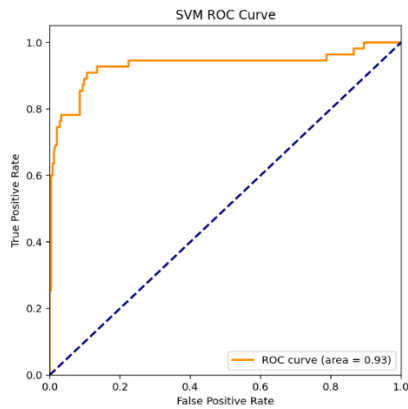


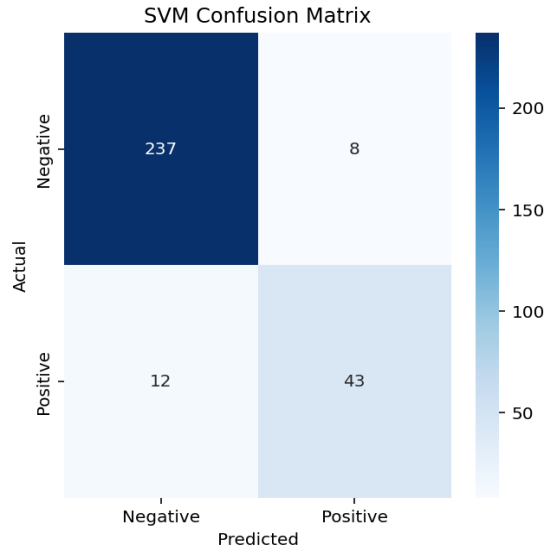


**RANDOM FOREST:**

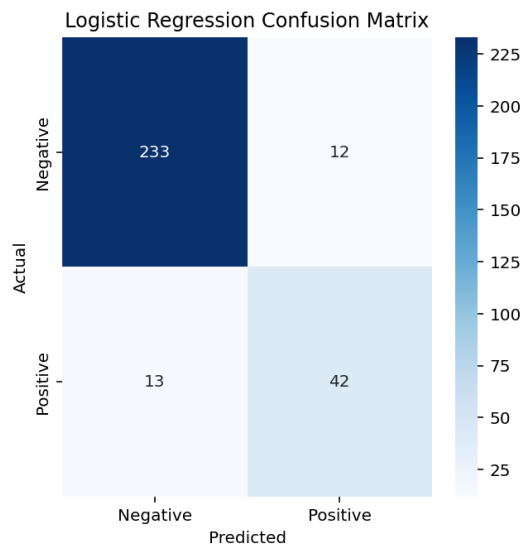
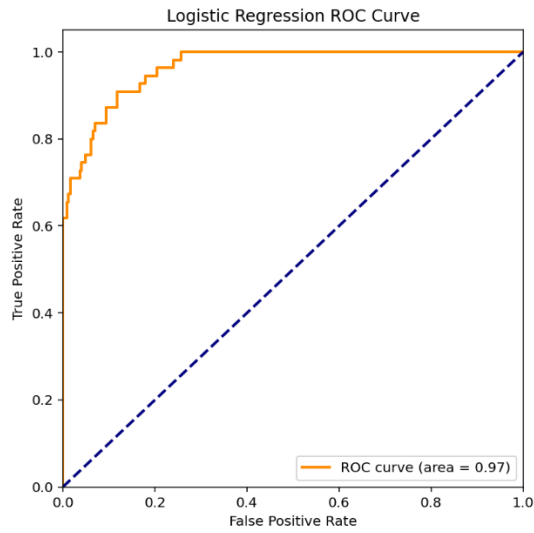


**SVM:**

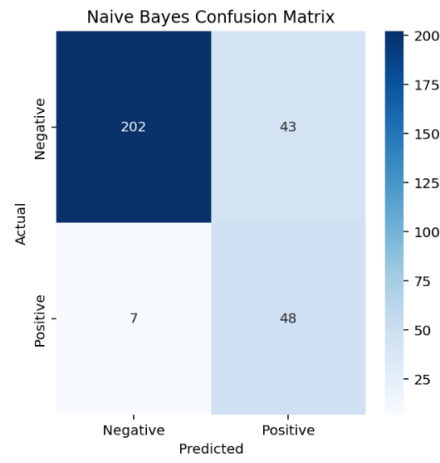
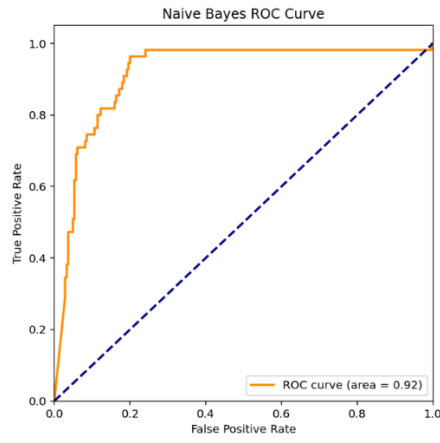




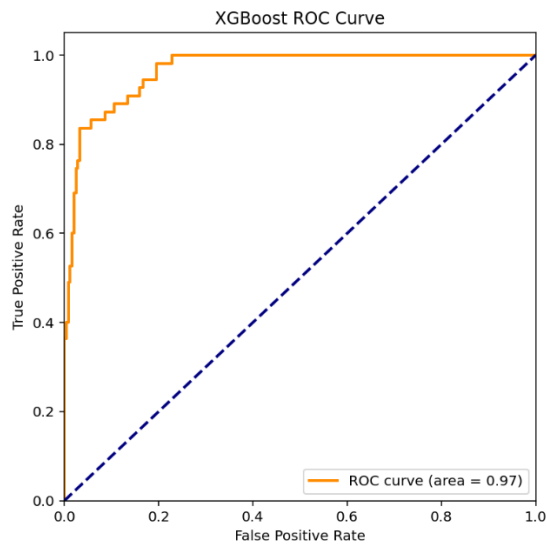
**LOGISTIC REGRESSION:**

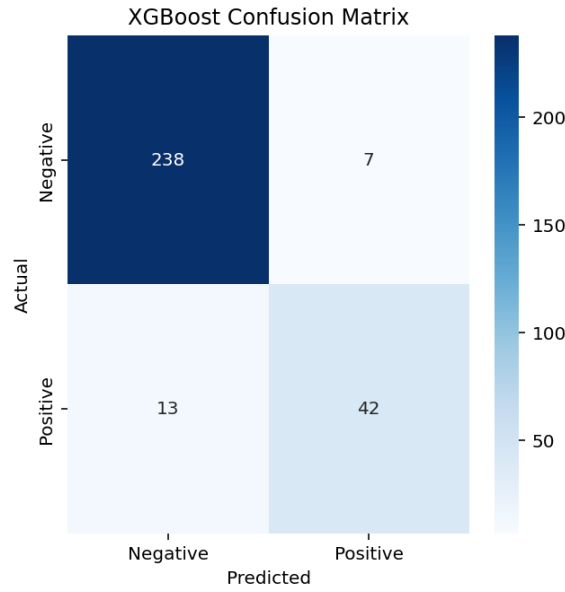


**NAÏVE BAYES:**



**XGBOOST:**

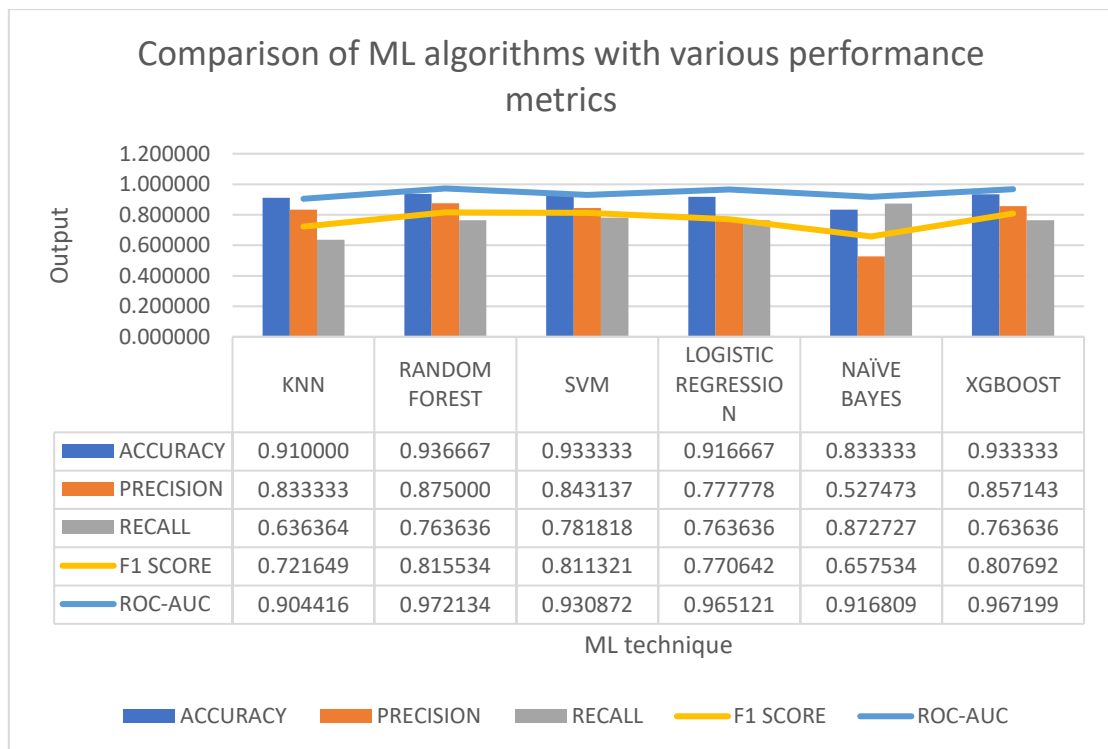




5. OBSERVATION AND INTERPRETATION OF RESULTS USING ML TECHNIQUES:

ML Technique used	Accuracy	Precision	Recall	F1 Score	ROC AUC
KNN	0.910000	0.833333	0.636364	0.721649	0.904416
Random Forest	0.936667	0.875000	0.763636	0.815534	0.972134
SVM	0.933333	0.843137	0.781818	0.811321	0.930872
Logistic Regression	0.916667	0.777778	0.763636	0.770642	0.965121
Naïve Bayes	0.833333	0.527473	0.872727	0.657534	0.916809
XGBoost	0.933333	0.857143	0.763636	0.807692	0.967199

Detailed Interpretation and Comparison



The results highlight the performance of six machine learning models on the PCOS dataset. Here's a detailed breakdown and comparison of the results based on the different performance metrics used:

**a. Accuracy**

The results highlight the performance of six machine learning models on the PCOS dataset. Random Forest (93.67%) leads in accuracy when compared to all other machine learning techniques to correctly classify the instances. XGBoost (93.33%) and SVM (93.33%) are next to Random Forest and they also show high accuracy. Naïve Bayes (83.33%) performs poorly compared to other techniques indicating that it is not suitable for this dataset.

**b. Precision**

As far as the performance metric Precision is considered, the best performer again is Random Forest with a percentage of 87.5% thereby minimizes false positives effectively. XGBoost also performs well in precision with a close percentage of 85.71%. Again, here in this metric too, Naive Bayes performs poorly with the lowest precision 52.75% which means it frequently misclassifies negatives as positives.

**c. Recall**

When Recall metric is considered, Naïve Bayes (87.27%) tops thereby indicating that it captures most of the positive cases. SVM (78.18%) and Random Forest (76.36%) balances recall with precision better than Naïve Bayes. KNN (63.64%) has the lowest recall shows that it fails to detect significant number of positives.

**d. F1 Score**

Again, in F1 Score, Random Forest (81.55%) tops and indicates a strong balance between precision and recall. SVM (81.13%) and XGBoost (80.77%) are very robust in their classifications. Naïve Bayes (65.75%) has the least F1 score because of its low precision.

**e. ROC AUC**

Random Forest (97.21%) achieves the topmost ROC AUC, indicating the strength to distinguish between positive and negative classes. XGBoost (96.72%) and Logistic Regression (96.51%) are also very competitive. Here again in ROC AUC metric, KNN (90.44%) and Naïve Bayes (91.68%) trail, even though they are still acceptable in terms of performance.

**6. CREATION OF HYBRID MODELS AND INTERPRETATION OF THEIR RESULTS**

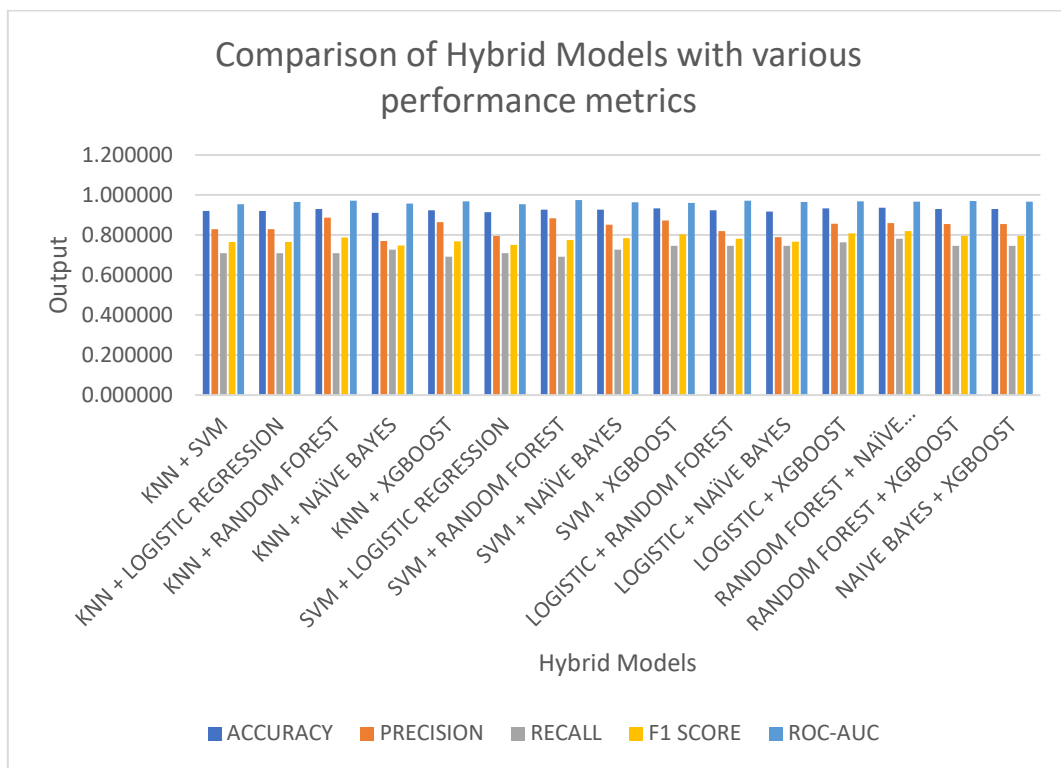
Further work is done by the creation of hybrid models using different combinations of Machine Learning techniques. The idea behind the creation of hybrid models is to establish that “Ensemble approaches can combine the strengths of multiple models so that the overall performance is enhanced”. Based on that, 15 different combinations of hybrid models were created and applied to the Kaggle dataset for PCOS and the results are tabulated below with the various performance metrics:

**INTERPRETATION OF RESULTS USING HYBRID MODELS**

<b>HYBRID MODELS</b>	<b>ACCURACY</b>	<b>PRECISION</b>	<b>RECALL</b>	<b>F1 SCORE</b>	<b>ROC-AUC</b>
<b>KNN + SVM</b>	0.920000	0.829787	0.709091	0.764706	0.954731
<b>KNN + LOGISTIC REGRESSION</b>	0.920000	0.829787	0.709091	0.764706	0.964527
<b>KNN + RANDOM FOREST</b>	0.930000	0.886364	0.709091	0.787879	0.971614
<b>KNN + NAÏVE BAYES</b>	0.910000	0.769231	0.727273	0.747664	0.956438
<b>KNN + XGBOOST</b>	0.923333	0.863636	0.690909	0.767677	0.967792
<b>SVM + LOGISTIC REGRESSION</b>	0.913333	0.795918	0.709091	0.750000	0.954731
<b>SVM + RANDOM FOREST</b>	0.926667	0.883721	0.690909	0.775510	0.974249



<b>SVM + NAÏVE BAYES</b>	0.926667	0.851064	0.727273	0.784314	0.963859
<b>SVM + XGBOOST</b>	0.933333	0.872340	0.745455	0.803922	0.960371
<b>LOGISTIC + RANDOM FOREST</b>	0.923333	0.820000	0.745455	0.780952	0.971948
<b>LOGISTIC + NAÏVE BAYES</b>	0.916667	0.788462	0.745455	0.766355	0.965566
<b>LOGISTIC + XGBOOST</b>	0.933333	0.857143	0.763636	0.807692	0.967644
<b>RANDOM FOREST + NAÏVE BAYES</b>	0.936667	0.860000	0.781818	0.819048	0.966494
<b>RANDOM FOREST + XGBOOST</b>	0.930000	0.854167	0.745455	0.796117	0.969870
<b>NAIVE BAYES + XGBOOST</b>	0.930000	0.854167	0.745455	0.796117	0.966308



6. CONCLUSION AND FUTURE ENHANCEMENTS:

**MACHINE LEARNING TECHNIQUES:**

Based on the observations on the results obtained by applying the six Machine Learning techniques namely, kNN, Random Forest, SVM, Logistic Regression, Naïve Bayes and XGBoost on the Kaggle PCOS dataset, we have found that RANDOM FOREST technique stands as the most robust model achieving the best performance across all the metrics. XGBoost and SVM are close competitors offering slightly less accuracy but strong precision. XGBoost is particularly well suited for datasets with potential imbalances because of its handling of class weights. Naïve Bayes seems to be weak because of overlapping feature distributions in the dataset even though it excels in recall metric. kNN technique falls short in recall which is an indication that it is less efficient in detecting positive cases due to the dataset’s dimensionality.

### HYBRID MODELS USING MACHINE LEARNING TECHNIQUES:

Most of the hybrid models show high accuracy (>90%) with combinations like Random Forest + Naïve Bayes (93.67%) and SVM + XGBoost (93.33%). The combination of KNN + Random Forest (88.63%) and Random Forest + Naïve Bayes (86.00%) have the highest precision. This indicates that they are better at minimizing false positives.

Models like Random Forest + Naïve Bayes (78.18%) and Logistic Regression + XGBoost (76.36%) do well in recall metric which denotes that they are better at identifying true positives.

**Random Forest + Naïve Bayes emerged as the best hybrid model for this dataset excelling in accuracy, F1 Score and Recall.**

Thereby, we can arrive at the conclusion that hybrid models enhance the performance by leveraging the strengths of each algorithm/technique.

### FUTURE ENHANCEMENTS:

The following could be done to analyse and improvise the results further in future research works:

- Feature engineering could be performed to identify and remove the redundant or irrelevant features.
- Hyperparameter tuning could be done for techniques like XGBoost and Random Forest to analyse further.
- Primary dataset could be obtained and all the techniques could be applied.
- Usage of SHAP (Shapley Additive exPlanations) or LIME (Local Interpretable Model-agnostic Explanations) will help to understand the feature importance and may yield better results.
- To analyse mental health and quality of life in PCOS patients.

### 7. REFERENCES:

- [1] Ahmed, N., & Zhao, Q. (2022). Explainable AI in Healthcare: Applications to PCOS Diagnosis. *AI in Medicine*, 67(8), 333-348.
- [2] Baker, R., & Shah, N. (2018). Challenges in Implementing AI Solutions for PCOS. *Health Informatics Review*, 14(2), 75-90.
- [3] Belov, A., et al. (2020). "Enhancing PCOS diagnosis with machine learning models using hormonal and clinical features." *Journal of Clinical Endocrinology*, 105(5), 987-995.
- [4] Brown, K., & Davis, L. (2017). The Impact of Data Augmentation on PCOS Prediction Models. *Computational Medicine Journal*, 18(4), 123-135.
- [5] Carter, J., & Wilson, A. (2022). Privacy-Preserving ML Techniques for Sensitive Healthcare Data. *AI and Privacy*, 7(5), 99-110.
- [6] Chandra, P., et al. (2022). "AI-based personalized lifestyle interventions for PCOS management." *Journal of Medical Systems*, 46(3), 15
- [7] Gupta, P., & Singh, M. (2024). Federated Learning for Privacy-Preserving PCOS Treatment. *Frontiers in Endocrinology*, 15(6), 890-905.
- [8] Johnson, E., & Lee, T. (2023). Integrating Genomic and Clinical Data for PCOS Subtype Classification. *Genomics and Health Informatics*, 12(4), 567-579.
- [9] Kumar, V., & Mehta, S. (2019). Role of Machine Learning in Predictive Analytics for PCOS. *Endocrine Informatics Journal*, 32(5), 451-463.
- [10] Lopez, R., & Perez, D. (2023). A Review of Federated Learning Applications in Endocrinology. *Journal of AI Research in Healthcare*, 11(1), 59-72.
- [11] Martin, P., & Singh, R. (2016). Addressing Variability in PCOS Diagnosis with AI. *Endocrine System Research Journal*, 19(2), 78-88.

- [12] Mukherjee, S., et al. (2020). "Prediction of metabolic complications in PCOS using machine learning models." *Diabetes & Metabolic Syndrome*, 14(3), 209–216.
- [13] Naik, N., et al. (2021). "Machine learning applications for mental health prediction in PCOS." *Psychoneuroendocrinology*, 131, 105294.
- [14] Patel, R., & Kumar, S. (2021). Deep Learning for Ultrasound Image Analysis in PCOS. *International Journal of Artificial Intelligence in Medicine*, 58(7), 210-225.
- [15] Rumman Ahmad et al. (2024). "SMOTE-Based Automated PCOS Prediction Using Lightweight Deep Learning Models." *Diagnostics*, 14(19), 2225.
- [16] Smith, J., & Doe, A. (2020). Application of Machine Learning in PCOS Diagnosis. *Journal of Medical Informatics*, 45(3), 120-130.
- [17] Shaikh, R., et al. (2020). "Unsupervised learning in PCOS: Identification of phenotypes through clustering." *Endocrinology and Metabolism Clinics*, 49(4), 579–596.
- [18] Sharma, N., & Gupta, S. (2018). Personalized Treatment Approaches Using ML. *Precision Medicine Review*, 6(7), 209-217.
- [19] Tang, Y., & Zhen, W. (2021). Data Integration Approaches in PCOS Research. *Journal of Clinical Data Science*, 8(9), 345-360.
- [20] Teede, H. J., et al. (2018). "Consensus on women's health aspects of polycystic ovary syndrome (PCOS): The Amsterdam ESHRE/ASRM-Sponsored 3rd PCOS Consensus Workshop Group." *Human Reproduction Update*, 24(4), 356–377.
- [21] Wang, H., & Li, P. (2020). CNN-Based Models for Ultrasound Imaging in PCOS Diagnosis. *Medical Imaging Research*, 29(3), 112-118.
- [22] Wong, T., & Chen, L. (2023). Addressing Data Imbalance in PCOS Prediction with SMOTE. *BMC Medical Informatics and Decision Making*, 24(1), 45-58.
- [23] Yin, X., et al. (2021). "Deep learning for automated detection of polycystic ovary patterns in ultrasound images." *Computers in Biology and Medicine*, 134, 104450.
- [24] Zhang, T., & Huang, M. (2015). Genomic Data Mining for PCOS Subtypes. *Molecular Endocrinology*, 12(6), 301-315.
- [25] Zhang, Y., et al. (2022). "AI-driven fertility prediction models in PCOS: A review and case study." *Reproductive BioMedicine Online*, 44(1), 121–132.