



e-ISSN: 2319-8753 | p-ISSN: 2347-6710

# IJRSET

International Journal of Innovative Research in  
**SCIENCE | ENGINEERING | TECHNOLOGY**



# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

Volume 13, Issue 8, August 2024

**ISSN** INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA

Impact Factor: 8.524

9940 572 462

6381 907 438

[ijirset@gmail.com](mailto:ijirset@gmail.com)

[www.ijirset.com](http://www.ijirset.com)

# Predicting PCOS : A Machine Learning Approach

Dr. D Kirubha, Bhoomika H K, Binetta P Biju, Divya J

Associate Professor, Department of Computer Engineering, Rajarajeswari College of Engineering, Bengaluru,  
Karnataka, India

U.G. Student, Department of Computer Engineering, Rajarajeswari College of Engineering, Bengaluru,  
Karnataka, India

**ABSTRACT:** Polycystic Ovary Syndrome (PCOS) is a widespread endocrine condition that affects women of reproductive age, characterized by hormonal imbalance and ovarian dysfunction. Early detection and accurate identification of PCOS are paramount for effective management and prevention of complications. In this proposal, we put forward a machine learning-based approach utilizing relevant clinical and demographic data to predict PCOS. The dataset includes features such as age, body mass index (BMI), hormonal levels (e.g., testosterone, luteinizing hormone), menstrual irregularities, and other parameter that are linked with PCOS. Gaussian Naive Bayes (GNB) algorithm is employed as the main algorithm, with feature selection performed using Chi-Square. The data is partitioned, and data visualizations including pie charts, count plots, scatter plots, lm plots, and histograms are utilized to explore correlations. Metrics for performance such as accuracy score, precision score, and recall score are assessed. Additionally, confusion matrices are generated to evaluate model performance. Our findings demonstrate the potential of the GNB algorithm in predicting PCOS with high accuracy and reliability. The developed models can aid healthcare practitioners in early identification of individuals at risk of PCOS, enabling timely interventions and personalized treatment strategies. Furthermore, findings obtained from this study contribute to a better understanding of the complex etiology of PCOS and guide future research in this domain.

**KEYWORDS:** Polycystic Ovary Syndrome (PCOS), Machine Learning, Predictive Modeling, Clinical Data Analysis, Feature Selection, Hormonal Imbalance, Feature Importance.

## I. INTRODUCTION

Polycystic Ovary Syndrome (PCOS) persists as a widespread and intricate medical condition, exerting significant impacts on the health and well-being of women during their reproductive years. Characterized by hormonal imbalances and disruptions in ovarian function, PCOS poses various challenges, from fertility issues to metabolic disturbances. Timely detection and accurate diagnosis are crucial in effectively managing PCOS and mitigating potential complications, highlighting the urgent need for innovative predictive approaches..

In response to this imperative, our study endeavors to harness the transformative potential of ML in predicting PCOS using a diverse range of health-related and demographic data. Recognizing that traditional diagnostic methods can be time-consuming and costly for both patients and healthcare providers, we aim to develop more efficient and effective predictive tools.

Central to our approach is the utilization of advanced ML algorithm, with GNB emerging as the mainstay of our predictive model. We employ a novel feature selection technique, dubbed CS-PCOS, which leverages the optimized chi-squared mechanism to identify the most salient predictors of PCOS from the dataset. By focusing on key features such as age, BMI (weight), hormonal levels, and menstrual irregularities, our model aims to distill complex clinical detail into awareness for early PCOS detection.

Furthermore, our study delves into the exploration of underlying data patterns through a comprehensive PCOS exploratory data analysis (PEDA). Through visualizations such as pie charts, count plots, scatter plots, and histograms, we seek to elucidate the intricate interplay of variables contributing to the manifestation of PCOS, thereby enhancing our awareness of the disease process.

Additionally, model development and data exploration, our research extends to the evaluation and validation of predictive performance. We apply ten advanced ML models, including stochastic gradient descent, random forest, support vector machine, and logistic regression, to compare and assess the efficacy of our proposed GNB model. Through k-fold cross-validation, we ensure the robustness and generalizability of our findings, providing confidence in the accuracy and reliability of our predictive framework.

By synthesizing insights from literature analysis, methodological approaches, and observed findings, our study intends to make meaningful contributions in the area of PCOS prediction. Through the integration of cutting-edge ML techniques and comprehensive data analysis, we seek to empower healthcare practitioners with the tools and knowledge necessary for early detection and personalized management of PCOS. Ultimately, our research underscores the transformative potential of machine learning in advancing healthcare and underscores the prominence of proactive measures in addressing the complexities of PCOS.

## II. RELATED WORK

In this section, we embark on an in-depth exploration of the literature related to PCOS prediction, elucidating past studies and pioneering techniques utilized in this realm. These studies illuminate the pressing nature and intricate nuances of PCOS diagnosis, providing invaluable insights into the evolving landscape of predictive modelling for this multifaceted condition.

One noteworthy study from the past employed a spectrum of ML methodologies, including Support Vector Machines (SVM), Random Forest (RF), CART, Logistic Regression, and Naive Bayes, to discern PCOS patients based on diverse medical indicators and signs. Notably, the RF algorithm showcased exceptional prowess, achieving an impressive 96% accuracy rate in diagnosing PCOS on a given dataset. This underscores the transformative potential of ML in enhancing the precision and efficacy of PCOS prediction. Similarly, another investigation harnessed a plethora of ML models on a dataset comprising 541 patients, employing sophisticated feature selection approaches to determine pivotal predictors of PCOS. The resultant RFLR algorithm, a fusion of Random Forest and Logistic Regression, yielded promising outcomes with a notable 90.01% accuracy score, showcasing the efficacy of ensemble learning paradigms in PCOS classification.

The emergence of sophisticated predictive modelling methods has sparked the development of advanced methods for early identification of PCOS. For example, a recent investigation introduced a model based on XGBRF and Cat Boost algorithms, supported by feature selection techniques, yielding impressive accuracy rates of 89% and 95%, correspondingly. These results highlight the value of utilizing state-of-the-art machine learning methods and feature selection approaches for accurate and timely prediction of PCOS.

Beyond conventional machine learning paradigms, researchers have ventured into innovative frameworks integrating genetic analysis and bioinformatics techniques for PCOS prediction. These inspection have delivered promising results, offering insights into the genetic underpinnings of PCOS and presenting novel avenues for predictive modelling in this domain. Furthermore, endeavors were made to develop automated systems for PCOS detection, leveraging a repertoire of ML models such as Gaussian Naive Bayes, SVM, k-nearest neighbours, Random Forest, and Logistic Regression. These studies underscore the potential of automated approaches in streamlining PCOS diagnosis and enhancing patient outcomes through early intervention. The corpus of literature surrounding PCOS prediction reflects the ongoing evolution in machine learning, genetic analysis, and bioinformatics, furnishing a diverse array of methodologies and techniques for accurate and efficient PCOS detection. By synthesizing insights from these studies, our research endeavors to contribute to this burgeoning field, offering innovative methods for early prediction and management of PCOS.

## III. METHODOLOGY

The workflow of our research methodology is meticulously examined to ensure coherence and effectiveness in each stage, from data preprocessing to model deployment.

This systematic approach underpins the reliability and validity of our research findings, facilitating meaningful contributions to the field of PCOS prediction and management, as in



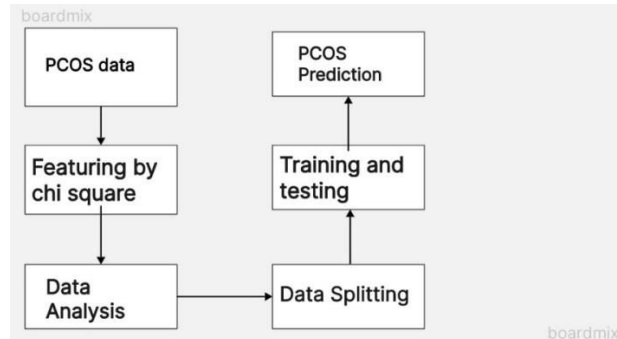


Figure 1: PCOS Prediction

A. Dataset Description and Preprocessing:

Our investigation employs a dataset containing pertinent clinical and physical attributes related to this condition to construct ML models. The dataset undergoes thorough feature engineering utilizing the innovative CS-PCOS method, with the aim of refining feature selection to enhance prediction accuracy. To unravel underlying data patterns contributing to PCOS, we conduct PCOS exploratory data analysis (PEDA), facilitating a deeper comprehension of the disease's etiology. The dataset is thoroughly pre-processed, including handling null values and dropping irrelevant columns, to ensure its suitability for model training and evaluation. Not use abbreviations in the title or heads unless they are unavoidable, as in figure 2

Sl. No	Patient File No.	PCOS (Y/N)	Age (yrs)	Weight (Kg)	Height(Cm)	BMI	Blood Group	Pulse rate(bpm)	RR (breaths/min)	Fast food (Y/N)	Reg.Exercise(Y/N)	BP_Systolic (mmHg)	BP_Diastolic (mmHg)	Follicle No. (L)	Follicle No. (R)	Avg. F size (mm)		
0	1	1	0	28	44.6	152.0	19.3	15	78	22	..	1.0	0	110	80	3	3	18.0
1	2	2	0	36	65.0	161.5	24.9	15	74	20	..	0.0	0	120	70	3	5	15.0
2	3	3	1	33	68.8	165.0	25.3	11	72	18	..	1.0	0	120	80	13	15	18.0
3	4	4	0	37	65.0	148.0	29.7	13	72	20	..	0.0	0	120	70	2	2	15.0
4	5	5	0	25	52.0	161.0	20.1	11	72	18	..	0.0	0	120	80	3	4	16.0

5 rows x 45 columns

Figure 2: Dataset

B. Novel CS-PCOS Feature Engineering Technique

In our methodology, we introduce a pioneering feature engineering technique, the CS-PCOS approach, tailored to transform dataset features into optimal representations for predictive modelling. This novel approach leverages the optimized chi-squared mechanism to assess feature importance and select the most salient predictors of PCOS. Through a systematic process, the CS-PCOS technique identifies vital feature statistics based on goodness of fit, thereby refining the dataset for predictive modelling as in figure 3

Feature	Chi-Square Score	P-value
23	AMH(ng/ml)	2.393024e-179
28	Follicle No. (f)	1.777709e-120
37	Follicle No. (l)	1.299632e-100
16	FSH(IU/ml)	1.850831e-96
2	Weight (kg)	2.820457e-72
1	Patient File No.	4.797859e-60
6	SI_No	4.797859e-60
18	FSH(LU)	7.186243e-55
14	I_beta-HCG(IU/ml)	2.713000e-31
20	hair_growth(Y/N)	1.322012e-15
30	Skin darkening (Y/N)	2.109234e-15
25	Wt_O3 (ng/ml)	1.273384e-14
28	Weight gain(Y/N)	7.459214e-13
9	Cycle(L/D)	9.728484e-13
2	Age (yrs)	9.702724e-09
10	Hip(Inch)	6.617673e-08
11	Hemoglobin (gms)	6.129372e-07
26	PRC(ng/ml)	8.492918e-07
33	Fast food (Y/N)	1.294884e-06
30	Avg. F size (L) (cm)	4.408648e-06
17	IM(In/ml)	4.868488e-05
20	Wt(Inch)	6.960125e-05
8	HB(g/dl)	6.947194e-05
40	Avg. F size (R) (cm)	1.837371e-04
32	Pimplast(Y/N)	2.527894e-04
24	PHI(ng/ml)	1.880389e-03
4	Height(Cm)	2.375544e-03
6	Pulse rate(bpm)	3.205859e-03
19	Cycle length(days)	1.658914e-01
22	THI (Cm/L)	1.525709e-02
31	Hair loss(Y/N)	4.199754e-02
41	Endometriosis (cm)	4.875094e-02
7	RR (breaths/min)	8.856861e-02
15	No. of abortions	9.559397e-02
21	Wt:Hip Ratio	1.924262e-01
27	HBS(ng/dl)	2.012348e-01
15	I_beta-HCG(IU/ml)	5.227191e-01
5	Blood Group	5.879928e-01
34	Reg.Exercise(Y/N)	7.250339e-01
12	Pregnant(Y/N)	7.281748e-01
36	BP_Diastolic (mmHg)	8.802300e-01
42	BP_Systolic (mmHg)	9.545395e-01
35	BP_Systolic (mmHg)	9.522061e-01

Figure 3: Feature selection

C. Visualization and Analysis of Data Patterns:

To delve into the dataset further, we utilize diverse visualization methods. Count plots are employed to illustrate the distribution of instances across PCOS categories, offering a glimpse into class imbalances. A pie chart provides a visual depiction of class proportions within the dataset, revealing the prevalence of PCOS. Furthermore, Bar plots and 3D scatter plots facilitate the visualization and examination of essential feature interactions, unveiling potential relationships influencing the occurrence of PCOS, as in figure 4

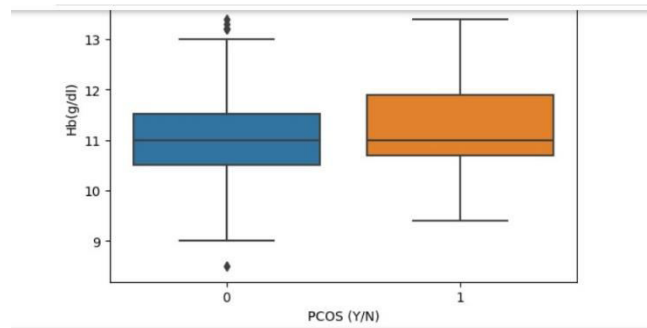


Figure 4: Bar plot

D. Regression Analysis and Visualization:

We employ lm plot, a versatile tool combining regression plots and Facet Grid, to analyse the relationship between high-value dataset features and PCOS. Through lm plot visualization, we uncover patterns and trends indicative of PCOS regression, offering valuable insights into the predictive power of individual features. Regression plots provide an in-depth comprehension of feature interactions and their impact on PCOS likelihood, as in figure 5

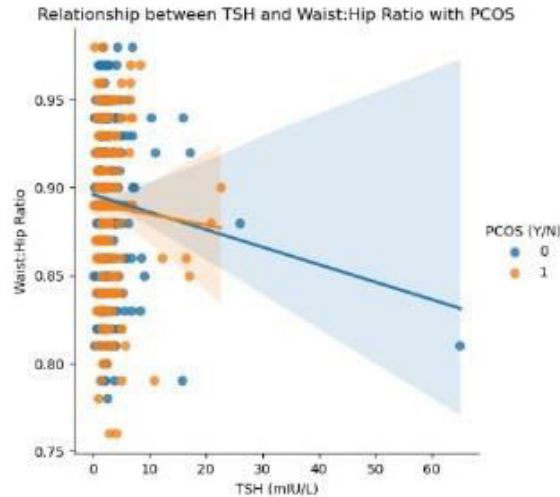


Figure 5: Im plot

E. Frequency Distribution Analysis:

Histograms are employed to analyze the frequency distribution of PCOS classes across imported features, facilitating a nuanced understanding of feature significance. By examining the frequency distribution of PCOS classes across key features such as Hip size, Haemoglobin levels, and Blood Pressure readings, we gain insights into the prevalence and distribution of PCOS within the dataset, as in figure 6

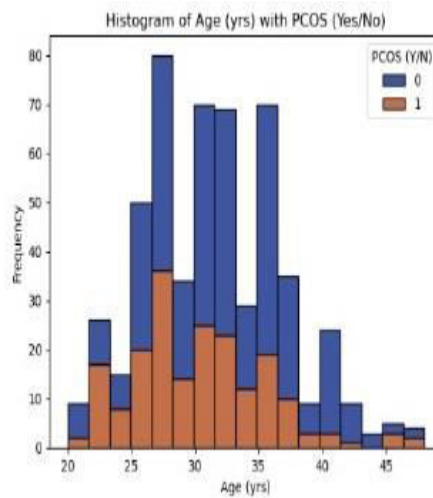


Fig 6: Histogram for Age with PCOS

F. Dataset Splitting for Model Evaluation:

For rigorous model evaluation and to mitigate overfitting, we splitting the data into training and testing sets, maintaining an 80:20 ratio. The training portion, comprising 80% of the dataset, is utilized for training the model, while the remaining 20% is reserved for evaluating model performance on unseen data. This rigorous approach ensures the reliability and generalizability of our machine learning models, yielding high accuracy results for PCOS prediction.

#### IV. INVESTIGATING MACHINE LEARNING METHODS FOR PREDICTING PCOS

This section delves into the utilization of diverse ML methodologies for the prediction of PCOS. We elucidate the operational mechanisms and mathematical formulations underlying these techniques, aiming to shed light on their efficacy in discerning PCOS occurrence.

Stochastic Gradient Descent (SGD) Classifier:

The SGD classifier is harnessed for its prowess in large-scale learning tasks, employing loss functions derived from the stochastic gradient descent routine. Leveraging the SGD's efficiency and ease of implementation, we optimize parameters to minimize the logistic cost function, as delineated in Equation 1.

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases} \quad \dots \quad (1)$$

Linear Regression (LIR)

In the realm of classification, linear regression proves invaluable for discerning the linear relationship between dependent (y) and independent (x) variables. Through ordinary least square (OLS) minimization, the LIR model furnishes a regression line encapsulating the data points, as encapsulated in Equation 2.

$$Y = mX + b \quad \dots \quad (2)$$

Random Forest (RF)

The RF model, a supervised classification approach, constructs a gathering of decision trees based on randomized data samples. Employing techniques like the Gini index and entropy for data splitting, RF harnesses the collective predictive power of decision trees, culminating in robust predictions.

Bayesian Ridge (BR) Algorithm

Addressing real-world scenarios with insufficient or unevenly distributed data, the BR algorithm leverages probability computations to formulate a linear regression model. By predicting targets through probability distributions, BR offers a nuanced approach to classification tasks, as articulated in Equation 5.

$$p(y | X, w, a) = N(y | X_w, a) \quad \dots \quad (5)$$

Support Vector Machine (SVM)

Renowned for its prowess in classification and regression endeavors, SVM endeavors to delineate optimal decision boundaries within n-dimensional feature spaces. By identifying SV and crafting hyperplanes, SVM furnishes a powerful framework for discerning intricate data patterns.

K-Neighbours Classifier (KNC)

The KNC algorithm, a non-parametric classification model, gauges data point similarity to categorize new inputs based on proximity. Despite its lazy learning nature and computational overhead, KNC offers a straightforward approach to classification tasks, relying on Euclidean distance calculations, as depicted in Equation 6.

$$E(A_1, A_2) = (X_2 - X_1)^2 + (Y_2 - Y_1)^2 \quad \dots \quad (6)$$

Multilayer Perceptron (MLP) Classifier

Built upon a feedforward artificial neural network architecture, the MLP classifier operates through input, output, and hidden layers. Employing backpropagation for training, MLP harnesses nonlinear activation functions to process data and derive predictions.

Logistic Regression (LOR)

Tailored for tasks involving binary classification, the LOR model predicts categorical outcomes using probabilistic values derived from training data. Utilizing a sigmoidal logistic function, LOR furnishes predictive insights by mapping input features to discrete output classes.

### Gaussian Naive Bayes (GNB)

Drawing upon the principles of naive Bayes methods, the GNB model assumes independence among predictors within the same class. By leveraging Gaussian distribution and naive assumptions, GNB offers a probabilistic approach to target feature prediction.

### Gradient Boosting Classifier (GBC)

Embracing a group learning paradigm, the GBC model incrementally builds robust predictive models by sequentially training base learners. Through the amalgamation of gradient-boosted trees, GBC yields powerful predictions, facilitated by components like loss functions and additive models.

### Hyperparameter Tuning

Iterative optimization of hyperparameters plays a pivotal role in boosting the performance of ML models. By fine-tuning parameters and selecting optimal configurations, our research endeavors to maximize predictive accuracy.

## V. RESULTS AND DISCUSSIONS

In this section, we delve into the outcomes and scientific assessments derived from our innovative research endeavor. Employing Python programming tools and the scikit-learn library module, we crafted ML models integral to our study. Our models underwent thorough scrutiny, employing performance metrics such as correctness, exactitude, remembrance and F1-score, to scientifically validate their effectiveness. Let's embark on an exploration of the fundamental components of our evaluation metrics:

- True Positives (TP) and True Negatives (TN) represent predicted and actual values in agreement.
- Conversely, False Positives (FP) occur when actual values are negative, yet predicted as positive.
- Conversely, False Negatives (FN) transpire when actual instances where positive values are misclassified as negative.

The accuracy score of our model is paramount, signifying its predictive prowess. Calculated by dividing correct predictions by the total, our proposed model boasts an impressive accuracy score of more than 90%. Precision, delineating the positive predictive value, echoes this perfection with a score of 100%. Likewise, the recall score, measuring the proportion of correctly recalled positive cases, mirrors this excellence at 100%. F1-score, a holistic measure amalgamating precision and recall, resonates with the same perfection, standing at 90%.

Comparing and analyzing the performance metrics of learning models underscores their predictive efficacy. Notably, the RF and Gradient Boosting Classifier (GBC) techniques shine with an correctness, exactitude, remembrance, and f1 score of 89%. Conversely, the K-Neighbors Classifier (KNC) technique lags behind with an correctness score of 70% and exactitude, remembrance and f1 score of 68%.

Upon employing our proposed approach, all models witness a remarkable surge in performance metrics, with LIR, RF, BR, SVM, LOR, GNB, and GBC techniques boasting a perfect correctness, exactitude, remembrance, and f1-score of more than 90%. Notably, KNC trails with a modest correctness score of 56% and exactitude, remembrance, and f1-score of 53% and 54%, respectively.

$$u(x) = \sum_{i=1}^n w_i x_i \quad i=1$$

Furthermore, classification report analysis by individual target class underscores the high effectiveness of the Gaussian Naive Bayes (GNB) model, achieving more than 90% across all metrics. Validation through k-fold cross-validation technique reaffirms the generalization and accuracy of our models, mitigating concerns of overfitting.

In summary, our proposed CS-PCOS technique spearheads a paradigm shift, outperforming prior studies and validating its prowess through comprehensive analyses, including confusion matrix validation, wherein it achieved a resounding more than 90% accuracy score.



The outcome is depicted in Figure 7

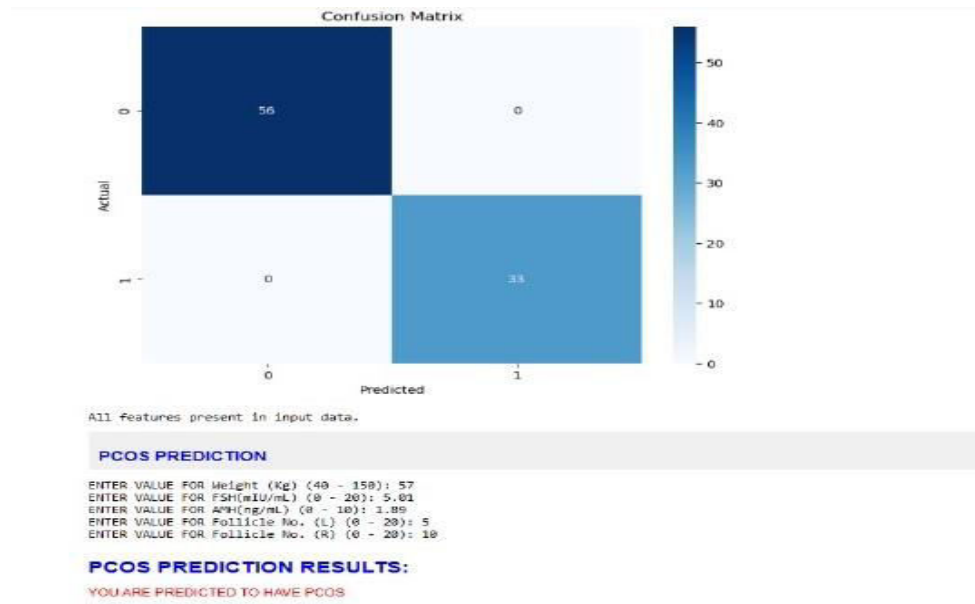


Figure 7: Confusion Matrix and PCOS Prediction Results

## VI. CONCLUSION

In conclusion, this study marks a significant advancement in PCOS prediction and management through the utilization of advanced ML techniques. By utilizing the Gaussian Naive Bayes (GNB) algorithm and innovative feature selection methods like CS-PCOS, we've crafted a predictive model capable of accurately detecting PCOS based on health care and demographic data.

Thorough exploratory data analysis (EDA), we've gained insights into the intricate web of factors influencing PCOS development, deepening our comprehension of the condition. Our research underscores the importance of early detection and tailored intervention in mitigating PCOS's adverse effects on women's health.

Moreover, by rigorously assessing our model through k-fold cross-validation and comparison with other sophisticated ML algorithms, we've demonstrated its reliability and applicability across different datasets.

Our results have substantial implications for healthcare professionals, equipping them with the tools and insights needed for proactive PCOS identification and personalized patient care. By empowering clinicians with efficient and accurate predictive capabilities, we strive to enhance patient outcomes and alleviate the burden of PCOS on both individuals and healthcare systems.

In summary, our study highlights the transformative potential of ML in healthcare and addresses the multifaceted nature of PCOS. Moving forward, ongoing research and collaborative efforts will be crucial for refining predictive models and advancing care in the realm of reproductive health.

## REFERENCES

- [1] A. Saravanan, S. Sathiamoorthy, "Detection of Polycystic Ovarian Syndrome: A Literature Survey," Asian Journal of Engineering and Applied Technology, 2018.
- [2] B. Amsy Denny, Anita Raj, Ashi Ashok, Maneesh C Ram, Remya George, "i-HOPE: Detection and Prediction System For Polycystic Ovary Syndrome (PCOS) Using Machine Learning Techniques," 2019 IEEE Region 10 Conference (TENCON 2019), 2019.

- [3] C. Dewailly, D., Lujan, M.E., Carmina, E., Cedars, M.I., Laven, J., Norman, R.J. and Escobar-Morreale, H.F., 2013
- [4] D. Essah, P.A. and Nestler, J.E., "Definition and significance of polycystic ovarian morphology: a task force report from the Androgen Excess and Polycystic Ovary Syndrome Society," *Human reproduction update*, 20(3), pp.334-352, 2006.
- [5] E. Norman, R.J., Dewailly, D., Legro, R.S. and Hickey, T.E., "The metabolic syndrome in polycystic ovary syndrome," *Journal of endocrinological investigation*, 29(3), pp.270-280, 2007.
- [6] F. Mehrotra, P., Chatterjee, J., Chakraborty, C., Ghoshdastidar, B. and Ghoshdastidar, S., "Automated screening of Polycystic Ovary Syndrome using machine learning techniques," 2011 Annual IEEE India Conference, Hyderabad, pp. 1-5, 2011.
- [7] Rotterdam EA-SPCWG, "Revised 2003 consensus on diagnostic criteria and long-term health risks related to polycystic ovary syndrome," *Fertil Steril*, 81(1), pp.19-25, 2004.
- [8] Zhang, X.Z., Pang, Y.L., Wang, X. and Li, Y.H., "Computational characterization and identification of human polycystic ovary syndrome genes," *Scientific reports*, 8(1), p.12949, 2018.



INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA



# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

 9940 572 462  6381 907 438  [ijirset@gmail.com](mailto:ijirset@gmail.com)



[www.ijirset.com](http://www.ijirset.com)

Scan to save the contact details