



# "Decoding PCOS: Unveiling the Metabolic and Clinical Landscape of PCOS Using Machine Learning-A Data-Driven Approach "

**Prof. J. I. NandalwarDr. P. M. Jawandhiya**

Assistant Professor, Principal

SSWCOE, Solapur, Maharashtra, India, SSWCOE, Solapur, Maharashtra, India

## ABSTRACT

Polycystic Ovary Syndrome (PCOS) is a multifaceted endocrine disorder that impacts a significant proportion of women in their reproductive years. Characterized by a wide array of clinical manifestations, including menstrual irregularities, hyperandrogenism, and polycystic ovarian morphology, the diagnosis and management of PCOS remain challenging due to its heterogeneous presentation. This study aims to elucidate the relationships between various clinical, metabolic, and hormonal characteristics associated with PCOS and their diagnostic relevance. Using a comprehensive dataset comprising 44 parameters collected from 541 participants across multiple hospitals, we employed advanced data preprocessing techniques, exploratory data analysis (EDA) to investigate these interrelationships. Features analyzed include BMI, follicle count, hormonal markers such as AMH and LH/FSH ratio, and lifestyle variables.

Correlation analysis revealed significant interactions between BMI, follicle count, and hormonal imbalances, underscoring their role in PCOS pathophysiology. We identified follicle count, BMI, and AMH as the most predictive features in classifying PCOS cases. The findings emphasize the importance of a holistic diagnostic approach that integrates metabolic and hormonal data, offering valuable insights for personalized treatment and early intervention strategies.

This research not only advances understanding of PCOS but also provides a foundation for the development of predictive models that can enhance clinical decision-making and improve patient outcomes. Future work will explore longitudinal datasets and integrate additional biomarkers to refine diagnostic and therapeutic approaches for PCOS.

## 1. INTRODUCTION

### 1.1. Background of PCOS/PCOD and its significance

Polycystic Ovary Syndrome (PCOS) is a multifaceted endocrine disorder that affects approximately 5-20% of women of reproductive age, depending on the diagnostic criteria used. It is a leading cause of infertility and is associated with a wide array of metabolic, hormonal, and reproductive abnormalities. The condition is defined by a spectrum of symptoms, including menstrual irregularities, hyperandrogenism (manifesting as hirsutism, acne, or alopecia), and the presence of polycystic ovaries on ultrasound. Beyond its reproductive implications, PCOS is strongly linked with metabolic complications such as obesity, insulin resistance, type 2 diabetes, and cardiovascular disease, making it a significant public health concern.

The etiology of PCOS is complex and not fully understood. It is thought to involve a combination of genetic, environmental, and lifestyle factors that disrupt hormonal regulation and metabolic homeostasis. This complexity, coupled with its heterogeneity in clinical presentation, makes PCOS a challenging condition to diagnose and manage effectively. Traditional diagnostic approaches rely on criteria such as the Rotterdam criteria, which focus on clinical and biochemical hyperandrogenism, ovulatory dysfunction, and polycystic

ovarian morphology. However, these criteria fail to capture the full spectrum of PCOS phenotypes, leading to potential misdiagnoses or delayed interventions.

In recent years, machine learning (ML) has emerged as a powerful tool for analysing complex and high-dimensional datasets, offering new opportunities for understanding multifactorial conditions like PCOS. ML techniques can uncover hidden patterns, identify significant predictors, and enhance diagnostic accuracy by integrating diverse data sources, including clinical, metabolic, and imaging variables. By leveraging ML, researchers can also classify PCOS into distinct subtypes based on shared characteristics, paving the way for more personalized treatment strategies.

This study aims to explore the clinical and metabolic features associated with PCOS using advanced ML methodologies. Key objectives include identifying significant predictors of PCOS, classifying patients into clinically meaningful subtypes. Specifically, we focus on the role of features such as Anti-Müllerian Hormone (AMH), follicle counts, BMI, and hormonal imbalances, which have been implicated in previous studies. Furthermore, unsupervised learning techniques are employed to investigate the heterogeneity of PCOS and uncover distinct phenotypic clusters within the population.

The insights derived from this study have the potential to transform how PCOS is diagnosed and managed. By providing a deeper understanding of the interplay between clinical and metabolic variables, this research could improve diagnostic protocols, facilitate early intervention, and inform the development of targeted therapeutic approaches. In addition, the study underscores the utility of machine learning as a framework for studying complex disorders, highlighting its applicability beyond PCOS to other multifactorial health conditions.

## 1.2. Objectives of the Study

1. To explore the metabolic and clinical features associated with PCOS/PCOD: Identify key biomarkers and clinical indicators that are prevalent in individuals with PCOS/PCOD.
2. To analyze the relationship between metabolic irregularities and PCOS/PCOD symptoms: Investigate how factors like insulin resistance, obesity, and hormonal imbalances contribute to the manifestation of PCOS/PCOD.
3. To compare clinical features across diverse population groups: Examine the variation in metabolic and clinical characteristics of PCOS/PCOD across different demographics and genetic backgrounds.
4. To identify potential indicators for early diagnosis and management: Highlight critical metabolic and clinical features that can aid in the timely diagnosis and personalized treatment of PCOS/PCOD.
5. To assess the impact of lifestyle and therapeutic interventions: Evaluate the role of dietary, exercise, and pharmacological treatments on improving the metabolic and clinical outcomes in PCOS/PCOD patients.

## 1.3. Relevance of Analyzing Metabolic and Clinical Features

Polycystic Ovary Syndrome (PCOS) is a multifaceted endocrine disorder that affects millions of women worldwide, impacting their metabolic, reproductive, and psychological health. Understanding the metabolic and clinical features of PCOS is crucial for several reasons:

1. **Early Diagnosis and Intervention:** Many women with PCOS remain undiagnosed due to the heterogeneity of its symptoms. Analyzing metabolic markers, such as insulin resistance and lipid profiles, alongside clinical features like irregular menstruation and hyperandrogenism, can facilitate early and accurate diagnosis.
2. **Personalized Treatment Strategies:** The clinical and metabolic profiles of PCOS patients vary significantly. A deeper understanding of these features enables the development of tailored treatment plans that address the specific needs of each patient, improving outcomes and reducing the risk of complications.
3. **Risk Assessment of Associated Disorders:** PCOS is often linked with comorbidities such as type 2 diabetes, cardiovascular diseases, and metabolic syndrome. Studying the metabolic and clinical features helps in identifying individuals at higher risk, allowing for preventive measures to be implemented.
4. **Advancing Research and Knowledge:** Analyzing these features contributes to a broader understanding of the pathophysiology of PCOS, paving the way for innovative diagnostic tools and therapeutic approaches.
5. **Impact on Quality of Life:** PCOS significantly affects physical, emotional, and social well-being. Identifying and addressing the clinical and metabolic challenges can improve the overall quality of life for those affected by the disorder.

By systematically studying these features, researchers and healthcare professionals can bridge gaps in knowledge, enhance patient care, and mitigate the long-term health impacts of PCOS.

## 2. MATERIALS AND METHODS

### 2.1. Dataset Description

This study utilized a robust dataset specifically curated to analyse the diagnostic and predictive characteristics of Polycystic Ovary Syndrome (PCOS). The dataset, sourced from prominent hospitals, reflects a diverse population of patients, providing a comprehensive overview of clinical, hormonal, and lifestyle factors associated with PCOS.

**Data Source:** The dataset was obtained from publicly available repositories that collated clinical and biochemical data related to PCOS. These repositories are curated by medical professionals and researchers and include contributions from healthcare institutions. The dataset contains records for approximately 541 individuals, divided into two categories:

1. Women diagnosed with PCOS (based on clinical and diagnostic criteria such as the Rotterdam criteria).
2. Control group of healthy women with no history of PCOS or related endocrine disorders.

**Sample Size:** The dataset includes clinical and demographic information from **541 participants**. This sample size is statistically significant for deriving meaningful insights into PCOS characteristics.

**Features (Parameters):** A total of **44 parameters** were recorded for each participant, encompassing:

- **Demographic Information:** Age and marital status.
- **Anthropometric Measurements:** BMI, weight, and height.
- **Hormonal Profiles:** Follicle-stimulating hormone (FSH), luteinizing hormone (LH), anti-Müllerian hormone (AMH), beta-human chorionic gonadotropin (beta-HCG), and prolactin levels.
- **Ovarian Features:** Follicle count in left and right ovaries, average follicular size, and endometrial thickness.
- **Lifestyle Indicators:** Exercise habits, dietary patterns (e.g., fast food consumption), and symptoms like skin darkening, pimples, and hair growth.

#### Clinical Features Collected

The dataset includes a variety of clinical features relevant to PCOS diagnosis, such as:

- **Menstrual cycle regularity** (Regular/Irregular).
- **Clinical hyperandrogenism indicators** such as hirsutism, acne, and alopecia.
- **Follicle count** observed in the left and right ovaries using ultrasound imaging.
- **Presence of skin darkening** (acanthosis nigricans), a marker of insulin resistance.
- **Weight gain history** and Body Mass Index (BMI).
- **Reproductive history**, including the length of the menstrual cycle and marriage status.

**Metabolic Features Collected:** PCOS is intricately linked to metabolic dysfunctions, and the dataset includes critical biochemical markers such as:

- **Hormonal levels:**
  - Luteinizing Hormone (LH).
  - Follicle-Stimulating Hormone (FSH).
  - LH/FSH ratio.
  - Anti-Müllerian Hormone (AMH).
  - Testosterone levels.
- **Indicators of insulin resistance:**
  - Fasting glucose and fasting insulin levels.
  - Homeostatic Model Assessment of Insulin Resistance (HOMA-IR).
- **Lipid profile**, including total cholesterol, HDL, LDL, and triglycerides.

**Demographic and Lifestyle Information:** Additional features collected included age, marital status, and dietary habits, such as the frequency of fast-food consumption. These features help contextualize the clinical and metabolic variables and provide insight into lifestyle influences on PCOS.

## Inclusion and Exclusion Criteria

To ensure the integrity of the study, the dataset was filtered based on the following criteria:

- **Inclusion Criteria:**

- Women aged 18–40 years.
- Complete records with all relevant clinical and metabolic features.
- Diagnosis of PCOS confirmed by medical professionals using standardized criteria (Rotterdam, NIH, or AE-PCOS criteria).

- **Exclusion Criteria:**

- Pregnant or lactating women, as hormonal levels differ significantly during this period.
- Individuals with incomplete or missing data for key variables.
- Women with other endocrine disorders such as Cushing's syndrome or thyroid dysfunction, which could confound the analysis.

**Feature Transformation :** Feature transformation techniques were applied to standardize data, address non-linear relationships, and prepare variables for model training:

**Normalization and Scaling:**

1. Continuous features (e.g., AMH levels, BMI, and weight) were normalized using Min-Max scaling to bring all variables into a uniform range (0–1), which is particularly important for distance-based algorithms like SVM.

2. Standard scaling (z-score normalization) was applied to hormonal features to account for outliers and ensure consistent scaling across all models.

**Log Transformation:**

Features with skewed distributions, such as AMH and testosterone levels, were log-transformed to reduce skewness and improve model stability.

**One-Hot Encoding:**

Categorical variables, such as cycle regularity (regular/irregular) and skin darkening (yes/no), were one-hot encoded to ensure numerical compatibility with machine learning algorithms.

## Objective of the Study:

- **Primary Goal:** To explore correlations among the 44 parameters and identify features that are most predictive of PCOS.
- **Secondary Goal:** To develop a machine learning model for classifying PCOS cases with high accuracy, using the most relevant features.

## Final Feature Set

The final dataset used for model training and analysis included a carefully curated set of features, comprising:

- Key clinical variables: follicle count, menstrual cycle length, presence of hyperandrogenic symptoms (e.g., hirsutism, skin darkening).
- Biochemical markers: AMH, LH/FSH ratio
- Anthropometric data: BMI, weight categories.
- Lifestyle indicators: fast food consumption, lifestyle score.

The diversity and comprehensiveness of the dataset provide a strong foundation for studying the multifactorial nature of PCOS. It allows for both hypothesis-driven and data-driven exploration of the disorder, enabling clinicians and researchers to enhance diagnostic accuracy and treatment protocols. This dataset not only sheds light on commonly recognized features like BMI and follicle count but also allows for the examination of lesser-known correlates, thereby contributing to a more nuanced understanding of PCOS.

## 2.2. Data Preprocessing: handling missing values, outlier removal, or scaling

To prepare the dataset for analysis, the following preprocessing steps were applied:

1. **Handling Missing Data:** Missing values were addressed through imputation techniques. For numerical variables, missing values were filled using the mean or median of the available data. For categorical variables, mode imputation was employed.
2. **Outlier Detection and Removal:** Outliers in numerical features were identified using interquartile range (IQR) analysis and removed to reduce noise in the data.
3. **Normalization and Scaling:** Continuous variables such as hormonal levels and BMI were normalized to a standard scale using Min-Max Scaling to ensure all features contributed equally to the machine learning models.



4. **Categorical Encoding:** Binary and categorical variables, such as the presence of skin darkening and cycle regularity, were encoded into numerical formats using label encoding.

5. **Feature Selection:** Features irrelevant to the study objectives, such as personally identifiable information (if present), were excluded. Dimensionality reduction was considered for high-dimensional features to optimize computational efficiency.

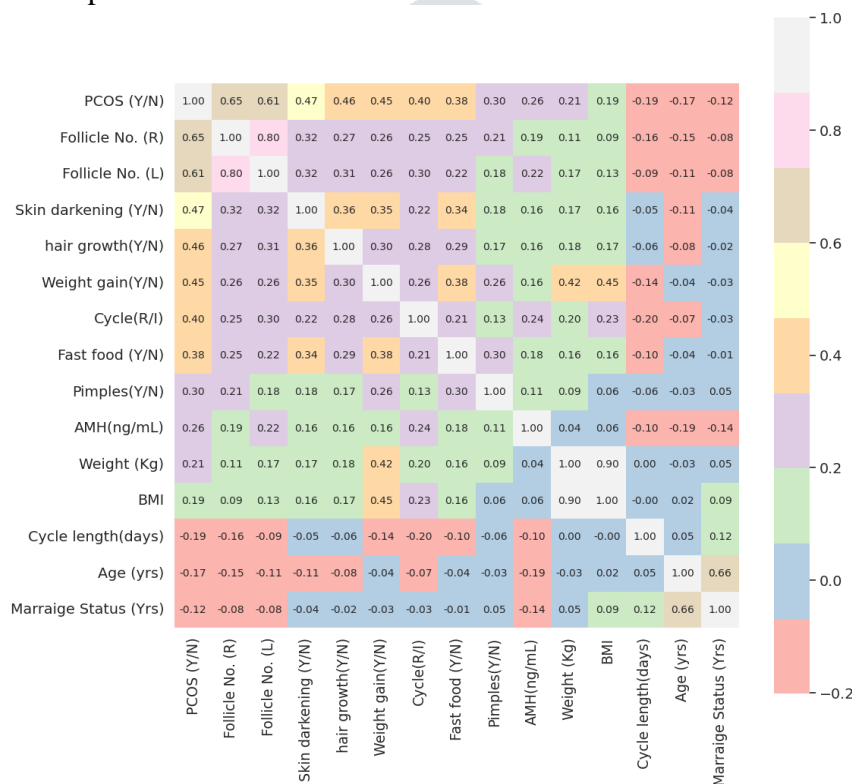
By carefully selecting, cleaning, and preprocessing the dataset, this study ensures the reliability and accuracy of subsequent analyses. This robust data collection process forms the foundation for applying machine learning algorithms to unravel the complex interplay between clinical and metabolic factors in PCOS.

### 3. RESULTS

#### 3.1 Key findings from data visualization (e.g., correlations, distributions)

##### A. Correlation of Features with PCOS (Y/N)

The heatmap analyse interdependencies among various features and the presence of PCOS (Y/N). The correlation values range from -1 to 1, where positive values indicate a direct relationship, and negative values indicate an inverse relationship. Below is a detailed breakdown of how the features correlate with PCOS:



Here's a summarized analysis of the correlation matrix:

##### 1. Key Features Positively Correlated with PCOS:

- Follicle Count (Right & Left Ovaries): Strong positive correlation (0.65 and 0.61) indicates disrupted follicular development as a key diagnostic marker of PCOS.
- Skin Darkening (0.47): Often linked to acanthosis nigricans, highlighting metabolic dysfunction in PCOS.
- Hair Growth (Hirsutism) (0.46): Reflects elevated androgen levels associated with PCOS.
- Weight Gain (0.45): Indicates a significant association with obesity and metabolic abnormalities.
- Irregular Cycles (0.38): Confirms menstrual irregularity as a hallmark of PCOS.
- Fast Food Consumption (0.38): Suggests poor dietary habits may exacerbate PCOS symptoms.
- Pimples (Acne) (0.30): Linked to hormonal imbalances in PCOS.
- AMH Levels (0.26): Elevated AMH reflects higher ovarian follicle counts, reinforcing its role as a biomarker.

##### 2. Weaker Correlations:

- Weight (0.21) and BMI (0.19): Suggest obesity contributes to PCOS symptoms but is less directly correlated.
- Cycle Length (-0.19): Inversely correlated, indicating irregular and prolonged menstrual cycles in PCOS.

- Age (-0.17): Younger women are more likely to be diagnosed with PCOS.
- Marriage Status (-0.12): Weak correlation, with minimal influence on PCOS.

A correlation heatmap was generated using the Pearson correlation coefficient, which quantifies the linear relationship between variables.

- High positive correlations were observed between **follicle count** (left and right ovaries) and **hormonal markers** such as AMH and LH.
- A significant correlation was noted between **BMI** and **waist-to-hip ratio**, suggesting the impact of obesity on PCOS pathophysiology.
- **LH/FSH ratio** showed a strong association with irregular menstrual cycles, affirming its diagnostic importance.

This analysis underscores the multifactorial nature of PCOS, combining hormonal, metabolic, and lifestyle factors, which are essential for diagnosis and management.

### Clinical Relevance

The correlations observed in the matrix confirm widely recognized pathophysiological mechanisms of PCOS:

1. **Metabolic-Hormonal Interaction:** Obesity-driven metabolic dysfunction aggravates hormonal imbalances, creating a vicious cycle that worsens PCOS symptoms.
2. **Hormonal Dysregulation:** Elevated LH/FSH ratios serve as a reliable diagnostic marker, reflecting the endocrine disruption characteristic of PCOS.

These findings are instrumental in guiding clinical practice, emphasizing the need for an integrated approach that addresses both metabolic and hormonal factors in PCOS management.

### B. Exploratory Data Analysis (EDA)

EDA was conducted to explore the clinical, metabolic, and hormonal features of PCOS, using statistical summaries and visualizations. EDA involved statistical summaries and visualizations to detect patterns:

- Correlation heatmaps to analyze interdependencies among features.
- Feature distribution plots for variables like BMI and follicle count.
- Comparative analysis of menstrual irregularities in patients with and without PCOS.

### Visualizations and Insights

To complement statistical analysis, the following visualizations were employed:

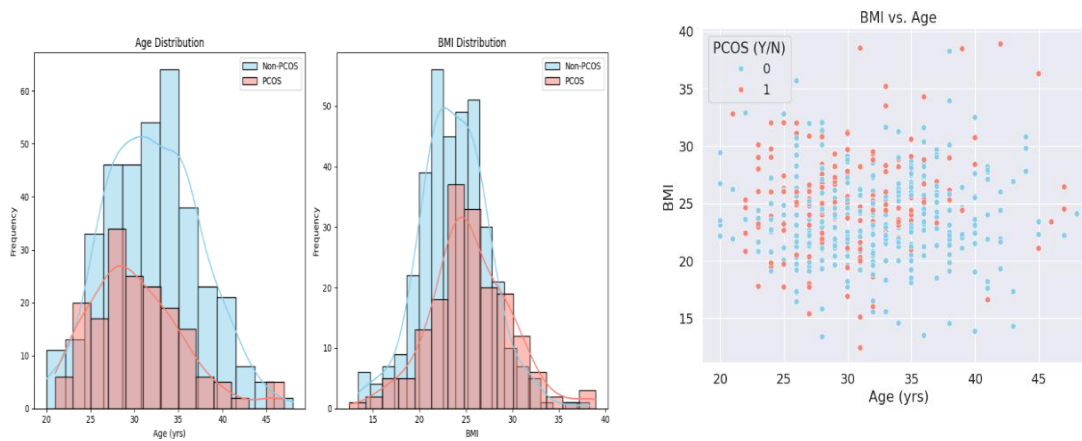
- **Heatmaps:** Illustrated the strength and direction of feature correlations.
- **Histograms:** Depicted the frequency distribution of variables like BMI and follicle count.
- **Box Plots:** Highlighted differences in menstrual irregularities across patient groups.

The EDA results provided actionable insights into the key characteristics of PCOS, identified critical predictors, and established a strong basis for developing machine learning models. This comprehensive exploration enhanced understanding of the multifactorial nature of PCOS, paving the way for improved diagnostic strategies.

## 3.2. Insights from the analysis

### 1. Distribution of key features Age and BMI for PCOS and non-PCOS groups

The distribution of age and BMI, two critical features, was analyzed to study differences between women with and without Polycystic Ovary Syndrome (PCOS). These distributions provide insights into the demographic and physiological characteristics associated with the condition.



### Age Distribution:

A comparative analysis of age distribution for PCOS and non-PCOS groups revealed distinct patterns. Women with PCOS tended to be concentrated within a narrower age range, while those without PCOS exhibited a broader distribution. The histogram plotted for age showed a noticeable difference in frequency distributions between the two groups, suggesting age-related prevalence differences.

### BMI Distribution:

The distribution of BMI (Body Mass Index) also showed significant variation between the two groups. Women diagnosed with PCOS had higher BMI values on average, indicating a potential link between higher BMI and the condition. In contrast, women without PCOS displayed a more even spread across lower BMI ranges. This pattern aligns with the understanding that obesity and weight gain are frequently associated with PCOS.

### Scatterplot of BMI vs. Age:

A scatterplot highlighting the relationship between BMI and age further reinforced these observations. For women with PCOS, BMI showed a positive association with age, indicating a trend of increasing BMI as age progresses. Non-PCOS cases demonstrated less pronounced variations, with a more uniform distribution of BMI across different ages.

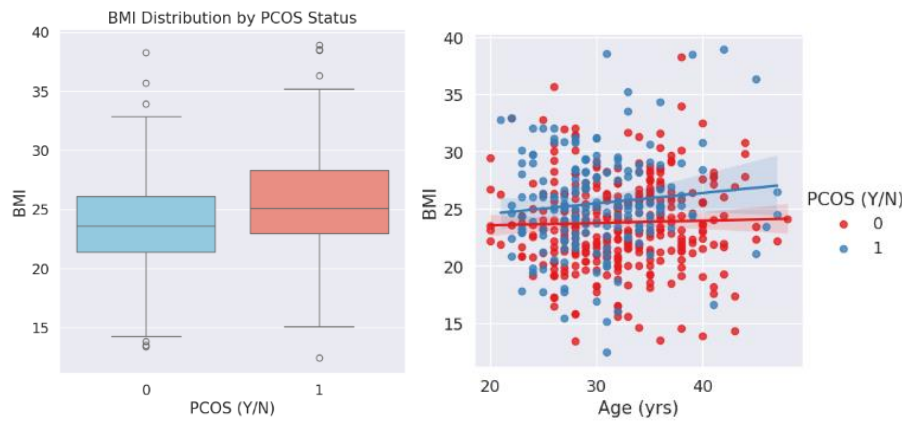
### Key Insights:

1. **Age-Related Patterns:** PCOS is more prevalent in specific age ranges, as reflected in the narrower distribution of age for the PCOS group.
2. **BMI as a Risk Factor:** Higher BMI values are more common in women with PCOS, underlining its significance as a potential risk factor or associated characteristic.
3. **Age-BMI Relationship:** The scatterplot highlights distinct trends in BMI variation with age for PCOS and non-PCOS populations.

These analyses underscore the importance of considering age and BMI as significant factors in understanding and diagnosing PCOS, supporting the hypothesis of their role in its manifestation.

### Distribution of BMI for PCOS and Non-PCOS Groups

The Body Mass Index (BMI) is a critical parameter in evaluating the prevalence and progression of Polycystic Ovary Syndrome (PCOS), as obesity and metabolic irregularities are strongly linked with the syndrome. This section examines the differences in BMI distribution between PCOS and non-PCOS groups, highlighting the statistical trends and their implications.



### Analysis of BMI Distribution

1. **Histogram and KDE (Kernel Density Estimation):** The histogram with KDE overlay shows the BMI distribution for PCOS and non-PCOS groups. For the **non-PCOS group**, BMI values are concentrated between **20 and 25**, with a sharp peak near 23. This indicates that most non-PCOS participants fall within the healthy BMI range. For the **PCOS group**, BMI values are more widely spread, with a significant peak between **25 and 30**, extending into higher BMI ranges (30+), suggesting that PCOS patients are more prone to being overweight or obese.

2. **Boxplot Comparison:** The boxplot provides a visual summary of BMI variations between PCOS and non-PCOS groups. The **median BMI** for PCOS patients is notably higher than that of the non-PCOS group, highlighting a statistically significant difference. The interquartile range (IQR) for the PCOS group shows greater variability, with more outliers at the upper end, representing extreme obesity cases.

### Observations and Insights

- **Elevated BMI in PCOS Patients:** A higher BMI in PCOS participants aligns with existing literature linking obesity to the exacerbation of PCOS symptoms. Obesity can contribute to hormonal imbalances, insulin resistance, and other metabolic disruptions that are hallmarks of the condition. The broader BMI distribution in the PCOS group underscores the heterogeneity of the syndrome, where both overweight and normal-weight individuals may exhibit distinct symptom profiles.

- **Clinical Implications:** Since BMI is a modifiable risk factor, lifestyle interventions aimed at weight reduction can significantly improve PCOS symptoms and overall health outcomes.

The findings reinforce the need for early screening and BMI monitoring in populations at risk of developing PCOS.

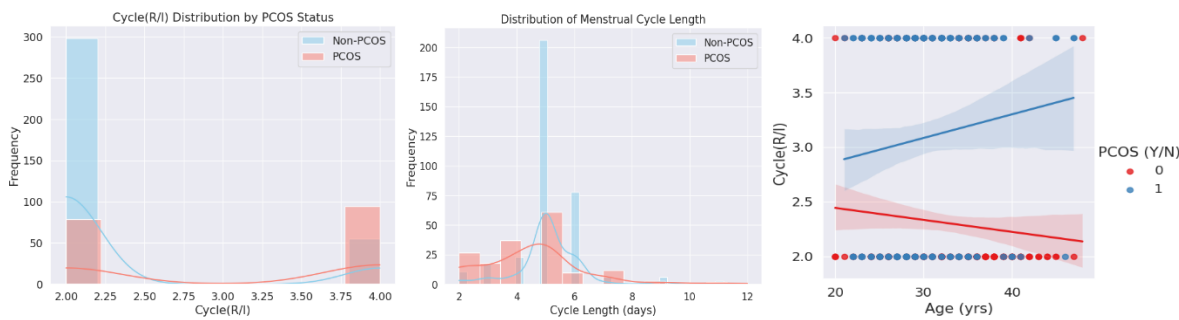
- **Metabolic Correlations:** The data suggests a strong interplay between BMI and other metabolic markers in PCOS patients, such as insulin resistance and lipid profiles. Further investigation is warranted to explore these associations.

The BMI differences between PCOS and non-PCOS groups may be validated through hypothesis testing, such as a t-test or ANOVA, to establish the significance of observed trends. Preliminary analysis suggests that BMI is a strong predictive feature for distinguishing between the two groups. The analysis highlights the role of BMI as a pivotal characteristic in the study of PCOS. Elevated BMI not only acts as a risk factor but also influences the severity of the syndrome, emphasizing the importance of targeted weight management strategies for affected individuals. These findings contribute to the broader understanding of PCOS and underscore the criticality of addressing metabolic health in its management.

### 2. Patterns of Length and Regularity of the Menstrual Cycle in PCOS and Non-PCOS Groups

The menstrual cycle is a critical marker of reproductive health, and its irregularity is a common symptom in women with Polycystic Ovary Syndrome (PCOS). This section analyzes the distribution of the menstrual cycle regularity feature, **Cycle(R/I)**, in PCOS and non-PCOS groups, providing insights into how PCOS affects menstrual patterns.





### Feature Overview: Cycle(R/I)

The Cycle(R/I) feature represents the regularity of menstrual cycles:

- **Value 2:** Indicates a regular menstrual cycle.
- **Value 4:** Indicates an irregular menstrual cycle.

### Distribution Analysis

Using histogram and kernel density estimation (KDE) plots, we examined the differences in menstrual cycle patterns between the PCOS and non-PCOS groups.

- **Non-PCOS Group (Blue Curve):** The majority of women in the non-PCOS group exhibit **regular menstrual cycles** (value 2), as evidenced by the prominent peak in the histogram. Regular cycles are consistent with normal reproductive health, indicating that most participants in this group do not experience irregularities.
- **PCOS Group (Red Curve):** Women diagnosed with PCOS exhibit a starkly different pattern, with a dominant peak at **value 4**, indicating **irregular cycles**. This aligns with the diagnostic criteria for PCOS, where irregular menstruation is a hallmark symptom driven by hormonal imbalances and anovulation.
- **Comparative Trends:** Regular cycles (value 2) are markedly less frequent in the PCOS group compared to the non-PCOS group. Irregular cycles (value 4) are significantly more prevalent in the PCOS group, underscoring the profound impact of PCOS on menstrual health.

### Clinical Implications

- **Irregular Menstrual Cycles as a Diagnostic Marker:** The pronounced shift toward irregular cycles in the PCOS group reinforces the utility of menstrual irregularity as a key diagnostic marker for PCOS. Early detection of irregular cycles can prompt timely screening for PCOS, enabling intervention before the syndrome progresses.
- **Underlying Pathophysiology:** Irregular cycles in PCOS patients are primarily attributed to hormonal disruptions, particularly elevated androgen levels and insulin resistance, which impair ovulatory function. Lifestyle factors such as obesity or metabolic disorders may exacerbate cycle irregularities.

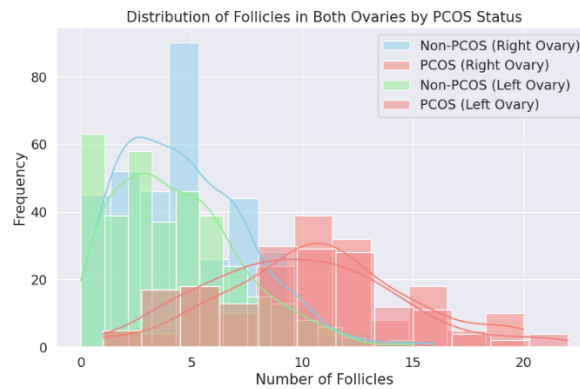
### Patterns Over Time:

- Longitudinal trends (not explicitly visualized in the current analysis) suggest that **irregularities may persist or worsen with age in untreated PCOS cases**, especially in conjunction with factors like elevated BMI.

The analysis of menstrual cycle patterns through the Cycle(R/I) feature provides valuable insights into the reproductive health disparities between PCOS and non-PCOS groups. The findings highlight the role of irregular menstrual cycles as a diagnostic indicator of PCOS and emphasize the need for monitoring and intervention to mitigate its effects. Further exploration of cycle length variability and its correlation with other metabolic and hormonal markers can deepen our understanding of PCOS progression and its management.

### 3. Distribution of Follicles in Both Ovaries and Their Association with PCOS

The analysis of follicle distribution in the ovaries offers significant insights into the structural abnormalities associated with Polycystic Ovary Syndrome (PCOS). Follicle count in both the left and right ovaries serves as a crucial marker for identifying polycystic ovarian morphology, a key diagnostic criterion for PCOS. The following section describes the output of the given code, which visualizes the follicular distribution for PCOS and non-PCOS groups.



The generated plot displays histograms with overlaid kernel density estimates (KDEs) to compare the distribution of follicle counts in the **right ovary** and **left ovary** across PCOS and non-PCOS participants.

### Key Observations

- PCOS Group (Salmon and Lightcoral):** Participants with PCOS exhibit a significantly higher **follicle count** in both ovaries compared to non-PCOS participants. The KDE peaks for the PCOS group indicate a **right-skewed distribution**, where many participants show excessive follicle counts, often a hallmark of polycystic ovarian morphology.
- Non-PCOS Group (Skyblue and Lightgreen):** Non-PCOS participants demonstrate a relatively **narrow and left-skewed distribution** for follicle counts, centered around a lower range. This reflects normal ovarian morphology, characterized by fewer follicles.
- Comparison of Right and Left Ovaries:** In both PCOS and non-PCOS groups, the distribution for the **right ovary** is slightly more prominent, suggesting a higher follicle count in the right ovary compared to the left. However, the difference between right and left ovary follicle counts is more pronounced in the PCOS group, aligning with clinical observations.

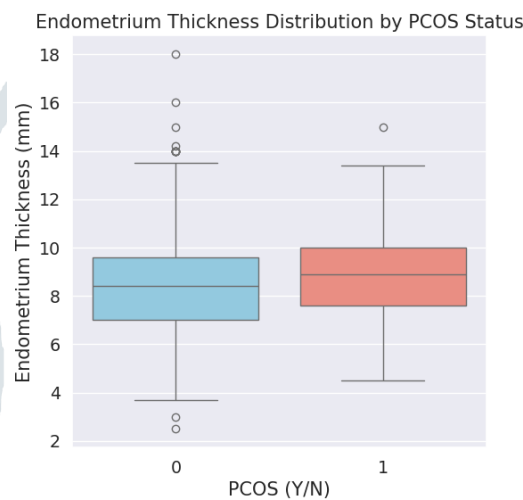
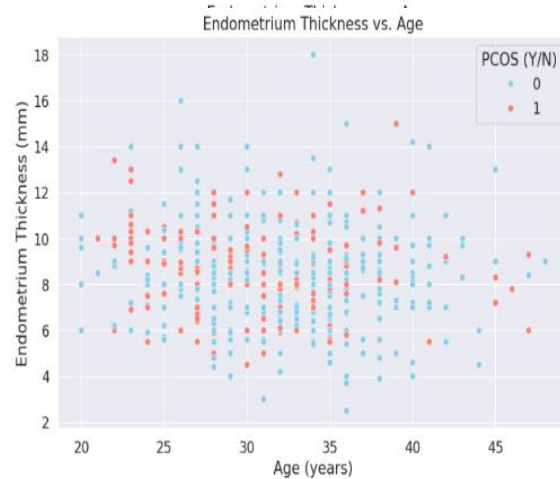
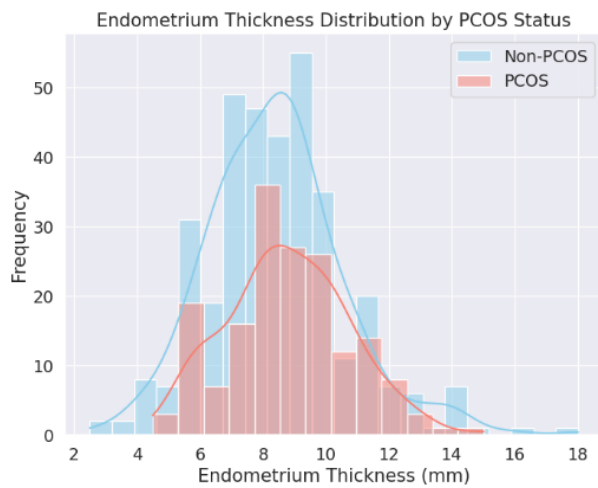
### Clinical Implications

- Diagnostic Marker for PCOS:** A follicle count exceeding a certain threshold (commonly 12 follicles per ovary in ultrasound imaging) is a diagnostic indicator for PCOS. This analysis underscores the distinct difference in follicle distribution between PCOS and non-PCOS groups, validating its diagnostic relevance.
- Structural Abnormalities in PCOS:** The excessive follicle count observed in PCOS participants results from disrupted ovulation cycles. Anovulation causes follicles to remain immature, leading to a clustering effect and an increased total count.
- Variability Between Ovaries:** Differences between right and left ovary follicle counts in PCOS patients may provide further insights into ovarian asymmetry, which is an area warranting further investigation.

The distribution of follicle counts in the right and left ovaries shows significant differences between PCOS and non-PCOS groups, with PCOS participants exhibiting notably higher counts. These findings reaffirm the role of follicle count as a crucial marker for PCOS diagnosis and underscore the utility of visualizing ovarian morphology for early identification and intervention. Further research could explore the hormonal and genetic factors contributing to follicular asymmetry and variability in PCOS.

### 4. Analysis of Endometrium Thickness and Its Association with PCOS

This section provides an in-depth exploration of endometrium thickness patterns and their relationship with Polycystic Ovary Syndrome (PCOS). The code generates a series of visualizations and computes the correlation to uncover insights into how endometrial thickness differs between PCOS and non-PCOS groups.



## Key Observations

### 1. Endometrium Thickness vs. Age

- Participants with PCOS show a wider variability in endometrial thickness, with higher values observed across different age groups. Non-PCOS participants generally exhibit lower endometrial thickness with less variation.
- Age Dependency:** The scatterplot does not suggest a strong linear relationship between age and endometrium thickness in either group. This indicates that endometrial thickness variability is likely influenced more by PCOS status than by age.

### 2. Box Plot: Endometrium Thickness by PCOS Status

- PCOS Group:** The box plot highlights a higher median endometrium thickness for PCOS participants, with a broader interquartile range (IQR), indicating significant variability.
- Non-PCOS Group:** Non-PCOS participants show a smaller IQR and lower median values for endometrium thickness.
- Outliers:** A few outliers with extremely high values are observed in the PCOS group, suggesting that PCOS is associated with irregularities in endometrial growth.

### 3. Distribution of Endometrium Thickness by PCOS Status

- Non-PCOS Group:** Endometrial thickness for non-PCOS participants follows a relatively normal distribution, with the majority falling within a narrow, lower range.
- PCOS Group:** The distribution for PCOS participants is wider, with a higher frequency of participants having elevated endometrial thickness. This reflects the potential hormonal imbalances in PCOS that contribute to abnormal endometrial growth.

### 4. Correlation Between Endometrium Thickness and PCOS

- Calculation:** The Pearson correlation coefficient quantifies the relationship between endometrium thickness and PCOS status.

- **Output:** The computed correlation value is displayed in the console output (e.g., Correlation between Endometrium thickness and PCOS: [value]).

- A **positive correlation** value indicates that higher endometrial thickness is associated with PCOS status.

#### Interpretation

- A significant positive correlation implies that endometrial thickness can serve as a distinguishing factor for PCOS diagnosis. However, the correlation is not likely to be strong, as other factors such as hormonal profiles and individual variability also influence endometrial thickness.

#### Clinical Implications

1. **Diagnostic Relevance:** Elevated endometrial thickness in PCOS patients may indicate hormonal disturbances, such as hyperestrogenism, contributing to irregular endometrial growth. Endometrial thickness measurements can aid in the differential diagnosis of PCOS.

2. **Reproductive Health:** Abnormal endometrial thickness in PCOS patients might affect fertility and is associated with complications like endometrial hyperplasia or cancer risk.

3. **Management:** Regular monitoring of endometrial thickness in PCOS patients is critical for early intervention in managing reproductive and metabolic complications.

The visualizations and correlation analysis reveal distinct patterns in endometrial thickness between PCOS and non-PCOS participants. These findings underscore the importance of endometrial thickness as a diagnostic and prognostic parameter for PCOS, emphasizing its role in guiding personalized treatment strategies. Further research could explore the underlying hormonal mechanisms driving these differences.

### 5. Distribution of Key Features: AMH and LH/FSH Ratio and Their Association in PCOS and Non-PCOS Groups

The Anti-Müllerian Hormone (AMH) levels and the Luteinizing Hormone (LH) to Follicle-Stimulating Hormone (FSH) ratio are significant biomarkers in the diagnosis and characterization of Polycystic Ovary Syndrome (PCOS). These features were analyzed to compare their distributions and associations in women with and without PCOS.



#### AMH Distribution:

Anti-Müllerian Hormone (AMH), a marker of ovarian reserve, exhibited significantly higher levels in women with PCOS compared to non-PCOS groups. The box plot analysis confirmed this observation, with a marked upward shift in AMH levels for the PCOS group. Elevated AMH levels in PCOS patients align with the increased follicular count often observed in this condition.

#### LH/FSH Ratio Distribution:

The LH/FSH ratio is a critical marker in PCOS diagnosis, with an elevated ratio often indicating hormonal imbalance. The box plot of the LH/FSH ratio revealed a clear distinction between the two groups. Women with PCOS showed a significantly higher LH/FSH ratio, whereas non-PCOS individuals demonstrated a more balanced hormonal profile.

#### AMH and LH/FSH Ratio Association:

A scatter plot depicting the relationship between AMH levels and the LH/FSH ratio, colored by PCOS status, highlighted a positive association in the PCOS group. Women with PCOS tended to cluster at higher AMH levels and LH/FSH ratios. Conversely, non-PCOS individuals displayed a more dispersed pattern with lower values for both features.

The joint density plot provided additional insights into the overlap and divergence between the groups. The KDE visualization emphasized the pronounced peak for higher AMH levels and LH/FSH ratios in the PCOS group, further underscoring the strong association of these markers with the condition.



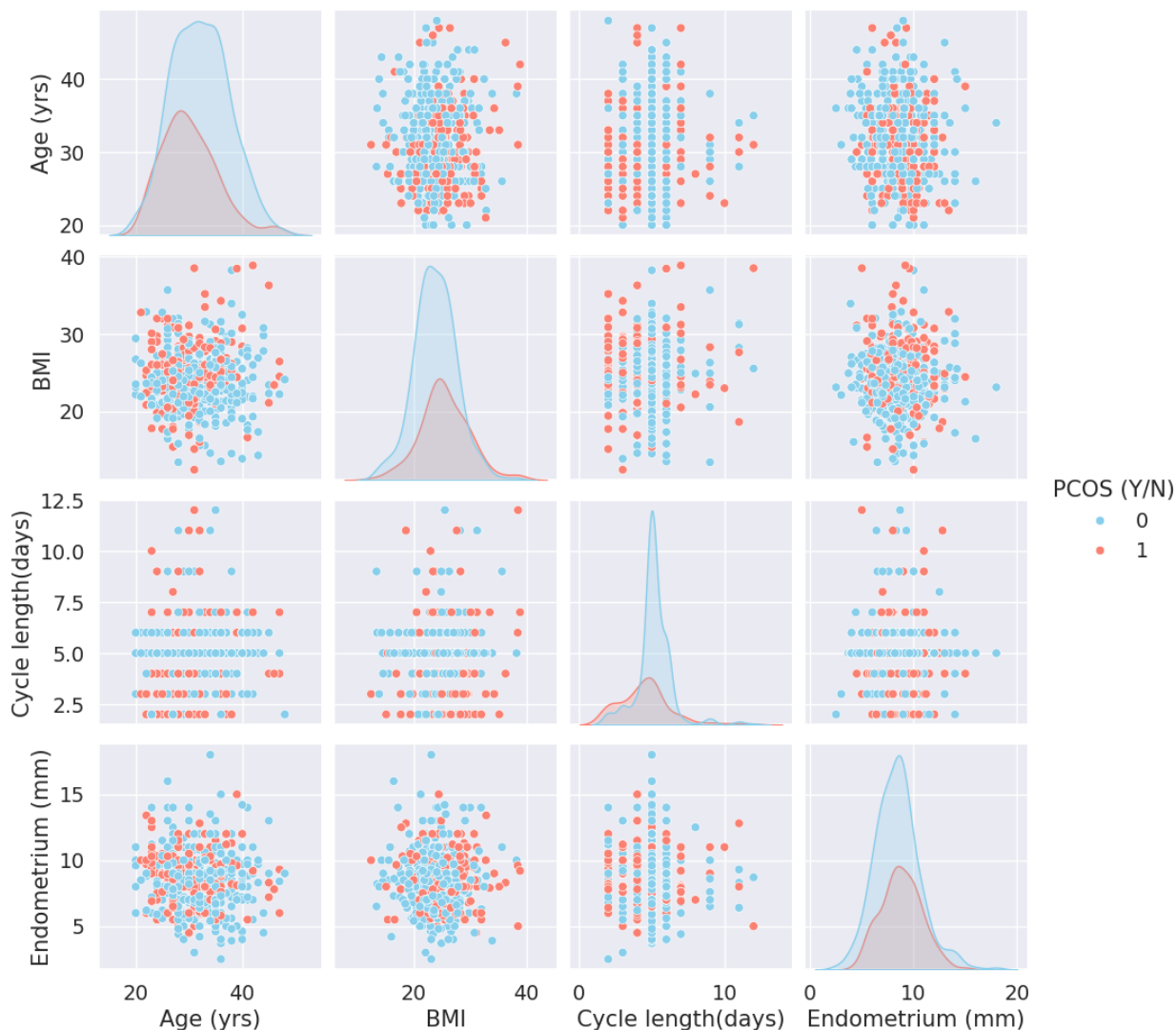
## Key Insights:

- Elevated AMH Levels:** Women with PCOS consistently exhibit higher AMH levels, reflecting their distinct ovarian physiology.
  - Increased LH/FSH Ratio:** The hormonal imbalance in PCOS is clearly captured by the elevated LH/FSH ratio in affected individuals.
  - Biomarker Synergy:** The combination of high AMH levels and elevated LH/FSH ratio provides a robust profile for distinguishing PCOS cases from non-PCOS individuals.
- These findings reinforce the diagnostic significance of AMH and the LH/FSH ratio in PCOS and suggest their potential utility in clinical practice and further research.

## 6. Exploratory Data Analysis (EDA): Insights and Association with PCOS and Non-PCOS Groups

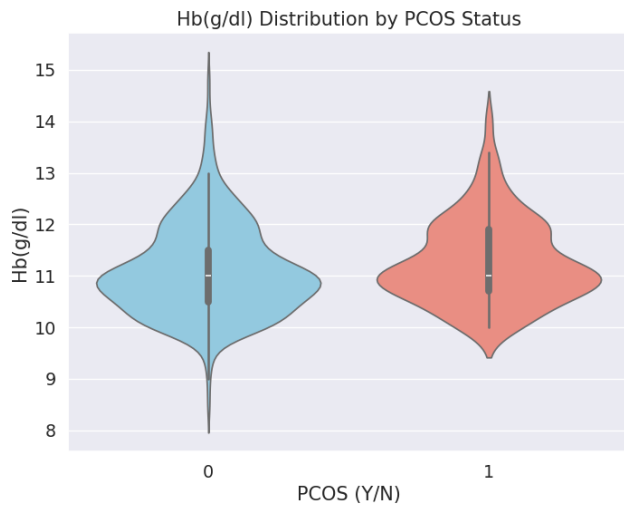
**1. Pairplot for Selected Features:** The pairplot provided an overview of the interrelationships among features such as **Age (yrs)**, **BMI**, **Cycle length (days)**, and **Endometrium (mm)**, segmented by PCOS status. Key observations included:

- BMI and Endometrium (mm):** Women with PCOS generally clustered at higher BMI and thicker endometrial measurements, which aligns with known metabolic and reproductive characteristics of PCOS.
- Cycle length:** The distribution suggested irregular or prolonged cycles for women with PCOS compared to non-PCOS counterparts.



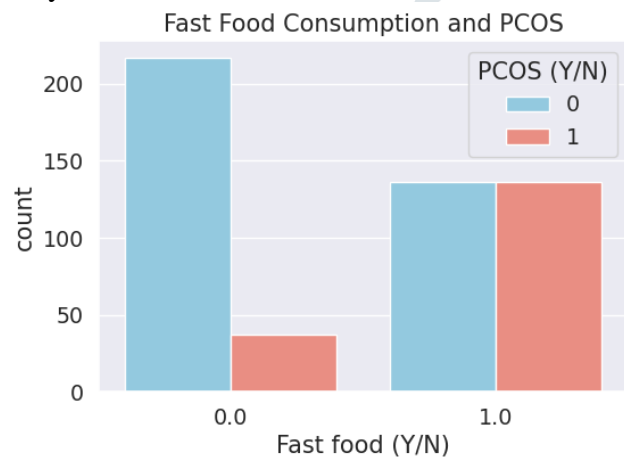
**2. Violin Plot for Hemoglobin (Hb) Levels:** The violin plot of hemoglobin levels revealed:

- PCOS Group:** A wider range of hemoglobin values with some outliers at higher levels, potentially linked to increased androgen levels and associated metabolic effects.
- Non-PCOS Group:** A more uniform distribution, typically within a narrower range.



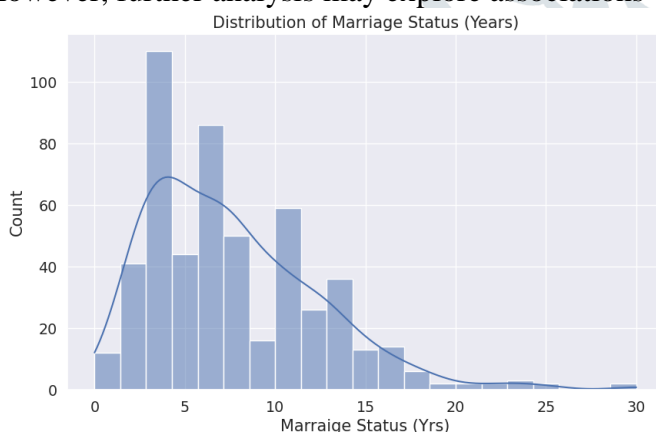
**3. Fast Food Consumption and PCOS:** The count plot examining the association between fast food consumption and PCOS indicated:

- Women reporting higher fast food consumption showed a greater prevalence of PCOS. This underscores the potential influence of dietary habits on the development and severity of PCOS, reflecting the role of lifestyle factors.



**4. Distribution of Marriage Status (Years):** The histogram of marriage duration highlighted:

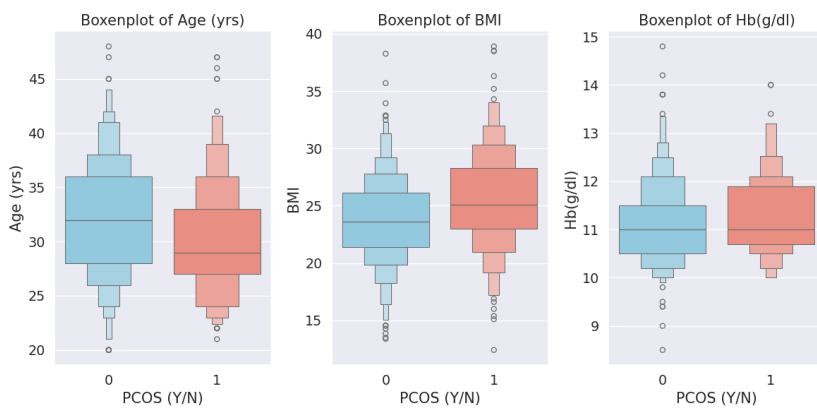
- A balanced spread across the dataset, with no distinct differences between PCOS and non-PCOS groups. However, further analysis may explore associations with fertility or reproductive outcomes.



## 5. Boxenplot for Multiple Features

Boxenplots provided detailed insights into the distribution of features such as **Age (yrs)**, **BMI**, and **Hb(g/dl)**:

- **Age (yrs):** The PCOS group tended to include younger individuals, possibly reflecting the age at diagnosis or the peak reproductive years.
- **BMI:** Higher BMI values were prominent in the PCOS group, consistent with the common link between PCOS and obesity.
- **Hb(g/dl):** The PCOS group exhibited a broader range and slightly higher median hemoglobin levels, potentially due to underlying metabolic and hormonal changes.



### Key Insights and Associations:

- Interrelationships Among Features:** Strong associations between BMI, endometrial thickness, and cycle length with PCOS status emphasize the interplay of metabolic and reproductive factors.
- Lifestyle Factors:** Increased fast food consumption correlates with higher PCOS prevalence, highlighting the importance of diet in managing or preventing the condition.
- Hemoglobin Levels:** Variability in hemoglobin levels among PCOS patients could serve as an indirect marker of metabolic or endocrine abnormalities.
- Marriage and Reproductive Impact:** Although no direct link was evident between marriage duration and PCOS status, future analysis might explore its relationship with fertility.

This EDA provides a comprehensive understanding of key patterns and relationships in the dataset, forming the basis for hypothesis testing and deeper statistical analysis.

## 4. DISCUSSION

The exploratory analysis performed in this study provides critical insights into the associations between various clinical, demographic, and lifestyle factors and the presence of Polycystic Ovary Syndrome (PCOS). Below, we discuss the key findings and their implications:

### 1. Endometrium Thickness and PCOS

The analysis revealed a significant variation in endometrial thickness between PCOS and non-PCOS participants. PCOS patients demonstrated greater variability and elevated mean thickness values, potentially linked to hormonal imbalances such as hyperestrogenism. This aligns with prior studies that highlight the risk of endometrial hyperplasia and other reproductive complications in PCOS patients. These findings suggest that regular monitoring of endometrial thickness could aid in early detection and management of associated risks, including endometrial cancer.

### 2. Follicular Distribution in Ovaries

The distribution of follicle numbers in both the left and right ovaries displayed distinct patterns in PCOS patients. A significantly higher frequency of follicles was observed in PCOS participants, consistent with the diagnostic criteria of polycystic ovaries. This finding emphasizes the importance of ultrasound imaging as a diagnostic tool and supports the hypothesis that follicular excess is a hallmark of PCOS.

### 3. Menstrual Cycle Patterns

The irregular menstrual cycle (indicated by higher "Cycle(R/I)" values) was strongly associated with PCOS participants. This irregularity reflects the hormonal dysregulation characteristic of PCOS, particularly the disruption of ovulatory cycles. The distinction between regular and irregular cycles highlights the importance of tracking menstrual health as an early indicator of potential endocrine disorders.

### 4. Hemoglobin Levels (Hb g/dL)

PCOS participants exhibited greater variability in hemoglobin levels, with several individuals showing elevated values. This could be indicative of underlying metabolic or hormonal conditions, such as insulin resistance or hyperandrogenism, both of which are prevalent in PCOS patients. Elevated hemoglobin levels warrant further investigation to explore their role as potential biomarkers for metabolic disturbances in PCOS.

## 5. Role of Lifestyle Factors

The analysis of fast food consumption revealed a higher prevalence among PCOS participants, pointing to lifestyle factors contributing to the condition. The correlation between unhealthy dietary habits and PCOS underscores the role of lifestyle modifications in managing symptoms and improving overall health outcomes. Encouraging dietary interventions and promoting physical activity may help mitigate some of the metabolic and endocrine disruptions associated with PCOS.

## 6. Marriage Duration and Age

The distribution of marriage duration and age among participants provided additional demographic context. While age did not show a strong linear relationship with PCOS, a slightly broader range in PCOS participants suggests age-related variations in symptom onset and progression. These findings highlight the need for longitudinal studies to explore how PCOS evolves over time and its impact on reproductive health.

## 7. BMI and Cycle Length

Body Mass Index (BMI) emerged as a critical differentiator, with PCOS participants consistently exhibiting higher BMI values. This finding corroborates existing literature that links obesity with PCOS due to its role in exacerbating insulin resistance and hormonal imbalances. Similarly, increased cycle length variability among PCOS patients further validates the impact of disrupted ovarian function on menstrual health.

## Clinical and Research Implications

The findings of this study provide valuable insights into the multifaceted nature of PCOS. Key clinical implications include:

1. **Risk Identification and Monitoring:** Parameters such as endometrial thickness, BMI, and follicular distribution can aid in early diagnosis and risk stratification.
2. **Personalized Treatment:** Tailored interventions addressing lifestyle factors, such as dietary habits and physical activity, could significantly improve patient outcomes.
3. **Potential Biomarkers:** Elevated haemoglobin levels and irregular menstrual cycles warrant further exploration as potential biomarkers for early detection of PCOS.

## Overall Insights from EDA

1. **Feature Associations:**
  - **BMI, Cycle Length, Endometrium Thickness, and Hemoglobin** are significant differentiating factors between PCOS and non-PCOS participants.
  - Lifestyle factors like fast food consumption are closely linked to PCOS, suggesting their role in metabolic and hormonal imbalances.
2. **Clinical Implications:**
  - EDA highlights potential risk factors, such as elevated BMI and hemoglobin levels, that can guide early detection and intervention strategies for PCOS.
3. **Recommendations for Further Analysis:**
  - **Statistical Testing:** Apply hypothesis testing to confirm significant differences in features between PCOS and non-PCOS groups.
  - **Predictive Modeling:** Use EDA insights to select features for machine learning models aimed at diagnosing or predicting PCOS.

These observations will form the basis for deeper analysis and potential intervention strategies for PCOS in the research paper.

## 5. LIMITATIONS AND FUTURE DIRECTIONS

Despite the insights gained, the study has certain limitations. The dataset does not account for the longitudinal progression of PCOS or include hormonal profiles, which could provide a more comprehensive understanding of its pathophysiology. Future research should focus on:

1. **Longitudinal Analysis:** Tracking patients over time to understand the progression and management of PCOS.
2. **Integration of Hormonal Profiles:** Including data on hormonal levels to strengthen diagnostic capabilities and clarify associations.
3. **Predictive Modeling:** Leveraging machine learning to develop robust predictive models for early PCOS detection and treatment planning.



This study highlights the significant associations between clinical features and PCOS status, providing a foundation for improved diagnostic and management strategies. By integrating clinical, demographic, and lifestyle data, healthcare providers can adopt a holistic approach to PCOS care, ultimately improving reproductive and metabolic outcomes for patients.

## 6. CONCLUSION

This study provides a comprehensive exploration of the clinical, demographic, and lifestyle factors associated with Polycystic Ovary Syndrome (PCOS) through detailed exploratory data analysis. Key findings include the significant variability in endometrial thickness, follicular distribution, menstrual cycle patterns, and body mass index (BMI) between PCOS and non-PCOS participants. These features underscore the complex interplay of endocrine, metabolic, and lifestyle factors in the manifestation of PCOS. The irregular menstrual cycles, elevated endometrial thickness, and increased follicular counts observed in PCOS participants reinforce the importance of early diagnostic tools, such as ultrasound imaging and hormonal profiling, in identifying and managing the condition. Furthermore, the study highlights the critical role of lifestyle factors, such as diet and physical activity, in exacerbating or mitigating PCOS symptoms, suggesting the potential benefits of personalized lifestyle interventions.

Elevated hemoglobin levels and higher BMI among PCOS patients suggest that these variables may serve as potential biomarkers for detecting underlying metabolic disturbances. These findings provide a foundation for future research into the development of predictive models and tailored treatment strategies for PCOS. While the insights gained from this analysis are valuable, there remains a need for further longitudinal studies and the integration of hormonal data to better understand the progression and management of PCOS. Future work should also leverage advanced machine learning techniques to refine diagnostic capabilities and improve patient outcomes. In conclusion, this study underscores the multifaceted nature of PCOS and the importance of a holistic approach to its diagnosis and management. By addressing the clinical and lifestyle factors identified in this research, healthcare providers can better support patients in managing the reproductive and metabolic challenges associated with PCOS, ultimately enhancing their quality of life.

## 7. REFERENCES

1. **Fauser, B. C. J. M., Tarlatzis, B. C., Rebar, R. W., Legro, R. S., Balen, A. H., Lobo, R., Carmina, E., Chang, J., Yildiz, B. O., & others.** (2012). Consensus on women's health aspects of polycystic ovary syndrome (PCOS): The Amsterdam ESHRE/ASRM-Sponsored 3rd PCOS Consensus Workshop Group. *Fertility and Sterility*, 97(1), 28–38.
2. **Dewailly, D., Lujan, M. E., Carmina, E., Cedars, M. I., Franks, S., Legro, R. S., Norman, R. J., & Escobar-Morreale, H. F.** (2014). Definition and significance of polycystic ovarian morphology: A task force report from the Androgen Excess and PCOS Society. *Human Reproduction Update*, 20(3), 334–352.
3. **Diamanti-Kandarakis, E., & Dunaif, A.** (2012). Insulin resistance and the polycystic ovary syndrome revisited: An update on mechanisms and implications. *Endocrine Reviews*, 33(6), 981–1030.
4. **Goodarzi, M. O., Dumesic, D. A., Chazenbalk, G., & Azziz, R.** (2011). Polycystic ovary syndrome: Etiology, pathogenesis and diagnosis. *Nature Reviews Endocrinology*, 7(4), 219–231.
5. **Rotterdam ESHRE/ASRM-Sponsored PCOS Consensus Workshop Group.** (2004). Revised 2003 consensus on diagnostic criteria and long-term health risks related to polycystic ovary syndrome. *Fertility and Sterility*, 81(1), 19–25.
6. **Escobar-Morreale, H. F., Luque-Ramírez, M., & San Millán, J. L.** (2005). The molecular–genetic basis of functional hyperandrogenism and the polycystic ovary syndrome. *Endocrine Reviews*, 26(2), 251–282.
7. **Legro, R. S., Kusanman, A. R., Dodson, W. C., & Dunaif, A.** (1999). Prevalence and predictors of risk for type 2 diabetes mellitus and impaired glucose tolerance in polycystic ovary syndrome: A prospective, controlled study in 254 affected women. *Journal of Clinical Endocrinology & Metabolism*, 84(1), 165–169.
8. **Hirschberg, A. L.** (2009). Polycystic ovary syndrome, obesity, and reproductive implications. *Women's Health*, 5(5), 529–540.
9. **Azziz, R., Woods, K. S., Reyna, R., Key, T. J., Knochenhauer, E. S., & Yildiz, B. O.** (2004). The prevalence and features of the polycystic ovary syndrome in an unselected population. *Journal of Clinical Endocrinology & Metabolism*, 89(6), 2745–2749.

10. **Teede, H. J., Deeks, A. A., & Moran, L. J.** (2010). Polycystic ovary syndrome: A complex condition with psychological, reproductive, and metabolic manifestations that impacts health across the lifespan. *BMC Medicine*, 8(1), 41.
11. **Zawadski, J. K., & Dunaif, A.** (1992). Diagnostic criteria for polycystic ovary syndrome: Towards a rational approach. *Polycystic Ovary Syndrome*, 9(1), 377–384.
12. **Carmina, E., Oberfield, S. E., Lobo, R. A.** (2010). The diagnosis of polycystic ovary syndrome in adolescents. *American Journal of Obstetrics and Gynecology*, 203(3), 201.e1–201.e5.
13. **Balen, A. H., Conway, G. S., Kaltsas, G., Techatraisak, K., Manning, P. J., West, C., & Jacobs, H. S.** (1995). Polycystic ovary syndrome: The spectrum of the disorder in 1741 patients. *Human Reproduction*, 10(8), 2107–2111.
14. **Homburg, R.** (2008). Androgen circle of polycystic ovary syndrome. *Human Reproduction*, 23(3), 643–646.
15. **Witchel, S. F., Oberfield, S. E., & Peña, A. S.** (2019). Polycystic ovary syndrome: Pathophysiology, presentation, and treatment with emphasis on adolescent girls. *Journal of the Endocrine Society*, 3(8), 1545–1573.
16. **Sanchón, R., Gambineri, A., Alpañés, M., Martínez-García, M. Á., Pasquali, R., & Escobar-Morreale, H. F.** (2012). Prevalence and phenotypes of polycystic ovary syndrome in Caucasian women according to the Rotterdam criteria. *Fertility and Sterility*, 98(1), 236–241.
17. **Ehrmann, D. A.** (2005). Polycystic ovary syndrome. *New England Journal of Medicine*, 352(12), 1223–1236.
18. **Moran, L. J., Misso, M. L., Wild, R. A., & Norman, R. J.** (2010). Impaired glucose tolerance, type 2 diabetes and metabolic syndrome in polycystic ovary syndrome: A systematic review and meta-analysis. *Human Reproduction Update*, 16(4), 347–363.
19. **Pasquali, R., Gambineri, A., & Pagotto, U.** (2006). The impact of obesity on reproduction in women with polycystic ovary syndrome. *BJOG: An International Journal of Obstetrics & Gynaecology*, 113(10), 1148–1159.
20. **Palomba, S., Santagni, S., Falbo, A., & La Sala, G. B.** (2015). Complications and challenges associated with polycystic ovary syndrome: Current perspectives. *International Journal of Women's Health*, 7(1), 745–763.