

2021

# DETECTION & CLASSIFICATION OF PCOD/PCOS USING MACHINE LEARNING ALGORITHMS

ATHULYA SHANTY  
2048030  
2MDS

CHRIST UNIVERSITY | Bangalore

# Detection and Classification of PCOD/ PCOS

## INTRODUCTION

The chosen domain is Healthcare in women. Ovaries are an important part in the reproductive system of females. The reproductive system of women is controlled by the complex interplay of primarily five reproductive hormones namely estrogen, gonadotropin-releasing hormone, follicle stimulating hormone, progesterone and luteinizing hormone. An imbalance within these hormones leads to a hormonal disorder called the polycystic ovary syndrome (PCOS) or polycystic ovarian disease (PCOD) among women of reproductive age. The signs and symptoms of this disease include anovulation, menstrual dysfunction and signs of hyperandrogenism. Other signs and symptoms include hirsutism, infertility, obesity, metabolic syndrome, and diabetes. Women with PCOS/PCOD have high chances of hypertension. PCOS/PCOD is clinically diagnosed with an ultrasound abdomen scanning.

My aim is to build a model using machine learning algorithms which provides at most accuracy to predict the presence of PCOD/PCOS in women. This project aims to detect the polycystic ovaries Disease/Syndrome in the ovaries of a woman in order to recommend an immediate attention and need of diagnosis through a set of questionnaires. The diagnosis of PCOS/PCOD is important as this disease can cause problems with menstrual periods and make it difficult for females to conceive. If not treated it can cause insulin resistant diabetes, obesity and high cholesterol leading to heart disease.

This project aims to give an alert or warning to the females with PCOS/PCOD at its earliest stage possible. Even though there are no certain medicines for this syndrome/disease, lifestyle modifications such as diet, exercise and weight loss are considered to be the first-line treatment for women with PCOS/PCOD.

## ABOUT THE DATA

The data was collected from a population of females living across the world. The source of data collection was primary and was carried out using Google forms which were distributed with the help of social media and messaging services. The google form contained 6 sections such as the main section, Personal details, Symptoms, Personal details & history, lifestyle and Medical history which contain sub questions that is curated to identify certain key criteria of PCOS/ PCOD which in turn help in developing a model to predict if the female has PCOD / PCOS or not. The collected data contains 882 rows and 27 observations.

## IMPLEMENTATION DETAILS

The raw dataset contained 882 observations and 27 variables. The lengthy column names were renamed to shorter ones. The duplicate values in the dataset were removed. Even though the desired population is women, the survey was also attended by some of the males whom must be removed for further analysis. Therefore there comes a need to filter the dataset such that the required dataset contains only the data from females. The same was done by removing the observations with 'other' and 'male' as gender. Also, there were some unwanted variables for our study and they were removed. There were missing data in the dataset. Dropping the missing values was not a solution as it dropped the entire dataset and those were handled by replacing the null values with the string 'unknown'. The column named 'kids' contained both numerical and string data which will result in errors in further processes. Therefore the numbers 0,1,2 were replaced to strings 'zero', 'one' and 'two'. After cleaning the data, the observations got reduced to 723 and the variables got reduced to 24.

The dataset had 511 observations which contained only independent variables whereas 212 observations contained both independent variables and dependent variable. These were separated into two different data frames for further analysis. There are no numerical variables found in the dataset whereas there were 24 categorical variables. These 24 categorical variables were converted to numerical variables by label encoding. The dataframe which contained both dependent and independent variables were divided into dependent variable and independent variables.

The dependent variable is 'diagnosis' and independent variables are 'age', 'body\_type', 'scanning', 'before\_period', 'irregular\_period', 'painful\_period', 'bleeding', 'period\_cycle', 'period\_duration', 'period\_pain', 'clots', 'alcohol', 'smoking', 'stress', 'exercise', 'Hereditary', 'diabetes', 'hypothyroidism', 'hair\_growth', 'acne', 'marital\_status', 'kids' and 'work'. Different bar plots were plotted with their frequency for unique category of the categorical variables for a better understanding of the dataset.

The correlation matrix gave the correlation between the variables and the heatmap of the same gave a better way of visualization. The relevant features for the study were chosen based on the correlation between variables. The variables with correlation more than 0.2 were considered to be the important variables.

The variable to be predicted was categorical and hence different algorithms such as Naive Bayes classification, Logistic regression, KNN classification, Random Forest classifier and Decision tree classifier were used for modelling the data. The data frame which have both the independent and

dependent variables were splitted to X\_train, y\_train, X\_test and y\_test. Each model was trained by inputting X\_train, y\_train values of the training dataset. Model evaluation and the confusion matrix were drawn for each algorithm.

The data frame which had only dependent variables were inputted to the two models with the highest accuracy for predicting the diagnosis. The predicted diagnosis by the models with highest accuracies were also compared.

## RESULTS AND DISCUSSION

- The lengthy column names was renamed to shorter ones.

```
1 #Reading the dataset
2 df=pd.read_excel(r'E:\files\Desktop\2MDS\MACHINE_LEARNING\projects\uncleaned_data.xlsx')
3 df.head()
```

Timestamp	What is your name?	Select your Gender	Pick your age limit	How would you describe your Body Physic ?	Have you done an ultrasound abdomen scanning and what does your report say?	Do you notice any of these right before your period begins?	Are you experiencing irregular or late periods?	Are you experiencing painful periods?	Are you experiencing Excessive bleeding?	...	Do you exercise regularly?	Is your mother diagnosed with PCOS/PCOD?	Do you suffer from diabetes?
-----------	--------------------	--------------------	---------------------	---	---	---	---	---------------------------------------	--	-----	----------------------------	--	------------------------------

```
1 df.rename(renamed, axis='columns', inplace=True)
2 df.head()
```

Timestamp	name	gender	age	body_type	scanning	before_period	irregular_period	painful_period	bleeding	...	exercise	Hereditary	diabetes	hypo
-----------	------	--------	-----	-----------	----------	---------------	------------------	----------------	----------	-----	----------	------------	----------	------

- Removing duplicates removed 113 duplicate records and resulted in 763 observations with 27 variables.

```
1 df.shape
(882, 27)
```

```
1 #Removing duplicates
2 df=df.drop_duplicates(subset=['name'])
3 df.shape
(763, 27)
```

- Dropping 'male' and 'other' from gender resulted in 723 observations with 27 variables.

```
1 #Removing males and other categories from gender variable.
2 df.drop(df[df['gender'] == 'Male'].index, inplace = True)
3 df.drop(df[df['gender'] == 'Other'].index, inplace = True)
```

```
1 df.shape
(723, 27)
```

- Removal of unwanted columns such as 'Timestamp', 'name', etc resulted in 723 observations with 24 variables.

```

1 df=df.drop(['Timestamp','name','gender'],axis=1)
2 df.shape

```

(723, 24)

- Handling missing data was done by filling the null values with the variable 'unknown'.

```

1 #Filling missing values with 'unknown'
2 df['marital_status']=df['marital_status'].fillna('unknown')
3 df['work']=df['work'].fillna('unknown')
4 df['kids']=df['kids'].fillna('unknown')
5 df['diagnosis']=df['diagnosis'].fillna('Didn't check')

```

- The dataset was splitted into 2 data frames. One for training and the other for testing. The observations which had both dependent and independent variables were added to 1<sup>st</sup> data frame for testing and the observations with only independent variables were added to 2<sup>nd</sup> data frame. The observations in which the response of 'Didn't check' has been added to data frame 1 for prediction.

```

1 df2=df1=df

```

```

1 #Data for prediction
2 df1=df1[df1['diagnosis'] == 'Didn't check']
3 #Dropping diagnosis column of predictive dataset as it contains only 'Didn't check'
4 df1=df1.drop(['diagnosis'],axis=1)
5 df1.shape

```

(511, 23)

```

1 #Data for modeling
2 df2.drop(df2[df2['diagnosis'] == 'Didn't check'].index, inplace = True)
3 df2.shape

```

(212, 24)

- The dataset contained 24 categorical variables and no numerical variables.

```

1 # List of numerical variables
2 numerical = [feature for feature in df.columns if ((df[feature].dtypes != 'O') & (feature not in ['deposit']))]
3 print('Number of numerical variables: ', len(numerical))
4
5 print(numerical)

```

Number of numerical variables: 0  
[]

```

1 #List of categorical variables
2 categorical=[]
3 for col in df.select_dtypes(include='object').columns:
4     categorical.append(col)
5 print('Number of categorical variables: ', len(categorical))
6 print(categorical)

```

Number of categorical variables: 24  
['age', 'body\_type', 'scanning', 'before\_period', 'irregular\_period', 'painful\_period', 'bleeding', 'period\_cycle', 'period\_duration', 'period\_pain', 'clots', 'alcohol', 'smoking', 'stress', 'exercise', 'Hereditary', 'diabetes', 'hypothyroidism', 'hair\_growth', 'acne', 'marital\_status', 'kids', 'work', 'diagnosis']

- The categorical variables of both the data frames have been converted to numerical by Label Encoding.

```

1 #Label Encoding for categorical variables in df1
2 from sklearn import preprocessing
3 for i in range(0,(len(categorical)-1)):
4     le=preprocessing.LabelEncoder()
5     le.fit(df1[categorical[i]])
6     df1[categorical[i]]=le.transform(df1[categorical[i]])

```

```

1 #Label Encoding for categorical variables in df2
2 from sklearn import preprocessing
3 for i in range(0,len(categorical)):
4     le=preprocessing.LabelEncoder()
5     le.fit(df2[categorical[i]])
6     df2[categorical[i]]=le.transform(df2[categorical[i]])

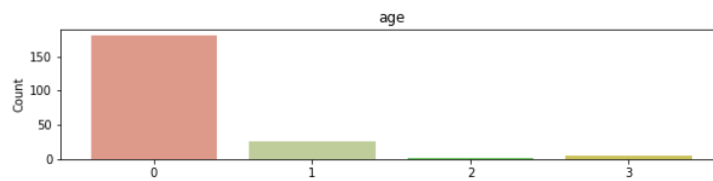
```

- The count of each value of every variables were plotted in a bar graph.

```

1 colors = ['#DD998A', '#BEC999', '#CAC158', '#58BF4B']
2 import matplotlib.pyplot as plt
3 def bar_plot(variable):
4     var = df[variable]
5     varValue = var.value_counts()
6     plt.figure(figsize = (10,2))
7     plt.bar(varValue.index,varValue,color=colors)
8     plt.xticks(varValue.index,varValue.index.values)
9     plt.ylabel("Count")
10    plt.title(variable)
11    plt.show()
12    print("{}: \n {}".format(variable,varValue))
13
14 for c in categorical:
15     bar_plot(c)

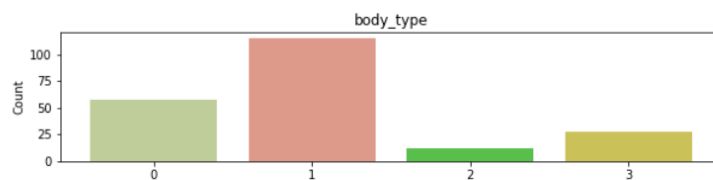
```



```

age:
0    180
1     26
3      4
2       2
Name: age, dtype: int64

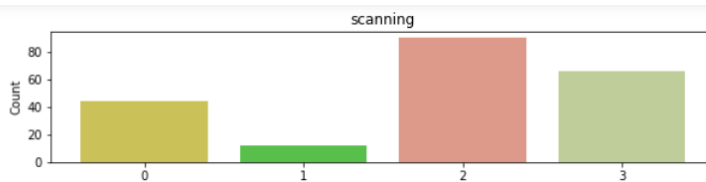
```



```

body_type:
1    115
0     58
3     27
2      12
Name: body_type, dtype: int64

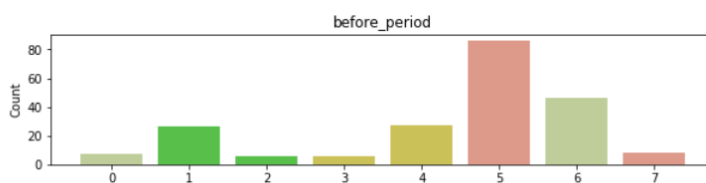
```



```

scanning:
2     90
3     66
0     44
1     12
Name: scanning, dtype: int64

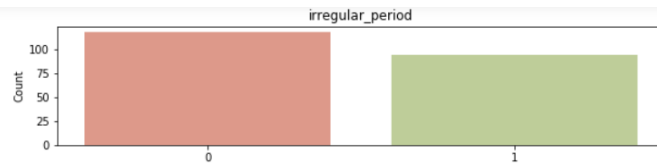
```



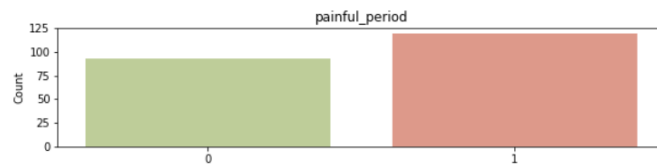
```

before_period:
5     86
6     46
4     27
1     26
7      8
0       7
3       6
2       6
Name: before_period, dtype: int64

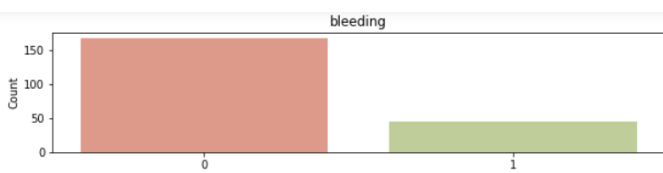
```



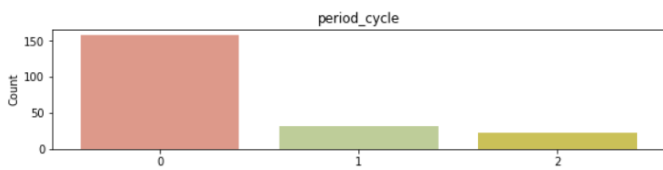
```
irregular_period:
0    118
1     94
Name: irregular_period, dtype: int64
```



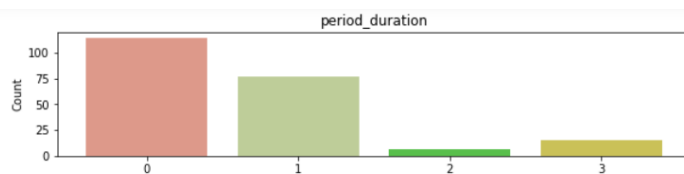
```
painful_period:
1    119
0     93
Name: painful_period, dtype: int64
```



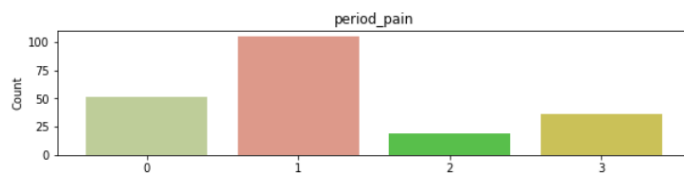
```
bleeding:
0    167
1     45
Name: bleeding, dtype: int64
```



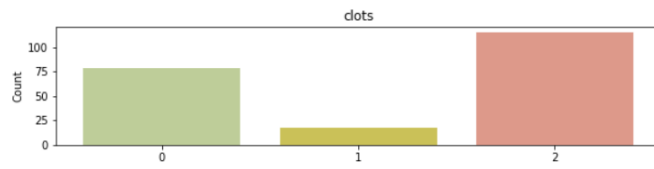
```
period_cycle:
0    158
1     31
2     23
Name: period_cycle, dtype: int64
```



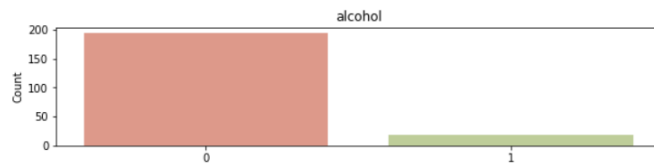
```
period_duration:
0    114
1     77
3     15
2      6
Name: period_duration, dtype: int64
```



```
period_pain:
1    105
0     52
3     36
2     19
Name: period_pain, dtype: int64
```



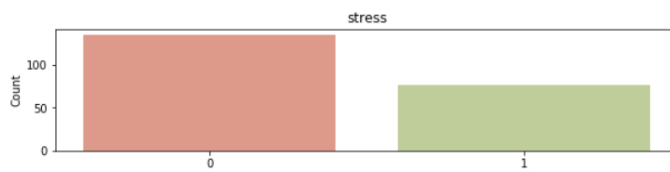
```
clots:
  2    115
  0     79
  1     18
Name: clots, dtype: int64
```



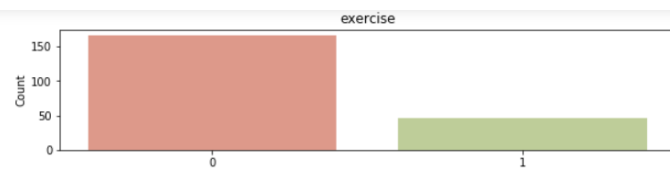
```
alcohol:
  0    194
  1     18
Name: alcohol, dtype: int64
```



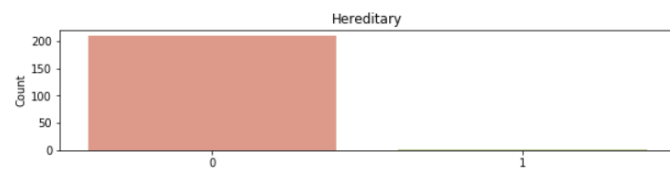
```
smoking:
  0    206
  1      6
Name: smoking, dtype: int64
```



```
stress:
  0    135
  1     77
Name: stress, dtype: int64
```



```
exercise:
  0    166
  1     46
Name: exercise, dtype: int64
```

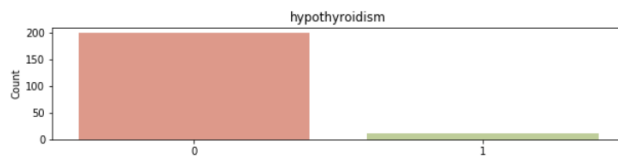


```
Hereditary:
  0    210
  1      2
Name: Hereditary, dtype: int64
```

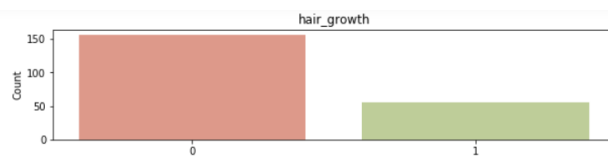




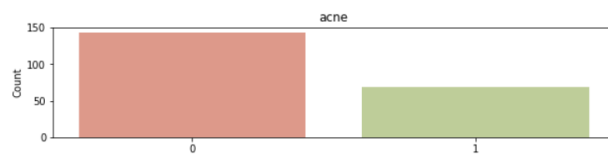
```
diabetes:
0    211
1         1
Name: diabetes, dtype: int64
```



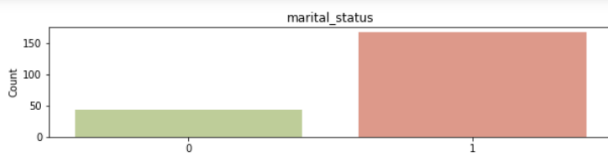
```
hypothyroidism:
0    200
1     12
Name: hypothyroidism, dtype: int64
```



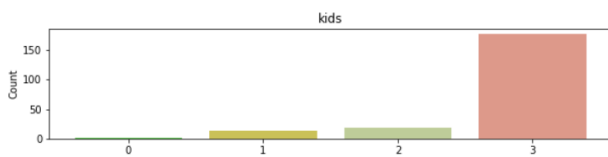
```
hair_growth:
0    156
1     56
Name: hair_growth, dtype: int64
```



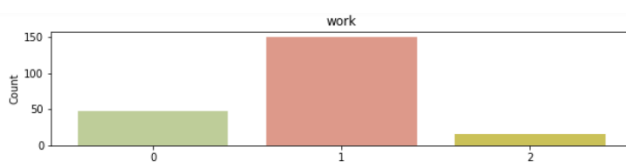
```
acne:
0    143
1     69
Name: acne, dtype: int64
```



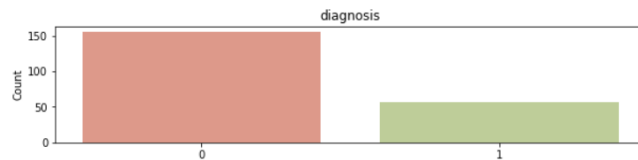
```
marital_status:
1    168
0     44
Name: marital_status, dtype: int64
```



```
kids:
3    177
2     19
1     14
0         2
Name: kids, dtype: int64
```



```
work:
1    150
0     47
2      15
Name: work, dtype: int64
```

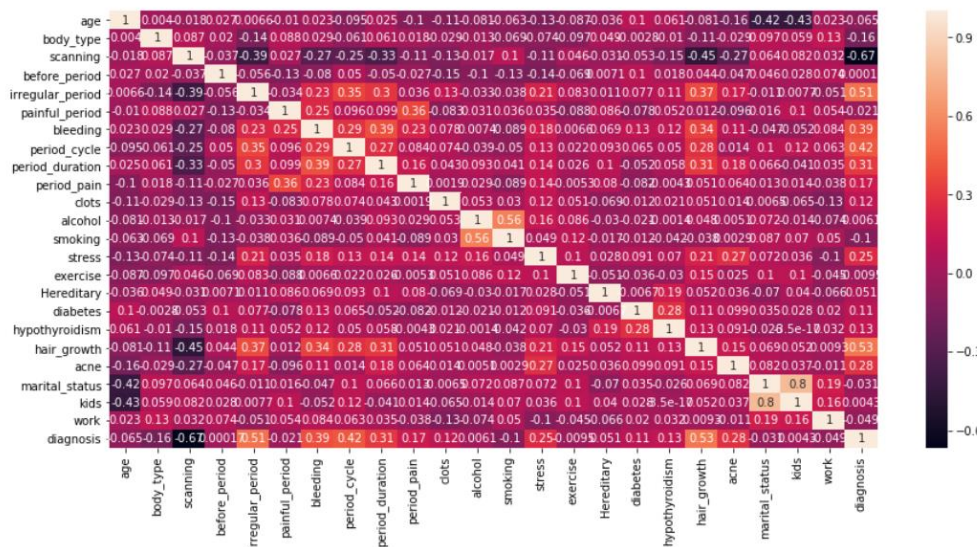


```
diagnosis:
0    155
1     57
Name: diagnosis, dtype: int64
```

- Correlation heat map has been plotted.

```
In [32]: 1 ## Checking for correlation
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 cor_mat=df2.corr()
5 fig = plt.figure(figsize=(15,7))
6 sns.heatmap(cor_mat,annot=True)
```

Out[32]: <matplotlib.axes.\_subplots.AxesSubplot at 0x28fb04da400>



- The relevant features was found to be 'scanning', 'irregular\_period', 'bleeding', 'period\_cycle', 'stress', 'hair\_growth' and 'acne' through the method of correlation with cut-off as 0.2.

```
1 #Correlation with output variable
2 cor=df2.corr()
3 cor_target = abs(cor['diagnosis'])
4 #Selecting highly correlated features
5 relevant_features = cor_target[cor_target>0.2]
6 print('There are ',(len(relevant_features)-1),'relevant features out of which diagnosis is the target variable')
7 relevant_features[0:(len(relevant_features)-1)]
8
```

There are 8 relevant features out of which diagnosis is the target variable

```
scanning          0.666752
irregular_period  0.508112
bleeding          0.387688
period_cycle      0.417865
period_duration   0.314478
stress            0.249914
hair_growth       0.529517
acne              0.282646
Name: diagnosis, dtype: float64
```

```
1 relevant_features=['scanning','irregular_period','bleeding','period_cycle','stress','hair_growth','acne']
2 print(relevant_features)
```

['scanning', 'irregular\_period', 'bleeding', 'period\_cycle', 'stress', 'hair\_growth', 'acne']

- The dataset which have both dependent and independent variable was splitted into Xtrain, ytrain, Xtest and ytest for modelling.

```
1 X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=0)
```

```
1 X_train.shape
```

(169, 23)

```
1 X_test.shape
```

(43, 23)

- There were a total of 5 models. They are Naïve Bayesian model, Logistic Regression model, KNN model, Random Forest Classification model, and Decision Tree model with an accuracy of 88.37%, 90.70%, 90.70%, 88.37% and 90.7% respectively. Confusion matrix for each matrix was plotted.

```
1 #Naive Bayesian Classifier
2 clf=GaussianNB()
3 model=clf.fit(X_train,y_train)
4 y_pred=model.predict(X_test)
5 ac=accuracy_score(y_test,y_pred,normalize=True)
6 a=np.round(ac*100,2)
7 print('Accuracy is ',a,'%')
```

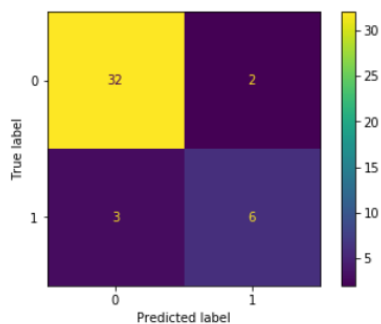
Accuracy is 88.37 %

```
1 #Model Evaluation of Naive Bayesian Classifier
2 print(f"Score in Test Data : {model.score(X_test,y_test)}")
3
4 cm=confusion_matrix(y_test, y_pred)
5 p_right=cm[0][0]+cm[1][1]
6 p_wrong=cm[0][1]+cm[1][0]
7
8 print(f"Right classification : {p_right}")
9 print(f"Wrong classification : {p_wrong}")
10
11 plot_confusion_matrix(clf, X_test, y_test)
12 plt.show()
```

Score in Test Data : 0.8837209302325582

Right classification : 38

Wrong classification : 5



```
1 #Logistic Regression
2 clf=LogisticRegression()
3 model=clf.fit(X_train,y_train)
4 y_pred=model.predict(X_test)
5 ac=accuracy_score(y_test,y_pred,normalize=True)
6 b=np.round(ac*100,2)
7 print('Accuracy is ',b,'%')
```

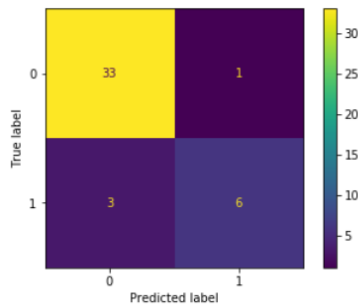
Accuracy is 90.7 %

```

1 #Model Evaluation of Logistic Regression
2 print(f"Score in Test Data : {model.score(X_test,y_test)}")
3
4 cm=confusion_matrix(y_test, y_pred)
5 p_right=cm[0][0]+cm[1][1]
6 p_wrong=cm[0][1]+cm[1][0]
7
8 print(f"Right classification : {p_right}")
9 print(f"Wrong classification : {p_wrong}")
10
11 plot_confusion_matrix(clf, X_test, y_test)
12 plt.show()

```

Score in Test Data : 0.9069767441860465  
 Right classification : 39  
 Wrong classification : 4



```

1 #KNN
2 clf=neighbors.KNeighborsClassifier(n_neighbors=3)
3 model=clf.fit(X_train,y_train)
4 y_pred=model.predict(X_test)
5 ac=accuracy_score(y_test,y_pred,normalize=True)
6 c=np.round(ac*100,2)
7 print('Accuracy is ',c,'%')

```

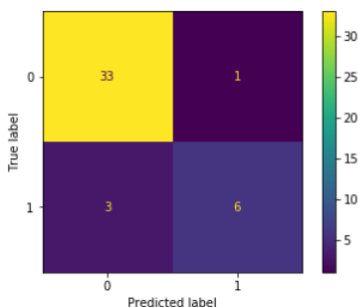
Accuracy is 90.7 %

```

1 #Model Evaluation of KNN
2 print(f"Score in Test Data : {model.score(X_test,y_test)}")
3
4 cm=confusion_matrix(y_test, y_pred)
5 p_right=cm[0][0]+cm[1][1]
6 p_wrong=cm[0][1]+cm[1][0]
7
8 print(f"Right classification : {p_right}")
9 print(f"Wrong classification : {p_wrong}")
10
11 plot_confusion_matrix(clf, X_test, y_test)
12 plt.show()

```

Score in Test Data : 0.9069767441860465  
 Right classification : 39  
 Wrong classification : 4



```

1 #Random Forest Classifier
2 clf = RandomForestClassifier(max_depth=2, random_state=0)
3 model=clf.fit(X_train, y_train)
4 y_pred=model.predict(X_test)
5 ac=accuracy_score(y_test,y_pred,normalize=True)
6 d=np.round(ac*100,2)
7 print('Accuracy is ',d,'%')

```

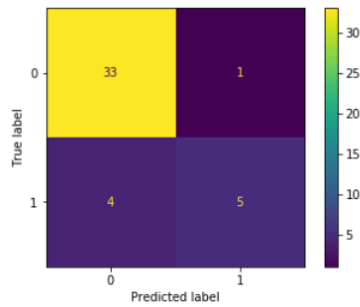
Accuracy is 88.37 %

```

1 #Model Evaluation of Random Forest Classifier
2 print(f"Score in Test Data : {model.score(X_test,y_test)}")
3
4 cm=confusion_matrix(y_test, y_pred)
5 p_right=cm[0][0]+cm[1][1]
6 p_wrong=cm[0][1]+cm[1][0]
7
8 print(f"Right classification : {p_right}")
9 print(f"Wrong classification : {p_wrong}")
10
11 plot_confusion_matrix(clf, X_test, y_test)
12 plt.show()

```

Score in Test Data : 0.8837209302325582  
 Right classification : 38  
 Wrong classification : 5



```

1 #Decision Tree
2 clf = tree.DecisionTreeClassifier(max_depth=4)
3 model = clf.fit(X_train, y_train)
4 y_pred=model.predict(X_test)
5 ac=accuracy_score(y_test,y_pred,normalize=True)
6 e=np.round(ac*100,2)
7 print('Accuracy is ',e,'%')

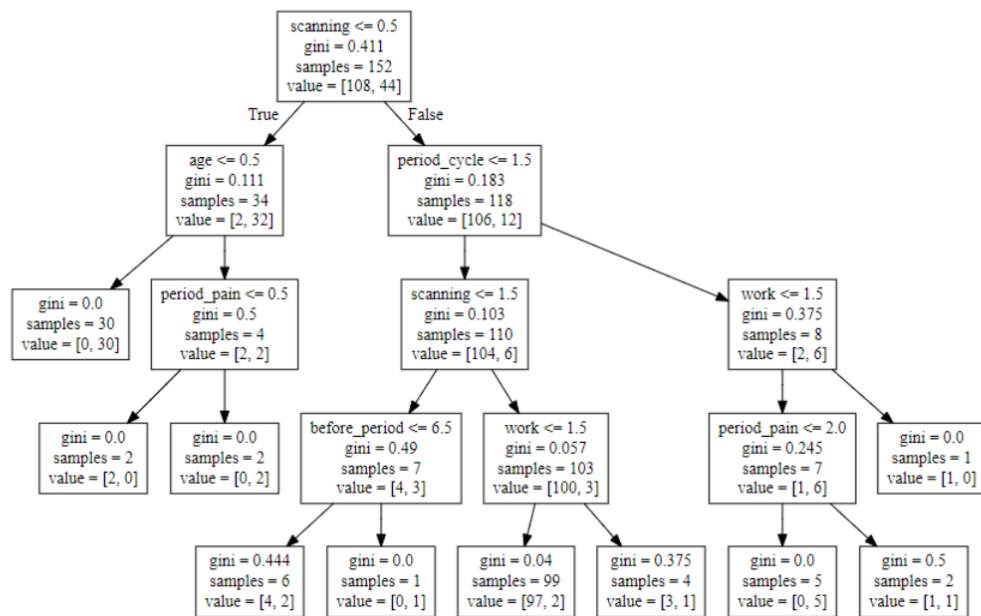
```

Accuracy is 88.37 %

```

1 import graphviz
2 with open("DecisionTree.dot", 'w') as f:
3     f=tree.export_graphviz(model, feature_names=X.columns, out_file=f);

```

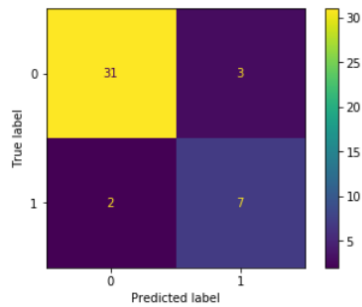


```

1 #Model Evaluation of Decision Tree
2 print(f"Score in Test Data : {model.score(X_test,y_test)}")
3
4 cm=confusion_matrix(y_test, y_pred)
5 p_right=cm[0][0]+cm[1][1]
6 p_wrong=cm[0][1]+cm[1][0]
7
8 print(f"Right classification : {p_right}")
9 print(f"Wrong classification : {p_wrong}")
10
11 plot_confusion_matrix(clf, X_test, y_test)
12 plt.show()

```

Score in Test Data : 0.8837209302325582  
 Right classification : 38  
 Wrong classification : 5



- The best models were found to be Logistic Regression and KNN with 90.70% of accuracy.

```

1 dictionary={'Algorithm':['Naive Bayesian Classifier','Logistic regression','KNN','Random Forest Classifier','Decision Tree Classifier'],
2             'Accuracy':[a,b,c,d,e]}
3 n=pd.DataFrame(dictionary)
4 n.sort_values(by=['Accuracy'],ascending=False)

```

	Accuracy	Algorithm
1	90.70	Logistic regression
2	90.70	KNN
0	88.37	Naive Bayesian Classifier
3	88.37	Random Forest Classifier
4	88.37	Decision Tree Classifier

- The diagnosis were predicted for the test data using the models with high accuracy.

```

In [56]: 1 #Prediction using Logistic Regression
2 clf=LogisticRegression()
3 model=clf.fit(X,y)
4 y_pred=model.predict(X2)
5 df3=df1
6 df4=pd.DataFrame(y_pred)
7 pred_diag={0:'Predicted Diagnosis'}
8 df4.rename(pred_diag, axis=1, inplace=True)
9 df4.head(10)

```

Out[56]:

	Predicted Diagnosis
0	0
1	0
2	0
3	0
4	0
5	0
6	0
7	0
8	1
9	1

```

1 #Prediction using KNN
2 clf=neighbors.KNeighborsClassifier(n_neighbors=3)
3 model=clf.fit(X, y)
4 y_pred=model.predict(X2)
5 df6=df1
6 df7=pd.DataFrame(y_pred)
7 pred_diag={0:'Predicted Diagnosis'}
8 df7.rename(pred_diag, axis=1, inplace=True)
9 df7.head(10)

```

Predicted Diagnosis	
0	0
1	0
2	0
3	0
4	0
5	0
6	0
7	0
8	1
9	0

- The diagnosis predicted by both the Logistic Regression model and KNN model were compared. 485 predictions were predicted the same by both the models and 26 predictions were predicted differently by both the models.

```

1 df9=df4==df7
2 df9['Predicted Diagnosis'].value_counts()

```

True 485  
False 26  
Name: Predicted Diagnosis, dtype: int64

## LIMITATIONS

- A minimum of 20 cysts in the ovaries were considered to be PCOD/PCOS. But we were not able to collect the number of cysts in each ovary which would have made the model much more accurate.
- The training data set was comparatively very less to the testing dataset. Reaching out to the ones with the disease should have been increased for a better accuracy and prediction.

## CONCLUSION

Hormone imbalances affect a woman's health in many different ways. PCOS/PCOD can increase the risk of infertility, metabolic syndrome, sleep apnea, endometrial cancer, and depression. A healthy lifestyle, eating a healthy diet and exercising regularly will help the women with PCOS/PCOD. A model with 90.70% of accuracy has been built for the prediction of PCOD/PCOS. The models Logistic Regression and KNN performed well than Naïve Bayesian, Random Forest, and Decision tree. This models with high accuracy helps in the diagnosis of the syndrome/disease at it earliest stage

possible and would be able to warn them beforehand. 485 predictions were same by both models whereas 26 predictions were not the same.

## **FUTURE WORK**

- I am planning to further continue this project by collecting more data which would increase the accuracy rate and improve the prediction.
- I am planning to perform unsupervised learning algorithm on the data frame where dependent variable is absent. I also wanted to compare the predictions of both supervised and unsupervised learning algorithms.
- I am planning to develop a mobile application or web page which will take all the variables as inputs and give the predictions as output with a warning message. The app will also display the methods to be performed for recovering from PCOS/PCOD.

## **BIBLIOGRAPHY & REFERENCES**

<https://www.healthline.com/health/polycystic-ovary-disease#health-effects>

<https://www.yashodahospitals.com/diseases-treatments/pcod-pcos-symptoms-causes-treatment/>