

Detection and Analysis of

Polycystic Ovary Syndrome

Data Science Boot Camp
Final Presentation

Ethan Bootehsaz, Kaiyuan Ma, Joy Wang

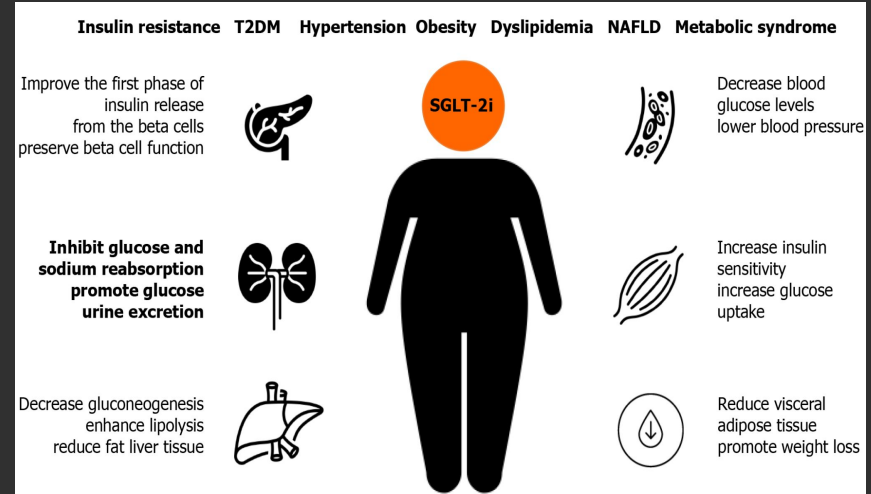
04/30/2024

Agenda

- **Goal**
- **Data Preprocessing Steps**
- **Exploration of Data**
 - More Dataset Features
- **Initial Hypothesis**
- **Model Selection**
 - Logistic Regression
 - Deep Learning
 - Deep learning with Cross Validation
- **Results & Discussion**
 - Model Performance
 - Future Improvements

What is PCOS?

PCOS affects 7–10% of women of and is the most common cause of infertility. In the United States, an estimated 5 to 6 million women have PCOS.



- Condition where ovaries produce an excessive amount of male hormones (androgens)
- May or may not develop cysts due to problems in menstrual cycle (anovulation)

Goal

Our goal was to develop a predictive model for the detection of Polycystic ovary syndrome (PCOS) by utilizing relevant physical and clinical parameters.

- We analyzed the dataset to identify patterns and trends.
- We developed a model that can predict the likelihood of PCOS.
- We evaluated the model's performance using various metrics.
- We created visualizations to support our findings.

Data Preprocessing Steps

- **Initial Data Cleaning**
 - Removed irrelevant columns: `Unnamed: 44`, `Sl. No`, and `Patient File No.` to focus on analytically useful data.
 - Renamed mislabeled columns to ensure consistency: Changed `Marraige Status (Yrs)` to `Marriage Status (Yrs)`.
- **Data Formatting**
 - Standardized column names by stripping leading/trailing whitespaces.
 - Cleaned data entries to remove trailing commas and periods for accuracy in analysis.
- **Handling Data Quality Issues**
 - Identified and removed rows where data had non-numeric characters in numeric fields to maintain data integrity.
- **Dataset Post-Processing**
 - Final dataset comprised of 539 rows ready for further analysis and modeling.

```
# Check for missing values in the dataset
```

```
missing_data = pcOS_data.isnull().sum()
```

```
missing_data[missing_data > 0]
```

```
Marraige Status (Yrs)      1
Fast food (Y/N)            1
Unnamed: 44                539
dtype: int64
```

```
Missing values in ../data/PCOS_processed.csv:
PCOS (Y/N)      0
Age (yrs)       0
Weight (Kg)     0
Height(Cm)      0
BMI             0
Blood Group     0
Pulse rate(bpm) 0
RR (breaths/min) 0
Hb(g/dl)        0
Cycle(R/I)      0
Cycle length(days) 0
Marriage Status (Yrs) 0
Pregnant(Y/N)   0
No. of abortions 0
I beta-HCG(mIU/mL) 0
II beta-HCG(mIU/mL) 0
FSH(mIU/mL)     0
LH(mIU/mL)      0
FSH/LH          0
Hip(inch)       0
Waist(inch)     0
Waist:Hip Ratio 0
TSH (mIU/L)     0
AMH(ng/mL)      0
...
```

```
#Dropping repeated/unnecessary columns
```

```
df = df.drop(['Unnamed: 44', 'Sl. No', 'Patient File No.'], axis=1, errors='ignore')
```

```
#Renaming column due to misspelling in original df
```

```
df.rename(columns={'Marraige Status (Yrs)': 'Marriage Status (Yrs)'}, inplace=True, errors='ignore')
```

```
# Fix column names - optional
```

```
df.columns = df.columns.str.strip() # .str.replace(' ', '_').str.lower()
```

```
df.columns = [re.sub(r'\s+', ' ', col).strip() for col in df.columns]
```

```
# Remove leading/trailing whitespaces from the data, remove trailing commas, periods
```

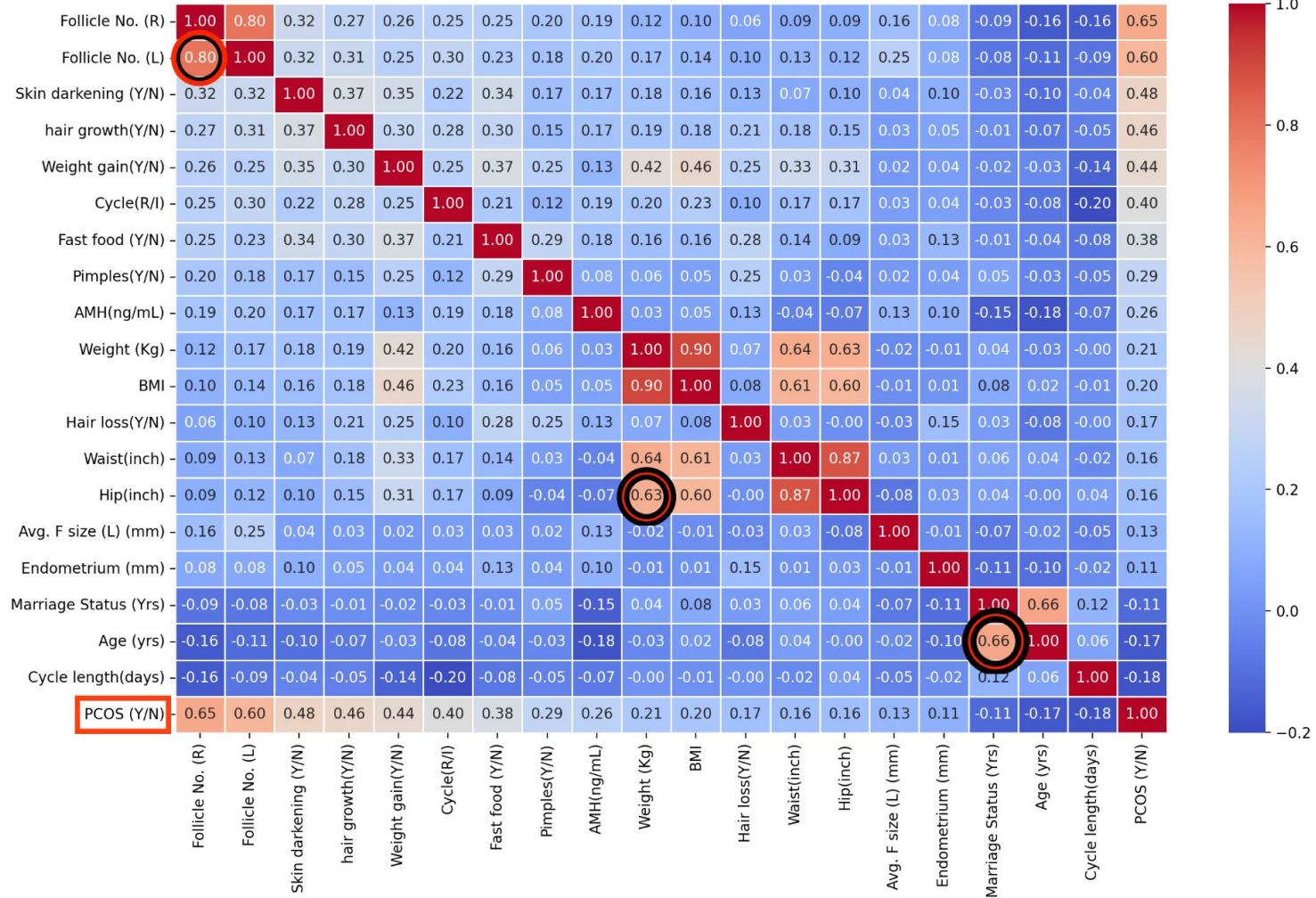
```
df = df.applymap(lambda x: x.strip() if isinstance(x, str) else x)
```

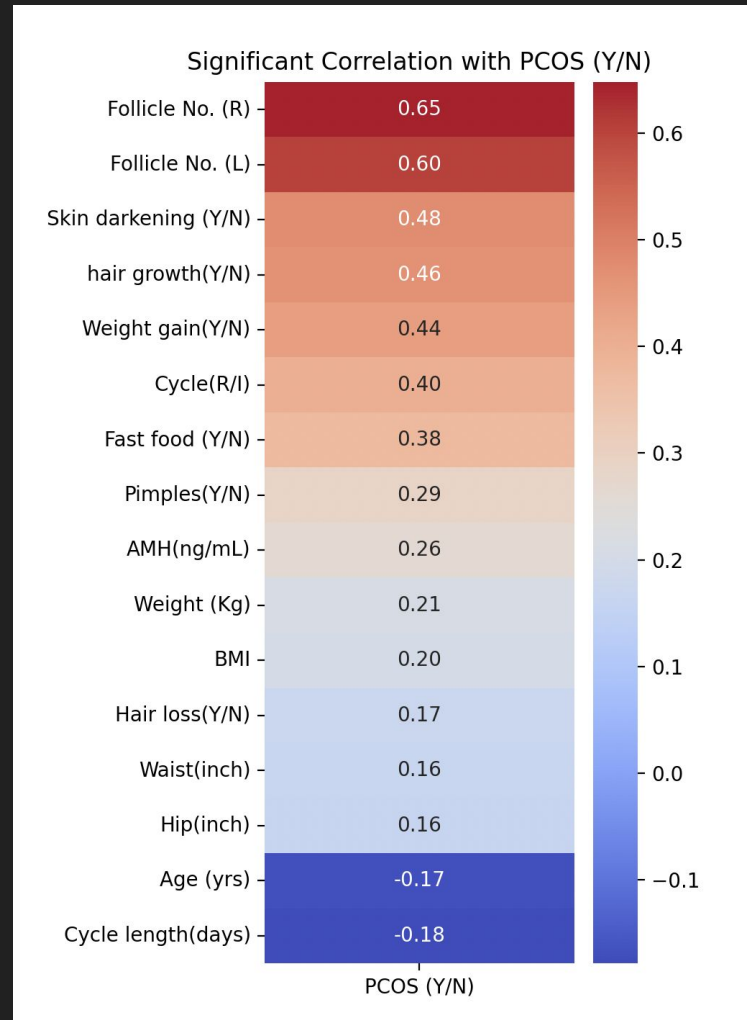
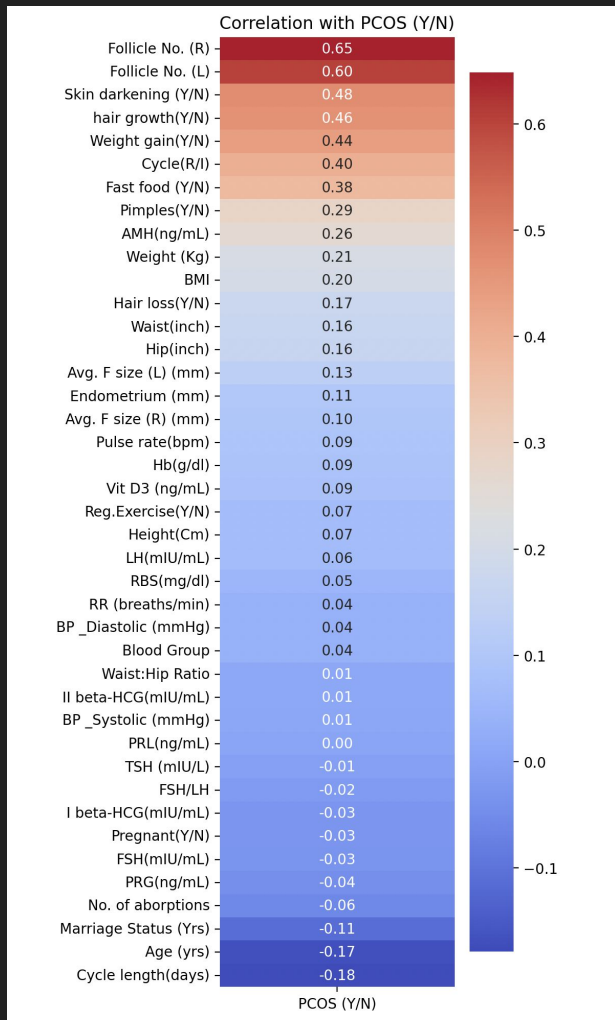
```
df = df.applymap(lambda x: x.rstrip('.') if isinstance(x, str) else x)
```

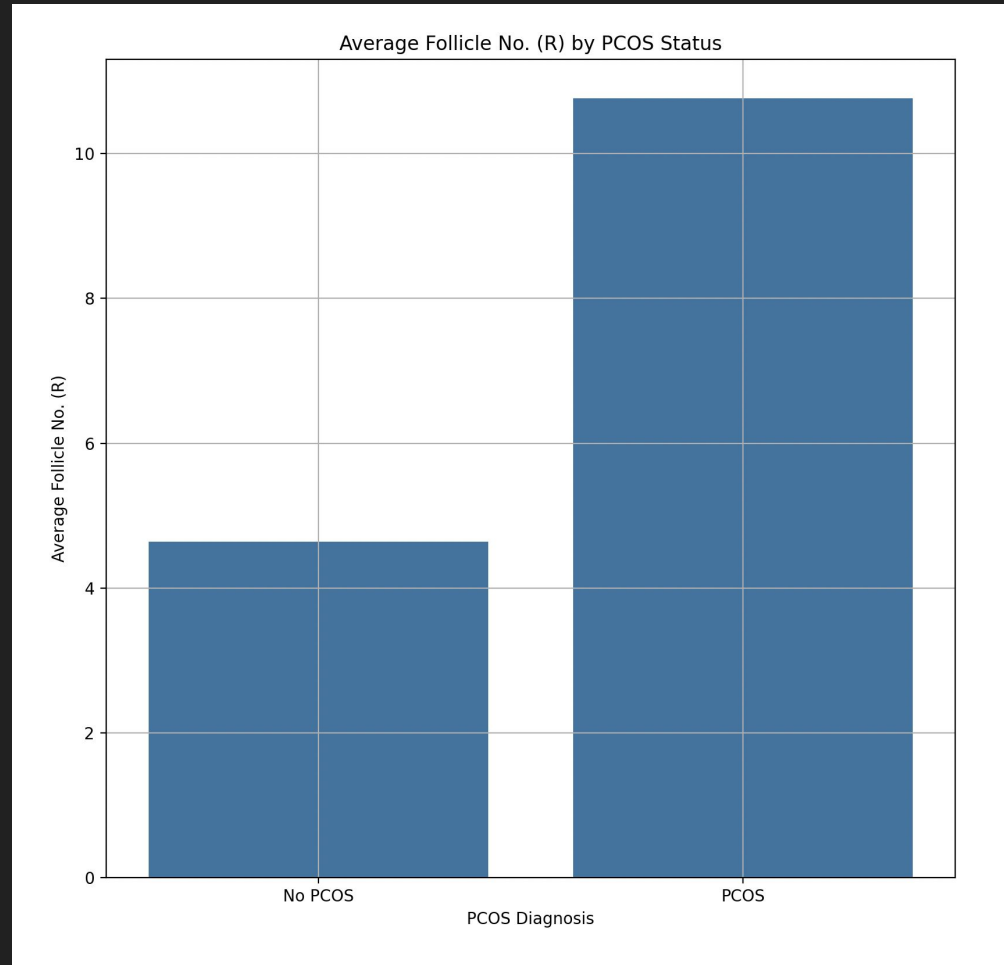
Exploration of Data

- Understanding Correlation Calculation in Data Analysis
 - `df.corr()['PCOS (Y/N)']`
- Correlation coefficients range from -1 to 1. A coefficient close to 1 indicates a strong positive correlation, meaning as one variable increases, so does the other. A coefficient close to -1 indicates a strong negative correlation, where one variable increases as the other decreases. A coefficient near 0 suggests no linear relationship.

Heatmap of Significantly Correlated Features with PCOS (Y/N)

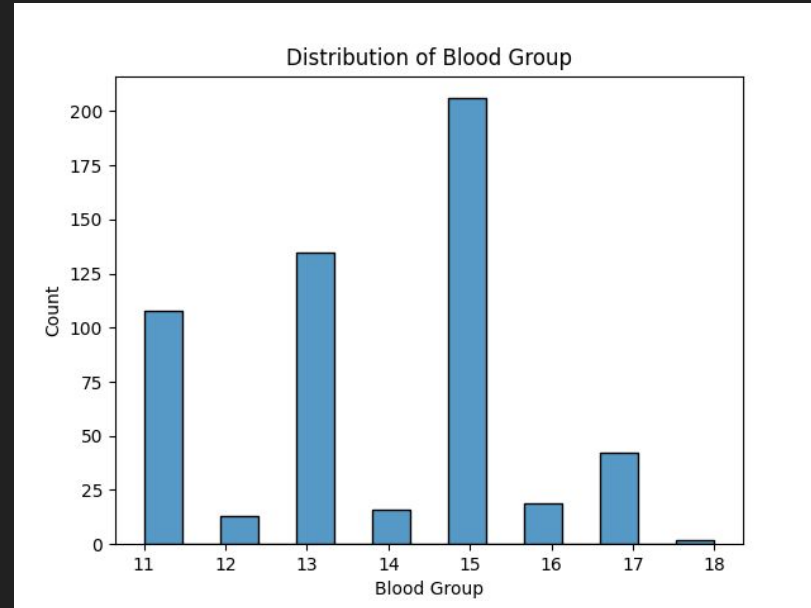
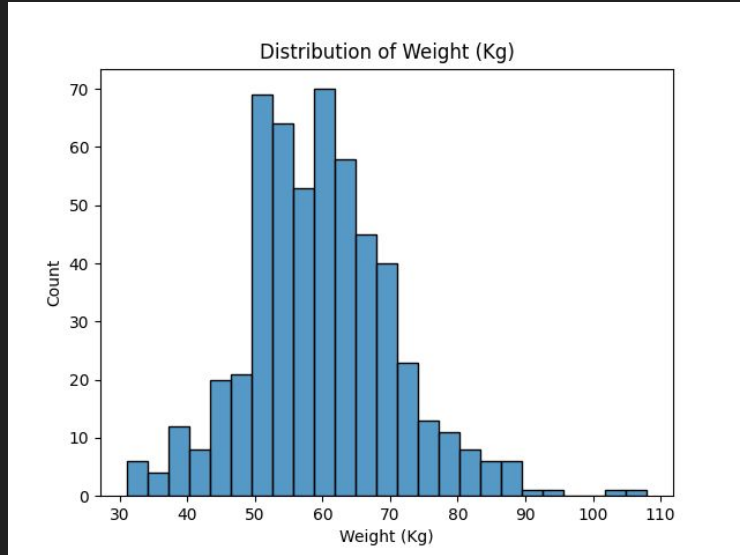






More Dataset Features

Distribution of numerical variables



mean

std

min

max

156.484835

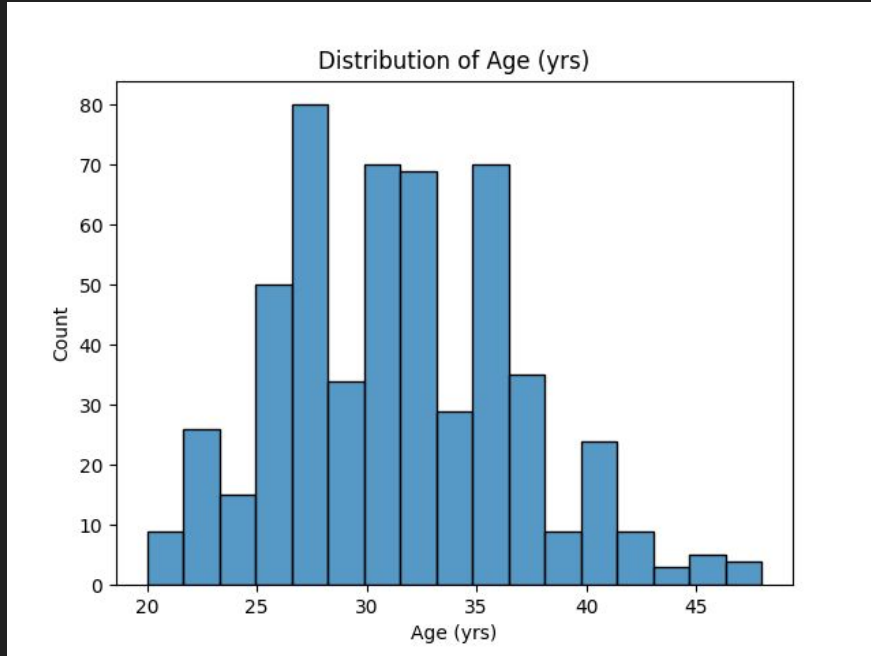
6.033545

137.000000

108.000000

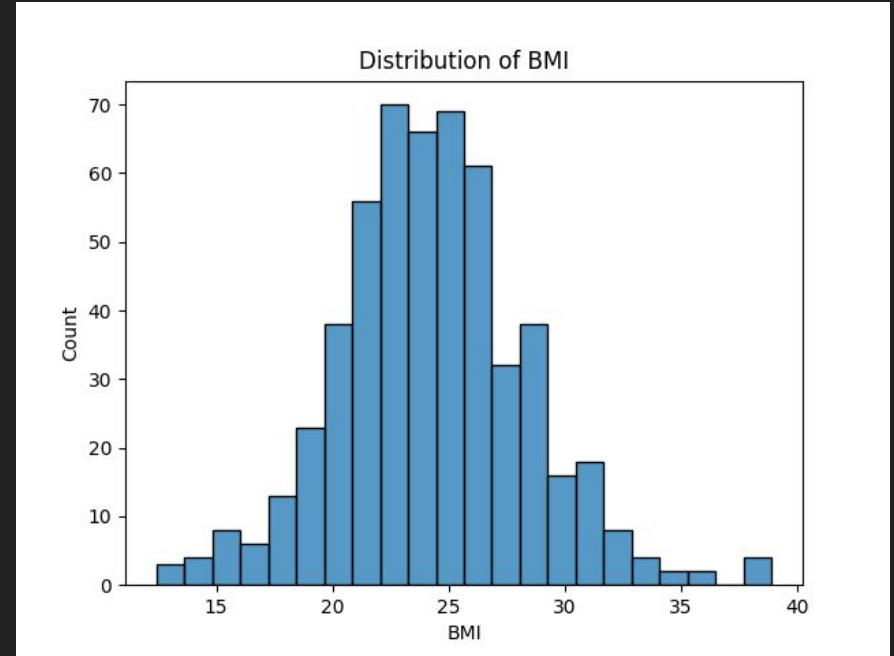
Blood Group of the patient A+ = 11, A- = 12, B+ = 13, B- = 14, O+ = 15, O- = 16, AB+ = 17, AB- = 18

Distribution of numerical variables cont.



Count mean std min max

541.0 31.430684 5.411006 20.000000 48.000000



Count mean std min max

541.0 24.311285 4.056399 12.417882 38.900000

Contextualize Dataset Info

The average patient is/has:

- 156.48 lbs
- 31.43 years old
- BMI of 24.311
- Been married 5 years
- More likely to be blood type (O+)

Cycle (R/I) + Cycle length (days)

Cycle R/I

There's no documentation on what the 2, 4, 5 mean >:(((

My guess is that 2 is regular, 4 is irregular, 5 is none but this is unusable

Cycle Length

Varies, but the average is 28 days

beta-HCG II (mIU/ML)

Used in fertility / pregnancy testing. Can also determine the age of the fetus based on gestation period.

FSH → Follicle-stimulating hormone

Levels vary based on puberty, menstruation, menopause

Waist/Hip ratio

Used with BMI to determine health risks

Health risk: obesity, type II diabetes, heart disease

- Low = 0.80 or below
- Moderate = 0.81 - 0.85
- High = 0.86 or higher

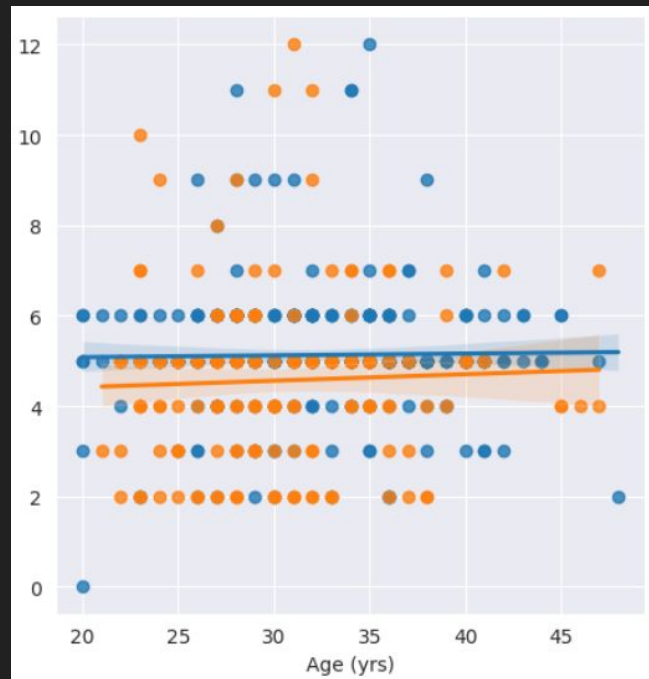
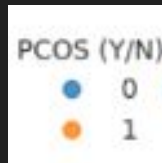
Proportion of diagnoses

	Age Group	PCOS (Y/N)	Proportion
0	(20, 25]	1	0.564516
1	(20, 25]	0	0.435484
2	(25, 30]	0	0.633880
3	(25, 30]	1	0.366120
4	(30, 35]	0	0.702381
5	(30, 35]	1	0.297619
6	(35, 40]	0	0.804348
7	(35, 40]	1	0.195652
8	(40, 45]	0	0.840000
9	(40, 45]	1	0.160000
10	(45, 50]	1	0.600000
11	(45, 50]	0	0.400000

```
# make a linear model plot for the menstrual length, PCOS vs NORMAL  
ax = sns.lmplot(data = df, x = "Age (yrs)", y = "Cycle length(days)", hue = "PCOS (Y/N)")  
plt.show(ax)
```

Y axis:

Cycle length



Higher incidence rate in younger & older pop.

Initial Hypothesis

- Age (yrs): Age can influence hormonal balance, which is crucial in PCOS.
- Weight (Kg) and BMI: There's a known correlation between body weight, fat distribution, and PCOS.
- Pulse rate (bpm): Although not a direct indicator, it might provide insights into metabolic rate and overall health.
- Follicle Number (L/R) and Average Follicle Size: Directly related to ovarian function, which is central to PCOS diagnosis.
- Endometrium (mm): Endometrial thickness can be affected by PCOS due to hormonal imbalances.
- Fast food (Y/N): Diet and lifestyle choices are relevant to symptoms and management of PCOS.

Model Selection

Logistic Regression

The logistic regression model performed quite well in predicting PCOS based on the selected features:

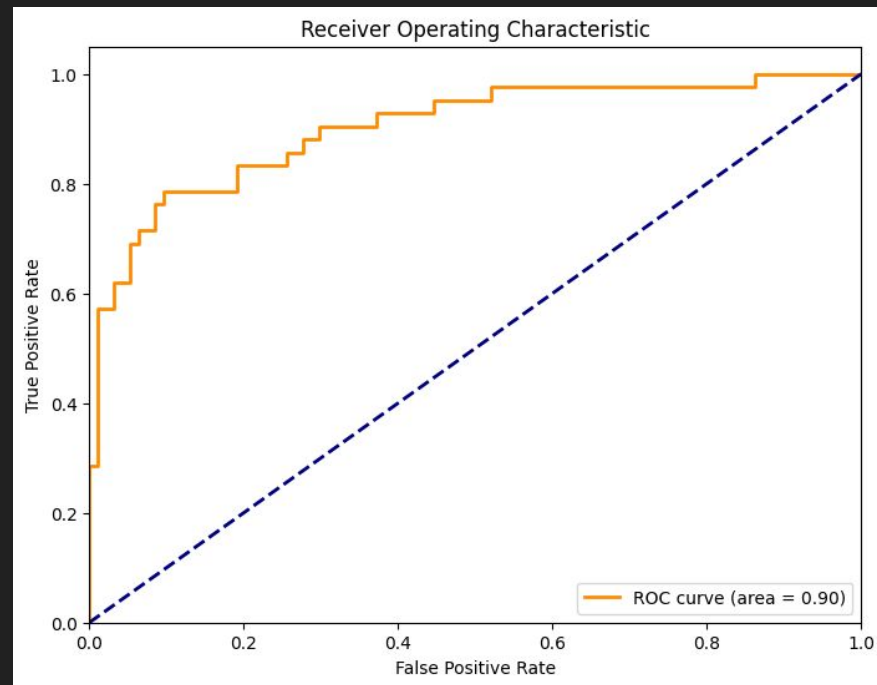
Accuracy: 86.03%

Precision for Class 1 (PCOS): 78%

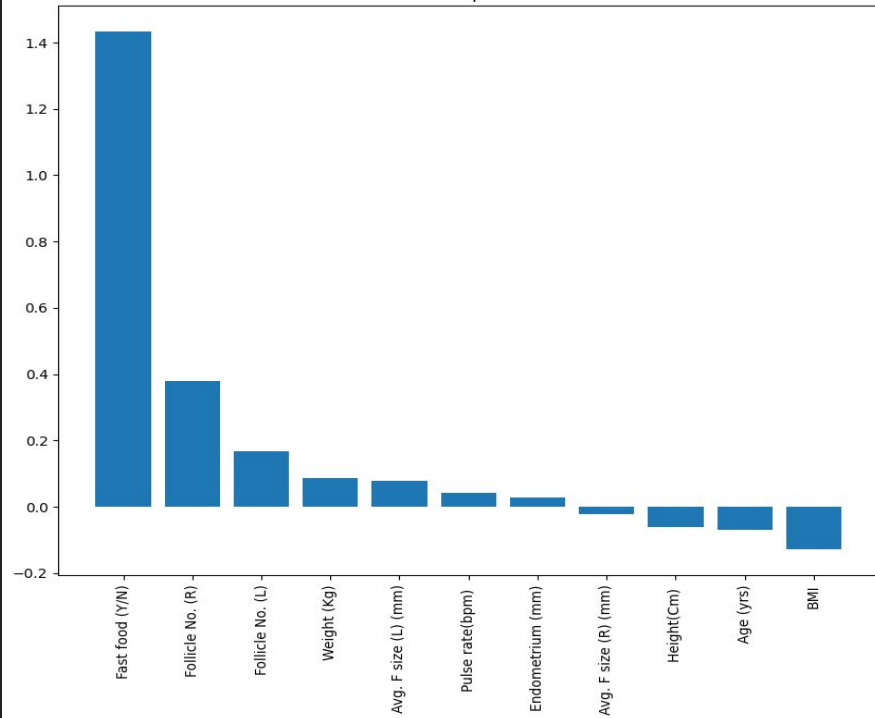
Recall for Class 1 (PCOS): 76%

F1-Score for Class 1 (PCOS): 77%

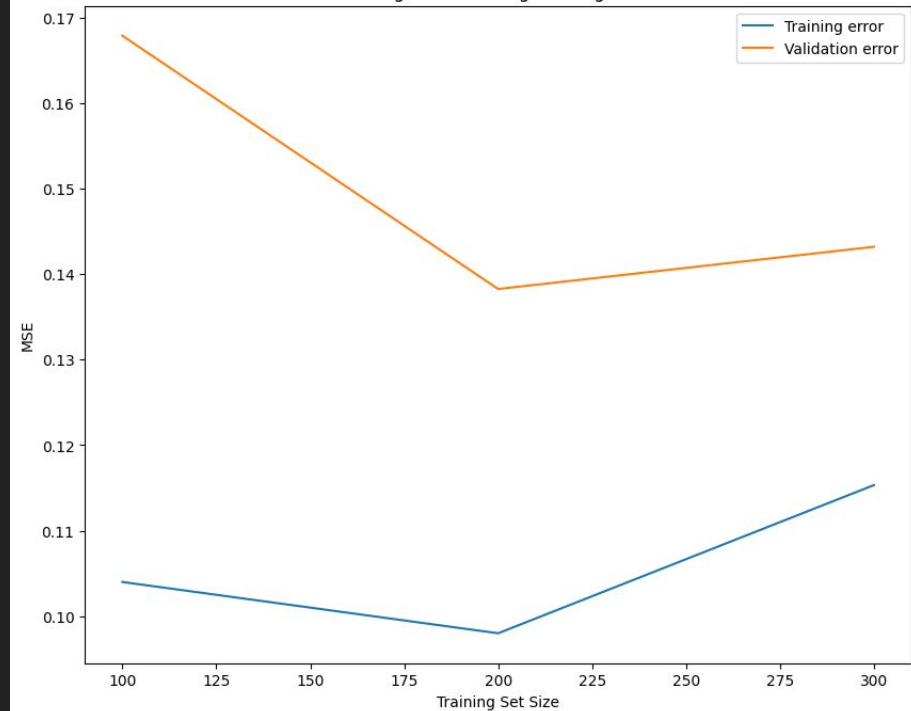
These results suggest that the model is quite effective, with good balance between precision and recall for detecting PCOS cases.



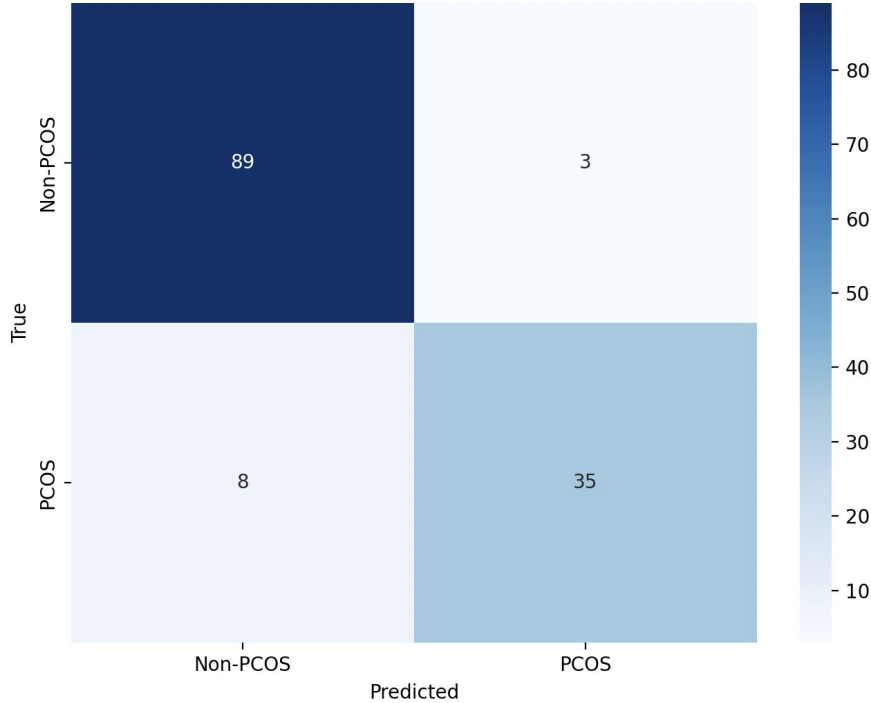
Feature Importances



Learning Curve for Logistic Regression



Confusion Matrix for Logistic Regression



Logistic Regression Results:

Accuracy:
0.9185

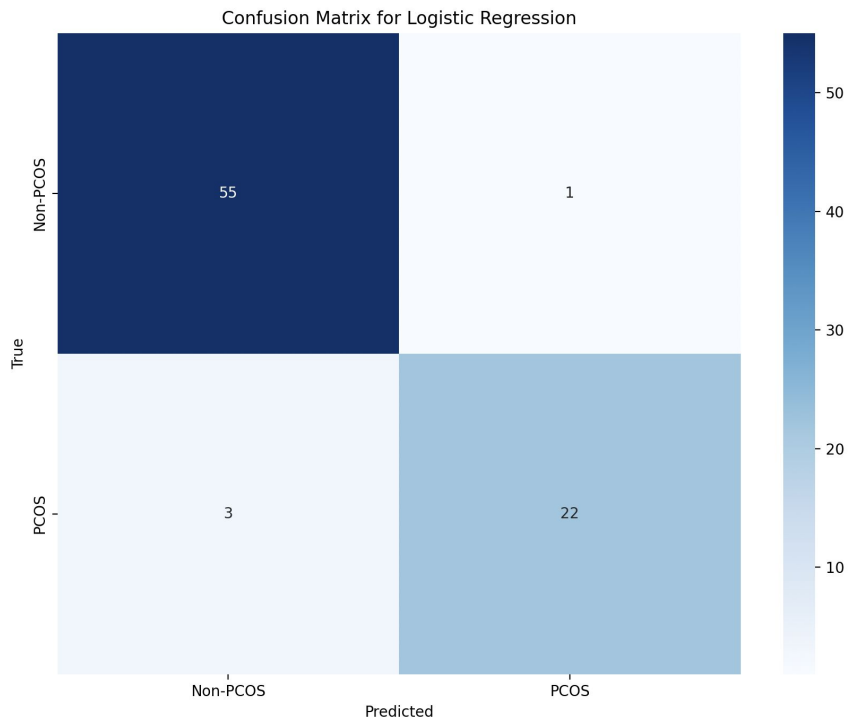
Classification_report:

	precision	recall	f1-score	support
0	0.92	0.97	0.94	92
1	0.92	0.81	0.86	43
accuracy			0.92	135
macro avg	0.92	0.89	0.90	135
weighted avg	0.92	0.92	0.92	135

Confusion_matrix:

Neural Network Results:

Accuracy:
0.9037



Logistic Regression Results:

Accuracy:
0.9506

Classification_report:

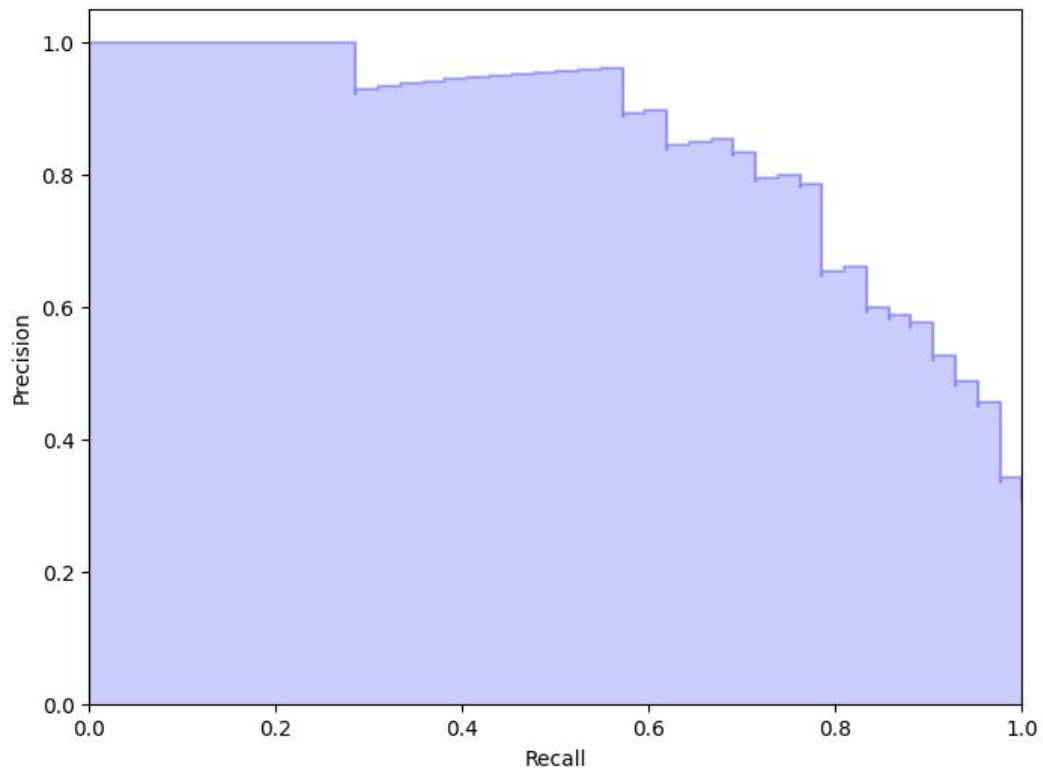
	precision	recall	f1-score	support
0	0.95	0.98	0.96	56
1	0.96	0.88	0.92	25
accuracy			0.95	81
macro avg	0.95	0.93	0.94	81
weighted avg	0.95	0.95	0.95	81

Confusion_matrix:

Neural Network Results:

Accuracy:
0.9383

2-class Precision-Recall curve: AP=0.85



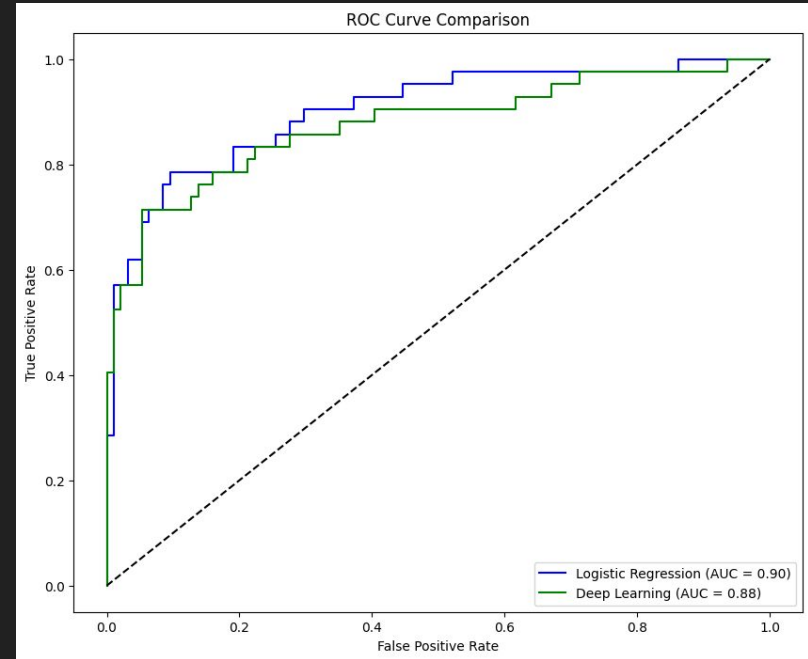
Deep Learning

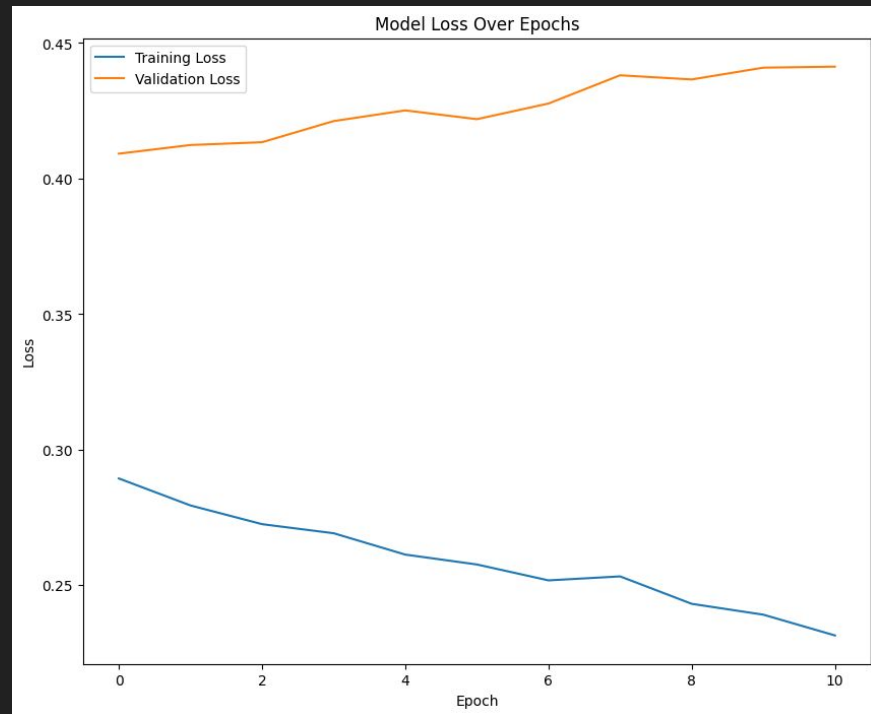
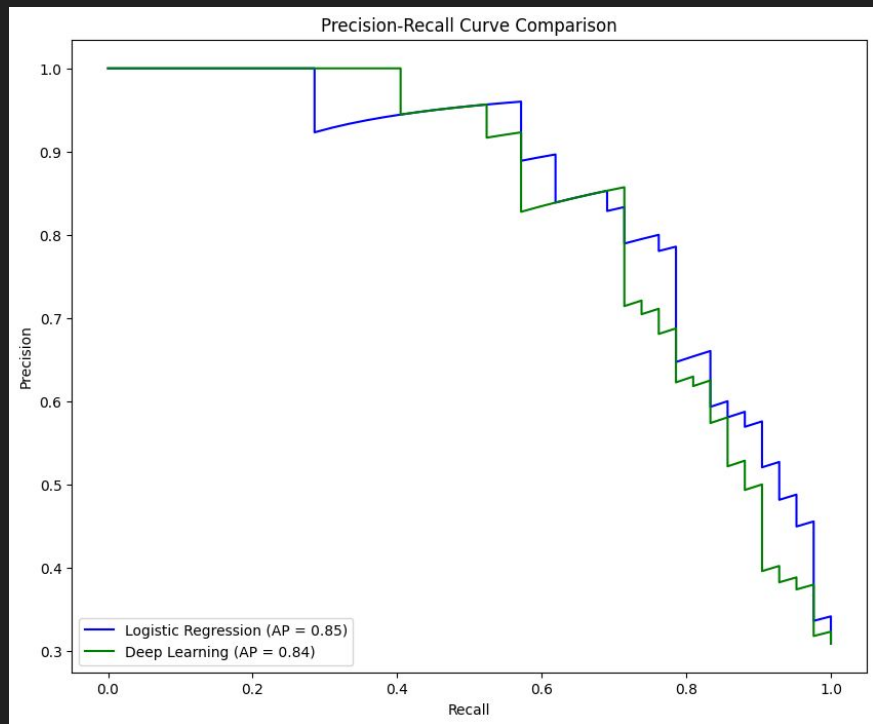
Deep Learning Test Accuracy: 84.56%

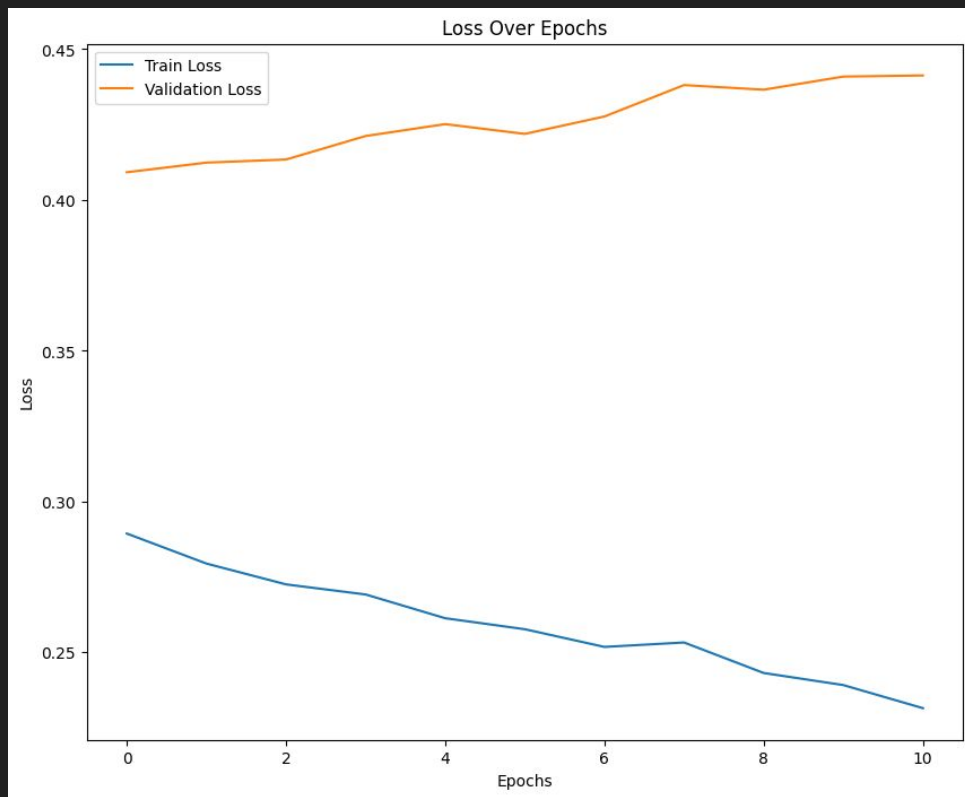
```
# Deep Learning Model
model = Sequential([
    Dense(64, activation='relu', input_shape=(X_train_scaled.shape[1],)),
    Dense(32, activation='relu'),
    Dense(1, activation='sigmoid')
])

model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
early_stopping = EarlyStopping(monitor='val_loss', patience=10, restore_best_weights=True)
model.fit(X_train_scaled, y_train, epochs=50, batch_size=10, validation_split=0.2, callbacks=[early_stopping])

# Evaluate the Deep Learning Model
loss, accuracy = model.evaluate(X_test_scaled, y_test)
print('Deep Learning Test Accuracy:', accuracy)
```







Deep learning with Cross Validation

Mean Accuracy over 5 folds:
0.879629623889923

```
# K-fold cross-validation
kfold = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
fold_no = 1
scores = []

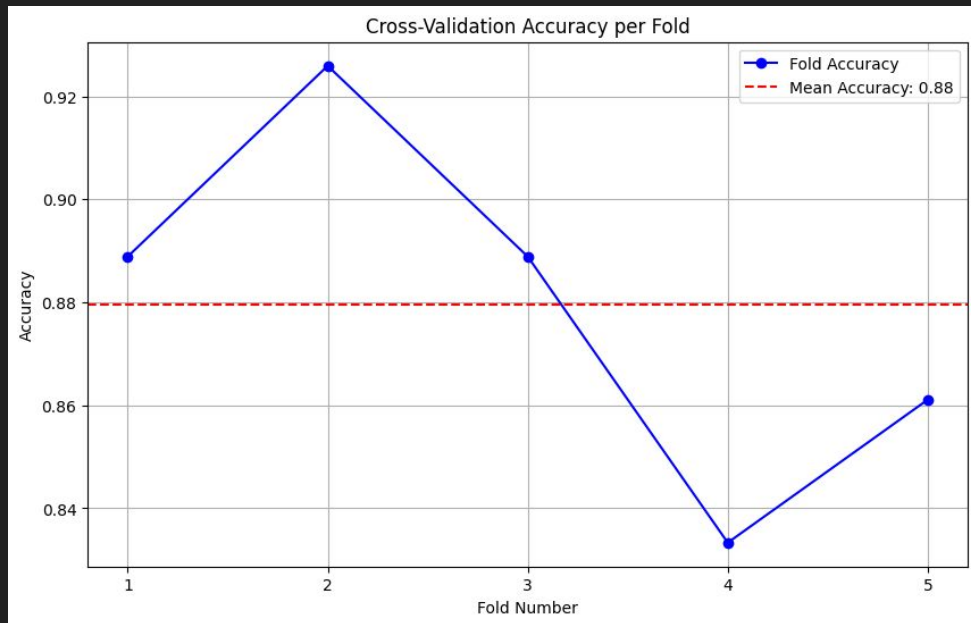
for train, test in kfold.split(X_resampled, y_resampled):
    # Neural network model
    model = Sequential([
        Dense(128, activation='relu', input_shape=(X_train_scaled.shape[1],)),
        BatchNormalization(),
        Dropout(0.5),
        Dense(64, activation='relu'),
        BatchNormalization(),
        Dropout(0.3),
        Dense(32, activation='relu'),
        Dense(1, activation='sigmoid')
    ])

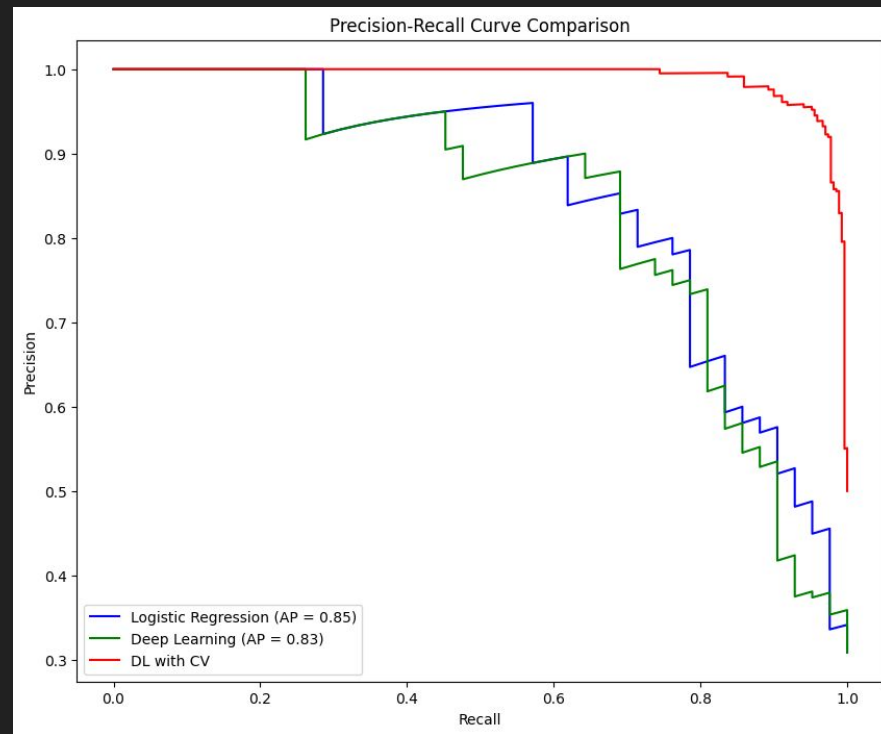
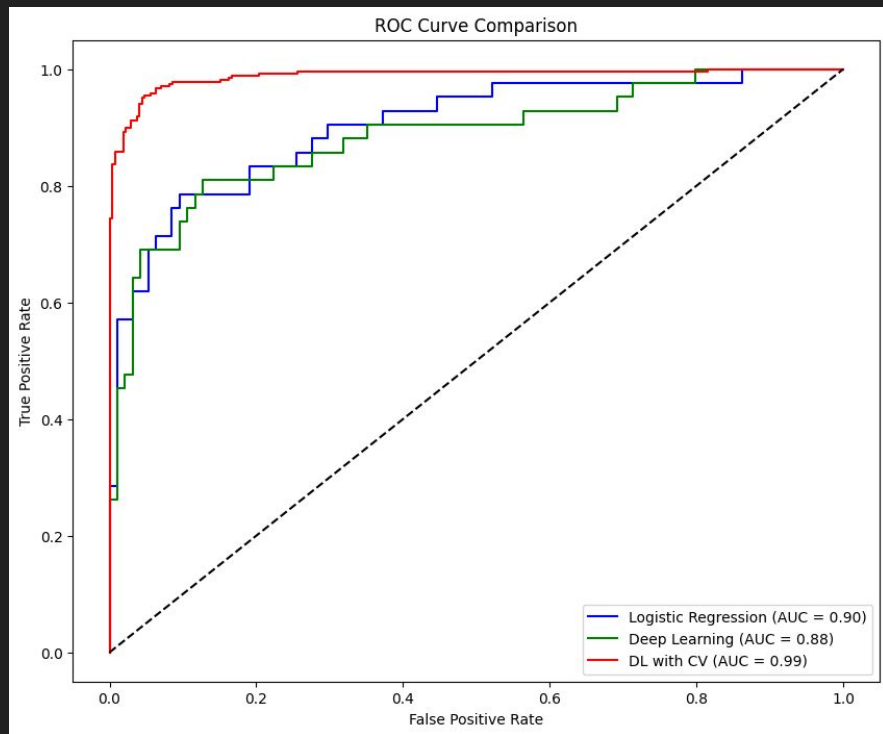
    # Compile the model
    model.compile(optimizer=Adam(learning_rate=0.001), loss='binary_crossentropy', metrics=['accuracy'])

    # Fit the model
    print(f'Training for fold {fold_no} ...')
    model.fit(X_resampled[train], y_resampled[train], epochs=50, batch_size=32, verbose=0)

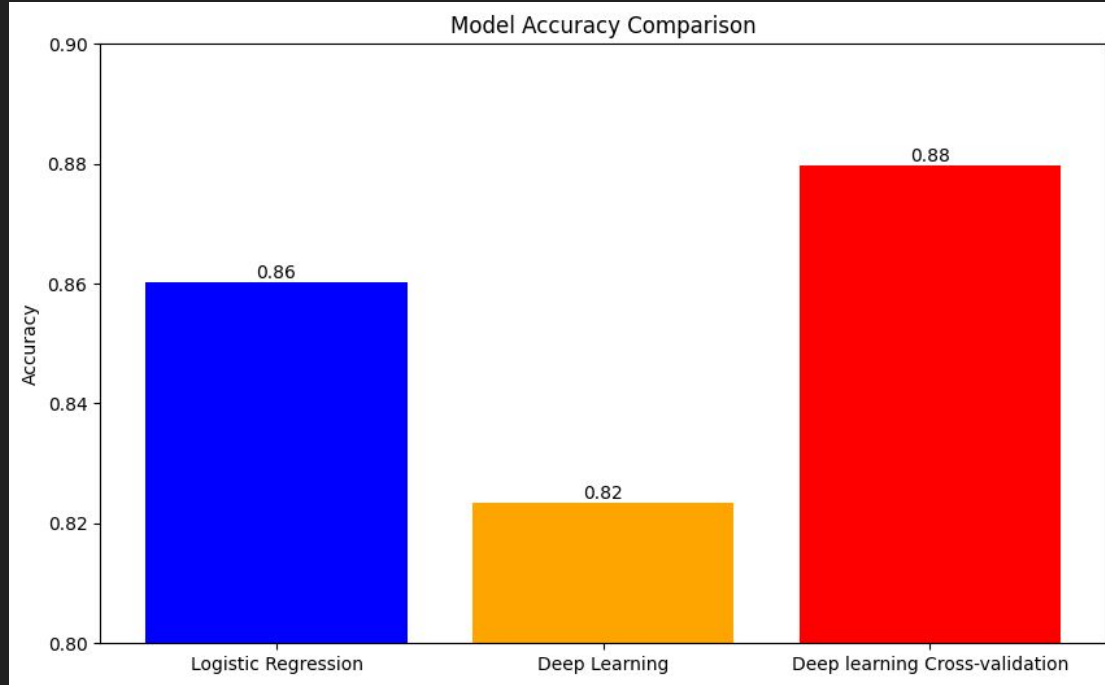
    # Evaluate the model
    scores.append(model.evaluate(X_resampled[test], y_resampled[test], verbose=0)[1])
    fold_no += 1

print(f'Mean Accuracy over {kfold.n_splits} folds: {np.mean(scores)}')
```





Model Comparisons



Results & Discussion

- The logistic regression model achieved an accuracy of 94%, with precision and recall at 96% and 93%, respectively. Feature importance analysis revealed that 'Fast Food', 'Follicle No. (R)', 'Weight' and 'Avg. F size (L) (mm)' were pivotal in predicting PCOS.
- Our deep learning model, particularly when enhanced with cross-validation, achieved an impressive accuracy of 88%, indicating robust generalizability.

Model Performance

High Accuracy Factors:

- **Advanced Capability:** Deep learning ability to capture non-linear relationships in data likely enhanced its performance.
- **Cross-Validation:** Rigorous cross-validation contributed to robustness, ensuring the model performs well on unseen data.
- **Preprocessing and Feature Selection:** Careful preprocessing and strategic feature selection based on domain expertise helped focus the model on the most relevant predictors, reducing noise.

Future Improvements

Considerations and Cautions:

- Risk of Overfitting: Deep learning models are prone to memorizing training data rather than generalizing from it, raising concerns about overfitting.
- Future Improvements:
 - Dropout Techniques: Incorporating dropout techniques and regularization could help mitigate overfitting risks.
 - External Validation: For clinical deployment, validating models on an independent external dataset is crucial to confirm their practical utility and reliability.

Exploration of data

relevant physical and clinical parameters for PCOS

- 1: Project presentation lacked any sort of depth and clarity, student(s) presenting did not clearly address issues from prompt
- 2: Project presentation included little detail and inconclusive results, difficult to understand
- 3: Project presentation addresses challenges stated in prompt clearly and concisely, easy for audience to follow and understand

PROJECT DESIGN

- 1: Project not designed well; yields no tangible results or applicability
- 2: Project designed to move forward, but disjointed approach
- 3: Project designed to move forward in systematic way

This is what the syllabus says:

https://docs.google.com/document/d/1mIJUfQBfq0tC8_9HiACewh-5BFG6YwUM-XBBFJKJt0k/edit

Coordination

Exploration of data

ADDRESSING STATED PROBLEMS

- 1: Results did not address stated problem at all
- 2: Results include missing details and inconclusive results
- 3: Results solves stated problem at face value

Fixing outliers
<decision point, challenges>

Choosing the models