# Detection and Analysis of
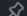# Polycystic Ovary Syndrome

## Data Science Boot Camp

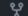Ethan Bootehsaz, Kaiyuan Ma, Joy Wang, Aaron Yang
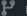
04/02/2024

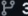# Created a Repository



Analysis-and-Detection-of-PCOS  Public

⚲ Unpin    👁 Unwatch  1  ▾    ⑂ Fork  0  ▾    ☆ Star  0  ▾

🎋 main ▾    ⑂ 3 Branches    ⬡ 0 Tags

🔍 Go to file    t    Add file ▾    <> Code ▾

About

Our goal is to develop a predictive model for the early detection of Polycystic ovary syndrome (PCOS) and infertility-related issues. We will evaluate the model's performance using standard metrics.

👤 ebootehsaz  update readme  ✕        055be89 · 9 minutes ago    🕐 21 Commits

📂 .github/workflows         Modded tests                    8 hours ago
📂 Kaggle_Data               Original data                   4 days ago
📂 data                      First clean processed data      4 days ago
📂 logs                      Changed data flow structure     8 hours ago
📂 src                       Changed data flow structure     8 hours ago
📄 .gitignore                modded .gitignore               8 hours ago
📄 LICENSE                   Initial commit                  4 days ago
📄 Makefile                  Changed data flow structure     8 hours ago
📄 PCOS dataset info.docx    First commit                    4 days ago
📄 PCOSData.ipynb            First commit                    4 days ago
📄 README.md                 update readme                   9 minutes ago
📄 requirements.txt          Changed data flow structure     8 hours ago
📄 start.sh                  precedent for logging           4 days ago

📖 Readme
⚖ MIT license
〰 Activity
☆ 0 stars
👁 1 watching
⑂ 0 forks

Releases

No releases published
Create a new release

Packages

No packages published
Publish your first package

Contributors  2

# Analysis-and-Detection-of-PCOS

Our goal is to develop a predictive model for the early detection of Polycystic ovary syndrome (PCOS) and infertility-related issues. PCOS is a complex hormonal disorder, and early diagnosis is crucial for effective management.

We will analyze the dataset to identify patterns and trends. We will develop a model that can predict the likelihood of PCOS and infertility-related issues. We will evaluate the model's performance using various metrics. We will present the results in a clear and concise manner.

## Table of Contents

# Presentation 1 (Progress Report)

- Explore the PCOS dataset to understand its structure and features.

- Identify missing values, outliers, and patterns in the data.

- Select relevant physical and clinical parameters for PCOS and infertility detection.

  - Assessment of Ovarian Reserve: AMH (Anti-Müllerian Hormone) levels give insight into the ovarian reserve, which is crucial for fertility.
  - Women with PCOS often have higher than normal AMH levels. This is because AMH is produced by the granulosa cells of ovarian follicles, and women with PCOS tend to have a higher number of small follicles in their ovaries.

- Perform any necessary preprocessing steps, such as handling missing values or encoding categorical variables.

- Choose appropriate machine learning models for binary classification (e.g., Logistic Regression)

- Formulate hypotheses related to PCOS and infertility based on the dataset.

- Create visualizations to support your findings.

- Generate tables/graphs to show the distribution of the target variable and key features.

# Created a Project Board

# Script to help get started

**Analysis-and-Detection-of-PCOS** / **start.sh**

ebootehsaz   precedent for logging

| Code | Blame |   Executable File · 15 lines (12 loc) · 453 Bytes

```bash
1    #!/bin/bash
2
3    # Check if the virtual environment already exists
4    if [ ! -d "dsbootcamp" ]; then
5        # Create a virtual environment named dsbootcamp if it doesn't exist
6        python3 -m venv dsbootcamp
7    fi
8
9    # Activate the virtual environment and install the required packages
10   source dsbootcamp/bin/activate
11   pip install -r requirements.txt >> logs/requirements.log 2>&1
12
13   # Pull updates from Git repository
14   echo "Pulling updates from Git repository..."
15   git pull
```

Code    Blame    88 lines (77 loc) · 2.89 KB

```python
import re
import pandas as pd
import os
from pandas import DataFrame

# This program defines utility functions to load, merge, and save data.

from constants import PCOS_woinf_filepath_page

def load_data(filepath: str) -> DataFrame:
    """
    Load data from a file path into a DataFrame.
    Args:
        filepath: str - file path to the data
    Returns:
        df - data loaded into a DataFrame
    """
    # Check if the file path exists
    if not os.path.exists(filepath):
        # remove ../ from the start of filepath
        filepath = re.sub(r'\.\./', '', filepath)
        if not os.path.exists(filepath):
            raise FileNotFoundError(f"File path {filepath} does not exist.")

    # They data is either in csv or excel format
    if filepath.endswith('.csv'):
        df = pd.read_csv(filepath)
        df.attrs['file_path'] = filepath  # Storing file path as an attribute
    elif filepath.endswith('.xlsx'):
        sheet_name = PCOS_woinf_filepath_page
        df = pd.read_excel(filepath, sheet_name)
        df.attrs['file_path'] = filepath  # Storing file path as an attribute
    else:
        raise ValueError(f"File path {filepath} is not a csv or excel file.")

    return df
```
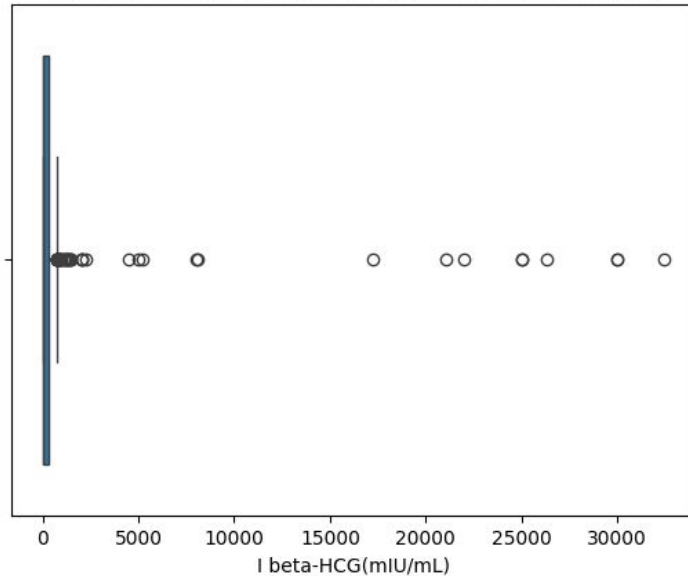
```python
import re
from pandas import DataFrame

from constants import PCOS_inf_filepath, PCOS_woinf_filepath, \
    PCOS_inf_processed_filepath, PCOS_woinf_processed_filepath, PCOS_merged_processed_filepath

from utils import load_data, merge_data, save_data_csv

# This program defines utility functions to process and clean data.

def process_data(df: DataFrame) -> DataFrame:
    # Checking column names, missing values, duplicates
    print(f"Columns in {df.attrs['file_path']}:", df.columns)
    print(f"Missing values in {df.attrs['file_path']}:\n{df.isnull().sum()}")
    print(f"Duplicates in {df.attrs['file_path']}: {df.duplicated().sum()}")

    #Dropping repeated/unnecessary columns
    df = df.drop(['Unnamed: 44','Sl. No_wo', 'PCOS (Y/N)_wo', '  I  beta-HCG(mIU/mL)_wo','II   beta-HCG(mIU/mL)_wo', 'AMH(ng/mL)_wo'], 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]

    # Fix missing values
    # Print out the first 5 missing rows for each column with missing values
    # Find rows with missing data across any column
    rows_with_missing_data = df[df.isnull().any(axis=1)]

    # Display the rows with missing data if any
    if not rows_with_missing_data.empty:
        print("Rows with missing data:")
        print(rows_with_missing_data)
    else:
        print("No missing data in any row.")

    df = df.fillna('None')

    # Drop duplicates
    df = df.drop_duplicates()

    # Take a random sample of the data
    print(f"Sample of the data in {df.attrs['file_path']}:", df.sample(5))

    return df
```

# Sample of data

# Boxplot of I beta-HCG (mIU/mL)

# Boxplot II beta-HCG (mIU/mL)

# Boxplot of AMH (ng/mL)

Percentage of Data Remaining vs. Removed Due to Outliers in Merged Dataset

Distribution of PCOS (Y/N)

Distribution of Weight (Kg)

Distribution of Blood Group

Distribution of Age (yrs)

Distribution of Height(Cm)

Distribution of BMI

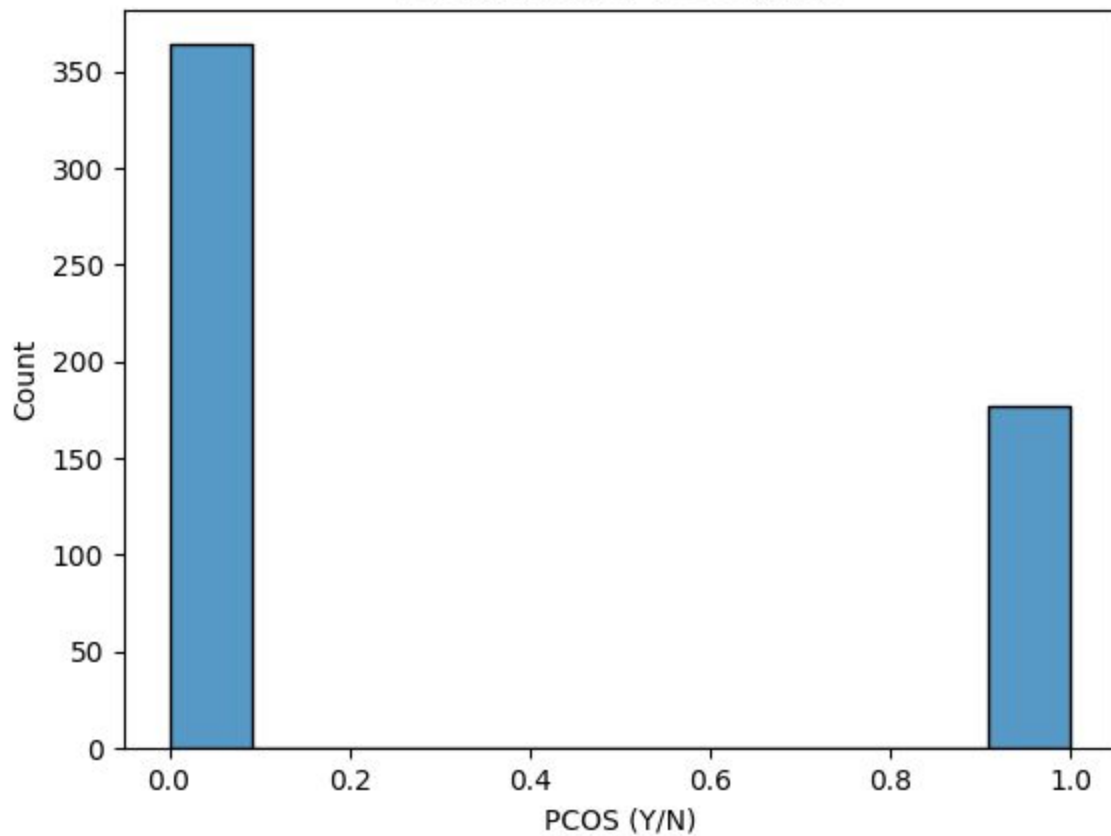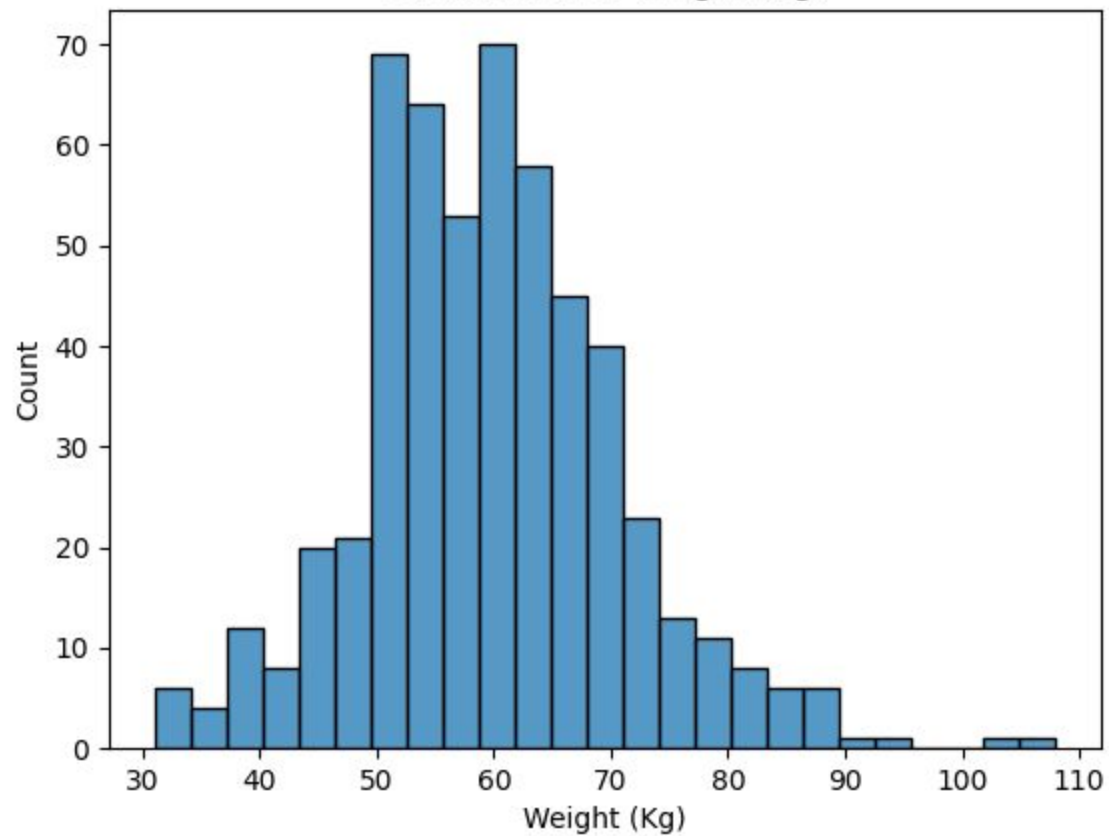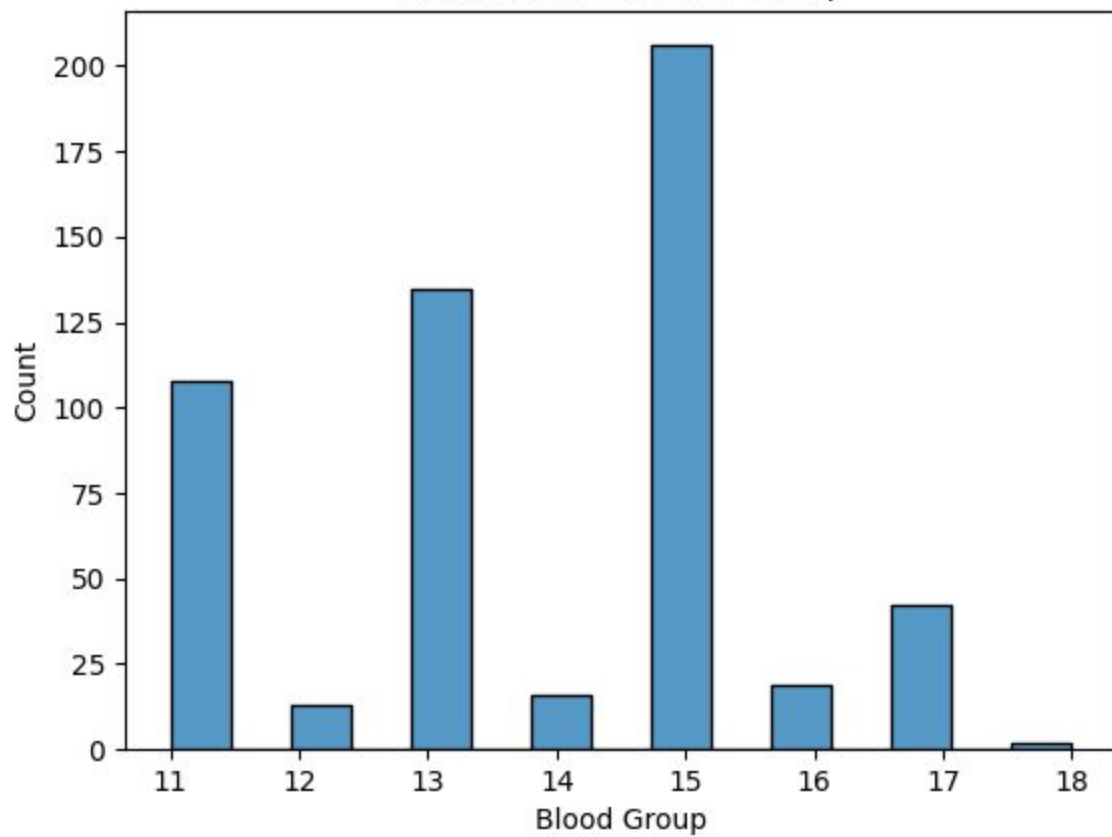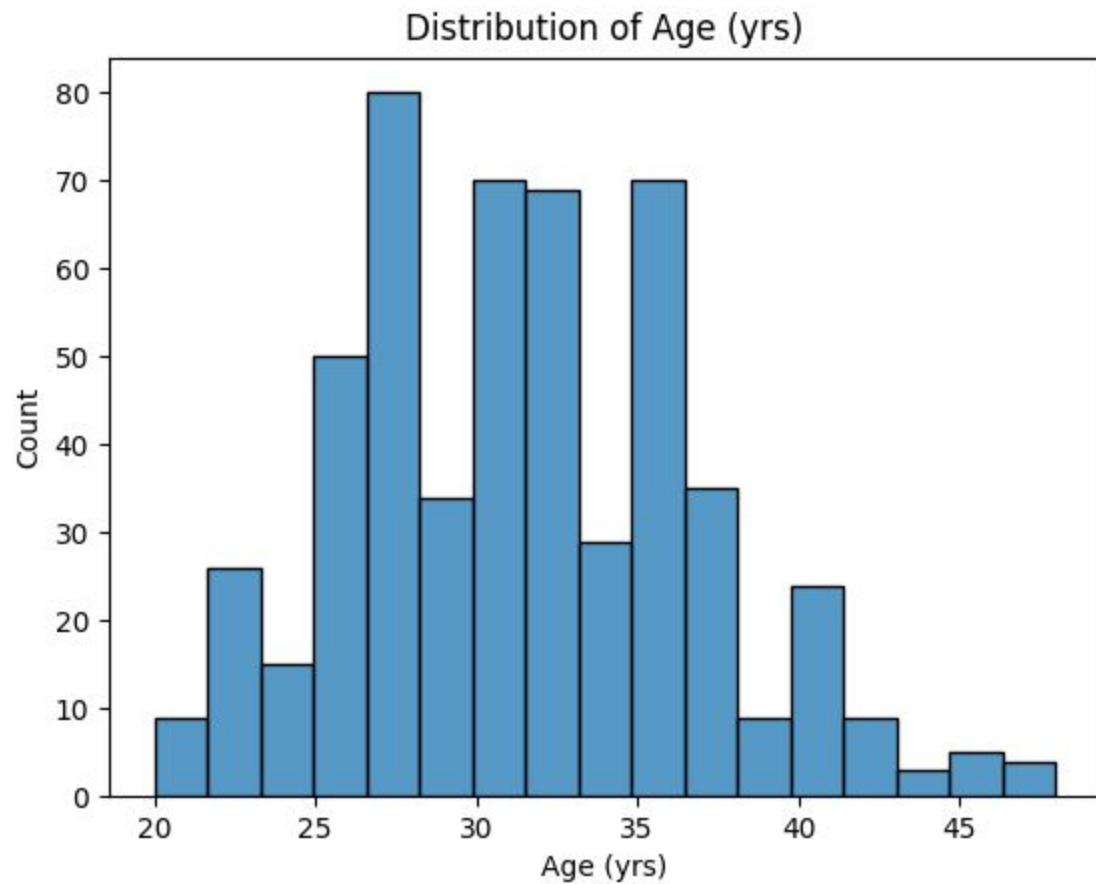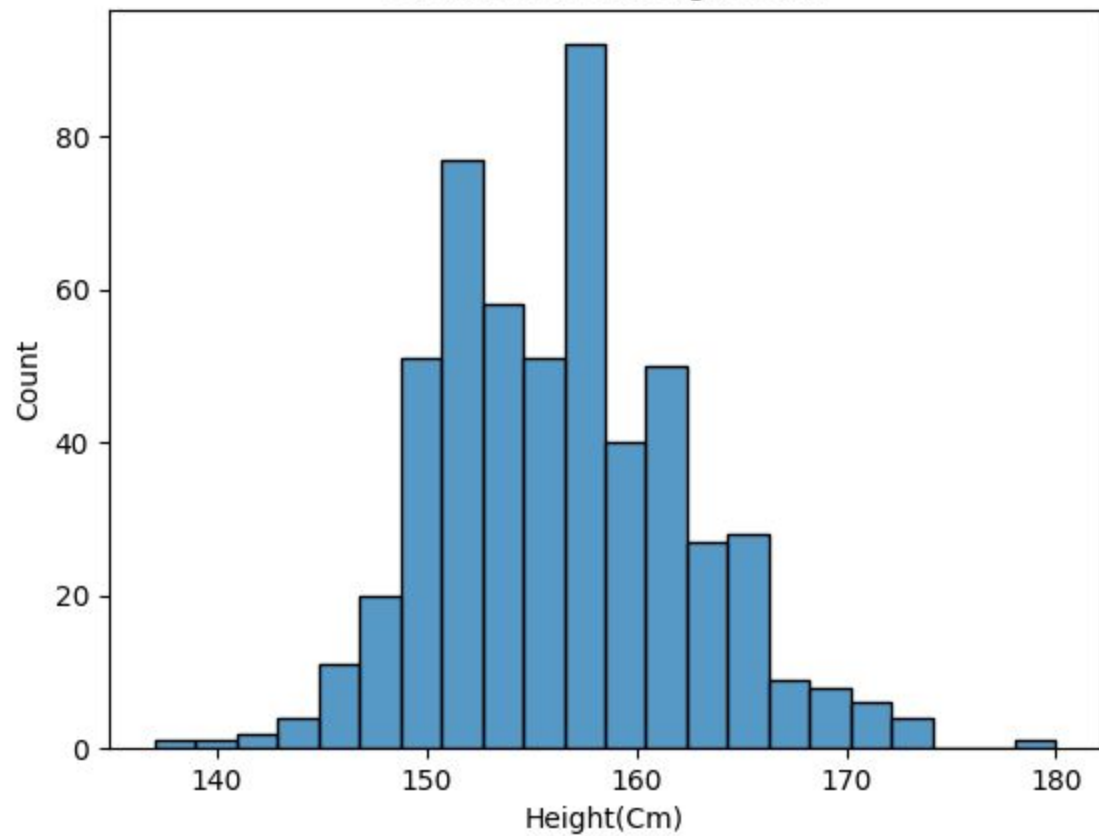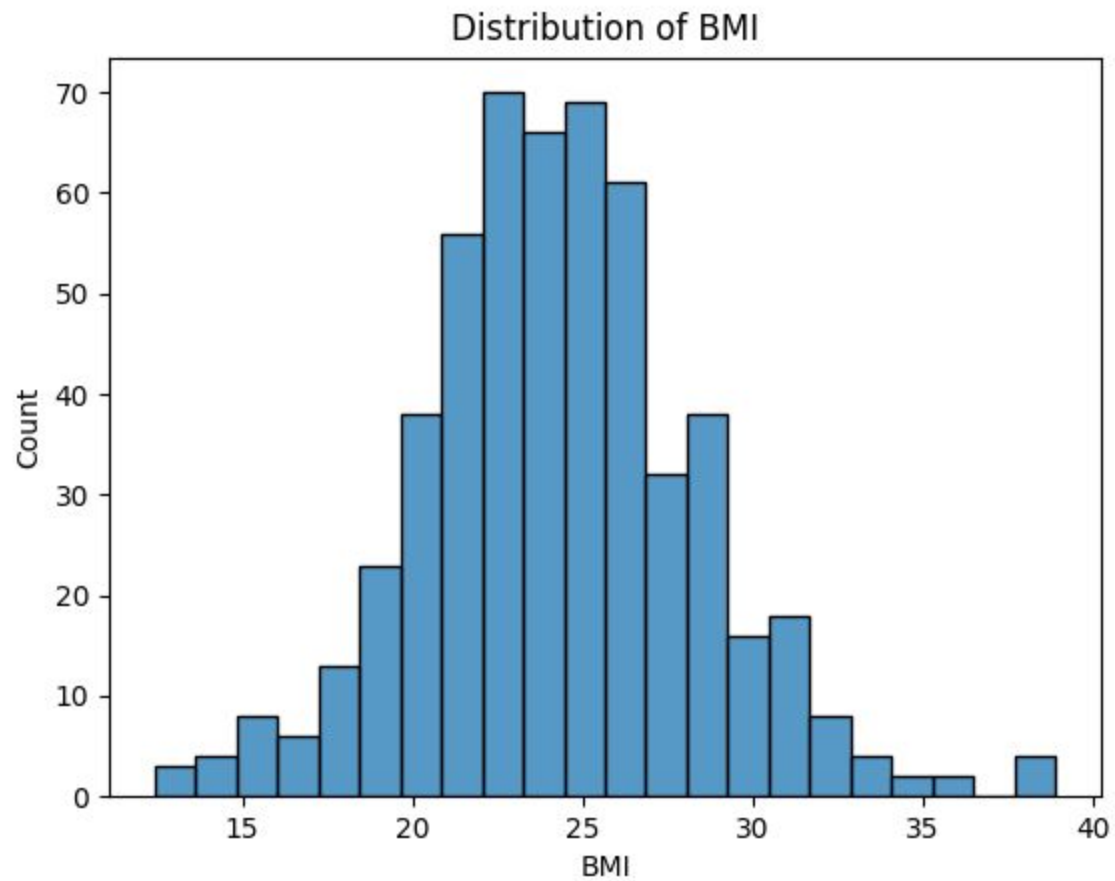| | |
|---|---|
| PCOS (Y/N) | 1.000000 |
| Follicle No. (R) | 0.648327 |
| Follicle No. (L) | 0.603346 |
| Skin darkening (Y/N) | 0.475733 |
| hair growth(Y/N) | 0.464667 |
| Weight gain(Y/N) | 0.441047 |
| Cycle(R/I) | 0.401644 |
| Fast food (Y/N) | 0.376183 |
| Pimples(Y/N) | 0.286077 |
| AMH(ng/mL) | 0.264141 |
| Weight (Kg) | 0.211938 |
| BMI | 0.199534 |
| Hair loss(Y/N) | 0.172879 |
| Waist(inch) | 0.164598 |
| Hip(inch) | 0.162297 |
| Avg. F size (L) (mm) | 0.132992 |
| Endometrium (mm) | 0.106648 |
| Avg. F size (R) (mm) | 0.097690 |
| Height(Cm) | 0.068254 |
| Reg.Exercise(Y/N) | 0.065337 |
| LH(mIU/mL) | 0.063879 |
| RBS(mg/dl) | 0.048922 |
| BP _Diastolic (mmHg) | 0.038032 |
| RR (breaths/min) | 0.036928 |
| Blood Group | 0.036433 |
| II    beta-HCG(mIU/mL) | 0.012760 |
| Waist:Hip Ratio | 0.012386 |
| BP _Systolic (mmHg) | 0.007942 |
| PRL(ng/mL) | 0.005143 |
| TSH (mIU/L) | -0.010140 |
| FSH/LH | -0.018336 |
| Pregnant(Y/N) | -0.027565 |
| I    beta-HCG(mIU/mL) | -0.027617 |
| FSH(mIU/mL) | -0.030319 |
| PRG(ng/mL) | -0.043834 |
| No. of aborptions | -0.057158 |
| Marraige Status (Yrs) | -0.113056 |
| Age (yrs) | -0.168513 |
| Cycle length(days) | -0.178480 |

# Assumptions

We first want to focus on the first three feature that is present on both data set which were 'I beta-HCG(mIU/mL)', 'II beta-HCG(mIU/mL)', 'AMH(ng/mL)' as features to train our dataset.

And our target variable will be PCOS positive/negative rate

# Future Goals

- Develop a predictive model for the early detection of PCOS and infertility-related issues.
- Analyze the dataset to identify patterns and trends.
- Develop a model that can predict the likelihood of PCOS and infertility-related issues.
- Evaluate the model's performance using various metrics.
- Present the results in a clear and concise manner.

# Feedback 4/2 12pm

- Take good traceability of every decision point where the experiment changed / was done / for what reason – Dani
- Oh that's all


- Need more stats stuff