# Report on Mini Project

# Subject: Machine Learning

# AY: 2024-25

## Air Passenger Data Analysis - Time Series Forecasting -

## SARIMAX

**Parth Puranik:**     **2103142**

**Mohib Abbas Sayed: 2103158**

**Hamza Sayyed:**     **2103159**

**Om Shete:**     **2103163**

Guided By

(Tanuja Sarode)

# CHAPTER 1: INTRODUCTION

In this project, we aim to forecast the number of air passengers for the next 24 months based on historical data from 1949 to 1960. This type of analysis is crucial in the aviation industry as it aids in capacity planning, demand forecasting, and ensuring smooth operations. The dataset consists of monthly totals of US airline passengers over 12 years. Time series forecasting is essential for capturing the data's underlying patterns such as seasonality, trends, and cycles.

We use the ARIMA (AutoRegressive Integrated Moving Average) model, particularly SARIMAX (Seasonal ARIMA with Exogenous Regressors), which accounts for seasonality and other external factors, to provide accurate forecasts. We aim to implement this model and generate estimates for the upcoming 24 months based on the data trends from 1949 to 1960.

# CHAPTER 2: DATA DESCRIPTION AND ANALYSIS

## 2.1 Dataset Overview

The dataset used in this project contains two key attributes:

- **Month**: The period in which the passengers travelled.

- **Passengers**: The total number of passengers who travelled in each month.

The dataset spans from January 1949 to December 1960, consisting of 144 observations. Each data point represents the total number of passengers for that month.

## 2.2 Data Preprocessing

Before proceeding with time series modelling, the following preprocessing steps were performed:

- **Parsing Dates**: The 'Month' attribute was converted into a DateTime object for time-based operations.

- **Handling Missing Data**: The dataset did not contain any missing values.
- **Stationarity Check**: The Augmented Dickey-Fuller (ADF) test was applied to check the stationarity of the series.
- **Seasonality and Trends**: Visualizations such as line plots and seasonal decomposition were used to identify trends, seasonality, and cycles in the data.

## 2.3 Exploratory Data Analysis (EDA)

Using libraries such as matplotlib and seaborn, we generated the following visual insights:

- **Trend Analysis**: A clear upward trend in the number of passengers over the years.

- **Seasonal Patterns**: A yearly recurring seasonal pattern was evident, with peak months typically being mid-year (summer travel).
- **Correlation**: We analyzed autocorrelations and partial autocorrelations to explore the relationship between monthly passenger numbers over time.

# CHAPTER 3: DESIGN OF DATA PIPELINE

The data pipeline for this project is designed to ensure a seamless flow from data ingestion to model deployment. The pipeline consists of the following steps:

## 3.1 Data Ingestion

The dataset was loaded using the pandas library, and date parsing was performed to ensure proper time series handling.

## 3.2 Data Preprocessing

- **Handling DateTime Format**: The 'Month' column was converted to a proper datetime index to facilitate time series analysis.

- **Differencing**: To achieve stationarity, differencing techniques were applied to remove trends and seasonality from the data.

- **Scaling**: Data normalization was applied to ensure consistency across variables, preparing the dataset for the ARIMA model.

## 3.3 Model Selection and Training

**ARIMA Model**: The ARIMA model was selected based on the dataset characteristics. We utilized SARIMAX (Seasonal ARIMA) due to the clear seasonality in the data. The model parameters (p, d, q) and seasonal components were tuned using grid search methods to optimize performance.

The model training process involved splitting the data into training and testing sets, followed by fitting the SARIMAX model on the training data.

## 3.4 Model Evaluation

- **Performance Metrics**: Metrics such as Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and AIC were used to assess model performance.

- **Cross-Validation**: The model's robustness was evaluated using time-series cross-validation techniques.

## 3.5 Forecasting

After model training and evaluation, we generated forecasts for the next 24 months. Visualizations were created to compare the forecasts with actual data, providing insights into expected trends and passenger growth.
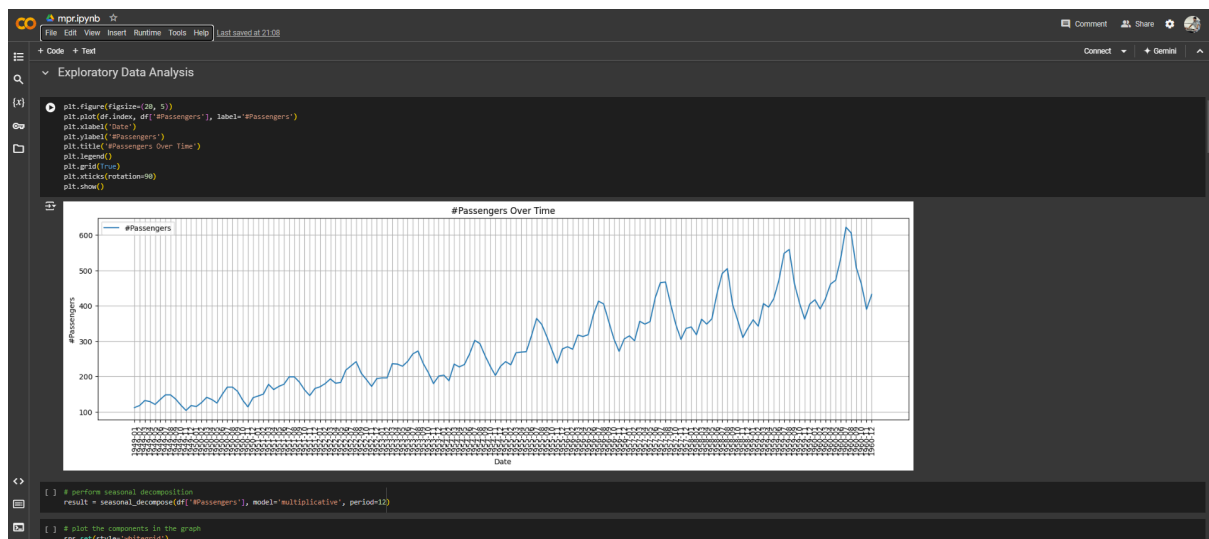
# CHAPTER 4: RESULT ANALYSIS

```python
seasonal_period = 12
```

```python
from statsmodels.tsa.stattools import adfuller # Augmented Dickey-Fuller Test

result = adfuller(df['#Passengers'], autolag='AIC') # Akaike Information Criterion
print('ADF Statistic:', result[0])
print('p-value:', result[1])
```

```
ADF Statistic: 0.815368079260441
p-value: 0.9918802434376489
```

```python
# first order differencing
result = adfuller(df['#Passengers'].diff().dropna(), autolag='AIC')
print('ADF Statistic:', result[0])
print('p-value:', result[1])
```

```
ADF Statistic: -2.829266824170034
p-value: 0.05421329028382497
```

```python
# second order differencing
result = adfuller(df['#Passengers'].diff().diff().dropna(), autolag='AIC')
print('ADF Statistic:', result[0])
print('p-value:', result[1])
```

```
ADF Statistic: -16.384231542460852
p-value: 2.732891850014085e-29
```

```python
# plot the differencing values
fig, (ax1, ax2, ax3) = plt.subplots(3)

ax1.plot(df)
ax1.set_title('Original Time Series')
ax1.axes.xaxis.set_visible(False)

ax2.plot(df.diff())
ax2.set_title('1st Order Differencing')
ax2.axes.xaxis.set_visible(False)

ax3.plot(df.diff().diff())
ax3.set_title('2nd Order Differencing')
ax3.axes.xaxis.set_visible(False)

plt.show()
```



```python
# the time series becomes stationary after first order differencing
```

## Define Parameters for ARIMA

```python
# p = 0 # MA - Moving Average - PACF
# d = 1 # order of differencing - I
# q = 0 # AR - Auto Regressive - ACF
```

```python
fig, ax = plt.subplots(2, 1, figsize=(12, 7))
sm.graphics.tsa.plot_acf(df.diff().dropna(), lags=40, ax=ax[0])
sm.graphics.tsa.plot_pacf(df.diff().dropna(), lags=40, ax=ax[1])
plt.show()
```



```python
p = 2 # pacf
d = 1 # 1st order difference
q = 1 # acf
```

```python
P = 1
D = 0
Q = 3
```

## Model Training

```python
# define the arima model
from statsmodels.tsa.statespace.sarimax import SARIMAX

model = SARIMAX(df['#Passengers'], order=(p, d, q), seasonal_order=(P, D, Q, seasonal_period))
fitted_model = model.fit()
print(fitted_model.summary())
```

```
                              SARIMAX Results
==============================================================================
Dep. Variable:                         #Passengers   No. Observations:                  144
Model:             SARIMAX(2, 1, 1)x(1, 0, [1, 2, 3], 12)   Log Likelihood             -563.224
Date:                            Mon, 30 Sep 2024   AIC                       1142.448
Time:                                    18:26:03   BIC                       1166.151
Sample:                                01-01-1949   HQIC                      1152.080
                                     - 12-01-1960
Covariance Type:                              opg
==============================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------
ar.L1          0.6244      0.101      6.168      0.000       0.426       0.823
ar.L2          0.1947      0.100      1.951      0.051      -0.001       0.390
ma.L1         -0.9675      0.039    -24.632      0.000      -1.045      -0.891
ar.S.L12       0.9619      0.036     26.615      0.000       0.891       1.033
ma.S.L12       0.1117      0.126     -0.898      0.369      -0.359       0.133
ma.S.L24       0.1355      0.129      1.053      0.292      -0.117       0.388
ma.S.L36       0.0049      0.147      0.033      0.973      -0.284       0.294
sigma2       124.2009     14.750      8.421      0.000      95.299     153.119
==============================================================================
Ljung-Box (L1) (Q):                   0.01   Jarque-Bera (JB):        16.14
Prob(Q):                              0.93   Prob(JB):                 0.00
Heteroskedasticity (H):               3.99   Skew:                     0.18
Prob(H) (two-sided):                  0.00   Kurtosis:                 4.61
==============================================================================

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).
```

## Forecasting

```python
# forecast for next 2 years
forecast_steps = 24
forecast = fitted_model.get_forecast(steps=forecast_steps)
```
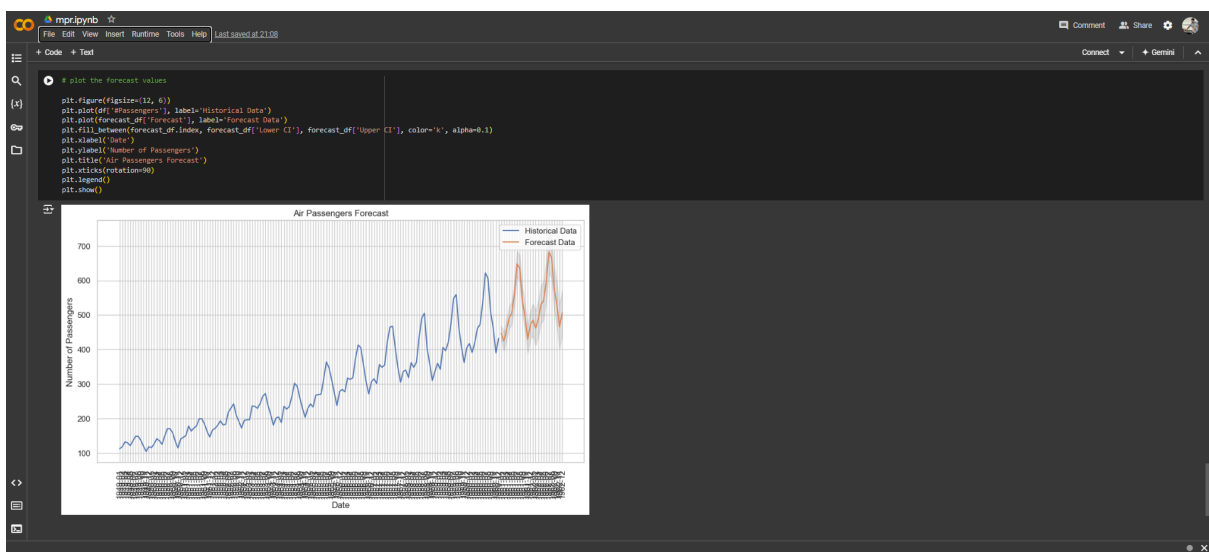
```python
# forecast for next 2 years
forecast_steps = 24
forecast = fitted_model.get_forecast(steps=forecast_steps)

# create the date range for the forecasted values
forecast_index = pd.date_range(start=df.index[-1], periods=forecast_steps+1, freq='M')[1:].strftime('%Y-%m') # remove start date
```

```python
# create a forecast dataframe
forecast_df = pd.DataFrame({
    "Forecast": list(forecast.predicted_mean),
    "Lower CI": list(forecast.conf_int().iloc[:, 0]),
    "Upper CI": list(forecast.conf_int().iloc[:, 1])
}, index=forecast_index)

forecast_df.head()
```

|         | Forecast   | Lower CI   | Upper CI   |
|---------|------------|------------|------------|
| 1961-01 | 446.711482 | 424.867844 | 468.555119 |
| 1961-02 | 423.325499 | 397.191175 | 449.459823 |
| 1961-03 | 456.418435 | 426.807576 | 486.029294 |
| 1961-04 | 491.562749 | 459.538919 | 523.586579 |
| 1961-05 | 505.131996 | 471.260404 | 539.003589 |

```python
# plot the forecast values
plt.figure(figsize=(12, 6))
plt.plot(df['#Passengers'], label='Historical Data')
plt.plot(forecast_df['Forecast'], label='Forecast Data')
plt.fill_between(forecast_df.index, forecast_df['Lower CI'], forecast_df['Upper CI'], color='k', alpha=0.1)
plt.xlabel('Date')
plt.ylabel('Number of Passengers')
plt.title('Air Passengers Forecast')
plt.xticks(rotation=90)
plt.legend()
plt.show()
```


Air Passengers Forecast

# CHAPTER 5: CONCLUSION AND FUTURE SCOPE

## 5.1 Conclusion

In this project, we successfully implemented time series forecasting on the Air Passenger dataset using the SARIMAX (Seasonal ARIMA with Exogenous Regressors) model. Through exploratory data analysis, we identified clear trends and seasonality in the dataset, which were crucial for accurate forecasting.

Our analysis revealed that air passenger traffic has been steadily increasing over time, with consistent seasonal peaks during the summer months. By applying the SARIMAX model, we were able to capture these patterns effectively and generate reliable forecasts for the next 24 months.

The accuracy of the model was evaluated using key performance metrics such as Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE). The results showed that the SARIMAX model was able to make reasonably accurate predictions, though there were some challenges with periods of higher variability in passenger numbers.

This project demonstrates the effectiveness of time series analysis and forecasting techniques in the airline industry, where accurate demand predictions can lead to better resource management, optimized flight scheduling, and improved customer service.

## 5.2 Future Scope

While the results of this analysis are promising, there are several areas for future improvement and exploration:

- **Model Refinement**: Future work can focus on further tuning the SARIMAX model or experimenting with more advanced models like Prophet or deep learning-based models (e.g., LSTM or GRU) to improve accuracy and better capture fluctuations in data.

- **Incorporating External Factors**: The current model only uses historical passenger data for predictions. Incorporating external factors such as fuel prices, economic conditions, or weather data could provide a more comprehensive forecast.

- **Real-Time Forecasting**: The current model is built on historical data up to 1960. A real-time data integration approach could allow for continuous updates to forecasts, enabling more dynamic and responsive models in a real-world setting.

- **Longer-Term Forecasting**: While this project focused on a 24-month forecast horizon, future studies could explore longer-term predictions, considering multi-year trends and global shifts in air travel demand.

- **Data Augmentation**: Exploring additional data sources, such as data from multiple airlines or other regions, could enhance model performance and provide broader insights into air passenger traffic patterns.

- **Forecast Validation in Modern Context**: While this dataset is historic, testing the forecasting model on more recent datasets would provide insights into how well these traditional models perform in today's airline industry context, which may have new patterns due to modern technologies and global factors like pandemics.

By expanding the scope of the model and incorporating more advanced techniques, this study can be extended to provide even more robust and actionable insights for the airline industry.