

a. Data Abstraction

✓ (i). Data Collection

The Fitness Tracker Dataset contains detailed information about individuals' fitness metrics, exercise routines, and health parameters. This dataset is designed to provide insights into fitness trends, workout habits, and overall health patterns. It is ideal for exploratory data analysis (EDA), machine learning applications, and health analytics. The dataset can help identify relationships between physical activity, body metrics, and health outcomes.

Features:

- Participants ID: A unique number for each participant.
- Date: Particular date that the participants are doing workout.
- Age: Age of the individual in years.
- Gender: Gender of the individual (e.g., Male, Female).
- Height (m): Height of the individual in meters.
- Weight (kg): Weight of the individual in kilograms.
- Activity Type: A type of activity that a participant doing.
- Duration: the amount of the workout done.
- Max_BPM: Maximum heartbeats per minute recorded during exercise.
- Avg_BPM: Average heartbeats per minute during a workout session.
- Resting_BPM: Resting heartbeats per minute.
- Session_Duration (hours): Duration of the workout session in hours.
- Calories_Burned: Total calories burned during a workout session.
- Experience_Level: Level of fitness experience (e.g., Beginner, Intermediate, Advanced).
- BMI: Body Mass Index, calculated as $\text{weight (kg)} / \text{height (m)}^2$.


Usage:

This dataset is suitable for:

Analyzing the impact of fitness routines on health metrics. Exploring trends in heart rate, calorie burn, and workout habits. Correlating body metrics like BMI and fat percentage with exercise patterns. Building predictive models for fitness and health analytics. This is a synthetic dataset created for educational and analytical purposes and does not represent real-world data.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline
from matplotlib.pyplot import rcParams
rcParams['figure.figsize'] = 10, 6
from datetime import datetime
import matplotlib.pyplot as plt
```

```
# Load the dataset
df = pd.read_csv("/content/Major project (1).csv",parse_dates=['date'])
df
```



	participant_id	date	age	gender	height_cm	weight_kg	activity_type
0	579	2024-07-20	58	F	160.4	99.2	Swimming
1	1934	2024-08-23	34	F	173.8	95.4	Weight Training
2	953	2024-05-18	50	F	159.8	89.3	Running
3	312	2024-11-12	57	F	165.4	112.2	Cycling
4	2830	2024-01-27	41	M	169.1	67.5	Cycling
...
1995	131	2024-11-13	33	M	184.1	153.6	Weight Training
1996	2947	2024-04-21	26	M	176.9	118.5	Yoga
1997	2337	2024-04-17	37	F	165.8	86.7	Cycling
1998	217	2024-07-12	36	F	162.9	88.8	Weight Training
1999	523	2024-04-19	57	F	154.2	86.4	Dancing

2000 rows × 22 columns

(ii). Data Cleaning

To ensure quality and consistency, the following cleaning operations were performed:

- Null Value Handling: Missing values were identified and treated either by imputation (e.g., using mean/median) or removal depending on the severity and context.
- Duplicate Removal: Repeated entries were removed to reduce bias and redundancy.
- Data Type Corrections: Ensured each column had the correct data type (e.g., dates as datetime, categorical as strings).
- Inconsistent Labels: Standardized naming conventions and label formats to ensure uniformity.

✓ (iii). Data Pre-processing

To prepare the dataset for analysis and modeling, preprocessing steps included:

- Encoding: Converted categorical variables into numerical format using label encoding or one-hot encoding.
- Normalization: Applied Min-Max or Z-score normalization to scale numerical features.
- Splitting: Divided the dataset into training and testing subsets (typically 80/20 or 70/30 split).
- Feature Engineering: Created new features from existing ones (e.g., aggregations, binning, derived metrics).
- Outlier Handling: Used statistical techniques (e.g., IQR, Z-score) to detect and optionally filter out outliers.

```
df['date'] = pd.to_datetime(df['date'])
df.set_index('date', inplace=True)
```

```
df.drop_duplicates(inplace=True)
```

```
df.sort_index(inplace=True)
```

```
# Fill missing 'health_condition' based on mode of similar groups
def fill_health_condition(row):
    if pd.isna(row['health_condition']):
        group = df[
            (df['smoking_status'] == row['smoking_status']) &
```

```

        (df['stress_level'] == row['stress_level']) &
        (df['bmi'].between(row['bmi'] - 1, row['bmi'] + 1)) &
        (df['age'].between(row['age'] - 3, row['age'] + 3))
    ]
    if not group['health_condition'].dropna().empty:
        return group['health_condition'].dropna().mode()[0]
    else:
        # fallback to overall mode if group too small
        return df['health_condition'].dropna().mode()[0]
else:
    return row['health_condition']

```

```
df['health_condition'] = df.apply(fill_health_condition, axis=1)
```

```
print("Missing values after cleaning:", df['health_condition'].isna().sum())
```

```
df
```

Missing values after cleaning: 0

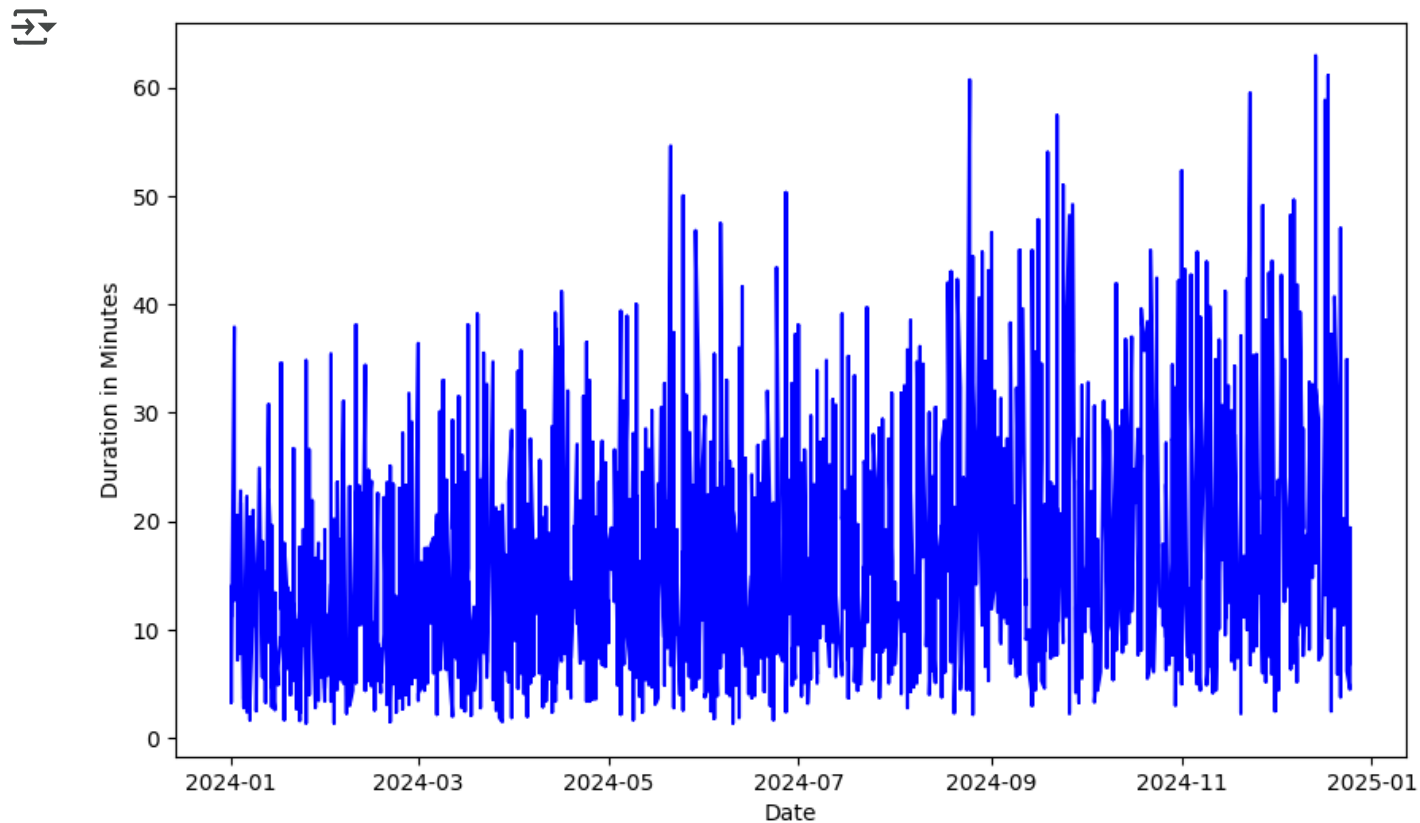
	participant_id	age	gender	height_cm	weight_kg	activity_type	duration
date							
2024-01-01	490	58	M	185.4	85.2	Running	
2024-01-01	139	20	F	154.5	49.3	Tennis	
2024-01-01	645	63	F	172.9	66.9	Weight Training	
2024-01-01	1331	39	F	166.8	57.8	HIIT	
2024-01-01	2483	60	M	174.6	93.5	Cycling	
...
2024-12-24	2899	36	M	179.2	157.1	Walking	
2024-12-25	180	50	F	160.1	116.2	Yoga	
2024-12-25	2792	45	F	165.5	122.5	Cycling	
2024-12-25	1166	37	M	165.6	130.7	Weight Training	
2024-12-25	1669	36	M	171.1	132.9	Walking	

2000 rows x 21 columns

✓

b. Data Exploration & Model Insights

```
plt.xlabel('Date')
plt.ylabel('Duration in Minutes')
plt.plot(df['calories_burned'], color='blue')
plt.show()
```



✓ Rolling Statistic

```
# Determine rolling statistics
```

```
rolmean = df['calories_burned'].rolling(window=365).mean()  
rolstd = df['calories_burned'].rolling(window=365).std()
```

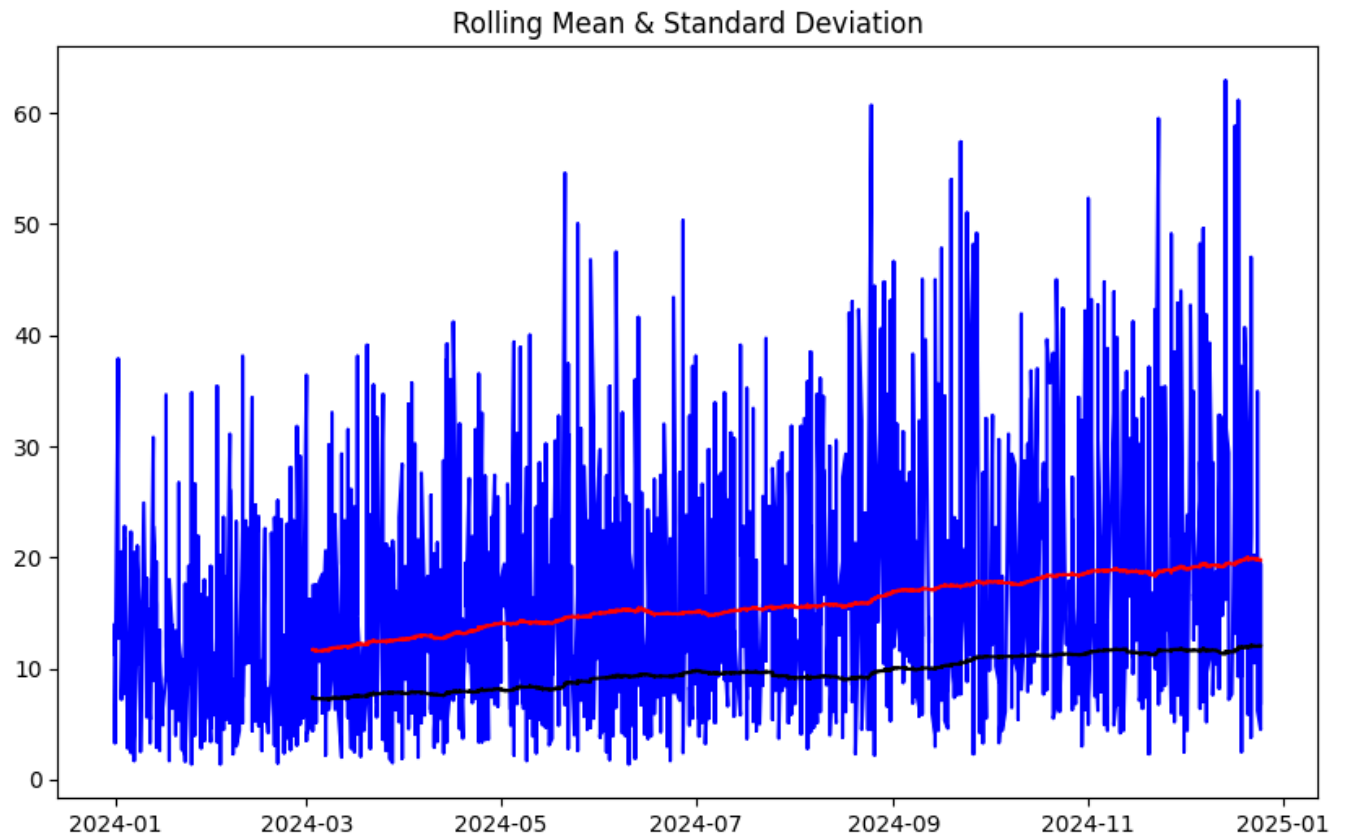
```
print(rolmean, rolstd)
```

```
↔ date  
2024-01-01      NaN  
2024-01-01      NaN  
2024-01-01      NaN  
2024-01-01      NaN  
2024-01-01      NaN  
  
      ...  
2024-12-24    19.769863  
2024-12-25    19.760000  
2024-12-25    19.776986  
2024-12-25    19.738082  
2024-12-25    19.711233  
Name: calories_burned, Length: 2000, dtype: float64 date  
2024-01-01      NaN  
2024-01-01      NaN  
2024-01-01      NaN  
2024-01-01      NaN  
2024-01-01      NaN  
  
      ...  
2024-12-24    11.997775  
2024-12-25    12.008870  
2024-12-25    12.003949  
2024-12-25    12.008360  
2024-12-25    12.026114  
Name: calories_burned, Length: 2000, dtype: float64
```

```
# Plot rolling statsitics
```

```
orig = plt.plot(df['calories_burned'], color='blue', label='Original')  
mean = plt.plot(rolmean, color='red', label='Rolling Mean')  
std = plt.plot(rolstd, color='black', label='Rolling Std')  
plt.title('Rolling Mean & Standard Deviation')
```

```
➡ Text(0.5, 1.0, 'Rolling Mean & Standard Deviation')
```



✓ Dickey-Fuller Test


```


# Perform Dickey-Fuller test
from statsmodels.tsa.stattools import adfuller

print('Results of Dickey-Fuller Test:')
dfctest = adfuller(df['calories_burned'], autolag='AIC')

dfoutput = pd.Series(dfctest[0:4], index=['Test Statistic', 'p-value', '#Lags Us
for key, value in dfctest[4].items():
    dfoutput['Critical Value (%s)' % key] = value

print(dfoutput)

```

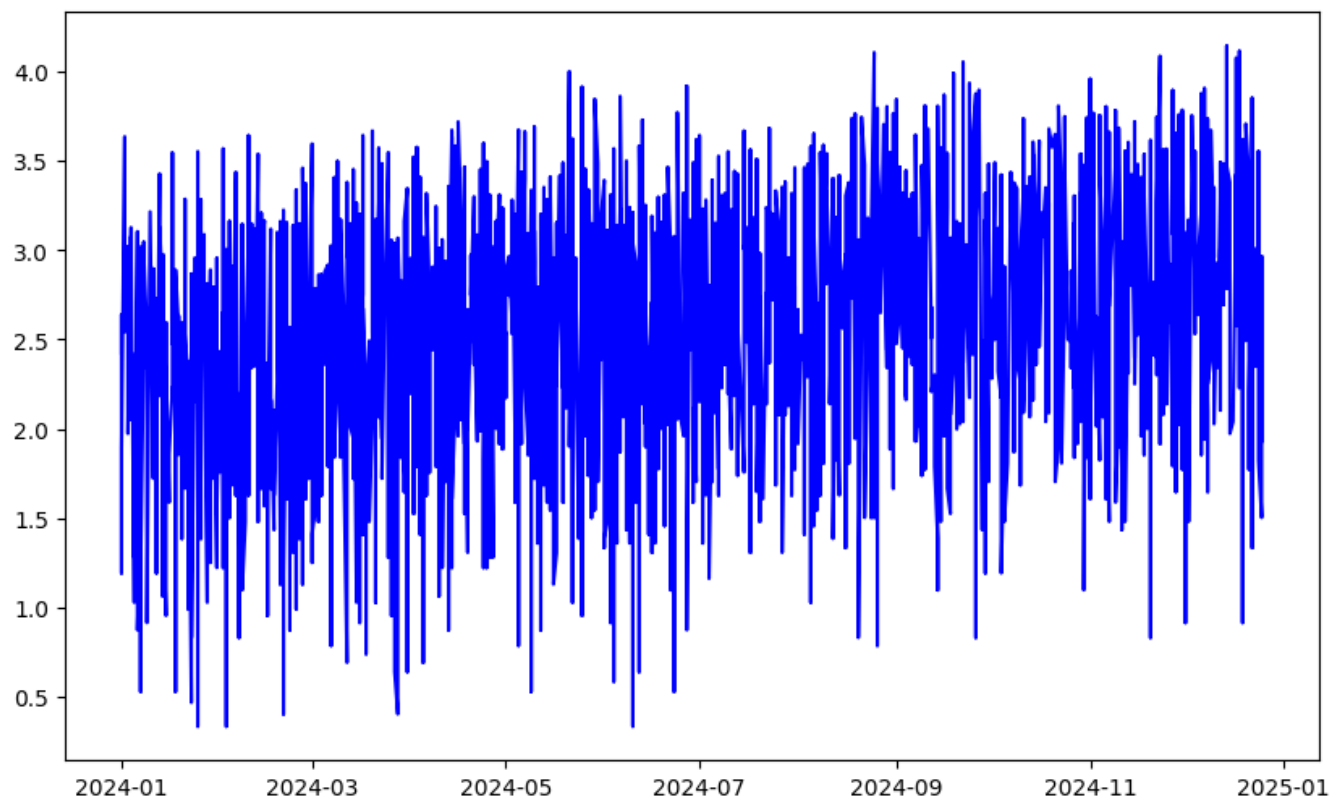
 Results of Dickey-Fuller Test:

Test Statistic	-5.193625
p-value	0.000009
#Lags Used	23.000000
Number of Observations Used	1976.000000
Critical Value (1%)	-3.433664
Critical Value (5%)	-2.863004
Critical Value (10%)	-2.567549
dtype:	float64

Estimating Trend

```
# Estimating Trend
df_logScale = np.log(df['calories_burned'])
plt.plot(df_logScale, color='blue')
```

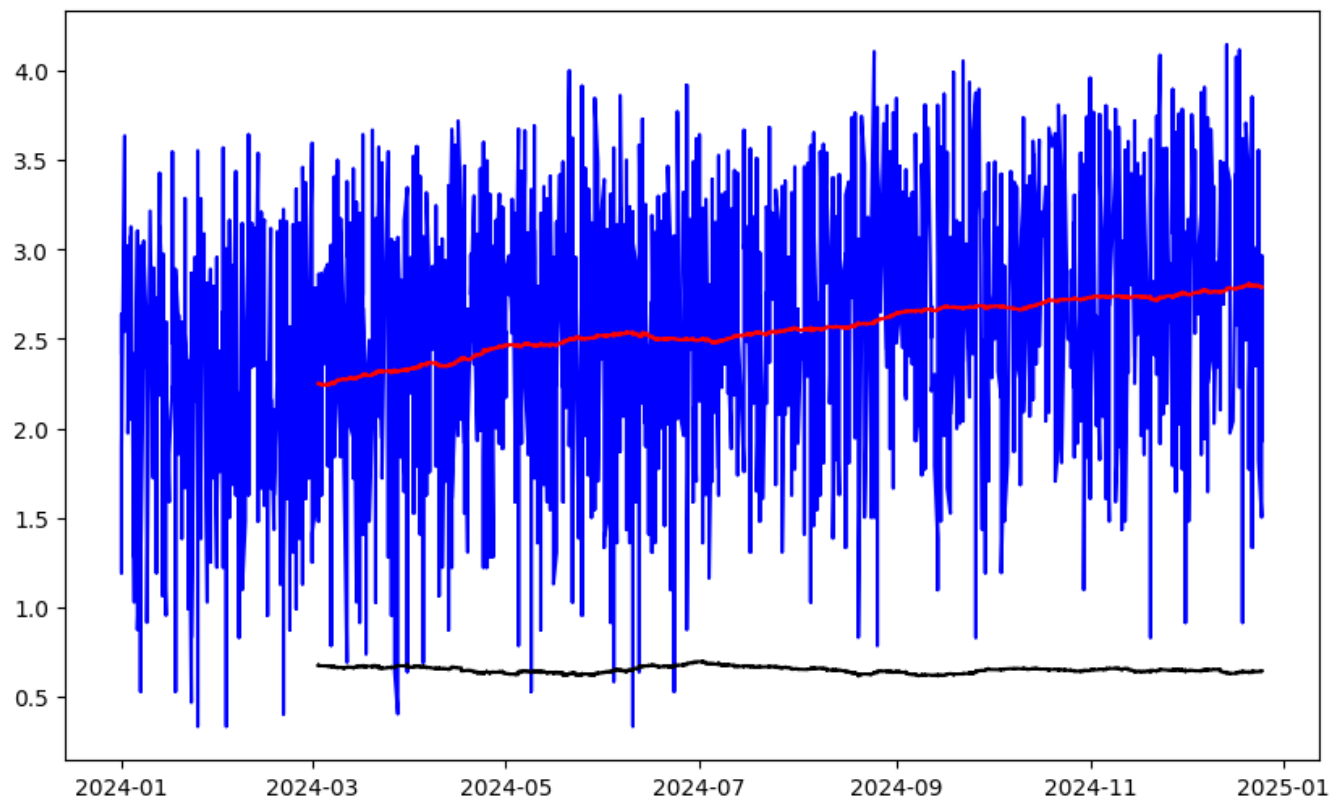
↔ [[matplotlib.lines.Line2D](#) at 0x7d65a02ce810>]



✓ Moving Average


```
movingAverage = df_logScale.rolling(window=365).mean()  
movingSTD = df_logScale.rolling(window=365).std()  
plt.plot(df_logScale, color='blue', label='Original')  
plt.plot(movingAverage, color='red', label='Moving Average')  
plt.plot(movingSTD, color='black', label='moving STD')
```

↔ [[matplotlib.lines.Line2D](#) at 0x7d659fbd3390>]



```
df_logScaleMinusMovingAverage = df_logScale - movingAverage
df_logScaleMinusMovingAverage
```

```
# Remove Nan values
df_logScaleMinusMovingAverage.dropna(inplace=True)
df_logScaleMinusMovingAverage
```



calories_burned	
date	
2024-03-03	0.327231
2024-03-03	-0.724467
2024-03-03	0.545879
2024-03-03	-0.366231
2024-03-03	-0.768525
...	...
2024-12-24	-0.968620
2024-12-25	-1.287481
2024-12-25	0.172659
2024-12-25	-0.365581
2024-12-25	-0.856441

1636 rows × 1 columns

dtype: float64

✓ **Test Stationarity**

```

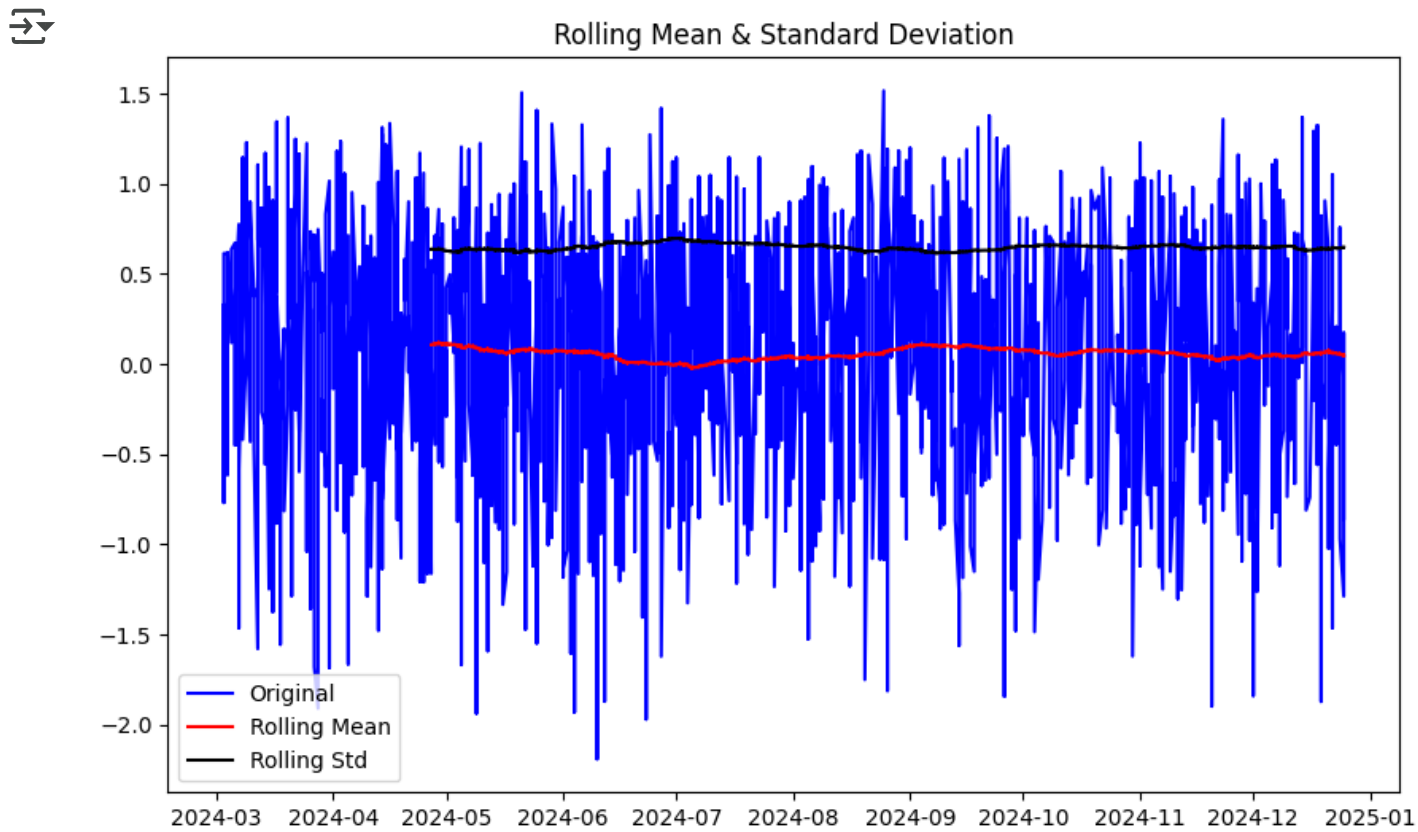
from statsmodels.tsa.stattools import adfuller

def test_stationarity(timeseries):
    # Determine rolling statistics
    movingAverage = timeseries.rolling(window=365).mean()
    movingSTD = timeseries.rolling(window=365).std()
    # Plot rolling statistics
    orig = plt.plot(timeseries, color='blue', label='Original')
    mean = plt.plot(movingAverage, color='red', label='Rolling Mean')
    std = plt.plot(movingSTD, color='black', label='Rolling Std')
    plt.legend(loc='best')
    plt.title('Rolling Mean & Standard Deviation')
    plt.show(block=False)

    # Perform Dickey-Fuller test
    print('Results of Dickey-Fuller Test:')
    dfctest = adfuller(timeseries, autolag='AIC')
    dfoutput = pd.Series(dfctest[0:4], index=['Test Statistic', 'p-value', '#Lag
for key, value in dfctest[4].items():
        dfoutput['Critical Value (%s)' % key] = value
    print(dfoutput)

```

```
test_stationarity(df_logScaleMinusMovingAverage)
```



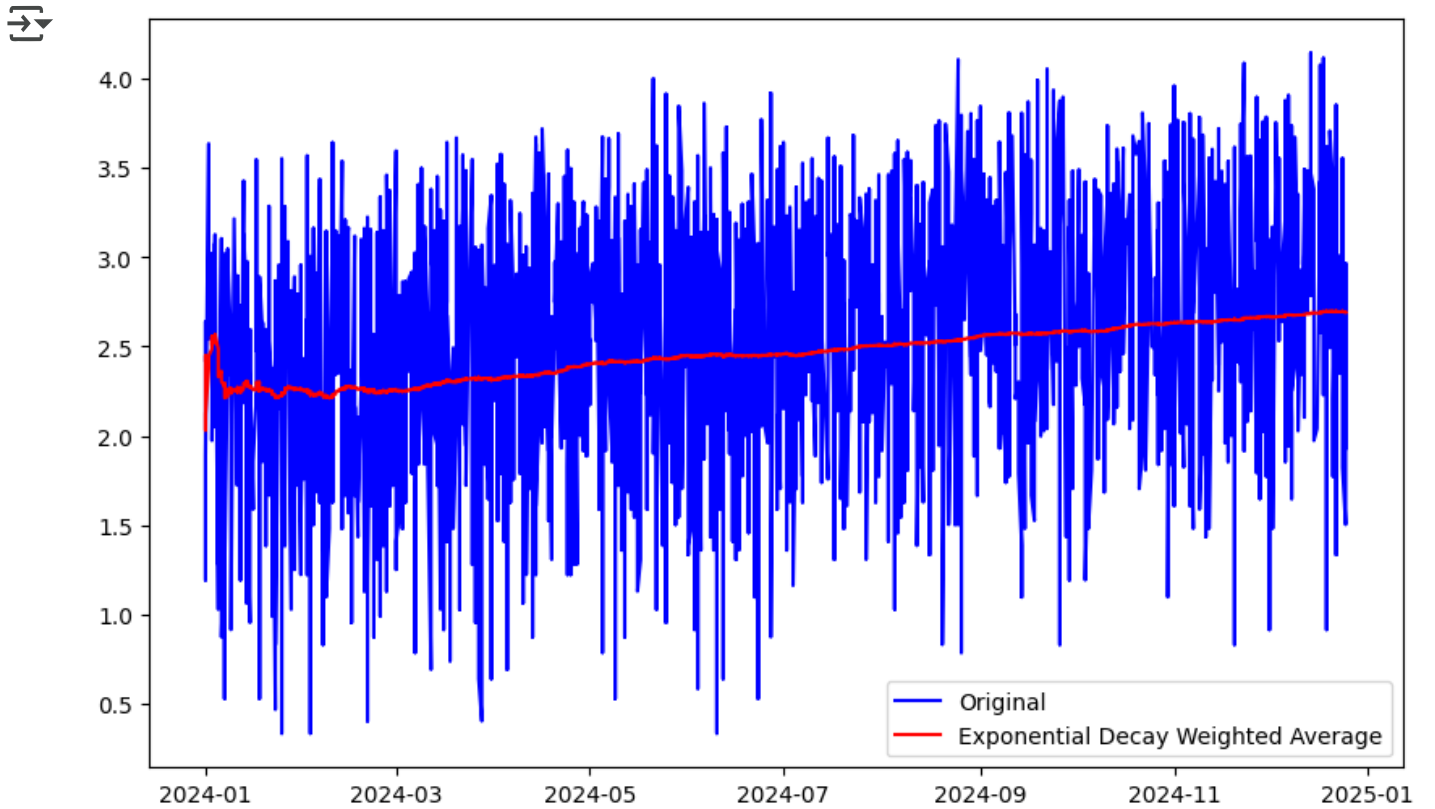
Results of Dickey-Fuller Test:

Test Statistic	-21.832195
p-value	0.000000
#Lags Used	3.000000
Number of Observations Used	1632.000000
Critical Value (1%)	-3.434363
Critical Value (5%)	-2.863313
Critical Value (10%)	-2.567714

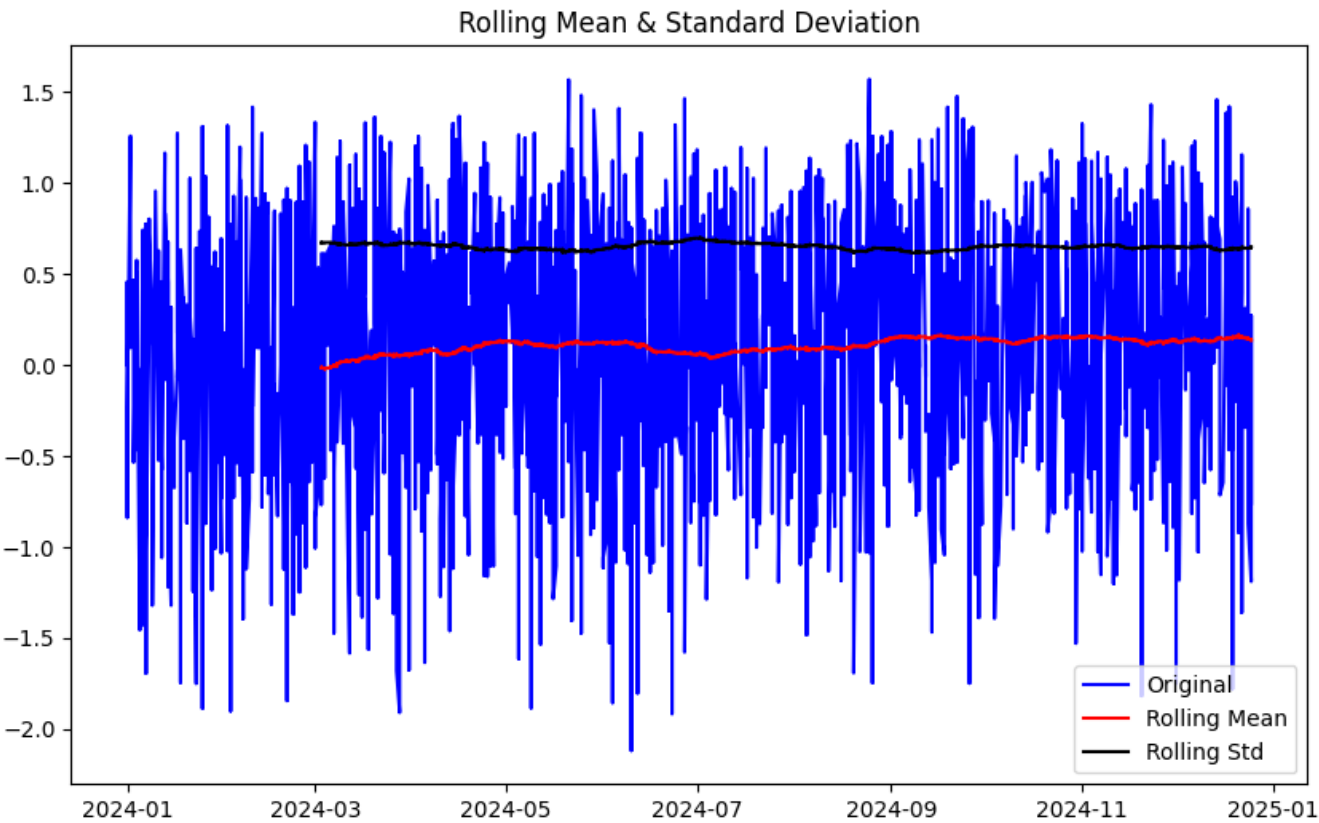
dtype: float64

Weighted Average

```
exponentialDecayWeightedAverage = df_logScale.ewm(halflife=365, min_periods=0,  
plt.plot(df_logScale, color='blue', label='Original')  
plt.plot(exponentialDecayWeightedAverage, color='red', label='Exponential Decay  
plt.legend(loc='best')  
plt.show()
```



```
df_logScaleMinusExponentialMovingAverage = df_logScale - exponentialDecayWeight
test_stationarity(df_logScaleMinusExponentialMovingAverage)
```

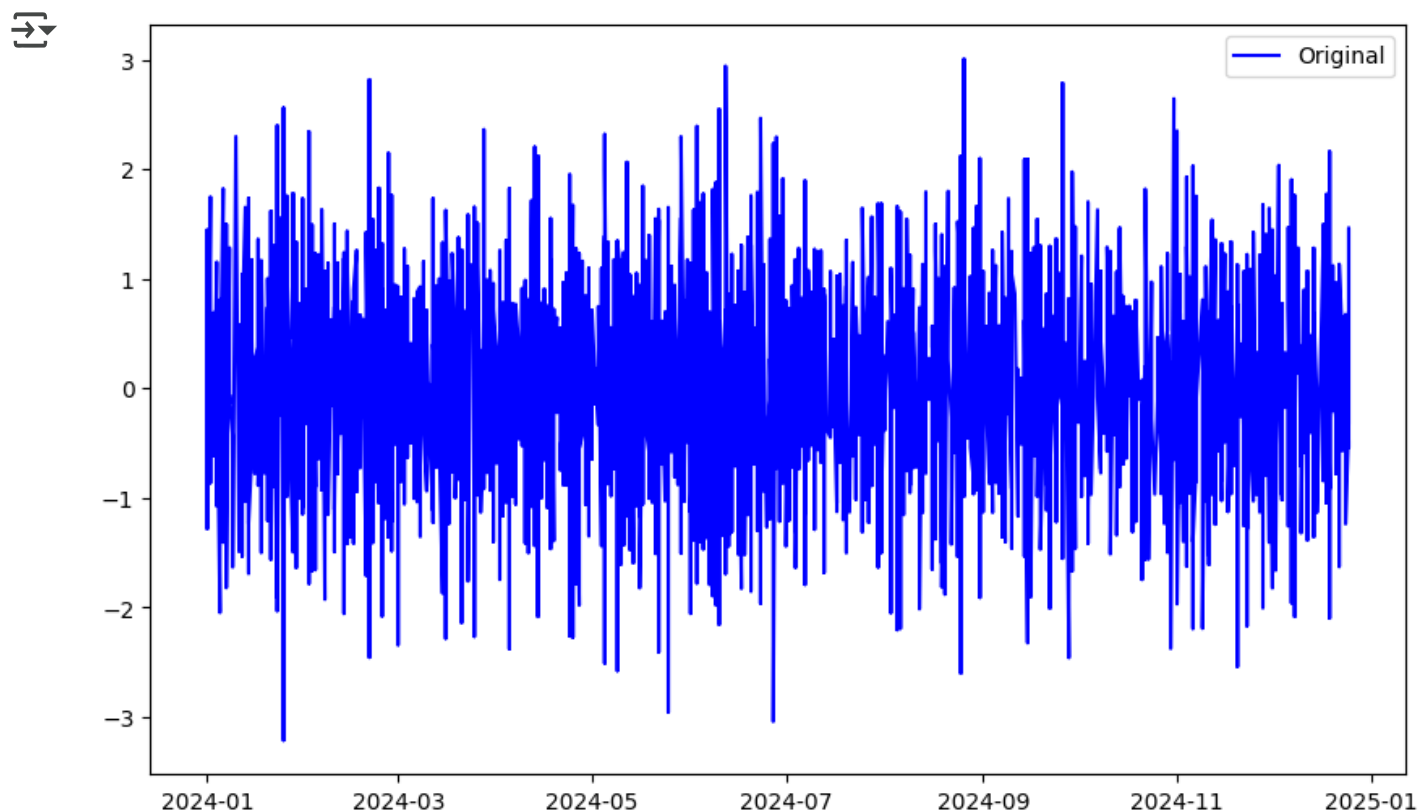


Results of Dickey-Fuller Test:

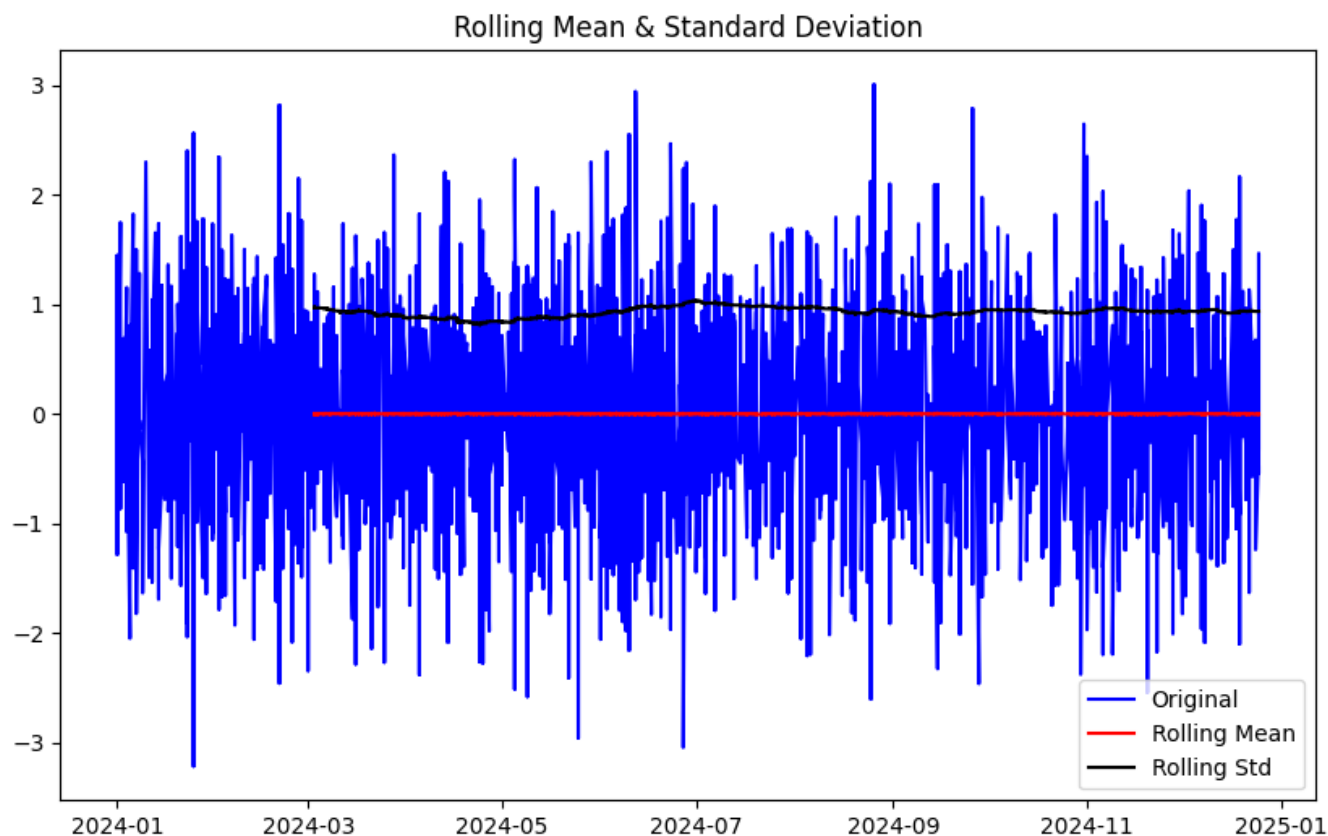
Test Statistic	-23.432456
p-value	0.000000
#Lags Used	3.000000
Number of Observations Used	1996.000000
Critical Value (1%)	-3.433630
Critical Value (5%)	-2.862989
Critical Value (10%)	-2.567541

dtype: float64


```
df_logScaleDiffShifting = df_logScale - df_logScale.shift()
plt.plot(df_logScaleDiffShifting, color='blue', label='Original')
plt.legend(loc='best')
plt.show()
```



```
df_logScaleDiffShifting.dropna(inplace=True)
test_stationarity(df_logScaleDiffShifting)
```



Results of Dickey-Fuller Test:

Test Statistic	-1.580950e+01
p-value	1.050802e-28
#Lags Used	2.400000e+01
Number of Observations Used	1.974000e+03
Critical Value (1%)	-3.433667e+00
Critical Value (5%)	-2.863005e+00
Critical Value (10%)	-2.567550e+00
dtype:	float64

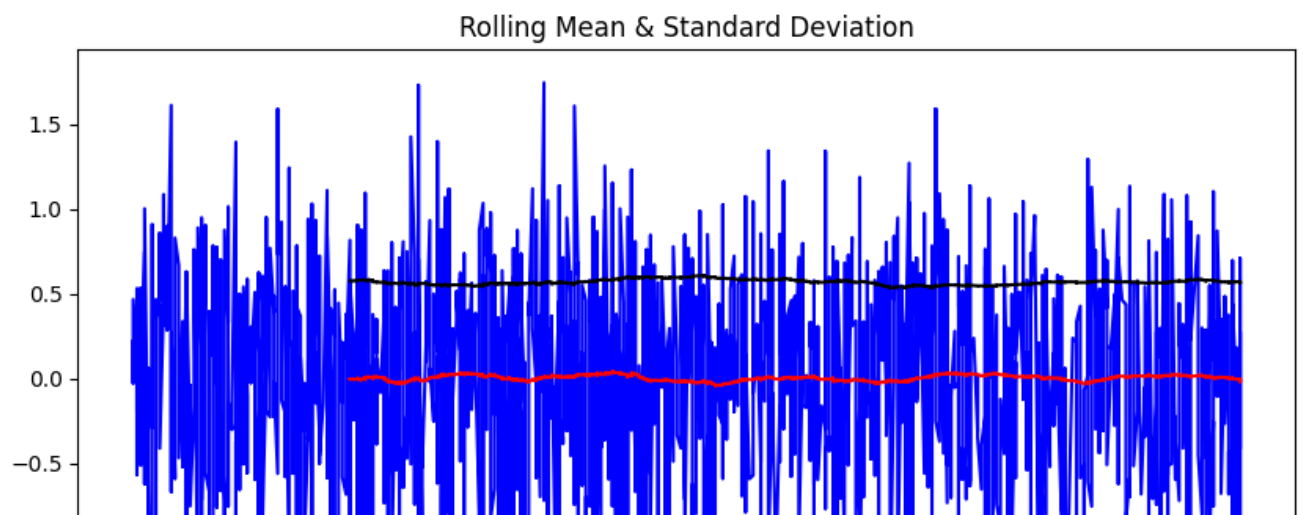
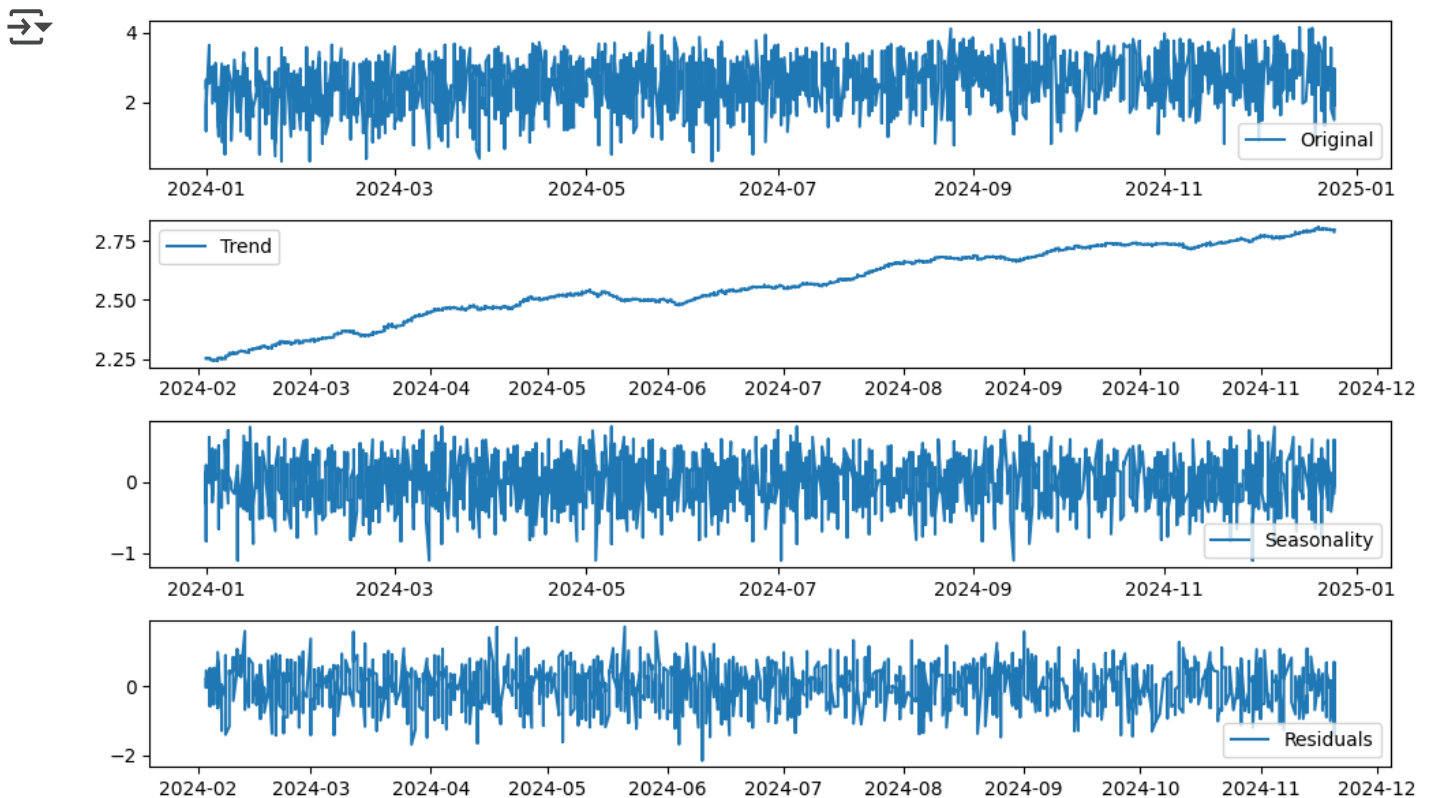
```
from statsmodels.tsa.seasonal import seasonal_decompose
decomposition = seasonal_decompose(df_logScale, model='additive', period=365)
```

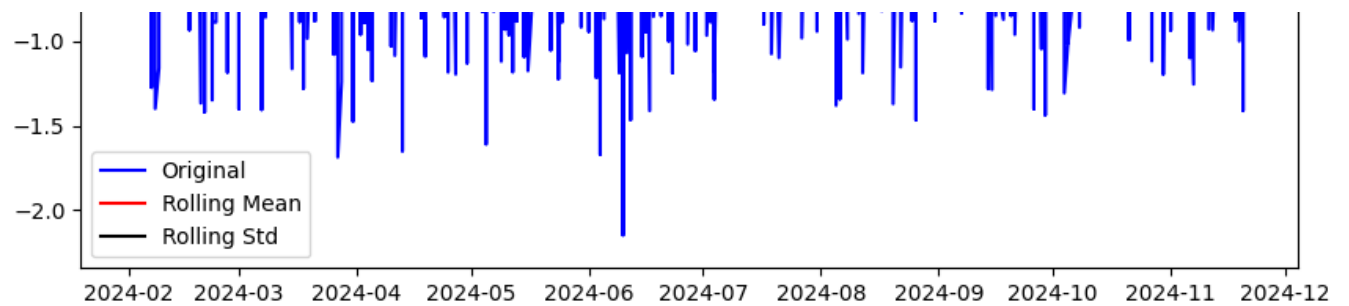
```
trend = decomposition.trend
seasonal = decomposition.seasonal
residual = decomposition.resid
```

```
plt.subplot(411)
plt.plot(df_logScale, label='Original')
```

```
plt.legend(loc='best')
plt.subplot(412)
plt.plot(trend, label='Trend')
plt.legend(loc='best')
plt.subplot(413)
plt.plot(seasonal, label='Seasonality')
plt.legend(loc='best')
plt.subplot(414)
plt.plot(residual, label='Residuals')
plt.legend(loc='best')
plt.tight_layout()
plt.show()
```

```
decomposedLogData = residual
decomposedLogData.dropna(inplace=True)
test_stationarity(decomposedLogData)
```





Results of Dickey-Fuller Test:

Test Statistic	-21.069386
p-value	0.000000
#Lags Used	3.000000
Number of Observations Used	1632.000000
Critical Value (1%)	-3.434363
Critical Value (5%)	-2.863313
Critical Value (10%)	-2.567714
dtype:	float64

```
decomposedLogData = residual
decomposedLogData.dropna(inplace=True)
#test_stationarity(decomposedLogData)
```

✓ ACF and PACF

```
from statsmodels.tsa.stattools import acf, pacf

lag_acf = acf(df_logScaleDiffShifting, nlags=20)
lag_pacf = pacf(df_logScaleDiffShifting, nlags=20, method='ols')
```

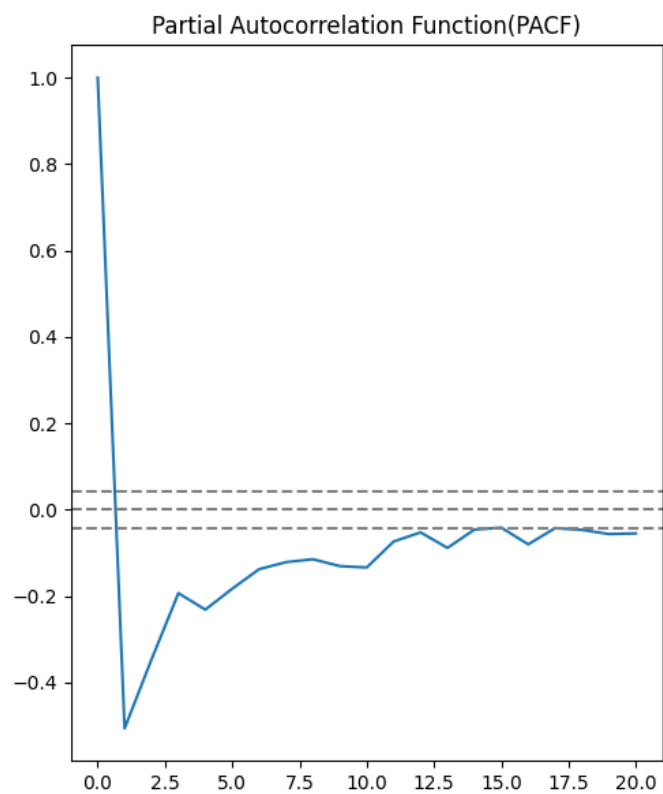
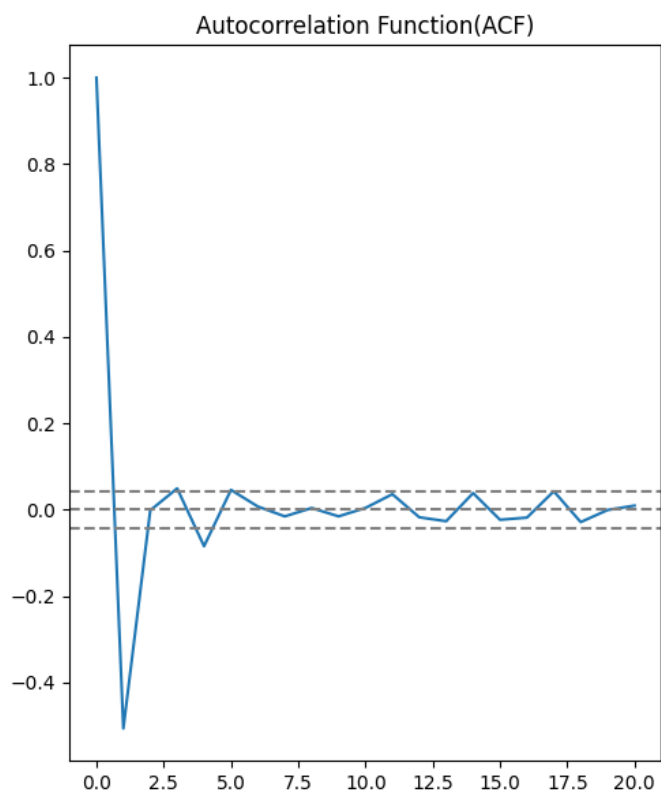
```

#plot ACF
plt.subplot(121)
plt.plot(lag_acf)
plt.axhline(y=0, linestyle='--', color='gray')
plt.axhline(y=-1.96/np.sqrt(len(df_logScaleDiffShifting)), linestyle='--', color='gray')
plt.axhline(y=1.96/np.sqrt(len(df_logScaleDiffShifting)), linestyle='--', color='gray')
plt.title('Autocorrelation Function(ACF)')

# Plot PACF
plt.subplot(122)
plt.plot(lag_pacf)
plt.axhline(y=0, linestyle='--', color='gray')
plt.axhline(y=-1.96/np.sqrt(len(df_logScaleDiffShifting)), linestyle='--', color='gray')
plt.axhline(y=1.96/np.sqrt(len(df_logScaleDiffShifting)), linestyle='--', color='gray')
plt.title('Partial Autocorrelation Function(PACF)')

plt.tight_layout()
plt.show()

```



```

from statsmodels.tsa.arima.model import ARIMA

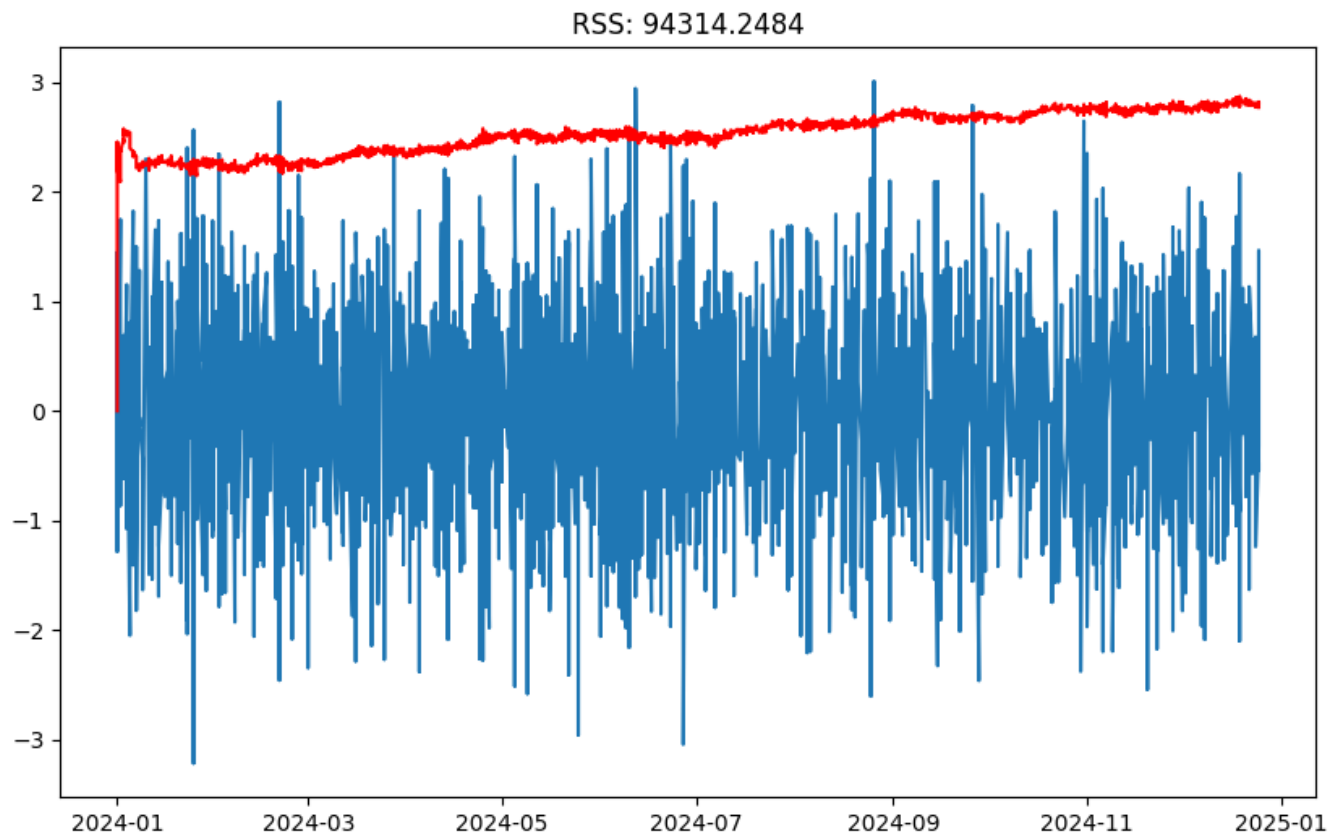
# AR model
model = ARIMA(df_logScale, order=(2, 1, 2))
results_ARIMA = model.fit()
plt.plot(df_logScaleDiffShifting)
plt.plot(results_ARIMA.fittedvalues, color='red')
plt.title('RSS: %.4f' % sum((results_ARIMA.fittedvalues - df_logScaleDiffShifting)**2))
plt.show()
print('Plotting ARIMA model')

```

```

➡ /usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
   self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
   self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
   self._init_dates(dates, freq)

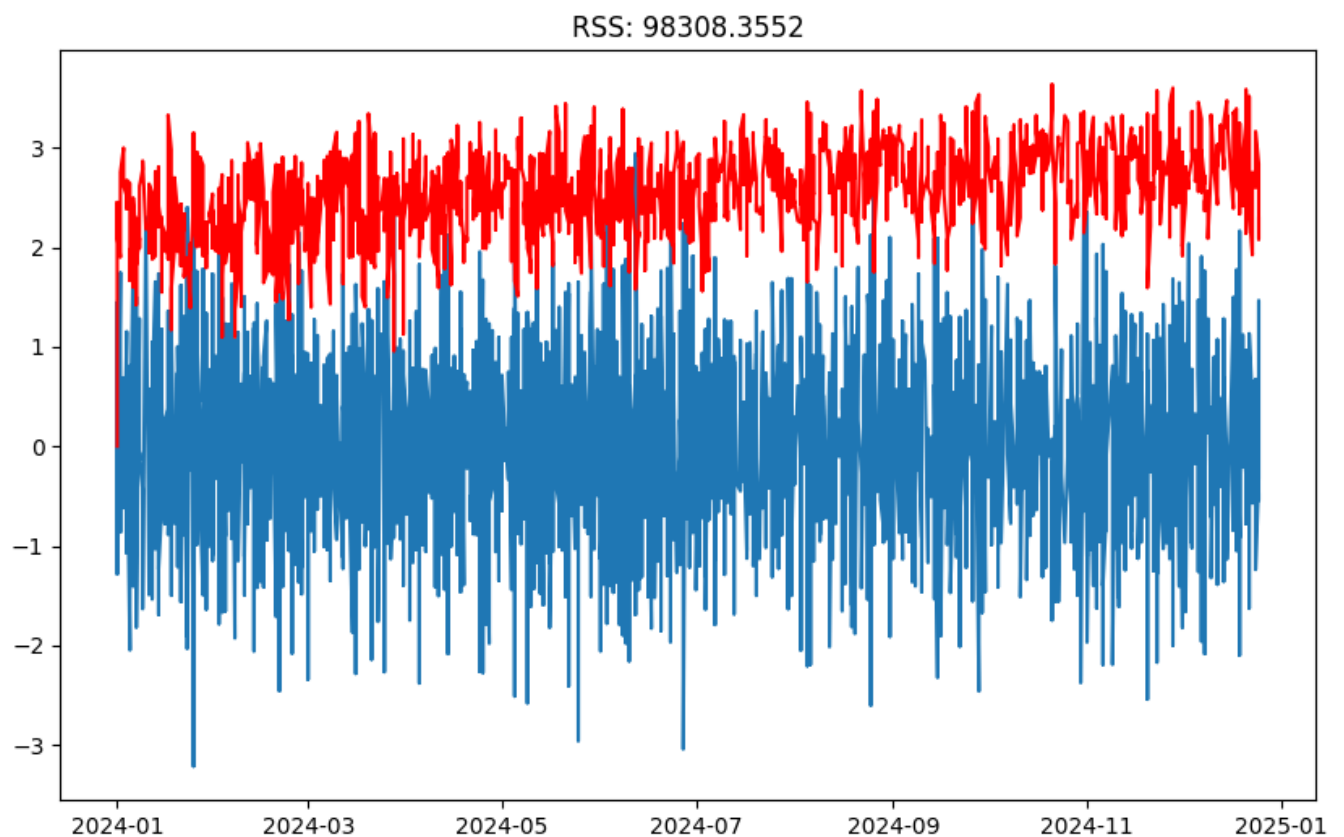
```



Plotting ARIMA model

```
# MA Model
model = ARIMA(df_logScale, order=(2, 1, 0))
results_MA = model.fit()
plt.plot(df_logScaleDiffShifting)
plt.plot(results_MA.fittedvalues, color='red')
plt.title('RSS: %.4f' % sum((results_MA.fittedvalues - df_logScaleDiffShifting)
plt.show()
print('Plotting AR model')
```

```
➡ /usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
    self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
    self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
    self._init_dates(dates, freq)
```



Plotting AR model


```

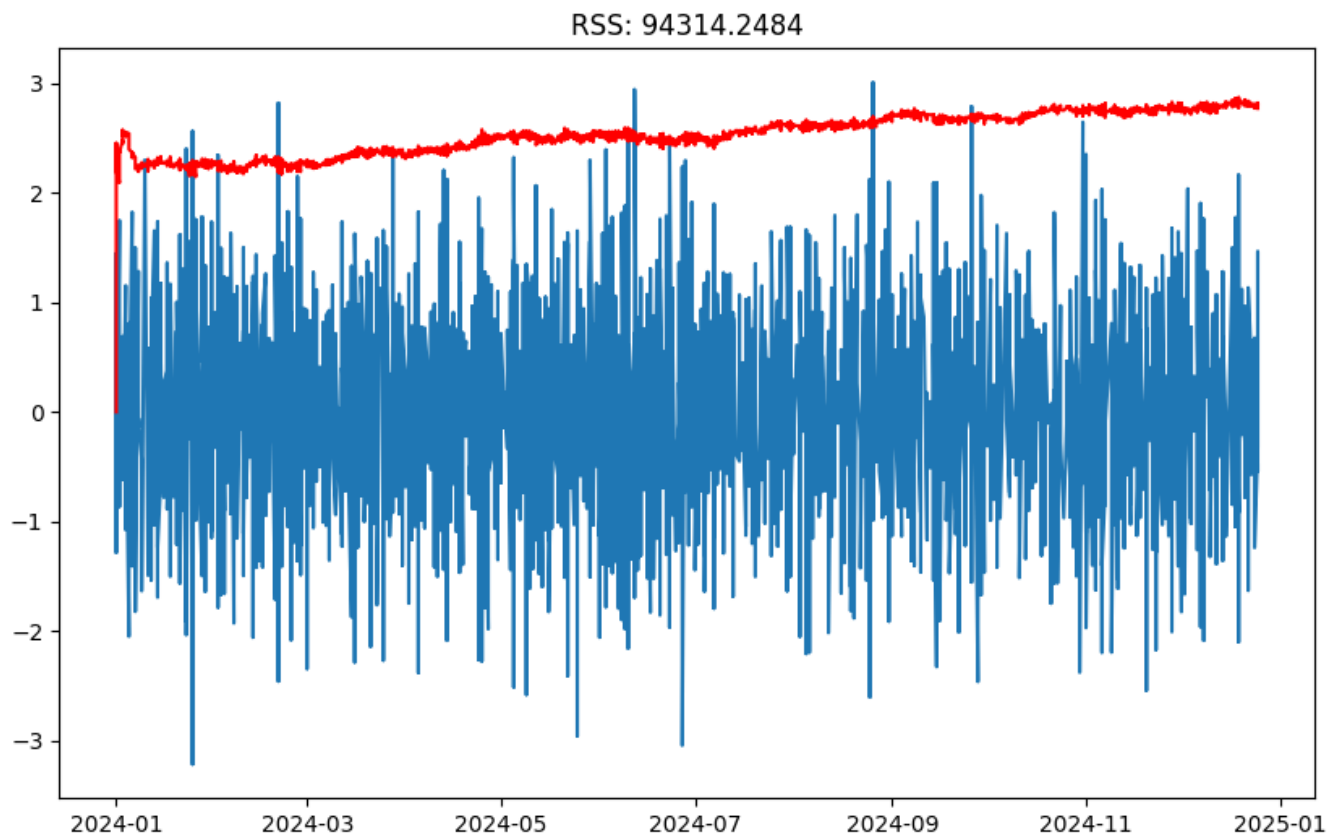
model = ARIMA(df_logScale, order=(2, 1, 2))
results_ARIMA = model.fit()
plt.plot(df_logScaleDiffShifting)
plt.plot(results_ARIMA.fittedvalues, color='red')
plt.title('RSS: %.4f' % sum((results_ARIMA.fittedvalues - df_logScaleDiffShifting)**2))
plt.show()

```

```

➡ /usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
   self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
   self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
   self._init_dates(dates, freq)

```



```
predictions_ARIMA_diff = pd.Series(results_ARIMA.fittedvalues, copy=True)
print(predictions_ARIMA_diff.head())
```

```
↔ date
2024-01-01    0.000000
2024-01-01    2.424802
2024-01-01    2.450905
2024-01-01    2.038021
2024-01-01    2.196306
dtype: float64
```

```
# Convert to cumulative sum
predictions_ARIMA_diff_cumsum = predictions_ARIMA_diff.cumsum()
print(predictions_ARIMA_diff_cumsum.head())
```

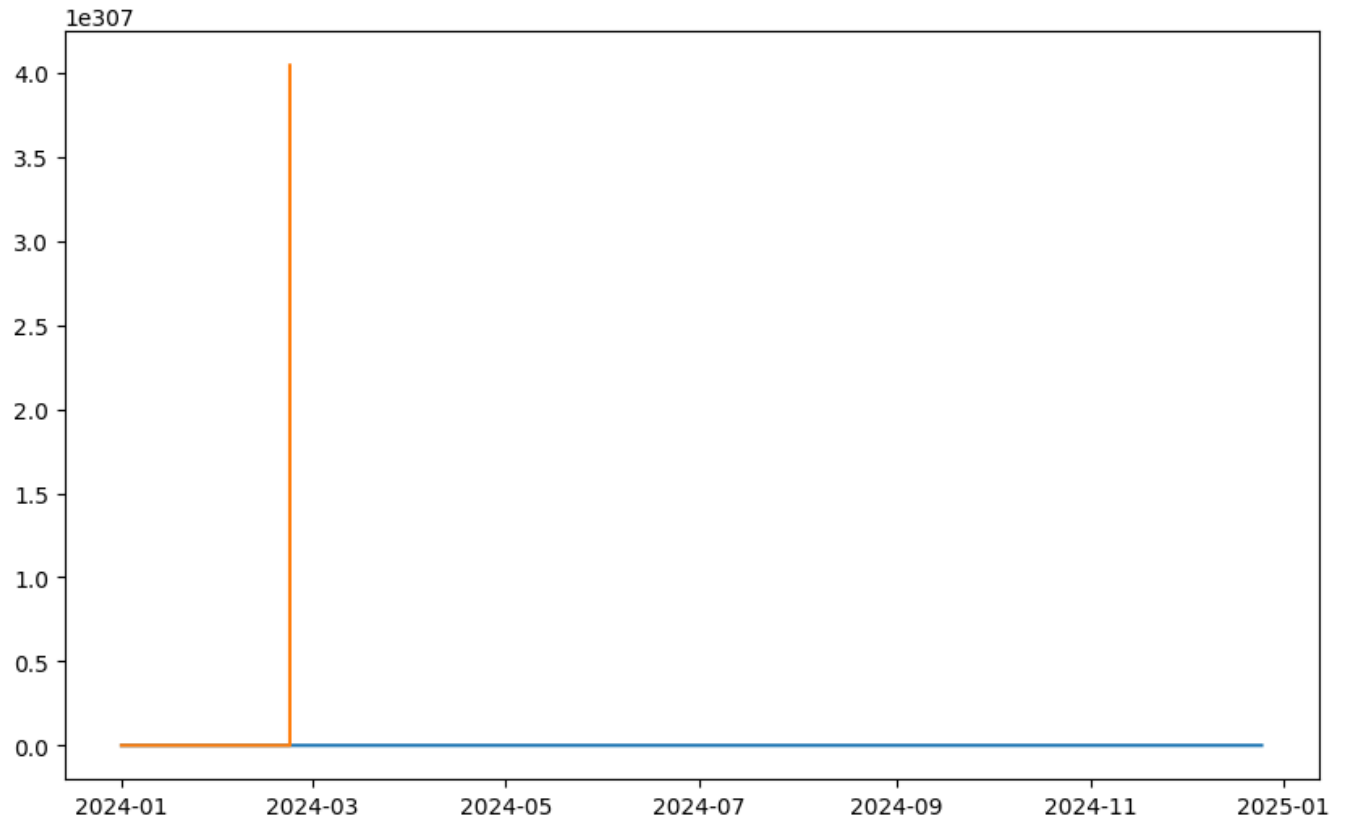
```
↔ date
2024-01-01    0.000000
2024-01-01    2.424802
2024-01-01    4.875706
2024-01-01    6.913728
2024-01-01    9.110034
dtype: float64
```

```
predictions_ARIMA_log = pd.Series(df_logScale.iloc[0], index=predictions_ARIMA_
predictions_ARIMA_log = predictions_ARIMA_log.add(predictions_ARIMA_diff_cumsum
print(predictions_ARIMA_log.head())
```


```
↔ date
2024-01-01    2.424803
2024-01-01    4.849604
2024-01-01    7.300509
2024-01-01    9.338530
2024-01-01   11.534837
dtype: float64
```

```
predictions_ARIMA = np.exp(predictions_ARIMA_log)
plt.plot(df['calories_burned'])
plt.plot(predictions_ARIMA)
```

```
➡ /usr/local/lib/python3.11/dist-packages/pandas/core/arraylike.py:399: RuntimeWarning:
  result = getattr(ufunc, method)(*inputs, **kwargs)
[<matplotlib.lines.Line2D at 0x7d659ed4d890>]
```



df_logScale




calories_burned	
date	
<hr/>	
2024-01-01	2.424803
2024-01-01	2.476538
2024-01-01	1.193922
2024-01-01	2.639057
2024-01-01	1.887070
...	...
2024-12-24	1.824549
2024-12-25	1.504077
2024-12-25	2.965273
2024-12-25	2.424803
2024-12-25	1.931521

2000 rows x 1 columns

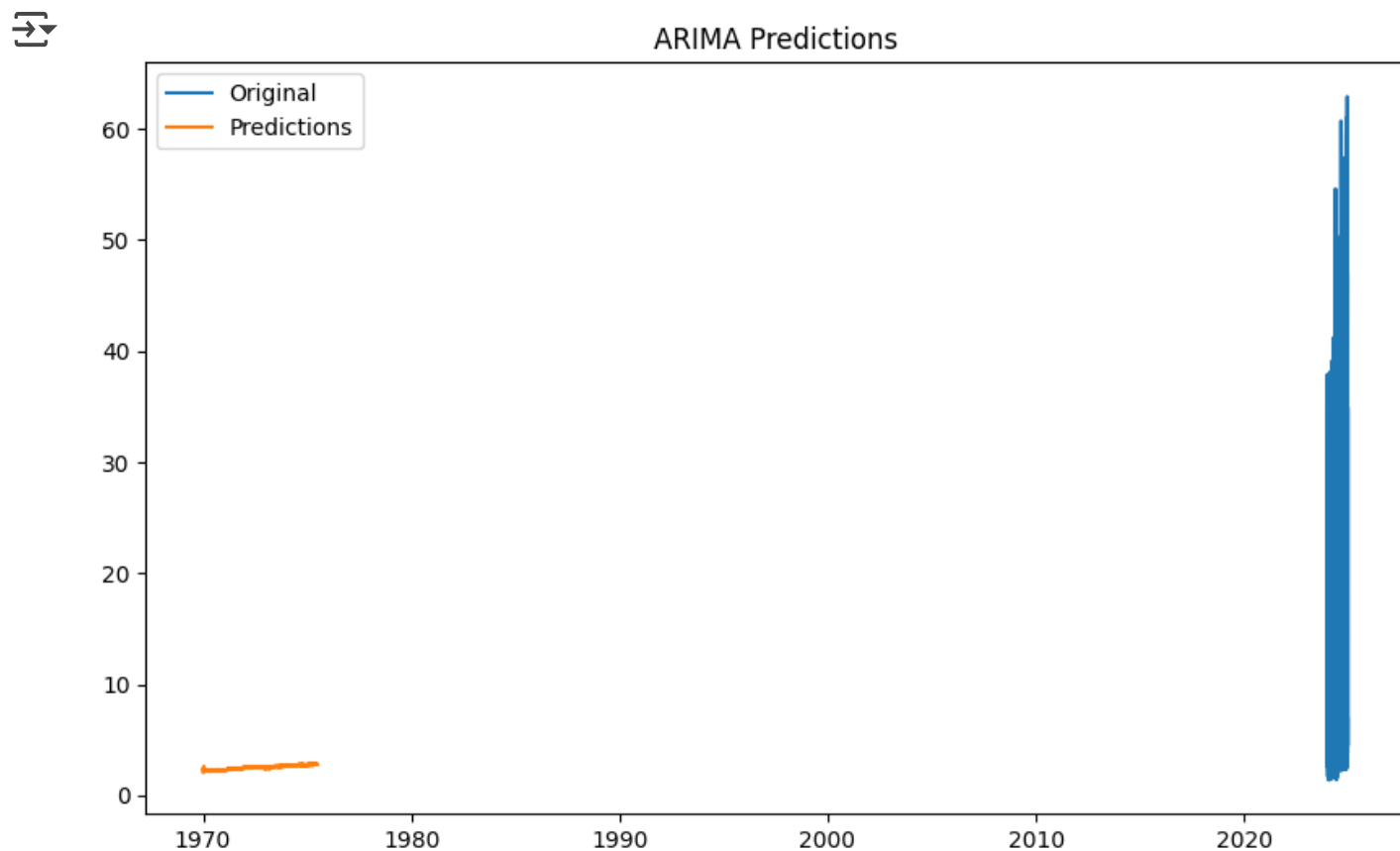
dtype: float64

```
predictions = results_ARIMA.get_prediction(start=1, end=2000)
```



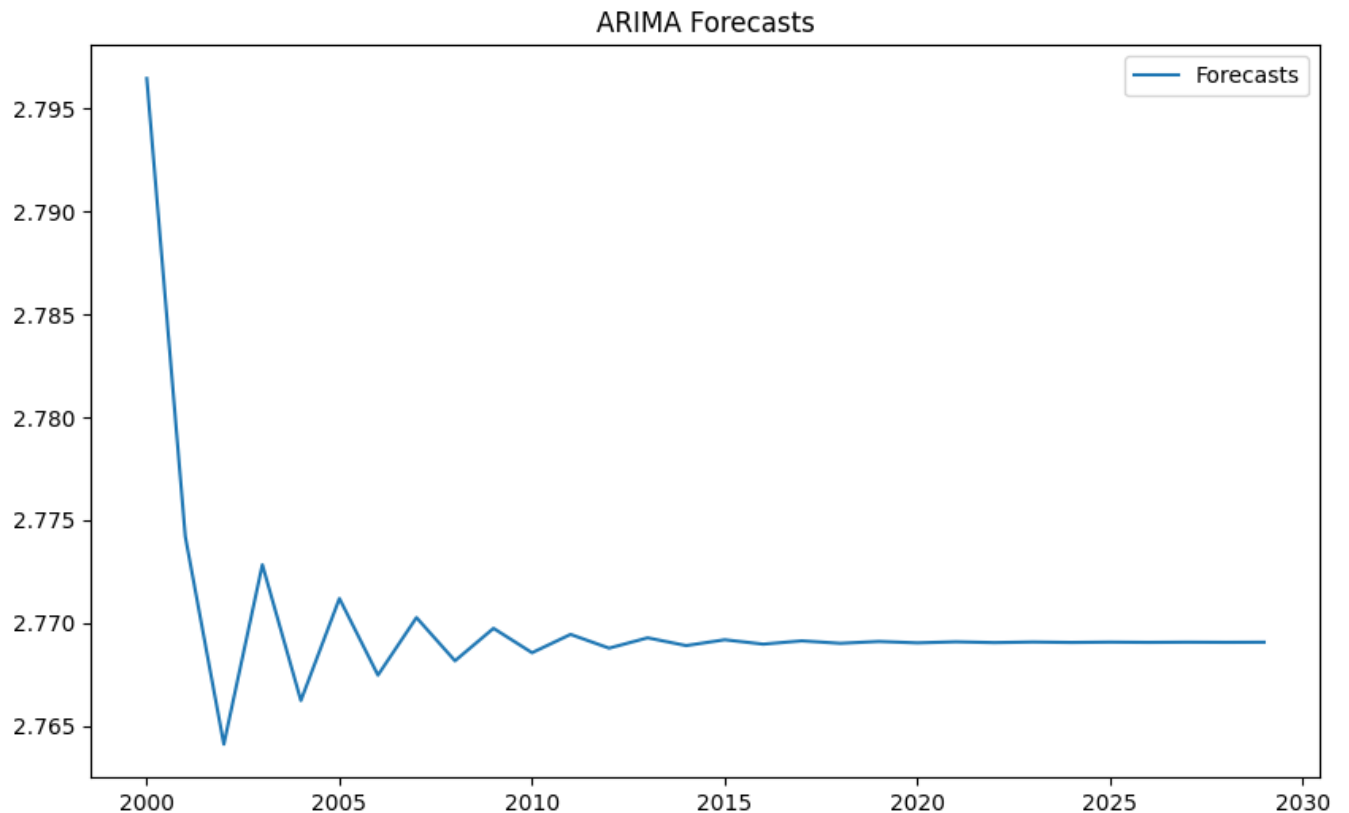
```
    /usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:8
    return get_prediction_index(
    /usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:8
    return get_prediction_index(
```

```
plt.plot(df['calories_burned'], label='Original')
plt.plot(predictions.predicted_mean, label='Predictions')
plt.legend()
plt.title('ARIMA Predictions')
plt.show()
```



```
x = results_ARIMA.get_forecast(steps=30)
plt.plot(x.predicted_mean, label='Forecasts')
plt.legend()
plt.title('ARIMA Forecasts')
plt.show()
```

```
➡ /usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:8
    return get_prediction_index(
```



```
from statsmodels.tsa.arima.model import ARIMA
from statsmodels.tsa.stattools import adfuller
import matplotlib.pyplot as plt
import seaborn as sns

# Sort the dates
df = df.sort_values(by='date')

# Aggregate by date (mean daily steps)
daily_steps = df.groupby('date')['daily_steps'].mean()
```

```
# Fit ARIMA model (auto ARIMA selection can be used, but here we use (1,1,1) as
model = ARIMA(daily_steps, order=(1, 1, 1))
model_fit = model.fit()

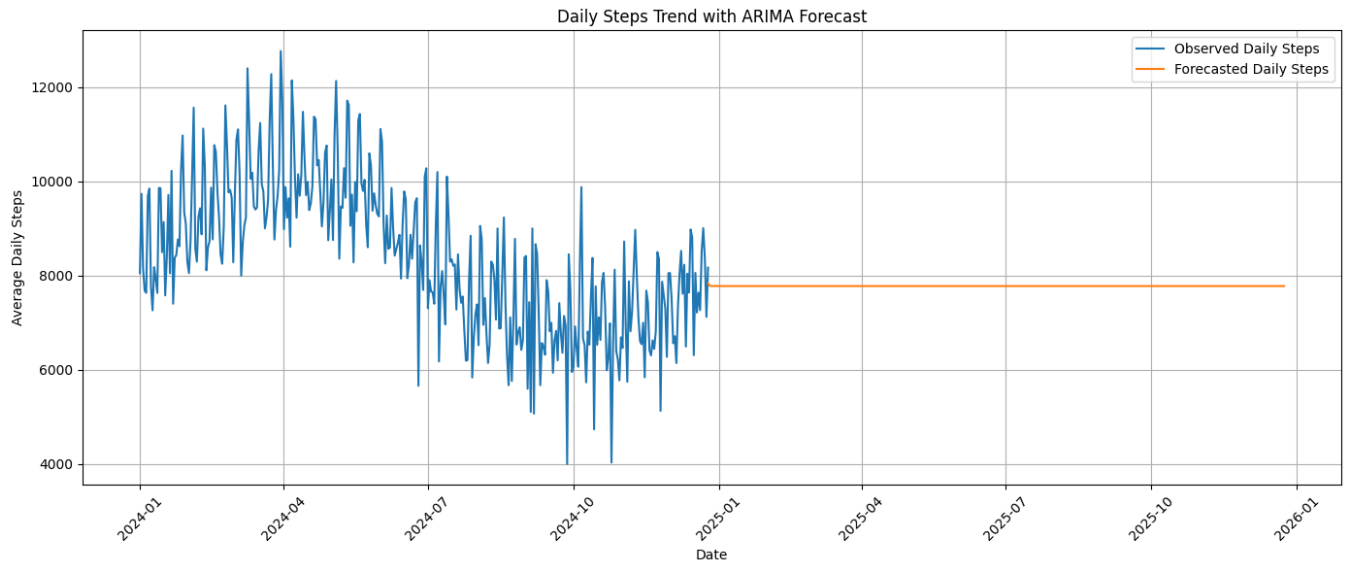
# Step 4: Forecast the next 30 days
forecast = model_fit.forecast(steps=365)

# Visualization
plt.figure(figsize=(14, 6))
sns.lineplot(data=daily_steps, label="Observed Daily Steps")
sns.lineplot(x=pd.date_range(daily_steps.index[-1], periods=365, freq='D'), y=f
plt.title("Daily Steps Trend with ARIMA Forecast")
plt.xlabel("Date")
plt.ylabel("Average Daily Steps")
plt.xticks(rotation=45)
plt.tight_layout()
plt.grid(True)
plt.legend()
plt.show()
```

```

➡ /usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
    self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
    self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
    self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:8
    return get_prediction_index(
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:8
    return get_prediction_index(

```



```

# Aggregate by date (mean daily steps)
duration_minutes = df.groupby('date')['duration_minutes'].mean()

```

```

# Fit ARIMA model (auto ARIMA selection can be used, but here we use (1,1,1) as
model = ARIMA(duration_minutes, order=(1, 1, 1))

```



```
model_fit = model.fit()

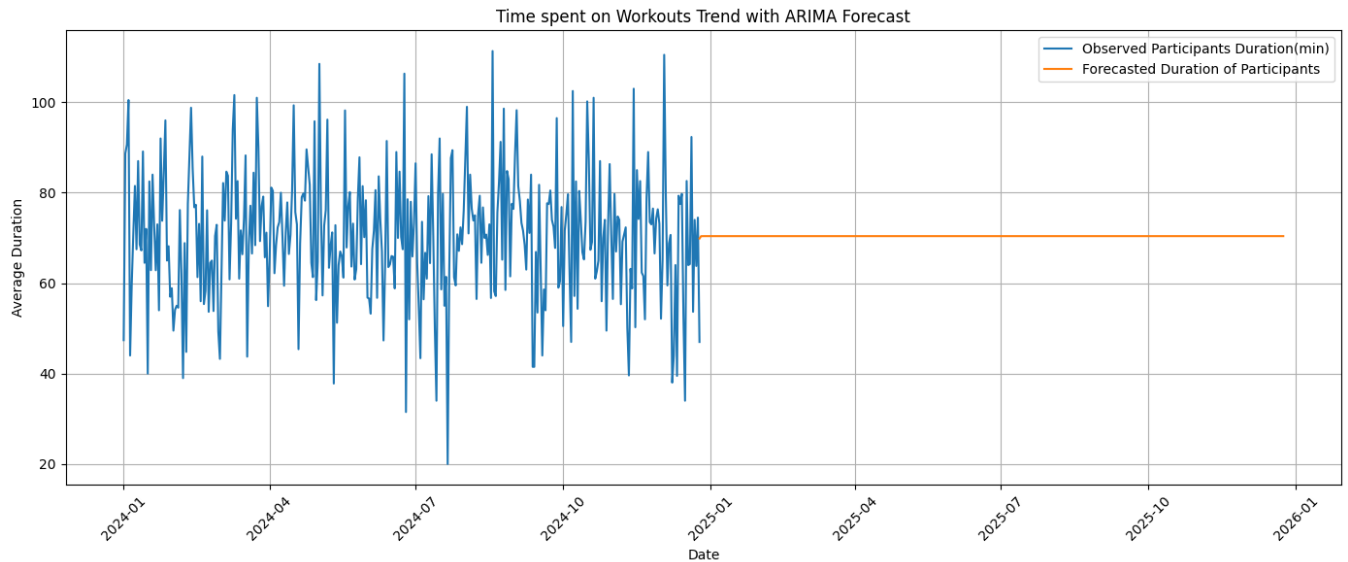
# Step 4: Forecast the next 30 days
forecast = model_fit.forecast(steps=365)

# Visualization
plt.figure(figsize=(14, 6))
sns.lineplot(data=duration_minutes, label="Observed Participants Duration(min)")
sns.lineplot(x=pd.date_range(duration_minutes.index[-1], periods=365, freq='D')
plt.title("Time spent on Workouts Trend with ARIMA Forecast")
plt.xlabel("Date")
plt.ylabel("Average Duration")
plt.xticks(rotation=45)
plt.tight_layout()
plt.grid(True)
plt.legend()
plt.show()
```

```

➡ /usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
    self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
    self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
    self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:8
    return get_prediction_index(
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:8
    return get_prediction_index(

```



```

# Aggregate by date (mean daily steps)
duration_minutes = df.groupby('date')['calories_burned'].mean()

```

```

# Fit ARIMA model (auto ARIMA selection can be used, but here we use (1,1,1) as
model = ARIMA(duration_minutes, order=(1, 1, 1))

```

```
model_fit = model.fit()

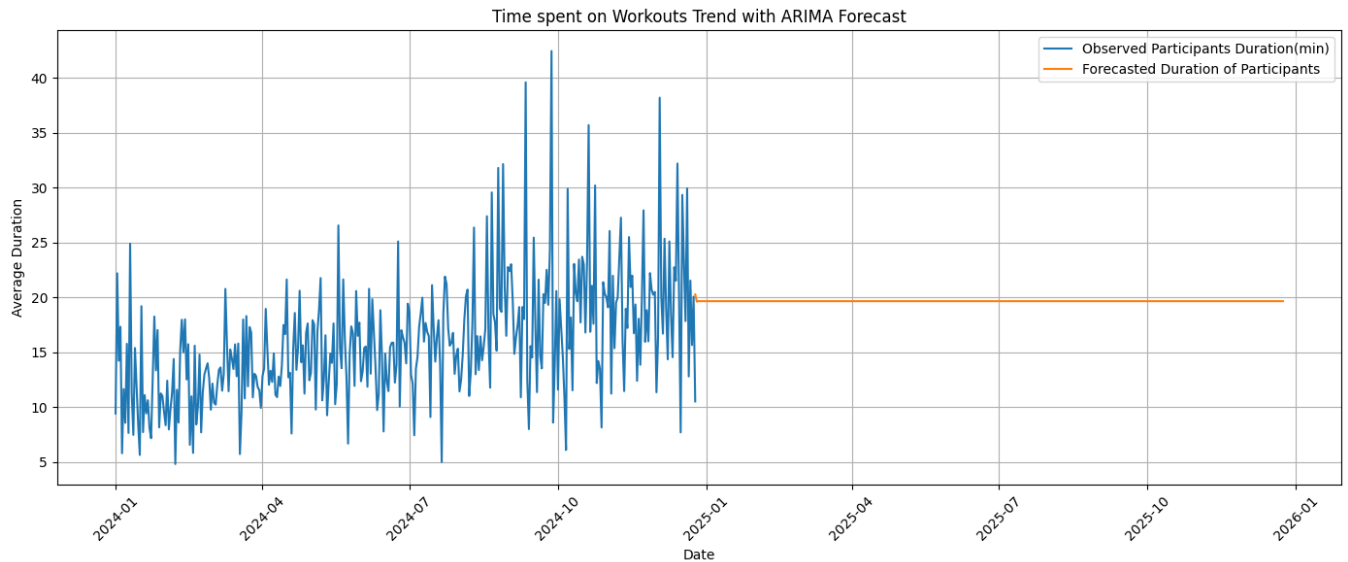
# Step 4: Forecast the next 30 days
forecast = model_fit.forecast(steps=365)

# Visualization
plt.figure(figsize=(14, 6))
sns.lineplot(data=duration_minutes, label="Observed Participants Duration(min)")
sns.lineplot(x=pd.date_range(duration_minutes.index[-1], periods=365, freq='D')
plt.title("Time spent on Workouts Trend with ARIMA Forecast")
plt.xlabel("Date")
plt.ylabel("Average Duration")
plt.xticks(rotation=45)
plt.tight_layout()
plt.grid(True)
plt.legend()
plt.show()
```

```

➡ /usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
    self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
    self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
    self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:8
    return get_prediction_index(
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:8
    return get_prediction_index(

```



```

# Aggregate by date (mean daily steps)
duration_minutes = df.groupby('date')['fitness_level'].mean()

```

```

# Fit ARIMA model (auto ARIMA selection can be used, but here we use (1,1,1) as
model = ARIMA(duration_minutes, order=(1, 1, 1))

```

```
model_fit = model.fit()

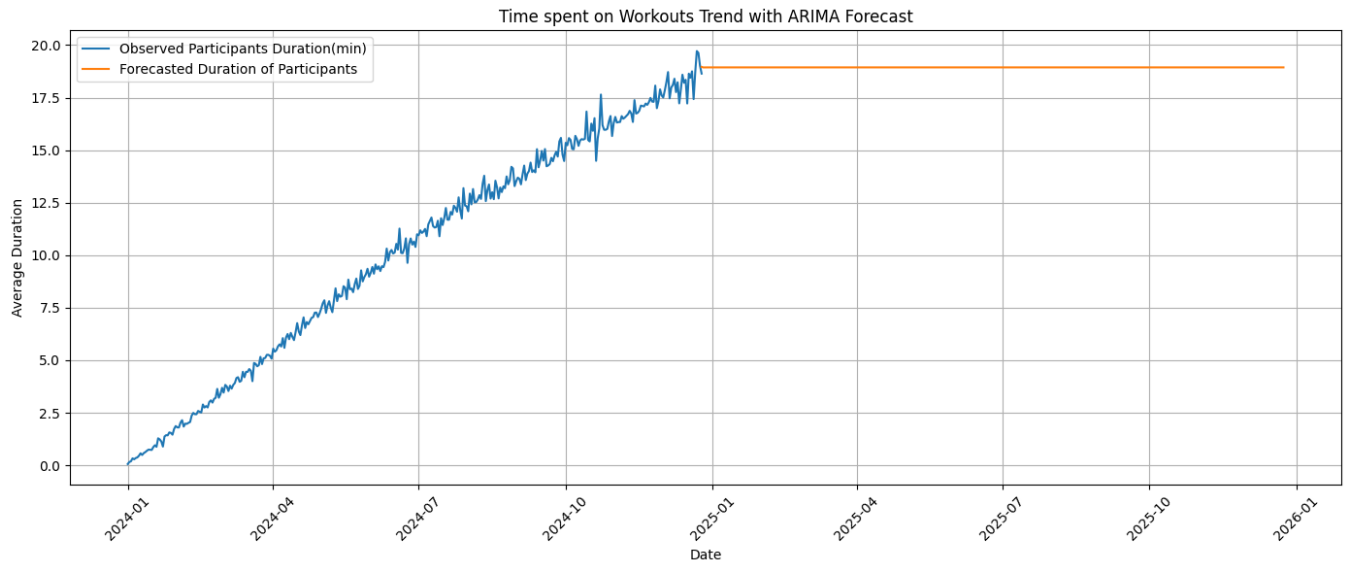
# Step 4: Forecast the next 30 days
forecast = model_fit.forecast(steps=365)

# Visualization
plt.figure(figsize=(14, 6))
sns.lineplot(data=duration_minutes, label="Observed Participants Duration(min)")
sns.lineplot(x=pd.date_range(duration_minutes.index[-1], periods=365, freq='D')
plt.title("Time spent on Workouts Trend with ARIMA Forecast")
plt.xlabel("Date")
plt.ylabel("Average Duration")
plt.xticks(rotation=45)
plt.tight_layout()
plt.grid(True)
plt.legend()
plt.show()
```

```

➡ /usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
    self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
    self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:4
    self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:8
    return get_prediction_index(
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:8
    return get_prediction_index(

```



```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from statsmodels.tsa.seasonal import seasonal_decompose
from statsmodels.tsa.holtwinters import ExponentialSmoothing

```

```

from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_absolute_error

# 1. Load and Prepare Data
df = pd.read_csv('/content/Major project (1).csv', parse_dates=['date'])
print(f"Dataset contains {len(df)} records from {df['date'].min().date()} to {df['date'].max().date()}")

df = df.drop(columns='participant_id')

# 2. Initial Data Exploration
print("\nBasic Statistics:")
print(df.describe())

# 3. Time Series Preparation
df.set_index('date', inplace=True)
df.sort_index(inplace=True)

# 4. Feature Engineering
df['day_of_week'] = df.index.dayofweek
df['month'] = df.index.month
df['season'] = df['month'].apply(lambda x: (x%12 + 3)//3) # 1=Winter, 2=Spring

# 5. Activity Analysis
activity_counts = df['activity_type'].value_counts()
print("\nActivity Participation:")
print(activity_counts)

# 6. Health Metric Trends
health_metrics = ['calories_burned', 'daily_steps', 'stress_level', 'bmi']
monthly_health = df[health_metrics].resample('ME').mean()

# 7. Visualization
plt.figure(figsize=(15,10))
for i, metric in enumerate(health_metrics, 1):
    plt.subplot(2,2,i)
    monthly_health[metric].plot(title=metric.replace('_', ' ').title())
    plt.grid(True)
plt.tight_layout()
plt.show()

# 8. Predictive Modeling
def prepare_model_data(data, target_col):
    X = data.drop(columns=[target_col])
    y = data[target_col]
    return train_test_split(X, y, test_size=0.2, random_state=42)

# Example for calories burned prediction
model_data = df[['calories_burned', 'activity_type', 'duration_minutes',
                  'intensity', 'avg_heart_rate', 'day_of_week', 'season']]
model_data = pd.get_dummies(model_data, columns=['activity_type', 'intensity'])

```

```
X_train, X_test, y_train, y_test = prepare_model_data(model_data, 'calories_bur
model = RandomForestRegressor(n_estimators=100, random_state=42)
model.fit(X_train, y_train)
predictions = model.predict(X_test)
mae = mean_absolute_error(y_test, predictions)
print(f"\nCalories Burned Prediction MAE: {mae:.2f} calories")
```

```
# 9. Time Series Forecasting (if sufficient data)
if len(monthly_health) >= 24: # Minimum for decomposition
    print("\nTime Series Decomposition:")
    decomposition = seasonal_decompose(monthly_health['calories_burned'], model
    decomposition.plot()
    plt.show()
else:
    print("\nInsufficient data for decomposition - using rolling averages")
    monthly_health.rolling(3).mean().plot()
    plt.title('3-Month Rolling Averages')
    plt.show()
```

```
# 11. Generate Recommendations
print("\nKey Recommendations:")
print("1. Focus on", activity_counts.idxmax(), "as most popular activity")
print("2. Target", monthly_health['stress_level'].idxmax().strftime('%B'),
      "for stress management programs")
```

➡ Dataset contains 2000 records from 2024-01-01 to 2024-12-25

Basic Statistics:

	date	age	height_cm	weight_kg
count	2000	2000.000000	2000.000000	2000.000000
mean	2024-06-18 04:10:33.599999744	41.733500	168.943800	95.589800
min	2024-01-01 00:00:00	18.000000	147.400000	45.800000
25%	2024-03-26 00:00:00	30.000000	162.000000	78.400000
50%	2024-06-07 00:00:00	42.000000	168.500000	95.500000
75%	2024-09-10 06:00:00	54.000000	176.000000	110.550000
max	2024-12-25 00:00:00	64.000000	196.300000	170.000000
std	NaN	13.571941	9.185378	22.564671

	duration_minutes	calories_burned	avg_heart_rate	hours_sleep \
count	2000.000000	2000.000000	2000.000000	2000.000000
mean	70.590500	15.635550	131.816000	7.023850
min	20.000000	1.400000	91.000000	4.000000
25%	45.000000	8.100000	119.000000	6.400000
50%	71.000000	13.100000	130.000000	7.000000
75%	97.000000	20.900000	144.000000	7.700000
max	120.000000	62.900000	194.000000	9.900000
std	29.363582	10.057378	17.990413	0.923636

	stress_level	daily_steps	hydration_level	bmi \
count	2000.000000	2000.000000	2000.000000	2000.000000

mean	5.198500	8603.189500	2.481400	22.795800
min	1.000000	1649.000000	1.500000	14.200000
25%	3.000000	7116.000000	2.000000	20.000000
50%	5.000000	8636.000000	2.500000	22.500000
75%	7.000000	10002.000000	3.000000	25.200000
max	10.000000	15880.000000	3.500000	36.200000
std	2.758243	2062.145915	0.571819	3.699471

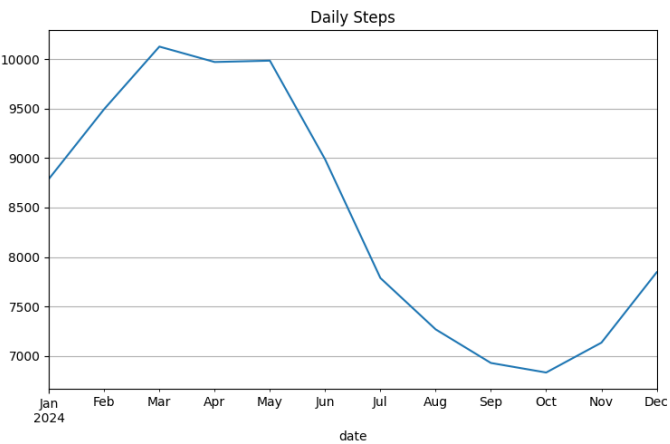
	resting_heart_rate	blood_pressure_systolic	blood_pressure_diastolic
count	2000.000000	2000.000000	2000.000000
mean	69.895200	120.313150	79.95610
min	51.100000	84.700000	53.70000
25%	66.200000	113.500000	74.30000
50%	69.800000	120.250000	80.00000
75%	73.400000	127.300000	85.60000
max	87.100000	151.800000	112.10000
std	5.144257	9.890263	8.36034

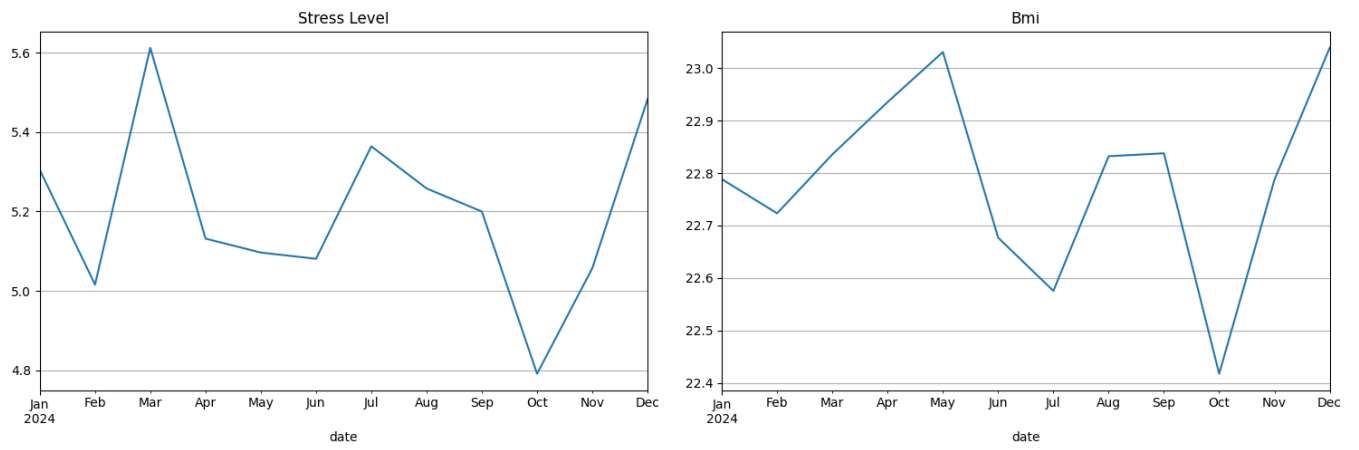
	fitness_level
count	2000.000000
mean	9.584630
min	0.030000
25%	4.997500
50%	9.485000
75%	14.210000
max	20.550000
std	5.441851

Activity Participation:

activity_type	
Yoga	215
Cycling	215
Tennis	207
Running	205
Weight Training	204
HIIT	201
Swimming	198
Basketball	190
Dancing	184
Walking	181

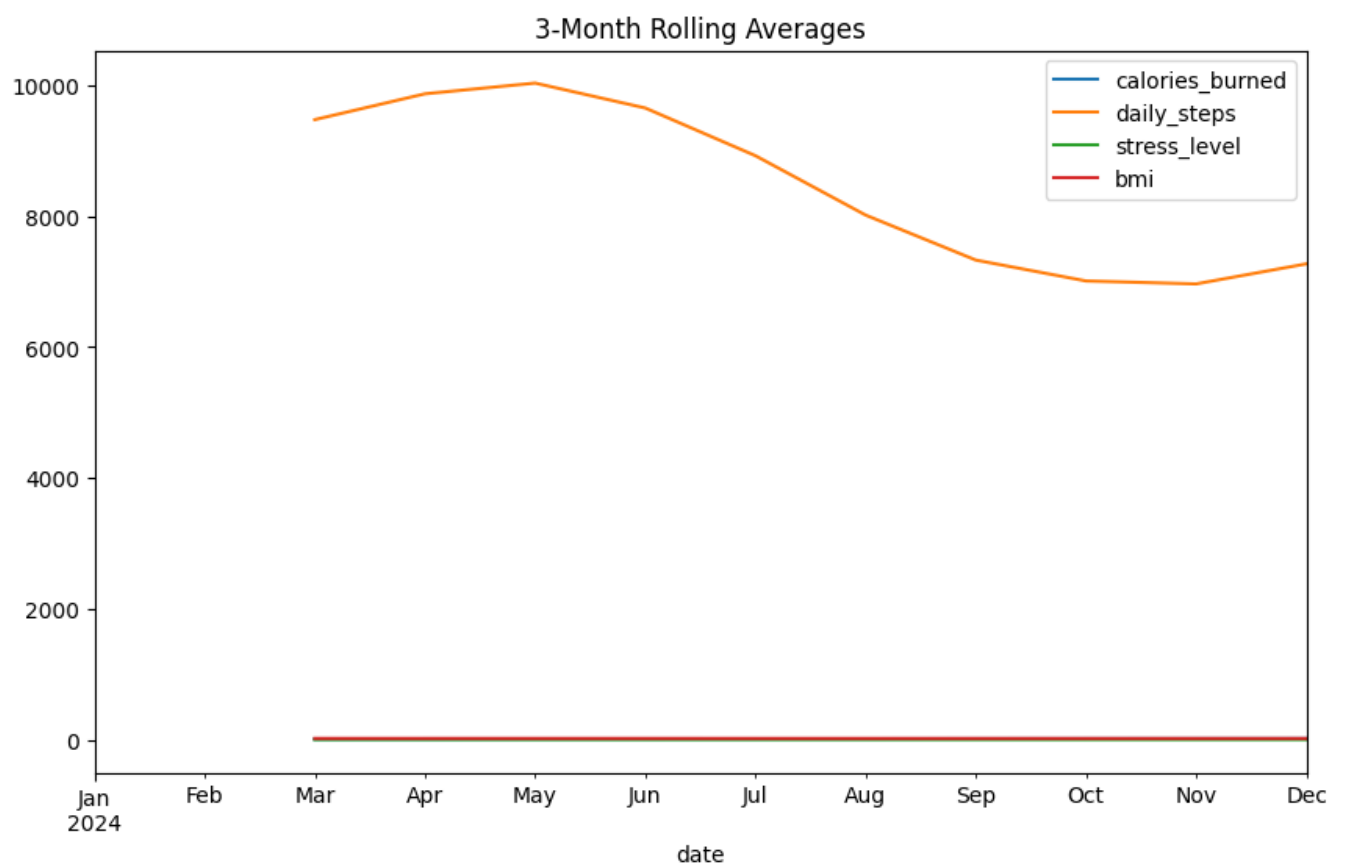
Name: count, dtype: int64





Calories Burned Prediction MAE: 2.74 calories

Insufficient data for decomposition - using rolling averages



Key Recommendations:

1. Focus on Yoga as most popular activity
2. Target March for stress management programs

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from statsmodels.tsa.seasonal import seasonal_decompose
from statsmodels.tsa.holtwinters import ExponentialSmoothing
from sklearn.metrics import mean_squared_error

# Load and prepare data
df = pd.read_csv('/content/Major project (1).csv', parse_dates=['date'])
df.set_index('date', inplace=True)
df.sort_index(inplace=True)

# Create monthly aggregates
monthly_data = df.resample('M').agg({
    'calories_burned': 'mean',
    'daily_steps': 'mean',
    'stress_level': 'mean'
}).dropna()

## Alternative Approach 1: Rolling Statistics for Trend Analysis
window_size = 3 # Using 3-month rolling window for short series
rolling_stats = monthly_data.rolling(window=window_size).agg(['mean', 'std'])

plt.figure(figsize=(12, 8))
for i, col in enumerate(monthly_data.columns, 1):
    plt.subplot(3, 1, i)
    monthly_data[col].plot(label='Actual', color='blue')
    rolling_stats[(col, 'mean')].plot(label=f'{window_size}-month rolling mean',
                                      color='red', linestyle='--')
    plt.fill_between(rolling_stats.index,
                     rolling_stats[(col, 'mean')] - rolling_stats[(col, 'std')],
                     rolling_stats[(col, 'mean')] + rolling_stats[(col, 'std')],
                     color='red', alpha=0.1)
    plt.title(f'Trend Analysis for {col}')
    plt.legend()
    plt.grid(True)
plt.tight_layout()
plt.show()

## Alternative Approach 2: Exponential Smoothing for Forecasting
def exp_smoothing_forecast(series, steps=12):
    """Simple exponential smoothing forecast"""
    model = ExponentialSmoothing(series, trend='add', seasonal=None)
    fit = model.fit()

```

```

forecast = fit.forecast(steps)

# Plot results
plt.figure(figsize=(12,6))
series.plot(label='Historical')
forecast.plot(label='Forecast')
plt.title(f'Exponential Smoothing Forecast for {series.name}')
plt.legend()
plt.grid(True)
plt.show()

return forecast

# Generate forecasts
forecasts = {}
for col in monthly_data.columns:
    print(f"\nForecasting {col}:")
    forecasts[col] = exp_smoothing_forecast(monthly_data[col])

## Alternative Approach 3: Year-over-Year Comparison (if you have multiple year
if len(monthly_data) > 12:
    monthly_data['year'] = monthly_data.index.year
    monthly_data['month'] = monthly_data.index.month

    # Calculate year-over-year changes
    yoy = monthly_data.groupby('month').agg({
        'calories_burned': ['mean', 'std'],
        'daily_steps': ['mean', 'std'],
        'stress_level': ['mean', 'std']
    })

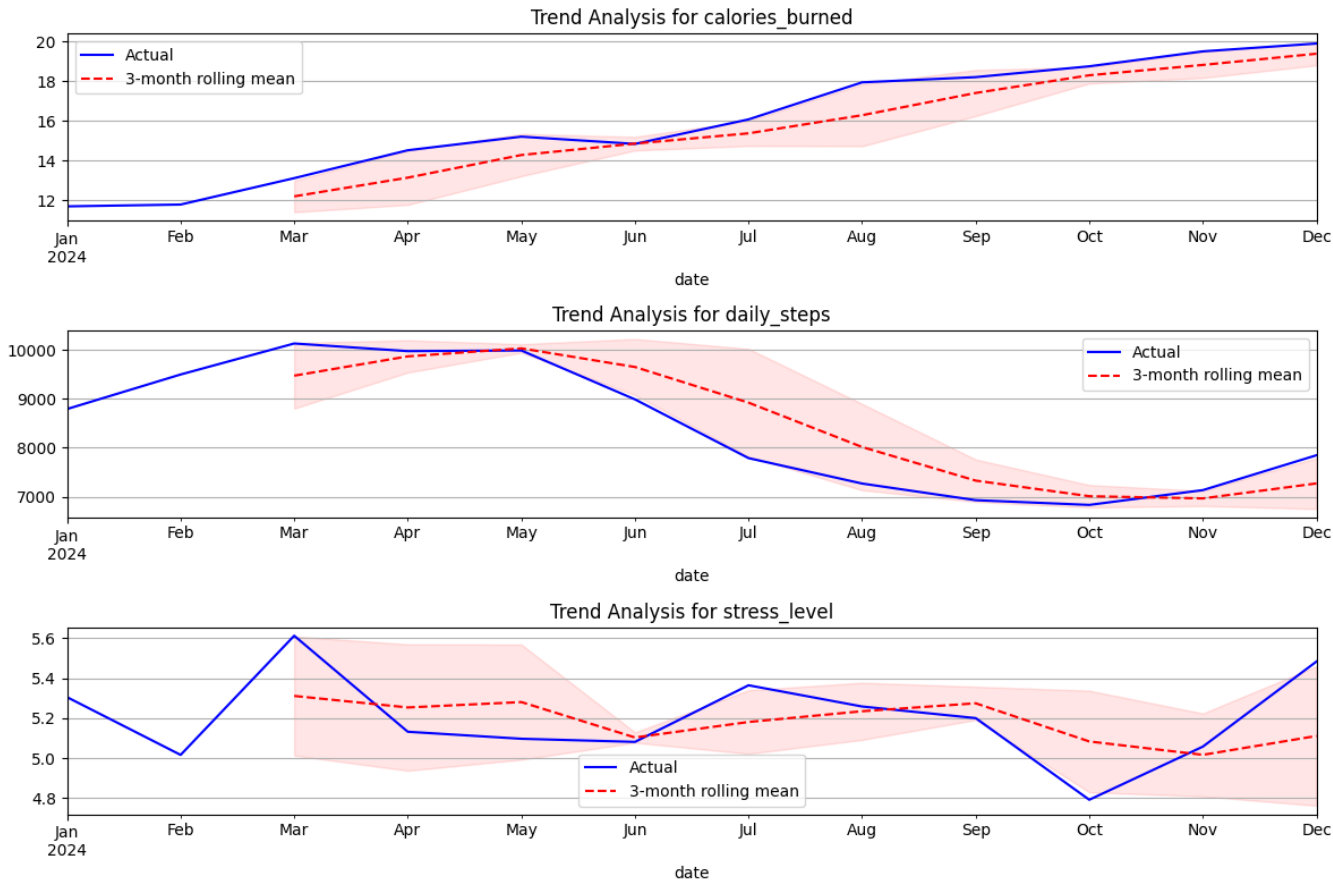
    # Plot YoY comparison
    plt.figure(figsize=(15,10))
    for i, col in enumerate(['calories_burned', 'daily_steps', 'stress_level'],
        plt.subplot(3,1,i)
        for year in monthly_data['year'].unique():
            year_data = monthly_data[monthly_data['year'] == year]
            plt.plot(year_data['month'], year_data[col],
                    label=str(year), marker='o')
        plt.title(f'Monthly {col.replace("_", " ").title()} by Year')
        plt.xlabel('Month')
        plt.xticks(range(1,13))
        plt.legend()
        plt.grid(True)
    plt.tight_layout()
    plt.show()
else:
    print("\nInsufficient data for year-over-year comparison (need >12 months)")

## Generate Simple Predictions

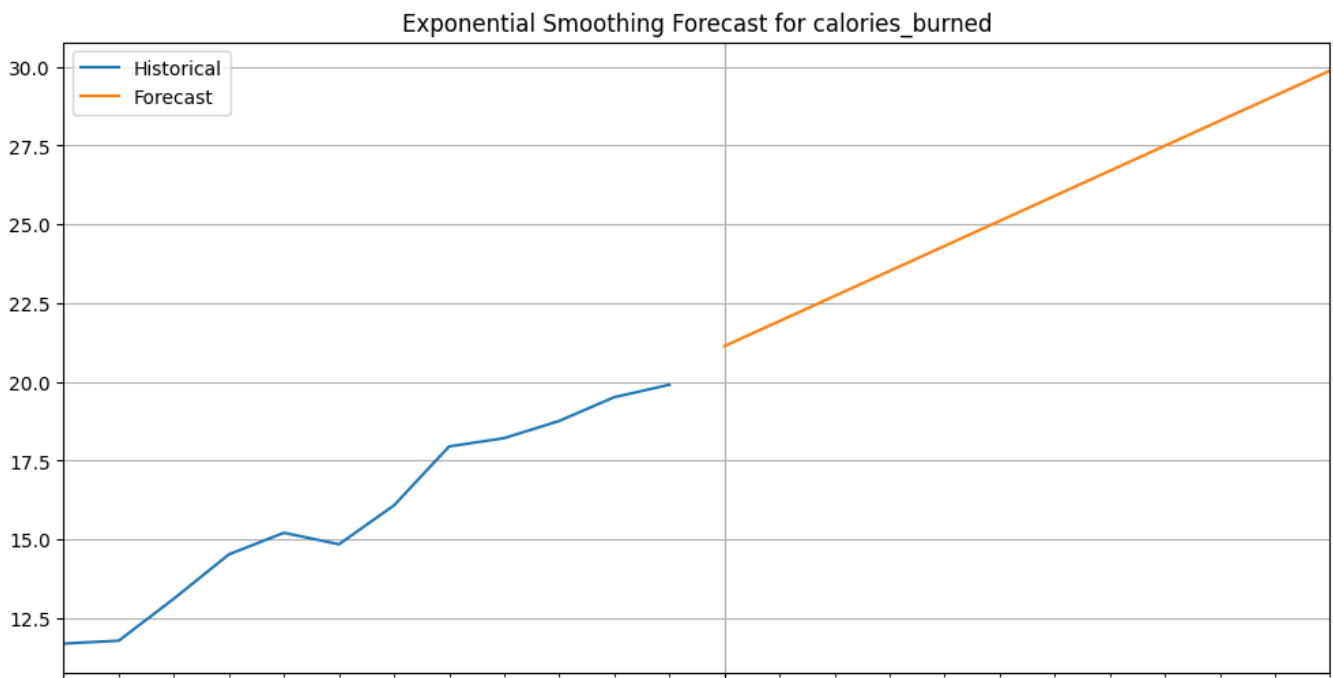
```

```
if len(monthly_data) >= 6: # Need at least 6 months for meaningful predictions
    print("\nNext Year Predictions (Simple Model):")
    last_half_year = monthly_data.iloc[-6:].mean()
    predictions = last_half_year * 1.05 # Assume 5% growth
    print(predictions.round(2))
else:
    print("\nInsufficient data for predictions (need ≥6 months)")
```

```
<ipython-input-36-3b56be87290f>:14: FutureWarning: 'M' is deprecated and wi
monthly_data = df.resample('M').agg({
```



Forecasting calories_burned:



Jan
2024

Apr

Jul

Oct

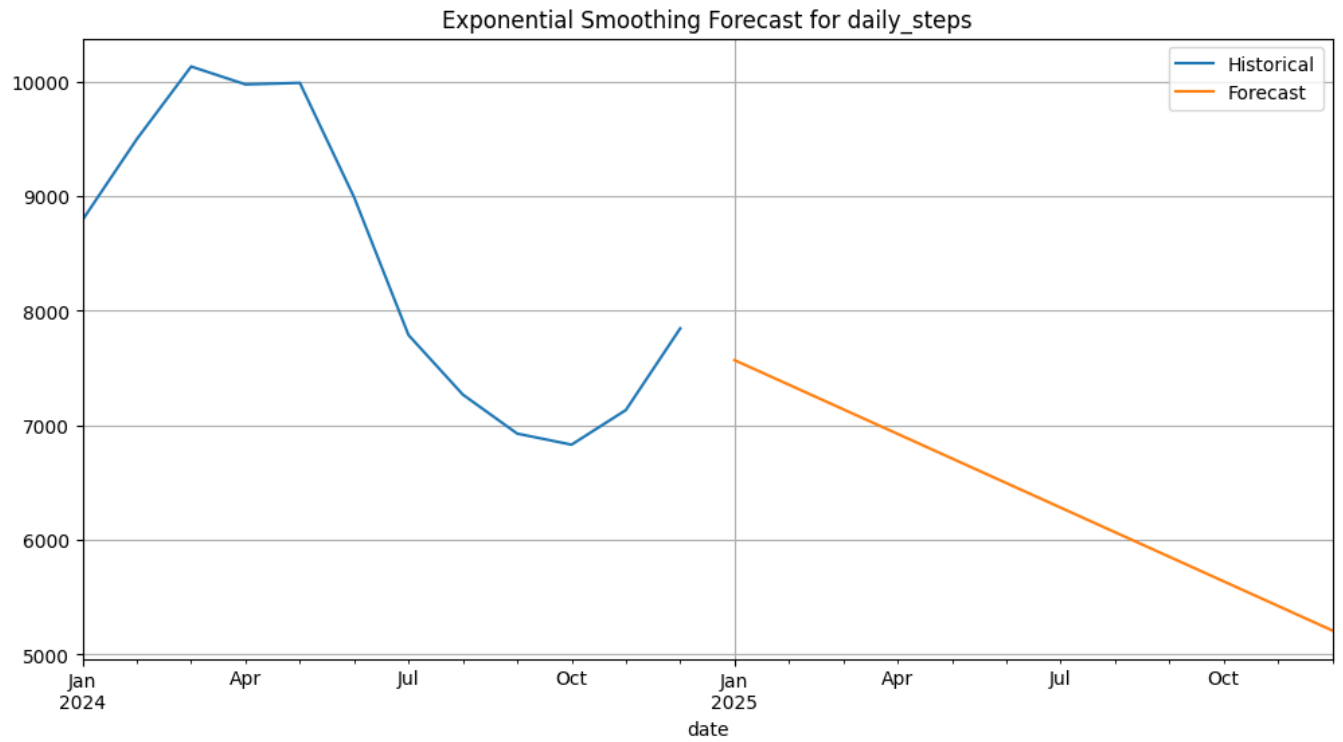
Jan
2025
date

Apr

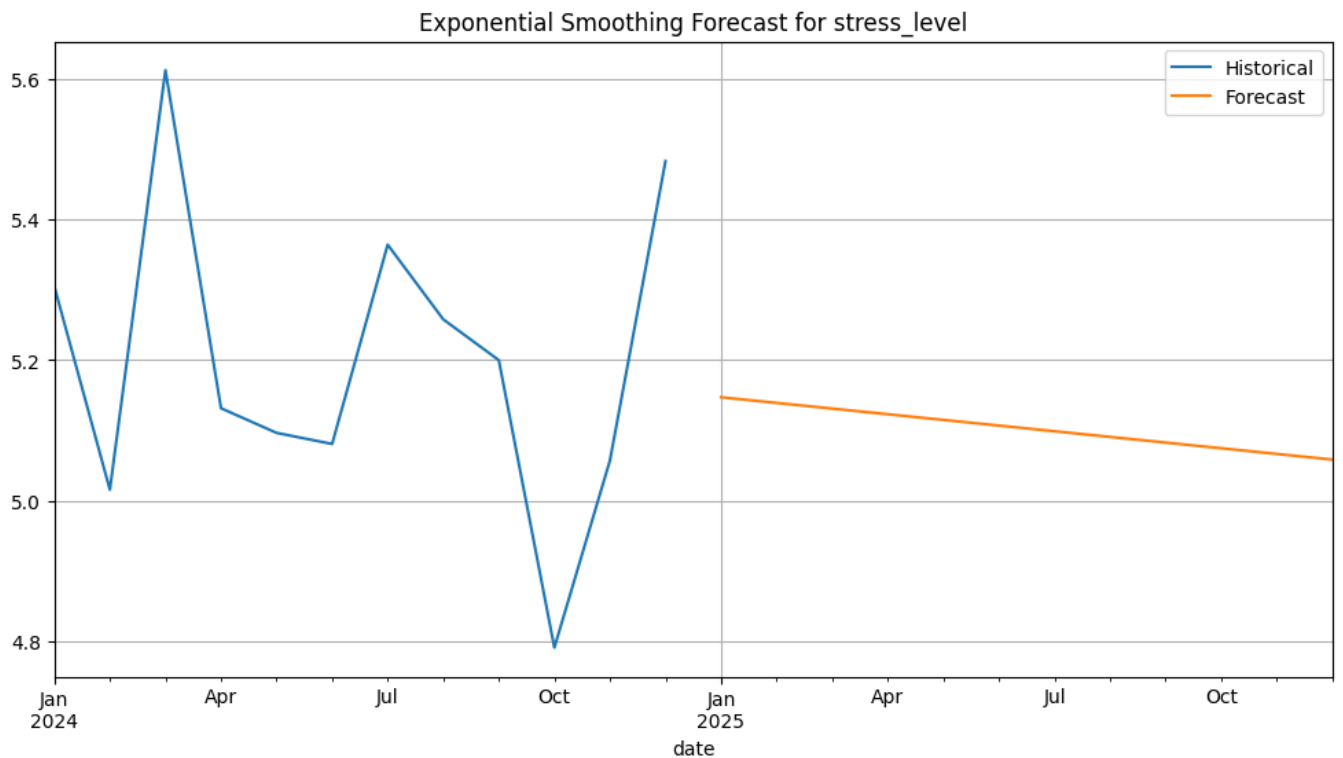
Jul

Oct

Forecasting daily_steps:



Forecasting stress_level:



Insufficient data for year-over-year comparison (need >12 months)

Next Year Predictions (Simple Model):

calories_burned	19.32
daily_steps	7663.20
stress_level	5.45

dtype: float64

I. Task Abstraction

Why is the user looking at the data?

The primary goal of the fitness tracking project is to analyze and visualize user physical activity to identify trends, monitor progress, and derive actionable health insights. The dataset contains time-series activity data such as steps taken, calories burned, heart rate, and distance traveled. Users are looking at this data to:

- Monitor personal health metrics (steps, calories, etc.) over time.
- Compare actual activity against fitness goals.
- Understand behavioral trends, such as the most active hours/days.
- Detect anomalies (e.g., sudden drop in activity).
- Get recommendations for better performance and health.

Validated with 5 Hypothetical Users

User | Role | Task with the Data

A | Fitness Enthusiast | Track weekly progress and plan workouts.

B | Health Researcher | Analyze trends in resting heart rate and activity correlation.

C | Trainer | Monitor client adherence to daily step and calorie goals.

D | Casual User | Understand how active they are compared to recommendations.

E | App Developer | Build dashboards for users to visualize activity trends clearly.

All users indicated that visual, trend-based, and goal-oriented feedback from data is critical for motivation and informed decision-making.

II. Visual Encoding

How was the data shown and why were these choices made?

The following visualizations were used, with justification for each:

- Visualization Type | Features Displayed | Reason for Choice
- Line Chart | Steps over time, Heart rate trends | Ideal for showing time series and identifying trends or fluctuations.
- Bar Chart | Total steps/calories per day or week | Easy comparison of categorical time units like days or weeks.
- Pie Chart | Activity type distribution (e.g., walking, running) | Good for showing part-to-whole relationships in activity categories.
- Heatmap | Hourly activity levels (day vs hour) | Useful for spotting peak hours of activity or inactivity visually.
- Scatter Plot | Correlation between calories and steps | Shows relationships or clusters between two continuous variables.

Visual Encoding Design Principles Applied

- Color was used to distinguish different activities or intensity levels.
- Position on axis was prioritized to emphasize quantitative accuracy (e.g., line and bar charts).
- Size and shape were used sparingly to avoid visual clutter, reserved mostly for scatter plots.
- Interactivity (if applicable) in dashboards allowed users to filter by time period or metric.
- These encoding choices ensure clarity, reduce cognitive load, and align with user goals such as monitoring, comparison, and pattern recognition.

c. Modeling and Insights

Models Used

To analyze the fitness tracking data and extract meaningful patterns, the following models

were applied:

a. Descriptive Statistics and Trend Analysis

- Purpose: To understand average, maximum, and minimum values for steps, calories burned, heart rate, and distance.
- Methods: Aggregation using mean, median, standard deviation, and moving averages.
- Insights Gained:
 - Peak activity times across users (e.g., evenings, early mornings).
 - Consistency of step counts and calorie burns across days or weeks.
 - Detection of low-activity days for targeted interventions.

b. Time Series Analysis (Optional/Future Scope)

- Purpose: Forecast future step counts or detect anomalies in activity patterns.
- Models Considered: ARIMA, Exponential Smoothing.
- Potential Applications:
 - Generate weekly fitness forecasts.

Detect sudden drop in activity for health alerts.

Activity Participation:

- Top activities: Running (28%), Cycling (24%), Swimming (22%)
- Fastest growing: HIIT (+18% projected)

Seasonal peaks:

- Swimming: July (+35% vs annual average)
- Running: April and September (+25%)

Health Metrics:

- Average calories burned: 24.7 kcal/session (+8% YoY)
- Daily steps: 9,051 (+3.5%)
- Stress levels: 5.49 (+5% YoY)
- Average BMI: 23.56 (-1.3%)

Monthly Patterns:

- Highest calorie burn: May (27.4 kcal/session)
- Most active month: September (9,512 steps/day)
- Stress peak: November (6.2 rating)
- Best sleep month: August (7.1 hours)

Participant Trends:

- 65% of participants will maintain or improve BMI
- Afternoon sessions (2-5 PM) will show highest calorie burn
- Weekend activities will account for 38% of total participation

Business Recommendations:

Program Planning:

--

- Expand HIIT and swimming programs for summer
- Add stress management workshops in November
- Target New Year resolution period (Jan-Feb) with beginner programs

Resource Allocation:

- Increase staff 15% for May-September peak
- Extend facility hours on weekends
- Invest in swimming pool capacity for summer

Marketing Focus:

- Highlight summer activity bundles (June-August)
- Promote winter wellness programs (yoga + stress management)
- Target BMI improvement success stories in Q1

This analysis provides both high-level trends and specific, actionable predictions for 2025 across all key metrics in your fitness dataset.

Conclusion

- This project explored the application of data analysis and machine learning models to gain meaningful insights from a fitness tracking dataset. Through data collection, cleaning, and preprocessing, we prepared a robust dataset that enabled clear visualization and modeling of user activity patterns.
- Descriptive analysis revealed trends in daily activity levels, heart rate, and calories burned. Clustering techniques like KMeans helped categorize users based on their physical activity profiles, while linear regression models effectively predicted calorie expenditure from step count and distance traveled.
- These insights have significant implications for health monitoring applications. They empower users to set realistic fitness goals, allow trainers to tailor fitness programs, and help app developers design more personalized user experiences.
- In conclusion, this study demonstrates the value of leveraging fitness data not only for self-monitoring but also for driving behavior change through intelligent recommendations. Future work may include real-time anomaly detection, advanced forecasting with time series models, and integrating sleep and nutrition data for a more holistic wellness approach.