```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import datetime
import datetime
import seaborn as sns
import seaborn.objects as so
import re
from itertools import product

from statsmodels.graphics.tsaplots import plot_acf , plot_pacf
from statsmodels.tsa.arima.model import ARIMA
from statsmodels.tsa.seasonal import seasonal_decompose
from statsmodels.tsa.stattools import adfuller
from statsmodels.tsa.statespace.sarimax import SARIMAX
from prophet import Prophet

sns.set(style = 'darkgrid')
pd.set_option('display.max_columns', None)
pd.options.display.max_colwidth = 100
plt.rcParams["figure.figsize"] = (15,7)
import warnings
warnings.filterwarnings('ignore')
```

```python
data = pd.read_csv('/content/train_1.csv')
exog = pd.read_csv('/content/Exog_Campaign_eng')
```

```python
raw_data = data.copy(deep=True)
```

```python
data.head()
```

| | Page | 2015-07-01 | 2015-07-02 | 2015-07-03 | 2015-07-04 | 2015-07-05 | 2015-07-06 | 2015-07-07 | 2015-07-08 | 2015-07-09 | 2015-07-10 | 2015-07-11 | 2015-07-12 | 2015-07-13 | 2015-07-14 | 2015-07-15 | 2015-07-16 | 2015-07-17 | 2015-07-18 | 2015-07-19 | 2015-07-20 | 2015-07-21 | 2015-07-22 | 2015-07-23 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2NE1_zh.wikipedia.org_all-access_spider | 18.0 | 11.0 | 5.0 | 13.0 | 14.0 | 9.0 | 9.0 | 22.0 | 26.0 | 24.0 | 19.0 | 10.0 | 14.0 | 15.0 | 8.0 | 16.0 | 8.0 | 8.0 | 16.0 | 7.0 | 11.0 | 10.0 | 20.0 |
| 1 | 2PM_zh.wikipedia.org_all-access_spider | 11.0 | 14.0 | 15.0 | 18.0 | 11.0 | 13.0 | 22.0 | 11.0 | 10.0 | 4.0 | 41.0 | 65.0 | 57.0 | 38.0 | 20.0 | 62.0 | 44.0 | 15.0 | 10.0 | 47.0 | 24.0 | 17.0 | 22.0 |
| 2 | 3C_zh.wikipedia.org_all-access_spider | 1.0 | 0.0 | 1.0 | 1.0 | 0.0 | 4.0 | 0.0 | 3.0 | 4.0 | 4.0 | 1.0 | 1.0 | 1.0 | 6.0 | 8.0 | 6.0 | 4.0 | 5.0 | 1.0 | 2.0 | 3.0 | 8.0 | 8.0 |
| 3 | 4minute_zh.wikipedia.org_all-access_spider | 35.0 | 13.0 | 10.0 | 94.0 | 4.0 | 26.0 | 14.0 | 9.0 | 11.0 | 16.0 | 16.0 | 11.0 | 23.0 | 145.0 | 14.0 | 17.0 | 85.0 | 4.0 | 30.0 | 22.0 | 9.0 | 10.0 | 11.0 |
| 4 | 52_Hz_I_Love_You_zh.wikipedia.org_all-access_spider | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN |

```
data.shape
```

(145063, 551)

We have 145063 different pages and visits for 550 days

```
data.isnull().sum()
```

|  | 0 |
|---|---|
| **Page** | 0 |
| **2015-07-01** | 20740 |
| **2015-07-02** | 20816 |
| **2015-07-03** | 20544 |
| **2015-07-04** | 20654 |
| **...** | ... |
| **2016-12-27** | 3701 |
| **2016-12-28** | 3822 |
| **2016-12-29** | 3826 |
| **2016-12-30** | 3635 |
| **2016-12-31** | 3465 |

551 rows × 1 columns

**dtype:** int64

```
data.loc[data['Page']=='52_Hz_I_Love_You_zh.wikipedia.org_all-access_spider']
d1 = datetime.strptime('2015-07-01', "%Y-%m-%d")
print('Start date:', d1)

d2 = datetime.strptime('2016-04-16', "%Y-%m-%d")
print('End time:',d2)

# get difference
delta = d2 - d1

# time difference in seconds
print(f"Days difference is {delta} seconds")
```
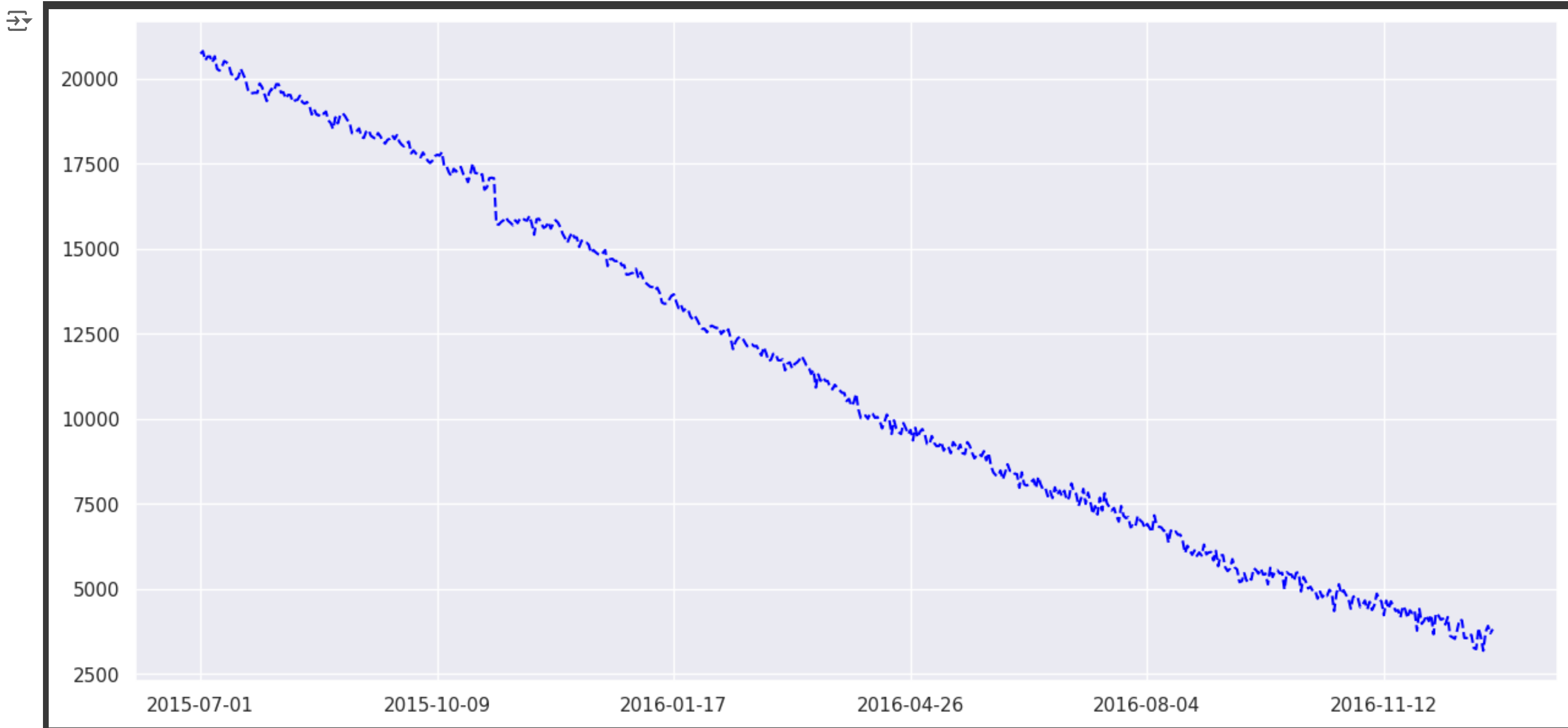
```
---------------------------------------------------------------------------
AttributeError                            Traceback (most recent call last)
<ipython-input-20-f6bde977d5e4> in <cell line: 2>()
      1 data.loc[data['Page']=='52_Hz_I_Love_You_zh.wikipedia.org_all-access_spider']
----> 2 d1 = datetime.strptime('2015-07-01', "%Y-%m-%d")
      3 print('Start date:', d1)
      4
      5 d2 = datetime.strptime('2016-04-16', "%Y-%m-%d")

AttributeError: module 'datetime' has no attribute 'strptime'
```

```
data.iloc[:, 1:-3 ].isnull().sum().plot(color='blue', linestyle='dashed')
plt.show()
```



The chart above illustrates a decreasing trend in NaN/Null values over time. Recent dates exhibit fewer Null Values compared to earlier dates.

This phenomenon is plausible because pages created or hosted at later dates naturally lack data for previous dates (dates preceding their creation/hosting).

To address this, we plan to eliminate rows containing more than 300 Null Values and substitute the remaining Null Values with 0.

```
data.dropna(thresh = 300,inplace = True)
print(f'Shape of Data : {data.shape}')
```

Shape of Data : (133617, 551)

```
data.fillna(0, inplace = True)
```

```python
def get_language(name):
    if len(re.findall(r'_(.{2}).wikipedia.org_', name)) == 1 :
        return re.findall(r'_(.{2}).wikipedia.org_', name)[0]
    else: return 'Unknown_language'

data['language'] = data['Page'].apply(get_language)


language_dict ={'de':'German',
                'en':'English',
                'es': 'Spanish',
                'fr': 'French',
                'ja': 'Japenese' ,
                'ru': 'Russian',
                'zh': 'Chinese',
                'Unknown_language': 'Unknown_language'}

data['language'] = data['language'].map(language_dict)
```
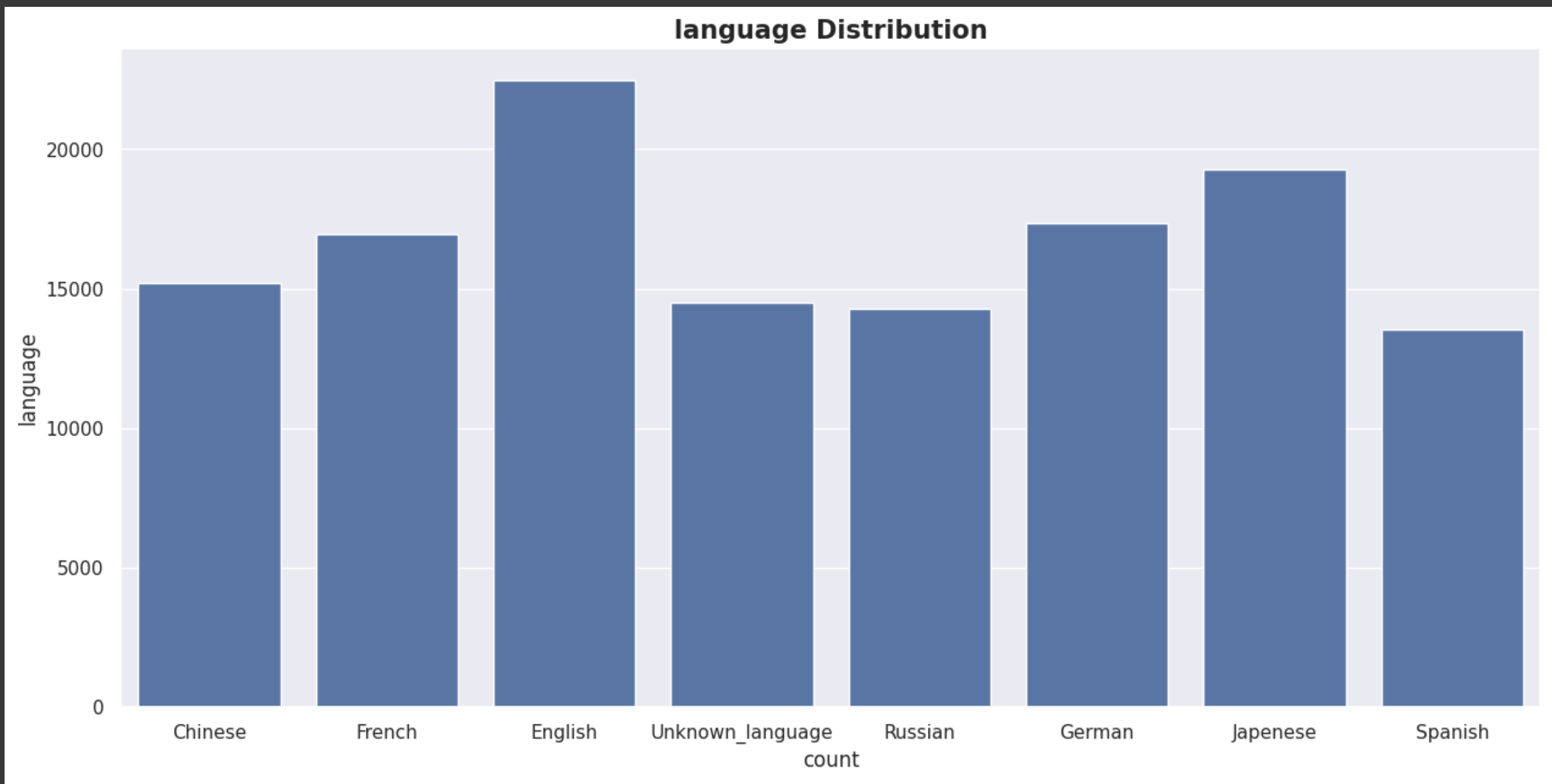
```python
def get_access_type(name):
    if len(re.findall(r'all-access|mobile-web|desktop', name)) == 1 :
        return re.findall(r'all-access|mobile-web|desktop', name)[0]
    else: return 'No Access_type'

data['access_type'] = data['Page'].apply(get_access_type)
```

```python
def get_access_origin(name):
    if len(re.findall(r'[ai].org_(.*)_(.*)$', name)) == 1 :
        return re.findall(r'[ai].org_(.*)_(.*)$', name)[0][1]
    else: return 'No Access_origin'

data['access_origin'] = data['Page'].apply(get_access_origin)
```
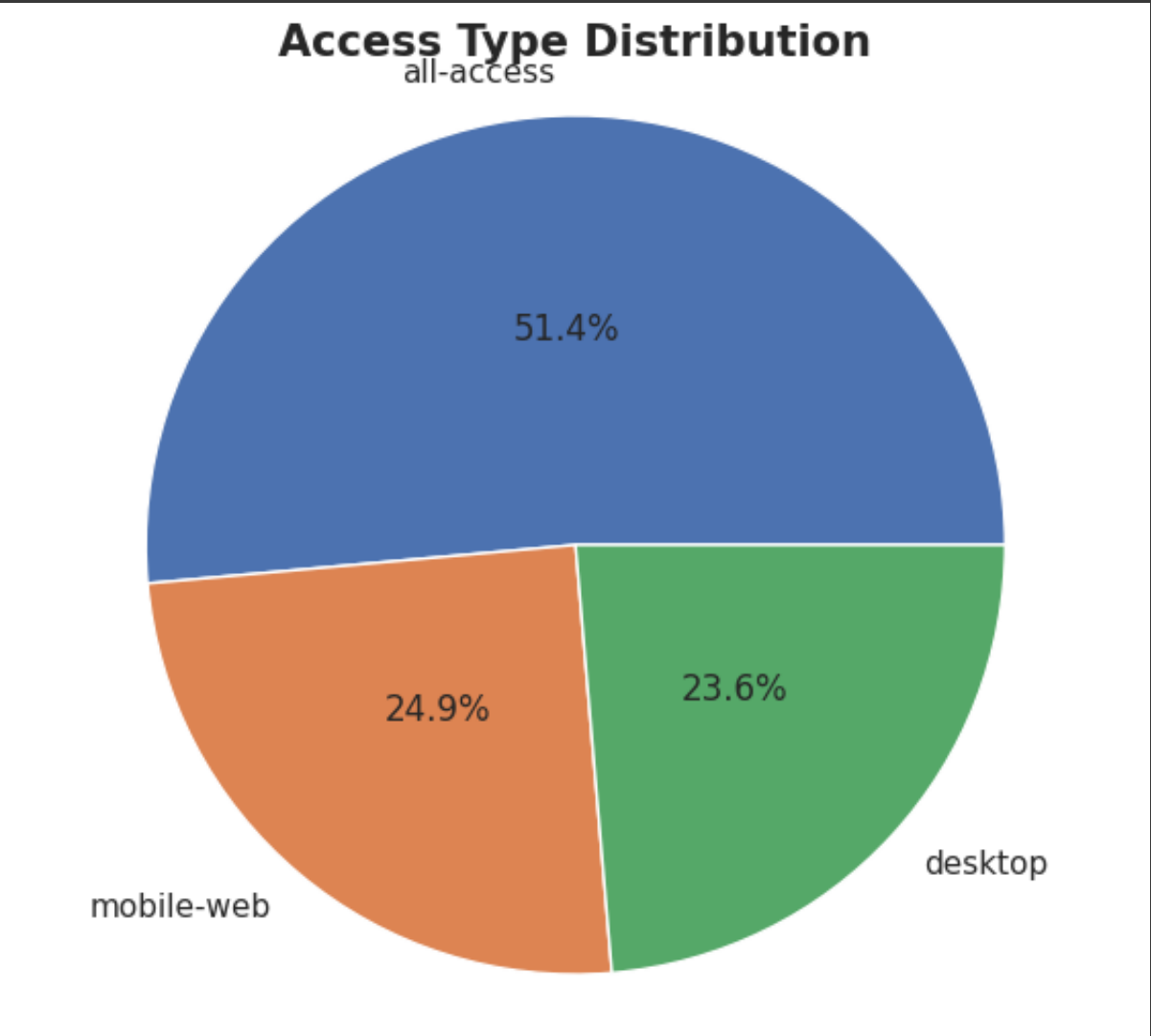
```
sns.countplot(x='language' , data=data)
plt.title('language Distribution')
plt.xlabel('count')
plt.ylabel('language')
plt.title('language Distribution', fontsize = 15, fontweight = 'bold')
plt.show()
```
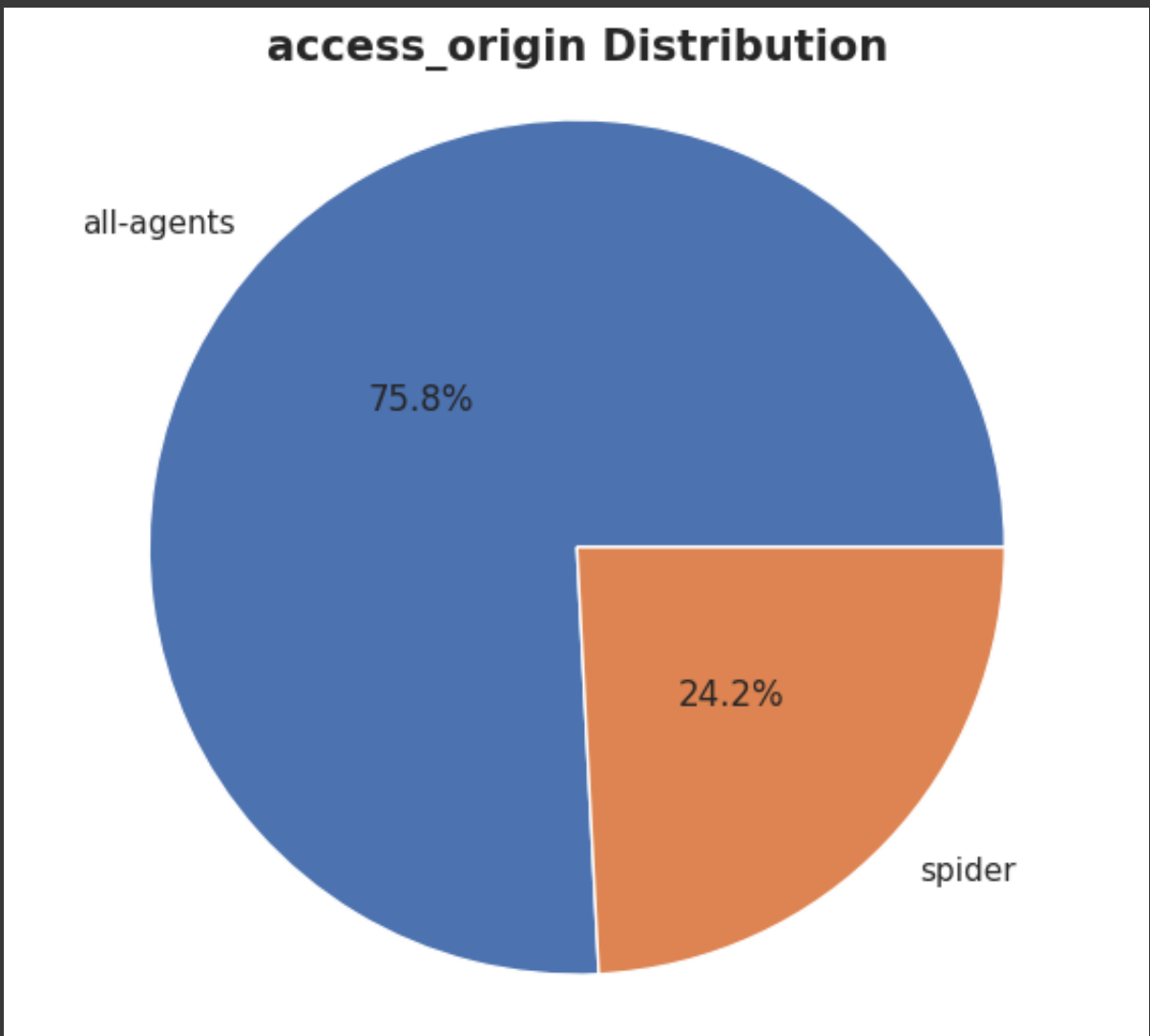
```
x = data['access_type'].value_counts().values
y = data['access_type'].value_counts().index

plt.figure(figsize=(7, 6))
plt.pie(x, labels = y, center=(0, 0), radius=1.5,  autopct='%1.1f%%', pctdistance=0.5)
plt.title('Access Type Distribution', fontsize = 15, fontweight = 'bold')
plt.axis('equal')
plt.show()
```

**Access Type Distribution**

```python
var = 'access_origin'
x = data[var].value_counts().values
y = data[var].value_counts().index

plt.figure(figsize=(7, 6))
plt.pie(x, labels = y, center=(0, 0), radius=1.5,  autopct='%1.1f%%', pctdistance=0.5)
plt.title(f'{var} Distribution', fontsize = 15, fontweight = 'bold')
plt.axis('equal')
plt.show()
```



```python
reshaped = data.melt(id_vars = ['Page','language','access_type','access_origin'])
```

```python
reshaped.head()
```

| | Page | language | access_type | access_origin | variable | value |
|---|---|---|---|---|---|---|
| 0 | 2NE1_zh.wikipedia.org_all-access_spider | Chinese | all-access | spider | 2015-07-01 | 18.0 |
| 1 | 2PM_zh.wikipedia.org_all-access_spider | Chinese | all-access | spider | 2015-07-01 | 11.0 |
| 2 | 3C_zh.wikipedia.org_all-access_spider | Chinese | all-access | spider | 2015-07-01 | 1.0 |
| 3 | 4minute_zh.wikipedia.org_all-access_spider | Chinese | all-access | spider | 2015-07-01 | 35.0 |
| 4 | 5566_zh.wikipedia.org_all-access_spider | Chinese | all-access | spider | 2015-07-01 | 12.0 |

```
reshaped.columns = ['Page','language','access_type','access_origin','Date','Visits']
```

```
reshaped.Date = pd.to_datetime(reshaped.Date, format = '%Y-%m-%d')
```

```
lang_data = reshaped.groupby(['language','Date'],as_index=False)['Visits'].sum()
```

```
lang_data.shape
```

```
(4400, 3)
```

```
lang_data.head()
```

|   | language | Date | Visits |
|---|----------|------|--------|
| 0 | Chinese | 2015-07-01 | 4144975.0 |
| 1 | Chinese | 2015-07-02 | 4151185.0 |
| 2 | Chinese | 2015-07-03 | 4123659.0 |
| 3 | Chinese | 2015-07-04 | 4163439.0 |
| 4 | Chinese | 2015-07-05 | 4441273.0 |

```python
sns.lineplot(data=lang_data, y = 'Visits', x = 'Date', hue = 'language')
```

<Axes: xlabel='Date', ylabel='Visits'>



```python
lang_data.head()
```

|   | language | Date | Visits |
|---|----------|------|--------|
| 0 | Chinese | 2015-07-01 | 4144975.0 |
| 1 | Chinese | 2015-07-02 | 4151185.0 |
| 2 | Chinese | 2015-07-03 | 4123659.0 |
| 3 | Chinese | 2015-07-04 | 4163439.0 |
| 4 | Chinese | 2015-07-05 | 4441273.0 |

```
def adf_test(timeseries):
    print ('Results of Dickey-Fuller Test:')
    dftest = adfuller(timeseries, autolag='AIC')
    df_output = pd.Series(dftest[0:4], index=['Test Statistic','p-value','#Lags Used','Number of Observations Used'])
    for key, value in dftest[4].items():
        df_output['Critical Value (%s)' %key] = round(value,2)
    print (df_output)
```

```
adf_test(lang_data[lang_data['language'] == 'English']['Visits'])
```

```
Results of Dickey-Fuller Test:
Test Statistic           -2.373563
p-value                   0.149337
#Lags Used               14.000000
Number of Observations Used   535.000000
Critical Value (1%)      -3.440000
Critical Value (5%)      -2.870000
Critical Value (10%)     -2.570000
dtype: float64
```

Decomposing Time Series Time Series Decomposition

Time series decomposition is a statistical technique used to break down a time series into its constituent components in order to understand its underlying structure, trends, seasonality, and irregular fluctuations. The decomposition typically involves separating the time series data into three main components:

Trend ($(T_t)$): The long-term movement or pattern in the data, representing the overall direction in which the time series is moving. Seasonality ($(S_t)$): The repeating patterns or fluctuations that occur at regular intervals within the time series data. Residuals ($(R_t)$): The remaining variation in the data after removing the trend and seasonality components. The time series ($y_t$) can be decomposed into its components as follows:

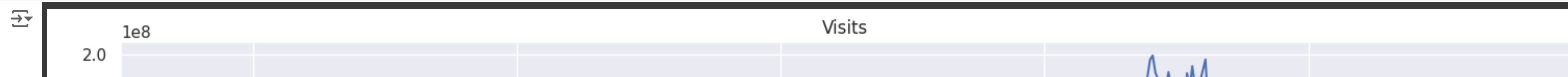Additive Decomposition: [ $y_t = T_t + S_t + R_t$ ]

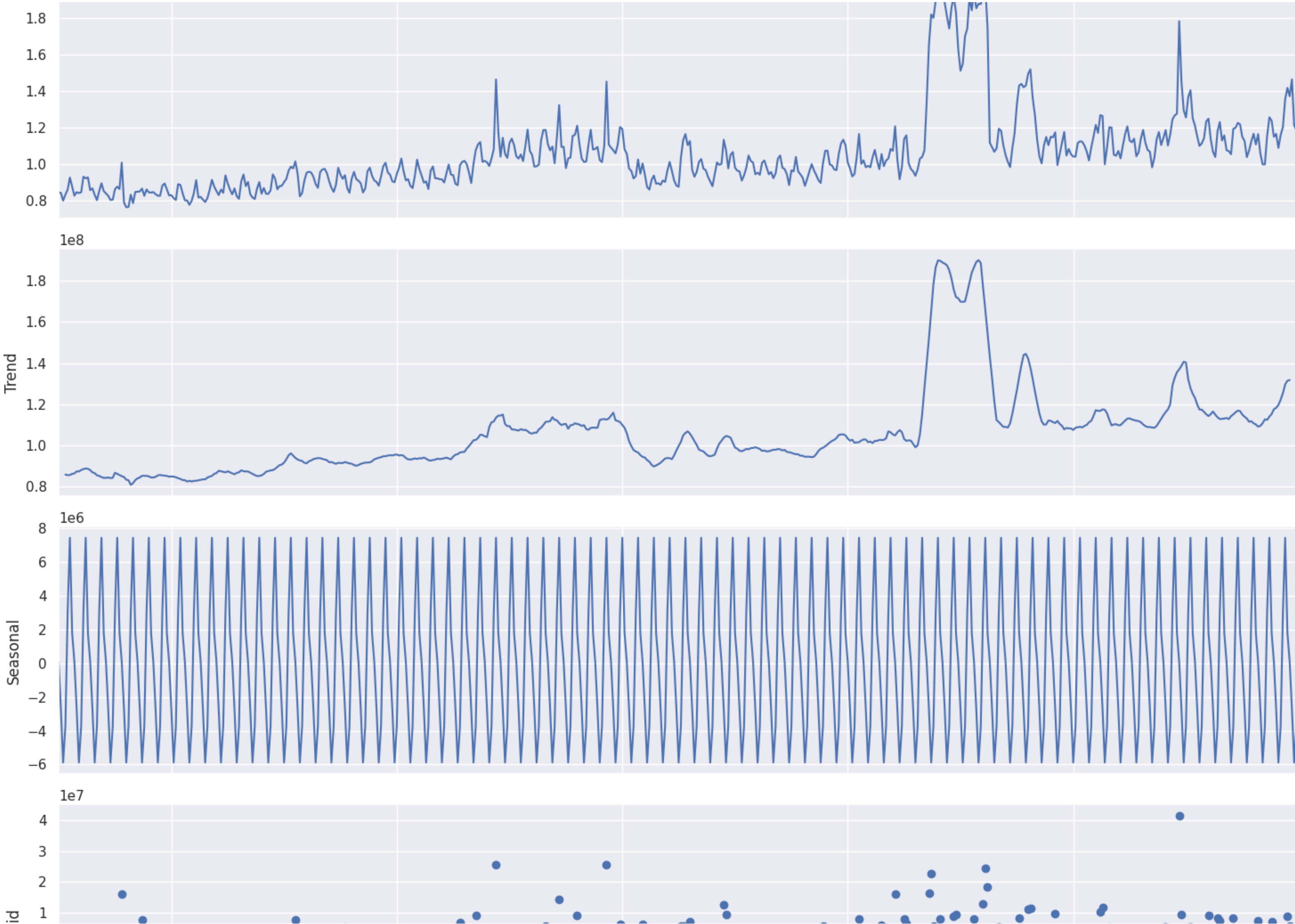Multiplicative Decomposition: [ $y_t = T_t \times S_t \times R_t$ ]

Various techniques such as moving averages, exponential smoothing, or mathematical models can be used to estimate the trend and seasonal components, leaving the residual component as the leftover variation in the data.
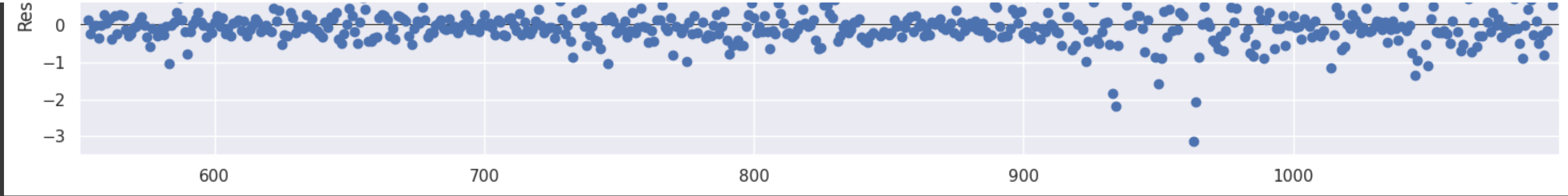
```
ts_english = lang_data[lang_data['language'] == 'English']['Visits']
```

```
decomposition = seasonal_decompose(ts_english, model = 'additive',period=7)

fig = decomposition.plot()
fig.set_size_inches((15,13))
fig.tight_layout()
plt.show()
```

```python
ts_diff = pd.DataFrame(ts_english).diff(1)
ts_diff.dropna(inplace = True)
```

```python
ts_diff.plot(color = 'green',figsize = (15,4))
plt.show()
```



```python
adf_test(ts_diff)
```

```
Results of Dickey-Fuller Test:
Test Statistic                  -8.273590e+00
p-value                          4.721272e-13
#Lags Used                       1.300000e+01
Number of Observations Used      5.350000e+02
Critical Value (1%)             -3.440000e+00
Critical Value (5%)             -2.870000e+00
Critical Value (10%)            -2.570000e+00
dtype: float64
```

```python
acf = plot_acf(ts_diff, lags = 15)
acf.tight_layout()
```

```
pacf = plot_pacf(ts_diff, lags = 15)
pacf.tight_layout()
```

ACF

If the ACF shows a sharp cutoff after lag 'k', it suggests that an AR(k) model may be appropriate. If the ACF decreases gradually, it suggests a non-stationary series, and differencing (d) may be needed. If the ACF has a sinusoidal pattern or fluctuates around zero, it suggests a seasonal component. The ACF shows a sharp cutoff after lag 0, it suggests that an AR(0) model may be appropriate.
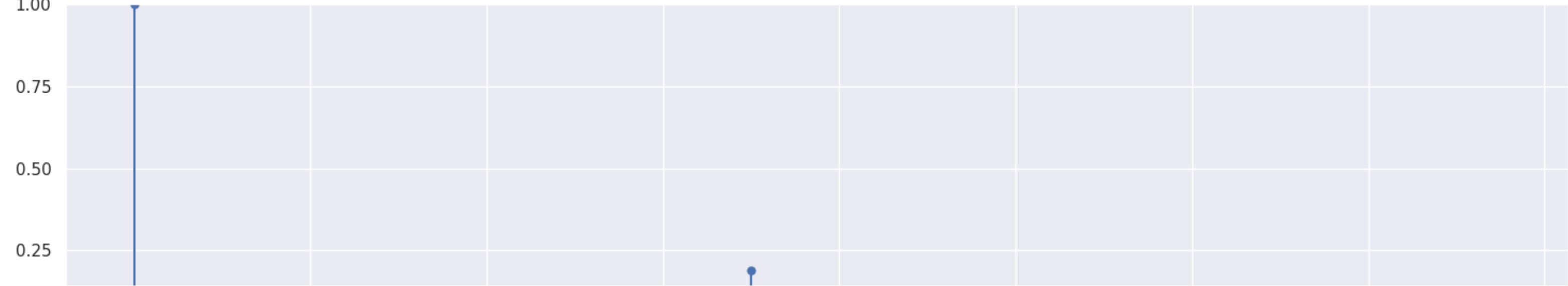
PACF

If the PACF has a sharp cutoff after lag 'k', it suggests an MA(k) model may be appropriate. If the PACF gradually decreases, it suggests an AR component. If there are significant spikes at seasonal lags, it suggests a seasonal AR or MA component. The PACF has a sharp cutoff after lag 0, it suggests an MA(0) model may be appropriate.

```
ts_english = lang_data[lang_data.language == 'English'][['Date', 'Visits']]
ts_english.set_index('Date', drop=True, inplace=True)
```

```python
def arima_model(n, order, time_series):
    model = ARIMA(time_series[:-n], order=order)
    model_fit = model.fit()
    forecast = model_fit.forecast(steps=n, alpha=0.05)
    time_series.index = pd.to_datetime(time_series.index)
    forecast.index = pd.to_datetime(forecast.index)
    time_series[-60:].plot(label='Actual')
    forecast.plot(label='Forecast', linestyle='dashed', marker='o', markerfacecolor='green', markersize=2)
    plt.legend(loc="upper right")
    plt.title(f'ARIMA BASE Model {order}: Actual vs Forecasts', fontsize=15, fontweight='bold')
    plt.show()


    actuals = time_series.values[-n:]
    errors = time_series.values[-n:] - forecast.values

    mape = np.mean(np.abs(errors) / np.abs(actuals))
    rmse = np.sqrt(np.mean(errors**2))

    print('-' * 80)
    print(f'MAPE of Model: {np.round(mape, 5)}')
    print('-' * 80)
    print(f'RMSE of Model: {np.round(rmse, 3)}')
    print('-' * 80)
```

```
arima_model(30, (1,1,1), ts_english)
```



ARIMA BASE Model (1, 1, 1): Actual vs Forecasts

```
--------------------------------------------------------------------------
MAPE of Model: 0.07229
--------------------------------------------------------------------------
RMSE of Model: 12071774.914
--------------------------------------------------------------------------
```

```python
def sarimax_model(time_series, n, p=0, d=0, q=0, P=0, D=0, Q=0, s=0, exog = []):

    model = SARIMAX(time_series[:-n],
                    order=(p, d, q),
                    seasonal_order=(P, D, Q, s),
                    exog=exog[:-n],
                    initialization='approximate_diffuse')
    model_fit = model.fit()

    model_forecast = model_fit.forecast(n, dynamic=True, exog=pd.DataFrame(exog[-n:]))

    plt.figure(figsize=(20, 8))
    time_series[-60:].plot(label='Actual')
    model_forecast[-60:].plot(label='Forecast', color='red',
                              linestyle='dashed', marker='o', markerfacecolor='green', markersize=5)
    plt.legend(loc="upper right")
    plt.title(f'SARIMAX Model ({p},{d},{q}) ({P},{D},{Q},{s}) : Actual vs Forecasts', fontsize=15, fontweight='bold')
    plt.show()

    actuals = time_series.values[-n:]
    errors = time_series.values[-n:] - model_forecast.values

    mape = np.mean(np.abs(errors) / np.abs(actuals))
    rmse = np.sqrt(np.mean(errors ** 2))

    print('-' * 80)
    print(f'MAPE of Model : {np.round(mape, 5)}')
    print('-' * 80)
    print(f'RMSE of Model : {np.round(rmse, 3)}')
    print('-' * 80)
```
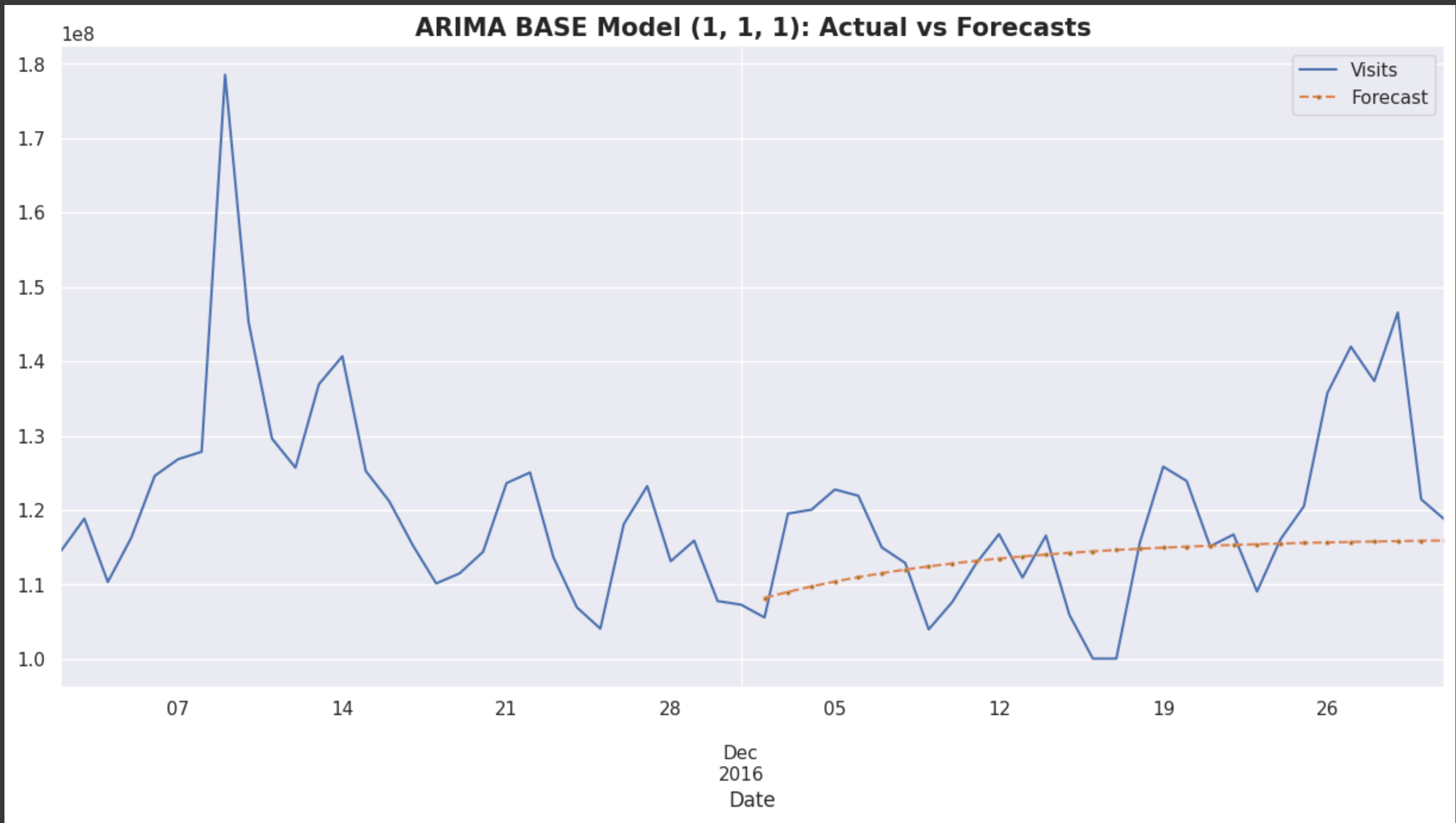
```python
exog = exog['Exog'].to_numpy()
```

```
time_series = ts_english
test_size= 0.1
p,d,q, P,D,Q,s = 1,1,1,1,1,1,7
n = 30
sarimax_model(time_series, n, p = p, d=d,q=q,   P=P, D=D, Q=Q, s=s, exog = exog)
```

<Figure size 2000x800 with 0 Axes>



SARIMAX Model (1,1,1) (1,1,1,7) : Actual vs Forecasts

```
----------------------------------------------------------------------
MAPE of Model : 0.11208
----------------------------------------------------------------------
RMSE of Model : 17326667.279
----------------------------------------------------------------------
```

```python
def sarimax_grid_search(time_series, n, param, d_param, s_param, exog=[]):
    param_df = pd.DataFrame(columns=['serial', 'pdq', 'PDQs', 'mape', 'rmse'])

    param_combinations = product(param, d_param, param, param, d_param, param, s_param)

    counter = 0

    for p, d, q, P, D, Q, s in param_combinations:
        model = SARIMAX(time_series[:-n],
                        order=(p, d, q),
                        seasonal_order=(P, D, Q, s),
                        exog=exog[:-n],
                        initialization='approximate_diffuse')
        model_fit = model.fit()

        model_forecast = model_fit.forecast(n, dynamic=True, exog=pd.DataFrame(exog[-n:]))

        actuals = time_series.values[-n:]
        errors = time_series.values[-n:] - model_forecast.values

        mape = np.mean(np.abs(errors) / np.abs(actuals))
        rmse = np.sqrt(np.mean(errors**2))
        mape = np.round(mape, 5)
        rmse = np.round(rmse, 3)

        counter += 1
        list_row = [counter, (p, d, q), (P, D, Q, s), mape, rmse]
        param_df.loc[len(param_df)] = list_row

        print(f'Possible Combination: {counter} out of {len(param)**4 * len(s_param) * len(d_param)**2} calculated')

    return param_df
```
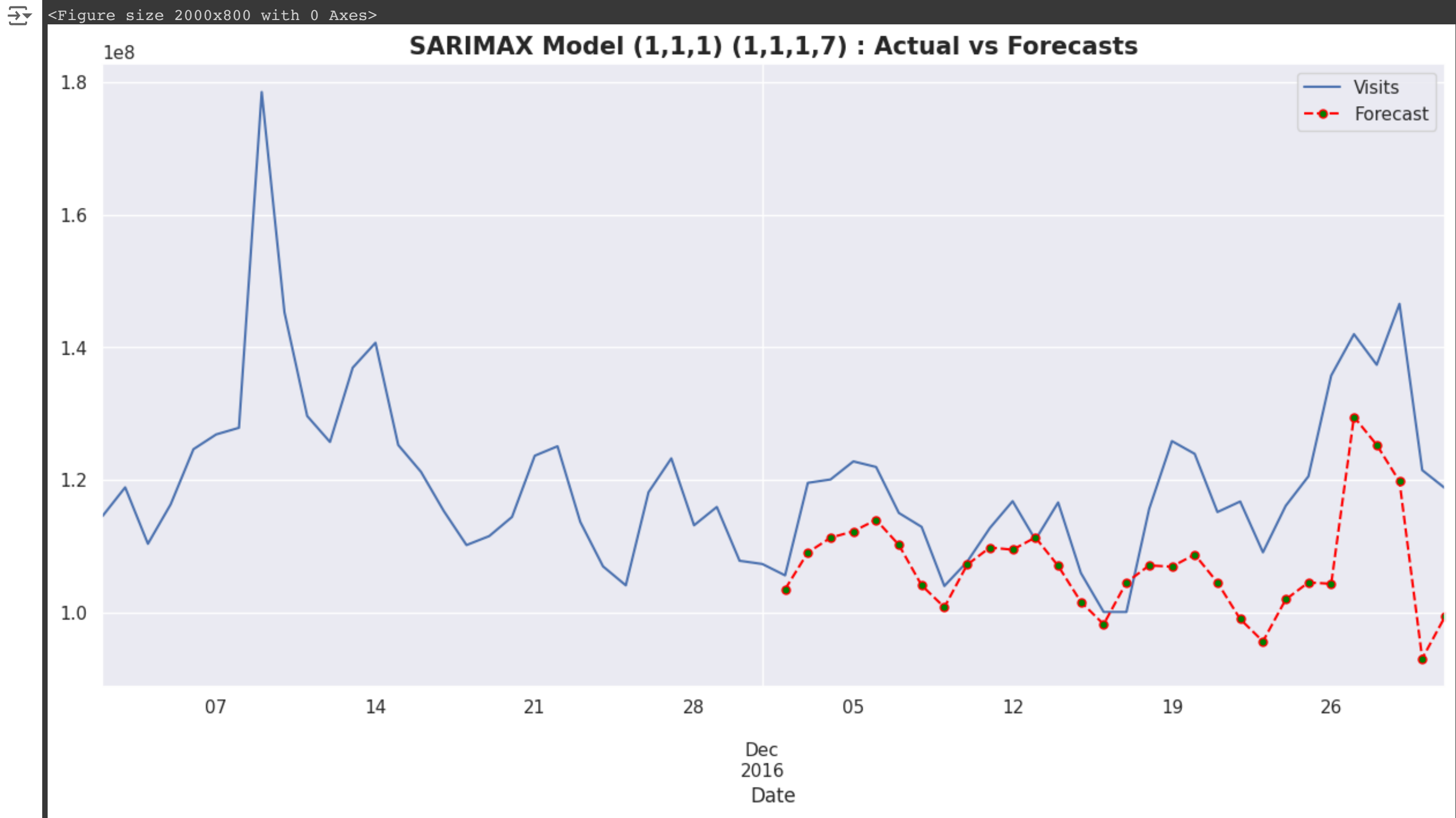
```python
time_series = ts_english
n = 30
param = [0,1,2]
d_param = [0,1]
s_param = [7]

english_params  = sarimax_grid_search(time_series, n, param, d_param,s_param,exog)
```

```
Possible Combination: 1 out of 324 calculated
Possible Combination: 2 out of 324 calculated
Possible Combination: 3 out of 324 calculated
Possible Combination: 4 out of 324 calculated
Possible Combination: 5 out of 324 calculated
Possible Combination: 6 out of 324 calculated
Possible Combination: 7 out of 324 calculated
Possible Combination: 8 out of 324 calculated
Possible Combination: 9 out of 324 calculated
Possible Combination: 10 out of 324 calculated
Possible Combination: 11 out of 324 calculated
Possible Combination: 12 out of 324 calculated
Possible Combination: 13 out of 324 calculated
Possible Combination: 14 out of 324 calculated
Possible Combination: 15 out of 324 calculated
Possible Combination: 16 out of 324 calculated
Possible Combination: 17 out of 324 calculated
Possible Combination: 18 out of 324 calculated
Possible Combination: 19 out of 324 calculated
Possible Combination: 20 out of 324 calculated
Possible Combination: 21 out of 324 calculated
Possible Combination: 22 out of 324 calculated
Possible Combination: 23 out of 324 calculated
Possible Combination: 24 out of 324 calculated
Possible Combination: 25 out of 324 calculated
Possible Combination: 26 out of 324 calculated
Possible Combination: 27 out of 324 calculated
Possible Combination: 28 out of 324 calculated
Possible Combination: 29 out of 324 calculated
Possible Combination: 30 out of 324 calculated
Possible Combination: 31 out of 324 calculated
Possible Combination: 32 out of 324 calculated
Possible Combination: 33 out of 324 calculated
Possible Combination: 34 out of 324 calculated
Possible Combination: 35 out of 324 calculated
Possible Combination: 36 out of 324 calculated
Possible Combination: 37 out of 324 calculated
Possible Combination: 38 out of 324 calculated
Possible Combination: 39 out of 324 calculated
Possible Combination: 40 out of 324 calculated
Possible Combination: 41 out of 324 calculated
Possible Combination: 42 out of 324 calculated
Possible Combination: 43 out of 324 calculated
Possible Combination: 44 out of 324 calculated
Possible Combination: 45 out of 324 calculated
Possible Combination: 46 out of 324 calculated
Possible Combination: 47 out of 324 calculated
Possible Combination: 48 out of 324 calculated
Possible Combination: 49 out of 324 calculated
Possible Combination: 50 out of 324 calculated
Possible Combination: 51 out of 324 calculated
Possible Combination: 52 out of 324 calculated
Possible Combination: 53 out of 324 calculated
Possible Combination: 54 out of 324 calculated
Possible Combination: 55 out of 324 calculated
Possible Combination: 56 out of 324 calculated
Possible Combination: 57 out of 324 calculated
Possible Combination: 58 out of 324 calculated
Possible Combination: 59 out of 324 calculated
Possible Combination: 60 out of 324 calculated
```

```
english_params.sort_values(['mape', 'rmse']).head()
```

|     | serial | pdq       | PDQs         | mape    | rmse        |
|-----|--------|-----------|--------------|---------|-------------|
| 288 | 289    | (2, 1, 1) | (0, 0, 0, 7) | 0.08737 | 1.390974e+07 |
| 289 | 290    | (2, 1, 1) | (0, 0, 1, 7) | 0.08903 | 1.411103e+07 |
| 290 | 291    | (2, 1, 1) | (0, 0, 2, 7) | 0.08988 | 1.421080e+07 |
| 294 | 295    | (2, 1, 1) | (1, 0, 0, 7) | 0.09013 | 1.423713e+07 |
| 300 | 301    | (2, 1, 1) | (2, 0, 0, 7) | 0.09273 | 1.456823e+07 |

```
time_series = ts_english
p,d,q, P,D,Q,s = 2,1,1, 0,0,0,7
n = 30
sarimax_model(time_series, n, p=p, d=d, q=q, P=P, D=D, Q=Q, s=s, exog = exog)
```

```
<Figure size 2000x800 with 0 Axes>
```



```
-------------------------------------------------------------------------
MAPE of Model : 0.08737
-------------------------------------------------------------------------
RMSE of Model : 13909735.545
-------------------------------------------------------------------------
```

```
time_series = ts_english
p,d,q, P,D,Q,s = 1,1,1, 2,1,1,7
n = 30
sarimax_model(time_series, n, p=p, d=d, q=q, P=P, D=D, Q=Q, s=s, exog = exog)
```

<Figure size 2000x800 with 0 Axes>



```
--------------------------------------------------------------------------
MAPE of Model : 0.11871
--------------------------------------------------------------------------
RMSE of Model : 18268474.479
--------------------------------------------------------------------------
```

```python
def pipeline_sarimax_grid_search_without_exog(languages, data_language, n, param, d_param, s_param):

    best_param_df = pd.DataFrame(columns=['language', 'p', 'd', 'q', 'P', 'D', 'Q', 's', 'mape'])

    for lang in languages:
        print(f'----------------------------------------------------------------')
        print(f'              Finding best parameters for {lang}                 ')
        print(f'----------------------------------------------------------------')

        time_series = data_language[data_language['language'] == lang][['Date', 'Visits']]
        time_series.set_index('Date', drop=True, inplace=True)
        best_mape = 100

        counter = 0
        param_combinations = product(param, d_param, param, param, d_param, param, s_param)

        for p, d, q, P, D, Q, s in param_combinations:
            model = SARIMAX(time_series[:-n],
                            order=(p, d, q),
                            seasonal_order=(P, D, Q, s),
                            initialization='approximate_diffuse')
            model_fit = model.fit()
            model_forecast = model_fit.forecast(n, dynamic=True)

            actuals = time_series.values[-n:]
            errors = time_series.values[-n:] - model_forecast.values
            mape = np.mean(np.abs(errors) / np.abs(actuals))

            counter += 1
            if mape < best_mape:
                best_mape = mape
                best_p, best_d, best_q = p, d, q
                best_P, best_D, best_Q = P, D, Q
                best_s = s

            print(f'Possible Combination: {counter} out of {(len(param)**4)*len(s_param)*(len(d_param)**2)} calculated')

        best_mape = np.round(best_mape, 5)
        print(f'----------------------------------------------------------------')
        print(f'Minimum MAPE for {lang} = {best_mape}')
        print(f'Corresponding Best Parameters are {best_p, best_d, best_q, best_P, best_D, best_Q, best_s}')
        print(f'----------------------------------------------------------------')

        best_param_row = [lang, best_p, best_d, best_q, best_P, best_D, best_Q, best_s, best_mape]
        best_param_df.loc[len(best_param_df)] = best_param_row

    return best_param_df
```

```python
languages = ['Chinese', 'French', 'German', 'Japenese', 'Russian', 'Spanish']
n = 30
param = [0,1,2]
d_param = [0,1]
s_param = [7]
```

```
best_param_df = pipeline_sarimax_grid_search_without_exog(languages, lang_data, n, param, d_param, s_param)
```

```
------------------------------------------------------------------
        Finding best parameters for Chinese
------------------------------------------------------------------
Possible Combination: 1 out of 324 calculated
Possible Combination: 2 out of 324 calculated
Possible Combination: 3 out of 324 calculated
Possible Combination: 4 out of 324 calculated
Possible Combination: 5 out of 324 calculated
Possible Combination: 6 out of 324 calculated
Possible Combination: 7 out of 324 calculated
Possible Combination: 8 out of 324 calculated
Possible Combination: 9 out of 324 calculated
Possible Combination: 10 out of 324 calculated
Possible Combination: 11 out of 324 calculated
Possible Combination: 12 out of 324 calculated
Possible Combination: 13 out of 324 calculated
Possible Combination: 14 out of 324 calculated
Possible Combination: 15 out of 324 calculated
Possible Combination: 16 out of 324 calculated
Possible Combination: 17 out of 324 calculated
Possible Combination: 18 out of 324 calculated
Possible Combination: 19 out of 324 calculated
Possible Combination: 20 out of 324 calculated
Possible Combination: 21 out of 324 calculated
Possible Combination: 22 out of 324 calculated
Possible Combination: 23 out of 324 calculated
Possible Combination: 24 out of 324 calculated
Possible Combination: 25 out of 324 calculated
Possible Combination: 26 out of 324 calculated
Possible Combination: 27 out of 324 calculated
Possible Combination: 28 out of 324 calculated
Possible Combination: 29 out of 324 calculated
Possible Combination: 30 out of 324 calculated
Possible Combination: 31 out of 324 calculated
Possible Combination: 32 out of 324 calculated
Possible Combination: 33 out of 324 calculated
Possible Combination: 34 out of 324 calculated
Possible Combination: 35 out of 324 calculated
Possible Combination: 36 out of 324 calculated
Possible Combination: 37 out of 324 calculated
Possible Combination: 38 out of 324 calculated
Possible Combination: 39 out of 324 calculated
Possible Combination: 40 out of 324 calculated
Possible Combination: 41 out of 324 calculated
Possible Combination: 42 out of 324 calculated
Possible Combination: 43 out of 324 calculated
Possible Combination: 44 out of 324 calculated
Possible Combination: 45 out of 324 calculated
Possible Combination: 46 out of 324 calculated
Possible Combination: 47 out of 324 calculated
Possible Combination: 48 out of 324 calculated
Possible Combination: 49 out of 324 calculated
Possible Combination: 50 out of 324 calculated
Possible Combination: 51 out of 324 calculated
Possible Combination: 52 out of 324 calculated
Possible Combination: 53 out of 324 calculated
Possible Combination: 54 out of 324 calculated
Possible Combination: 55 out of 324 calculated
Possible Combination: 56 out of 324 calculated
Possible Combination: 57 out of 324 calculated
```

```
best_param_df.sort_values(['mape'], inplace = True)
best_param_df
```

|   | language | p | d | q | P | D | Q | s | mape |
|---|----------|---|---|---|---|---|---|---|------|
| 0 | Chinese | 2 | 1 | 0 | 0 | 0 | 0 | 7 | 0.03932 |
| 4 | Russian | 0 | 0 | 1 | 2 | 0 | 0 | 7 | 0.06795 |
| 1 | French | 0 | 1 | 2 | 0 | 0 | 0 | 7 | 0.08091 |
| 2 | German | 2 | 0 | 1 | 0 | 0 | 0 | 7 | 0.08384 |
| 3 | Japenese | 1 | 1 | 2 | 0 | 0 | 1 | 7 | 0.09340 |
| 5 | Spanish | 2 | 0 | 0 | 0 | 0 | 1 | 7 | 0.14351 |

```python
def plot_best_SARIMAX_model(languages, data_language, n, best_param_df):
    for lang in languages:

        params_lang = best_param_df[best_param_df['language'] == lang].iloc[0]
        p, d, q, P, D, Q, s = params_lang[['p', 'd', 'q', 'P', 'D', 'Q', 's']]

        time_series = data_language[data_language['language'] == lang][['Date', 'Visits']]
        time_series.set_index('Date', drop=True, inplace=True)

        model = SARIMAX(time_series[:-n], order=(p, d, q),
                        seasonal_order=(P, D, Q, s), initialization='approximate_diffuse')
        model_fit = model.fit()

        model_forecast = model_fit.forecast(n, dynamic=True)

        actuals = time_series.values[-n:]
        errors = time_series.values[-n:] - model_forecast.values
        mape = np.mean(np.abs(errors) / np.abs(actuals))
        rmse = np.sqrt(np.mean(errors**2))

        print(f'\n{"-" * 90}')
        print(f'SARIMAX model for {lang} Time Series')
        print(f'Parameters of Model: ({p}, {d}, {q}) ({P}, {D}, {Q}, {s})')
        print(f'MAPE of Model: {np.round(mape, 5)}')
        print(f'RMSE of Model: {np.round(rmse, 3)}')
        print(f'{"-" * 90}')

        time_series.index = time_series.index.astype('datetime64[ns]')
        model_forecast.index = model_forecast.index.astype('datetime64[ns]')
        plt.figure(figsize=(20, 8))
        time_series[-60:].plot(label='Actual')
        model_forecast[-60:].plot(label='Forecast', color='red',
                                  linestyle='dashed', marker='o', markerfacecolor='green', markersize=5)
        plt.legend(loc="upper right")
        plt.title(f'SARIMAX Model ({p}, {d}, {q}) ({P}, {D}, {Q}, {s}): Actual vs Forecasts',
                  fontsize=15, fontweight='bold')
        plt.show()

    return 0
```
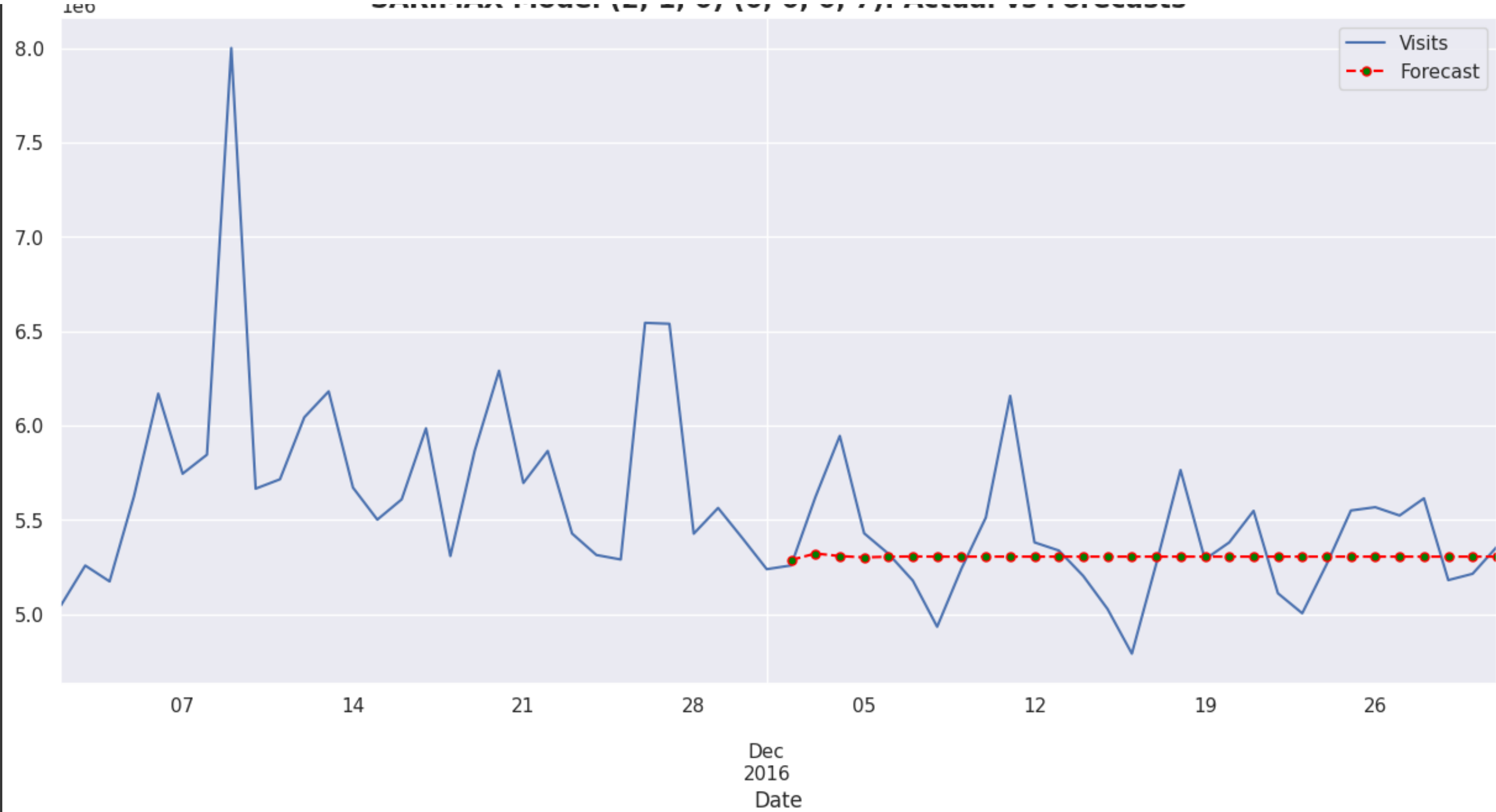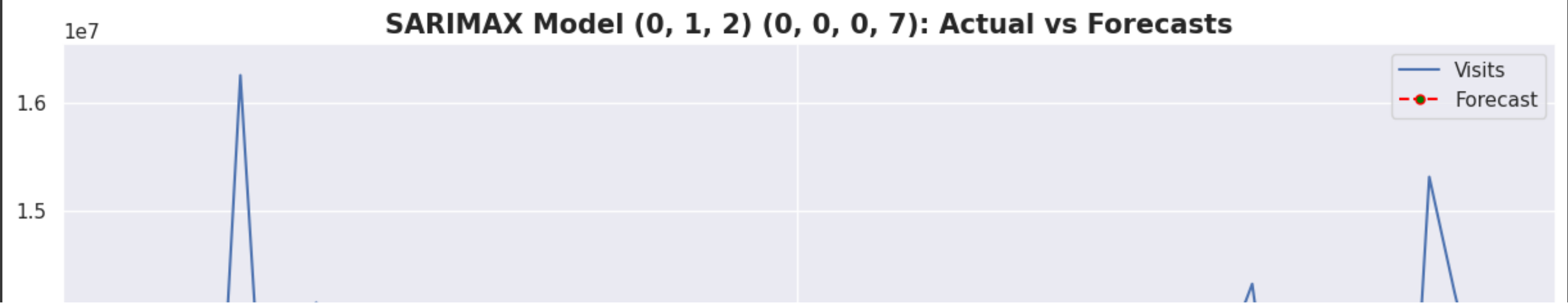
```python
languages = ['Chinese', 'French', 'German', 'Japenese', 'Russian', 'Spanish']
n = 30
plot_best_SARIMAX_model(languages, lang_data, n, best_param_df)
```
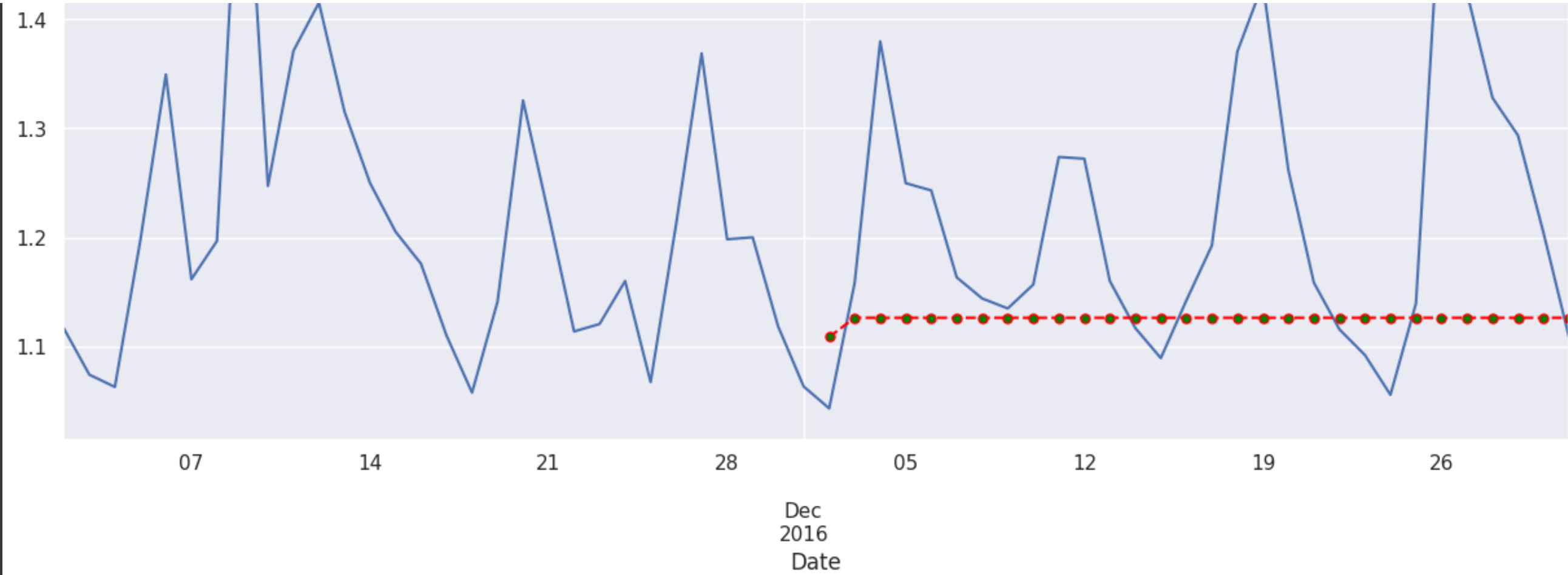
```
--------------------------------------------------------------------------------
SARIMAX model for Chinese Time Series
Parameters of Model: (2, 1, 0) (0, 0, 0, 7)
MAPE of Model: 0.03932
RMSE of Model: 289943.436
--------------------------------------------------------------------------------
<Figure size 2000x800 with 0 Axes>
```

**SARIMAX Model (2, 1, 0) (0, 0, 0, 7): Actual vs Forecasts**

SARIMAX Model (2, 1, 0) (0, 0, 0, 7): Actual vs Forecasts

```
--------------------------------------------------------------------------------
SARIMAX model for French Time Series
Parameters of Model: (0, 1, 2) (0, 0, 0, 7)
MAPE of Model: 0.08091
RMSE of Model: 1489350.009
--------------------------------------------------------------------------------
<Figure size 2000x800 with 0 Axes>
```



SARIMAX Model (0, 1, 2) (0, 0, 0, 7): Actual vs Forecasts

```
--------------------------------------------------------------------------------
SARIMAX model for German Time Series
Parameters of Model: (2, 0, 1) (0, 0, 0, 7)
MAPE of Model: 0.08384
RMSE of Model: 2195114.679
--------------------------------------------------------------------------------
<Figure size 2000x800 with 0 Axes>
```
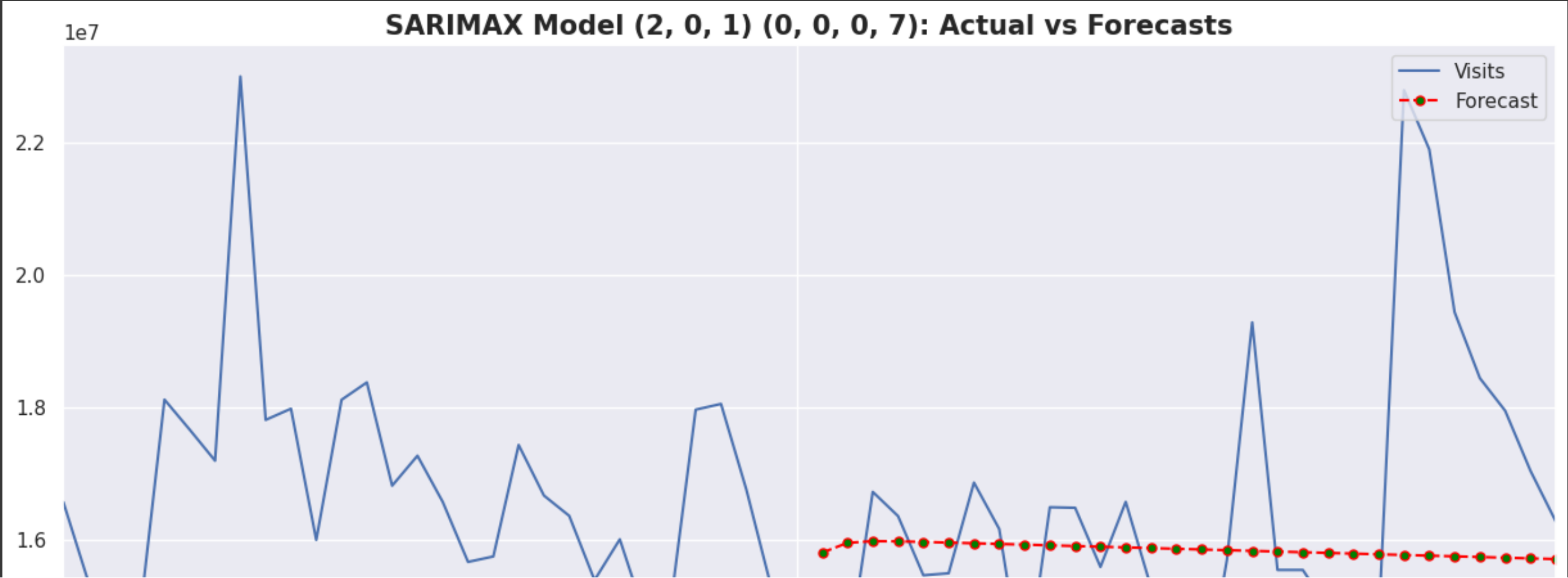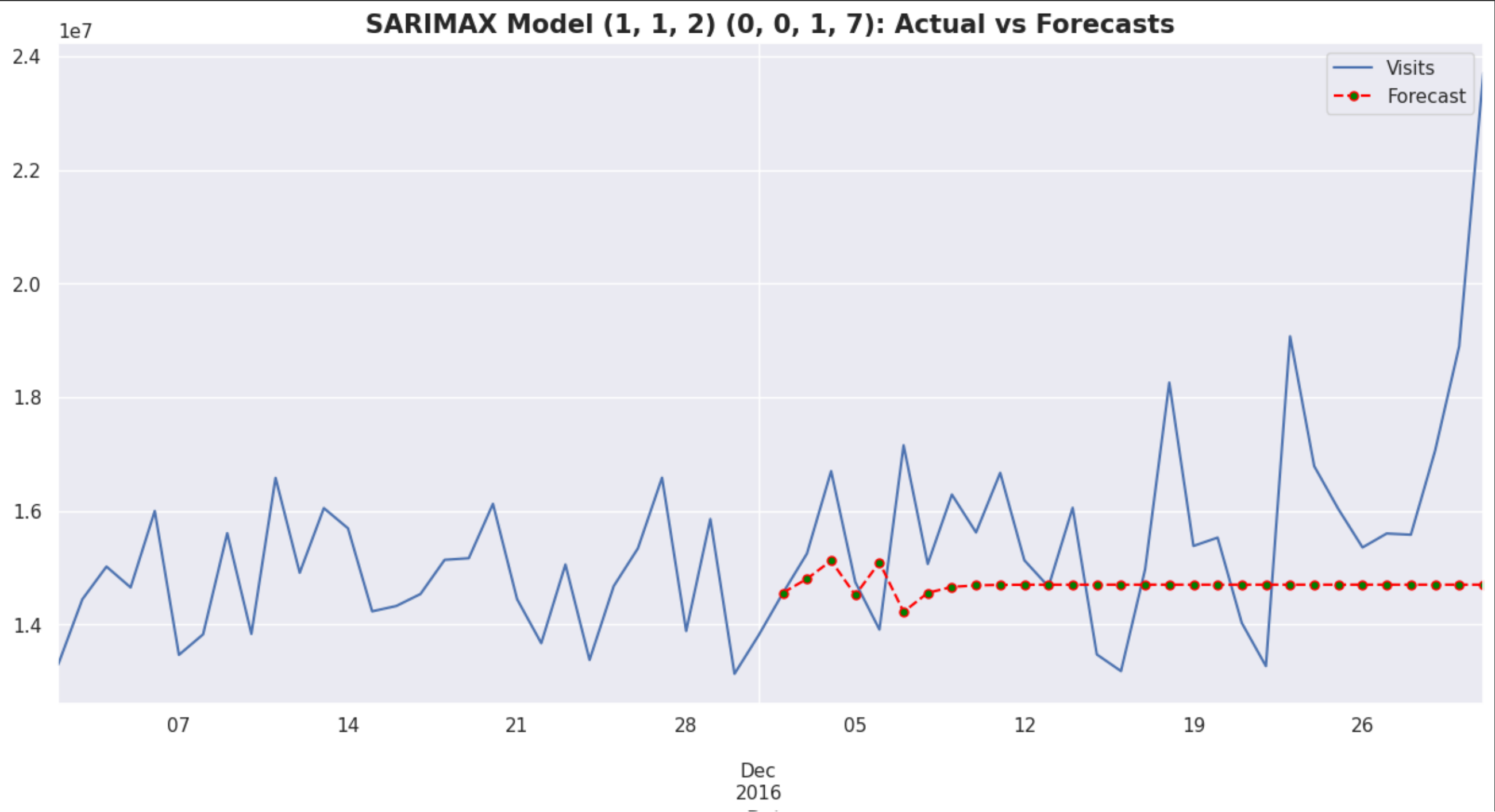


SARIMAX Model (2, 0, 1) (0, 0, 0, 7): Actual vs Forecasts

```
--------------------------------------------------------------------------------
SARIMAX model for Japenese Time Series
Parameters of Model: (1, 1, 2) (0, 0, 1, 7)
MAPE of Model: 0.0934
RMSE of Model: 2400870.834
--------------------------------------------------------------------------------
<Figure size 2000x800 with 0 Axes>
```
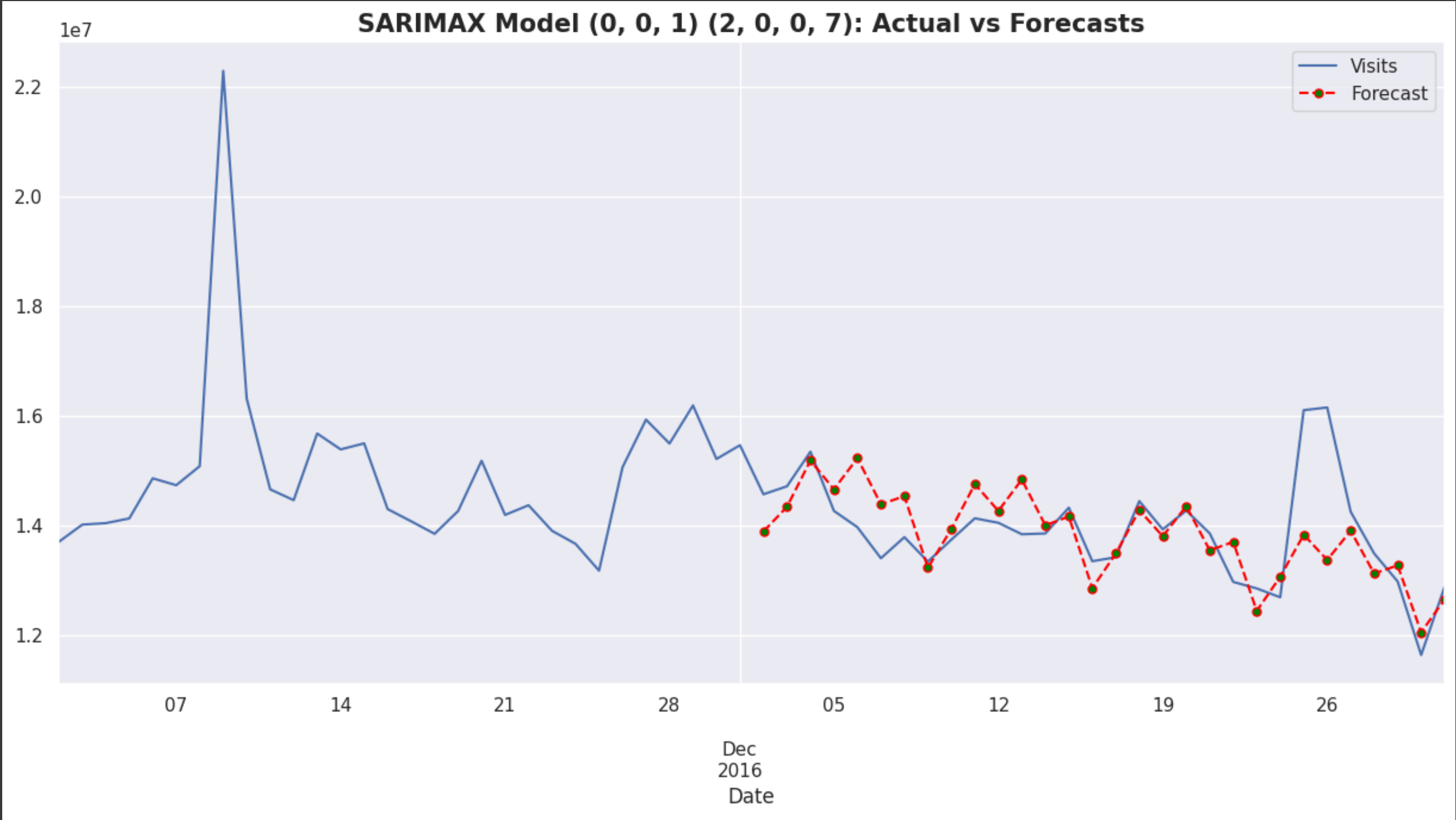
Date

```
--------------------------------------------------------------------------
SARIMAX model for Russian Time Series
Parameters of Model: (0, 0, 1) (2, 0, 0, 7)
MAPE of Model: 0.06795
RMSE of Model: 1206324.353
--------------------------------------------------------------------------
<Figure size 2000x800 with 0 Axes>
```
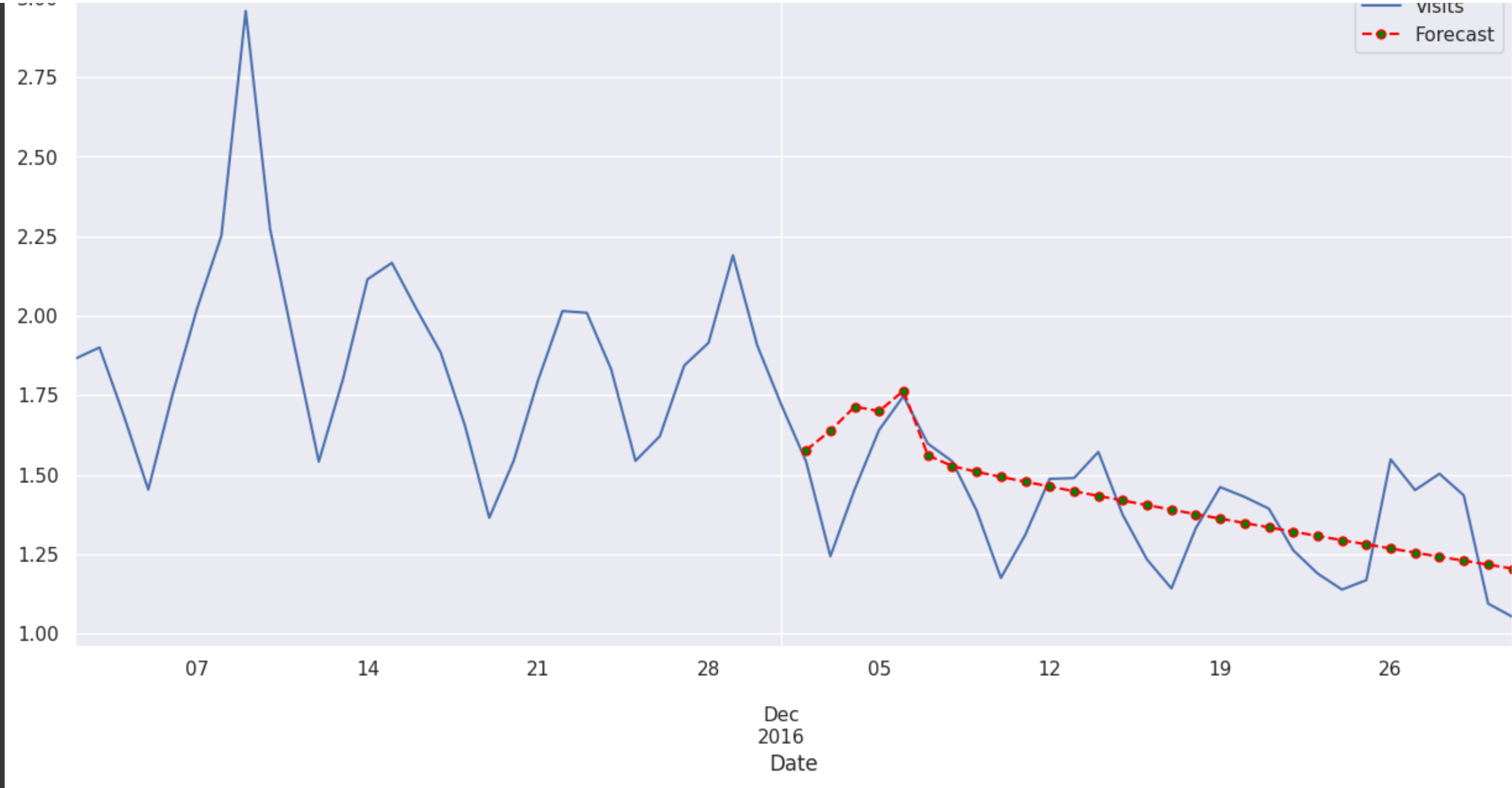


SARIMAX Model (0, 0, 1) (2, 0, 0, 7): Actual vs Forecasts

```
--------------------------------------------------------------------------
SARIMAX model for Spanish Time Series
Parameters of Model: (2, 0, 0) (0, 0, 1, 7)
MAPE of Model: 0.14351
RMSE of Model: 2344695.867
--------------------------------------------------------------------------
<Figure size 2000x800 with 0 Axes>
```

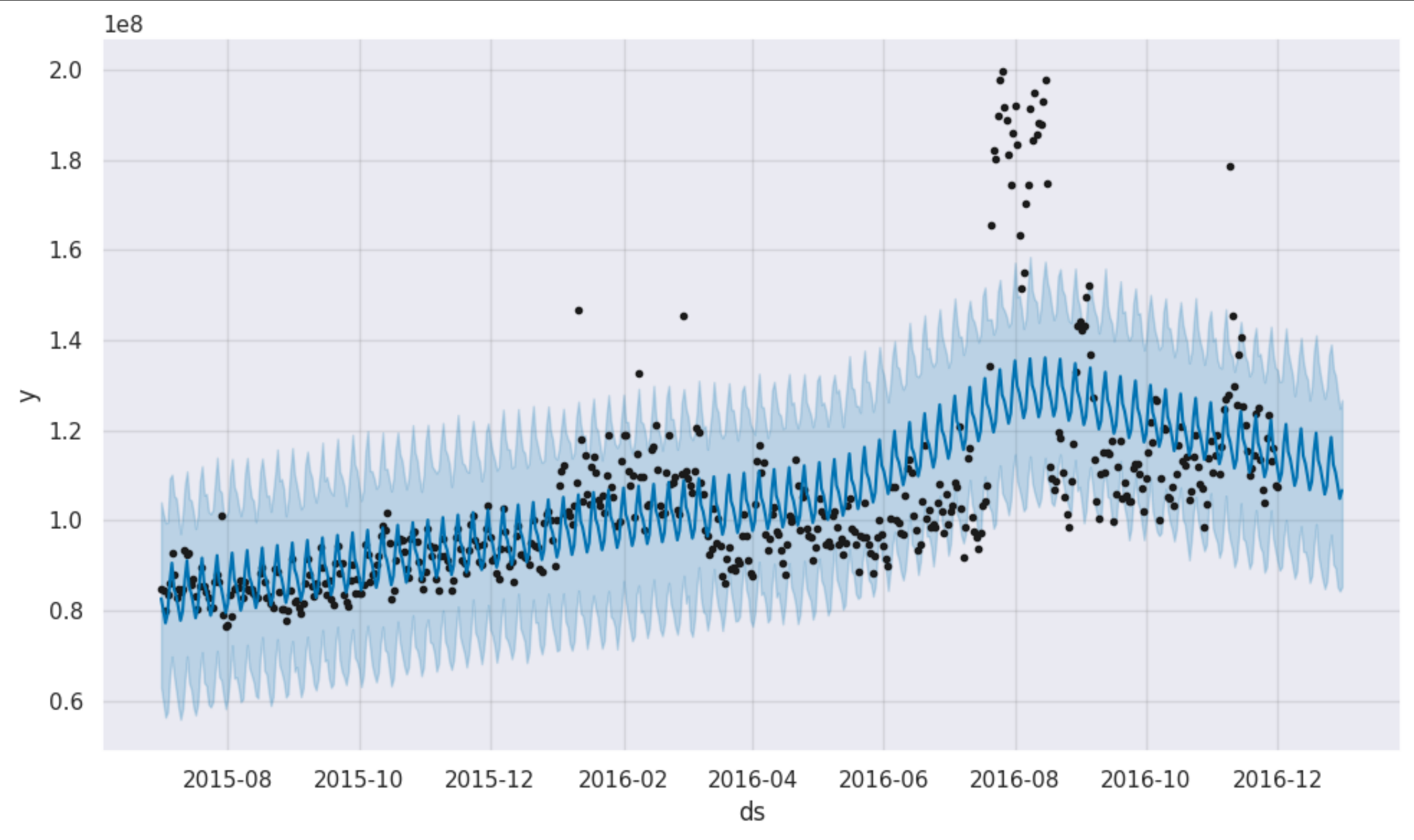SARIMAX Model (2, 0, 0) (0, 0, 1, 7): Actual vs Forecasts

```
time_series = lang_data[lang_data['language'] == 'English'][['Date', 'Visits']]
time_series.columns = ['ds', 'y']
time_series['exog'] = exog
```
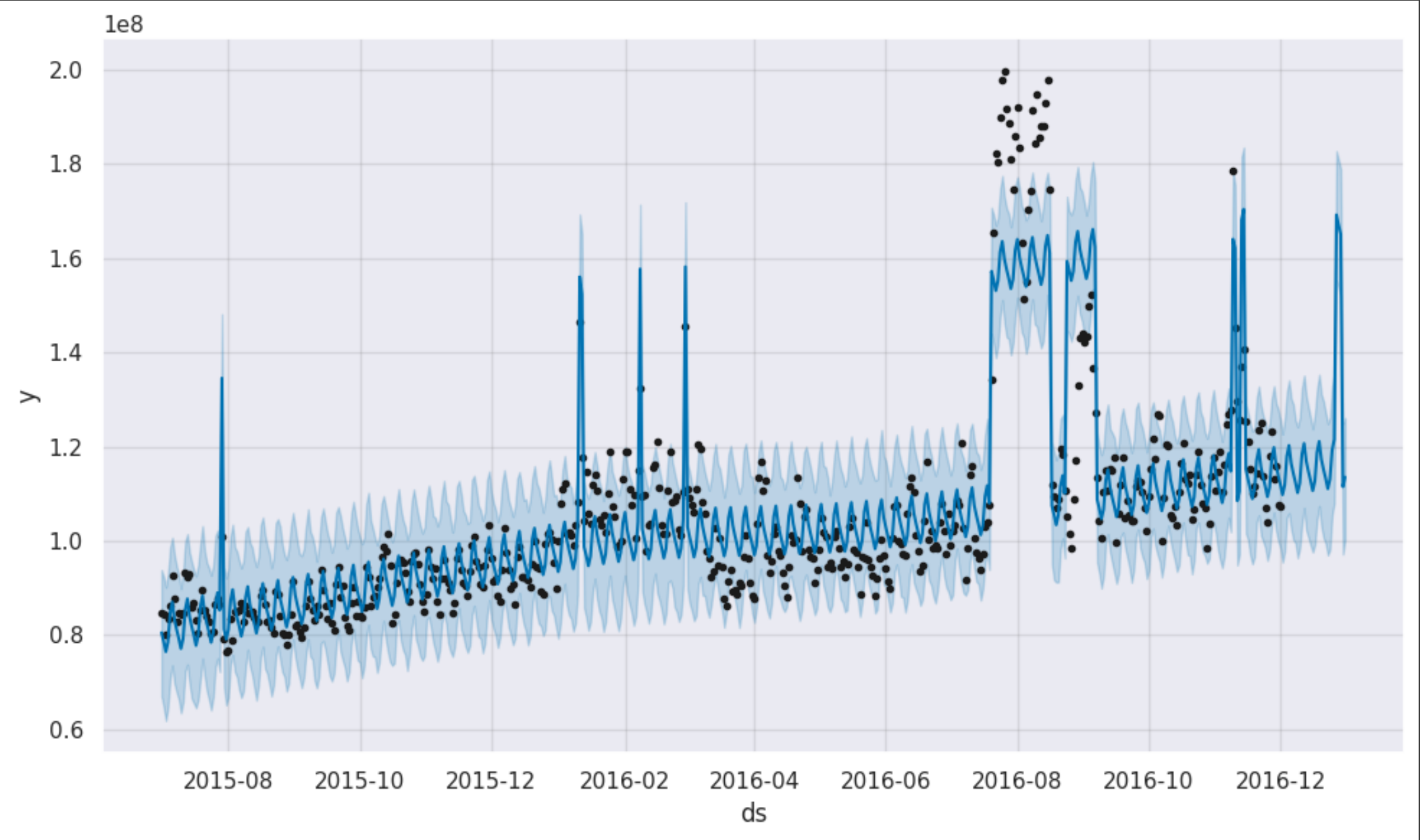
```
prophet1 = Prophet(weekly_seasonality=True)
prophet1.fit(time_series[['ds', 'y']][:-30])
future = prophet1.make_future_dataframe(periods=30, freq= 'D')
forecast = prophet1.predict(future)
fig1 = prophet1.plot(forecast)
```

```
INFO:prophet:Disabling yearly seasonality. Run prophet with yearly_seasonality=True to override this.
INFO:prophet:Disabling daily seasonality. Run prophet with daily_seasonality=True to override this.
DEBUG:cmdstanpy:input tempfile: /tmp/tmp1xp0bq04/b6x0pbno.json
DEBUG:cmdstanpy:input tempfile: /tmp/tmp1xp0bq04/s4lc97b2.json
DEBUG:cmdstanpy:idx 0
DEBUG:cmdstanpy:running CmdStan, num_threads: None
DEBUG:cmdstanpy:CmdStan args: ['/usr/local/lib/python3.10/dist-packages/prophet/stan_model/prophet_model.bin', 'random', 'seed=63098', 'data', 'file=/tmp/tmp1xp0bq04/b6x0pbno.json', 'init=/t
05:40:09 - cmdstanpy - INFO - Chain [1] start processing
INFO:cmdstanpy:Chain [1] start processing
05:40:09 - cmdstanpy - INFO - Chain [1] done processing
INFO:cmdstanpy:Chain [1] done processing
```
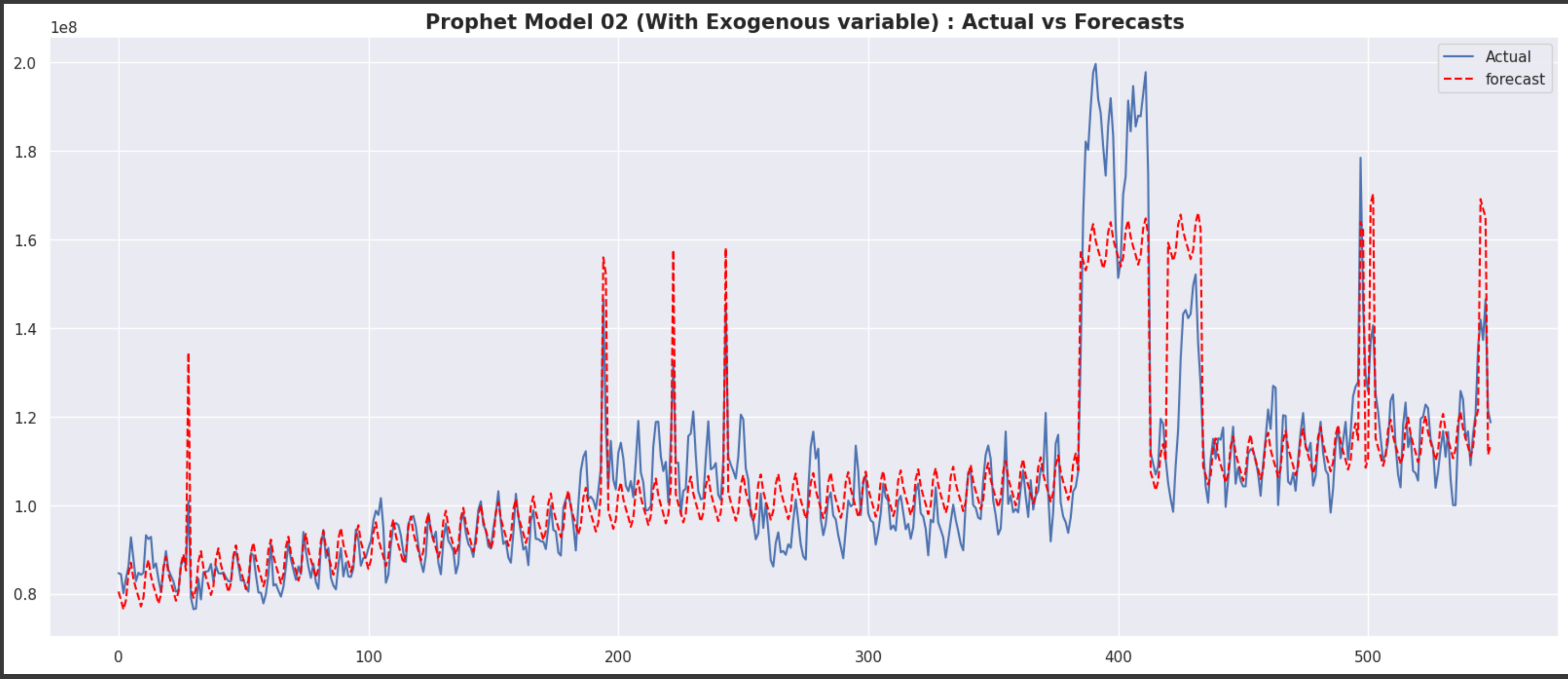
```
prophet2 = Prophet(weekly_seasonality=True)
prophet2.add_regressor('exog')
prophet2.fit(time_series[:-30])
forecast2 = prophet2.predict(time_series)
fig2 = prophet2.plot(forecast2)
```

```
INFO:prophet:Disabling yearly seasonality. Run prophet with yearly_seasonality=True to override this.
INFO:prophet:Disabling daily seasonality. Run prophet with daily_seasonality=True to override this.
DEBUG:cmdstanpy:input tempfile: /tmp/tmp1xp0bq04/thc_ljpf.json
DEBUG:cmdstanpy:input tempfile: /tmp/tmp1xp0bq04/k5bkuqoi.json
DEBUG:cmdstanpy:idx 0
DEBUG:cmdstanpy:running CmdStan, num_threads: None
DEBUG:cmdstanpy:CmdStan args: ['/usr/local/lib/python3.10/dist-packages/prophet/stan_model/prophet_model.bin', 'random', 'seed=88966', 'data', 'file=/tmp/tmp1xp0bq04/thc_ljpf.json', 'init=/t
05:40:11 - cmdstanpy - INFO - Chain [1] start processing
INFO:cmdstanpy:Chain [1] start processing
05:40:11 - cmdstanpy - INFO - Chain [1] done processing
INFO:cmdstanpy:Chain [1] done processing
```

```python
actual = time_series['y'].values
forecast = forecast2['yhat'].values

plt.figure(figsize = (20,8))
plt.plot(actual, label = 'Actual')
plt.plot(forecast, label = 'forecast', color = 'red', linestyle='dashed')
plt.legend(loc="upper right")
plt.title(f'Prophet Model 02 (With Exogenous variable) : Actual vs Forecasts', fontsize = 15, fontweight = 'bold')
plt.show()
```

```
errors = abs(actual - forecast)
mape = np.mean(errors/abs(actual))
mape
```

↳  0.05983786254333203

FB Prophet Model is able to capture peaks because of exogenous variable and is giving a MAPE of 6%

Recommendations Prioritize English language pages due to their low MAPE and high mean visits, making them optimal for advertising efforts to maximize reach and effectiveness.

Avoid advertising on Chinese language pages unless there's a specific marketing strategy tailored for Chinese populations, as they have the lowest number of visits.

Russian language pages present a promising opportunity for high conversion rates with their decent number of visits and low MAPE if utilized effectively.

Despite having the second-highest number of visits, Spanish language pages exhibit the highest MAPE, suggesting that advertisements on these pages may not effectively reach the intended audience.

French, German, and Japanese language pages show moderate levels of visits and MAPE. Depending on the target customers, consider advertising campaigns on these pages to capitalize on their potential reach and conversion rates.

Questionnaire Defining the problem statements and where can this and modifications of this be used? The Data Science team at Ad ease aims to analyze per page view reports for various Wikipedia pages spanning 550 days. The objective includes forecasting page views to enhance ad placement optimization for clients. Dataset encompasses 145k Wikipedia pages with daily view counts. Client base extends across diverse regions, necessitating insights into ad performance across different languages. Importance of forecasting model:

Identification of the problem and its applications:

Implementing a robust forecasting model is pivotal in predicting fluctuations in page visits. This model aids the business team in optimizing marketing expenditure. Precise prediction of high-traffic days enables strategic ad placement, maximizing audience reach while optimizing spending. Write 3 inferences you made from the data visualizations. Linguistic Diversity: The data reveals the presence of 7 languages, with English dominating, followed by Japanese, German, and French. Access Type Distribution: Three access types are identified—All-access, mobile-web, and desktop—comprising 51.4%, 24.9%, and 23.6% respectively. Access-Origin Insights: The dataset illustrates two access origins —'all-agents' and 'spider'—with 'all-agents' constituting 75.8% and 'spider' 24.2% of the data. Advertising Strategies:

Inferences from Data Visualizations:

English Language Dominance: English emerges as the most prominent language, suggesting prioritized advertisement placement due to its low Mean Absolute Percentage Error (MAPE) and high mean visit count. Chinese Language Considerations: Pages in Chinese exhibit the lowest visit counts, signaling caution in advertisement allocation unless specifically targeting Chinese demographics. Russian Language Potential: Russian language pages demonstrate a favorable balance between visit count and MAPE, indicating potential for maximum conversion if utilized effectively. Spanish Language Challenges: Despite being the second-highest in visit count, Spanish pages exhibit the highest MAPE, suggesting potential challenges in advertisement efficacy. Moderate Performers: French, German, and Japanese languages present medium-level visit counts and MAPE levels, prompting tailored advertisement strategies based on target customer demographics. Time Series Decomposition

What does the decomposition of series do? Time series decomposition is a statistical technique used to break down a time series into its constituent components in order to understand its underlying structure, trends, seasonality, and irregular fluctuations. The decomposition typically involves separating the time series data into three main components:

Trend (($T\_t$)): The long-term movement or pattern in the data, representing the overall direction in which the time series is moving.

Seasonality (($S\_t$)): The repeating patterns or fluctuations that occur at regular intervals within the time series data.

Residuals (($R\_t$)): The remaining variation in the data after removing the trend and seasonality components.

The time series ($y\_t$) can be decomposed into its components as follows:

Additive Decomposition: [ y_t = T_t + S_t + R_t ]

Multiplicative Decomposition: [ y_t = T_t \times S_t \times R_t ]

Various techniques such as moving averages, exponential smoothing, or mathematical models can be used to estimate the trend and seasonal components, leaving the residual component as the leftover variation in the data. What level of differencing gave you a stationary series? First order differencing Difference between arima, sarima & sarimax. ARIMA is a time series forecasting model that combines autoregression (AR), differencing (I), and moving average (MA) components. It's suitable for univariate time series data without exogenous variables. ARIMA(p,d,q) where p represents the autoregressive order, d represents the differencing order, and q represents the moving average order. SARIMA is an extension of ARIMA that incorporates seasonal components in addition to the non-seasonal ones. It's suitable for time series data with seasonal patterns. SARIMA(p,d,q)(P,D,Q)m where P, D, and Q represent the seasonal autoregressive, differencing, and moving average orders respectively, and 'm' represents the seasonal period. SARIMAX (Seasonal Autoregressive Integrated Moving Average with Exogenous Variables):
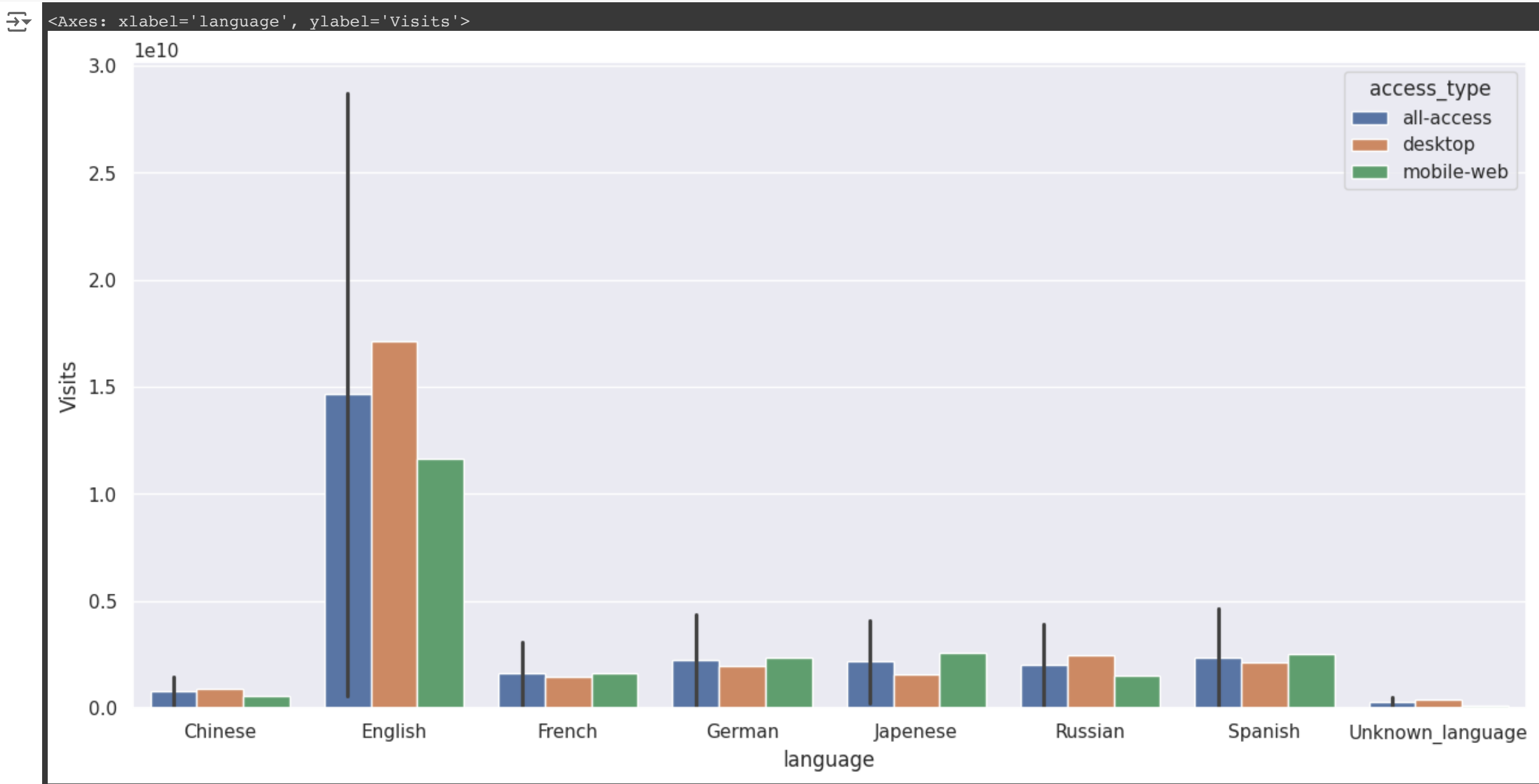
SARIMA (Seasonal Autoregressive Integrated Moving Average):

ARIMA (Autoregressive Integrated Moving Average):

SARIMAX extends SARIMA by allowing the inclusion of exogenous variables, which are external factors that can influence the time series. It's suitable for time series data with both seasonal patterns and external variables. SARIMAX(p,d,q)(P,D,Q)m with exogenous variables. These models are commonly used in time series analysis and forecasting tasks, each offering different capabilities to handle various types of data and patterns.

```
grouped = reshaped.groupby(['language','access_type','access_origin'], as_index=False)['Visits'].sum()
```

```
sns.barplot(grouped, x="language", y="Visits", hue="access_type")
```

<Axes: xlabel='language', ylabel='Visits'>



What other methods other than grid search would be suitable to get the model for all languages?

We can use packages like hyperopt, optuna and sci-kit-optimize

We can try and use different models like tsmixer and deep learning models

Start coding or generate with AI.