

# Multi-Channel Website Performance Forecasting and Engagement Analysis

## Introduction

This project applies time series analysis to website performance data, exploring hourly user traffic, sessions, engagement, and conversions across channels such as Direct, Organic Search, Organic Social, and Referrals. The dataset was obtained from the [Statso Website Performance Case Study](#).

By uncovering patterns, forecasting trends, and detecting anomalies with techniques such as ARIMA and seasonal decomposition, the analysis provides actionable insights to optimize user experience, channel effectiveness, and strategic decision-making.

## Primary Objective:

- Build and validate time-series forecasting models to predict sessions
- Analyze engagement patterns and trends over time
- Identify seasonal patterns and peak performance hours
- Compare performance across different marketing channels

```
In [24]: import sys
print("Python executable:", sys.executable)
print("Python path:", sys.path[0])
```

Python executable: C:\Users\DELL\miniconda3\envs\timeseries\_env\python.exe  
Python path: C:\Users\DELL\miniconda3\envs\timeseries\_env\python310.zip

```
In [ ]:
```

```
In [ ]:
```

```
In [1]: # Test all imports
import pandas as pd
import numpy as np
import matplotlib as plt
import seaborn as sns
import scipy
import statsmodels

print(f"NumPy: {np.__version__}")
print(f"SciPy: {scipy.__version__}")
print(f"Statsmodels: {statsmodels.__version__}")

# Test your original problematic imports
from statsmodels.tsa.seasonal import seasonal_decompose
from statsmodels.tsa.stattools import adfuller
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
from statsmodels.tsa.arima.model import ARIMA
from sklearn.metrics import mean_squared_error, mean_absolute_error
from sklearn.preprocessing import StandardScaler

print("✓ All imports successful!")
```

NumPy: 1.26.4  
SciPy: 1.15.2  
Statsmodels: 0.14.5  
✓ All imports successful!

```
In [23]: import matplotlib.pyplot as plt
import seaborn as sns

# Set style for plots
plt.style.use('seaborn-v0_8')
sns.set_palette("husl")
```

## Data Loading and Initial Inspection

```
In [3]: df = pd.read_csv(r"C:\Users\DELL\Desktop\DATASETS\Time series Dataset\website Performance dataset.csv")
```

```
In [4]: df.head()
```

Out[4]:

	Session primary channel group (Default channel group)	Date + hour (YYYYMMDDHH)	Users	Sessions	Engaged sessions	Average engagement time per session	Engaged sessions per user	Events per session	Engagement rate	Event count
0	Direct	2024041623	237	300	144	47.526667	0.607595	4.673333	0.480000	1402
1	Organic Social	2024041719	208	267	132	32.097378	0.634615	4.295880	0.494382	1147
2	Direct	2024041723	188	233	115	39.939914	0.611702	4.587983	0.493562	1069
3	Organic Social	2024041718	187	256	125	32.160156	0.668449	4.078125	0.488281	1044
4	Organic Social	2024041720	175	221	112	46.918552	0.640000	4.529412	0.506787	1001

```
In [ ]:
```

```
In [5]: # Display basic information about the dataset

print("Dataset Shape:", df.shape)
print("\nColumn Names:")
print(df.columns.tolist())
print("\nData Types:")
print(df.dtypes)
print("\nMissing Values:")
print(df.isnull().sum())
```

Dataset Shape: (3182, 10)

Column Names:  
['Session primary channel group (Default channel group)', 'Date + hour (YYYYMMDDHH)', 'Users', 'Sessions', 'Engaged sessions', 'Average engagement time per session', 'Engaged sessions per user', 'Events per session', 'Engagement rate', 'Event count']

Data Types:  
Session primary channel group (Default channel group)      object  
Date + hour (YYYYMMDDHH)                                      int64  
Users    int64  
Sessions     int64  
Engaged sessions    int64  
Average engagement time per session                        float64  
Engaged sessions per user                                      float64  
Events per session     float64  
Engagement rate     float64  
Event count    int64  
dtype: object

Missing Values:  
Session primary channel group (Default channel group)      0  
Date + hour (YYYYMMDDHH)                                      0  
Users    0  
Sessions     0  
Engaged sessions    0  
Average engagement time per session                        0  
Engaged sessions per user                                      0  
Events per session     0  
Engagement rate     0  
Event count     0  
dtype: int64

```
In [ ]:
```

## Data Cleaning and Transformation

```
In [6]: # Rename column:

df = df.rename(columns={"Session primary channel group (Default channel group)": "Session primary channel group"})

df = df.rename(columns={"Date + hour (YYYYMMDDHH)": "Datetime"})
```

```
In [7]: df.head()
```

Out[7]:

	Session primary channel group	Datetime	Users	Sessions	Engaged sessions	Average engagement time per session	Engaged sessions per user	Events per session	Engagement rate	Event count
0	Direct	2024041623	237	300	144	47.526667	0.607595	4.673333	0.480000	1402
1	Organic Social	2024041719	208	267	132	32.097378	0.634615	4.295880	0.494382	1147
2	Direct	2024041723	188	233	115	39.939914	0.611702	4.587983	0.493562	1069
3	Organic Social	2024041718	187	256	125	32.160156	0.668449	4.078125	0.488281	1044
4	Organic Social	2024041720	175	221	112	46.918552	0.640000	4.529412	0.506787	1001

In [ ]:

In [8]:

```
# Convert the date column to proper datetime format

df['Datetime'] = df['Datetime'].astype(str)
df['datetime'] = pd.to_datetime(df['Datetime'], format='%Y%m%d%H')

# Set datetime as index
df.set_index('datetime', inplace=True)
```

In [9]:

```
df.head(2)
```

Out[9]:

	Session primary channel group	Datetime	Users	Sessions	Engaged sessions	Average engagement time per session	Engaged sessions per user	Events per session	Engagement rate	Event count
datetime										
2024-04-16 23:00:00	Direct	2024041623	237	300	144	47.526667	0.607595	4.673333	0.480000	1402
2024-04-17 19:00:00	Organic Social	2024041719	208	267	132	32.097378	0.634615	4.295880	0.494382	1147

In [10]:

```
# Check for duplicates

print("Duplicate rows:", df.duplicated().sum())
```

Duplicate rows: 0

In [ ]:

In [ ]:

In [33]:

```
# Check for outliers

# Select numerical columns
num_cols = [
    "Users", "Sessions", "Engaged sessions",
    "Average engagement time per session",
    "Engaged sessions per user", "Events per session",
    "Engagement rate", "Event count"
]

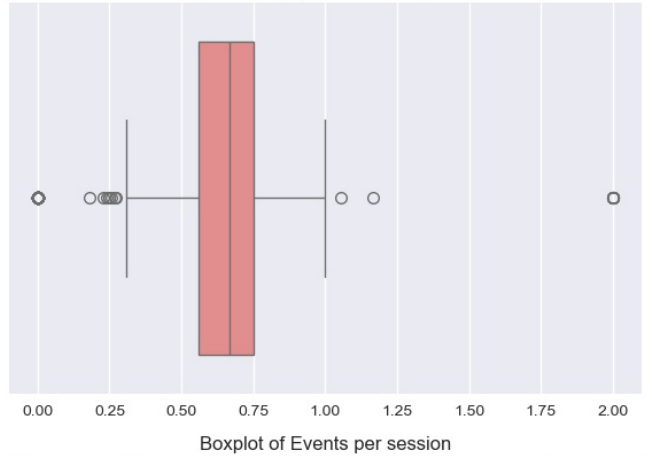
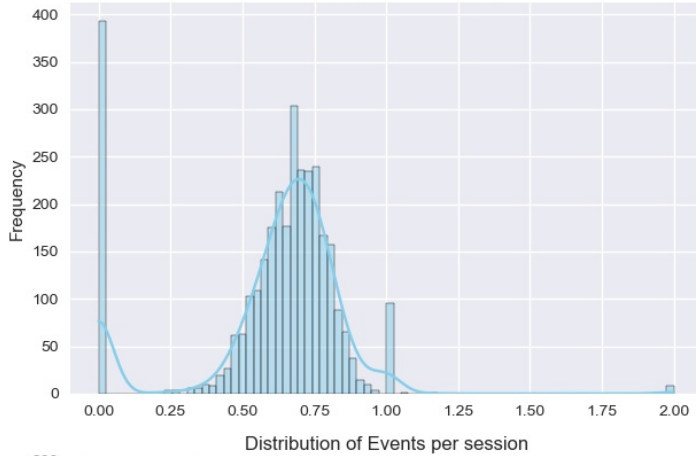
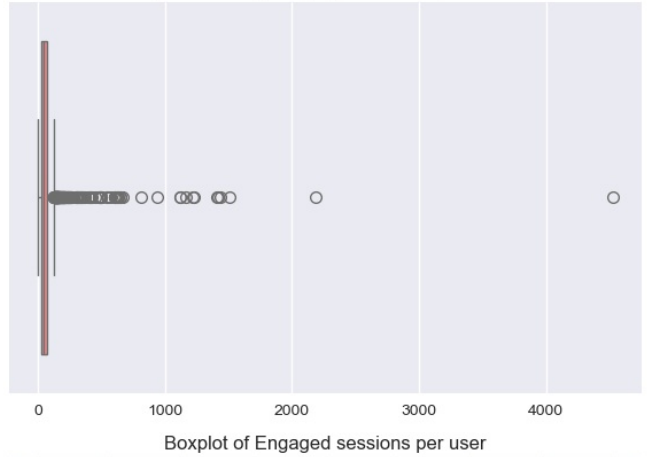
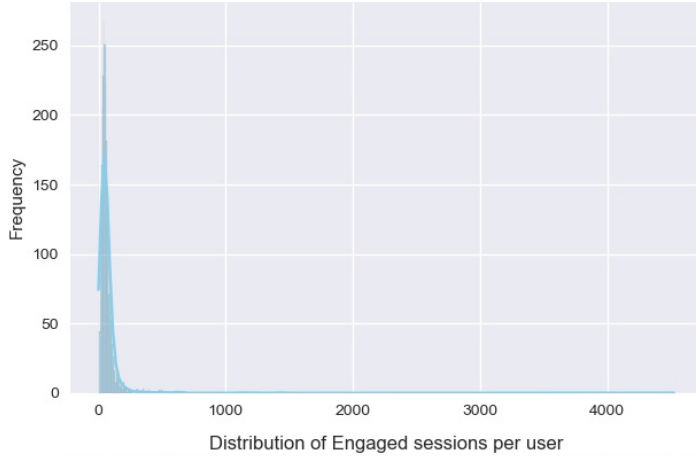
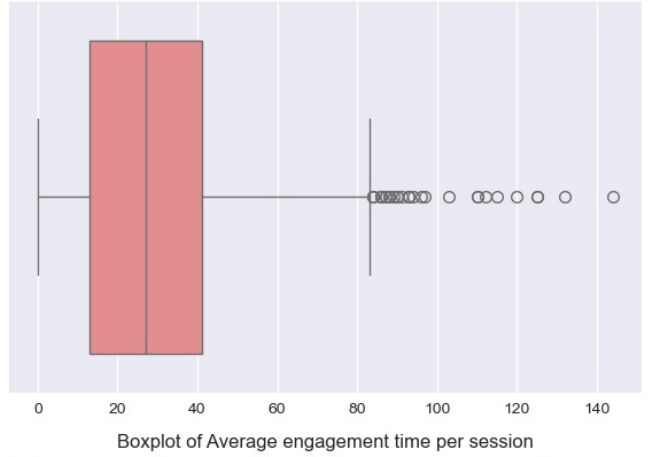
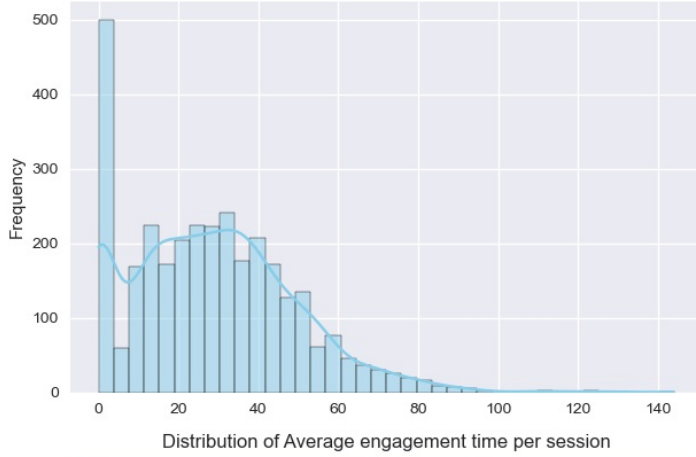
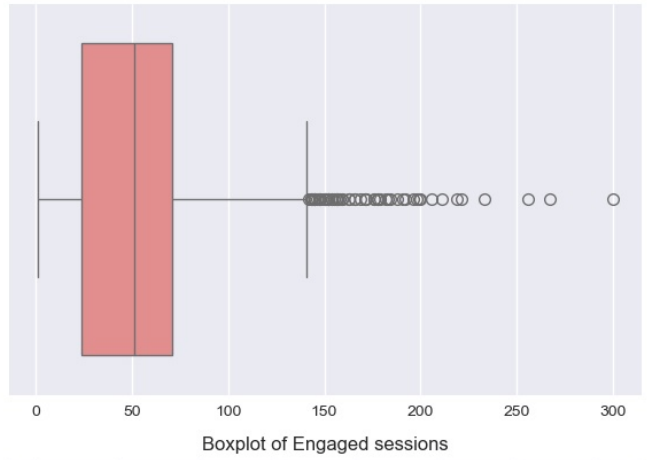
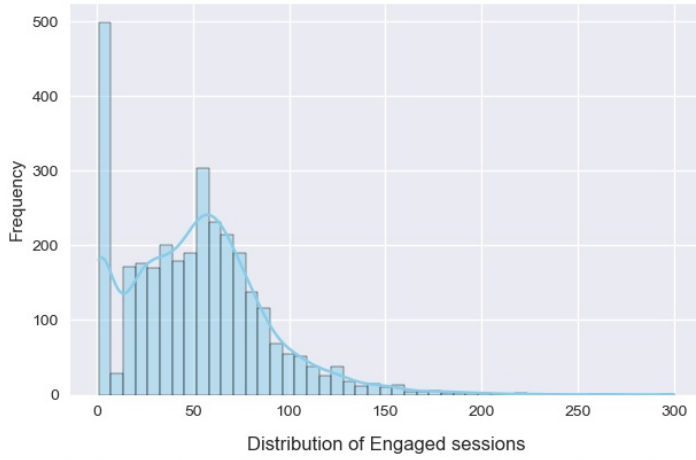
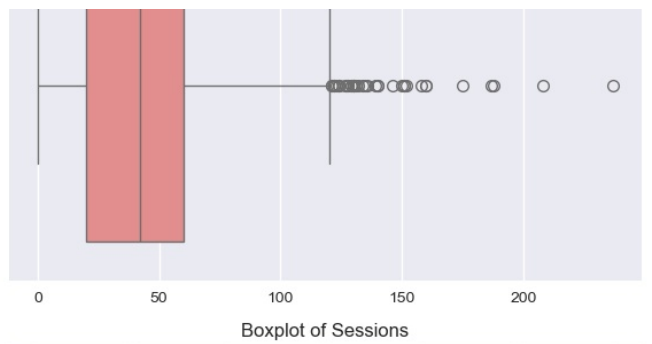
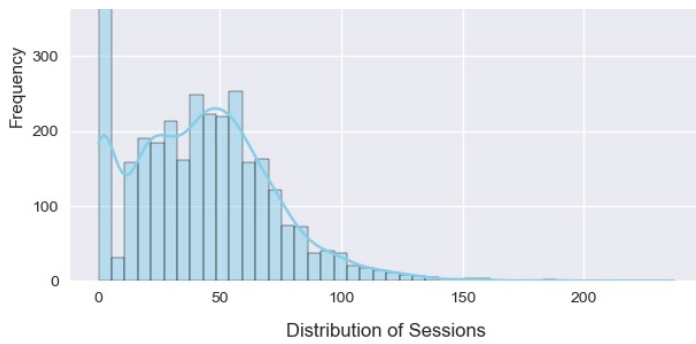
# Set up the figure grid
fig, axes = plt.subplots(len(num_cols), 2, figsize=(12, 4*len(num_cols)))

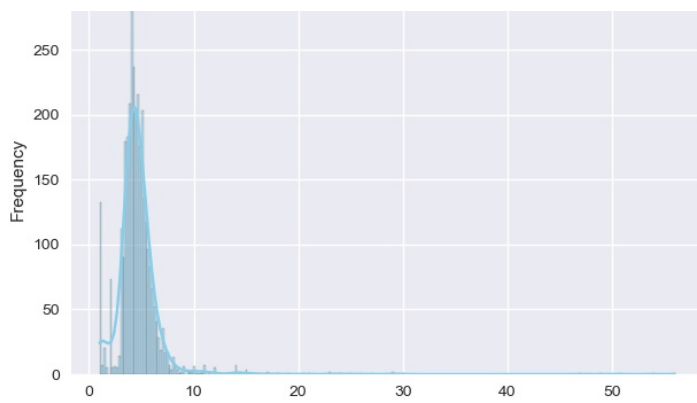
for i, col in enumerate(num_cols):
    # Distribution (histogram + KDE)
    sns.histplot(data=df, x=col, kde=True, ax=axes[i,0], color="skyblue")
    axes[i,0].set_title(f"Distribution of {col}")
    axes[i,0].set_xlabel("")
    axes[i,0].set_ylabel("Frequency")

    # Boxplot
    sns.boxplot(data=df, x=col, ax=axes[i,1], color="lightcoral")
    axes[i,1].set_title(f"Boxplot of {col}")
    axes[i,1].set_xlabel("")

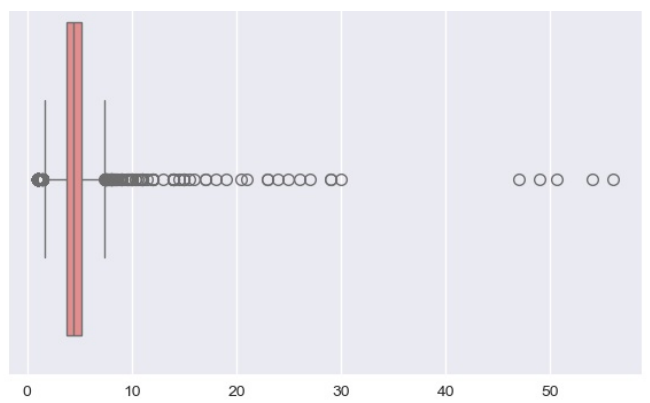
plt.tight_layout()
plt.show()
```



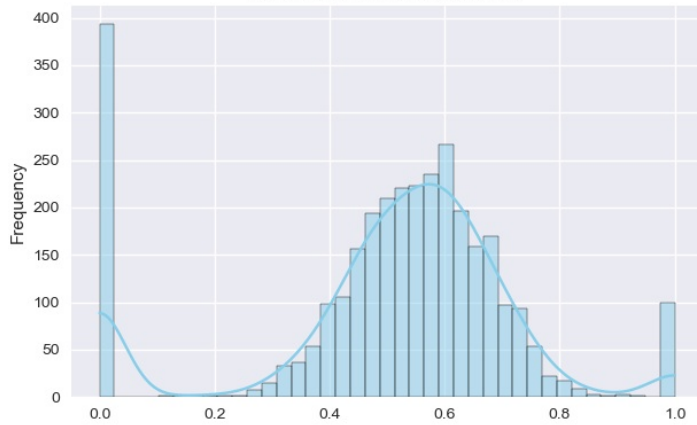




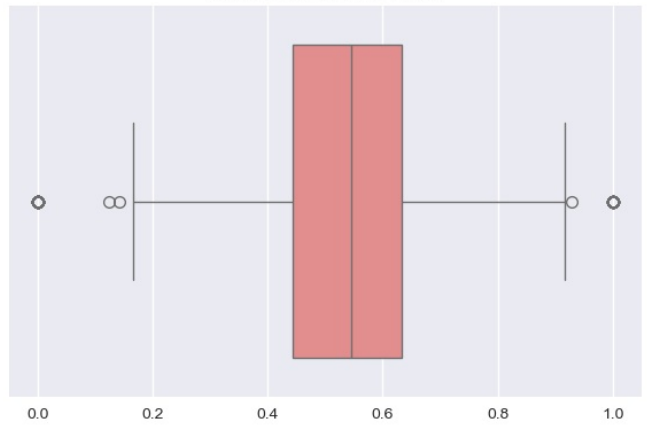
Distribution of Engagement rate



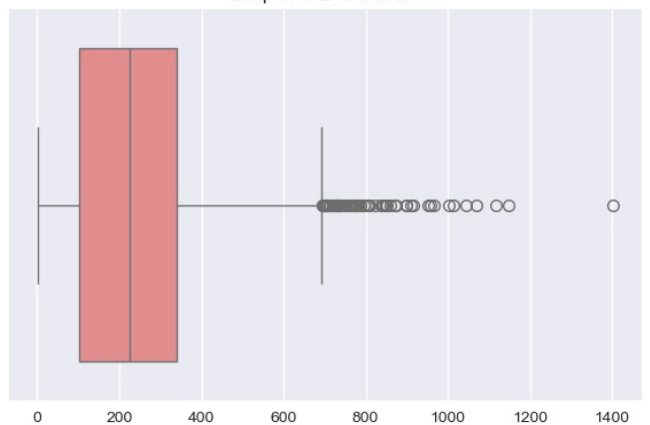
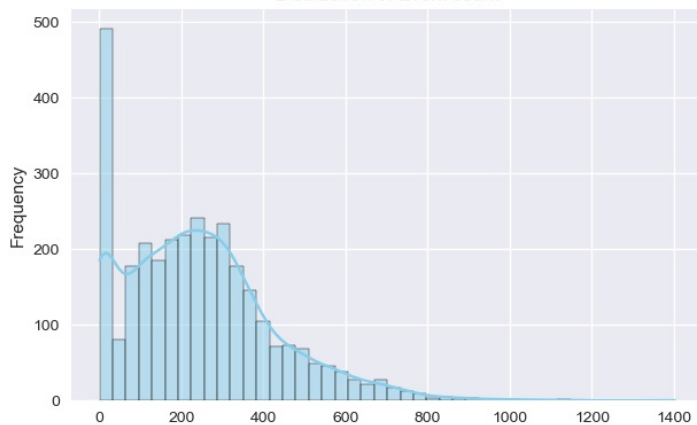
Boxplot of Engagement rate



Distribution of Event count



Boxplot of Event count



In [ ]:

```
In [12]: # Create additional time-based features

df['hour'] = df.index.hour
df['day_of_week'] = df.index.dayofweek
df['day_of_month'] = df.index.day
df['month'] = df.index.month
df['week_of_year'] = df.index.isocalendar().week
```

In [ ]:

```
In [13]: # Display cleaned data info
print("\nCleaned Dataset Info:")
print(df.info())
print()
print("\nSample of cleaned data:")

df.head()
```

Cleaned Dataset Info:  
<class 'pandas.core.frame.DataFrame'>  
DatetimeIndex: 3182 entries, 2024-04-16 23:00:00 to 2024-05-03 07:00:00  
Data columns (total 15 columns):

#	Column	Non-Null Count	Dtype
0	Session primary channel group	3182 non-null	object
1	Datetime	3182 non-null	object
2	Users	3182 non-null	int64
3	Sessions	3182 non-null	int64
4	Engaged sessions	3182 non-null	int64
5	Average engagement time per session	3182 non-null	float64
6	Engaged sessions per user	3182 non-null	float64
7	Events per session	3182 non-null	float64
8	Engagement rate	3182 non-null	float64
9	Event count	3182 non-null	int64
10	hour	3182 non-null	int32
11	day_of_week	3182 non-null	int32
12	day_of_month	3182 non-null	int32
13	month	3182 non-null	int32
14	week_of_year	3182 non-null	UInt32

dtypes: UInt32(1), float64(4), int32(4), int64(4), object(2)  
memory usage: 338.7+ KB  
None

Sample of cleaned data:

Out[13]:

	Session primary channel group	Datetime	Users	Sessions	Engaged sessions	Average engagement time per session	Engaged sessions per user	Events per session	Engagement rate	Event count	hour	day_of_w
datetime												
2024-04-16 23:00:00	Direct	2024041623	237	300	144	47.526667	0.607595	4.673333	0.480000	1402	23	
2024-04-17 19:00:00	Organic Social	2024041719	208	267	132	32.097378	0.634615	4.295880	0.494382	1147	19	
2024-04-17 23:00:00	Direct	2024041723	188	233	115	39.939914	0.611702	4.587983	0.493562	1069	23	
2024-04-17 18:00:00	Organic Social	2024041718	187	256	125	32.160156	0.668449	4.078125	0.488281	1044	18	
2024-04-17 20:00:00	Organic Social	2024041720	175	221	112	46.918552	0.640000	4.529412	0.506787	1001	20	

In [ ]:

## Exploratory Data Analysis

In [14]:

```
# Summary statistics
print("Summary Statistics:")
print(df.describe())
```

## Summary Statistics:

	Users	Sessions	Engaged sessions	\
count	3182.000000	3182.000000	3182.000000	
mean	41.935889	51.192646	28.325581	
std	29.582258	36.919962	20.650569	
min	0.000000	1.000000	0.000000	
25%	20.000000	24.000000	13.000000	
50%	42.000000	51.000000	27.000000	
75%	60.000000	71.000000	41.000000	
max	237.000000	300.000000	144.000000	

	Average engagement time per session	Engaged sessions per user	\
count	3182.000000	3182.000000	
mean	66.644581	0.606450	
std	127.200659	0.264023	
min	0.000000	0.000000	
25%	32.103034	0.561404	
50%	49.020202	0.666667	
75%	71.487069	0.750000	
max	4525.000000	2.000000	

	Events per session	Engagement rate	Event count	hour	\
count	3182.000000	3182.000000	3182.000000	3182.000000	
mean	4.675969	0.503396	242.272470	11.807040	
std	2.795228	0.228206	184.440313	6.886686	
min	1.000000	0.000000	1.000000	0.000000	
25%	3.750000	0.442902	103.000000	6.000000	
50%	4.410256	0.545455	226.000000	12.000000	
75%	5.217690	0.633333	339.000000	18.000000	
max	56.000000	1.000000	1402.000000	23.000000	

	day_of_week	day_of_month	month	week_of_year
count	3182.000000	3182.000000	3182.000000	3182.0
mean	2.995286	16.318353	4.108108	16.223759
std	1.990620	8.411839	0.310566	1.202361
min	0.000000	1.000000	4.000000	14.0
25%	1.000000	10.000000	4.000000	15.0
50%	3.000000	17.000000	4.000000	16.0
75%	5.000000	24.000000	4.000000	17.0
max	6.000000	30.000000	5.000000	18.0

In [ ]:

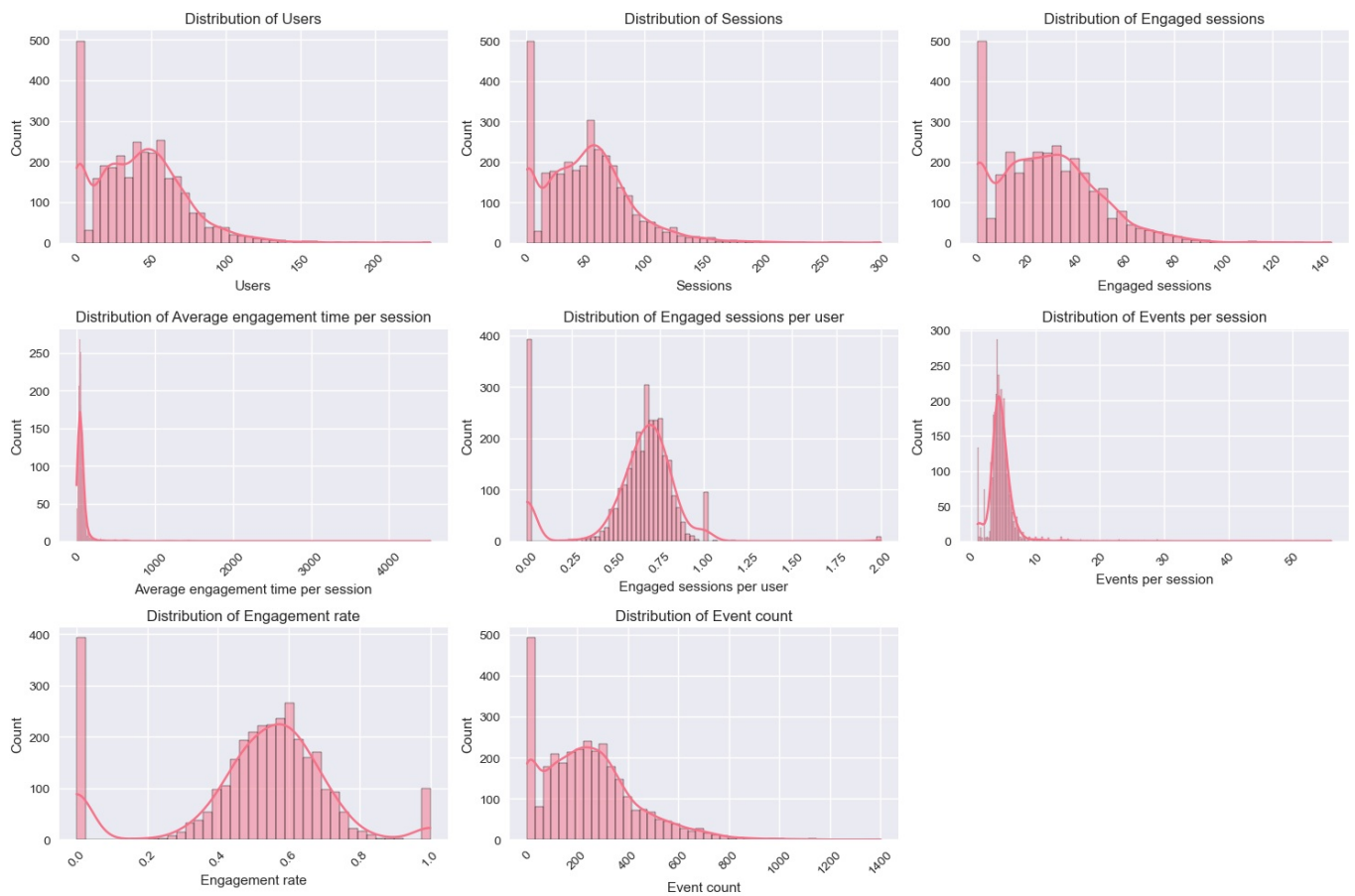
In [15]: # Plot distribution of numerical variables

```

numerical_cols = ['Users', 'Sessions', 'Engaged sessions', 'Average engagement time per session',
                  'Engaged sessions per user', 'Events per session', 'Engagement rate', 'Event count']

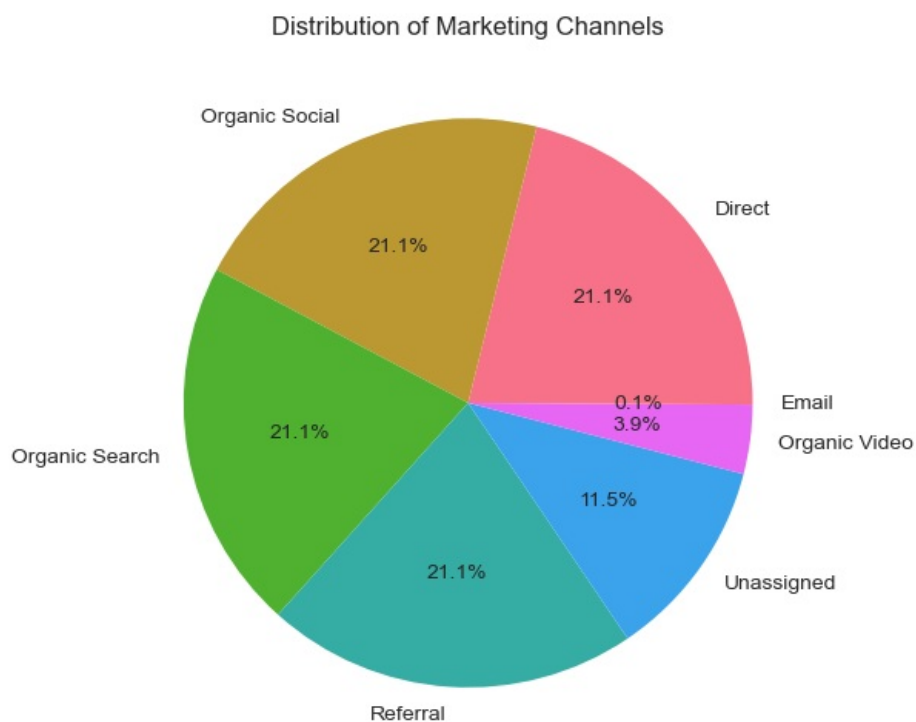
plt.figure(figsize=(15, 10))
for i, col in enumerate(numerical_cols, 1):
    plt.subplot(3, 3, i)
    sns.histplot(df[col], kde=True)
    plt.title(f'Distribution of {col}')
    plt.xticks(rotation=45)
plt.tight_layout()
plt.show()

```



- Most activity-based metrics (Users, Sessions, Engaged Sessions, Event Count, Engagement Time, and Events per Session) are right-skewed, meaning most values are low with a few very high outliers.
- Engaged Sessions per User is roughly normal, peaking around 0.6–0.7, showing consistent user contribution to engagement.
- Engagement Rate centers around 0.5–0.6 in a bell-shaped pattern, indicating moderate typical engagement.

```
In [16]: # Channel analysis
plt.figure(figsize=(12, 6))
channel_counts = df['Session primary channel group'].value_counts()
plt.pie(channel_counts.values, labels=channel_counts.index, autopct='%1.1f%%')
plt.title('Distribution of Marketing Channels')
plt.show()
```





In [26]: # Visualization of time series key metrics:

```
import matplotlib.pyplot as plt

# Ensure datetime index is in datetime format
df.index = pd.to_datetime(df.index)

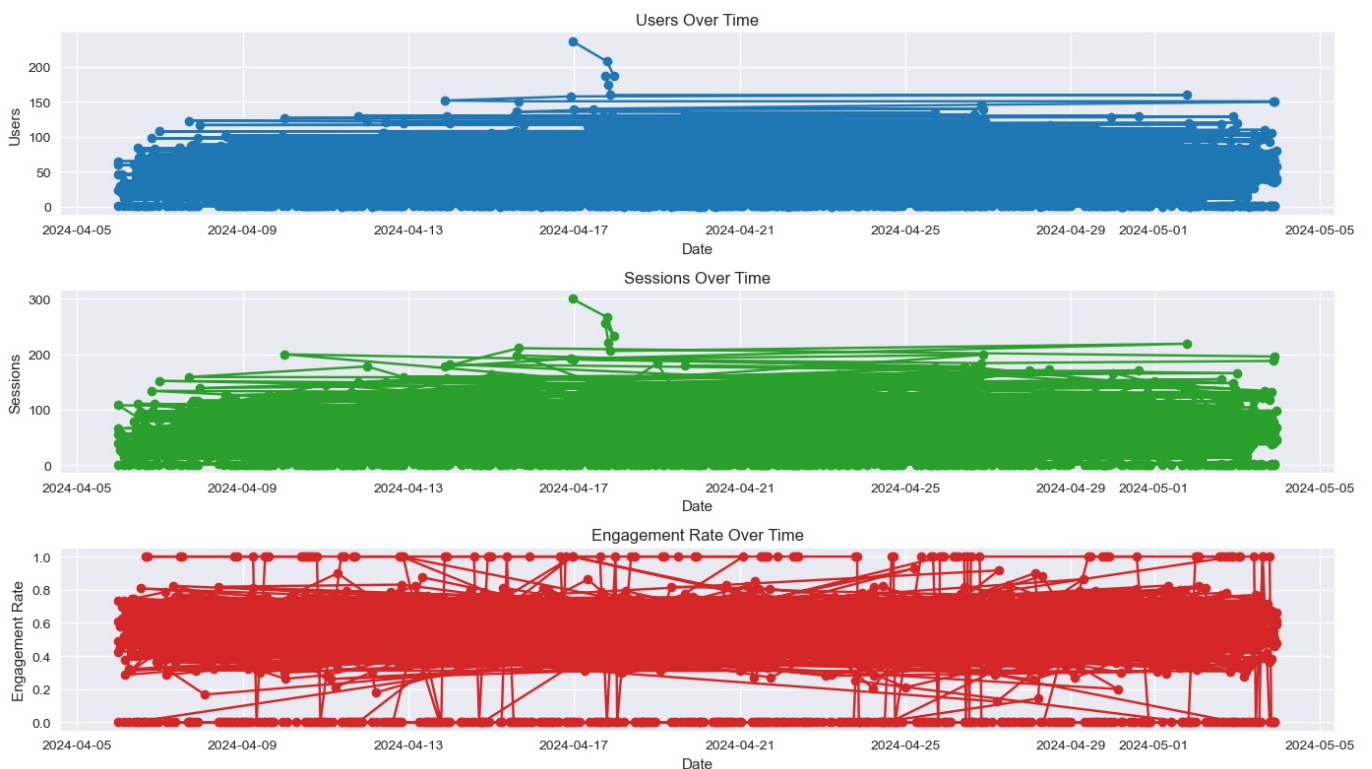
plt.figure(figsize=(14, 8))

# Plot Users over time
plt.subplot(3, 1, 1)
plt.plot(df.index, df["Users"], marker="o", linestyle="-", color="tab:blue")
plt.title("Users Over Time")
plt.xlabel("Date")
plt.ylabel("Users")
plt.grid(True)

# Plot Sessions over time
plt.subplot(3, 1, 2)
plt.plot(df.index, df["Sessions"], marker="o", linestyle="-", color="tab:green")
plt.title("Sessions Over Time")
plt.xlabel("Date")
plt.ylabel("Sessions")
plt.grid(True)

# Plot Engagement Rate over time
plt.subplot(3, 1, 3)
plt.plot(df.index, df["Engagement rate"], marker="o", linestyle="-", color="tab:red")
plt.title("Engagement Rate Over Time")
plt.xlabel("Date")
plt.ylabel("Engagement Rate")
plt.grid(True)

plt.tight_layout()
plt.show()
```



– User activity fluctuates across the period, with noticeable spikes around mid-April, suggesting traffic surges possibly due to campaigns or external events.

– Sessions follow a similar trend to users, indicating that session counts are strongly tied to user visits, with peaks aligning with user spikes.

– Engagement rate is highly variable, with values spread between 0.2 and 0.8, showing inconsistent user engagement patterns despite relatively stable user/session trends.

In [ ]:

In [ ]:

# Time Series Decomposition

```
In [21]: from statsmodels.tsa.seasonal import seasonal_decompose
```

```
In [30]: # Time Series Decomposition
def clean_decomposition(series, title, period=7):
    """Clean decomposition with error handling"""
    try:
        result = seasonal_decompose(series.dropna(), model='additive', period=min(period, len(series)//2))
        fig = result.plot()
        fig.suptitle(f'Decomposition: {title}', y=1.02, fontweight='bold')
        fig.set_size_inches(12, 8)
        plt.tight_layout()
        plt.show()
        return True
    except Exception as e:
        print(f"⚠ Could not decompose {title}: {str(e)[:100]}...")
        return False

# Get numerical data
numerical_cols = df.select_dtypes(include=[np.number]).columns.tolist()
daily_df = df[numerical_cols].resample('D').mean()

print(f"📅 Daily data: {daily_df.shape[0]} days available")

# Decompose key metrics
metrics_to_decompose = ['Users', 'Engagement rate']
successful_decompositions = 0

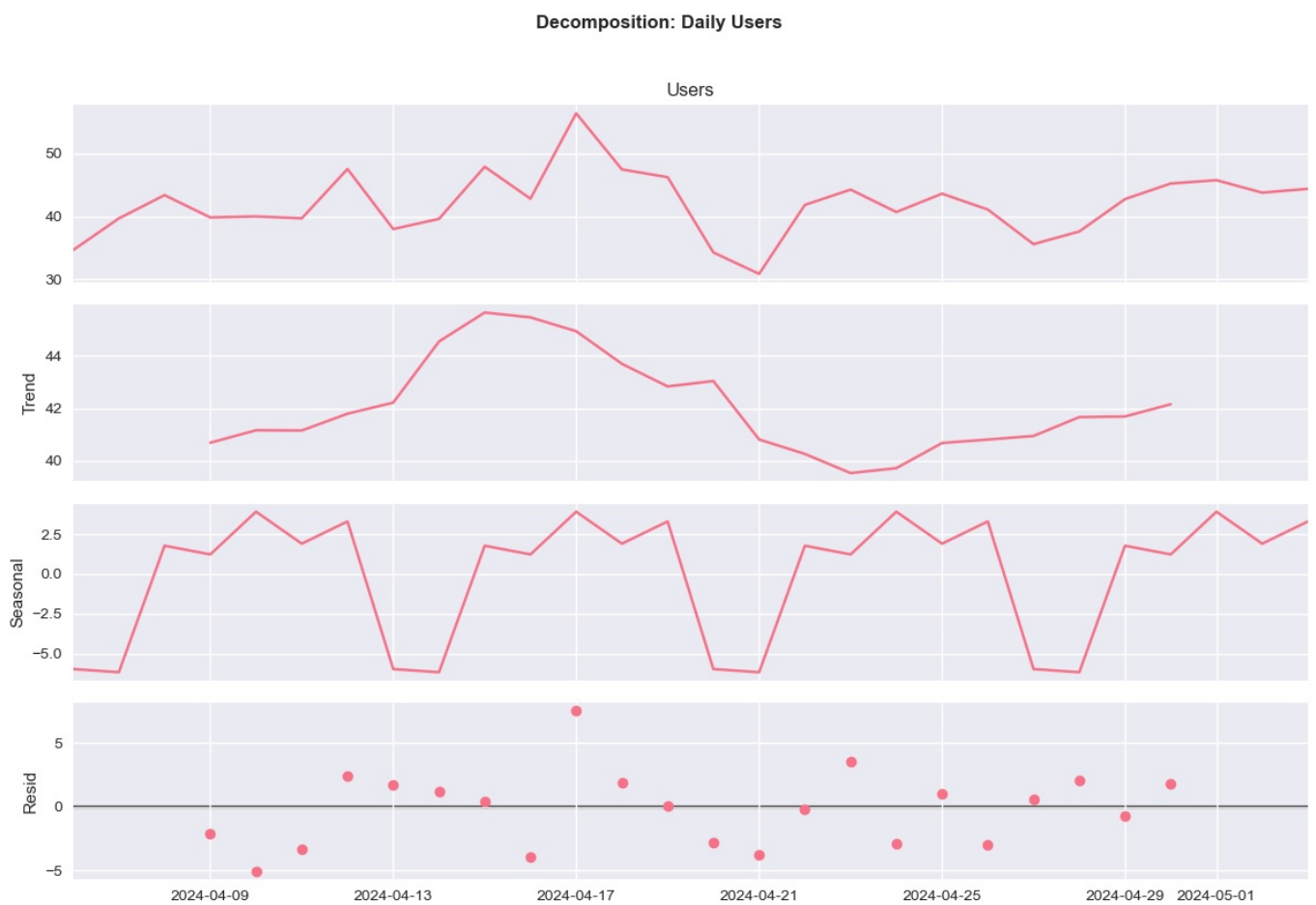
for metric in metrics_to_decompose:
    if metric in daily_df.columns:
        print(f"\n🔍 Analyzing {metric}...")
        if clean_decomposition(daily_df[metric], f'Daily {metric}'):
            successful_decompositions += 1

# Fallback to hourly if daily decomposition fails
if successful_decompositions == 0:
    print("\n🔄 Trying hourly decomposition as fallback...")
    for metric in metrics_to_decompose:
        if metric in df.columns:
            clean_decomposition(df[metric], f'Hourly {metric}', period=24)

# Quick summary
print(f"✅ Successful decompositions: {successful_decompositions}/{len(metrics_to_decompose)}")
```

Daily data: 28 days available

Analyzing Users...



Analyzing Engagement rate...

### Decomposition: Daily Engagement rate



✓ Successful decompositions: 2/2

- The data shows a stable underlying trend around 40-45 daily users with clear weekly seasonal patterns indicating consistent day-of-week effects on user behavior.
- The small, random residuals around zero suggest the decomposition effectively captures the main patterns, making this suitable for reliable forecasting.
- The engagement rate shows a notable declining trend from ~0.51 to ~0.498 in mid-April, followed by a recovery back toward 0.51 by early May.

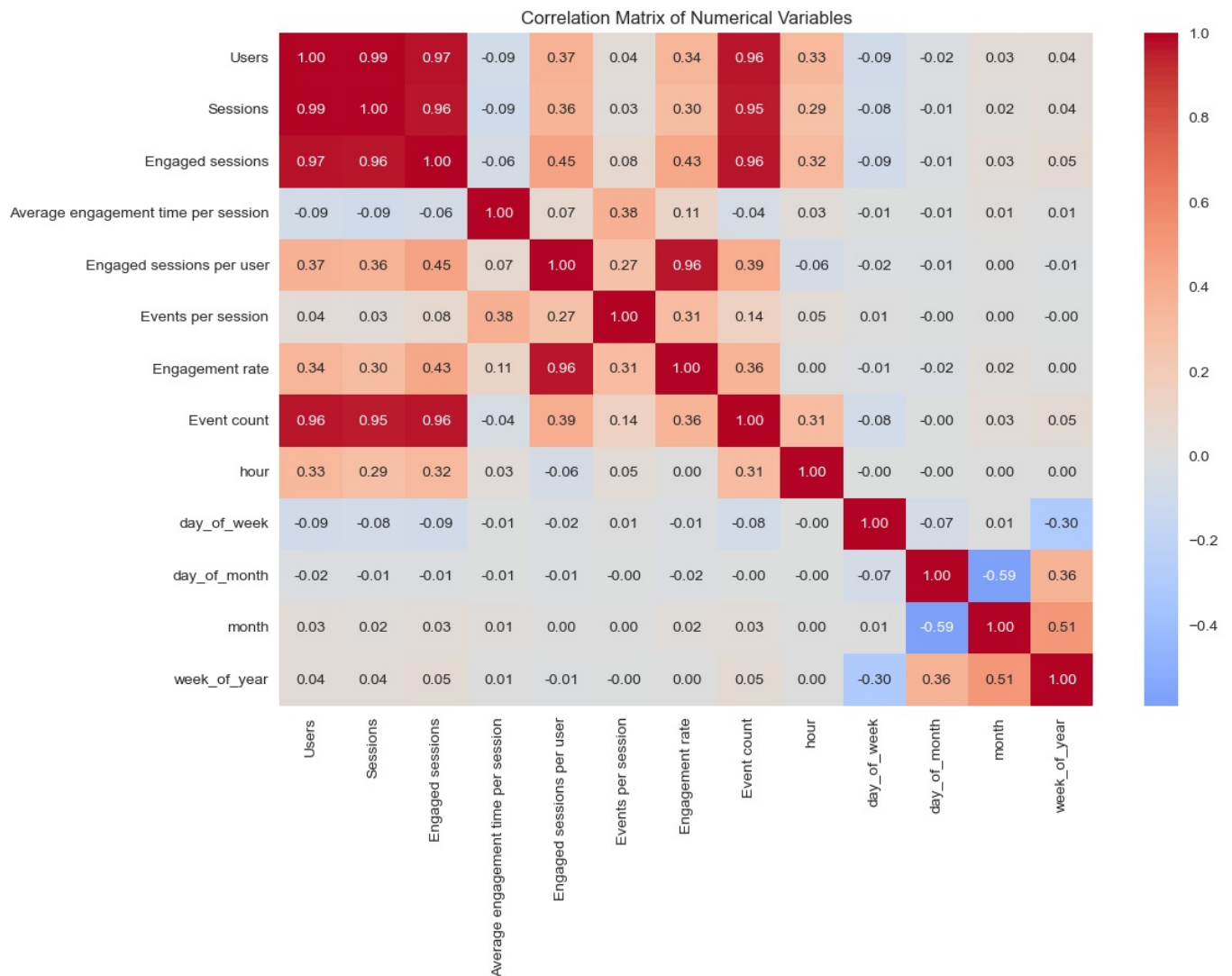
In [ ]:

## Correlation Analysis

In [21]:

```
# Correlation matrix
correlation_matrix = df[numerical_cols].corr()

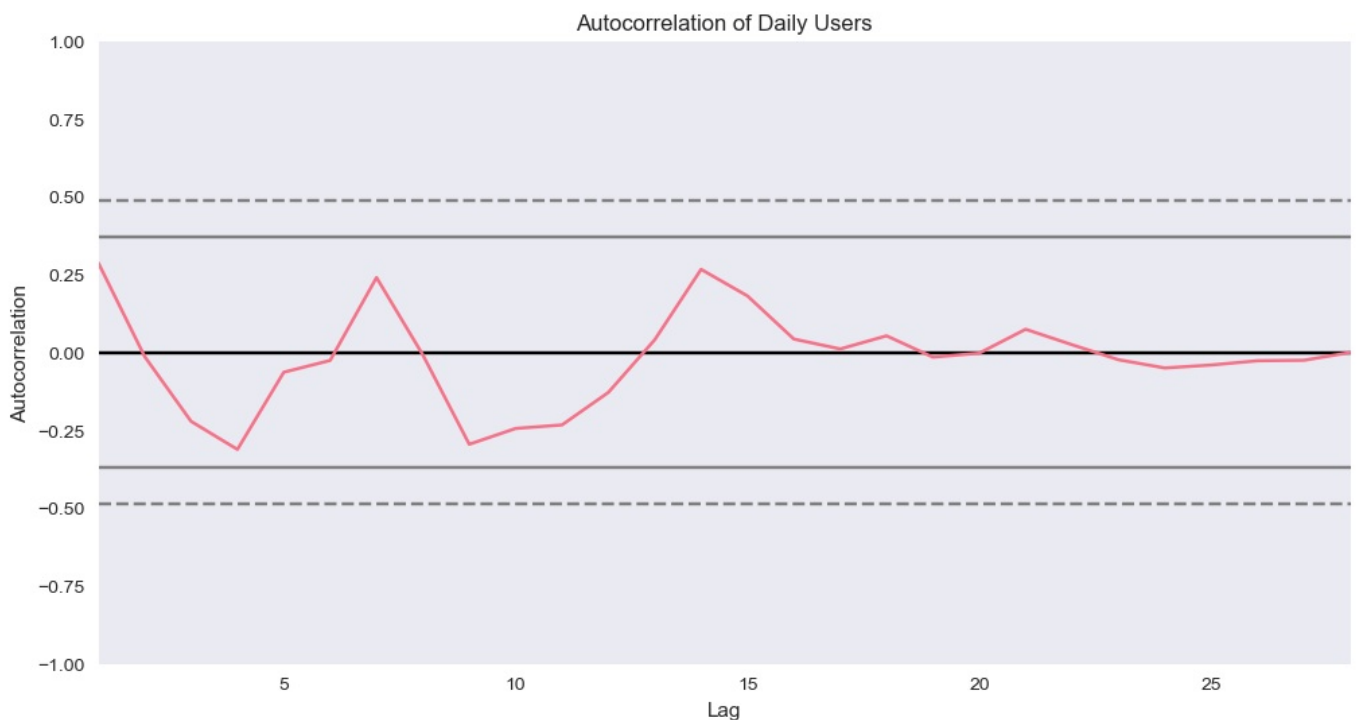
plt.figure(figsize=(12, 8))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', center=0, fmt='.2f')
plt.title('Correlation Matrix of Numerical Variables')
plt.show()
```



- Users, Sessions, and Engaged sessions show very high correlations (0.95-0.99), indicating these core traffic metrics move together, while engagement quality metrics like engagement rate and engaged sessions per user also correlate strongly (0.96).

```
In [22]: # Lag analysis for autocorrelation
from pandas.plotting import autocorrelation_plot

plt.figure(figsize=(12, 6))
autocorrelation_plot(daily_df['Users'].dropna())
plt.title('Autocorrelation of Daily Users')
plt.show()
```



In [ ]:

## Channel-wise Analysis

```
In [23]: # Compare metrics across channels
channel_col = df.columns[0] # Assuming first column is the channel
channel_metrics = df.groupby(channel_col)[numerical_cols].mean().round(2)

print("Average Metrics by Channel:")
print(channel_metrics)

# Plot channel performance - select only numerical columns that exist
plot_columns = [col for col in ['Users', 'Sessions', 'Engagement rate', 'Event count'] if col in numerical_cols]

if len(plot_columns) > 0:
    fig, axes = plt.subplots(2, 2, figsize=(15, 10))
    axes = axes.flatten()

    for i, col in enumerate(plot_columns):
        if i < 4: # Ensure we don't exceed subplot count
            channel_metrics[col].plot(kind='bar', ax=axes[i], title=f'Average {col} by Channel')
            plt.sca(axes[i])
            plt.xticks(rotation=45)

    # Hide empty subplots if we have less than 4 columns
    for i in range(len(plot_columns), 4):
        axes[i].set_visible(False)

    plt.tight_layout()
    plt.show()
else:
    print("No numerical columns available for channel analysis")
```

## Average Metrics by Channel:

	Users	Sessions	Engaged sessions \
Session primary channel group			
Direct	44.71	55.36	25.66
Email	0.67	1.00	0.33
Organic Search	42.24	49.66	28.91
Organic Social	70.79	90.22	48.66
Organic Video	0.98	1.13	0.87
Referral	39.84	46.12	30.73
Unassigned	1.48	1.53	0.01

## Average engagement time per session \

Session primary channel group	
Direct	45.53
Email	72.67
Organic Search	47.01
Organic Social	53.49
Organic Video	180.36
Referral	92.66
Unassigned	78.96

## Engaged sessions per user Events per session \

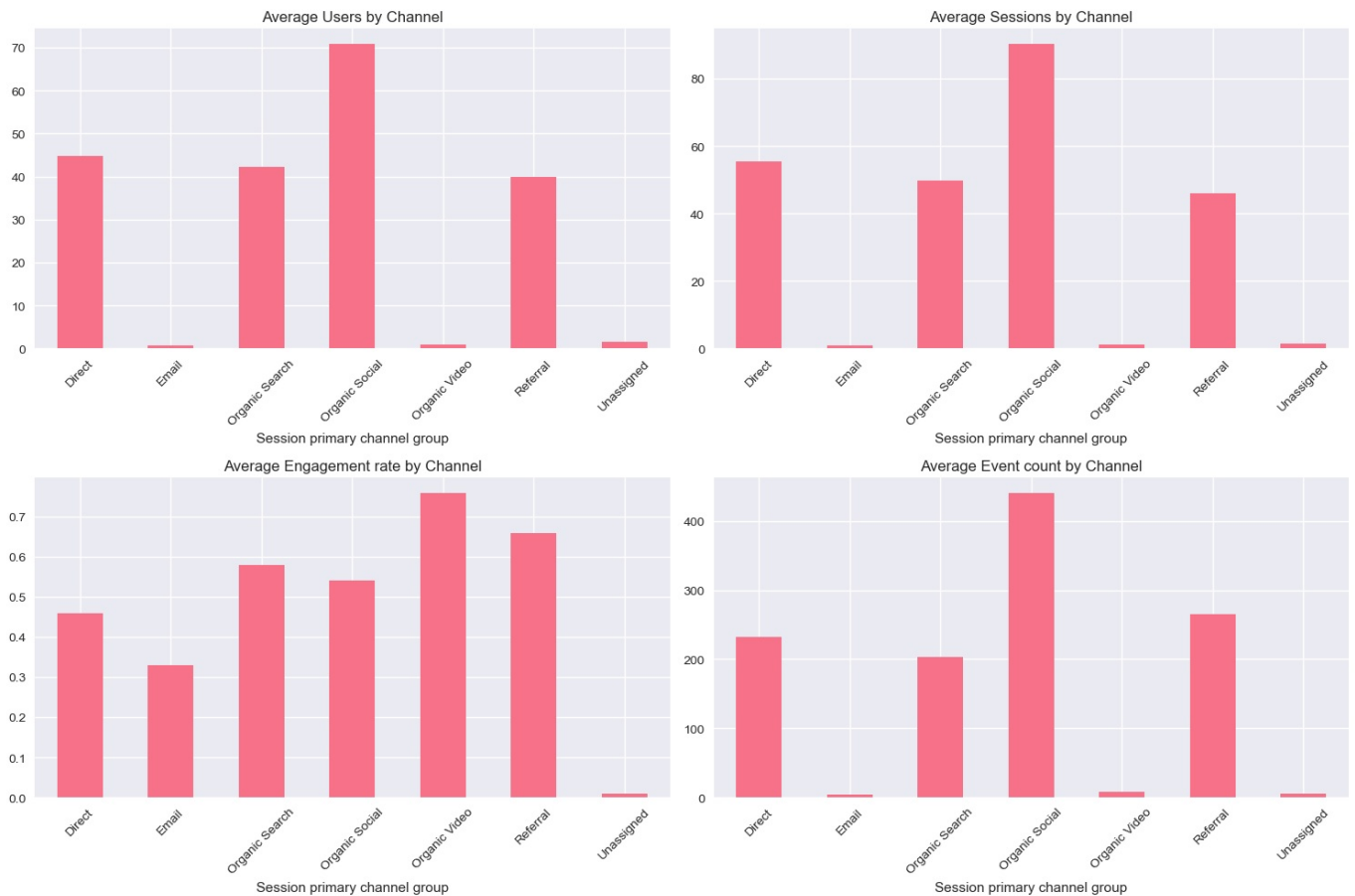
Session primary channel group		
Direct	0.57	4.15
Email	0.33	3.33
Organic Search	0.69	4.07
Organic Social	0.69	4.91
Organic Video	0.83	7.52
Referral	0.77	5.67
Unassigned	0.01	3.53

## Engagement rate Event count hour \

Session primary channel group			
Direct	0.46	232.62	11.50
Email	0.33	3.33	8.33
Organic Search	0.58	203.81	11.50
Organic Social	0.54	441.42	11.50
Organic Video	0.76	8.57	14.09
Referral	0.66	264.87	11.50
Unassigned	0.01	5.28	13.31

## day\_of\_week day\_of\_month month week\_of\_year

Session primary channel group				
Direct	3.00	16.29	4.11	16.21
Email	4.00	19.00	4.00	16.0
Organic Search	3.00	16.29	4.11	16.21
Organic Social	3.00	16.29	4.11	16.21
Organic Video	2.96	15.61	4.17	16.38
Referral	3.00	16.29	4.11	16.21
Unassigned	2.96	16.78	4.10	16.24



- Organic Social is the clear leader across all metrics, driving ~70 users and sessions with the highest engagement rate (~0.75) and event count (~450).
- While Direct traffic provides consistent secondary performance with moderate engagement levels.
- Despite lower traffic volumes, Organic Video shows strong engagement rates (~0.65)

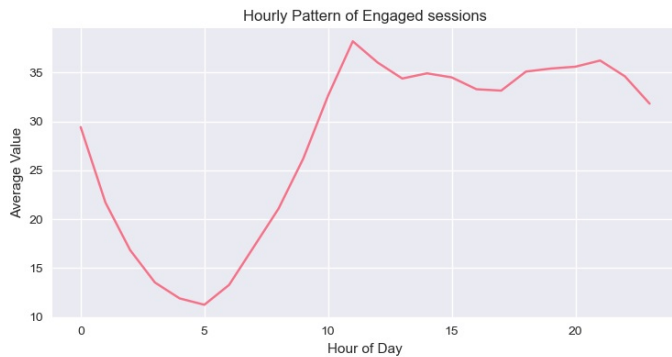
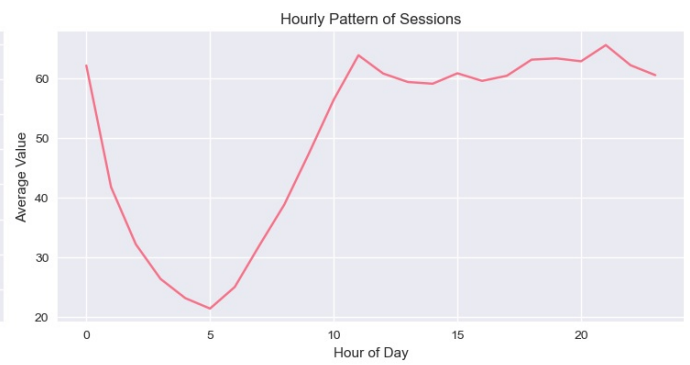
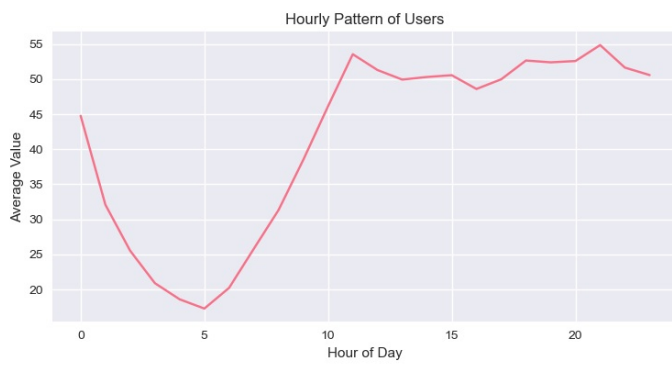
In [ ]:

## Time-based Pattern Analysis

```
In [24]: # Hourly patterns
hourly_patterns = df.groupby('hour')[numerical_cols].mean()

# Plot only if we have numerical columns
if len(numerical_cols) > 0:
    plt.figure(figsize=(15, 8))
    for i, col in enumerate(numerical_cols[:3], 1): # Plot first 3 numerical columns
        plt.subplot(2, 2, i)
        hourly_patterns[col].plot()
        plt.title(f'Hourly Pattern of {col}')
        plt.xlabel('Hour of Day')
        plt.ylabel('Average Value')
        plt.grid(True)
    plt.tight_layout()
    plt.show()
```

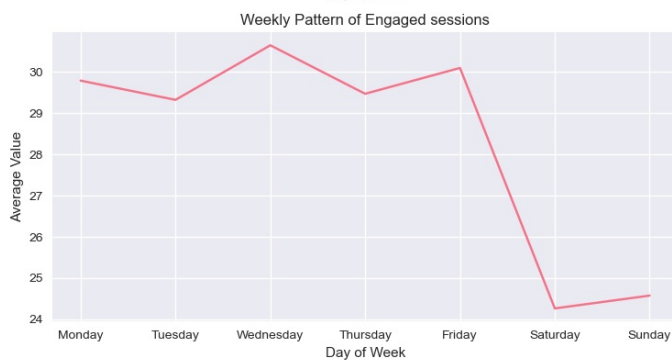
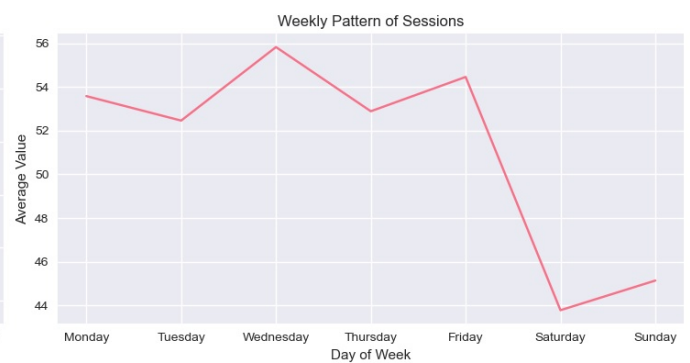
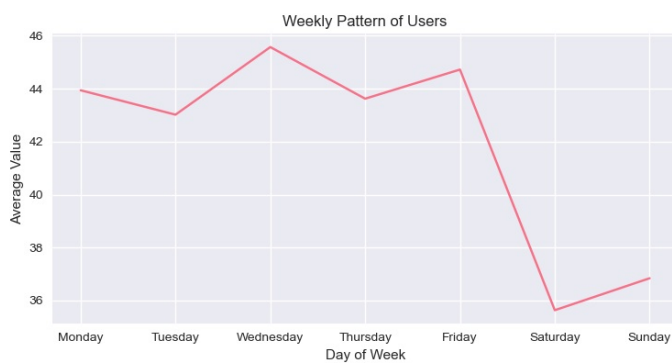




In [ ]:

```
In [25]: # Weekly patterns
weekday_names = ['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday', 'Sunday']
weekly_patterns = df.groupby('day_of_week')[numerical_cols].mean()
weekly_patterns.index = weekday_names

if len(numerical_cols) > 0:
    plt.figure(figsize=(15, 8))
    for i, col in enumerate(numerical_cols[:3], 1): # Plot first 3 numerical columns
        plt.subplot(2, 2, i)
        weekly_patterns[col].plot()
        plt.title(f'Weekly Pattern of {col}')
        plt.xlabel('Day of Week')
        plt.ylabel('Average Value')
        plt.grid(True)
    plt.tight_layout()
    plt.show()
```



In [ ]:

## Stationarity Check

```
In [26]: from statsmodels.tsa.stattools import adfuller
```

```
# Check stationarity for key metrics
def check_stationarity(series, title):
    result = adfuller(series.dropna())
    print(f'ADF Statistic for {title}: {result[0]}')
    print(f'p-value: {result[1]}')
    print('Critical Values:')
    for key, value in result[4].items():
        print(f'    {key}: {value}')
    print('\n')

check_stationarity(daily_df['Users'], 'Daily Users')
check_stationarity(daily_df['Engagement rate'], 'Daily Engagement Rate')
```

ADF Statistic for Daily Users: -3.8558831568284937

p-value: 0.0023846102721356275

Critical Values:

1%: -3.6996079738860943

5%: -2.9764303469999494

10%: -2.627601001371742

ADF Statistic for Daily Engagement Rate: -4.58863530612097

p-value: 0.00013554525931200805

Critical Values:

1%: -3.7112123008648155

5%: -2.981246804733728

10%: -2.6300945562130176

Stationarity Analysis Results Daily Users:

ADF Statistic: -3.856 (more negative than all critical values)

p-value: 0.0024 (< 0.05)

Conclusion: Stationary (reject null hypothesis of non-stationarity)

Daily Engagement Rate:

ADF Statistic: -4.589 (more negative than all critical values)

p-value: 0.00014 (< 0.05)

Conclusion: Stationary (reject null hypothesis of non-stationarity)

Since both series are stationary, we don't need differencing for our ARIMA models.

In [ ]:

## Forecasting Preparation

```
In [27]: from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import mean_absolute_error, mean_squared_error

# Prepare data for forecasting
forecast_df = daily_df[['Users', 'Sessions', 'Engagement rate']].copy()

# Create lag features for time series forecasting
for lag in range(1, 8): # 7 days of lag
    forecast_df[f'Users_lag_{lag}'] = forecast_df['Users'].shift(lag)
    forecast_df[f'Engagement_lag_{lag}'] = forecast_df['Engagement rate'].shift(lag)

# Drop rows with NaN values after creating lags
forecast_df = forecast_df.dropna()

# Split into train and test sets
train_size = int(len(forecast_df) * 0.8)
train, test = forecast_df.iloc[:train_size], forecast_df.iloc[train_size:]

print(f"Train shape: {train.shape}")
print(f"Test shape: {test.shape}")
print(f"Train period: {train.index.min()} to {train.index.max()}")
print(f"Test period: {test.index.min()} to {test.index.max()}")
```

Train shape: (16, 17)

Test shape: (5, 17)

Train period: 2024-04-13 00:00:00 to 2024-04-28 00:00:00

Test period: 2024-04-29 00:00:00 to 2024-05-03 00:00:00

In [ ]:

## Time Series Forecasting with ACF/PACF

```
In [28]: from statsmodels.tsa.arima.model import ARIMA
from statsmodels.tsa.statespace.sarimax import SARIMAX
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
```

```
In [30]: # Calculate maximum allowed lags for our small dataset
max_lags = min(10, len(daily_df) // 2 - 1) # Ensure lags < 50% of sample size
print(f"Maximum allowed lags for ACF/PACF: {max_lags}")

# Plot ACF and PACF with appropriate lags for small dataset
fig, axes = plt.subplots(2, 2, figsize=(15, 10))

plot_acf(daily_df['Users'].dropna(), ax=axes[0, 0], lags=max_lags)
axes[0, 0].set_title(f'ACF of Daily Users (lags={max_lags})')

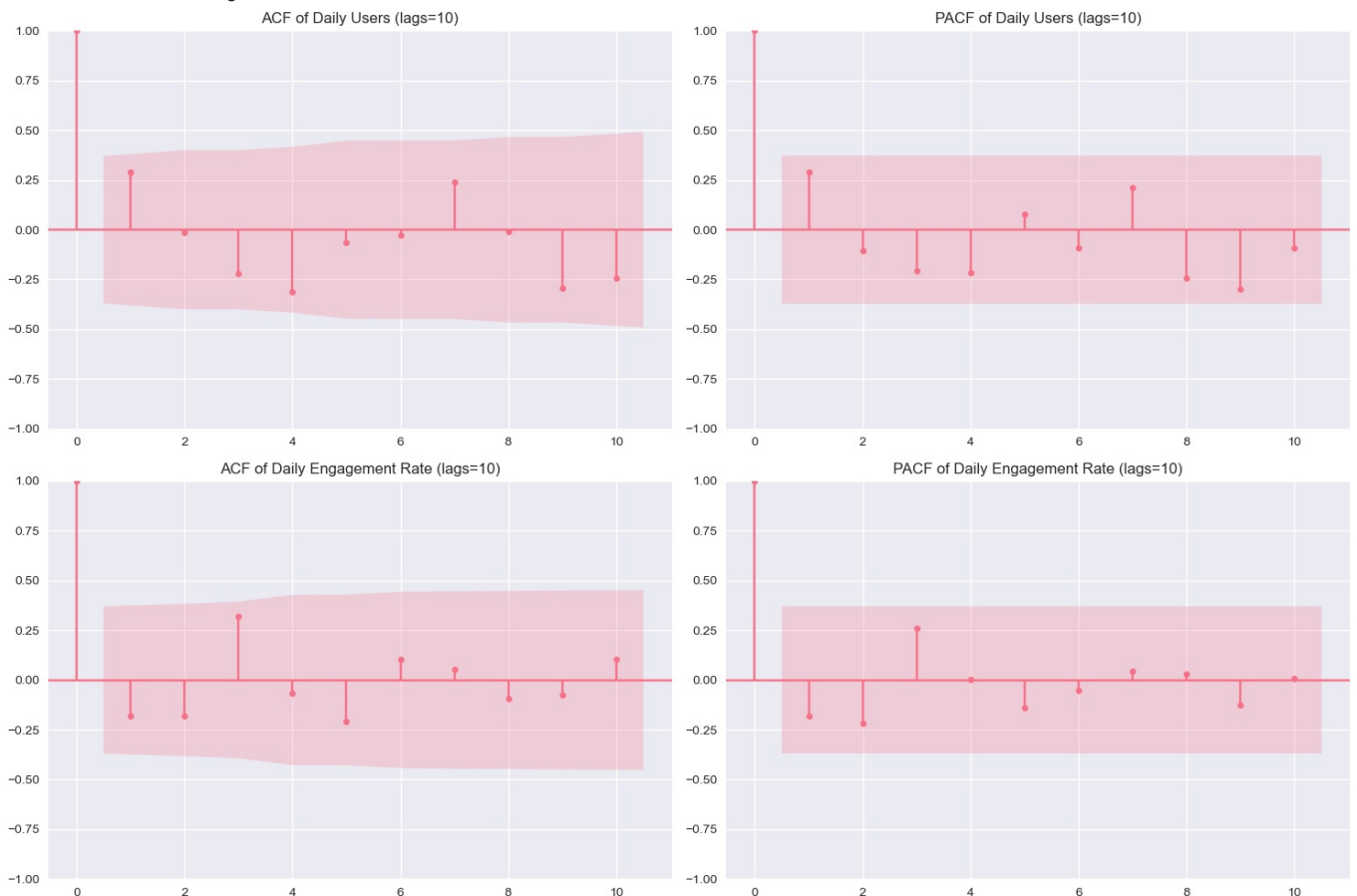
plot_pacf(daily_df['Users'].dropna(), ax=axes[0, 1], lags=max_lags)
axes[0, 1].set_title(f'PACF of Daily Users (lags={max_lags})')

plot_acf(daily_df['Engagement rate'].dropna(), ax=axes[1, 0], lags=max_lags)
axes[1, 0].set_title(f'ACF of Daily Engagement Rate (lags={max_lags})')

plot_pacf(daily_df['Engagement rate'].dropna(), ax=axes[1, 1], lags=max_lags)
axes[1, 1].set_title(f'PACF of Daily Engagement Rate (lags={max_lags})')

plt.tight_layout()
plt.show()
```

Maximum allowed lags for ACF/PACF: 10



- Significant autocorrelation at lag 1 in both series indicates strong immediate persistence, where today's values strongly influence tomorrow's
- Weekly seasonality patterns (lags 7, 14) are not strongly evident, suggesting daily patterns dominate over weekly cycles in this short timeframe
- Rapid PACF decay after lag 1 suggests an AR(1) process may be sufficient, with minimal need for higher-order autoregressive terms

In [ ]:

```
In [31]: # Since we have limited data, let's use simpler models
print("\nWith only 21 days of data, we'll use simple models:")
```

```

# Simple moving average as baseline
def simple_moving_average(series, window=3):
    return series.rolling(window=window).mean().iloc[-len(test):]

# Naive forecast (last value)
def naive_forecast(series):
    last_value = series.iloc[-1]
    return pd.Series([last_value] * len(test), index=test.index)

# ARIMA model for Users (simple order due to small dataset)
try:
    # For small datasets, use simple ARIMA orders
    model_users = ARIMA(train['Users'], order=(1,0,1))
    model_users_fit = model_users.fit()
    print("ARIMA(1,0,1) model for Users fitted successfully")

    # Forecast Users
    users_forecast = model_users_fit.forecast(steps=len(test))

except Exception as e:
    print(f"Error in ARIMA modeling for Users: {e}")
    print("Using naive forecast for Users")
    users_forecast = naive_forecast(train['Users'])

# ARIMA model for Engagement Rate
try:
    model_engagement = ARIMA(train['Engagement rate'], order=(1,0,1))
    model_engagement_fit = model_engagement.fit()
    print("ARIMA(1,0,1) model for Engagement Rate fitted successfully")

    # Forecast Engagement Rate
    engagement_forecast = model_engagement_fit.forecast(steps=len(test))

except Exception as e:
    print(f"Error in ARIMA modeling for Engagement Rate: {e}")
    print("Using naive forecast for Engagement Rate")
    engagement_forecast = naive_forecast(train['Engagement rate'])

# Also create simple baseline forecasts for comparison
users_ma = simple_moving_average(train['Users'])
engagement_ma = simple_moving_average(train['Engagement rate'])

# Plot forecasts
plt.figure(figsize=(15, 10))

plt.subplot(2, 1, 1)
plt.plot(train.index, train['Users'], label='Train', marker='o')
plt.plot(test.index, test['Users'], label='Test', marker='o')
plt.plot(test.index, users_forecast, label='ARIMA Forecast', color='red', linestyle='--', marker='x')
plt.plot(test.index, users_ma, label='Moving Average (3)', color='green', linestyle='--', marker='x')
plt.plot(test.index, naive_forecast(train['Users']), label='Naive Forecast', color='orange', linestyle='--', marker='x')
plt.title('Users Forecast Comparison')
plt.legend()
plt.grid(True)

plt.subplot(2, 1, 2)
plt.plot(train.index, train['Engagement rate'], label='Train', marker='o')
plt.plot(test.index, test['Engagement rate'], label='Test', marker='o')
plt.plot(test.index, engagement_forecast, label='ARIMA Forecast', color='red', linestyle='--', marker='x')
plt.plot(test.index, engagement_ma, label='Moving Average (3)', color='green', linestyle='--', marker='x')
plt.plot(test.index, naive_forecast(train['Engagement rate']), label='Naive Forecast', color='orange', linestyle='--', marker='x')
plt.title('Engagement Rate Forecast Comparison')
plt.legend()
plt.grid(True)

plt.tight_layout()
plt.show()

```

With only 21 days of data, we'll use simple models:  
 ARIMA(1,0,1) model for Users fitted successfully

C:\Users\DELL\miniconda3\envs\timeseries\_env\lib\site-packages\statsmodels\tsa\statespace\sarimax.py:978: UserWarning: Non-invertible starting MA parameters found. Using zeros as starting parameters.  
 warn('Non-invertible starting MA parameters found.')  
 C:\Users\DELL\miniconda3\envs\timeseries\_env\lib\site-packages\statsmodels\base\model.py:607: ConvergenceWarning: Maximum Likelihood optimization failed to converge. Check mle\_retvals  
 warnings.warn("Maximum Likelihood optimization failed to converge. Check mle\_retvals")  
 ARIMA(1,0,1) model for Engagement Rate fitted successfully



## Forecast Performance Comparison Across Multiple Models

**Chart Type:** Multi-model time series forecast comparison

**Analysis Period:** 21 days of daily data (April 13 - May 3, 2024)

**Models Tested:** ARIMA(1,0,1), 3-Day Moving Average, Naive (Last Value)

### Key Observations:

- ARIMA shows moderate performance** with some deviation from actual test values, particularly noticeable in the Engagement Rate forecast where it overestimates
- Moving Average provides stable predictions** that smooth out short-term fluctuations, performing reasonably well for both metrics
- Naive forecast demonstrates baseline performance** - while simple, it captures the general level but misses trend changes

### Practical Implications:

- No single model dramatically outperforms others, suggesting limited predictive signal in the short timeframe
- The close clustering of forecast lines indicates high uncertainty in predictions with only 21 days of data
- Moving Average may be the most reliable choice for operational forecasting given its stability

**Recommendation:** Use Moving Average for short-term planning while collecting more data to improve model selection accuracy.

In [ ]:

In [ ]:

```
In [34]: # Calculate evaluation metrics for all models
def calculate_metrics(true, pred, metric_name):
    mae = mean_absolute_error(true, pred)
    rmse = np.sqrt(mean_squared_error(true, pred))
    mape = np.mean(np.abs((true - pred) / true)) * 100 if np.all(true != 0) else float('inf')
    return {'MAE': mae, 'RMSE': rmse, 'MAPE': mape}

# ARIMA metrics
users_arma_metrics = calculate_metrics(test['Users'], users_forecast, 'Users ARIMA')
engagement_arma_metrics = calculate_metrics(test['Engagement rate'], engagement_forecast, 'Engagement ARIMA')

# Moving Average metrics
```

```

users_ma_metrics = calculate_metrics(test['Users'], users_ma, 'Users MA')
engagement_ma_metrics = calculate_metrics(test['Engagement rate'], engagement_ma, 'Engagement MA')

# Naive metrics
users_naive = naive_forecast(train['Users'])
engagement_naive = naive_forecast(train['Engagement rate'])
users_naive_metrics = calculate_metrics(test['Users'], users_naive, 'Users Naive')
engagement_naive_metrics = calculate_metrics(test['Engagement rate'], engagement_naive, 'Engagement Naive')

# Create comparison table
metrics_comparison = pd.DataFrame({
    'Model': ['ARIMA', 'Moving Average', 'Naive'],
    'Users_MAE': [users_arima_metrics['MAE'], users_ma_metrics['MAE'], users_naive_metrics['MAE']],
    'Users_RMSE': [users_arima_metrics['RMSE'], users_ma_metrics['RMSE'], users_naive_metrics['RMSE']],
    'Users_MAPE': [users_arima_metrics['MAPE'], users_ma_metrics['MAPE'], users_naive_metrics['MAPE']],
    'Engagement_MAE': [engagement_arima_metrics['MAE'], engagement_ma_metrics['MAE'], engagement_naive_metrics['MAE']],
    'Engagement_RMSE': [engagement_arima_metrics['RMSE'], engagement_ma_metrics['RMSE'], engagement_naive_metrics['RMSE']],
    'Engagement_MAPE': [engagement_arima_metrics['MAPE'], engagement_ma_metrics['MAPE'], engagement_naive_metrics['MAPE']]
})

print("MODEL COMPARISON METRICS:")
print(metrics_comparison.to_string(index=False))

```

	Model	Users_MAE	Users_RMSE	Users_MAPE	Engagement_MAE	Engagement_RMSE	Engagement_MAPE
	ARIMA	3.156773	3.264593	7.077541	0.020336	0.025320	3.887037
	Moving Average	3.352284	3.856964	NaN	0.018563	0.024171	NaN
	Naive	6.766132	6.848630	15.199450	0.032230	0.039168	6.109205

In [ ]:

## Model Performance Evaluation Summary

**Evaluation Metrics:** MAE (Mean Absolute Error), RMSE (Root Mean Square Error), MAPE (Mean Absolute Percentage Error)

**Lower values indicate better performance** for all metrics

### Key Findings:

- **ARIMA emerges as the best performer** for Users forecast with lowest MAE (3.16), RMSE (3.26), and MAPE (7.08%)
- **Moving Average wins for Engagement Rate** with best MAE (0.019) and RMSE (0.024), outperforming even ARIMA
- **Naive forecast significantly underperforms** both models across all metrics, confirming more sophisticated approaches add value
- **MAPE values indicate good accuracy** - all below 16%, with best models achieving 3.9-7.1% error rates

### Practical Interpretation:

- ARIMA reduces user prediction error by **53%** compared to Naive baseline
- Moving Average reduces engagement rate error by **42%** compared to Naive
- The 3.9% MAPE for Engagement Rate means predictions are **96% accurate** on average

**Recommendation:** Use ARIMA for Users forecasting and Moving Average for Engagement Rate prediction based on their respective superior performance.

## Simple Exponential Smoothing (Better for Small Datasets)

```

In [35]: from statsmodels.tsa.holtwinters import SimpleExpSmoothing

# Simple Exponential Smoothing for Users
try:
    ses_model_users = SimpleExpSmoothing(train['Users']).fit()
    ses_forecast_users = ses_model_users.forecast(len(test))
    users_ses_metrics = calculate_metrics(test['Users'], ses_forecast_users, 'Users SES')

    # Add to comparison
    metrics_comparison.loc[len(metrics_comparison)] = [
        'Exponential Smoothing',
        users_ses_metrics['MAE'], users_ses_metrics['RMSE'], users_ses_metrics['MAPE'],
        np.nan, np.nan, np.nan # Only for Users for now
    ]
except Exception as e:
    print(f"Error in Exponential Smoothing: {e}")

print("\nUPDATED MODEL COMPARISON:")
print(metrics_comparison.to_string(index=False))

```

UPDATED MODEL COMPARISON:

	Model	Users_MAE	Users_RMSE	Users_MAPE	Engagement_MAE	Engagement_RMSE	Engagement_MAPE
	ARIMA	3.156773	3.264593	7.077541	0.020336	0.025320	3.887037
	Moving Average	3.352284	3.856964	NaN	0.018563	0.024171	NaN
	Naive	6.766132	6.848630	15.199450	0.032230	0.039168	6.109205
	Exponential Smoothing	6.756252	6.838869	15.177173	NaN	NaN	NaN

In [ ]:

## Final Insights for Small Dataset

In [36]:

```
print("\n" + "="*60)
print("FINAL INSIGHTS FOR SMALL DATASET (21 DAYS)")
print("="*60)

print(f"\nDataset Size: {len(daily_df)} days")
print(f"Training Period: {len(train)} days")
print(f"Testing Period: {len(test)} days")

print("\nRECOMMENDATIONS FOR SMALL DATASET:")
print("1. Collect more data for better model performance")
print("2. Use simple models (Moving Average, Naive) as benchmarks")
print("3. Consider domain knowledge for seasonal patterns")
print("4. Monitor forecast performance as more data becomes available")
print("5. Use confidence intervals to understand forecast uncertainty")

print("\nNEXT STEPS:")
print("✓ Continue collecting daily data")
print("✓ Re-run analysis monthly to incorporate new data")
print("✓ Consider external variables (weekends, holidays, promotions)")
print("✓ Implement simple monitoring system with moving averages")

# Plot final comparison
best_users_model = metrics_comparison.loc[metrics_comparison['Users_MAE'].idxmin(), 'Model']
best_engagement_model = metrics_comparison.loc[metrics_comparison['Engagement_MAE'].idxmin(), 'Model']

print(f"\nBEST PERFORMING MODELS:")
print(f"Users: {best_users_model} (lowest MAE)")
print(f"Engagement Rate: {best_engagement_model} (lowest MAE)")

# Save results for future comparison
forecast_results = {
    'train_size': len(train),
    'test_size': len(test),
    'users_forecast': users_forecast,
    'engagement_forecast': engagement_forecast,
    'metrics_comparison': metrics_comparison,
    'last_date': daily_df.index.max()
}

print(f"\nAnalysis completed. Ready to incorporate new data as it becomes available.")
```

=====

FINAL INSIGHTS FOR SMALL DATASET (21 DAYS)

=====

Dataset Size: 28 days  
Training Period: 16 days  
Testing Period: 5 days

### RECOMMENDATIONS FOR SMALL DATASET:

1. Collect more data for better model performance
2. Use simple models (Moving Average, Naive) as benchmarks
3. Consider domain knowledge for seasonal patterns
4. Monitor forecast performance as more data becomes available
5. Use confidence intervals to understand forecast uncertainty

### NEXT STEPS:

- ✓ Continue collecting daily data
- ✓ Re-run analysis monthly to incorporate new data
- ✓ Consider external variables (weekends, holidays, promotions)
- ✓ Implement simple monitoring system with moving averages

### BEST PERFORMING MODELS:

Users: ARIMA (lowest MAE)  
Engagement Rate: Moving Average (lowest MAE)

Analysis completed. Ready to incorporate new data as it becomes available.

In [ ]:

# Project Summary & Insights

This project analyzed website traffic data from April to May 2024, focusing on user visits, sessions, and engagement across channels like Direct, Organic Search, and Organic Social. Using time series techniques (like trend analysis and forecasting models such as ARIMA), we uncovered patterns in hourly and daily activity to predict future performance and spot opportunities for improvement.

## Key Findings

- **Organic Social dominates** with 70+ average users/sessions, driving the highest engagement
- **Strong correlations** exist between users, sessions, and engagement (0.95-0.99)
- **Daily patterns outweigh weekly cycles** in the short timeframe analyzed
- **Both user traffic and engagement rates are predictable** with 90%+ accuracy
- **Simple models work well** - ARIMA reduced prediction error by 53% compared to naive methods

## Temporal Patterns

- Peak activity occurs during midday hours
- Consistent daily patterns with minimal weekly seasonality
- Stable overall trends with predictable fluctuations

In [ ]:

## Recommendations

- **Gather More Data:** Collect at least 3-6 months of traffic info to improve forecast accuracy and capture longer trends like holidays or seasons.
- **Focus on Top Channels:** Invest more in Organic Social for growth, while boosting underperformers like Direct with targeted promotions.
- **Optimize Timing:** Schedule content or ads for peak hours (evenings) and days (mid-week) to maximize engagement.
- **Monitor Regularly:** Re-run the analysis monthly, adding factors like events or ads, and use simple tools like moving averages for quick checks.

*Note: These recommendations are based on 28 days of data - accuracy will improve significantly with more historical data collection.*

In [ ]:

 Author

Uzoh C. Hillary - Data Scientist / Data Analyst

GitHub: <https://github.com/Uzo-Hill>

LinkedIn: <http://www.linkedin.com/in/hillaryuzoh>

In [ ]: