



MOBILE USERS BEHAVIOR DATA MANAGEMENT AND ANALYSIS

BY SOAD ATEF AND HAGER AHMED



www.reallygreatsite.com





INTRODUCTION

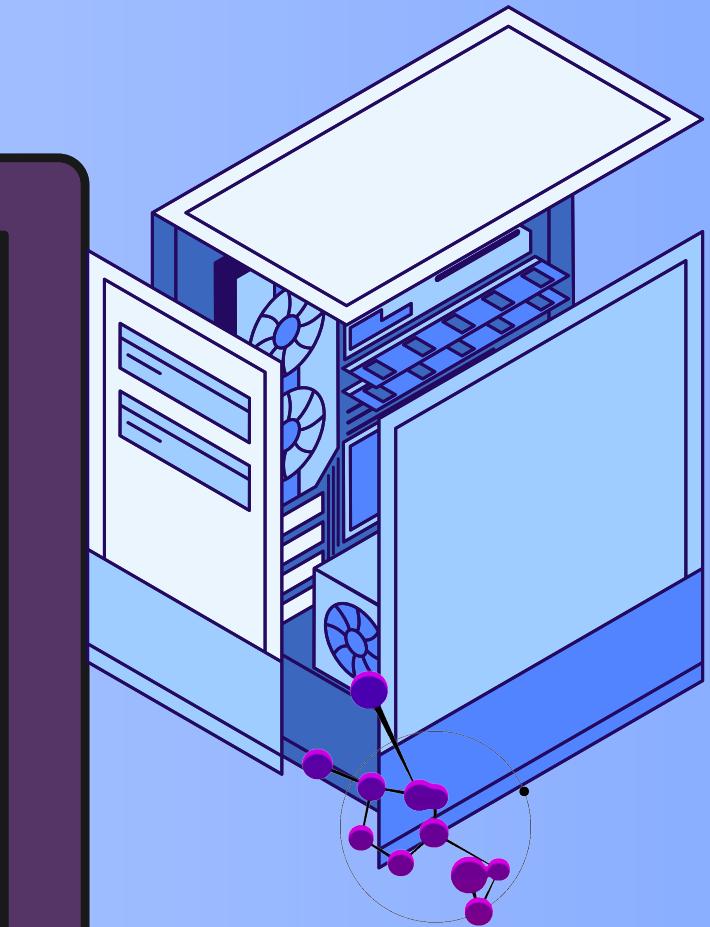
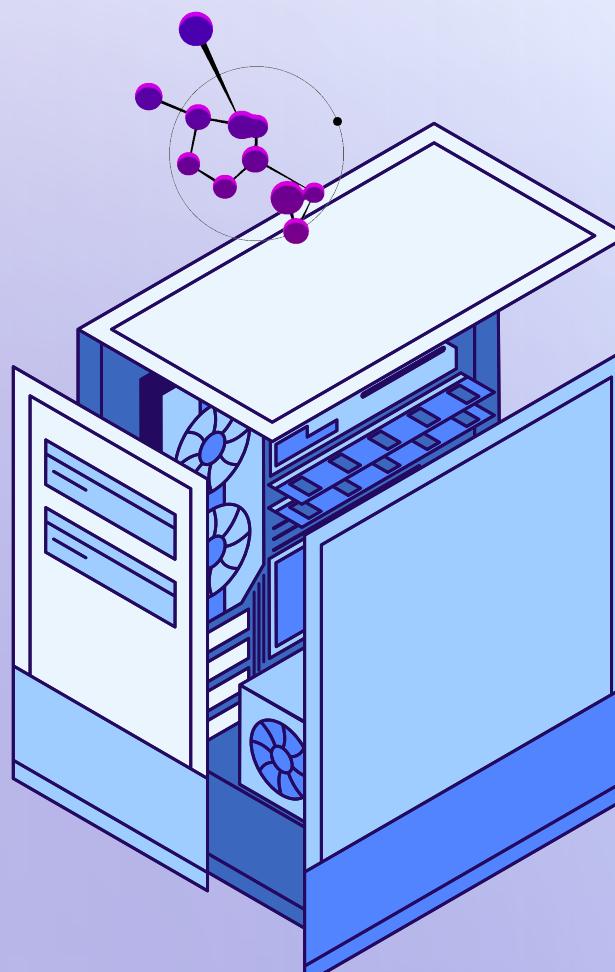
-

What is the project about?

- The project revolves around managing, cleaning, and analyzing mobile users' behavioral data. It combines SQL, SSIS (SQL Server Integration Services), and Python to handle data efficiently.

What tools were used?

- SQL for database setup and management.
- SSIS for data extraction, transformation, and loading (ETL).
- Python for data analysis and model building.





GOAL AND IMPORTANCE

- **Goal:** The goal of this project is to collect, clean, analyze, and store mobile user behavior data to gain insights that help improve services, marketing strategies, and user satisfaction.
- **Importance:**
 - Informed decision-making based on data.
 - Insights into user trends, preferences, and patterns.
 - Support for product improvement and customer satisfaction strategies.
 - Data-driven decisions are critical for competitive business strategies.
 - The ability to handle large datasets is essential for mobile user behavior analysis.
 - Combining SQL, SSIS, and Python provides an efficient solution for comprehensive data management and analysis





KEY CONCEPTS



- SQL Database to manage user data efficiently.
- SSIS for loading and transforming large datasets.
- Python for further data cleaning and advanced analysis.



DATA OVERVIEW



- **Dataset:** Mobile Users Behavior Dataset
- **Source:** The dataset was sourced from user interactions across various mobile platforms.
- **Number of Records:** Over 700 user behavior records.
- **Main Features:**
 - User ID
 - Device Model
 - App Usage
 - Operating_System
 - Device ID
 - App_Usage_Time
 - Screen_On_Time
 - Battery_Drain
 - Number_of_Apps_Installed
 - Data_Usage
 - Age
 - Gender
 - User_Behavior_Class
 - OS_ID



DATA CLEANING

Handled Missing Data:

- Identified and managed missing values in columns such as app usage and interaction type. Methods like mean imputation and removal of records were used.

Removed Duplicates:

- Ensured data integrity by identifying and removing duplicate records.



```
1 import pandas as pd
2
3 data = pd.read_csv("user_behavior.csv")
4 data.columns = data.columns.str.replace(' ', '_')
5
6 DeviceModel = data["Device_Model"]
7 OS = data["Operating_System"]
8 Gender = data["Gender"]
9 DeviceModelNumeric = {
10     "Google Pixel 5": 1,
11     "OnePlus 9": 2,
12     "iPhone 12": 3,
13     "Xiaomi Mi 11": 4,
14     "Samsung Galaxy S21": 5,
15 }
16 OSNumeric = {
17     "Android": 0,
18     "iOS": 1,
19 }
20 GenderNumeric = {
21     "Male": 0,
22     "Female": 1,
23 }
24
25
26 data["Device_ID"] = data["Device_Model"].replace(DeviceModelNumeric)
27 data["OS_ID"] = data["Operating_System"].replace(OSNumeric)
28 data["Gender"] = data["Gender"].replace(GenderNumeric)
29
30
31 output_path = "analuser_behavior_dataset.csv"
32 data.to_csv(output_path, index=False)
```

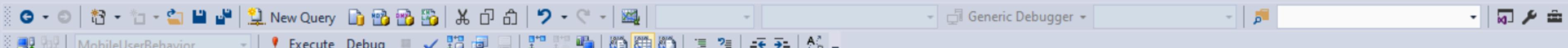
In [6] Down Mean linear lasso Ridge Logistic C:\Us 678: number regular mod C:\Us Conve STOP: Incre h Please h n_i To [7]

```
1 import pandas as pd
2
3 # Assuming you have a DataFrame df with 'User_ID' and 'OS_ID'
4 df = pd.read_csv("analuser_behavior_dataset.csv")
5
6 # Create a new composite key column by concatenating 'User_ID' and 'OS_ID'
7 df['Composite_Key1'] = (df['User_ID'].astype(str) + df['OS_ID'].astype(str)).astype(int)
8 df['Composite_Key2'] = df['User_ID'] * 1000
9
10 df = df.drop(columns=['Device_ID', 'OS_ID'])
11
12 df = df.rename(columns={'Composite_Key2': 'Device_ID'})
13 df = df.rename(columns={'Composite_Key1': 'OS_ID'})
14 # Now 'Composite_Key' is a combination of both columns
15 print(df)
16 df.to_excel('analuser_behavior_dataset.xls', index=False, engine='openpyxl')
```

Console 1/A X

In [1]:

New Menu File Edit View Query Project Debug Tools Window Help



Object Explorer

Connect ▾

- LAPTOP-9M36E8AP (SQL Server 16.0) ▾
 - Databases
 - + System Databases
 - + Database Snapshots
 - + AdventureWorks
 - + db
 - + dd
 - MobileUserBehavior
 - + Database Diagrams
 - + Tables
 - + System Tables
 - + FileTables
 - + External Tables
 - + dbo.Dim_Device
 - + dbo.Dim_Operating_System
 - + dbo.Dim_User
 - + dbo.Fact_User_Device_Usage
 - + Views
 - + External Resources
 - + Synonyms
 - + Programmability
 - + Query Store
 - + Service Broker
 - + Storage
 - + Security
 - + populationHealthDB
 - + TSQL
 - + US_Accidents
 - + work
 - + Security
 - + Server Objects
 - + Replication
 - + PolyBase
 - + AlwaysOn High Availability

SQLQuery2.sql - not connected*

```
1 CREATE TABLE Dim_User (
2     User_ID INT PRIMARY KEY,
3     Age INT,
4     Gender INT,
5     User_Behavior_Class INT
6 );
7 CREATE TABLE Dim_Device (
8     Device_ID INT PRIMARY KEY,
9     Device_Model VARCHAR(50)
10 );
11 CREATE TABLE Dim_Operating_System (
12     OS_ID INT PRIMARY KEY,
13     Operating_System VARCHAR(50)
14 );
15 CREATE TABLE Fact_User_Device_Usage (
16     User_ID INT,
17     Device_ID INT,
18     OS_ID INT,
19     App_Usage_Time FLOAT,
20     Screen_On_Time FLOAT,
21     Battery_Drain FLOAT,
22     Number_of_Apps_Installed INT,
23     Data_Usage FLOAT,
24     FOREIGN KEY (User_ID) REFERENCES Dim_User(User_ID),
25     FOREIGN KEY (Device_ID) REFERENCES Dim_Device(Device_ID),
26     FOREIGN KEY (OS_ID) REFERENCES Dim_Operating_System(OS_ID)
27 );
28
29
```

110 %

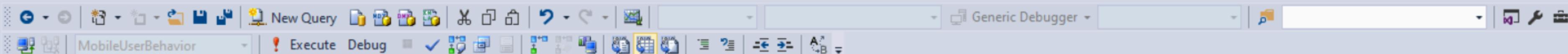
Messages

Command(s) completed successfully.

110 %

Disconnected.

New Menu File Edit View Query Project Debug Tools Window Help



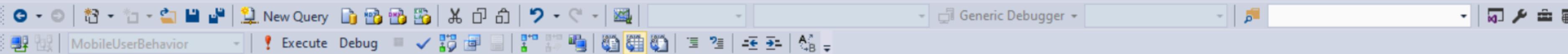
MobileUserBehavior | ! Execute Debug | SQLQuery11.sql - not connected* | SQLQuery17.sql - not connected | model.sql - not connected | SQLQuery16.sql - not connected* | SQLQuery15.sql - not connected* | SQLQuery12.sql - not connected*

```
1 CREATE TABLE temp_Dim_User (
2     User_ID INT PRIMARY KEY,
3     Age INT,
4     Gender INT,
5     User_Behavior_Class INT
6 );
7 INSERT INTO temp_Dim_User
8 SELECT DISTINCT *
9 FROM Dim_User |
10 DELETE FROM Dim_User
11 INSERT INTO Dim_User
12 SELECT * FROM temp_Dim_User
13 DROP TABLE temp_Dim_User
14 CREATE TABLE temp_Dim_Device (
15     Device_ID INT PRIMARY KEY,
16     Device_Model VARCHAR(50),
17     Brand VARCHAR(50)
18 );
19 INSERT INTO temp_Dim_Device
20 SELECT DISTINCT *
21 FROM Dim_Device
22 DELETE FROM Dim_Device
23 INSERT INTO Dim_Device
24 SELECT * FROM temp_Dim_Device
25 DROP TABLE temp_Dim_Device
26 CREATE TABLE temp_Dim_Operating_System (
27     OS_ID INT PRIMARY KEY,
28     Operating_System VARCHAR(50)
29 );
30 INSERT INTO temp_Dim_Operating_System
31 SELECT DISTINCT *
32 FROM Dim_Operating_System
33 DELETE FROM Dim_Operating_System
34 INSERT INTO Dim_User
```

110 %

Disconnected.

New Menu File Edit View Query Project Debug Tools Window Help



MobileUserBehavior | Execute Debug

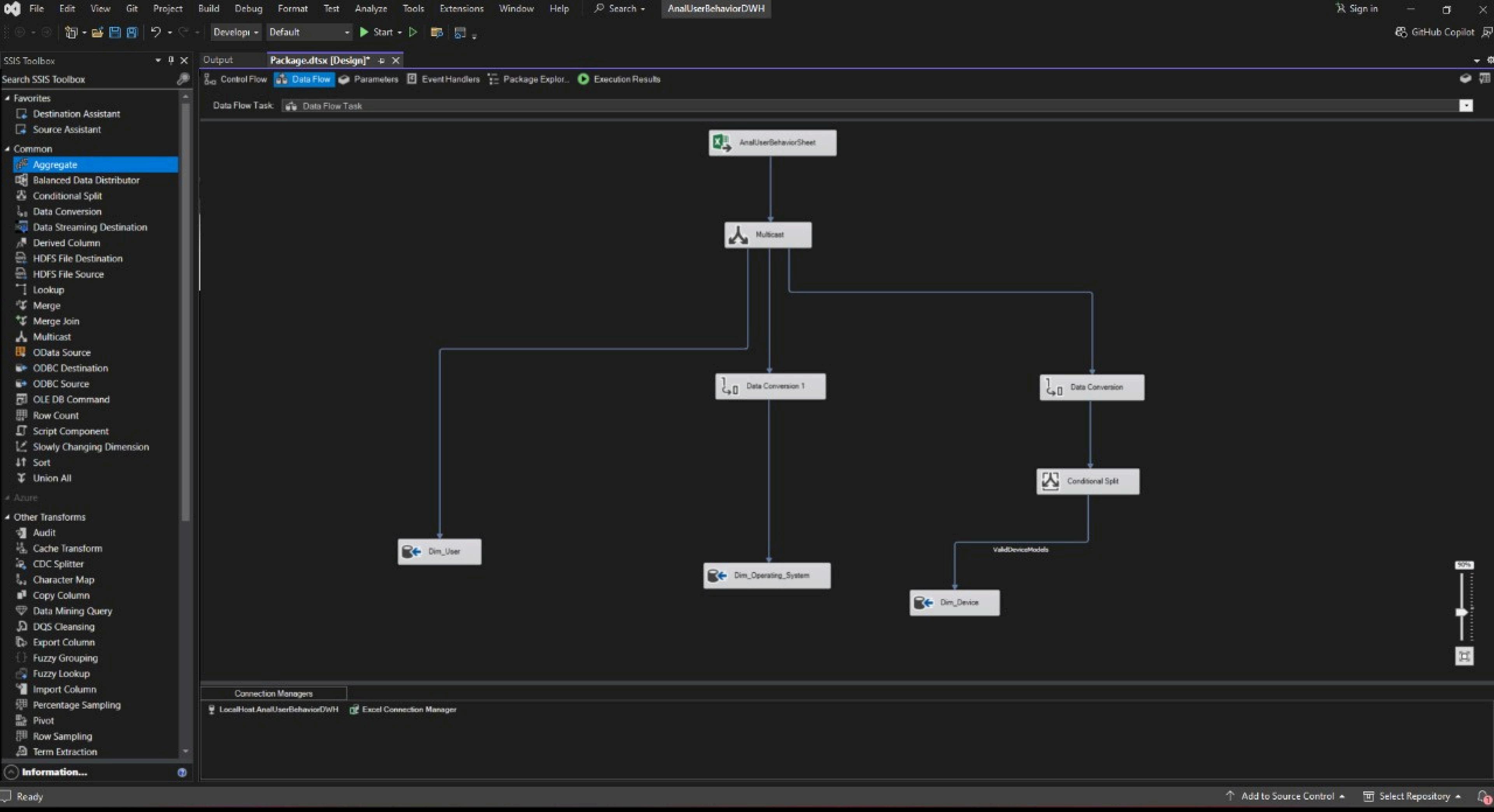
SQLQuery11.sql - not connected* SQLQuery17.sql - not connected model.sql - not connected SQLQuery16.sql - not connected* SQLQuery15.sql - not connected* SQLQuery12.sql - not connected*

```
31  SELECT DISTINCT *
32  FROM Dim_Operating_System
33  DELETE FROM Dim_Operating_System
34  INSERT INTO Dim_User
35  SELECT * FROM temp_Dim_Operating_System
36  DROP TABLE temp_Dim_Operating_System
37  CREATE TABLE temp_Fact_User_Device_Usage (
38      User_ID INT,
39      Device_ID INT,
40      OS_ID INT,
41      App_Usage_Time FLOAT,
42      Screen_On_Time FLOAT,
43      Battery_Drain FLOAT,
44      Number_of_Apps_Installed INT,
45      Data_Usage FLOAT
46      FOREIGN KEY (User_ID) REFERENCES Dim_User(User_ID),
47      FOREIGN KEY (Device_ID) REFERENCES Dim_Device(Device_ID),
48      FOREIGN KEY (OS_ID) REFERENCES Dim_Operating_System(OS_ID)
49 );
50  INSERT INTO temp_Fact_User_Device_Usage
51  SELECT DISTINCT *
52  FROM Fact_User_Device_Usage
53  DELETE FROM Fact_User_Device_Usage
54  INSERT INTO Fact_User_Device_Usage
55  SELECT * FROM temp_Fact_User_Device_Usage
56  DROP TABLE temp_Fact_User_Device_Usage
```



110 %

Disconnected.





رواد مصر الرقمية

FEATURE SELECTION

- **Why Feature Selection?** To remove irrelevant data and focus on features that would significantly impact the analysis.
- **Key Features Selected:**
 - Number of apps installed
 - Battery Drain
 - Data Usage



```
1 from sklearn.feature_selection import SelectKBest ,mutual_info_regression,chi2
2 from sklearn.preprocessing import LabelEncoder
3 import pandas as pd
4
5 file_path = 'cleaned_user_behavior_dataset.csv'
6 df = pd.read_csv(file_path)
7
8 X= df.drop('User Behavior Class' , axis=1)
9 Y= df[ 'User Behavior Class']
10
11 le = LabelEncoder()
12 for columns in X.columns:
13     X[columns]= le.fit_transform(X[columns])
14
15 Y = le.fit_transform(Y)
16
17 selector = SelectKBest(score_func= mutual_info_regression, k=3)#score_fun -> method of selection , k -> number of feature
18 selector.fit_transform(X,Y)
19 selected_feature = X.columns[selector.get_support()].values#-> select columns with True values
20 print(X[selected_feature])
```

History Help Variable Explorer Plots Files

Console 1/A X

In [7]: runfile('C:/Users/soada/Downloads/untitled23.py', wdir='C:/Users/soada/Downloads')

	Battery Drain (mAh/day)	Number of Apps Installed	Data Usage (MB/day)
0	391	54	
376	275	30	
1	158	21	
316	350	44	
2	280	46	
108	
3	220	15	
292	407	55	
4	190	11	
332	47	3	
..	266	37	
695	
131	
696	
397	
697	
162	
698	
58	
699	
279	

[700 rows x 3 columns]



MODEL SELECTION AND TRAINING

- **What kind of models?** Machine learning models were developed using Python for predictive analytics. **Models tested included:**
 - Logistic Regression for user retention prediction.
 - Linear Regression
 - Lasso
 - Ridge
- **Training Process:** The dataset was split into training and test sets. The models were trained using the training data and validated against the test data.



```

1 import pandas as pd
2 from sklearn.linear_model import LogisticRegression, LinearRegression, Ridge, Lasso
3 from sklearn.metrics import mean_squared_error
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import PolynomialFeatures
6 from sklearn.pipeline import make_pipeline
7
8 data = pd.read_csv("analuser_behavior_dataset.csv")
9
10 x = data.drop(['Device_Model', 'Operating_System', 'User_Behavior_Class'], axis=1)
11 y = data["User_Behavior_Class"]
12 xTrain, xTest, yTest = train_test_split(x, y, test_size=0.2, random_state=42)
13 lr = LinearRegression()
14 lr.fit(xTrain, yTrain)
15 yPredict = lr.predict(xTest)
16 mse = mean_squared_error(yTest, yPredict)
17 print("Mean Squared Error : ", mse)
18 print("linear regression accuracy : ", lr.score(xTest, yTest))
19 # ----- Lasso -----
20 lasso = make_pipeline(PolynomialFeatures(degree=2), Lasso(alpha=0.3))
21 lasso.fit(xTrain, yTrain)
22 print("Lasso accuracy : ", lasso.score(xTest, yTest))
23 # ----- Rigid -----
24 rg = Ridge(alpha=0.3)
25 rg.fit(xTrain, yTrain)
26 yPredict = rg.predict(xTest)
27 print("Ridge accuracy : ", rg.score(xTest, yTest))
28 # ----- Logistic Regression -----
29 from sklearn.linear_model import LogisticRegression
30 from sklearn.metrics import accuracy_score
31
32 lr = LogisticRegression()
33 lr.fit(xTrain,yTrain)
34 Y_predict = lr.predict(xTest)
35 print('Logistic Regression accuracy : ',accuracy_score(yTest, Y_predict))
36

```

Console 1/A

```

In [6]: runfile('C:/Users/soada/Downloads/untitled19.py', wdir='C:/Users/soada/Downloads')
Mean Squared Error :  0.0318563757754207
linear regression accuracy :  0.9825297995187956
lasso accuracy :  0.9902959564068873
Ridge accuracy :  0.982529917472532
Logistic Regression accuracy :  0.5785714285714286
C:\Users\soada\anaconda3\Lib\site-packages\sklearn\linear_model\_coordinate_descent.py:678: ConvergenceWarning: Objective did not converge. You might want to increase the number of iterations, check the scale of the features or consider increasing regularisation. Duality gap: 5.801e+00, tolerance: 1.117e-01
    model = cd_fast.enet_coordinate_descent(
C:\Users\soada\anaconda3\Lib\site-packages\sklearn\linear_model\_logistic.py:469: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
n_iter_i = _check_optimize_result()

In [7]:

```



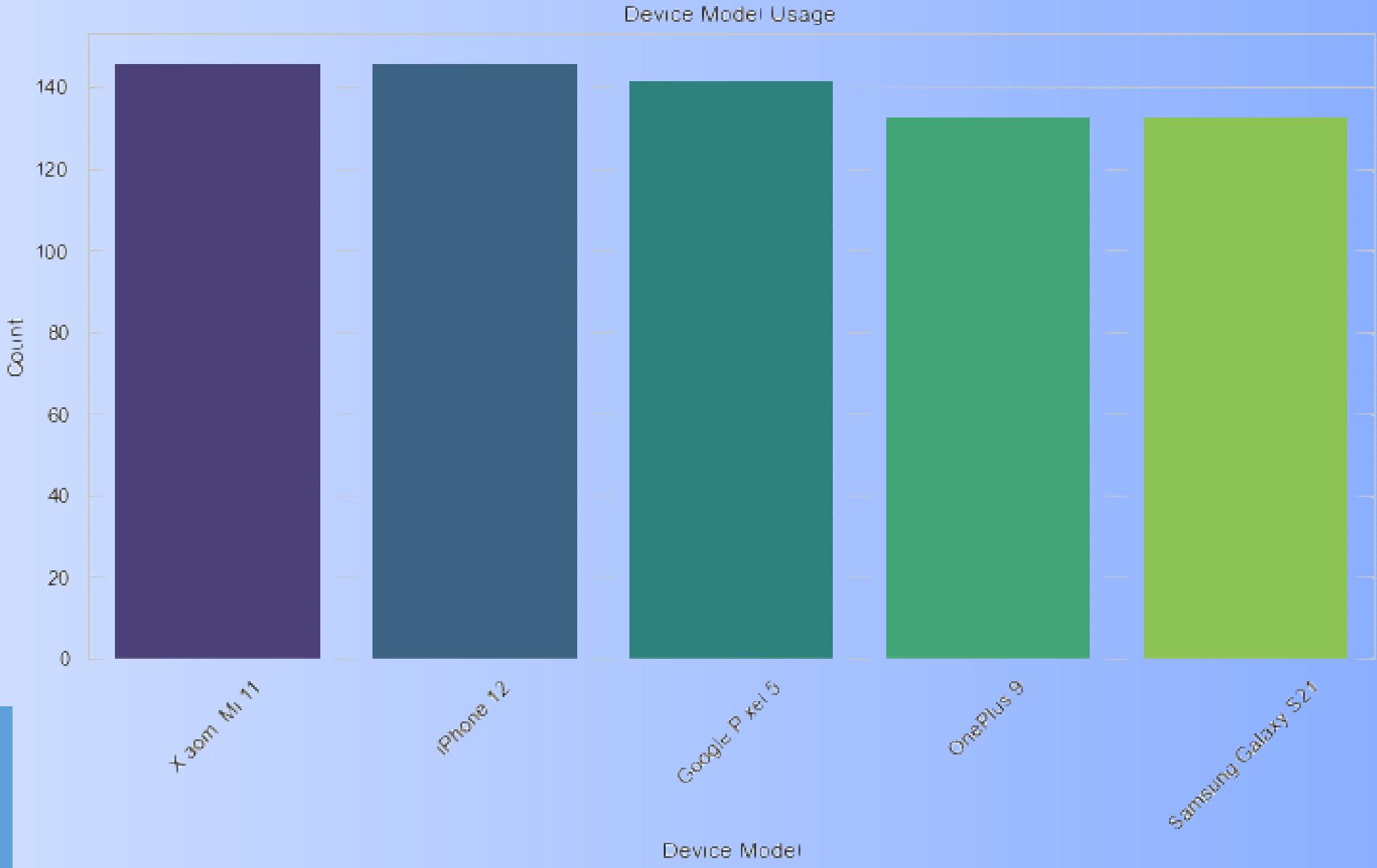
DATA VISUALIZATION

- **Why Visualization?** To gain better insights from the data through graphical representation.
- **Visualizations Created:**
- Bar Charts
- HeatMap
- Box Plot
- Histogram
- Line Graph





DATA VISUALIZATION



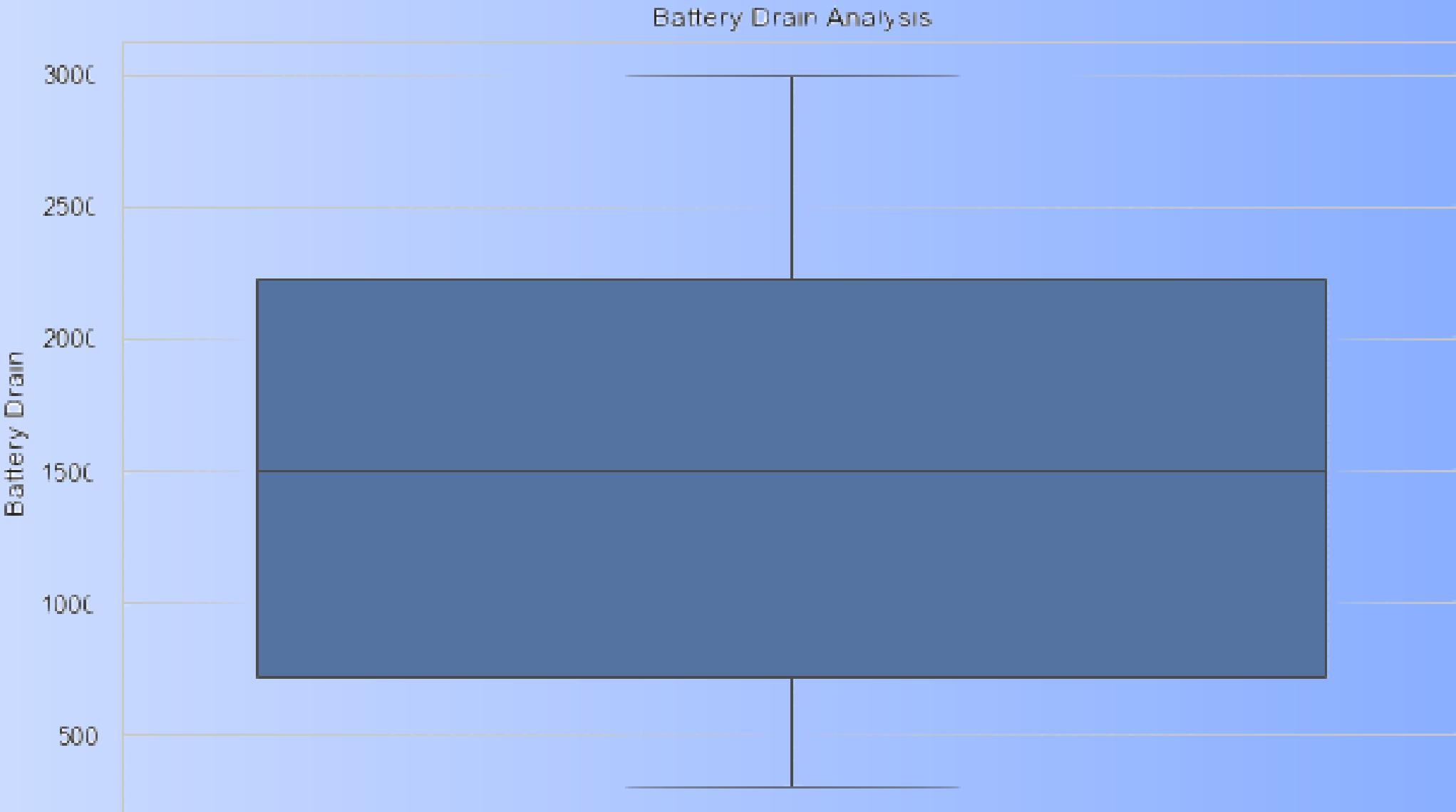
Device Model Usage

Type: Bar Chart

Purpose: Display the number of users for each device model to identify popular devices among users



DATA VISUALIZATION



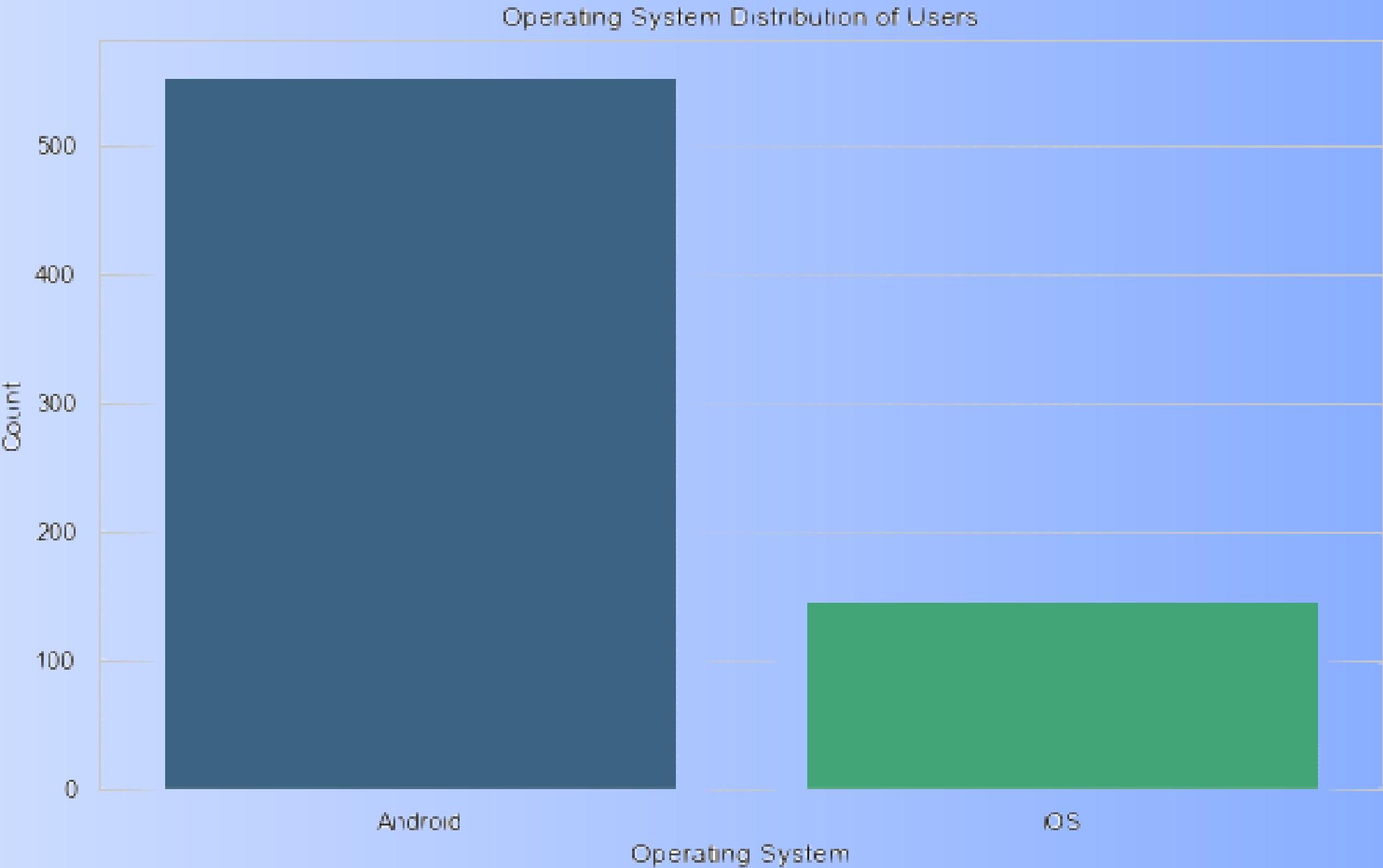
Battery Drain Analysis

Type: Box Plot

Purpose: Show the distribution of battery drain values to identify outliers and the general trend of battery usage among users.



DATA VISUALIZATION



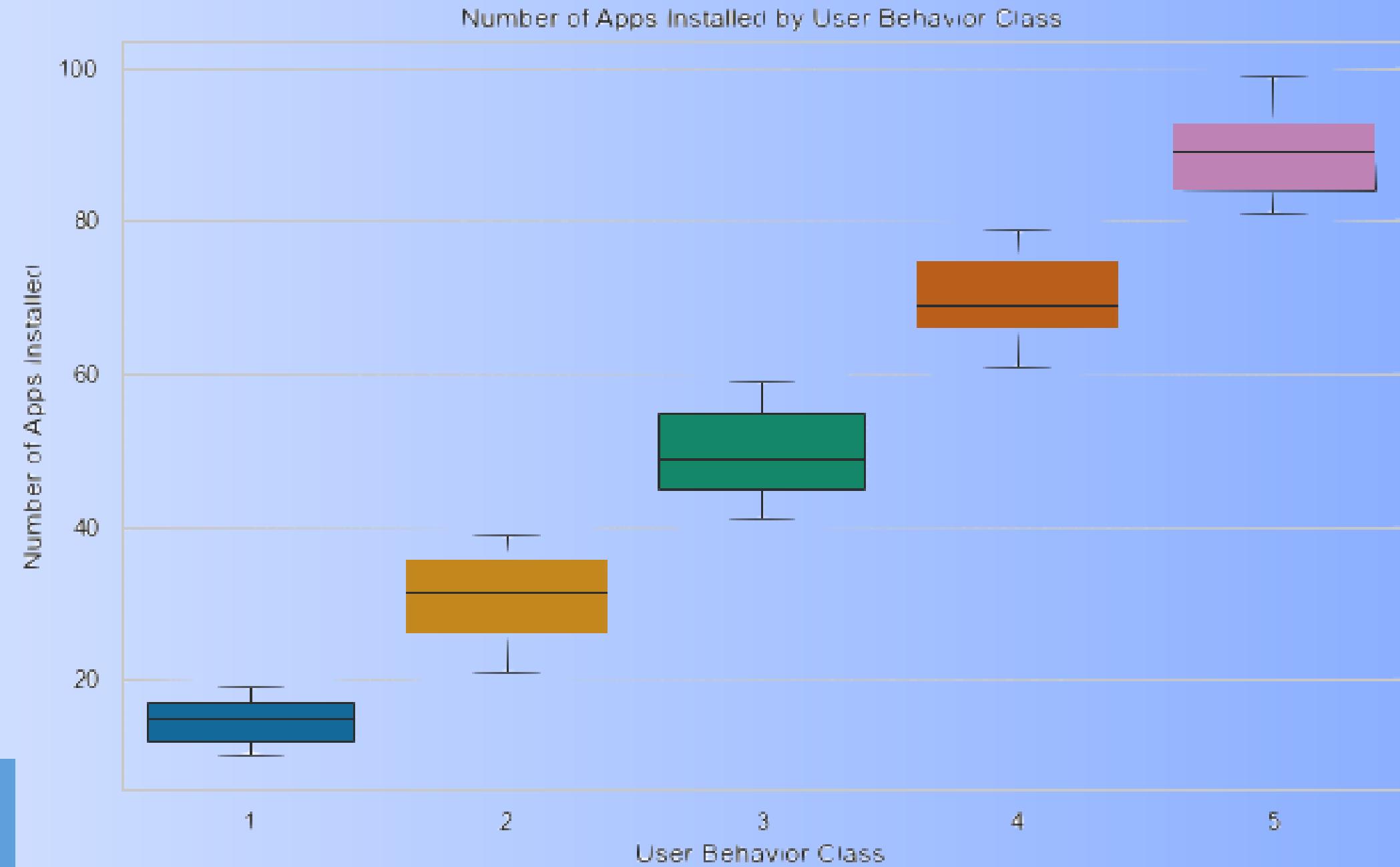
Operating System Distribution

Type: Bar Chart

Purpose: Illustrate the distribution of users across different operating systems (e.g., iOS, Android).



DATA VISUALIZATION



Number of Apps Installed vs. User Behavior Class

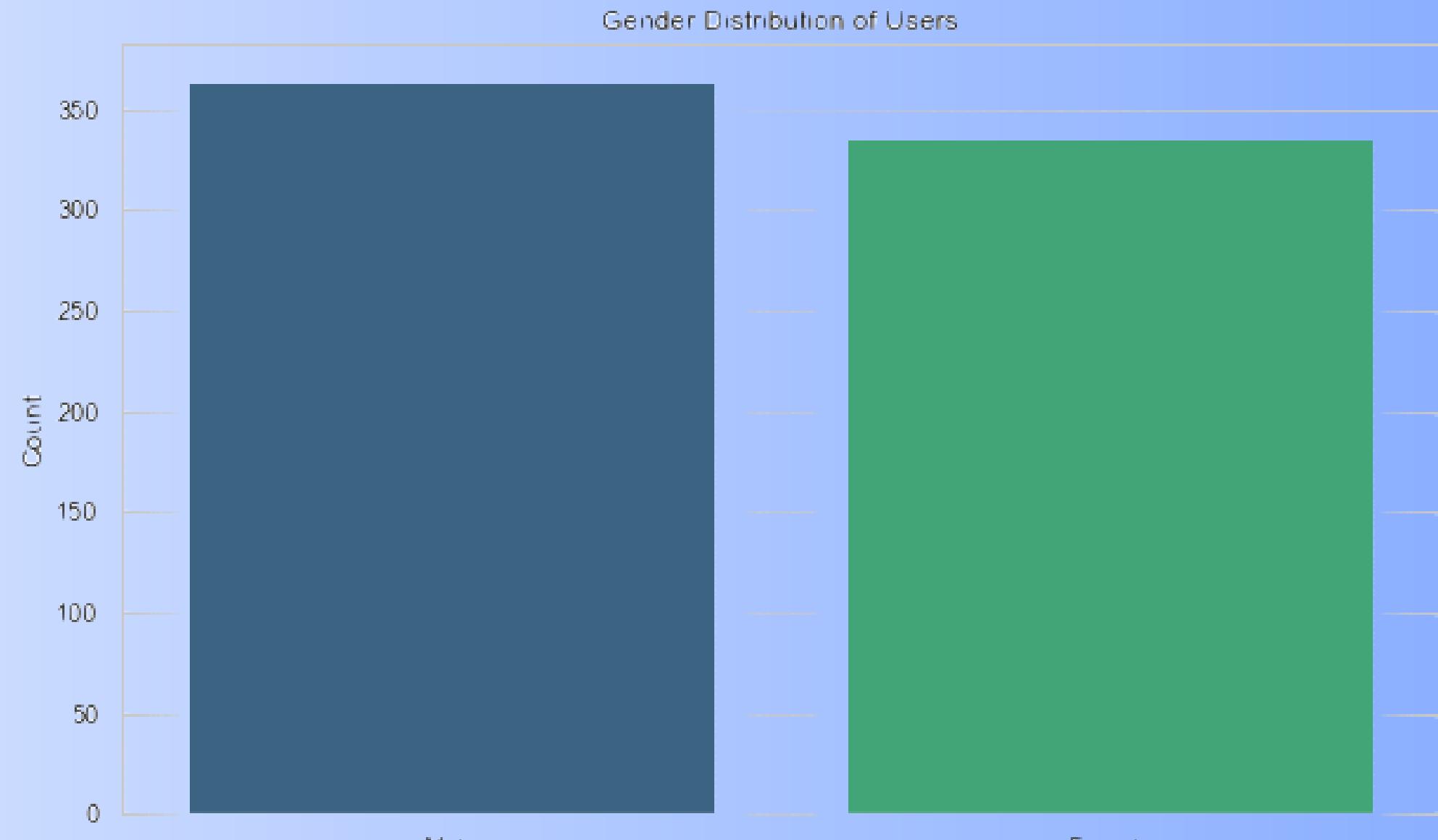
Type: Box Plot

Purpose: Compare the number of apps installed across different user behavior classes to identify any trends.



رواد مصر الرقمية

DATA VISUALIZATION



Gender Distribution

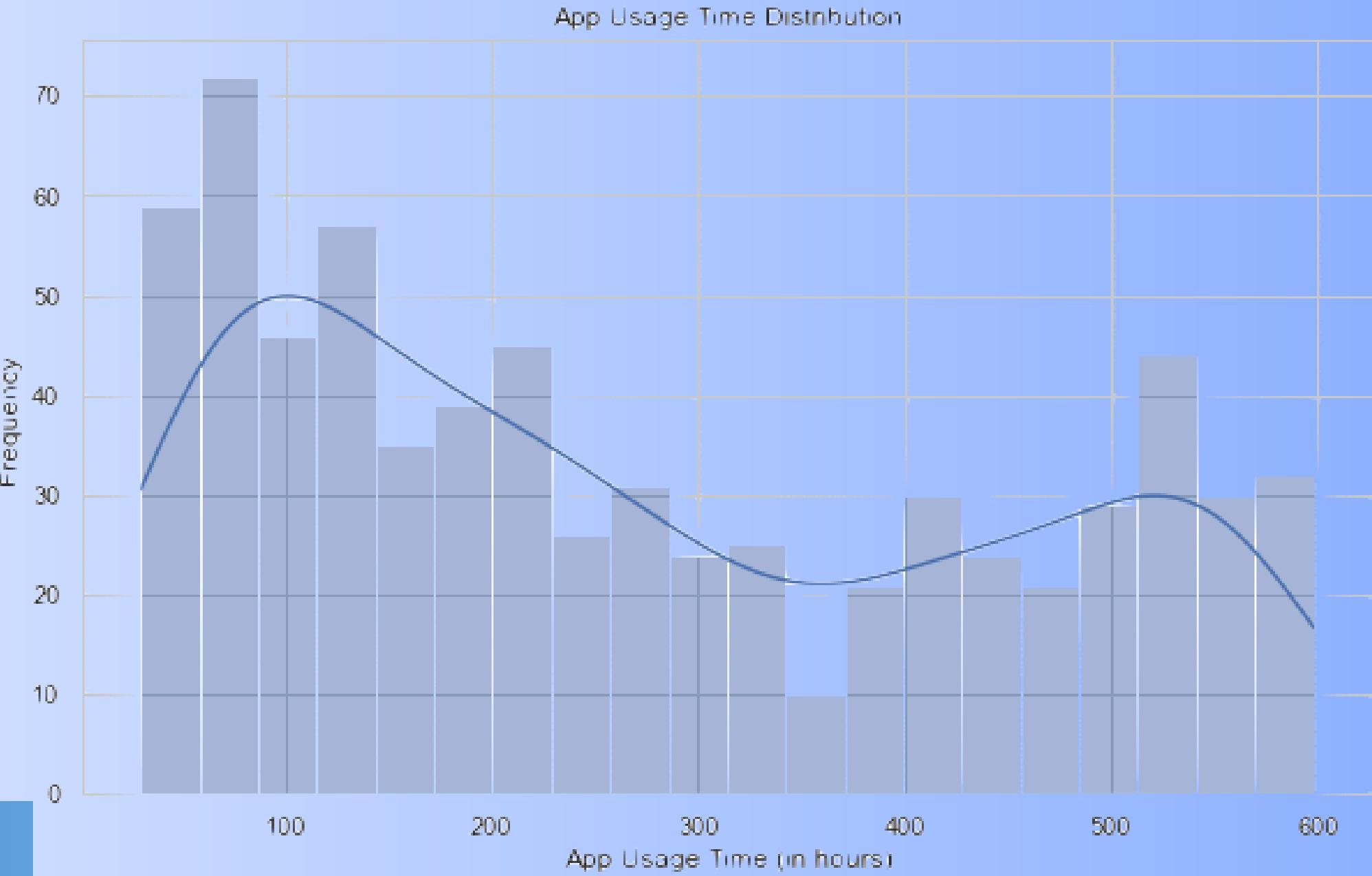
Type: Bar Chart

Purpose: Visualize the proportion of male and female users to assess gender representation in the dataset.



DATA

VISUALIZATION



App Usage Time distribution

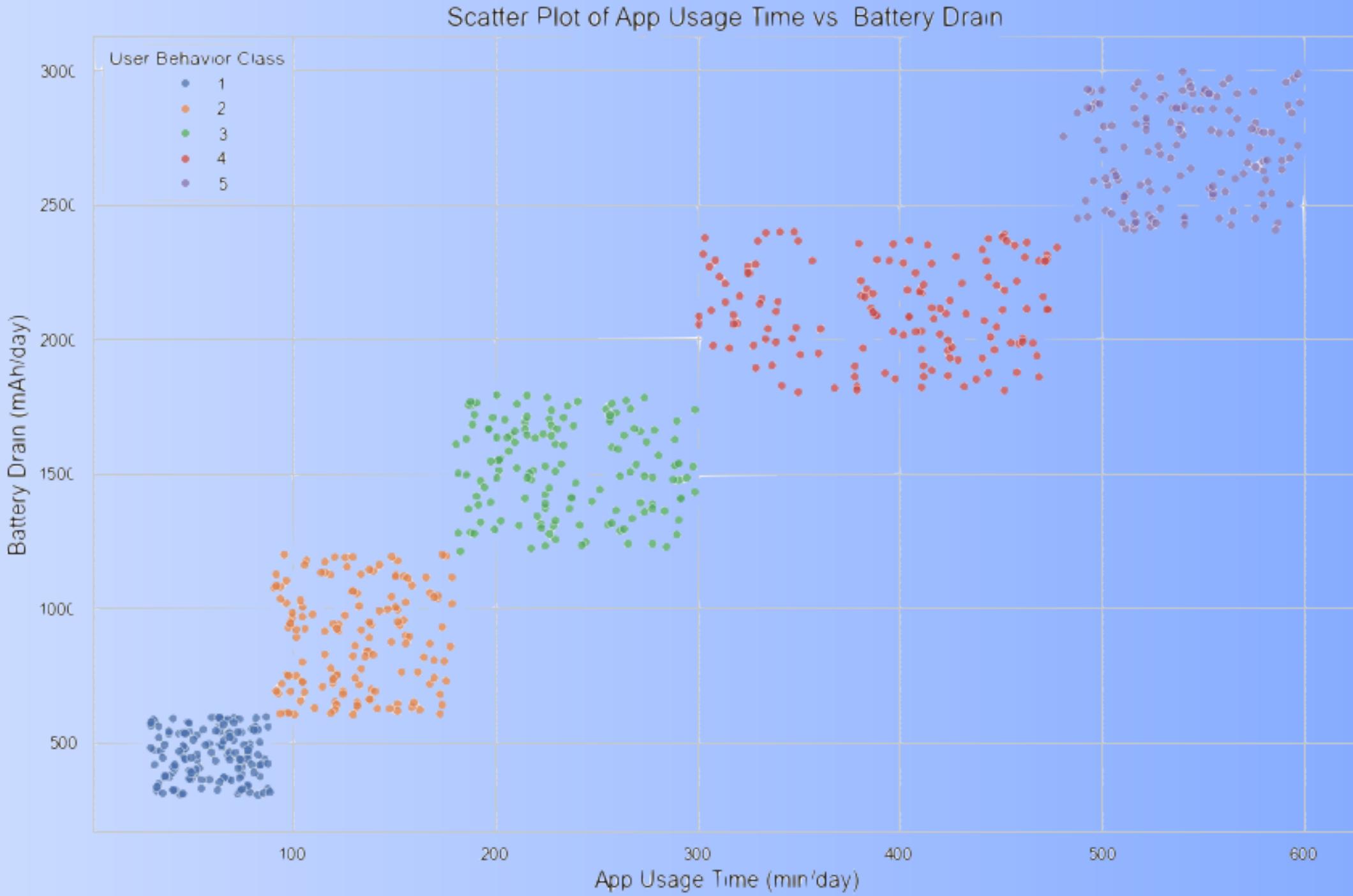
Type: Histogram

Purpose: Analyze the distribution of app usage time to see how much time users spend on their devices.



DATA

VISUALIZATION

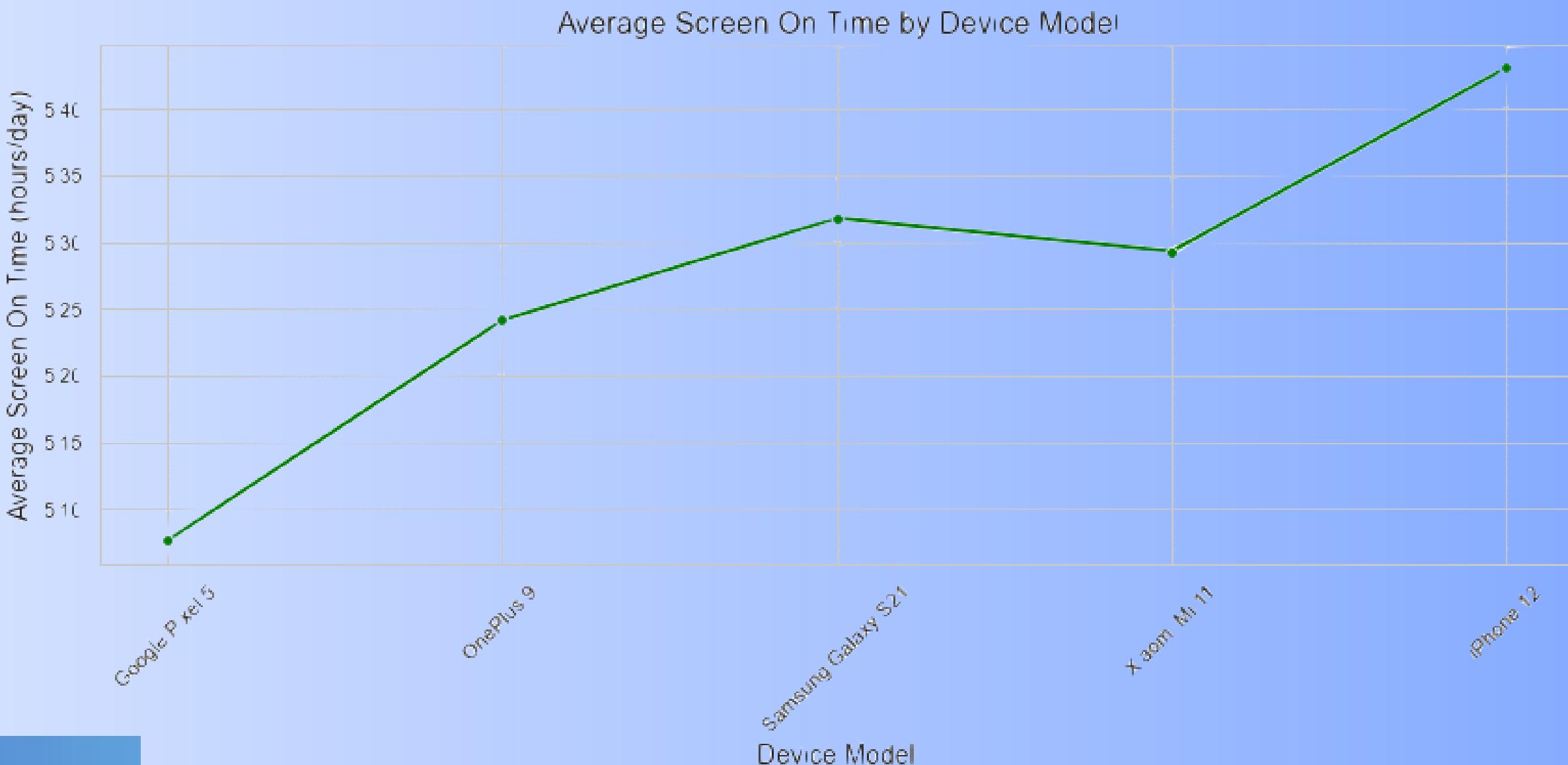


App Usage Time VS Battery Drain
Type: Scatter Plot

Purpose: Analyze the distribution of app usage time to see how much battery is drained



DATA VISUALIZATION



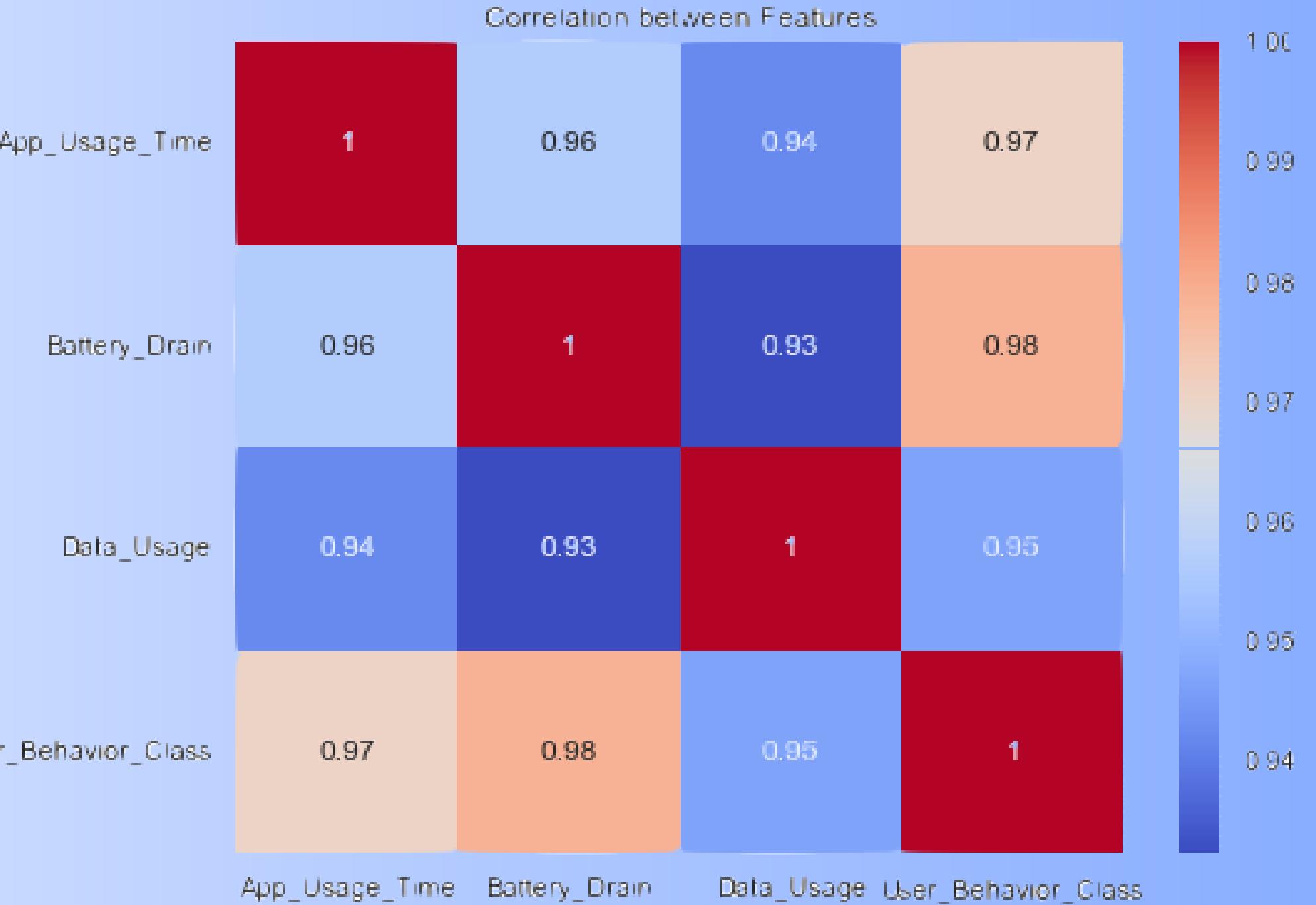
Average Screen On Time by Device Model

Type: Line Graph

Purpose: Analyze the average screen time



DATA VISUALIZATION



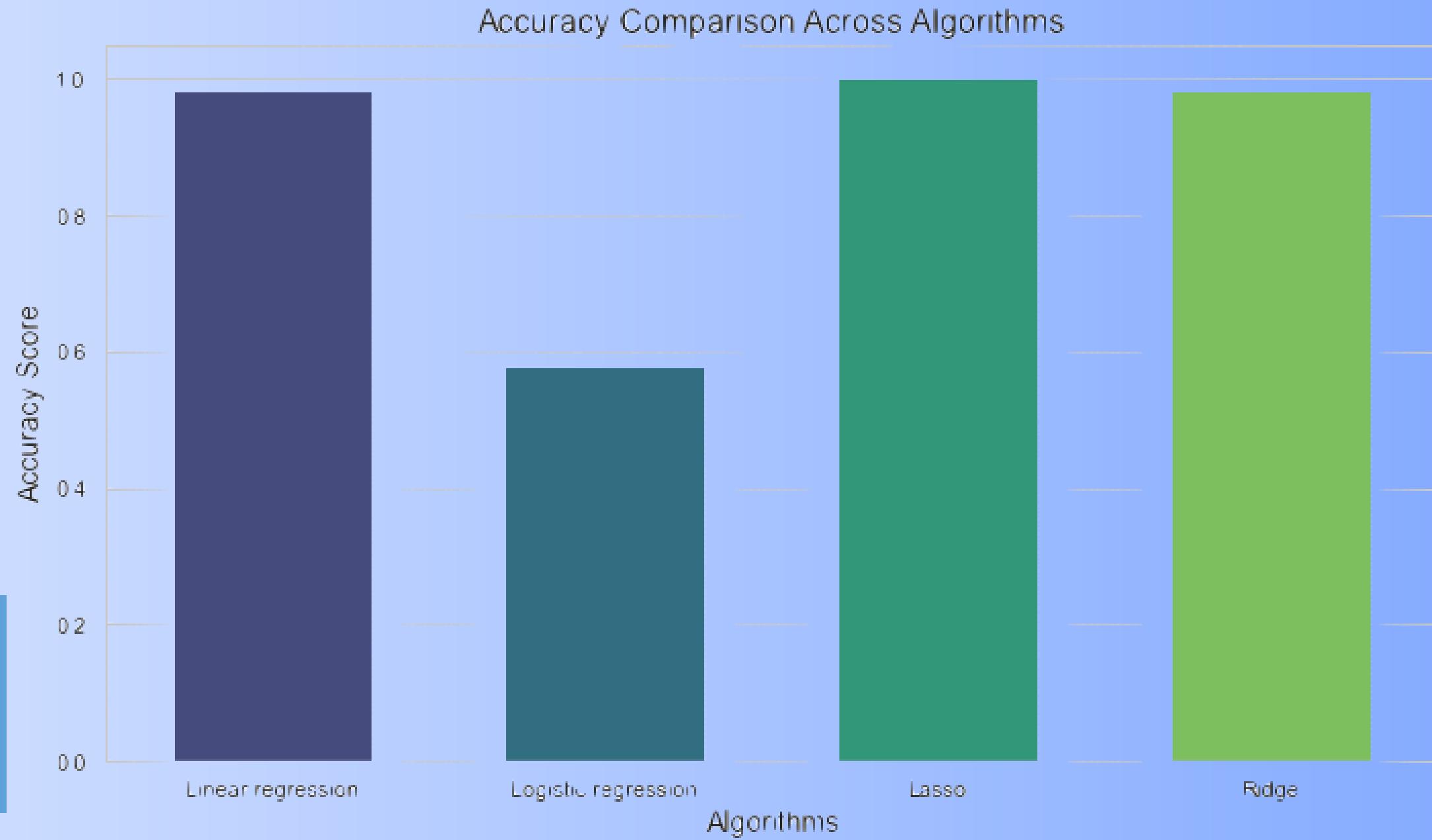
Correlation Between Features

Type: Table



رواد مصر الرقمية

DATA VISUALIZATION



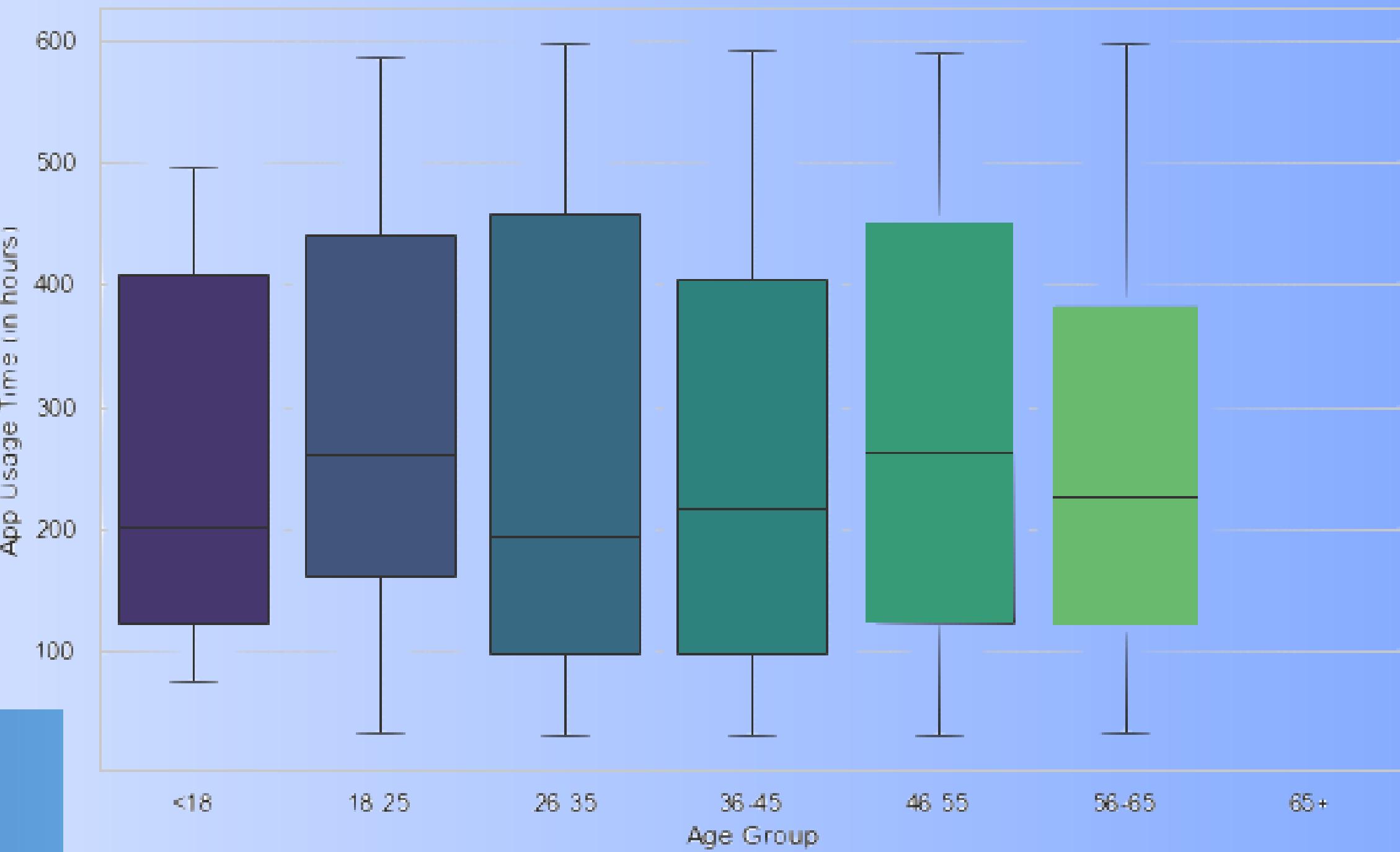
Accuracy Comparison Across Algorithms
Type: Bar Graph



DATA VISUALIZATION



App Usage Time by Age Group



Type: Box plot



RESULTS

Insights Gained:

The most used apps are entertainment-based and are used by men with the age between 26-35

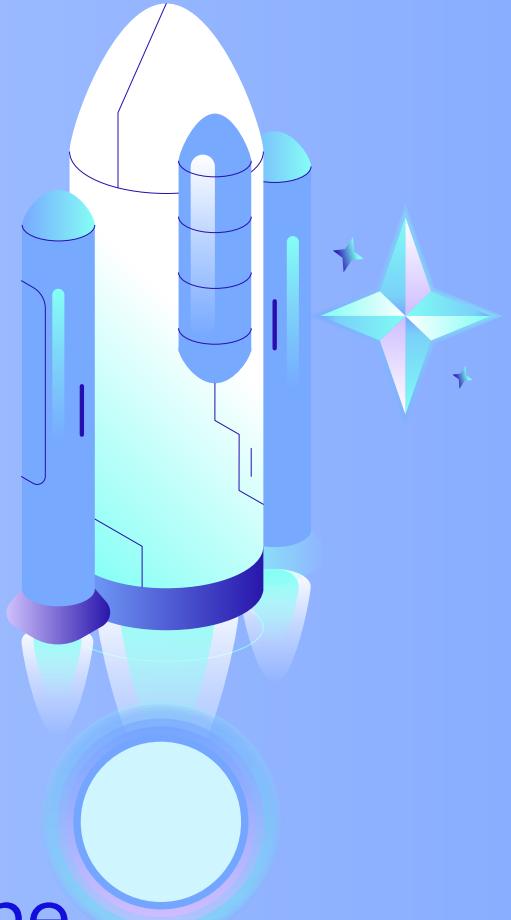
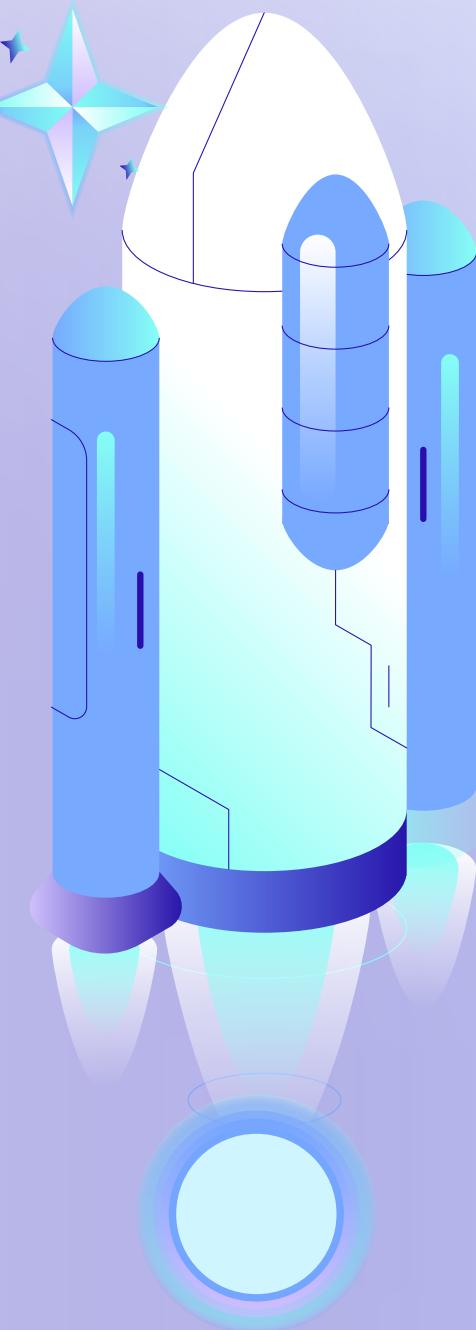
Device type significantly influences app usage patterns.

Model Performance:

Lasso achieved 99% accuracy in predicting user retention.

Logistic Regression had an accuracy of 57% in classifying app behavior





CONCLUSION

Key Takeaways:

- Data management and cleaning are crucial steps that determine the quality of analysis.
- Machine learning models provide actionable insights into user behavior.
- Visualizations are an effective tool to communicate findings.

Future Work: Implement additional models to predict user behavior more accurately and optimize the database for faster processing.





THANK YOU!



FROM: HAGER AHMED AND SOAD ATEF