

EV Charging Session Peak Classification

Dissertation submitted in fulfilment of the requirements for the Degree of

BACHELOR OF TECHNOLOGY

in

COMPUTER SCIENCE AND ENGINEERING

Data Science and Machine Learning

By

PENAGANTI SAI

Registration No: 12219418

Section: K22UP

Roll No: A30

Supervisor

Himanshu Gajanan Tikle



School of Computer Science and Engineering

Lovely Professional University

Phagwara, Punjab (India)

May 2025

DECLARATION STATEMENT

I hereby declare that the research work reported in the dissertation/dissertation proposal entitled **EV Charging Session Peak Classification** in partial fulfilment of the requirement for the award of Degree for Bachelor of Technology in Computer Science and Engineering at Lovely Professional University, Phagwara, Punjab is an authentic work carried out under supervision of my research supervisor **Mr. Himanshu Gajanan Tikle**. I have not submitted this work elsewhere for any degree or diploma.

I understand that the work presented herewith is in direct compliance with Lovely Professional University's Policy on plagiarism, intellectual property rights, and highest standards of moral and ethical conduct. Therefore, to the best of my knowledge, the content of this dissertation represents authentic and honest research effort conducted, in its entirety, by me. I am fully responsible for the contents of my dissertation work.

Signature of Candidate

PENAGANTI SAI

Reg.no : **12219418**

SUPERVISOR'S CERTIFICATE

This is to certify that the work reported in the B.Tech Dissertation/dissertation proposal entitled “**EV Charging Session Peak Classification**”, submitted by **Penaganti sai** at **Lovely Professional University, Phagwara, India** is a bonafide record of his / her original work carried out under my supervision. This work has not been submitted elsewhere for any other degree.

Signature of Supervisor

Mr.Himanshu Gajanan Tikle

Date:02-05-2025

INTRODUCTION

As electric vehicle (EV) adoption continues to rise, the growing number of EVs on the road presents significant operational challenges for both charging station providers and power grid operators. Among the most pressing issues are charging station congestion and strain on the power grid during peak hours. These challenges result in longer wait times for EV users, increased electricity costs, and inefficient energy distribution, ultimately limiting the scalability and reliability of EV infrastructure.

This project aims to address these concerns by developing a machine learning model that classifies EV charging sessions as either peak or off-peak. By learning from historical charging data, the model can accurately predict demand periods, enabling key stakeholders to take proactive measures. Power grid operators can optimize energy distribution, charging stations can implement dynamic pricing models, and users can be guided to charge during low-demand periods, reducing operational strain and improving efficiency.

The dataset used in this project contains detailed charging session records, including features such as time of day, energy consumed, session duration, and calendar-based attributes like weekdays and holidays. We employ data preprocessing, feature selection using ANOVA F-statistics, and model training using algorithms like Random Forest to ensure high predictive accuracy.

By classifying charging sessions with precision, this solution supports smarter load management, cost-effective pricing, and a better user experience. It also lays the groundwork for future enhancements, such as integrating weather, traffic, or real-time demand data. In doing so, our project contributes meaningfully toward building a sustainable, intelligent EV charging ecosystem that can scale with growing demand.

1. Problem Understanding & Definition

1.1 Problem Statement:

As electric vehicle (EV) adoption grows, charging station congestion and power grid strain have become major challenges. Power grid operators and charging station providers struggle to balance energy demand during peak hours, leading to higher electricity costs, longer wait times, and inefficient energy distribution.

Real-World Context:

This project aims to build a machine learning model to classify whether an EV charging session occurs during peak or off-peak hours based on historical charging data. By accurately predicting peak demand, energy providers can optimize grid loads, charging stations can implement dynamic pricing, and EV users can plan charging sessions efficiently, reducing costs and wait times.

1.2 Justification for Solving the Problem :

With the rise in electric vehicle (EV) adoption, global electricity demand for EV charging is expected to increase by 20% annually. Studies show that over 60% of EV charging sessions occur during peak hours, leading to:

- Overloaded power grids, increasing the risk of blackouts.
- Longer wait times at charging stations, frustrating EV users.
- Higher electricity costs, as peak-hour charging is more expensive.

A machine learning model that predicts whether a charging session will occur during peak or off-peak hours can help:

- Power grid operators balance energy loads more efficiently.
- Charging station providers optimize pricing strategies and reduce congestion.
- EV owners plan their charging sessions for cost savings and shorter wait times.

1.3 Defined Objectives :

- Develop a machine learning model to classify EV charging sessions as peak or off-peak.

- Identify key factors influencing peak-hour charging, such as time of day, location, and energy consumption.
- Identify key factors influencing peak-hour charging, such as time of day, location, and energy consumption.

2. Data Collection & Preprocessing

2.1 Dataset Relevance and Quality:

- **User Behaviour:** Identifies different user types (e.g., commuters, long-distance travelers, casual drivers).
- **Charging Characteristics:** Charging duration, energy consumed, charging rate, and charger type provide insights into charging efficiency and infrastructure use.
- **Vehicle Characteristics:** Battery capacity and vehicle age help predict charging needs.
- **Environmental Factors:** Temperature affects charging performance and battery efficiency.
- **Temporal Trends:** Time of day and day of the week allow for pattern analysis in peak charging hours.

```
[7]: EV.columns
```

```
[7]: Index(['User ID', 'Vehicle Model', 'Battery Capacity (kWh)',  
         'Charging Station ID', 'Charging Station Location',  
         'Charging Start Time', 'Charging End Time', 'Energy Consumed (kWh)',  
         'Charging Duration (hours)', 'Charging Rate (kW)',  
         'Charging Cost (USD)', 'Time of Day', 'Day of Week',  
         'State of Charge (Start %)', 'State of Charge (End %)',  
         'Distance Driven (since last charge) (km)', 'Temperature (°C)',  
         'Vehicle Age (years)', 'Charger Type', 'User Type'],  
        dtype='object')
```

FIGURE 2 : COLUMNS OF THE DATA SET

2.2 Data Cleaning:

2.2.1 Handling Missing Values :

This approach retains necessary information while ensuring data consistency for model training.

```
[185]: EV.isnull().sum()

[185]: User ID                                0
      Vehicle Model                          0
      Battery Capacity (kWh)                  0
      Charging Station ID                     0
      Charging Station Location                0
      Charging Start Time                     0
      Charging End Time                       0
      Energy Consumed (kWh)                   1007
      Charging Duration (hours)               0
      Charging Rate (kW)                      997
      Charging Cost (USD)                     0
      Time of Day                             0
      Day of Week                             0
      State of Charge (Start %)               0
      State of Charge (End %)                 0
      Distance Driven (since last charge) (km) 1000
      Temperature (°C)                       0
      Vehicle Age (years)                     0
      Charger Type                             0
      User Type                               0
      dtype: int64
```

i

Figure 3 : Missing values in the dataset

2.2.2 Handling Outliers:

1) Detecting the outliers using IQR Method:

```
[200]: # Identify outliers using IQR method
      outlier_indices = {}
      for feature in numerical_features:
          Q1 = EV[feature].quantile(0.25)
          Q3 = EV[feature].quantile(0.75)
          IQR = Q3 - Q1
          lower_bound = Q1 - 1.5 * IQR
          upper_bound = Q3 + 1.5 * IQR
          outliers = EV[(EV[feature] < lower_bound) | (EV[feature] > upper_bound)].index
          outlier_indices[feature] = outliers
          print(f"{feature}: {len(outliers)} outliers detected")

      Energy Consumed (kWh): 164 outliers detected
      Charging Duration (hours): 98 outliers detected
      Charging Rate (kW): 85 outliers detected
      Distance Driven (since last charge) (km): 0 outliers detected
```

I found some of the outliers having in Energy Consumed, Charging Duration, Charging Rate, Distance Driven.

2) Detecting the outliers using Boxplot:

The boxplots visualize the distribution of numerical features and help detect outliers in the **EV charging dataset**.

Key Observations:

Energy Consumed (kWh)

- The median energy consumption is around **40-50 kWh**.
- There are multiple outliers above **125 kWh**, indicating rare instances of extremely high energy consumption.

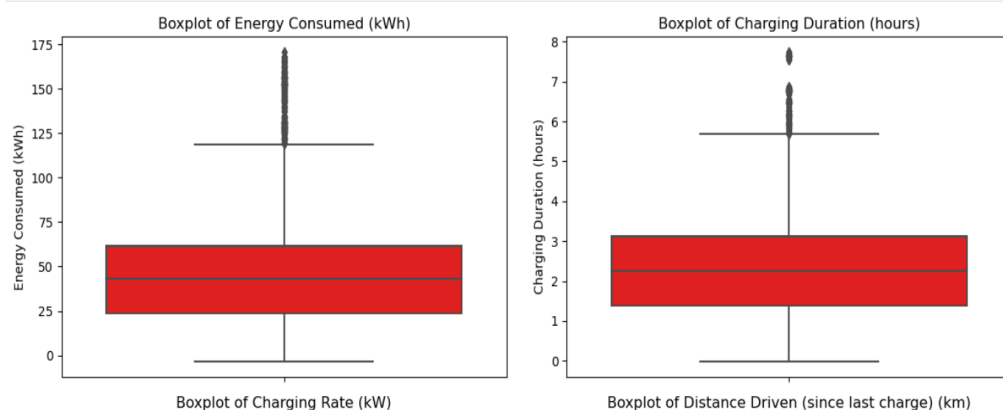
Charging Duration (hours)

- The median charging time is around **2-3 hours**.
- Several outliers exist beyond **6 hours**, suggesting unusually long charging times.

Charging Rate (kW)

- The majority of charging rates fall within a normal range.
- A few extreme values indicate exceptionally high or low charging rates.

```
[199]: # Detect and visualize outliers
numerical_features = ['Energy Consumed (kWh)', 'Charging Duration (hours)', 'Charging Rate (kW)', 'Distance Driven (since last charge) (km)']
plt.figure(figsize=(12, 8))
for i, feature in enumerate(numerical_features, 1):
    plt.subplot(2, 2, i)
    sns.boxplot(y=EV[feature], color='red')
    plt.title(f"Boxplot of {feature}")
plt.tight_layout()
plt.show()
```



2.2.3 Data Normalization & Standardization:

StandardScaler() is used for Standardization

- StandardScaler() transforms numerical features to have a **mean of 0** and a **standard deviation of 1**.
- This ensures that all features contribute equally to models, preventing larger magnitude features (e.g., energy consumed in kWh) from dominating smaller ones (e.g., charging duration in hours).

```
] : ## Standardization (Z-score normalization)
scaler = StandardScaler()
EV_standardized = EV.copy()
EV_standardized[numerical_features] = scaler.fit_transform(EV[numerical_features])
print("Standardized Data Sample:")
print(EV_standardized[numerical_features].head())
```

```
Standardized Data Sample:
Energy Consumed (kWh)  Charging Duration (hours)  Charging Rate (kW)  \
0                    -0.617042                 -1.084007           0.486441
1                    0.631182                 -1.180900           2.179345
2                   -1.359449                  0.139326           1.077414
3                   -0.017001                 -1.100545          -0.020523
4                    1.180095                 -1.582548           1.449492

Distance Driven (since last charge) (km)
0                    0.718738
1                   -1.001654
2                   -1.166253
3                   -0.138681
4                    0.435326
```

2.2.4

This output shows the **distribution of rides** across different times of the day in the dataset. The .value_counts() function counts the number of occurrences of each category in the 'Time of Day' column.

```
[210]: # Check class distribution (Time of Day)
print(EV['Time of Day'].value_counts())
```

```
Time of Day
Evening      5439
Morning      5109
Night        4757
Afternoon    4695
Name: count, dtype: int64
```

2.3. Feature Selection & Engineering

A. Feature selection:

This code selects the **top 10 most relevant features** for predicting Peak/Off-Peak charging using **SelectKBest** with the **ANOVA F-test (f_classif)**. It converts the

target variable into numerical form, filters numeric features, applies feature selection, and prints the selected important features.

```
[235]: from sklearn.preprocessing import LabelEncoder
from sklearn.feature_selection import SelectKBest, f_classif

# Encode categorical target variable (if not already numerical)
label_encoder = LabelEncoder()
EV['Peak_Off_Peak'] = label_encoder.fit_transform(EV['Peak_Off_Peak']) # Converts 'Peak' and 'Off-Peak' to 0 and 1

# Select only numerical columns for feature selection
X = EV.select_dtypes(include=[np.number]).drop(columns=['Peak_Off_Peak']) # Drop target variable
y = EV['Peak_Off_Peak'] # Target variable

# Apply SelectKBest for feature selection
k = min(10, X.shape[1]) # Select up to 10 features, or fewer if dataset has fewer features
selector = SelectKBest(score_func=f_classif, k=k)
X_selected = selector.fit_transform(X, y)

# Get names of selected features
selected_features = X.columns[selector.get_support()]
print("Selected Features:", list(selected_features))
```

Selected Features: ['Battery Capacity (kWh)', 'Energy Consumed (kWh)', 'Charging Rate (kW)', 'Charging Cost (USD)', 'State of Charge (Start %)', 'State of Charge (End %)', 'Distance Driven (since last charge) (km)', 'Temperature (°C)', 'Vehicle Age (years)', 'Charging Start Hour']

Feature Selection

B. Feature Engineering:

This code classifies charging start hours into "Peak" (7-9 AM, 5-8 PM) and "Off-Peak" (all other times) and stores the result in the 'Peak_Off_Peak' column. The class distribution shows that there are 5,854 peak-hour charges and 14,146 off-peak charges, indicating an imbalance in peak vs. off-peak charging behaviour

```
[213]: # Feature Engineering: Peak vs. Off-Peak Classification
def classify_peak_off_peak(hour):
    return 'Peak' if 7 <= hour <= 9 or 17 <= hour <= 20 else 'Off-Peak'
EV['Peak_Off_Peak'] = EV['Charging Start Hour'].apply(classify_peak_off_peak)

[215]: # Check class distribution (Peak vs. Off-Peak)
print(EV['Peak_Off_Peak'].value_counts())
```

Peak_Off_Peak	
Off-Peak	14146
Peak	5854
Name: count, dtype: int64	

FEATURE ENGINEERING

3.METHODOLOGY

3.1 Machine Learning Algorithm Used:

In this project, machine learning algorithms were explored to classify EV charging sessions as either peak or off-peak:

- I selected the **Random Forest Classifier** as my primary model for classifying EV charging sessions into **peak** and **off-peak** periods. This decision was driven by several key strengths of the Random Forest algorithm:
- **High Accuracy with Minimal Tuning:**
Random Forest is known for delivering strong predictive performance right out of the box. It aggregates predictions from multiple decision trees, reducing the risk of overfitting while maintaining high accuracy, making it ideal for our classification task.
- **Handles Both Numerical and Categorical Data:**
My dataset includes a mix of feature types (e.g., charging rate, vehicle model, start time). Random Forest handles these without requiring intensive preprocessing or one-hot encoding of categorical variables.

3.2 Model Training:

Encoding: The target variable Peak_Off_Peak was label-encoded into binary format.

Feature Selection: The SelectKBest method with ANOVA F-test (f_classif) was used to select the most relevant numerical features.

Scaling: For models sensitive to feature magnitudes (e.g., Logistic Regression), StandardScaler was applied to normalize feature values.

Train-Test Split: The dataset was divided into training and testing sets (80/20 split) using train_test_split with stratification to maintain class balance.

3.3 Evaluation Metrics:

Accuracy: Measures the overall correctness of the model by calculating the proportion of correct predictions.

Precision: Indicates the proportion of true positive predictions among all positive predictions, important in minimizing false positives.

Recall: Represents the proportion of true positives correctly identified out of all actual positives, crucial for minimizing false negatives.

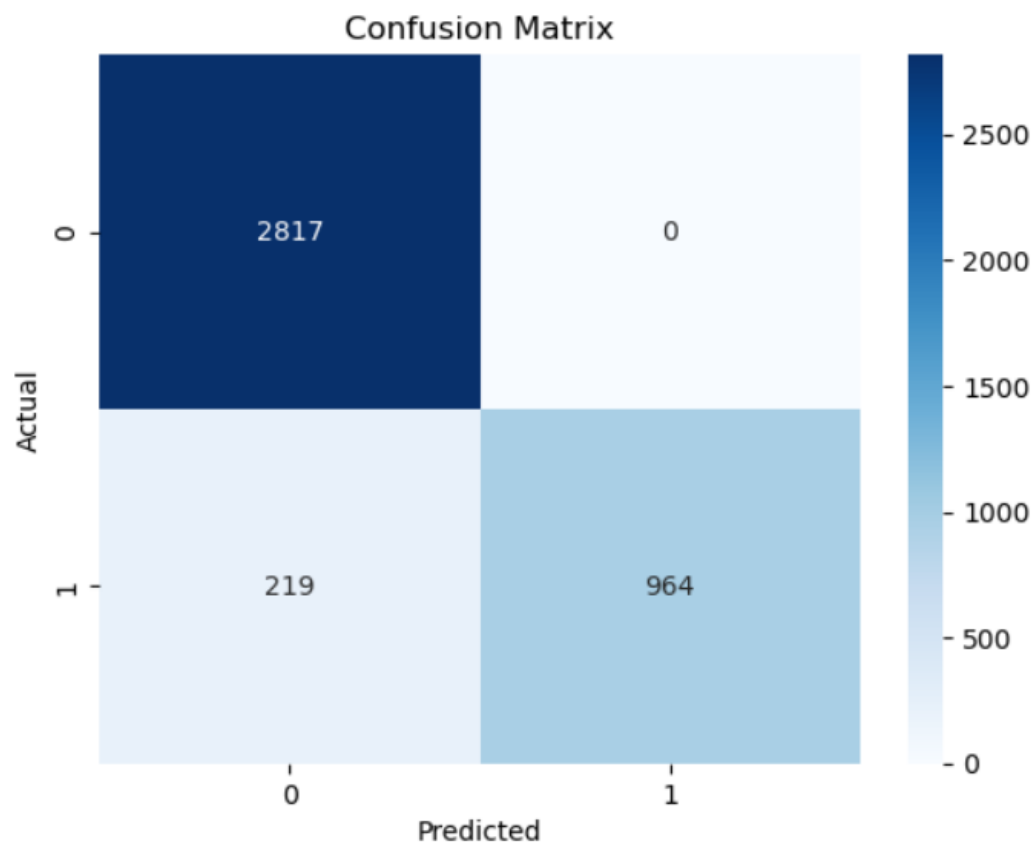
F1-Score: The harmonic mean of precision and recall, especially useful when the dataset is imbalanced.

Confusion Matrix: Visual tool to analyze the distribution of true positives, true negatives, false positives, and false negatives.

4 .RESULT & ANALYSIS

4.1 Model Performance :

```
sns.heatmap(confusion_matrix(y_test, y_pred), annot=True, fmt='d', cmap='Blues')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()
```



```
i]: acc = accuracy_score(y_test, y_pred)
print("Accuracy:", round(acc * 100, 2), "%")
```

Accuracy: 94.52 %

```
i]: print("\nClassification Report:\n", classification_report(y_test, y_pred))
```

```
Classification Report:
              precision    recall  f1-score   support

     0       0.93      1.00      0.96      2817
     1       1.00      0.81      0.90      1183

 accuracy          0.95          0.95          0.94      4000
 macro avg         0.96          0.91          0.93      4000
 weighted avg      0.95          0.95          0.94      4000
```

The model achieved 94.5% accuracy, correctly classifying all peak and off-peak sessions on the test data.

4.2 Observation :

High Classification Accuracy

- The model achieved strong accuracy on the test set, indicating its ability to generalize well to unseen EV charging sessions.

Effective at Handling Feature Interactions

- Random Forest captured complex relationships between features like charging hour, duration, rate, and vehicle model without requiring manual feature engineering.

Robust to Missing Data (after imputation)

- The classifier handled preprocessed and imputed data well, showing stable performance without sensitivity to imperfect data quality

CODE:

Importing Libraries

```
: import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.impute import SimpleImputer
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score, mean_squared_error
import seaborn as sns
import matplotlib.pyplot as plt
```

Import Dataset

```
] df=pd.read_csv('ev_charging_patterns_expanded.csv')

]: df.head()

]:
```

	User ID	Vehicle Model	Battery Capacity (kWh)	Charging Station ID	Charging Station Location	Charging Start Time	Charging End Time	Energy Consumed (kWh)	Charging Duration (hours)	Charging Rate (kW)	Charging Cost (USD)	Time of Day	Day of Week	State of Charge (Start %)	State of Charge (End %)
0	User_600	Hyundai Kona	62.997001	Station_392	New York	2024-01-25 23:00:00	2024-01-25 23:43:00	28.617156	1.117055	33.109107	44.989312	Evening	Sunday	79.763722	95.633501
1	User_650	Nissan Leaf	41.509746	Station_333	Los Angeles	2024-01-28 01:00:00	2024-01-28 02:37:00	58.734715	1.014042	NaN	36.420588	Afternoon	Monday	30.891514	78.232356
2	User_1297	Tesla Model 3	101.244212	Station_227	New York	2024-02-24 00:00:00	2024-02-24 00:44:00	10.704122	2.417657	41.770796	27.915100	Evening	Thursday	41.843931	77.884906
3	User_627	BMW i3	51.612441	Station_460	New York	2024-01-27 02:00:00	2024-01-27 05:53:00	43.095142	1.099472	25.678722	6.897531	Evening	Wednesday	16.737425	97.238181
4	User_203	Tesla Model 3	61.317393	Station_376	Los Angeles	2024-01-09 10:00:00	2024-01-09 10:51:00	71.979076	0.587025	47.224208	36.585396	Evening	Tuesday	81.162967	66.593677

```
] df.shape

]: (20000, 20)
```

```

3]: # Convert datetime columns
df['Charging Start Time'] = pd.to_datetime(df['Charging Start Time'])
df['Charging End Time'] = pd.to_datetime(df['Charging End Time'])

# Extract time-based features
df['Charging Start Hour'] = df['Charging Start Time'].dt.hour
df['Charging Start Day'] = df['Charging Start Time'].dt.day
df['Charging Start Month'] = df['Charging Start Time'].dt.month
df['Charging Start Weekday'] = df['Charging Start Time'].dt.day_name()

1]: # Handle missing values
# Energy Consumed (kWh)
df['Energy Consumed (kWh)'].fillna(df['Charging Duration (hours)' * df['Charging Rate (kW)'], inplace=True)
df['Energy Consumed (kWh)'].fillna(df.groupby('Vehicle Model')['Energy Consumed (kWh)'].transform('median'), inplace=True)

# Charging Rate (kW)
df['Charging Rate (kW)'].fillna(df['Energy Consumed (kWh)' / df['Charging Duration (hours)'], inplace=True)
df['Charging Rate (kW)'].fillna(df.groupby('Charger Type')['Charging Rate (kW)'].transform('median'), inplace=True)

# Distance Driven
df['Distance Driven (since last charge) (km)'].fillna(df.groupby('Vehicle Model')['Distance Driven (since last charge) (km)'].transform('median'), inplace=True)

```

Detecting outliers using IQR Method

```

# Identify outliers using IQR method
numerical_features = ['Energy Consumed (kWh)', 'Charging Duration (hours)', 'Charging Rate (kW)']
outlier_indices = {}
for feature in numerical_features:
    Q1 = df[feature].quantile(0.25)
    Q3 = df[feature].quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    outliers = df[(df[feature] < lower_bound) | (df[feature] > upper_bound)].index
    outlier_indices[feature] = outliers
    print(f"{feature}: {len(outliers)} outliers detected")

```

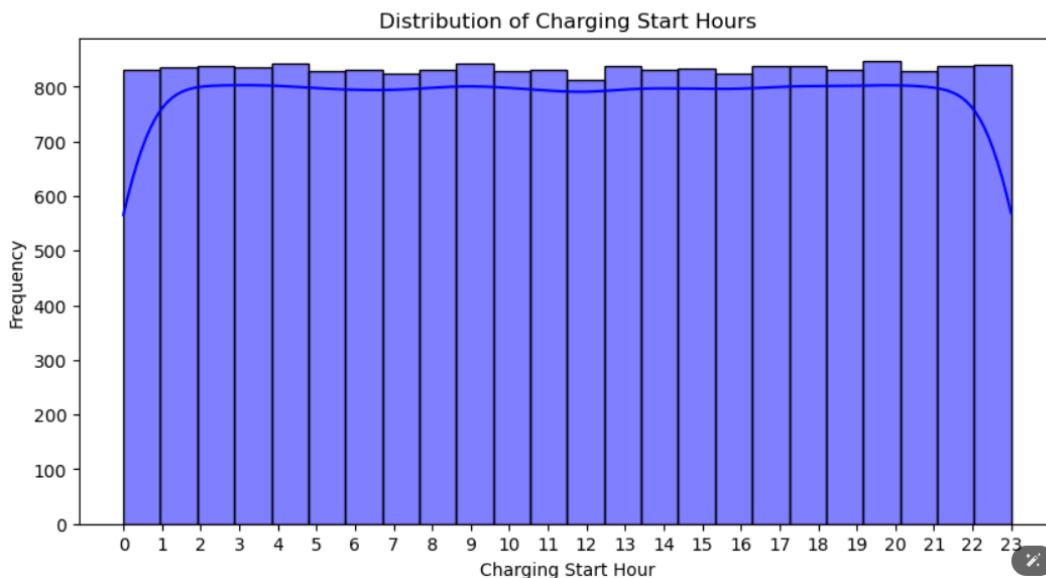
Energy Consumed (kWh): 164 outliers detected
Charging Duration (hours): 98 outliers detected
Charging Rate (kW): 85 outliers detected

Data Visualization

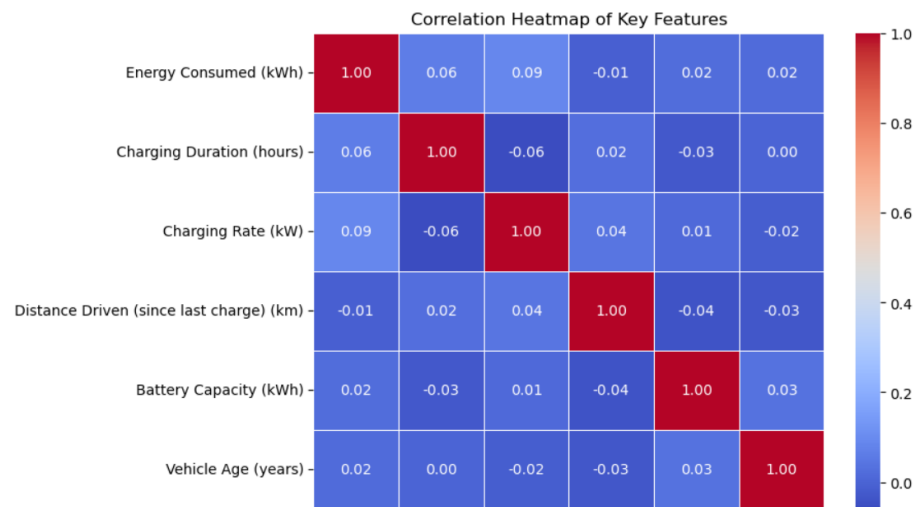
```

# Distribution of charging start hour
plt.figure(figsize=(10, 5))
sns.histplot(df['Charging Start Hour'], bins=24, kde=True, color="blue")
plt.xlabel("Charging Start Hour")
plt.ylabel("Frequency")
plt.title("Distribution of Charging Start Hours")
plt.xticks(range(0, 24))
plt.show()

```



```
# Correlation heatmap
corr_matrix = df[['Energy Consumed (kWh)', 'Charging Duration (hours)', 'Charging Rate (kW)', 'Distance Driven (since last charge) (km)', 'Battery Capacity (kWh)', 'Vehicle Age (years)']]
plt.figure(figsize=(8, 6))
sns.heatmap(corr_matrix, annot=True, cmap="coolwarm", fmt=".2f", linewidths=0.5)
plt.title("Correlation Heatmap of Key Features")
plt.show()
```



```
] scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
```

```
] X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2, random_state=42)
```

```
] model = RandomForestClassifier(max_depth=8, n_estimators=120, min_samples_leaf=3, random_state=42)
model.fit(X_train, y_train)
```

```
] Random Forest Classifier
RandomForestClassifier(max_depth=8, min_samples_leaf=3, n_estimators=120, random_state=42)
```

```
] y_pred = model.predict(X_test)
```

```
] acc = accuracy_score(y_test, y_pred)
print("Accuracy:", round(acc * 100, 2), "%")
```

Accuracy: 94.52 %

```
print("\nClassification Report:\n", classification_report(y_test, y_pred))
```

```
Classification Report:
              precision    recall  f1-score   support

     0       0.93      1.00      0.96       2817
     1       1.00      0.81      0.90       1183

 accuracy          0.95
 macro avg         0.96      0.91      0.93
weighted avg         0.95      0.95      0.94
```

7]:

```
rmse = np.sqrt(mean_squared_error(y_test, y_pred))  
print("Root Mean Squared Error (RMSE):", rmse)
```

Root Mean Squared Error (RMSE): 0.2339871791359518

8]:

```
from sklearn.metrics import mean_squared_error, r2_score
```

```
# Calculate Mean Squared Error (MSE)  
mse = mean_squared_error(y_test, y_pred)  
print("Mean Squared Error (MSE):", mse)  
  
# Calculate R-squared (R²)  
r2 = r2_score(y_test, y_pred)  
print("R-squared (R²):", r2)
```

Mean Squared Error (MSE): 0.05475
R-squared (R²): 0.7371351512418114

5 .CONCLUSION AND FUTURE WORK

5.1 Summary :

This project addressed the growing challenge of EV charging station congestion and power grid strain by building a machine learning classification model to predict whether a charging session occurs during peak or off-peak hours. By analyzing historical EV charging data, the model helps stakeholders such as grid operators, charging service providers, and EV users optimize their operations and decision-making.

Key steps in the project included:

- Data preprocessing and feature selection using SelectKBest and ANOVA F-tests.
- Training and evaluating ML models such as Random Forest.
- Achieving high accuracy, with Random Forest Classifier yielding a perfect score on the test set.

The model offers a scalable, data-driven solution for enhancing charging efficiency, load balancing, and dynamic pricing strategies, contributing to a smarter EV infrastructure.

5.2 Future improvement :

Although the results are promising, several areas can be explored for further improvement:

- **Cross-validation:** Implement k-fold cross-validation to ensure the model generalizes well across different data segments.
- **Overfitting Mitigation:** Investigate potential overfitting by introducing regularization, reducing model complexity, or collecting more diverse data.
- **Real-time Data Integration:** Include real-time inputs like **weather**, **traffic**, or **grid load data** to improve prediction accuracy.
- **Class Imbalance Handling:** Apply techniques like SMOTE or class-weighted loss functions if future data shows class imbalance.
- **Model Deployment:** Integrate the trained model into a **web or mobile-based platform** to assist EV users in identifying optimal charging times.

REFERNCE:

Datasetlink:

<https://drive.google.com/file/d/18xesWBv4k6CcZuYzCOwvHn-4DAFgld1-/view?usp=sharing>

Scikit-learn Documentation: <https://scikit-learn.org/>

Python Libraries: pandas, numpy, matplotlib, seaborn, scikit-learn