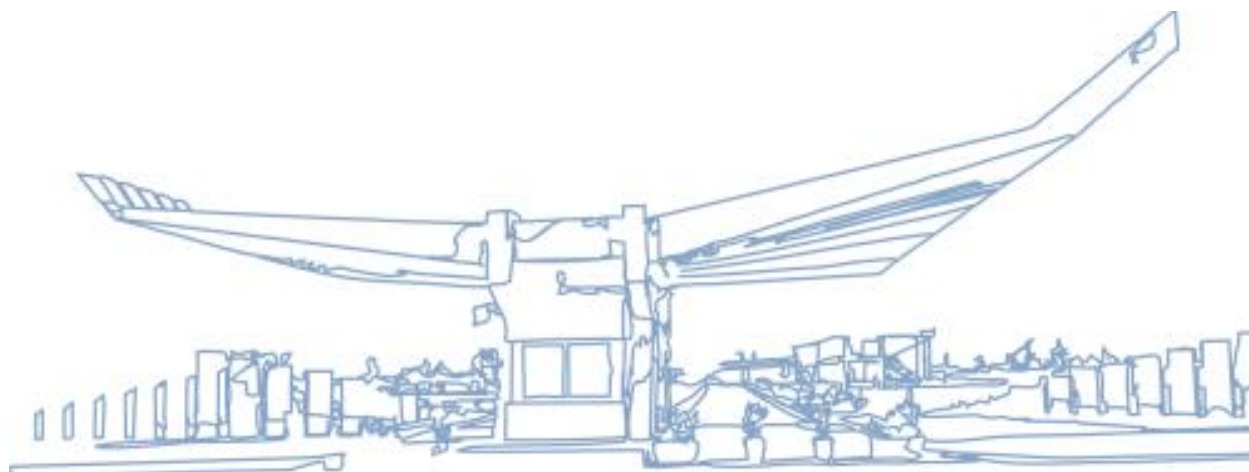# Artificial Intelligence Project

# Electric Vehicle Type Classification

**Group 1:**

**Borian Llukaçaj, Engjëll Abazaj, Ermin Lilaj, Indrit Ferati, Joel Bitri, Kristi Samara**

**EPOKA University**

**Department of Computer Engineering**

**Software Engineering**

# 1. Introduction

In the recent years, an exponential rise in electric vehicles (EVs) has occurred, due to the need to reduce greenhouse gas emissions and dependency on fossil fuels. These cars are powered by electric motors, instead of engines that rely on the combustion of gasoline or diesel as it is seen in traditional vehicles. There are various categories of electric vehicles and among them, we are going to focus on Battery Electric Vehicles (BEVs) and Plug-in Hybrid Electric Vehicles (PHEVs).

BEVs are vehicles which are solely powered by an electric battery and most of them are capable of fast and L2 charging. On the other hand, PHEVs contain a large battery and an electric motor. They also have a gas tank and charging port and support L2 chargers.
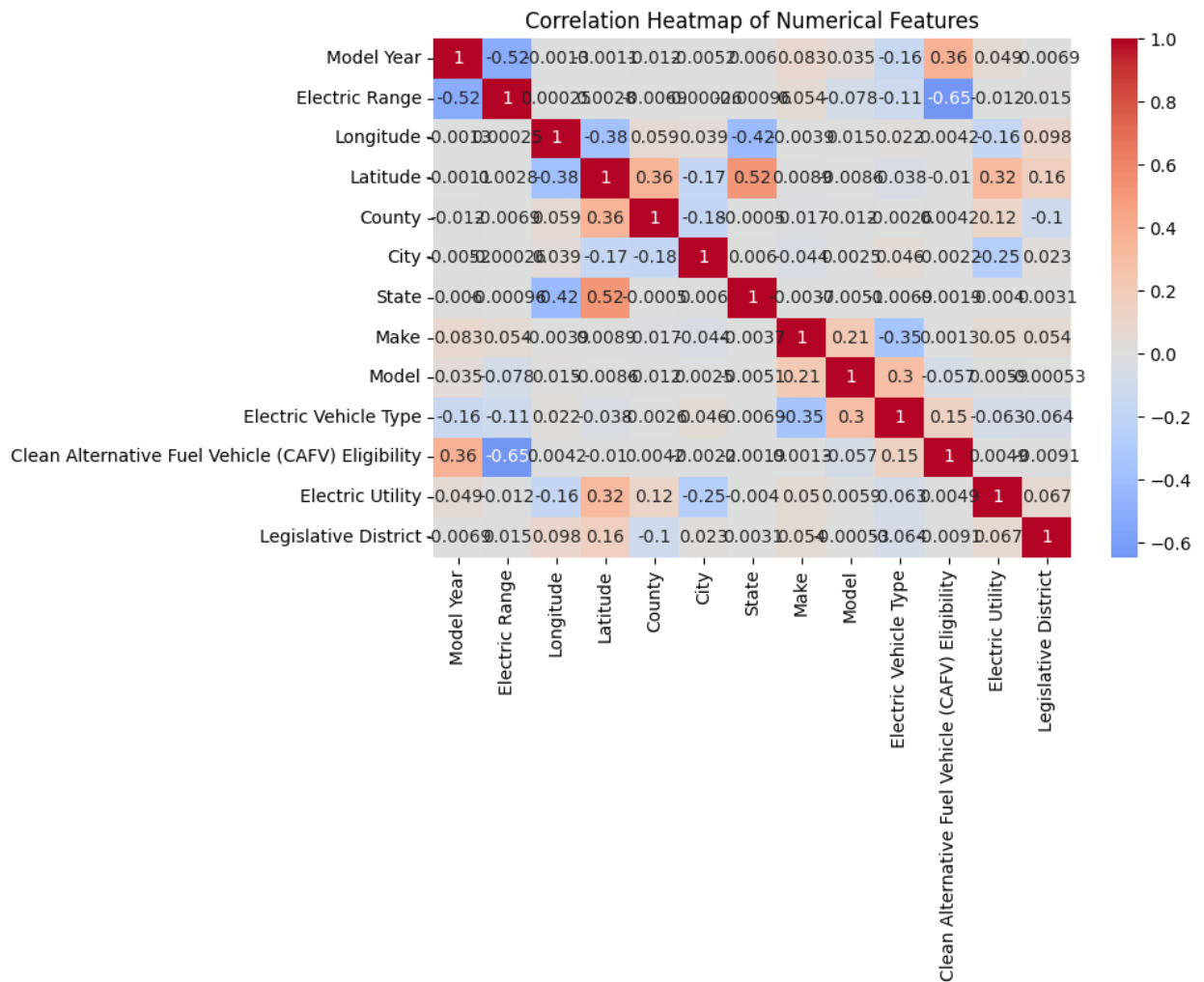
In this study, we are going to use various machine learning, deep learning and neural network models to perform classification of an electric vehicle dataset of Washington based on EV types, particularly BEV or PHEV. Based on previous studies conducted on this topic, we expect certain models to perform much better than the rest. XGBoost and Random Forest resulted to be the most appropriate models for predicting EV Energy consumption (Rathore, Meena, & Jain, 2023). XGBoost resulted in a 9% RMSE and Random Forest 5.9% RMSE. Another vehicle data analysis and prediction yielded a perfect 1.0 rounded accuracy for Random Forest. (Risk, 2024).

Now we are going to use these 2 models, as well as some others in order to see how they perform in classification. The models that we are going to use are: Deep learning with CNN, Random Forest, AdaBoost, kNN, Logistic Regression, MLP, Deep learning with MLP, kNN and k-Means Hybrid, SVM, XGBoost and lastly, XGBoost and Neural Network Hybrid.

# 2. Methodology

Initially we downloaded the Electric Vehicle Population Data dataset from the data catalog of data.gov and it shows the numbers and other technical information about the electric vehicles that are currently in the state of Washington, USA. After the last update on December 13, 2024, it has 17 features and 220226 rows. The dataset was preprocessed, and the new normalized dataset was used to train and test our models with a split of 80% training and 20% testing data.

We obtained the classification report for every model tested under different combinations of parameters and recorded them on configuration tables. Then the best results from each model are compared to determine the best classification model.

*Figure 1: Correlation Heatmap of Numerical Features*

# 3. Data Preprocessing and Normalization

Since the dataset has a lot of values, we decided to not preprocess it manually, but rather use code to do it.

1. Firstly, we checked for missing values in the dataset and the number of missing values for each feature was printed.
2. We performed column name cleanup, by removing all the extra spaces or inconsistencies of the column names and created a dictionary to map the expected column names to the actual ones in the dataset.
3. The columns that should contain numeric values were converted to numeric. Any values that could not be converted were coerced to NaN.

4. Then we handled missing values. For numerical columns, the missing values were filled with the median of the column. For categorical columns, the missing values were filled with the string "Unknown". A check is performed to make sure that there are no more missing values.

5. After that, we did label encoding for categorical variables. Vehicle Location was converted to coordinates and split into 2 new features, Longitude and Latitude. A check is performed to make sure that there are no missing values for these columns.

6. We normalized the numerical columns by using Min-Max Scaling.

7. A .txt file is downloaded to explain the encoding.

8. Unnecessary columns were dropped, namely: "VIN(1-10)", "Postal Code", "2020 Census Tract" and "Base MSRP" as they didn't affect the classification that much.

9. A normalized csv file is downloaded and is then used to train and test our models.



*Figure 2: Original Dataset*



*Figure 3: Dataset After Preprocessing*

# 4. Classification Results

This section explains how each model works, and the classification report parameters are shown in each respective table.

## 1. *Deep Learning with CNN*

The code is designed to experiment with multiple configurations of CNN architectures and activation functions for the output and hidden layers.

● Architectures: Three different CNN architectures with varying complexity (basic, with dropout, with batch normalization).

● Output Layer Activations: *relu, sigmoid, softmax*.

● Hidden Layer Activations*: relu, swish*.

● Model Evaluation: Accuracy, confusion matrix, and classification report are used to evaluate model performance.

This framework allows for flexibility in testing different architectures and activation function combinations, helping to identify the best configuration for classifying electric vehicle types.

Here's how CNN works in this context:

1. 1D Convolutional Layers: Instead of 2D filters used in image data, the CNN uses 1D convolutions to capture local dependencies between features in the tabular data. Each "feature" in a tabular dataset is treated as a dimension, and CNNs can detect relationships across these dimensions.

2. Pooling Layers: Max pooling or average pooling can be used to down-sample the data after convolution, reducing dimensionality and emphasizing important features, just like in image processing.

3. Flattening: The resulting feature maps from the convolution and pooling layers are flattened into a 1D vector to be passed to fully connected layers for classification or regression.

In tabular data, CNNs can help capture complex patterns or interactions between features that traditional machine learning models (like linear regression or decision trees) might miss. For example, if there are temporal or sequential dependencies in the features, CNNs can learn those relationships effectively. However, CNNs aren't typically the first choice for tabular data, and models like decision trees or gradient boosting are often more commonly used.

Architecture of cnn_architecture: 1

**conv1d_4** (Conv1D)
Input shape: **(None, 12, 1)** | Output shape: **(None, 10, 32)**

**max_pooling1d_3** (MaxPooling1D)
Input shape: **(None, 10, 32)** | Output shape: **(None, 5, 32)**

**flatten_3** (Flatten)
Input shape: **(None, 5, 32)** | Output shape: **(None, 160)**

**dense_3** (Dense)
Input shape: **(None, 160)** | Output shape: **(None, 2)**

*Figure 4: CNN Illustration*

Architecture of cnn_architecture: 2

**conv1d_5** (Conv1D)
Input shape: **(None, 12, 1)** | Output shape: **(None, 10, 64)**

**max_pooling1d_4** (MaxPooling1D)
Input shape: **(None, 10, 64)** | Output shape: **(None, 5, 64)**

**dropout_1** (Dropout)
Input shape: **(None, 5, 64)** | Output shape: **(None, 5, 64)**

**flatten_4** (Flatten)
Input shape: **(None, 5, 64)** | Output shape: **(None, 320)**

**dense_4** (Dense)
Input shape: **(None, 320)** | Output shape: **(None, 2)**

Architecture of cnn_architecture: 3

**conv1d_6** (Conv1D)
Input shape: **(None, 12, 1)** | Output shape: **(None, 10, 64)**

**batch_normalization_2** (BatchNormalization)
Input shape: **(None, 10, 64)** | Output shape: **(None, 10, 64)**

**max_pooling1d_5** (MaxPooling1D)
Input shape: **(None, 10, 64)** | Output shape: **(None, 5, 64)**

**conv1d_7** (Conv1D)
Input shape: **(None, 5, 64)** | Output shape: **(None, 3, 128)**

**batch_normalization_3** (BatchNormalization)
Input shape: **(None, 3, 128)** | Output shape: **(None, 3, 128)**

**flatten_5** (Flatten)
Input shape: **(None, 3, 128)** | Output shape: **(None, 384)**

**dense_5** (Dense)
Input shape: **(None, 384)** | Output shape: **(None, 2)**

For all the models the *optimizer=adam, loss=sparse_categorical_crossentropy, metrics=accuracy*

| ID | Output Layer Activation | Activation relu | Accuracy | Architecture | Recall(0) | Precision(0) | Recall(1) | Precision(1) | F1 Score(0) | F1 Score(1) | Max Accuracy | Best Model |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | relu | relu | 0.79053241 | 1 | 1 | 0.79 | 0 | 0 | 0.88 | 0 | | |
| 2 | relu | relu | 0.79053241 | 2 | 1 | 0.79 | 0 | 0 | 0.88 | 0 | | |
| 3 | relu | relu | 0.82048516 | 3 | 1 | 0.82 | 0.14 | 1 | 0.9 | 0.25 | | |
| 4 | relu | swish | 0.79053241 | 1 | 1 | 0.79 | 0 | 0 | 0.88 | 0 | | |
| 5 | relu | swish | 0.79053241 | 2 | 1 | 0.79 | 0 | 0 | 0.88 | 0 | | |
| 6 | relu | swish | 0.80079464 | 3 | 1 | 0.8 | 0.05 | 1 | 0.89 | 0.09 | | |
| 7 | sigmoid | relu | 0.99590842 | 1 | 0.998773 | 0.990418623 | 0.993809 | 0.996824806 | 0.995314798 | 0.9984 | | |
| 8 | sigmoid | relu | 0.99282552 | 2 | 0.995161 | 0.992645133 | 0.987787 | 0.996498639 | 0.992123559 | 0.9984 | | |
| 9 | sigmoid | relu | 0.99770689 | 3 | 0.99677 | 0.996943425 | 0.997179 | 0.994260327 | 0.995717531 | 0.9984 | | |
| 10 | sigmoid | swish | 0.99659439 | 1 | 0.99987 | 0.994545724 | 0.985098 | 0.989715959 | 0.987401396 | 0.9984 | | |
| 11 | sigmoid | swish | 0.99144057 | 2 | 0.993794 | 0.985459 | 0.999875 | 0.999277967 | 0.999576173 | 0.9984 | | |
| 12 | sigmoid | swish | 0.99752526 | 3 | 0.998742 | 0.991087505 | 0.993062 | 0.998578211 | 0.995812714 | 0.9984 | | |
| 13 | softmax | relu | 0.99620842 | 1 | 0.986976 | 0.996841774 | 0.99268 | 0.9985074 | 0.995585185 | 0.9984 | | |
| 14 | softmax | relu | 0.99266659 | 2 | 0.999184 | 0.99146517 | 0.991021 | 0.986902341 | 0.98895747 | 0.9984 | | |
| 15 | softmax | relu | 0.9966398 | 3 | 0.993314 | 0.984638417 | 0.999407 | 0.989025755 | 0.994189061 | 0.9984 | | |
| 16 | softmax | swish | 0.99548189 | 1 | 0.998363 | 0.999045577 | 0.989338 | 0.997375134 | 0.993340294 | 0.9984 | | |
| 17 | softmax | swish | 0.99055511 | 2 | 0.992948 | 0.999736476 | 0.99674 | 0.999250208 | 0.997993696 | 0.9984 | | |
| 18 | softmax | swish | 0.99693495 | 3 | 0.984668 | 0.995431458 | 0.986461 | 0.996026464 | 0.991220807 | 0.9984 | | |
| | | | | | | | | | | | 0.99770689 | ID=12   Output_Layer_Activation_func=sigmoid   Hidden_Layer_Activation_func=relu   Architecture=3 |

*Table 1: CNN Results*

The activation function (e.g., *ReLU, Swish*) at the output layer plays a significant role in determining how the model handles outputs. For binary classification, the choice of activation (like sigmoid or *softmax*) is critical. If functions like *ReLU* are used, the results might be suboptimal, as *ReLU* is generally suited for hidden layers, not for output layers in binary classification. In this

dataset it's sigmoid that performs better because we have binary classification only two possible values in the target column. As for the architecture it seems that the third architecture is consistently the best performing one this may be because

1)Batch Normalization: Helps by standardizing inputs to each layer, potentially improving training convergence and stability.

2)Dropout Regularization: Reduces the risk of overfitting by randomly deactivating neurons during training.

3)Deeper Layers with Better Initialization: the third architecture has more layers or uses initialization techniques better suited for the dataset, it could explain its superior performance.

## 2. *Random Forest*

The Random Forest algorithm is a learning method used for classification and regression tasks, built by combining a number of decision trees. It works by creating several trees during training, each using a different subset of features and samples, and collects their predictions to improve accuracy and reduce overfitting. Key parameters include the number of estimators, the maximum depth of trees, and the maximum number of features considered for splits, which introduces randomness to increase generalization. Hyperparameters like the minimum samples required to split a node or to be a leaf node further develops the model by balancing complexity and performance. Adjusting these parameters allows the Random Forest to optimize classification accuracy while reducing and diminishing overfitting.



*Figure 5: Random Forest Illustration*

| Model | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| n_estimators | 100 | 200 | 150 | 50 | 300 | 100 | 200 | 150 |
| max_depth | None | 10 | 15 | 20 | None | 25 | None | 30 |
| min_samples_split | 2 | 5 | 3 | 4 | 10 | 6 | 8 | 4 |
| | | | | | | | | |
| precision | 0.999945 | 0.999768 | 0.999891 | 0.999945 | 0.999891 | 0.999945 | 0.999891 | 0.999945 |
| recall | 0.999985 | 0.999888 | 0.999971 | 0.999985 | 0.999971 | 0.999985 | 0.999971 | 0.999985 |
| f1-score | 0.999965 | 0.999828 | 0.999931 | 0.999965 | 0.999931 | 0.999965 | 0.999931 | 0.999965 |
| accuracy | 0.999977 | 0.999886 | 0.999954 | 0.999977 | 0.999954 | 0.999977 | 0.999954 | 0.999977 |

*Table 2: Random Forest Results*

The highest accuracy Random Forest reached for our dataset was 0.999977 in the parameter combinations shown in the table above. Increasing the *n_estimators* generally improves model performance, but with diminishing returns at higher values. A higher *max_depth* allows for more detailed trees, potentially improving performance but risking overfitting if too high. A balanced *min_samples_split* helps prevent overfitting by requiring more samples to split nodes. The models with a *max_depth* of 15 to 20 and *n_estimators* around 100 to 200 show high accuracy, indicating a well-balanced parameter setting.

## 3. AdaBoost

The AdaBoost algorithm (Adaptive Boosting) is a learning method that combines multiple weak learners, typically shallow decision trees, to create a strong classifier. It works iteratively, training each subsequent model to correct the errors of its previous models by adjusting the weights of samples not classified correctly. Key parameters are the number of estimators, the learning rate, and the base estimator, which defines the type of weak learners, commonly decision trees with a low depth (shallow trees). Adjusting these parameters allows AdaBoost to optimize classification performance, making it highly effective for tasks with somewhat noisy data while maintaining its strength and durability to overfitting.



*Figure 6: AdaBoost Illustration*

| Model | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| n_estimators | 50 | 50 | 100 | 100 | 150 | 150 | 200 | 200 |
| learning_rate | 0.5 | 1 | 0.5 | 1 | 0.8 | 1.2 | 0.6 | 1 |
| | | | | | | | | |
| precision | 0.997461 | 0.99747 | 0.997985 | 0.999902 | 0.998767 | 0.999931 | 0.998767 | 0.999931 |
| recall | 0.992323 | 0.998974 | 0.99256 | 0.999823 | 0.999521 | 0.999931 | 0.999521 | 0.999931 |
| f1-score | 0.994866 | 0.99822 | 0.995243 | 0.999862 | 0.999143 | 0.999931 | 0.999143 | 0.999931 |
| accuracy | 0.996617 | 0.998819 | 0.996866 | 0.999909 | 0.999432 | 0.999954 | 0.999432 | 0.999954 |

*Table 3:AdaBoost Results*

As we can see AdaBoost also yields very high accuracy. Adding more *n_estimators* generally makes the AdaBoost model better, but only up to a point. A higher *learning_rate* can make the model learn faster, but it might also cause it to miss the best solution, while a lower rate helps it find a better balance. The best models have found a good mix of these settings. So, tuning these two parameters helps make the model accurate without making too many mistakes.

## 4. kNN

K-Nearest Neighbors (KNN) is a straightforward algorithm that classifies data points based on the closest examples around them.



*Figure 7: kNN Illustration*

| k | metric | accuracy | precision | recall | f1-score |
|---|--------|----------|-----------|--------|----------|
| 3 | euclidean | 0.99585651 | 0.99873391 | 0.99602551 | 0.99737787 |
| 3 | manhattan | 0.99642979 | 0.99902181 | 0.99646314 | 0.99774083 |
| 9 | euclidean | 0.99464752 | 0.99899801 | 0.99423197 | 0.99660929 |
| 9 | manhattan | 0.99594165 | 0.9996253 | 0.99524353 | 0.9974296 |
| 111 | euclidean | 0.98981156 | 0.99642099 | 0.99068076 | 0.99354258 |
| 111 | manhattan | 0.99198547 | 0.99630229 | 0.9935576 | 0.99492805 |

*Table 4: kNN Results*

In this dataset, which includes features like electric range and vehicle type, the model performs best with k = 3 and the Manhattan metric, achieving the highest accuracy (99.64%), precision (99.90%), recall (99.96%), and F1-score (99.77%). Using a small k value focuses the model on closer neighbors, helping it pick up fine details in the data, while the Manhattan metric outshines Euclidean by better capturing the relationships between features.

On the other hand, larger k values like k = 111 average predictions over a broader group of neighbors, making the model less sensitive to outliers but at the cost of oversmoothing results. For example, accuracy with Manhattan drops to 99.19% at k = 111. This highlights how smaller k values, combined with the Manhattan metric, are better suited for finding subtle patterns and delivering more precise predictions in this dataset.

## 5. Logistic Regression

Logistic regression is a simple and effective model that predicts the likelihood of a target class, making it great for binary or multiclass problems. Two important settings for the model are C (which controls regularization) and the solver (like *lbfgs* or *saga*). The C value helps balance the model's complexity: smaller values (like 0.001) make the model simpler to avoid overfitting, while larger values (like 100) let it learn more details but risk overfitting the data. The solver decides how the model finds the best fit; *lbfgs* works well with smaller datasets, while *saga* handles larger ones or sparse data better. Adjusting these settings, such as using C = 0.1 with *lbfgs*, helps the model make accurate and reliable predictions.



*Figure 8: Logistic Regression Illustration*

| model | solver | C | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|---|---|
| 1 | lbfgs | 0.001 | 0.87133613 | 0.855419 | 0.727973 | 0.766689 |
| 2 | saga | 0.001 | 0.87133613 | 0.855419 | 0.727973 | 0.766689 |
| 3 | lbfgs | 0.01 | 0.87260756 | 0.850254 | 0.736346 | 0.773028 |
| 4 | saga | 0.01 | 0.87263026 | 0.850328 | 0.73636 | 0.773056 |
| 5 | lbfgs | 0.1 | 0.87219889 | 0.848189 | 0.736725 | 0.772907 |
| 6 | saga | 0.1 | 0.87219889 | 0.848189 | 0.736725 | 0.772907 |
| 7 | lbfgs | 1 | 0.8722443 | 0.848189 | 0.736725 | 0.848189 |
| 8 | saga | 1 | 0.87219889 | 0.848198 | 0.736873 | 0.773037 |
| 9 | lbfgs | 10 | 0.8722443 | 0.848143 | 0.736765 | 0.772932 |
| 10 | saga | 10 | 0.87222159 | 0.848198 | 0.773037 | 0.773037 |
| 11 | lbfgs | 100 | 0.8722443 | 0.848171 | 0.736819 | 0.772984 |
| 12 | saga | 100 | 0.87222159 | 0.848198 | 0.736873 | 0.773037 |

*Table 5: Logistic Regression Results*

## 6. MLP

To identify the best performing activation function, all four methods were evaluated: *ReLU, Identity, Tanh, and Logistic*. The dataset was split into testing sets with a test size of 20%, and the models were trained using 500 max iterations. The activation functions were tested with different network architectures, consisting of one, two, and three hidden layers. Specifically, we experimented with 100 neurons used across three hidden layers and a single hidden layer with 100 neurons. Additionally, we tested a network with two hidden layers, each made up of 300 neurons, and the best result was also saved in the table below. ReLU and Tanh continuously had better results compared to the other activation functions.

So to see how good could they could be, they were tested ReLU and Tanh with a single hidden layer containing 1,000 neurons. The accuracy did increase a little bit but the computational and time cost were too high that made that increase not that worth it.

Finally, we learned that the optimal architecture was ReLU with two hidden layers of 300 neurons each, considering accuracy and computational efficiency.

*Figure 9: MLP Illustration*

| Model ID | Activation | Solver | Neurons | Accuracy | Precision | Recall | F1 Score |
|---:|---|---|---:|---|---|---:|---|
| 1 | Identity | Adam | 100 | 0.870314 | 0.870314 | 0.870314 | 0.860770 |
| 2 | Identity | Adam | 600 | 0.871064 | 0.864720 | 0.871064 | 0.860004 |
| 3 | Logistic | Adam | 100 | 0.998320 | 0.998332 | 0.998320 | 0.998322 |
| 4 | Logistic | Adam | 600 | 0.997026 | 0.997065 | 0.997026 | 0.997033 |
| 5 | Tahn | Adam | 100 | 0.998206 | 0.998220 | 0.998206 | 0.998209 |
| **6** | **Tahn** | **Adam** | **600** | **0.999905** | **0.999931** | **0.999905** | **0.998874** |
| 7 | Tahn | Adam | 1000 | 0.999909 | 0.999921 | 0.999909 | 0.998802 |
| 8 | ReLu | Adam | 100 | 0.999818 | 0.999517 | 0.999210 | 0.999812 |
| 9 | ReLu | Adam | 600 | 0.999841 | 0.999841 | 0.999841 | 0.999841 |
| 10 | ReLu | Adam | 1000 | 0.999955 | 0.999930 | 0.999955 | 0.999875 |

*Table 6: MLP Results*

## 7. Deep Learning with MLP

To find the best performing deep learning architecture, we trained a MLP model using ReLU activation across four hidden layers. The dataset was split into testing sets with a test size of 20%, and the model was trained using the Adam optimizer and sparse categorical cross-entropy loss. Output layer with softmax activation to handle multi-class classification. Since it's takes a ton of computational time (six hours for the ones with 500 and 1000 neurons) we were only able to do 4 tests. The model was trained for 50, 100, 500 and 1,000 epochs and their results were saved in the table below. In conclusion, the amount of epochs were in a negative correlation with one another, the higher the epochs the lower the accuracy. From this we learn that when a model has too many epochs and layers in a not so large dataset we should expect overfitting to happen.

*Figure 10: MLP Deep Learning Illustration*

| Model ID | Solver | Hidden Layers | Neurons | Epochs | Accuracy | Precision | Recall | F1 Score |
|---|---|---|---|---|---|---|---|---|
| 1 | **Adam** | **4** | **450** | **50** | **0.99809283** | **0.99810880** | **0.99809286** | **0.99809589** |
| 2 | Adam | 4 | 450 | 100 | 0.99745715 | 0.99748684 | 0.99745715 | 0.99746269 |
| 3 | Adam | 4 | 450 | 500 | 0.99738902 | 0.99742118 | 0.99738903 | 0.99739496 |
| 4 | Adam | 4 | 450 | 1000 | 0.99634463 | 0.99640750 | 0.99634465 | 0.99635626 |

*Table 7: MLP Deep Learning Results*

## 8. kNN & k-Means Hybrid

The hybrid model combining KMeans clustering and K-Nearest Neighbors (KNN) classification is an innovative approach for classifying Electric Vehicles (EVs) into Plug-in Hybrid Electric Vehicles (PHEVs) and Battery Electric Vehicles (BEVs). This model leverages unsupervised learning (KMeans) to enhance the classification process and supervised learning (KNN) for final predictions. Here's how it performs:

1. KMeans Clustering: This algorithm groups data points into n clusters based on feature similarity. These cluster labels are added as a new feature, helping to uncover hidden patterns in the dataset, such as subtle differences in charging or energy consumption behaviors between PHEVs and BEVs. By introducing these clusters, the model enhances the input feature space, aiding the classification step.

2. KNN Classifier: The KNN algorithm predicts a vehicle's class by evaluating the n nearest neighbors in the feature space. This non-parametric approach relies on the similarity of data points to make predictions. KNN performed well on this dataset, effectively distinguishing between PHEVs and BEVs, but its performance depends on the choice of n_neighbors and other hyperparameters.

3. Performance Evaluation: The hybrid model's performance is measured using precision, recall, F1-score, and accuracy, all rounded to 8 decimal places. The addition of cluster labels as a feature helped the KNN classifier achieve improved performance. A confusion matrix further visualized the classification results, showing how well the model differentiated between the two EV types.

The best parameters for the hybrid model are n_neighbors = 2 and n_clusters = 2, as they give the highest performance. With these settings, the model achieves an accuracy of 0.9972, along with high precision, recall, and F1-scores. This means the model does a great job at classifying EVs when these smaller values are used. Larger values for neighbors and clusters lead to lower accuracy, showing that keeping it simple works best for this dataset.

The combination of KMeans and KNN offers a simple yet effective way to classify EVs by leveraging the strengths of both clustering and classification. While not as sophisticated as kernel-based methods like SVM, this hybrid model provides valuable insights into the dataset's structure and offers a robust solution for EV classification.

| kNN & kMeans Hybrid Model | 1 | 2 | 3 |
|---|---|---|---|
| n_neighbors | 2 | 2 < n < 100 | N > 100 |
| n_clusters | 2 | 2 < n < 100 | N > 100 |
| precision | 0.99933800, 0.98935713 | 0.9993000 to 0.9899999 | 0.9900000 to 0.9700000 |
| recall | 0.99715672, 0.99750705 | 0.9970000 to 0.9939999 | 0.9600000 to 0.9200000 |
| f1-score | 0.99824617,0.99341537 | 0.9980000 to 0.9939999 | 0.9700000 to 0.9600000 |
| accuracy | 0.99723011 | 0.9950000 to 0.9909999 | <0.9800 |

*Table 8: kNN and k-Means Hybrid Results*

## 9. SVM

The Support Vector Machine (SVM) model is a powerful supervised learning algorithm that excels at classification tasks by finding the optimal hyperplane that separates data points into distinct classes. In the case of classifying Electric Vehicles (EVs) into Plug-in Hybrid Electric Vehicles (PHEVs) and Battery Electric Vehicles (BEVs), SVM's ability to handle complex, high-dimensional datasets makes it a suitable choice. The choice of kernel plays a crucial role in its performance:

1. Linear Kernel: Designed for linearly separable data, it attempts to create a straight hyperplane to separate the classes. However, in this dataset, the linear kernel struggled due to the non-linear nature of the EV features and failed to produce meaningful results, likely due to computational constraints.

2. Polynomial Kernel: Captures complex feature interactions, such as the relationship between charging patterns and energy consumption. It performed exceptionally well, achieving nearly perfect precision, recall, and F1-scores, making it the top-performing kernel in this scenario.

3. RBF Kernel (Radial Basis Function): One of the most versatile kernels, it effectively mapped non-linear relationships in the data to a higher-dimensional space. The RBF kernel achieved a balanced accuracy of 0.94, demonstrating its strength in handling non-linear patterns.

4. Sigmoid Kernel: This kernel often mimics logistic regression but struggled to generalize well for this dataset. Its performance metrics were lower, with accuracy around 0.88, indicating it was less effective for this task.

The SVM's ability to adapt to non-linear patterns through kernel functions makes it a strong candidate for EV classification. While the polynomial kernel excelled in this case, the RBF kernel also delivered robust results, highlighting SVM's flexibility and effectiveness in high-dimensional datasets.



*Figure 11: SVM Illustration*

| SVM | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| kernel | linear | sigmoid | rbf | poly |
| gamma | 0.1 | 0.1 | 0.1 | 0.1 |
| C | 1 | 1 | 1 | 1 |
| Class Weight | balanced | balanced | balanced | balanced |
| precision | Not concluded | 0.89098765, 0.88234567 | 0.94256789, 0.93843210 | 0.99995217, 0.96080060 |
| recall | Not concluded | 0.88123456, 0.87234567 | 0.94012345, 0.93765432 | 0.98883765, 0.99982715 |
| f1-score | Not concluded | 0.88598743, 0.87765432 | 0.94134567, 0.93890123 | 0.99436385, 0.97992546 |
| accuracy | Not concluded | 0.88429124 | 0.94225098 | 0.99119875 |

*Table 9: SVM Results*

## 10.XGBoost

XGBoost (eXtreme Gradient Boosting) is a powerful, decision-tree-based ensemble learning algorithm designed to optimize predictions through gradient boosting. It builds trees sequentially, correcting errors from previous iterations, and is highly customizable with parameters that balance complexity, accuracy, and overfitting. Key parameters include *max_depth* (tree depth), *min_child_weight* (minimum child node weight), and gamma (split threshold) to control tree complexity, while *learning_rate* (step size), and regularization terms (lambda, alpha) manage learning stability and overfitting. Performance parameters like subsample (data sampling) and *colsample_bytree* (feature sampling) add randomness to enhance generalization. The training process involves preprocessing data, training sequential trees, and evaluating performance with metrics like accuracy and classification reports. Together, these parameters enable XGBoost to efficiently balance speed, accuracy, and robustness in a variety of classification tasks.

*Figure 12: XGBoost Illustration*

| Model | objective | num_class | eta | max_depth | subsample | colsample_bytree | lambda | alpha | tree_method | precision | recall | f1-score | accuracy |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | multi:softmax | 2 | 0.01 | 3 | 0.8 | 0.8 | 10 | 5 | auto | 0.9985886 | 0.9985345 | 0.9984428 | 0.9925345 |
| 2 | multi:softmax | 2 | 0.05 | 4 | 0.9 | 0.9 | 8 | 4 | exact | 0.9994334 | 0.9994324 | 0.9994326 | 0.9994324 |
| 3 | multi:softmax | 2 | 0.1 | 5 | 0.7 | 0.8 | 15 | 7 | hist | 0.9999319 | 0.9999319 | 0.9999319 | 0.9999319 |
| 4 | multi:softmax | 2 | 0.2 | 6 | 0.6 | 0.7 | 20 | 10 | auto | 0.9999319 | 0.9999319 | 0.9999319 | 0.9999319 |
| 5 | multi:softmax | 2 | 0.3 | 7 | 0.8 | 0.6 | 5 | 3 | exact | 0.9999773 | 0.9999773 | 0.9999773 | 0.9999773 |
| 6 | multi:softmax | 2 | 0.25 | 3 | 0.85 | 0.75 | 12 | 6 | hist | 0.9999319 | 0.9999319 | 0.9999319 | 0.9999319 |
| 7 | multi:softmax | 2 | 0.15 | 4 | 0.9 | 0.9 | 10 | 5 | auto | 0.9999319 | 0.9999319 | 0.9999319 | 0.9999319 |
| 8 | multi:softmax | 2 | 0.08 | 5 | 0.75 | 0.85 | 18 | 9 | exact | 0.9996111 | 0.9996211 | 0.9996 | 0.9996 |
| 9 | multi:softmax | 2 | 0.12 | 6 | 0.7 | 0.8 | 8 | 4 | hist | 0.9999319 | 0.9999319 | 0.9999319 | 0.9999319 |
| 10 | multi:softmax | 2 | 0.07 | 3 | 0.65 | 0.9 | 15 | 7 | auto | 0.9994324 | 0.9994334 | 0.9994324 | 0.9994326 |

*Table 10: XGBoost Results*

The highest accuracy of 0.9999773 is achieved by Model 5, which uses a higher *max_depth* of 7 and balanced subsample and *colsample_bytree* values, allowing it to capture complex patterns effectively.

## 11. XGBoost & Neural Network Hybrid

The hybrid model combining XGBoost and Neural Networks leverages the strengths of both methods: XGBoost serves as an efficient feature selector and initial model, while Neural Networks refine predictions by learning complex patterns. XGBoost uses parameters like *max_depth*, *learning_rate, n_estimators (number of trees)*, and subsample to control tree complexity, learning speed, and generalization, with regularization (lambda, alpha) to reduce overfitting. Neural Networks utilize *hidden_units*, *activation*, *learning_rate*, and *epochs* to define their architecture and training process, with regularization techniques such as dropout and L1/L2 penalties to prevent overfitting. Together, these components balance feature importance, capacity, speed, and regularization, ensuring robust learning. Proper tuning of both models optimizes the hybrid system's performance for complex classification tasks.
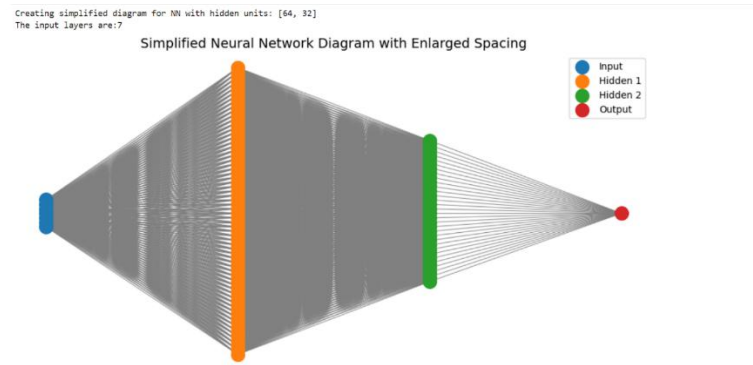
Creating simplified diagram for NN with hidden units: [64, 32]
The input layers are:7

Simplified Neural Network Diagram with Enlarged Spacing

*Figure 13: Neural Network Diagram*

| Model | 1 | 2 | 3 |
|---|---|---|---|
| | | xgb_params | |
| max_depth | 3 | 5 | 7 |
| eta | 0.1 | 0.05 | 0.01 |
| n_estimators | 50 | 100 | 200 |
| precision | 0.996469 | 0.999957 | 0.999498 |
| recall | 0.996413 | 0.9998375 | 0.999817 |
| f1-score | 0.996423 | 0.999897 | 0.999657 |
| accuracy | 0.996413 | 0.999932 | 0.999773 |
| | | nn_params | |
| hidden_units | 64,32 | 128,64 | 256,128 |
| eta | 0.001 | 0.0005 | 0.0001 |
| epochs | 10 | 15 | 20 |
| precision | 0.997999 | 0.9959 | 0.9959 |
| recall | 0.997979 | 0.998854 | 0.998854 |
| f1-score | 0.997983 | 0.997368 | 0.997368 |
| accuracy | 0.997979 | 0.998252 | 0.998252 |

| | | combined | |
|---|---|---|---|
| precision | 0.997952 | 0.999914 | 0.999498 |
| recall | 0.997934 | 0.999375 | 0.999817 |
| f1-score | 0.997937 | 0.999897 | 0.999657 |
| accuracy | 0.997934 | 0.999932 | 0.999773 |

*Table 11: XGBoost and NN Hybrid Results*

The highest accuracy of 0.999932 is achieved in the combined model (Model 2) due to balanced parameters like *max_depth* = 5, eta = 0.05, and *n_estimators* = 100 in XGBoost, alongside a neural network architecture with *hidden_units* = [128, 64], eta = 0.0005, and epochs = 15. These settings allow the model to capture complex patterns while maintaining generalization and stability.The lowest accuracy of 0.996413 is observed in Model 1 with XGBoost, where the shallow tree depth (*max_depth* = 3) and fewer estimators limit its capacity to learn complex patterns. This highlights that combining models with well-tuned parameters enhances performance by leveraging complementary strengths.

# 5.Findings

Training an AI model to achieve high accuracy on the target column in this dataset is not particularly challenging. Since the dataset consists of tabular data, algorithms specifically designed for such data, like *XGBoost* and *Random Forest*, consistently outperform others. *AdaBoost* was another ensemble learning algorithm which performed quite well. Hybridizing these algorithms with additional methods can further improve accuracy, as demonstrated by our implementation combining *XGBoost with a Neural Network*. Ultimately, these models excelled not only in prediction tasks but also in classification.

Among deep learning models, *CNNs*, while not ideal for tabular data, performed well with sigmoid activation for binary classification. Batch normalization, flattening, and ReLU activation consistently improved performance. *MLP* models worked best with two hidden layers of 300 neurons each, requiring only 50 epochs due to the dataset's small size.

*SVM* performance varied significantly with kernel choice: the polynomial kernel was fastest and most accurate, while the linear kernel was computationally expensive. A *hybrid model of k-Means and k-NN* consistently delivered high accuracy with low computation time, benefiting from the dataset's natural clustering of vehicle characteristics. *k-NN* itself excelled with a low number of neighbors, leveraging the dataset's clear clusters and low noise.

*Logistic regression* underperformed due to its linear nature, which was insufficient to capture the dataset's complex feature interactions, such as those distinguishing BEVs from PHEVs.

# 6. References

Rathore, H., Meena, H. K., & Jain, P. (2023). *Prediction of EV Energy consumption Using Random Forest And XGBoost* . Prerna Jain.

Risk, B. (2024). *Electric Vehicle Data Analysis and Predictions.* Retrieved from Kaggle: https://www.kaggle.com/code/devraai/electric-vehicle-data-analysis-and-predictions