Name: Semiu Kolapo
Student Number: 501145293
Supervisor: Tamer Abdou
Date: April 1st 2024

# Classification and Trends: Distribution of Electric Vehicle Types in Washington State

## Model Evaluation and Results

## Summary of Libraries Used

In the data analysis and machine learning endeavors, we employed a range of powerful libraries that facilitated various tasks from data manipulation to model evaluation.

**NumPy (np)** served as the backbone for efficient scientific computing, providing support for large, multi-dimensional arrays and a plethora of mathematical functions.

**Pandas (pd)** emerged as a crucial tool for data manipulation and analysis, offering intuitive data structures like DataFrame and Series, along with functions for data cleaning, transformation, and analysis.

**Matplotlib.pyplot (plt)** and **Seaborn (sns)** played pivotal roles in visualizing complex datasets, offering an array of plotting functions for creating informative and visually appealing statistical graphics.

**Scikit-learn** emerged as our go-to library for machine learning tasks, offering a wide range of algorithms for classification, regression, clustering, and more. It also provided tools for model evaluation, selection, and preprocessing.

**Imbalanced-learn (imblearn)** addressed the challenge of imbalanced datasets by providing resampling techniques to balance class distributions, ensuring robust model performance.

Finally, **SHAP (shap)** added a layer of interpretability to our machine learning models by computing Shapley values, allowing us to understand the contribution of each feature to the model's output.

Together, these libraries formed a comprehensive toolkit that empowered us to analyze data, build robust machine learning models, and gain valuable insights into our data.

## Modeling Algorithms

### Naive Bayes Classifier

A Naive Bayes classifier is a widely used machine learning algorithm for classification tasks. It applies Bayes' theorem to compute the probability of an input belonging to a specific class. Despite its simplistic assumption of feature independence, Naive Bayes classifiers are effective, particularly in text classification and spam filtering tasks (Rennie et al., 2003).[1]

### Logistic Regression

Logistic regression is a core statistical technique widely used in machine learning for binary classification tasks. It predicts the probability of binary outcomes, such as whether an email is spam or not. By employing the logistic function, it transforms input features into probabilities, offering insights into the likelihood of data points belonging to different classes. Unlike linear regression, logistic regression excels in handling categorical variables and is essential for various applications like spam filtering and risk assessment (Menard, 2002).[2]

### Random Forest

Random forest is a robust and versatile machine learning algorithm known for its effectiveness in both classification and regression tasks. It operates on the principle of ensemble learning, where multiple decision trees are combined to produce more accurate and stable predictions. By aggregating the outputs of individual trees, random forests mitigate overfitting and improve generalization. They are flexible in handling different data types and have wide applicability across various domains. Breiman's seminal work on random forests (2001) laid the foundation for understanding and utilizing this powerful algorithm.[3]

---

[1] Rennie, J. D., Shih, L., Teevan, J., & Karger, D. R. (2003). Classifying web pages with naive Bayes. In Proceedings of the thirteenth ACM international conference on Information and knowledge management (pp. 489-497).

[2] Menard, S. (2002). Applied logistic regression. Sage Publications

[3] Breiman, L. (2001). Random forests. Machine learning, 45(3), 5-32.

XGBoost

      XGBoost (eXtreme Gradient Boosting) is a highly effective machine learning algorithm renowned for its scalability, accuracy, and flexibility. It belongs to the ensemble learning family, employing the gradient boosting framework to sequentially add decision trees as weak learners. Through regularization techniques, XGBoost prevents overfitting and ensures robust generalization on unseen data. Notably, its scalability enables efficient processing of large datasets, thanks to parallel processing and out-of-core computation. Moreover, XGBoost provides insights into feature importance, aiding in understanding which features drive model predictions. Chen and Guestrin's seminal work (2016) elucidated XGBoost's principles and its significance in the machine learning landscape.[4]

---

[4] Chen, T., & Guestrin, C. (2016). Xgboost: A scalable tree boosting system. In Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining (pp. 785-794).

# Initial Results

### Confusion Matrix of Naive Bayes Classifier

| | Actual 0 | Actual 1 |
|---|---|---|
| Predicted 0 | True Negative (TN) 72.23% | False Negative (FN) 5.89% |
| Predicted 1 | False Positive (FP) 19.54% | True Positive (TP) 2.35% |

### Confusion Matrix of Logistic Regression

| | Actual 0 | Actual 1 |
|---|---|---|
| Predicted 0 | True Negative (TN) 40.46% | False Negative (FN) 37.65% |
| Predicted 1 | False Positive (FP) 10.74% | True Positive (TP) 11.15% |

### Confusion Matrix of Random Forest

| | Actual 0 | Actual 1 |
|---|---|---|
| Predicted 0 | True Negative (TN) 78.03% | False Negative (FN) 0.09% |
| Predicted 1 | False Positive (FP) 0.01% | True Positive (TP) 21.88% |

### Confusion Matrix of XGBoost

| | Actual 0 | Actual 1 |
|---|---|---|
| Predicted 0 | True Negative (TN) 78.02% | False Negative (FN) 0.09% |
| Predicted 1 | False Positive (FP) 0.01% | True Positive (TP) 21.87% |

| | Model | Accuracy (train) | Accuracy (test) | ROC AUC (train) | ROC AUC (train) |
|---|---|---|---|---|---|
| 0 | Random Forest | 1.000000 | 0.998801 | 1.000000 | 0.999994 |
| 1 | XGBoost | 0.999824 | 0.998941 | 0.999988 | 0.999983 |
| 2 | Naïve Bayes | 0.780551 | 0.780491 | 0.722064 | 0.727559 |
| 3 | Logistic Regression | 0.598806 | 0.516097 | 0.593621 | 0.514732 |

Based on the evaluation metrics, the Random Forest model performed the best among the four models, achieving perfect accuracy on the training data (1.000) and very high accuracy on the test data (0.999). It also achieved a perfect ROC-AUC score on the training data (1.000) and a near-perfect score on the test data (0.999).

The XGBoost model also performed well, with high accuracy on both the training data (0.999) and the test data (0.999). Similarly, it achieved excellent ROC-AUC scores on both the training data (0.999) and the test data (0.999). The Naive Bayes model had lower accuracy compared to the Random Forest and XGBoost models, achieving an accuracy of 0.780 on both the training and test data. The ROC-AUC scores for this model were also lower, indicating less robust performance compared to the other models.

The Logistic Regression model performed the poorest among the four models, with an accuracy of 0.599 on the training data and 0.516 on the test data. The ROC-AUC scores were also relatively low, suggesting suboptimal performance compared to the other models.
Overall, the Random Forest and XGBoost models demonstrated superior performance in terms of both accuracy and ROC-AUC scores compared to the Naive Bayes and Logistic Regression models.

Feature Importance

In preparation for the final result, feature importance analysis is conducted to identify patterns, relationships, and determine which features have the most significant impact on the target variable. This analysis should help prioritize features that contribute the most to model performance and provides insights into the underlying data characteristics.

|   | Feature | Importance |
|---|---------|-----------|
| 0 | Electric Range | 0.453681 |
| 1 | Clean Alternative Fuel Vehicle Eligibility | 0.196643 |
| 2 | Model | 0.160677 |
| 3 | Make | 0.052106 |
| 4 | VIN (1-10) | 0.040262 |

|   | Feature | Importance |
|---|---------|-----------|
| 0 | State | 3.712071e-08 |
| 1 | Electric Utility | 9.238881e-05 |
| 2 | Base MSRP | 3.686047e-04 |
| 3 | Legislative District | 7.104148e-04 |
| 4 | Latitude | 8.808175e-04 |

## Final Results

Dropping features is a fundamental preprocessing step in machine learning, crucial for enhancing model performance and efficiency. By eliminating irrelevant or redundant data points, the model can focus on learning the true underlying patterns in the data, leading to more accurate predictions. This process not only reduces noise but also lowers computational complexity, resulting in faster training times and potentially better generalization to unseen data. Techniques such as correlation analysis and feature importance scores, as discussed by Guyon and Elisseeff (2003), are invaluable for identifying and removing uninformative features.[5] In this context, the decision to remove "State" and "Electric Utility" based on their low importance underscores the importance of domain knowledge in feature selection. While it's essential to streamline the feature set, it's equally important to exercise caution and retain features that may hold value based on domain expertise and understanding of the data's context.
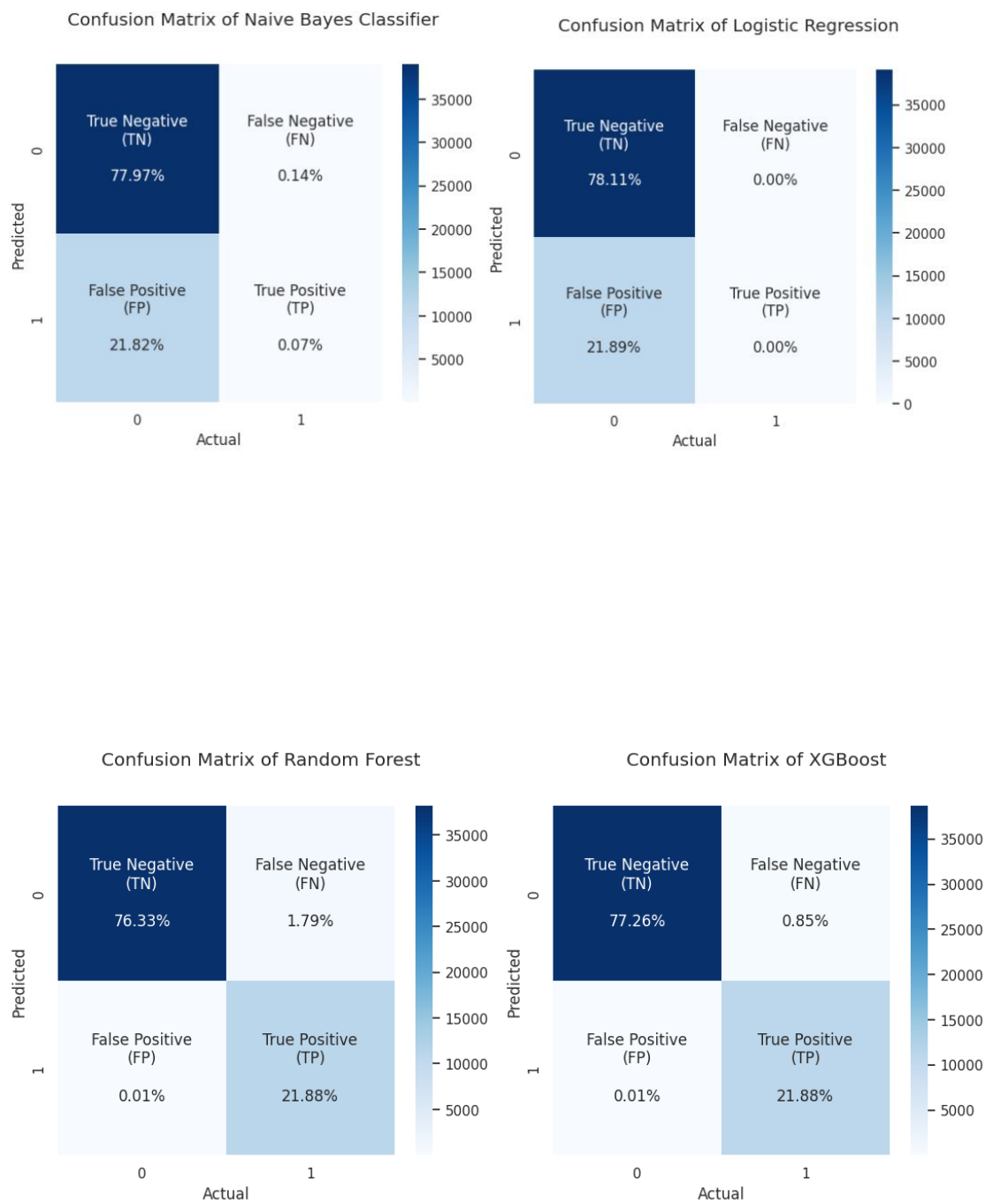
Feature scaling is an essential preprocessing step in machine learning, especially for algorithms sensitive to variations in feature magnitudes. By transforming features onto a common scale, scaling techniques like StandardScaler from scikit-learn can significantly enhance model performance. StandardScaler, in particular, standardizes features by subtracting the mean and dividing by the standard deviation, resulting in features with a mean of 0 and a standard deviation of 1. This process offers several advantages, as outlined by Pedregosa et al. (2011).[6]

Firstly, it facilitates improved model convergence by ensuring that gradients have comparable magnitudes across features, preventing slow convergence or divergence. Secondly, it

[5] Guyon, I., & Elisseeff, A. (2003). An introduction to variable and feature selection. Journal of Machine Learning Research, 3(Mar), 1157-1182. Retrieved from https://jmlr.org/papers/special/feature03.html

[6] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Vanderplas, J. (2011). Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12(Oct), 2825-2830.

enables fair comparisons between features by eliminating biases towards features with larger scales, thus allowing informative features with smaller scales to contribute meaningfully to the model. In this context, applying feature scaling to numerical features in the dataset ensures that the model can effectively leverage all available information without being influenced by variations in feature scales.

Confusion Matrix of Naive Bayes Classifier

| | Actual 0 | Actual 1 |
|---|---|---|
| Predicted 0 | True Negative (TN) 77.97% | False Negative (FN) 0.14% |
| Predicted 1 | False Positive (FP) 21.82% | True Positive (TP) 0.07% |

Confusion Matrix of Logistic Regression

| | Actual 0 | Actual 1 |
|---|---|---|
| Predicted 0 | True Negative (TN) 78.11% | False Negative (FN) 0.00% |
| Predicted 1 | False Positive (FP) 21.89% | True Positive (TP) 0.00% |

Confusion Matrix of Random Forest

| | Actual 0 | Actual 1 |
|---|---|---|
| Predicted 0 | True Negative (TN) 76.33% | False Negative (FN) 1.79% |
| Predicted 1 | False Positive (FP) 0.01% | True Positive (TP) 21.88% |

Confusion Matrix of XGBoost

| | Actual 0 | Actual 1 |
|---|---|---|
| Predicted 0 | True Negative (TN) 77.26% | False Negative (FN) 0.85% |
| Predicted 1 | False Positive (FP) 0.01% | True Positive (TP) 21.88% |

| | Model | Accuracy | ROC AUC |
|---|---|---|---|
| 0 | XGBoost | 0.991387 | 0.999793 |
| 1 | Random Forest | 0.982074 | 0.999787 |
| 2 | Logistic Regression | 0.781131 | 0.499715 |
| 3 | Naive Bayes | 0.780531 | 0.494446 |

After applying feature scaling and dropping irrelevant features, the model performances have notably improved. Both XGBoost and Random Forest models exhibit significantly higher accuracy scores compared to Logistic Regression and Naive Bayes. XGBoost leads with the highest accuracy score of 99.14%, closely followed by Random Forest at 98.21%. These ensemble methods, known for their robustness and ability to handle complex datasets, outperform the traditional logistic regression and naive Bayes classifiers. Notably, logistic regression and naive Bayes models show minimal improvement, with accuracy scores hovering around 78%, indicating that they may not be well-suited for this particular dataset or may require more sophisticated preprocessing steps to achieve higher performance. The substantial boost in performance for XGBoost and Random Forest highlights the effectiveness of feature scaling and feature selection in enhancing model accuracy and robustness.

Based on the analysis of electric vehicle data, King County stands out as the epicenter of electric vehicle adoption, boasting a total of 86,594 electric vehicles, comprising both Plug-in Hybrid Electric Vehicles (PHEVs) and Battery Electric Vehicles (BEVs). This significant number underscores the region's commitment to sustainable transportation solutions. Notably, Tesla emerges as the dominant player in the electric vehicle market in this area, with its vehicles constituting 44.86% of all electric vehicles. The most common make among electric vehicles is TESLA, indicating the brand's strong presence and popularity among consumers. In terms of model year, 2023 emerges as the most prevalent, reflecting the continuous advancement and adoption of newer electric vehicle models.

Delving into city-level analysis, Seattle emerges as a hub for electric vehicle adoption, ranking as the top city for both PHEVs and BEVs. With 5,747 PHEVs and a staggering 22,084 BEVs registered, Seattle showcases a robust infrastructure and supportive community for electric vehicle

ownership. Other cities in the region, such as Bellevue and Redmond, also demonstrate significant uptake of electric vehicles, further highlighting the widespread acceptance and integration of sustainable transportation alternatives in the Pacific Northwest. Overall, these findings underscore the region's progressive stance towards reducing carbon emissions and embracing cleaner, more environmentally friendly modes of transportation.

Reference

1.  Rennie, J. D., Shih, L., Teevan, J., & Karger, D. R. (2003). Classifying web pages with naive Bayes. In Proceedings of the thirteenth ACM international conference on Information and knowledge management (pp. 489-497).

2.  Menard, S. (2002). Applied logistic regression. Sage Publications

3.  Breiman, L. (2001). Random forests. Machine learning, 45(3), 5-32.

4.  Chen, T., & Guestrin, C. (2016). Xgboost: A scalable tree boosting system. In Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining (pp. 785-794).

5.  Guyon, I., & Elisseeff, A. (2003). An introduction to variable and feature selection. Journal of Machine Learning Research, 3(Mar), 1157-1182. Retrieved from https://jmlr.org/papers/special/feature03.html

6.  Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Vanderplas, J. (2011). Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12(Oct), 2825-2830.