

MILES AND EMISSIONS:

How cars contribute to CO₂ levels



3 GOOD HEALTH AND WELL-BEING



11 SUSTAINABLE CITIES AND COMMUNITIES



12 RESPONSIBLE CONSUMPTION AND PRODUCTION



13 CLIMATE ACTION



ADRIAN MARC FEDIER PURNAMA - 2702247210

GERVASIUS RUSSELL - 2702247450

HANS ARDIANTA - 2702241063

JONATHAN WILLIAM GUNAWAN - 2702251794

DATA SCIENCE, BINUS UNIVERSITY



Intro: What is this about?

As we all know, air pollution has become a significant global challenge in this modern era. Vehicular carbon dioxide (CO₂) emissions has become one of the major contributor to the polluted air condition. So by this project, we are aiming to create vehicular CO₂ emission prediction based on vehicle characteristics and behaviors by implementing machine learning techniques. The dataset is obtained from Kaggle. The dataset comprises car characteristics, fuel consumption behavior, and the target variable, carbon dioxide emissions (g/km). Since the target variable in this dataset is continuous, we employ regression machine learning models for prediction. The ultimate goal of this project is to contribute to the identification of carbon dioxide emissions produced by cars, providing valuable insights that can support efforts to reduce environmental impact.

Data Understanding

Model	Vehicle Class	Engine Size (cc)	Cylinders	Power (bhp)	Fuel Type	Fuel Consumption City (l/100 km)	Fuel Consumption Hwy (l/100 km)	Fuel Consumption City (l/100 km)	Fuel Consumption Hwy (l/100 km)	CO2 (g/km)
0	ALFA ROMEO	1.3	4	105	P	7.7	5.7	7.7	5.7	160

This is the independent variable (Y) or the target variable

6 numerical variables

5 categorical variables

Data Quality issues:

This data has no null/missing values, the only quality issue in this dataset is outliers, but it is also only in small quantity.

Handling Outliers

To handle outliers, we use the **Winsorizer** technique with caps at the 10th and 90th percentiles. Why Winsorizer? Winsorization is more robust against outliers compared to other methods like IQR or z-score because it **limits extreme values instead of removing or fully transforming them**. Winsorizer adjusts values while preserving the overall data structure and minimizing the loss of information. This makes it particularly effective for datasets with skewed distributions or extreme outliers.

We can also see that every numerical feature in the dataset is **highly correlated** with the target variable

Machine Learning Model

20% test data
80% train data

Models that are used

Decision Tree

Linear Regression

Random Forest

XGBRegressor

ensemble learning

Key takeouts

- Decision Tree - **best performer** - records highest R² Score and the lowest MAE and MSE values, indicating **strong accuracy**.
- Linear Regression - **weakest performer**
- Both Random Forest and XGBRegressor deliver similar results, performing well but **slightly below** the Decision Tree.

Feature Engineering

Label Encoding

Normalization (RobustScaler)

Feature Selection with PCA

Data ready to be split into train and test data

Model Performance Evaluation

Label Encoder converts categorical values into numerical labels, allowing machine learning models to process and interpret categorical data.

(for example we use "Model" variable from this dataset)

Before	After
F-150 FV	0
MUSTANG	1
CAMARO	2
SIERRA	3
ACCORD	4
SONIC	5
SIERRA	6
and so on...	and so on...

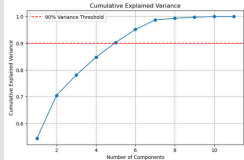
Encoding creates categorical variables into numerical formats, preserving their information and making them compatible with machine learning models.

Why do we use RobustScaler?

RobustScaler is chosen because it is well-suited for skewed data and less sensitive to outliers. Unlike Min-Max Scaler or Z-Scaler, which rely on the mean and standard deviation—making them prone to the influence of outliers and non-normally distributed data—RobustScaler uses the median and interquartile range (IQR). This makes it more effective in scaling data while maintaining robustness against extreme values, ensuring better performance for datasets with skewed distributions or outliers.

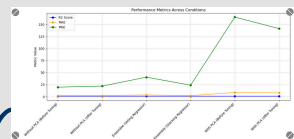
RobustScaler formula:

$$X_{\text{scale}} = \frac{x_i - x_{\text{med}}}{x_{75} - x_{25}}$$



Number of components for 90% explained variance

We use PCA with 5 components, determined through the Cumulative Explained Variance (CEV) technique with a 90% threshold. This means we select the smallest number of principal components that collectively explain at least 90% of the variance in the data.



The plot shows that without PCA, the model performs well and improves further with tuning and ensemble methods. In contrast, applying PCA significantly worsens performance, with higher errors (MSE and MAE) and near-zero R² scores, even after tuning, indicating PCA is not beneficial for this dataset.

Conclusion

By implementing various machine learning techniques, we managed to make car CO₂ emission prediction model and get the most out of it. The results show that the model's accuracy is higher without PCA, indicating that the original features are already sufficiently informative and we can use Decision Tree or Random Forest model to achieve high prediction accuracy. This project is beneficial for supporting air pollution control efforts by providing a predictive tool to identify high-emission vehicles, design strategies to reduce emissions, and support data-driven policy development. By leveraging machine learning technology, this project contributes to creating more efficient, data-based solutions to improve air quality and promote environmental sustainability.

To see the full code, you can scan this QR code.



The code is available and provided in Google Colab.