**Predictive Modeling Analysis: Used Vehicle Price Estimation**

# 1. Introduction

### 1.1 Background

The used vehicle market is influenced by various factors, including mileage, condition, age, and market trends. Traditional pricing models often assume that vehicle condition is the dominant factor, but our analysis challenges this assumption. By leveraging machine learning, we aimed to uncover the true determinants of used car prices.

### 1.2 Objective

The primary goal of this study was to build a predictive model that accurately estimates used vehicle prices based on key attributes. We explored multiple modeling techniques to identify the most effective approach and determine the most impactful features.

# 2. Data Preparation

### 2.1 Data Source

The dataset used for this study was obtained from Kaggle's 'Vehicle Sales Data' and contained information on vehicle characteristics, sales prices, and transaction details.

### 2.2 Data Cleaning and Feature Engineering

To ensure high-quality predictions, we performed the following preprocessing steps:

- **Data Type Conversion:** Transformed 'saledate' into a DateTime format.
- **Feature Creation:** Introduced 'price per mile' as an additional predictor.
- **Encoding Categorical Variables:** Converted categorical data (e.g., 'body', 'make', 'model') into numerical values using one-hot encoding.
- **Multicollinearity Check:** Used the Variance Inflation Factor (VIF) to remove redundant features.
- **Handling Missing Data:** Removed null values and standardized missing entries.

# 3. Model Selection and Training

### 3.1 Understanding R² Score

The **R² score** (coefficient of determination) measures how well a model explains the variance in the target variable (price). A score of 1 indicates perfect prediction, while 0 means the model explains no variance. Higher $R^2$ values suggest better model performance.

### 3.2 Baseline Models

We first experimented with traditional models to establish a baseline for performance:

- **Linear Regression**: Provided a basic understanding of variable relationships but lacked accuracy due to its linear assumptions.
- **Random Forest Regressor**: Improved performance by capturing nonlinear relationships but still had some variance.

### 3.3 Advanced Model: XGBoost

To enhance prediction accuracy, we implemented **XGBoost**, a powerful gradient– boosting model. After hyperparameter tuning, the final model achieved an **R² score of 0.9934**, indicating exceptional predictive power.

### 3.4 Model Performance Metrics

| Model | R² Score |
|---|---|
| Linear Regression | 0.3624 |
| Random Forest | 0.6256 |
| XGBoost (Tuned) | **0.9934** |

# 4. Key Insights

### 4.1 Feature Importance

Our analysis revealed surprising insights:

- **Market valuation (MMR) was the most significant factor** influencing price, indicating that industry trends shape vehicle pricing more than individual attributes.

- **Mileage (odometer) had a stronger effect than reported condition**, challenging the assumption that physical condition is the primary determinant.
- **Price per mile was a crucial indicator**, reflecting how buyers assess a car's remaining value.

## 4.2 Business Implications

- Buyers should prioritise **market trends and mileage** over condition ratings when evaluating vehicle prices.
- Dealerships can adjust pricing strategies by incorporating data-driven insights rather than relying solely on traditional valuation methods.
- Predictive modeling can be used to optimise pricing recommendations and inventory management in the used vehicle market.

# 5. Conclusion

This study demonstrated that predictive modeling is a powerful tool for estimating used vehicle prices. While condition remains a factor, **market valuation and mileage are stronger predictors**. Our findings suggest a shift in how consumers and dealerships should assess vehicle pricing, leveraging data-driven approaches for more accurate evaluations.

# 6. Future Work

To further refine our model, we recommend:

- **Expanding the dataset** to include regional variations and economic factors.
- **Testing deep learning techniques** for additional accuracy gains.
- **Developing a real-time pricing tool** for dealerships and consumers.