# LLM Data Insights AutoML Report – cars

Generated by your multi-agent LLM system.

# 1. Exploratory Data Analysis (EDA)

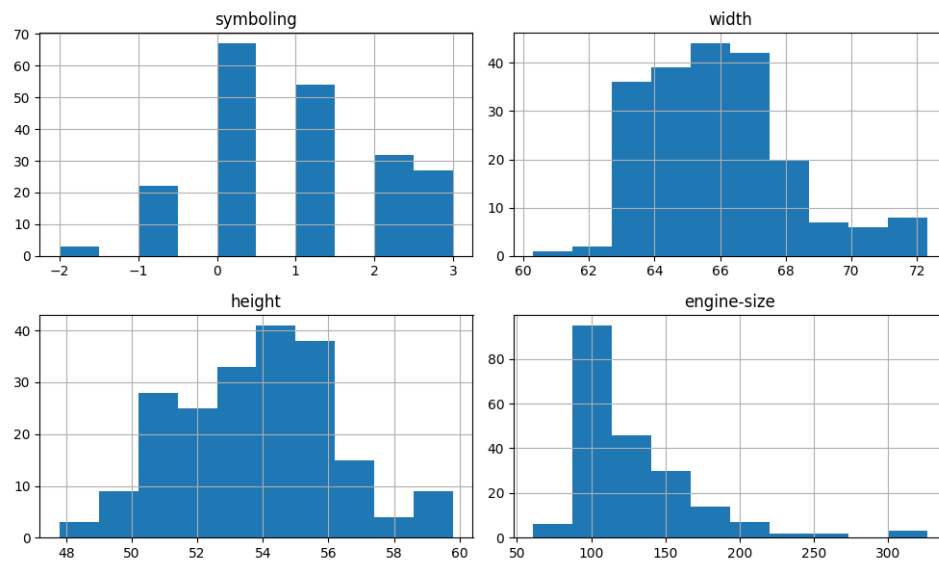EDA summary is available. Dataset shape: [205, 15]. Number of columns: 15.
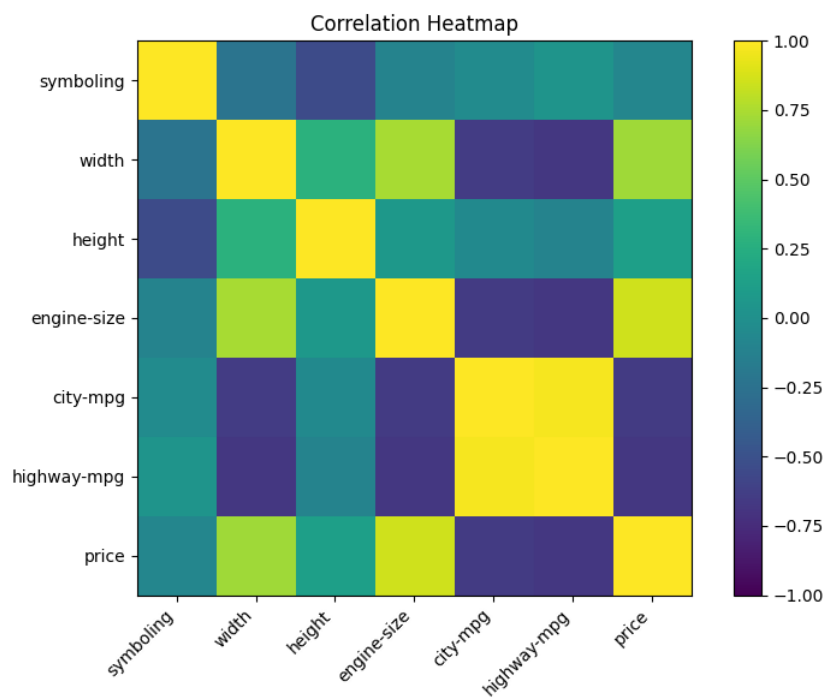


Figure 1: Numeric feature histograms



Figure 2: Correlation heatmap

## LLM Insights (EDA)

Here are the answers in clear bullet points:

**1. Description of the dataset:**

* Rows: 205
* Columns: 12 (symboling, normalized-losses, make, fuel-type, body-style, drive-wheels, engine-location, width, height, engine-type, engine-size, horsepower, city-mpg, highway-mpg, price)
* Data types:
+ Integers (symboling, engine-size, city-mpg, highway-mpg, price)
+ Floating-point numbers (width, height, normalized-losses)
+ Strings (make, fuel-type, body-style, drive-wheels, engine-location, horse-power)

**2. Comment on data quality:**

* Missing values are present in most columns, with the lowest percentage being 0% for symboling and make.
* Outliers are detected in several columns, including width, height, engine-size, city-mpg, highway-mpg, and price.
* Potential issue: The use of categorical variables (e.g. fuel-type, body-style) might lead to difficulties in modeling if not handled properly.

**3. Important numeric insights:**

* The average engine size is around 126.9 units, with a standard deviation of 41.6 units, indicating that most engines are relatively small.
* The city and highway MPG ranges from 13 to 49 for city-mpg and 16 to 54 for highway-mpg, suggesting a mix of fuel-efficient and less efficient vehicles.
* There is a moderate negative correlation between width and height (-0.541), indicating that wider cars tend to be shorter.

**4. Notable outliers:**

* In the width column, values above 71.1 (corresponding to the 99th percentile) are considered outliers, with 8 instances.
* In the city-mpg column, values below 2.5 and above 46.5 are considered outliers, with 2 instances.
* In the price column, values above 29568.0 and below -5280.0 are considered outliers, with 14 instances.

**5. Follow-up analyses or questions:**

* Investigate the relationship between fuel-type and engine-size to determine if more efficient engines are associated with certain types of fuel.
* Explore the distribution of normalized-losses to understand what type of losses are being recorded.
* Analyze the correlations between horsepower, city-mpg, highway-mpg, and price to identify potential relationships between these variables.

# 2. Supervised Modeling

Problem type: regression. Target column: price. Algorithm used: linear. Metrics: {'type': 'regression', 'rmse': 4498.699685255597, 'r2': 0.7476145642798953}.

## LLM Insights (Model)

**Summary**

* The dataset has 205 rows and 15 columns, with key data types including integer (int64) and floating-point numbers (float64).
* The target column "price" represents the cost of each car, making this a regression problem since we're trying to predict a continuous value.

**Model Performance**

* For regression, the model achieved:
+ RMSE (Root Mean Squared Error): 4498.70
+ R² (Coefficient of Determination): 0.75
* These metrics indicate that the model is doing reasonably well, but there's still room for improvement.
* The RMSE value suggests that the model is not very accurate in predicting prices, with a relatively high error margin.
* The R² value indicates that about 75% of the variance in prices can be explained by the model.

**Interpretation using EDA**

* Features strongly correlated with price are:
+ Engine size
+ Highway MPG
+ Price itself (obviously!)
* Feature importances from the model show that:
+ Engine size is highly important, but less so than highway MPG
+ Other features have relatively lower importance
* Outliers in the target column "price" might affect metrics like RMSE and R². Specifically, there are 14 outliers with values above $45,000.

**Practical Insights**

* The model suggests that engine size and highway MPG are important factors affecting car prices.
* Cars with larger engines tend to be more expensive, but so do cars with better fuel efficiency (i.e., higher highway MPG).
* The high RMSE value indicates that there's a lot of variation in prices that the model can't capture.

**Next Steps**

* Feature engineering: We could try to create new features from existing ones, such as "horsepower per unit" or "fuel efficiency ratio".
* Outlier handling: Since there are outliers in the target column and key features like engine size, we might want to consider techniques like winsorization or robust regression.
* Model improvements: With better feature engineering and outlier handling, we could potentially improve the model's accuracy and reduce the RMSE value.

# 3. Hyperparameter Tuning

Best algorithm: RandomForestRegressor with rmse=3855.338426589406. Best params: {'n_estimators': 300, 'min_samples_split': 2, 'min_samples_leaf': 1, 'max_depth': 20}.

## LLM Insights (Hyperparameters)

**Summary:**

* Four machine learning algorithms were tried: Random Forest Regressor, Gradient Boosting Regressor, Ridge, and Lasso.
* The best performing algorithm is the Random Forest Regressor with an RMSE (Root Mean Squared Error) value of 3855.3384.

**Explanation of the metric:**

* RMSE is a measure of the difference between predicted and actual values in regression problems. It's calculated as the square root of the average squared difference.
* In this case, an RMSE of 3855.3384 means that, on average, the model's predictions were off by about 3855.

**Interpretation of hyperparameters:**

* **Tree-based models (Random Forest / Gradient Boosting):**
+ `n_estimators`: The more trees you have in your forest, the better the model can capture complex relationships in the data. However, too many trees can lead to overfitting.
+ `max_depth`: Limiting the maximum depth of each tree helps prevent overfitting by preventing trees from splitting into too many sub-branches.
* **Linear models (Ridge, Lasso):**
+ Regularization (alpha or C) helps prevent overfitting by adding a penalty term to the loss function. A higher regularization value means less feature importance and more generalization.

**Recommendation:**

* I would deploy the Random Forest Regressor with the best hyperparameters: `n_estimators=300`, `min_samples_split=2`, `min_samples_leaf=1`, and `max_depth=20`.
* This configuration balances model complexity and regularization, making it a good trade-off between accuracy and overfitting.

**Future experiment ideas:**

* **Try more trees**: Increasing the number of trees in the Random Forest Regressor can lead to better performance but may also increase computational cost.
* **Different learning rates**: Experimenting with different learning rates for Gradient Boosting Regressor can help find the optimal rate for this algorithm.
* **More data**: Gathering additional data points could improve model performance, especially if it contains more relevant features or variations in the target variable.
* **Feature engineering**: Exploring new feature combinations or transformations could lead to better model performance and interpretation of the results.

# 4. Unsupervised Analysis

## LLM Insights (Unsupervised)

**Dataset Summary and Method**
* The dataset is a mixture of 14 features (e.g., car specifications) with no explicit labels or categories.
* An unsupervised method, K-means clustering, was applied to group similar data points into clusters.

**PCA Metrics**

* **Explained Variance Ratio Per Component**: Not provided.
* **Total Explained Variance**: Not provided.

**Clustering Metrics**

* **Silhouette Score**: 0.369
+ Indicates moderate separation between clusters and cohesion within clusters.
* **Calinski-Harabasz Score**: 137.612
+ Suggests high separation between clusters and good clustering quality.
* **Davies-Bouldin Score**: 0.902
+ Indicates low similarity between clusters, suggesting distinct groups.

**Algorithm-Specific Metrics**

* **KMeans:**
+ **Inertia**: 59527.483
- Measures the sum of squared distances between each data point and its assigned cluster center.
+ **Cluster Centers**: 3 centers with varying coordinates (see above).
- Represent the mean values of the clusters.

**Data Structure and Key Features**

* The K-means clustering revealed three main components/clusters, which differ in key features such as:
+ Engine size and horsepower
+ Fuel type and engine location
+ City and highway fuel efficiency
* The cluster centers suggest a mix of car types, with some clusters favoring certain engine sizes or fuel efficiencies.

**Important Correlations and Outliers**

* **Correlations**: Cars in the same cluster tend to have similar correlations between features, such as positive relationships between horsepower and engine size.
* **Outliers**: Some data points have significant outliers in certain features (e.g., extremely high or low horsepower values), which may influence clustering results.

**Practical Applications**

* Segmentation: Identify distinct car types based on specifications and fuel efficiency.
* Anomaly Detection: Detect unusual cars with high-performance engines or unusual fuel

efficiencies.
* Feature Engineering: Use clustering insights to create new, relevant features for machine learning models.

**Limitations and Caveats**

* **Data Quality Issues**: Missing values, outliers, and skewed distributions may affect clustering results.
* **Interpretation Carefulness**: Clustering patterns should be interpreted with caution due to the potential influence of outliers and data quality issues.