

Hastings DIRECT

# INSURANCE CLAIM ANALYSIS

*Presented by: Sarah Moin*





# PROJECT OVERVIEW

- **Objective:** Analyze insurance claims and predict claim severity.
- **Importance:** Improve risk assessment, reduce fraud, and optimize operations.
- **Tools Used:** Python, Pandas, Matplotlib, Scikit-learn.

# PROBLEM STATEMENT

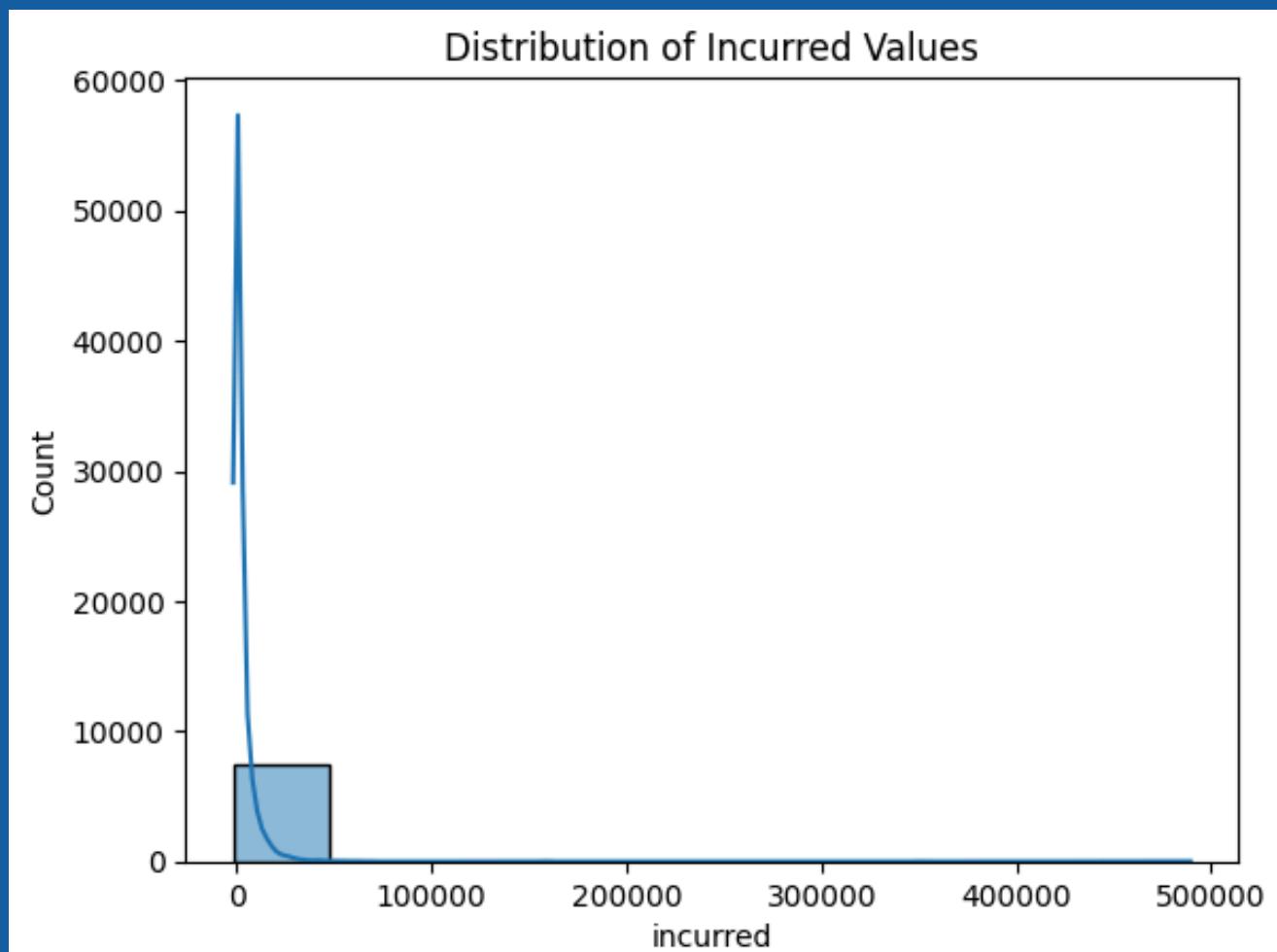
- Challenges in understanding claim behaviour.
- Need for identifying predictive features.
- Ensuring data quality and readiness for weekly updates.



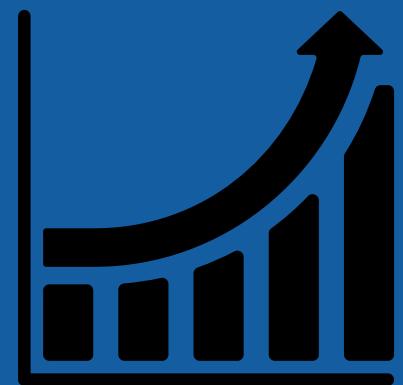
# Data Exploration

## Understanding the Dataset

- Dataset contains important columns: claim\_count, incurred, vehicle\_make, net\_earned\_premium, etc.
- Initial Observations:
  - a. Missing values in Incurred.
  - b. Outliers in Incurred (negative values).
- Visualization: Visualizations for some of the columns are shown in the next pages.



# Data Statistics



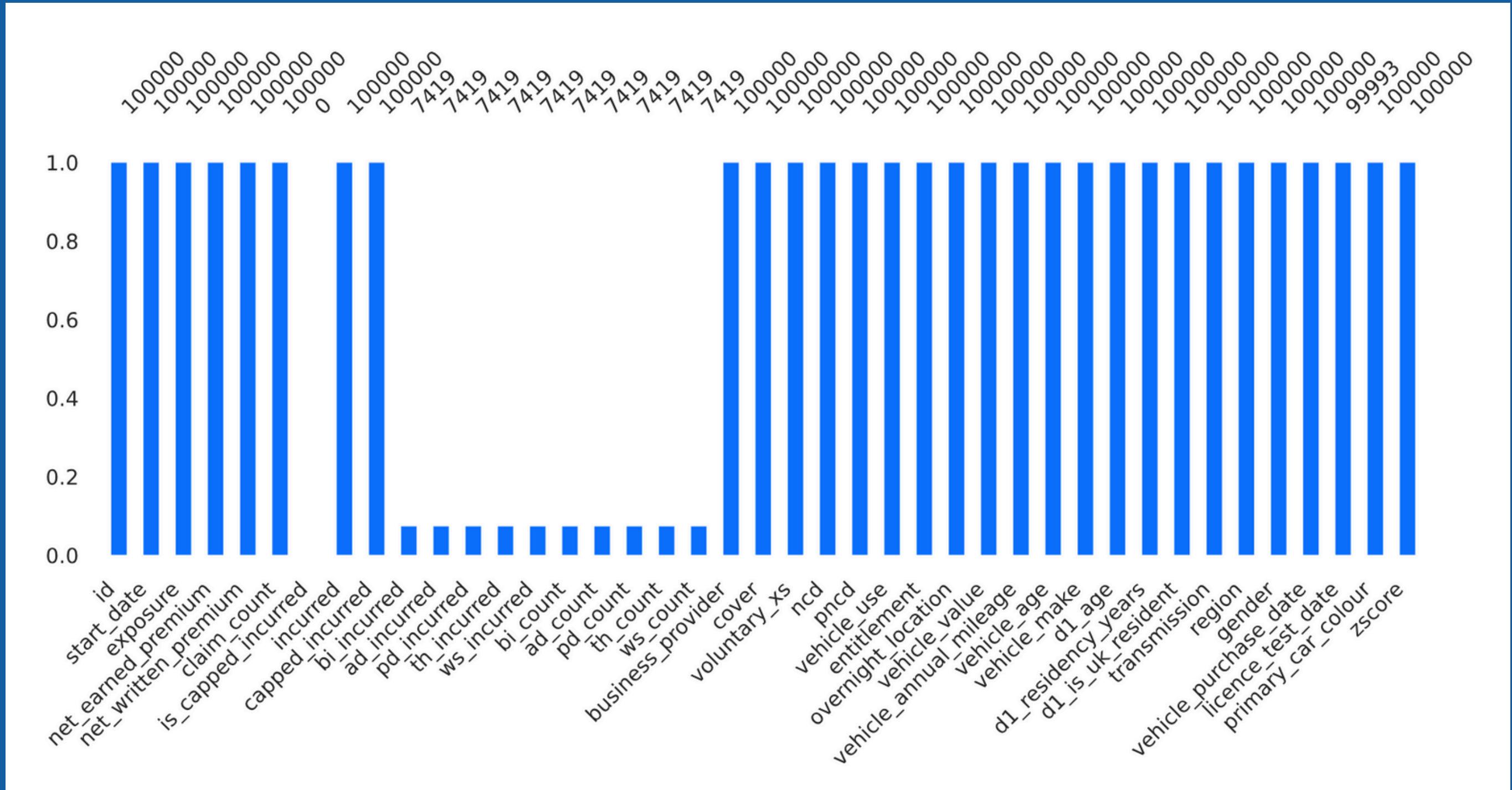
## Dataset statistics

Number of variables	41
Number of observations	100000
Missing cells	1025817
Missing cells (%)	25.0%
Duplicate rows	0
Duplicate rows (%)	0.0%
Total size in memory	88.6 MiB
Average record size in memory	928.6 B

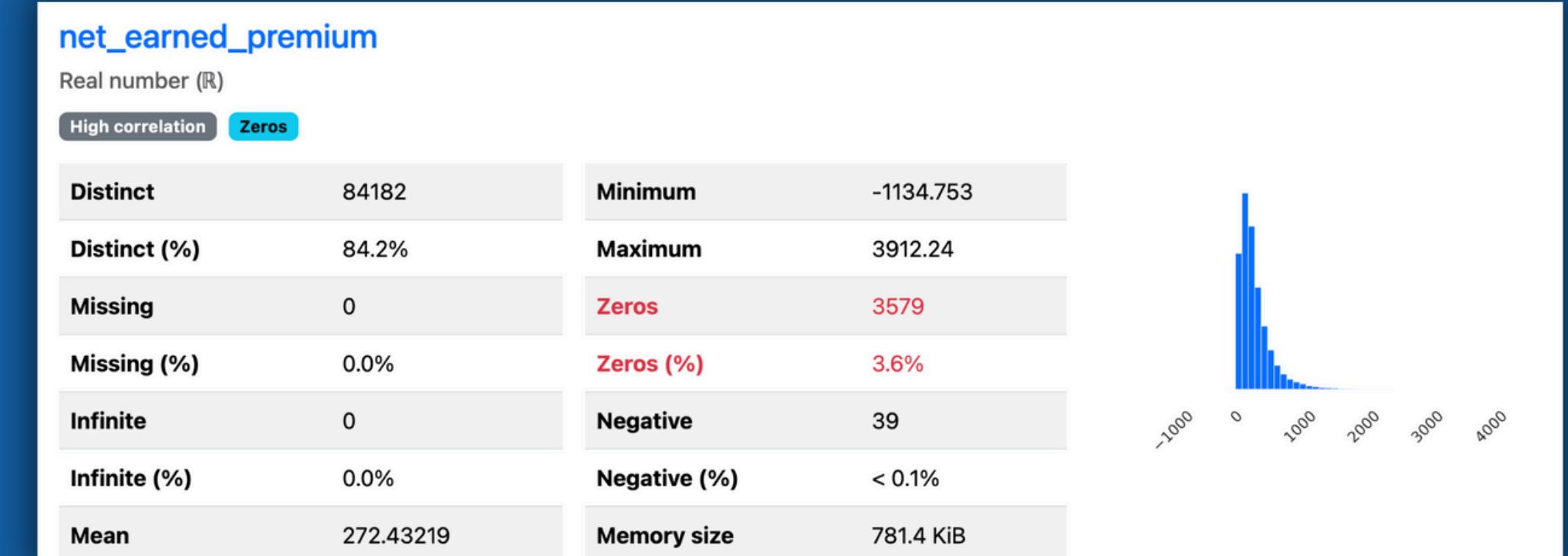
## Variable types

Numeric	19
DateTime	3
Categorical	15
Unsupported	1
Text	2
Boolean	1

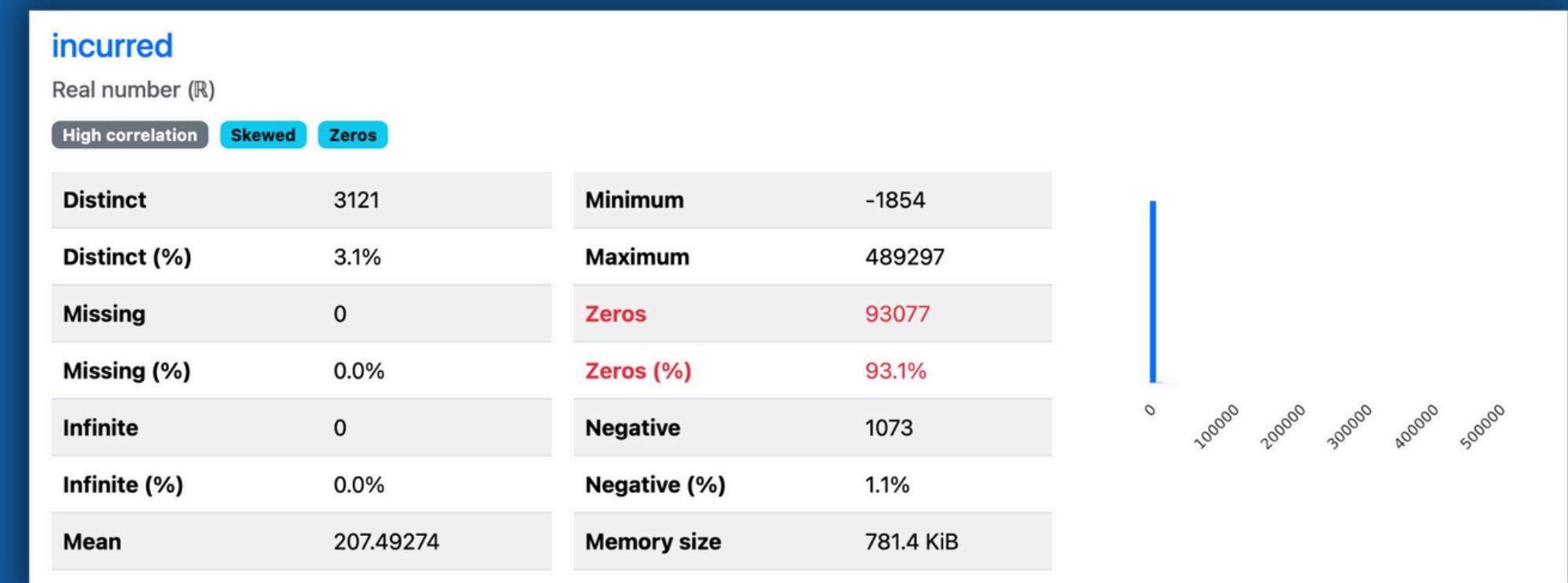
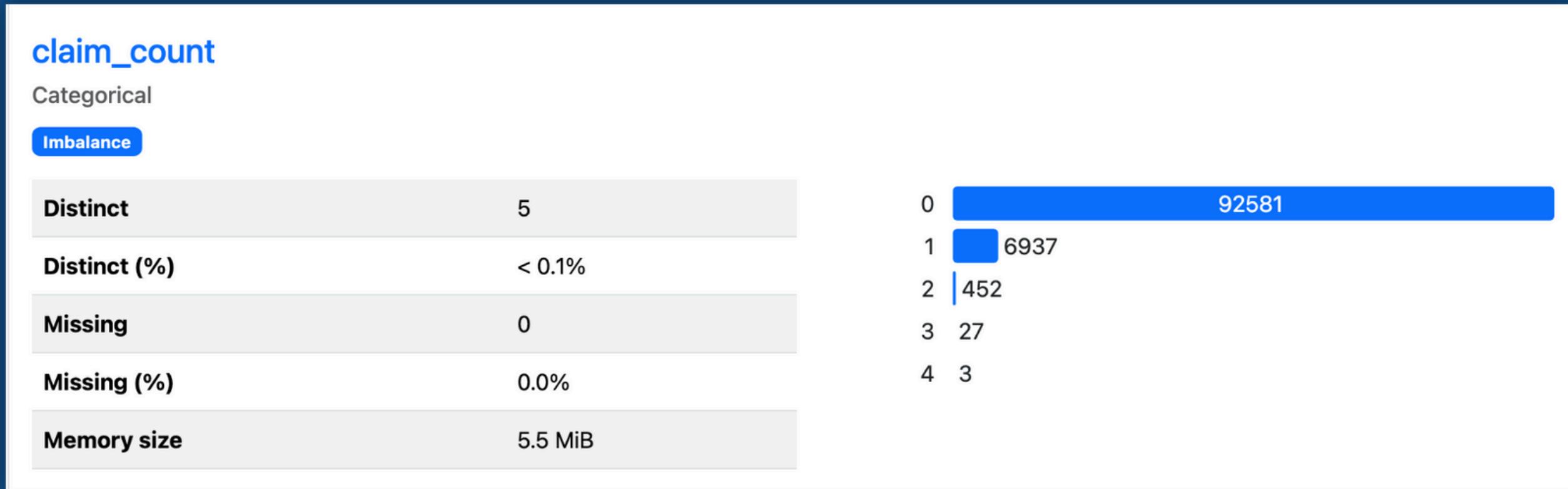
# Missing Value Plot



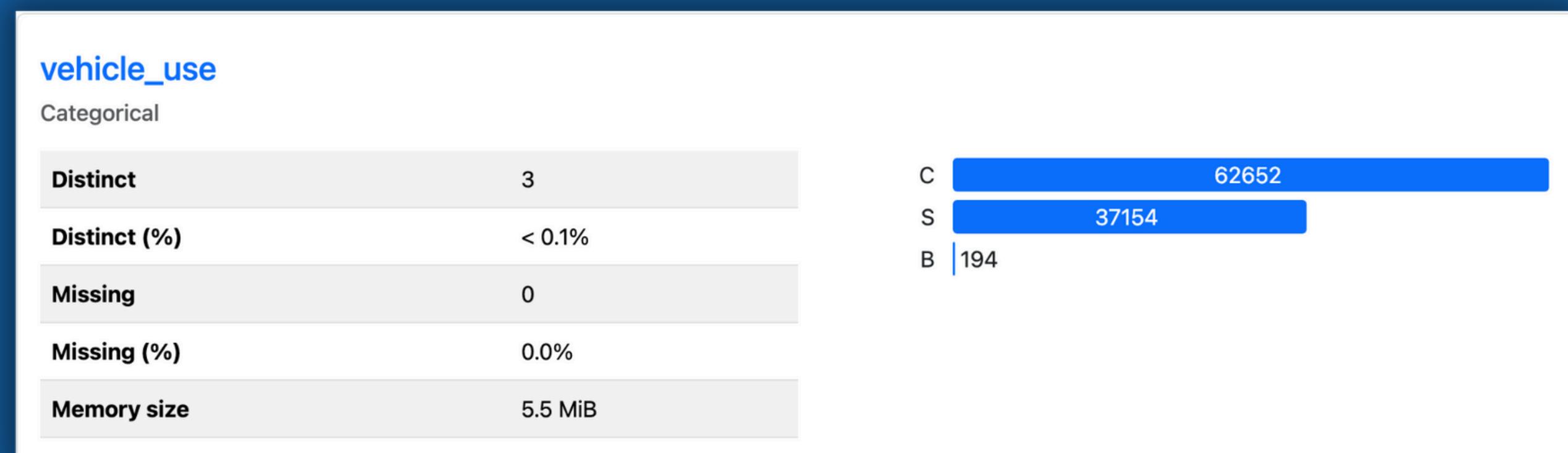
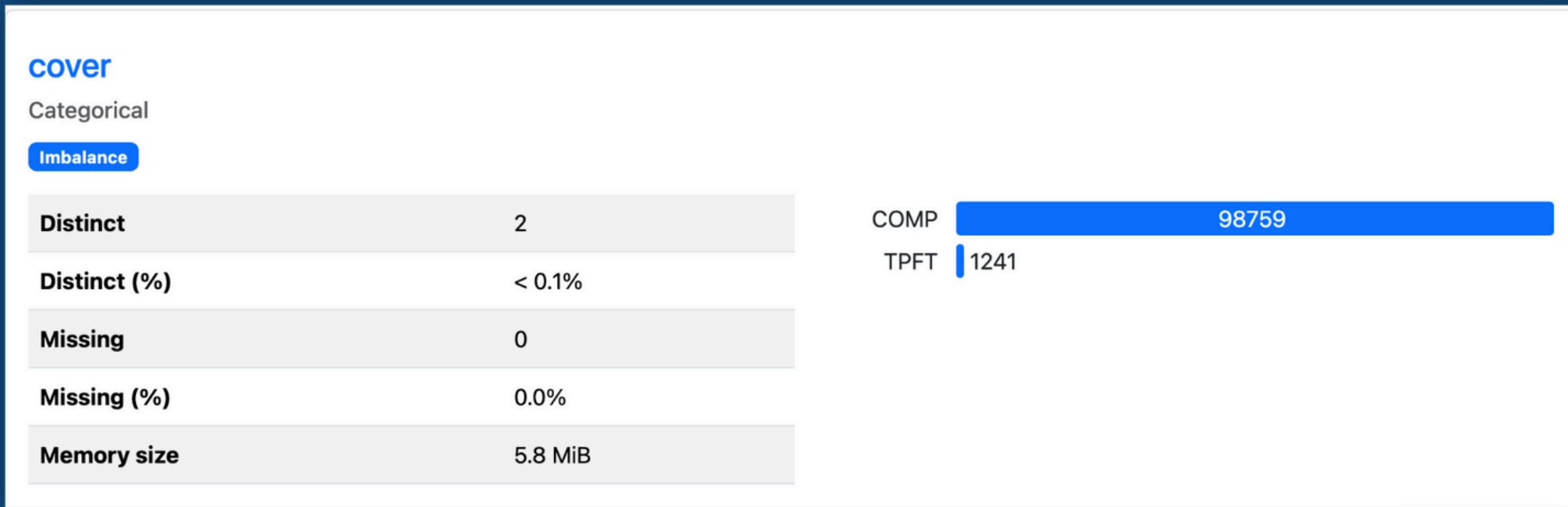
# Exposure & Net premium earned



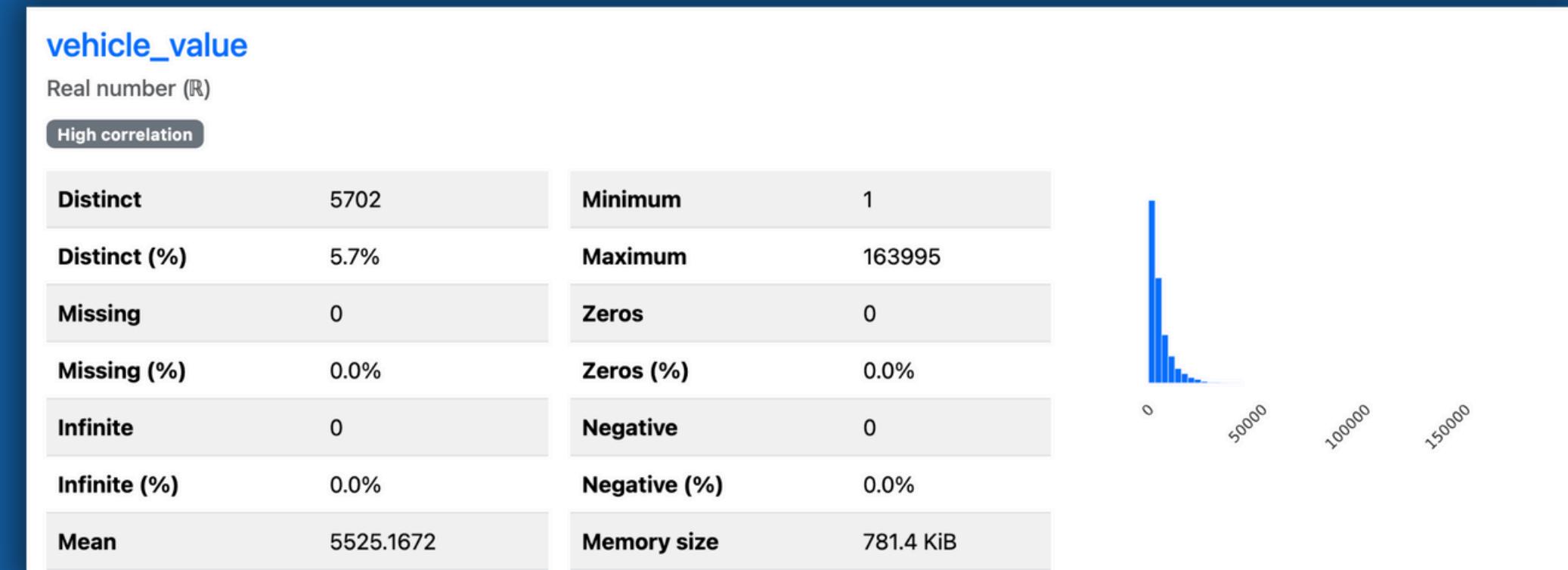
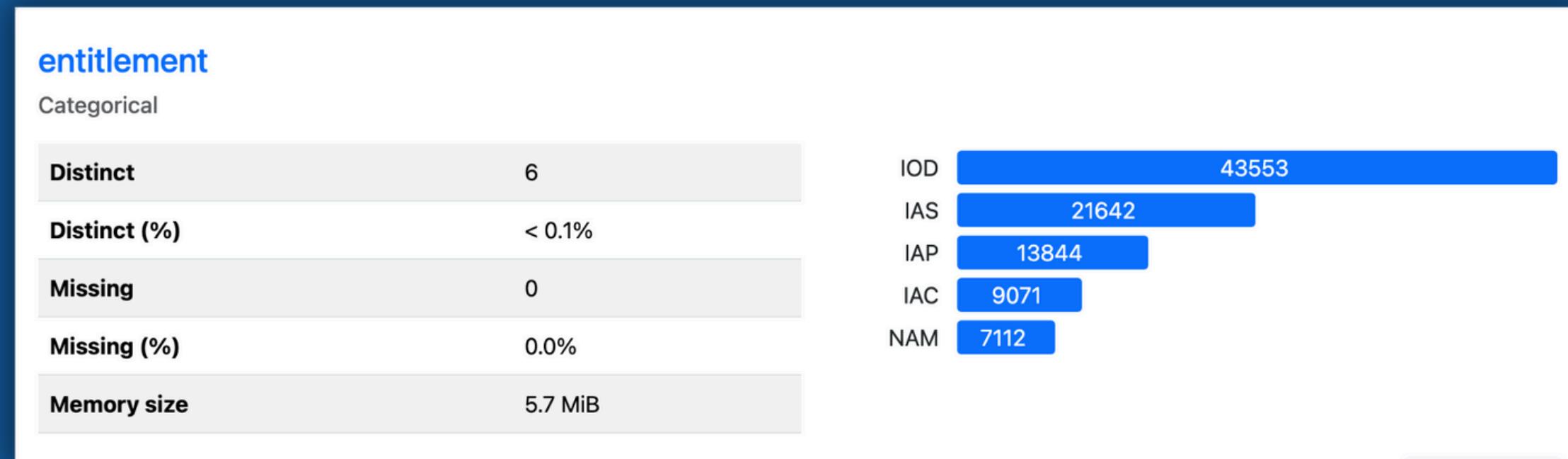
# Claim count & Incurred



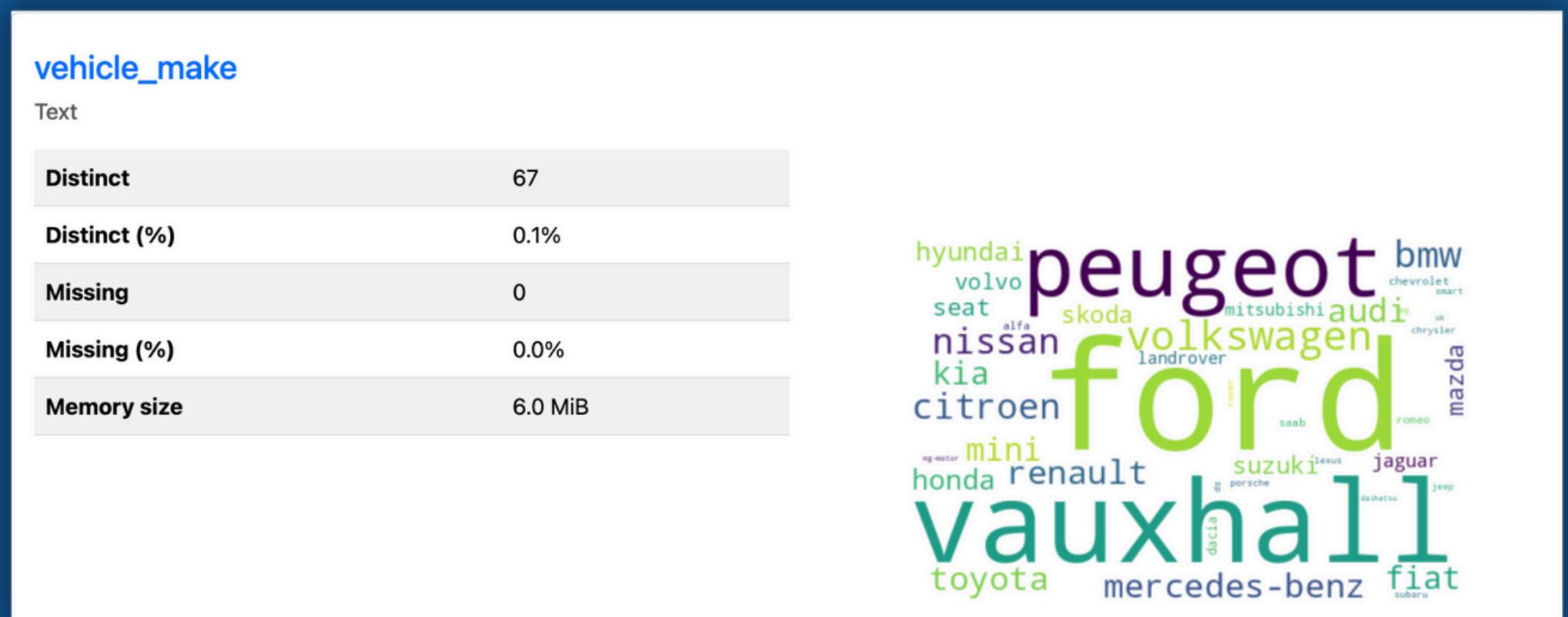
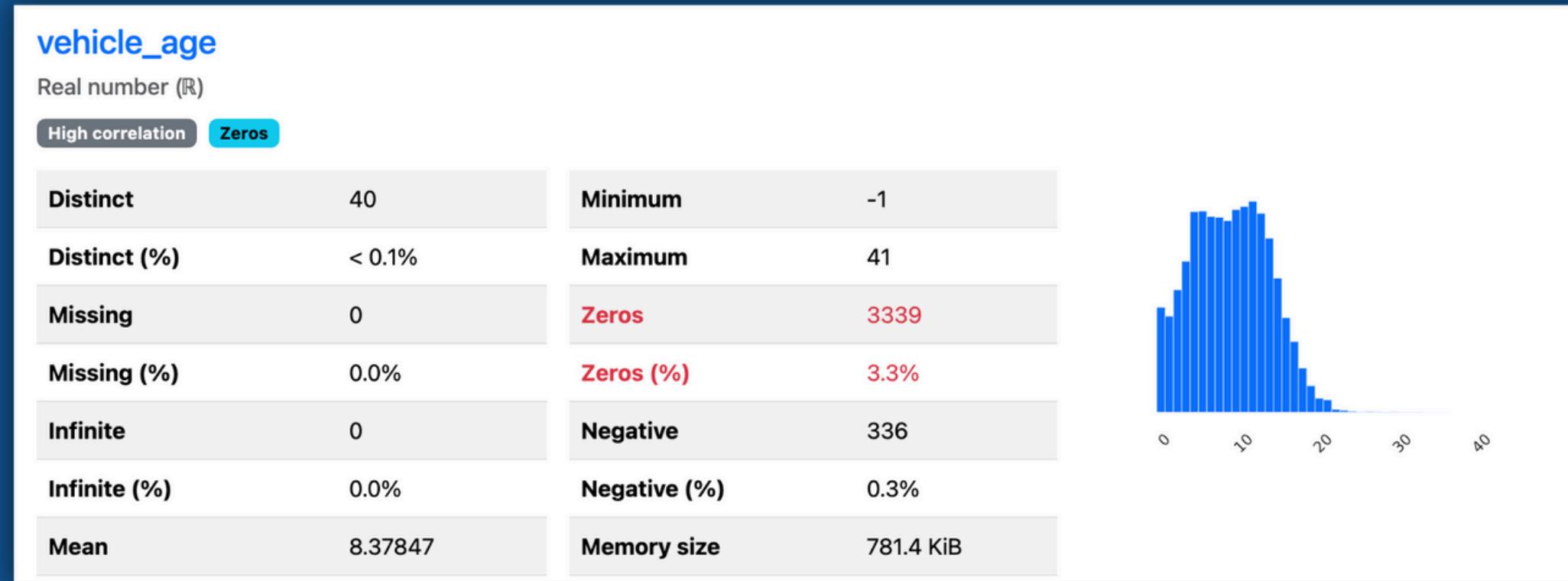
# Cover & Vehicle use



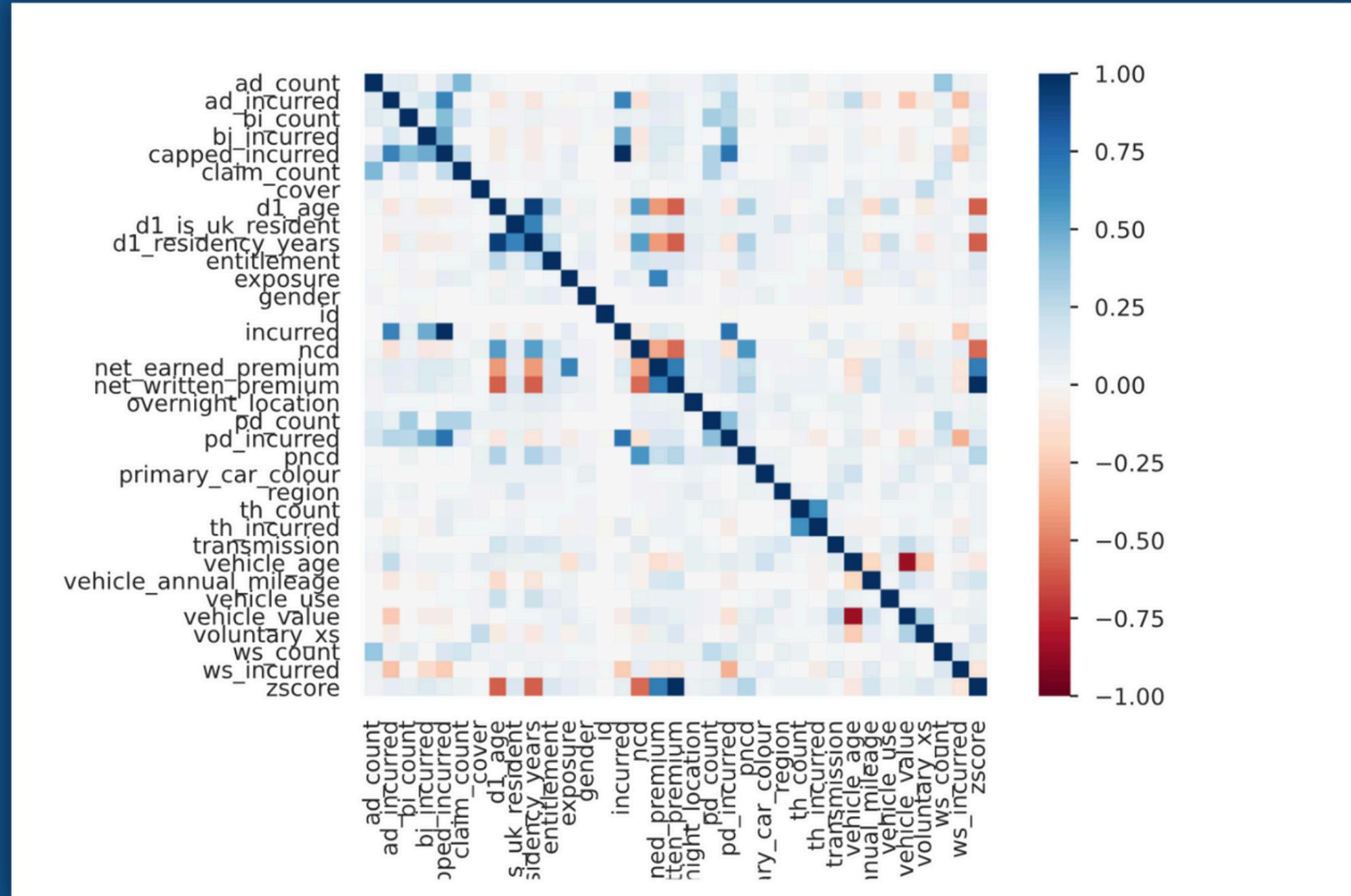
# Entitlement & Vehicle value



# Vehicle age & Vehicle make



# Correlation Matrix



# Data Preprocessing

## Cleaning

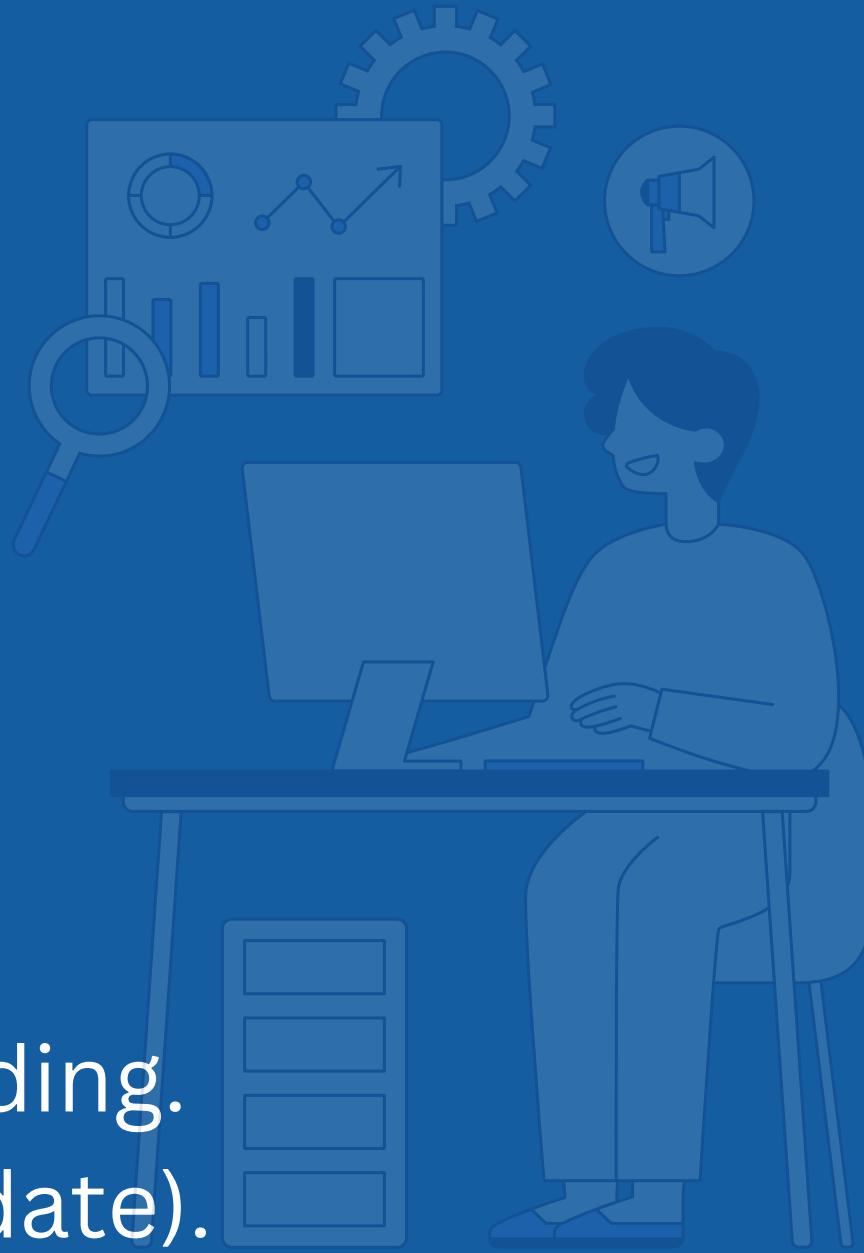
- 1. Addressed missing values:**
  - Categorical: Replaced with mode.
  - Numerical: Imputed with median.
- 2. Handled outliers (e.g., negative incurred values set to 0).**
- 3. Transformed data for consistency:**
  - Binned vehicle\_age.
  - Binned vehicle value.
- 4. Feature Engineering: Created few extra columns based on understanding of the dataset.**
  - policy duration
  - Claim severity
  - Driver experience

## Feature Engineering

# Feature Engineering

## New Features Added:

- Claim Severity = incurred / claim\_count
- Premium per Exposure = net\_earned\_premium / exposure
- Vehicle make categories transformed using one-hot encoding.
- Date features (e.g., policy\_duration extracted from start\_date).



# Top Predictive Features



## Key Insights from Analysis:

### Top Features:

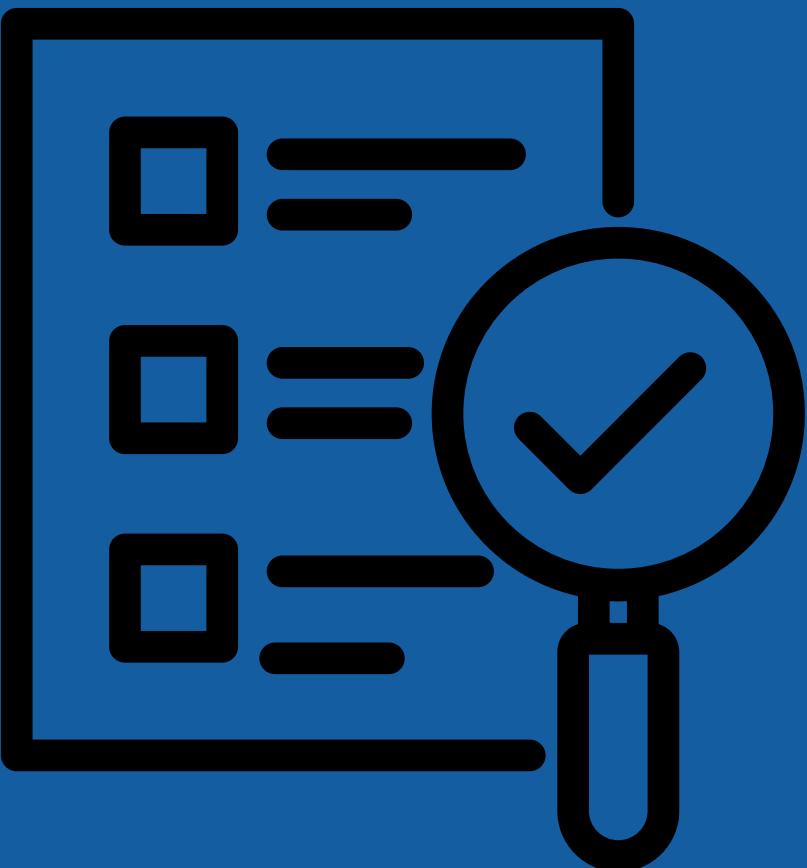
1. Vehicle Value
2. Vehicle Age Group
3. Vehicle Make



# Data Quality Checks

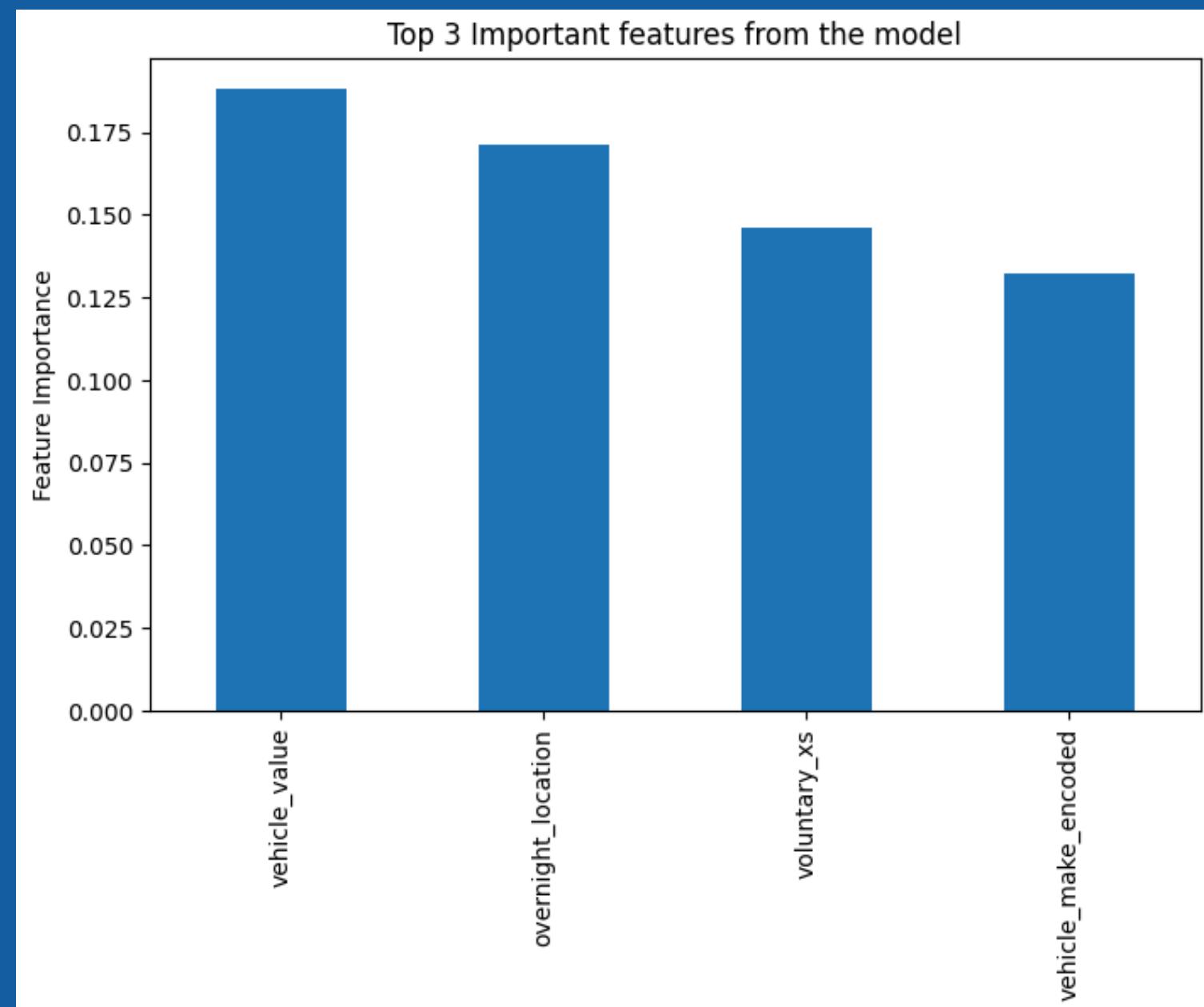
Implemented tests for:

- Missing values.
- Outliers.
- Duplicates.
- Automated flags for invalid data (e.g., negative incurred values).

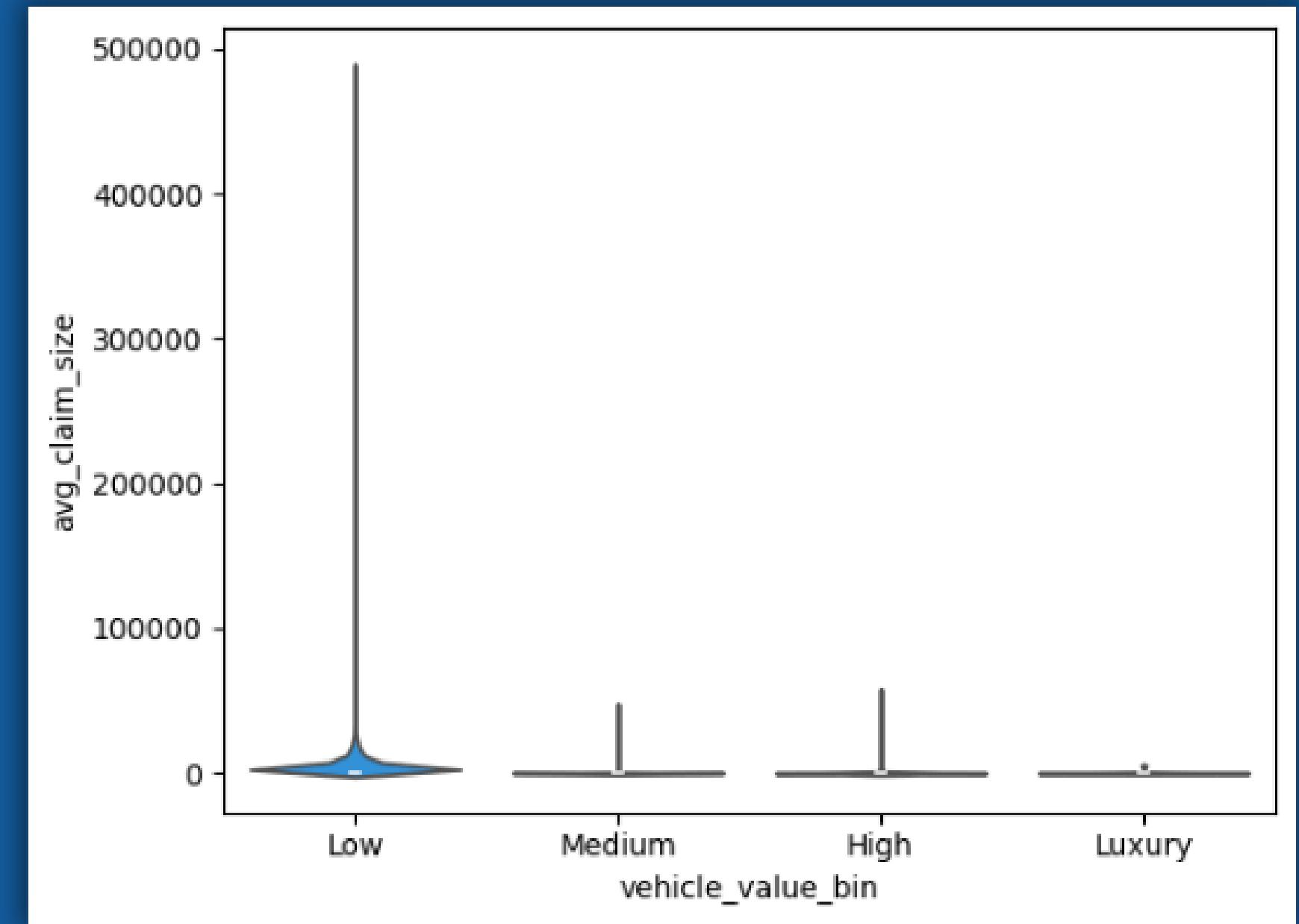
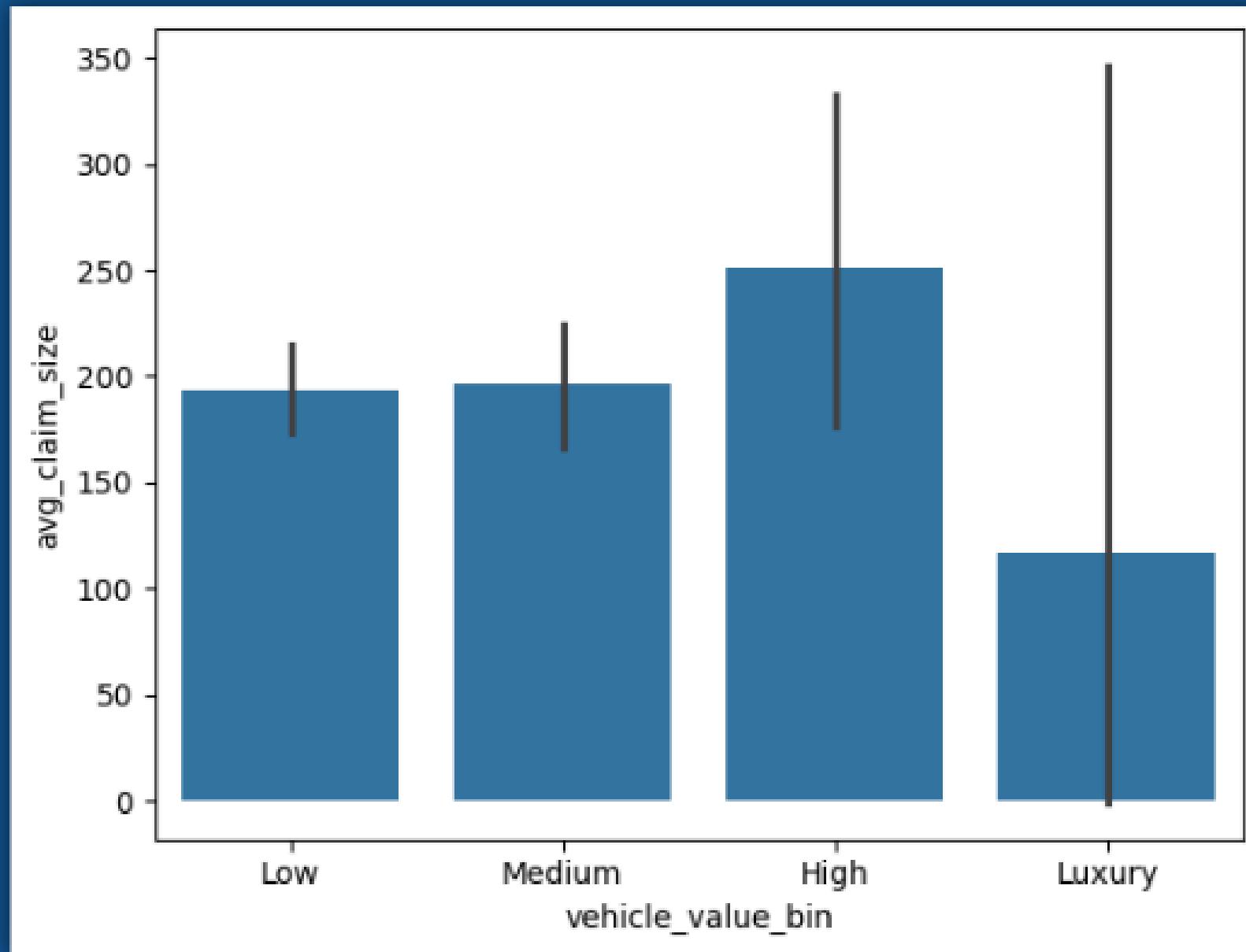


# Model Building

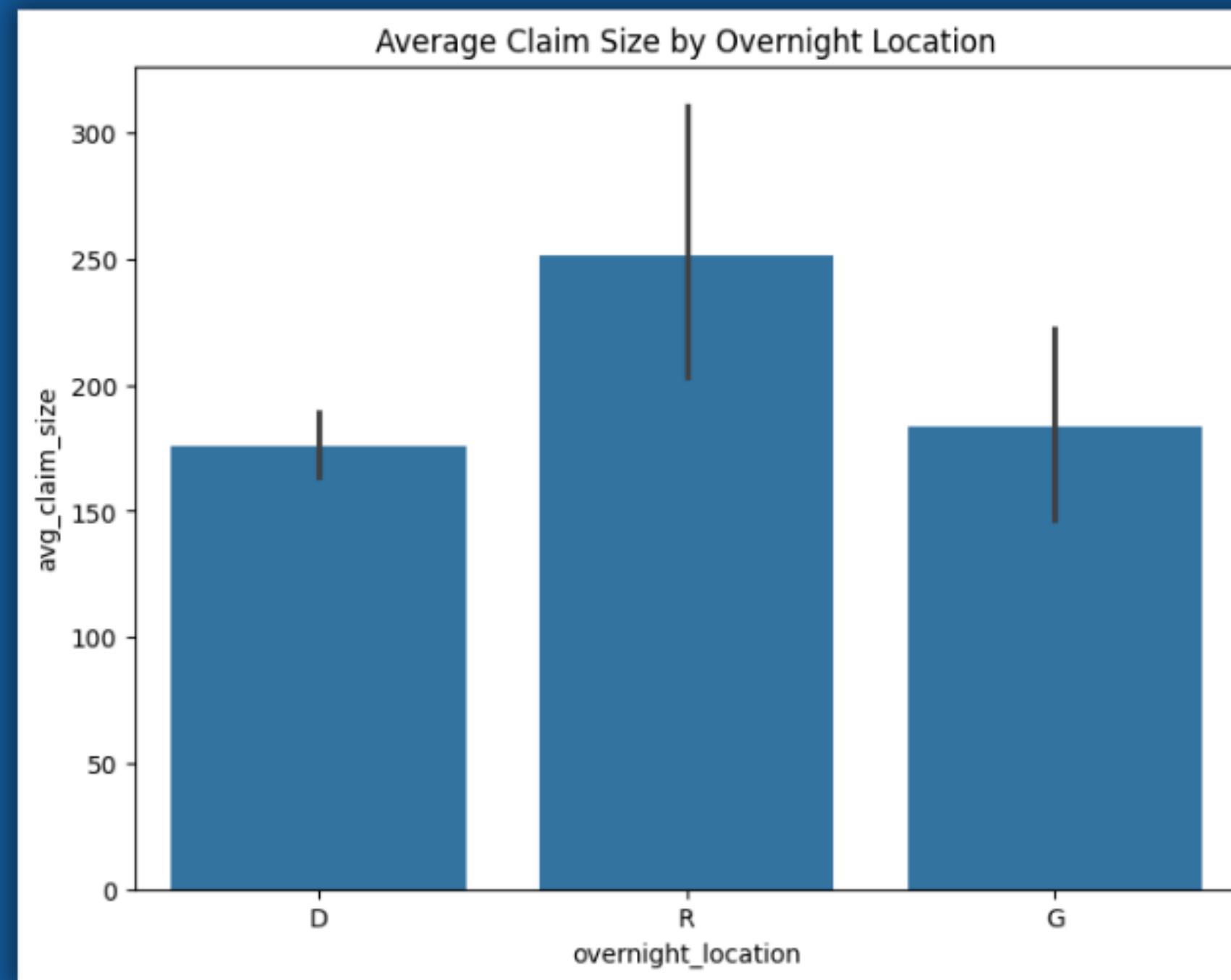
- Random Forest used for feature importance.
- Identified top features.



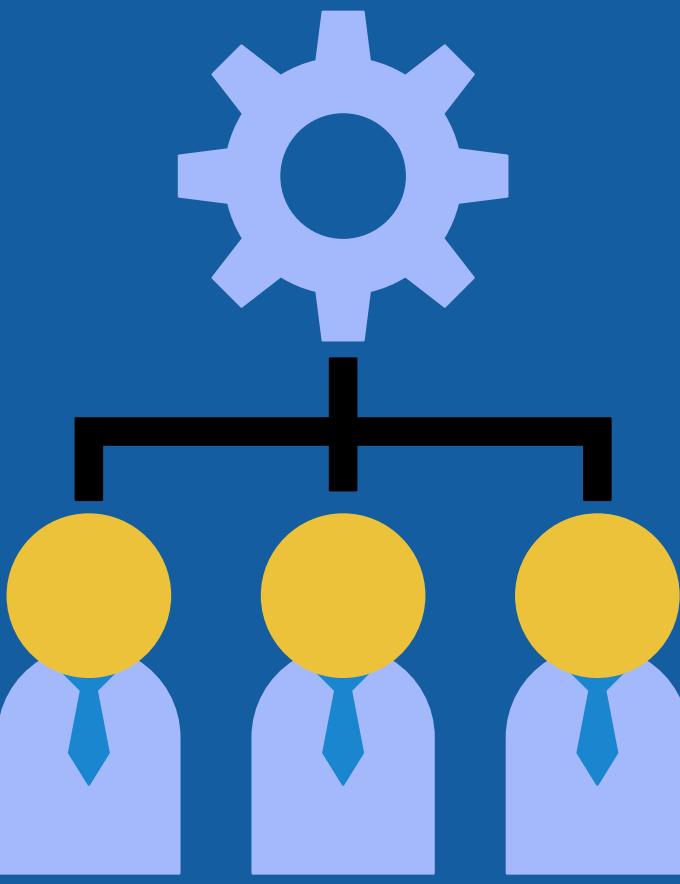
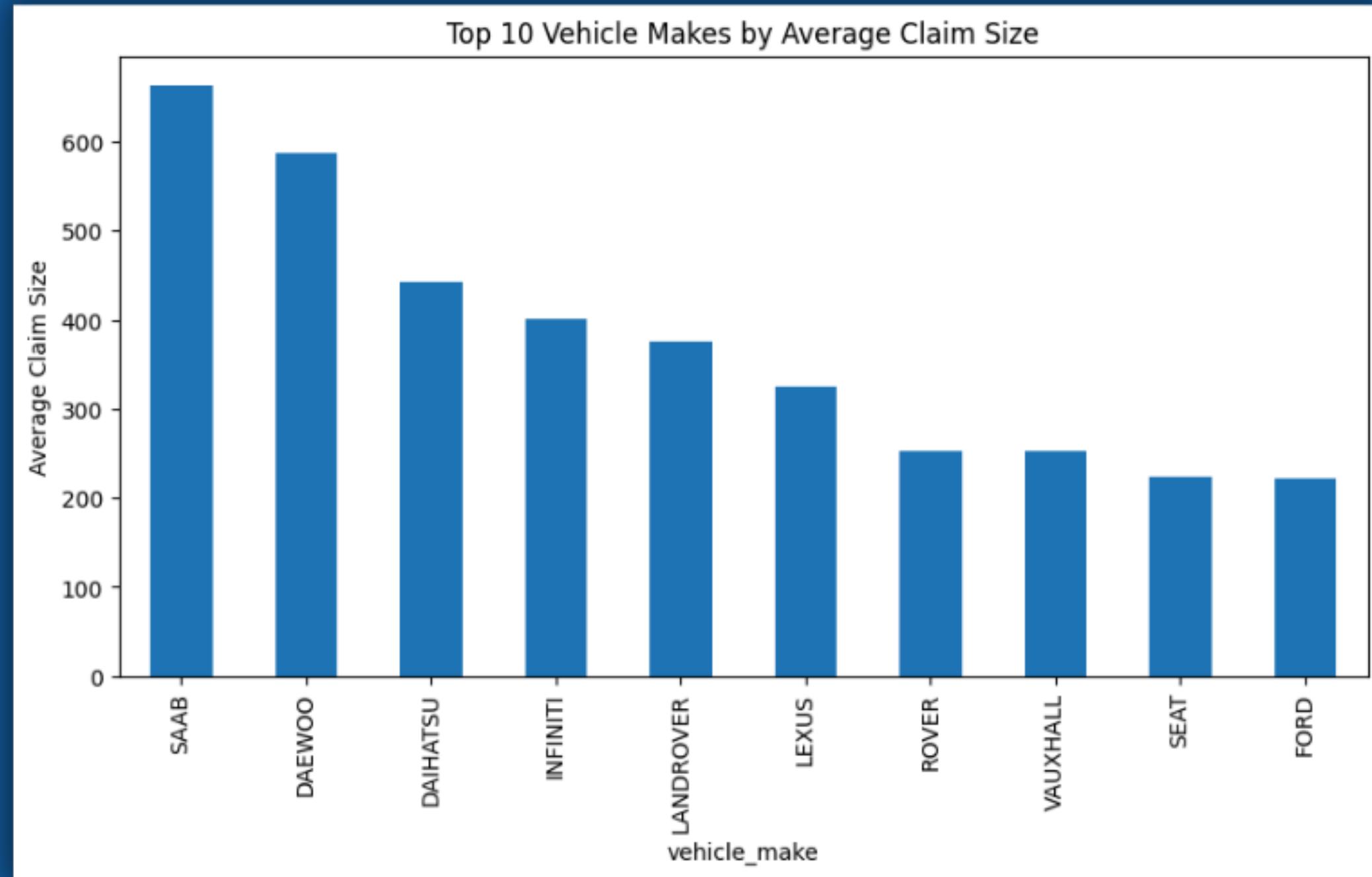
# Vehicle Value



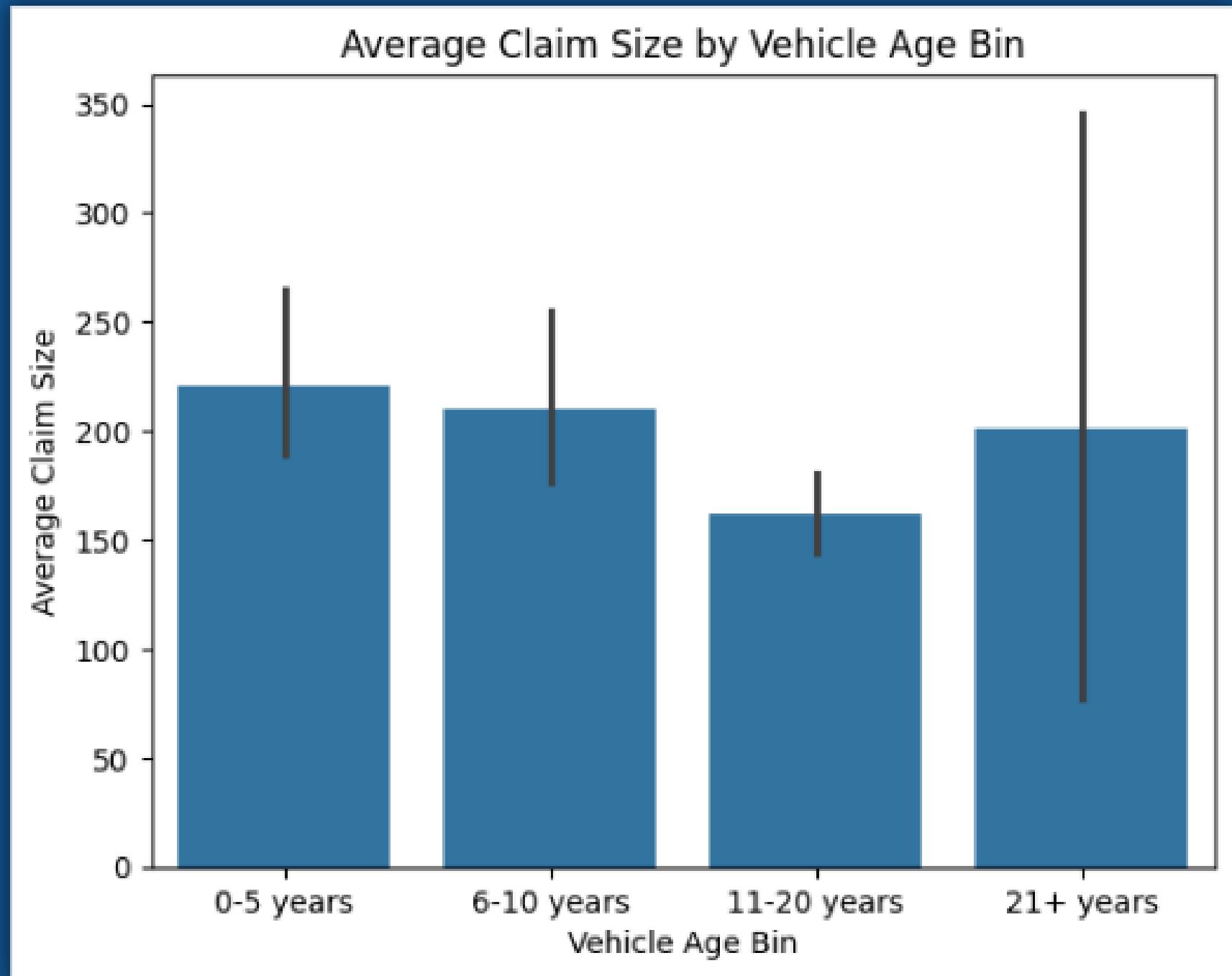
# Overnight Location



# Vehicle Make



# Vehicle Age



# Challenges

- Handling skewed data distributions.
- Dealing with high-cardinality categorical features.
- Balancing predictive accuracy with interpretability.

# Future Work

- Implement real-time data pipelines.
- Use advanced models (e.g., Gradient Boosting).
- Explore external data sources for enrichment.



# Conclusion

- Data transformation and feature engineering improved the data quality and hence in future will be used for model building purpose also.
- Top 3 feature were explored in details.