

Brief Overview of the Project

Context of the Project:

The objective of the project was to build a logistic regression model to predict whether a customer would file a claim on their car insurance based on various features. The dataset provided included multiple attributes of clients such as age, gender, driving experience, education, income, credit score, vehicle ownership, vehicle year, marital status, number of children, annual mileage, vehicle type, number of speeding violations, DUIs, and past accidents.

Project Instructions:

The main task was to identify the single feature that best predicts whether a customer will file a claim (the "outcome" column). This involved:

1. Reading and exploring the dataset to understand the data types, missing values, and distributions.
2. Filling any missing values appropriately.
3. Preparing the data for modeling by converting categorical variables to numeric and creating lists for storing models and results.
4. Building logistic regression models for each feature.
5. Measuring the performance of each model by calculating their accuracy.
6. Identifying the best-performing model with the highest accuracy.

How the Code Was Executed:

Step 1: Reading and Exploring the Dataset

The dataset was loaded into a pandas DataFrame. The data types, missing values, and distributions were examined to understand the initial state of the dataset.

Step 2: Filling Missing Values

Any missing values in the dataset were filled using the median to ensure no gaps in the data for modeling.

Step 3: Preparing for Modeling

Categorical variables were converted to numeric using LabelEncoder. This step ensured that all features were in a suitable format for modeling.

Step 4: Building and Storing the Models

Logistic regression models were built for each feature. Each feature was used as the predictor variable in separate models to identify its predictive power.

Step 5: Measuring Performance

The accuracy of each model was calculated by making predictions on the test set and comparing them with the actual outcomes. The accuracies were stored for comparison.

Step 6: Finding the Best Performing Model

The feature with the highest accuracy was identified as the best predictor for whether a customer would file a claim.

Detailed Explanation of the Solution

The solution part of the project involved identifying the best feature for predicting whether a customer will file a claim on their car insurance using logistic regression models. After processing the dataset and building models for each feature, the feature with the highest accuracy was identified.

Best Feature: Driving Experience

Best Accuracy: 0.7771

What Does This Mean?

1. Best Feature: Driving Experience

The feature identified as the best predictor for whether a customer will file a claim is "driving experience." This indicates that among all the features in the dataset, the number of years a client has been driving is the most significant factor in determining the likelihood of filing a claim.

2. Best Accuracy: 0.7771

The best accuracy achieved by the logistic regression model using the "driving experience" feature is 0.7771. This means that the model correctly predicted whether a customer would file a claim approximately 77.71% of the time based on their driving experience.

Implications of the Solution:

- **Significance of Driving Experience:** The result suggests that customers with different levels of driving experience have varying probabilities of filing a claim. For instance, newer drivers may be more likely to file claims compared to those with more years of driving experience.
- **Model Performance:** An accuracy of 77.71% indicates that the model is fairly reliable in predicting the outcome based on the driving experience feature alone. While not perfect, this level of accuracy is useful for risk assessment and decision-making processes in the insurance industry.
- **Risk Assessment:** Insurance companies can use this insight to refine their risk models. Understanding that driving experience is a key factor can help in tailoring insurance policies, pricing strategies, and marketing efforts to different customer segments.
- **Resource Allocation:** By focusing on the most predictive feature, insurers can allocate resources more efficiently. For example, they could develop targeted educational programs for newer drivers to reduce the likelihood of claims.

Conclusion:

The solution involved building logistic regression models for each feature and identifying "driving experience" as the most predictive feature with an accuracy of 0.7771. This result highlights the importance of driving experience in predicting insurance claims and provides valuable insights for improving risk assessment and resource allocation in the insurance industry.