

In this assignment, we focused on building a logistic regression model for lead scoring, with the goal of predicting the likelihood of lead conversion based on various features from the dataset. Below is a summary of the steps taken and the key learnings we gathered throughout the process.

1. Data Cleaning and Preprocessing

We began by loading the dataset and performing extensive data cleaning. We handled missing values by identifying columns with more than 50% missing data and removing them. Numerical columns with missing values were filled using the median, while categorical columns were filled using the mode. The 'Tags' column, which was deemed irrelevant, was also removed.

A specific challenge arose with the 'Lead Profile' column, where some rows contained the value 'Select,' which was treated as missing. We replaced these 'Select' values with 'Unknown' to maintain consistency. We then applied label encoding to all categorical columns, including 'Lead Profile,' to transform them into numeric values suitable for modeling.

2. Feature Engineering and Scaling

Next, we focused on scaling the numerical features. We removed irrelevant columns, such as 'Lead Number' and 'Converted,' before applying StandardScaler to normalize the numerical features. Scaling ensures that features with larger magnitudes do not dominate the model. After scaling, we split the data into training and testing sets to ensure the model could be validated properly.

3. Model Development and Evaluation

We trained a logistic regression model using the scaled data. The model was evaluated using various performance metrics: accuracy, precision, recall, F1 score, and ROC-AUC. These metrics gave us a comprehensive view of the model's performance, balancing the need for both false positives and false negatives in lead conversion prediction.

Through these evaluations, we learned that accuracy alone may not provide a complete picture of model performance, especially when dealing with imbalanced datasets. Precision, recall, and the ROC-AUC score were essential in understanding the model's ability to discriminate between converted and non-converted leads.

4. PCA for Dimensionality Reduction

To explore dimensionality reduction, we applied Principal Component Analysis (PCA). By retaining 95% of the variance, we reduced the number of features while preserving most of the information. This allowed us to visualize how the reduced data impacted the model's performance. We trained a logistic regression model on the PCA-transformed data and achieved similar performance to the model trained on the original features. PCA helped reduce the complexity of the model and may be beneficial for future analyses involving more features.

5. Confusion Matrix and Feature Importance

We visualized the confusion matrix to assess the types of errors our model was making (false positives and false negatives). This helped us understand where the model might need improvement.

Additionally, we analyzed the coefficients of the logistic regression model to identify which features were most influential in predicting lead conversion. Features like 'Total Time Spent on Website' and 'Lead Quality' emerged as having significant coefficients, highlighting their importance in the conversion process.

6. Key Learnings

From this assignment, we learned several important lessons:

- **Data Cleaning:** Handling missing data and preprocessing steps like encoding categorical variables and scaling numerical features are critical for building effective models.
- **Feature Selection:** Feature engineering plays a significant role in model performance. Identifying and removing irrelevant or redundant features can improve accuracy and reduce model complexity.
- **Model Evaluation:** Multiple evaluation metrics, including precision, recall, and ROC-AUC, provide a more balanced view of model performance, especially when working with imbalanced datasets.
- **Dimensionality Reduction:** PCA can be a useful technique to reduce the number of features, which simplifies the model while retaining most of the variance.

In conclusion, we successfully built a logistic regression model, applied dimensionality reduction, and evaluated its performance. The process underscored the importance of data cleaning, feature selection, and the use of multiple evaluation metrics to assess model performance comprehensively.