

$\frac{1}{4}$

Dataset

Source: The dataset is obtained from the City of Chicago's Traffic Crashes dataset.

Description: The dataset includes information about vehicle crashes, such as crash location, time, weather conditions, and the primary and secondary contributory causes.

Size: Approximately X million records.

Features:

- CRASH_DATE
- TRAFFIC_CONTROL_DEVICE
- WEATHER_CONDITION
- LIGHTING_CONDITION
- ROADWAY_SURFACE_COND
- PRIM_CONTRIBUTORY_CAUSE
- MOST_SEVERE_INJURY

Additional engineered features like Is_Weekend, Speed_Weather_Interaction, etc.

Installation

To run this project locally, follow these steps:

Clone the repository:

```
'git clone https://github.com/yourusername/Traffic-Crash-Analysis.git,'
```



Navigate to the project directory:

```
'cd Traffic-Crash-Analysis'
```



Create a virtual environment:

```
'python -m venv venv'
```



Activate the virtual environment:

```
'venv\Scripts\activate'
```



On macOS/Linux:

```
'source venv/bin/activate'
```



Install the required packages:

```
'pip install -r requirements.txt'
```



Usage

Preprocessing and Feature Engineering:

Run preprocessing.ipynb to clean and prepare the data.

Key steps include handling missing values, feature engineering, and encoding categorical variables. Model Training:

Run model_training.ipynb to train various classification models.

Includes Logistic Regression, Ridge Classifier, Decision Trees, Random Forest, and Gradient Boosting. Model Tuning and Evaluation:

Use `model_tuning.ipynb` to perform hyperparameter tuning using GridSearchCV and RandomizedSearchCV.

Evaluate models using metrics such as accuracy, precision, recall, and F1-score.

Final Testing and Results:

Run `final_testing.ipynb` to test the best-performing model on the test set.

Generate performance reports and confusion matrices. Visualization:

Use `visualization.ipynb` for exploratory data analysis (EDA) and to create plots like heatmaps, pair plots, and bar charts comparing model performances.

Modeling Process

Baseline Model: Logistic Regression was initially used as the baseline.

Feature Selection: Features were selected based on correlation and domain knowledge.

Model Comparison: Multiple models were trained, including Logistic Regression, Ridge Classifier, Random Forest, and Gradient Boosting.

Hyperparameter Tuning: RandomizedSearchCV was used for tuning hyperparameters to improve model performance.

Final Model: Gradient Boosting was selected as the best-performing model based on accuracy and other evaluation metrics.

Results

Best Model: Gradient Boosting

Accuracy: 1.00%

Key Features: The most important features contributing to the model's predictions were `Weather_Condition_Mode`, `Speed_Weather_Interaction`, and `Is_Weekend`.

Insights and Recommendations

Traffic Management: The model suggests that weather conditions and speed limits are significant contributors to crashes.

Implementing stricter speed regulations during adverse weather could reduce accidents.

Policy Changes: Focus on improving road conditions and traffic control devices in areas with high accident rates.

Further Analysis: Recommend analyzing the impact of time of day and specific locations (using latitude and longitude) on crash rates for targeted interventions.

Contributing

Contributions are welcome! Please create a pull request or open an issue for any suggestions or improvements.

License

This project is licensed under the GPL-2.0 License.

Acknowledgements

Special thanks to the City of Chicago for providing the open dataset used in this analysis.

Tools used include Python, Colab Notebook, Jupyter Notebook, Pandas, Scikit-learn, Matplotlib, Seaborn.

For visualizations and coding procedures, feel free to open the notebook and presentation files above

Gladwell Chepkorir.



Releases

No releases published
[Create a new release](#)

Packages

No packages published
[Publish your first package](#)

Languages

