# Clustering vehicle crash data

## <u>Problem Description :</u>

**Background:** Road safety is a crucial concern worldwide, with traffic accidents resulting in significant human and economic losses. Understanding the factors that contribute to accidents and their severity is essential for implementing effective safety measures. In this project, we aim to cluster vehicle crash data to identify patterns and factors associated with different levels of accident severity. By doing so, we can gain insights into the key risk factors and develop targeted strategies to reduce road accidents and their impact.

**Dataset Information:** The dataset used in this project contains information about vehicle crashes and their severity. It includes various attributes related to the accident, the involved vehicles, and the casualties. Some of the key features in the dataset are:

1. **Number_of_vehicles_involved:** The number of vehicles involved in the accident.
2. **Number_of_casualties:** The number of casualties in the accident.
3. **Road_type:** The type of road where the accident occurred (e.g., highway, local road).
4. **Speed_limit:** The speed limit on the road at the accident site.
5. Weather_conditions: The weather conditions during the accident (e.g., clear, rainy, snowy).
6. **Light_conditions:** The light conditions at the time of the accident (e.g., daylight, darkness with street lighting).
7. **Accident_severity:** The severity of the accident (the target variable).
8. **Problem Description:** The main objective of this project is to cluster the vehicle crash data into groups based on the patterns and characteristics of accidents. The clustering analysis will help identify common factors associated with different levels of accident severity. This will provide valuable insights to policymakers,

road safety authorities, and insurance companies to implement targeted interventions and preventive measures.

**Approach: To achieve the project's goal, we will follow the following steps:**

1. **Data Loading and Exploration:**
- Load the vehicle crash dataset and explore its structure, size, and data types.
- Analyze the distribution of accident severity levels to understand the class imbalance.

2. **Data Preprocessing and Feature Engineering:**
- Handle missing values and outliers appropriately to ensure data quality.
- Engineer new features if necessary to capture relevant information that can improve clustering.

3. **Exploratory Data Analysis (EDA):**
- Perform exploratory data analysis to visualize the distribution of accidents, causalities, and vehicles involved.
- Examine the relationships between various factors like weather conditions, road types, and light conditions with accident severity.

4. **Feature Selection and Dimensionality Reduction:**
- Use feature selection techniques to identify the most relevant attributes for clustering.
- Apply dimensionality reduction methods to reduce the feature space while preserving key information.

5. **Handling Imbalanced Data:**
- Address the class imbalance issue, which is common in accident severity datasets, using techniques like oversampling or undersampling.

6. **Clustering Algorithms:**
- Implement different clustering algorithms such as K-Means, DBSCAN, and hierarchical clustering.
- Evaluate the performance of each algorithm using appropriate metrics.

7. **Model Evaluation and Interpretation:**
- Evaluate the clustering models and choose the best-performing one based on validation metrics.
- Interpret the clusters and identify the key factors contributing to different accident severity levels.

8. **Conclusion and Recommendations:**

- Summarize the findings from the clustering analysis.
- Provide actionable recommendations to improve road safety based on the identified patterns and risk factors.

# Possible Framework :

## I. Introduction

- Overview of the project's goal and objective.
- Brief explanation of road safety importance and the need for analyzing crash data.

## II. Data Loading and Exploration

- Import necessary libraries.
- Load the vehicle crash dataset.
- Explore the dataset's structure, size, and data types.
- Analyze the distribution of accident severity levels.

## III. Data Preprocessing and Feature Engineering

- Handle missing values and outliers in the data.
- Engineer new features if required.
- Convert categorical features to numerical using encoding techniques.

## IV. Exploratory Data Analysis (EDA)

- Visualize the distribution of accidents, causalities, and vehicles involved.
- Examine the relationship between various factors and accident severity.
- Use graphs and plots for clear visualization.

## V. Feature Selection and Dimensionality Reduction

- Select relevant features for clustering using techniques like RFE or LASSO regression.
- Implement dimensionality reduction methods like PCA or t-SNE if needed.

## VI. Handling Imbalanced Data

- Address the class imbalance issue using oversampling or undersampling.

## VII. Clustering Algorithms

- Implement K-Means, DBSCAN, and hierarchical clustering algorithms.
- Evaluate each algorithm's performance using appropriate metrics.

**VIII. Model Evaluation and Interpretation**

- Compare the clustering results from different algorithms.
- Choose the best-performing model based on validation metrics.
- Interpret the clusters and identify key factors contributing to different accident severity levels.

**IX. Conclusion and Recommendations**

- Summarize the findings from the clustering analysis.
- Provide actionable recommendations to improve road safety based on the identified patterns and risk factors.

**X. Future Work**

- Suggest potential future work to enhance the analysis and gain deeper insights.
- Discuss areas for further research and improvements in the clustering process.

**XI. References**

- List any external sources or datasets used for the project.

# Code Explanation :

**\*If this section is empty, the explanation is provided in the .ipynb file itself.

The provided code is aimed at clustering vehicle crash data for safety analysis. It performs various steps, such as data loading, data preprocessing, data visualization, and model creation using different algorithms like K-Nearest Neighbors (KNN), Naive Bayes, Support Vector Machine (SVM), Decision Tree, and Random Forest. Let's break down the code and understand each part:

1. **Data Loading and Exploration:**
   - The code starts by importing necessary libraries like NumPy, Pandas, Seaborn, Matplotlib, and warnings.
   - It reads the vehicle crash data from a CSV file using the Pandas read_csv() function and stores it in a DataFrame called df.
   - It prints the shape and size of the dataset to get an idea of the number of rows and columns.
   - It uses describe() to get the statistical summary of the numerical columns in the dataset.
   - The code also checks the data types of each column using the info() function.

2. **Exploratory Data Analysis (EDA):**
   - The code checks for duplicate values in the dataset using duplicated().sum().
   - It displays the distribution of accident severity levels using value_counts() and a count plot.
   - The focus is on understanding the class imbalance, i.e., whether one severity level is dominating the dataset.

3. **Handling Missing Values:**
   - The code identifies columns with a significant number of missing values (>2500) and drops them from the dataset using drop().
   - For categorical columns, missing values are filled with the mode (most frequent value) of each column using a for loop.

4. **Data Visualization:**
   - The code uses Seaborn to create scatter plots and joint plots to visualize the relationship between the number of casualties and the number of vehicles involved in accidents.

- It generates a heatmap to visualize the correlation between numerical columns in the dataset.

5. **Feature Selection and Dimensionality Reduction:**

- The code encodes categorical columns using Label Encoding to convert them into numerical format for machine learning models.
- It performs a chi-square test (chi2) to select the most relevant features for clustering.

6. **Clustering Algorithms:**

- The code uses multiple clustering algorithms like KNN, Naive Bayes, SVM, Decision Tree, and Random Forest to create models for accident severity prediction.
- It splits the dataset into training and testing sets using train_test_split() from Scikit-learn.
- Each model is trained using the training data and then used to predict accident severity on the test data.

7. **Model Evaluation and Interpretation:**

- The code evaluates the performance of each model using confusion matrices, accuracy scores, and classification reports.
- The models are compared based on their accuracy scores to identify the best-performing one.

# Future Work :

**Step 1: Data Collection and Preparation**

- Obtain more extensive and updated vehicle crash data from various sources to improve the model's accuracy and generalizability.
- Perform rigorous data cleaning, handling missing values, and dealing with outliers more effectively to enhance data quality.

**Step 2: Feature Engineering**

- Experiment with feature engineering techniques to create more meaningful features from the existing ones.
- Consider incorporating external data, such as road conditions, traffic patterns, and vehicle safety features, to enrich the dataset and improve model performance.

**Step 3: Exploratory Data Analysis (EDA)**

- Conduct a more in-depth exploratory data analysis to gain deeper insights into the relationships between different factors and accident severity.
- Use interactive visualizations and geographical mapping to understand the spatial distribution of accidents and severity levels.

**Step 4: Advanced Clustering Algorithms**

- Implement advanced clustering algorithms like Density-Based Spatial Clustering of Applications with Noise (DBSCAN), Gaussian Mixture Models (GMM), and Spectral Clustering to compare their performance with the existing models.

**Step 5: Handling Imbalanced Data**

- Explore advanced techniques for handling imbalanced data, such as Synthetic Minority Over-sampling Technique (SMOTE), Adaptive Synthetic (ADASYN), and class-weighted learning.

**Step 6: Model Tuning and Ensemble Methods**

- Fine-tune hyperparameters of the best-performing model using techniques like GridSearchCV or RandomizedSearchCV for better accuracy.

- Consider using ensemble methods like Voting Classifier or Stacking to combine multiple models and improve prediction accuracy.

**Step 7: Interpretability and Explainability**

- Implement techniques like LIME or SHAP (SHapley Additive exPlanations) to explain the model's predictions and understand the factors contributing to different severity levels.
- Generate easy-to-understand visual explanations to aid decision-making.

**Step 8: Deployment and Real-Time Prediction**

- Deploy the final model as a web application or API for real-time accident severity prediction.
- Implement a user-friendly interface to input accident details and receive immediate predictions.

**Step 9: Regular Updates and Maintenance**

- Regularly update the model with new data to ensure its accuracy and relevance.
- Monitor the model's performance and retrain it periodically to account for any changes in traffic patterns or road conditions.

**Step-By-Step Implementation Guide**

1. Gather vehicle crash data from reliable sources in CSV format.
2. Load the data into a Pandas DataFrame and perform basic data exploration to understand its structure.
3. Preprocess the data by handling missing values, encoding categorical variables, and removing irrelevant columns.
4. Conduct exploratory data analysis (EDA) using Seaborn and Matplotlib to visualize data distributions and relationships.
5. Use chi-square tests or other feature selection techniques to identify essential features for clustering.
6. Split the dataset into training and testing sets using Scikit-learn's train_test_split.
7. Implement different clustering algorithms like KNN, Naive Bayes, SVM, Decision Tree, and Random Forest.
8. Train each model using the training data and evaluate its performance using confusion matrices, accuracy scores, and classification reports.

9. Compare the models based on accuracy and choose the best-performing one.
10. Conduct further analysis to interpret the clusters and identify key factors contributing to different accident severity levels.
11. Prepare a comprehensive report summarizing the findings and actionable recommendations for road safety improvement.
12. For future work, follow the outlined steps to collect more extensive and updated data, perform feature engineering, and implement advanced clustering algorithms and ensemble methods.
13. Deploy the final model as a web application or API for real-time accident severity prediction.
14. Monitor the model's performance, update it regularly with new data, and maintain its accuracy over time.

# Concept Explanation :

**Explanation of the Clustering Algorithm (K-Nearest Neighbors - KNN)**

Hey there, my curious friend! Let me tell you about the magical world of K-Nearest Neighbors, also known as KNN! Imagine you have a super cool secret spy gadget, and your mission is to figure out which neighborhood in a city is the safest.

**Step 1:** The Secret Spy Gadget Your secret spy gadget, KNN, is a straightforward algorithm that classifies things based on their neighbors. Yes, just like how you decide to wear a fancy hat because your best friend wears one too!

**Step 2:** Finding the Neighbors So, the first thing KNN does is find the "k" closest neighbors to the thing it wants to classify. Let's say you want to know if a neighborhood is safe or not. KNN will check which k neighborhoods are most similar to the one you want to classify.

**Step 3:** The Voting Game Now, here comes the fun part! KNN gathers all the k neighbors and asks them, "Hey, neighbors, are you safe or not?" Each neighbor raises their hand, shouting "Safe!" or "Not Safe!" KNN counts the votes to see which one has more. Whichever has the most votes wins! It's like deciding where to have lunch with your friends by going where the majority wants to go.

**Step 4:** Making the Decision Once all the neighbors have voted, KNN announces the result. "Hey, everyone! Our super spy gadget says this neighborhood is safe because most of its neighbors are safe!" It's like a neighborhood popularity contest!

**Step 5:** Be Cautious! But, hey, my friend, you must be careful! Choosing the right "k" is vital. If you pick too few neighbors, you might end up with some noisy or mischievous neighbors who give wrong answers. And if you choose too many neighbors, your gadget might get confused because of conflicting opinions. So, finding the right "k" value is like finding the perfect balance in your spy gadget!

**Example:** Imagine you have a map of different neighborhoods in a city. Each neighborhood has attributes like the number of schools, crime rates, parks, and shops. You want to use KNN to classify whether a new neighborhood is safe or not based on these attributes. KNN will look at the k closest neighborhoods and ask them if they are

safe or not. If most of them say "safe," then KNN will classify the new neighborhood as safe too.

So, my dear friend, K-Nearest Neighbors is a friendly algorithm that relies on the support of its neighbors to make decisions. It's like having a group of friends who help you decide where to go and what to do. Just be mindful of picking the right friends (neighbors) and finding the right balance (k value) for your secret spy gadget to work its magic!

# Exercise Questions :

**1. How would you explain the purpose of this vehicle crash data analysis project to a non-technical audience?**

**Answer:** The purpose of this project is to use data from vehicle crash incidents to understand patterns and trends related to accident severity. By analyzing various factors like the number of casualties, types of vehicles involved, road conditions, and weather, we aim to identify potential risk factors and improve safety measures on the roads.

**2. What are the steps involved in Exploratory Data Analysis (EDA) for this project?**

**Answer:** The steps in EDA include data cleaning to handle missing values and duplicates, exploring the distribution of accident severity levels, visualizing numerical statistics and correlations, and analyzing categorical features through count plots and cross-tabulations.

**3. Explain the concept of handling missing values in the dataset. How did you decide to fill the missing values for categorical features in this project?**

**Answer:** Handling missing values involves deciding how to replace or impute the missing data in the dataset. For categorical features, we often use the mode (most frequent value) to fill the missing values because it is a common and simple approach.

**4. What is the significance of using the K-Nearest Neighbors (KNN) algorithm in this project? How does KNN classify vehicle crash incidents based on their neighbors?**

**Answer:** KNN is used in this project to classify the severity of vehicle crash incidents based on the similarity of their features to the neighboring incidents. It finds the k closest incidents to the one we want to classify and uses majority voting to determine the incident's severity based on the severity of its neighbors.

**5. How did you handle class imbalance in the dataset? Why is it important to address class imbalance before building a classification model?**

**Answer:** Class imbalance occurs when certain classes have significantly more or fewer instances than others. In this project, we addressed class imbalance using the Synthetic Minority Over-sampling Technique (SMOTE), which creates synthetic samples of the minority class to balance the dataset. Addressing class imbalance is crucial because an imbalanced dataset can lead to biased and inaccurate classification results.

**6. What are some potential evaluation metrics used to assess the performance of classification models in this project? Which metric would you prioritize for this vehicle crash analysis?**

**Answer:** Some evaluation metrics include accuracy, precision, recall, F1-score, and the confusion matrix. For this vehicle crash analysis, we would prioritize accuracy as it gives an overall measure of the model's correctness in classifying incidents based on their severity.

**7. Explain the difference between Decision Tree and Random Forest algorithms used in this project. Why did you choose Random Forest over Decision Tree for classification?**

**Answer:** Decision Tree is a single decision-making tree, while Random Forest is an ensemble of multiple Decision Trees. Random Forest combines the predictions of multiple trees to improve accuracy and reduce overfitting. We chose Random Forest because it usually outperforms Decision Tree in complex datasets and provides a more robust classification.

**8. Can you explain the concept of feature importance in Random Forest? How did you use it to identify important features in this project?**

**Answer:** Feature importance in Random Forest measures the contribution of each feature in predicting the target variable. We used it to identify which features (e.g., number of casualties, weather conditions) have the most significant impact on predicting the severity of vehicle crash incidents.

**9. What are some potential future improvements for this vehicle crash data analysis project? How can we enhance the model's performance and make it more practical for real-world applications?**

**Answer:** Some potential improvements include collecting more diverse and detailed data, exploring advanced ensemble models like Gradient Boosting, and incorporating

real-time data sources to make the model more dynamic and useful for real-time accident severity prediction.

**10. How would you deploy the final classification model to help improve road safety in a city? What are the challenges and ethical considerations in implementing such a system?**

**Answer:** To deploy the model, we could integrate it with existing traffic management systems or mobile apps to provide real-time risk assessments for different areas. Challenges include ensuring data privacy, handling potential biases in the model, and avoiding over-reliance on the model, as human judgment and context are equally essential for road safety decision-making. Ethical considerations involve transparently communicating the limitations and potential biases of the model to users and continuously monitoring its performance for fairness and accuracy.