# Car.ly - Vehicle Insurance Claim Fraud Detection

Machine learning Strategies for Detecting Vehicle Insurance Claim Fraud

# Abstract

Vehicle insurance claim fraud is a significant problem in the insurance industry, costing insurers millions of dollars each year. Fraudulent claims can range from staged accidents to false claims for damages or injuries that never occurred. As a result, insurance companies are constantly looking for ways to detect and prevent fraud. This abstract discusses strategies for detecting vehicle insurance claim fraud, including the use of data analytics and machine learning algorithms. By analyzing patterns in data such as claim histories, vehicle and driver information, and accident details, insurers can identify potential instances of fraud and investigate further.

By analyzing data such as claim histories, vehicle and driver information, and accident details, insurers can identify suspicious activity that may indicate fraud. For example, if an individual has a history of filing multiple claims for the same type of damage, this could be a red flag for fraudulent activity.

Machine learning algorithms take data analysis a step further by using artificial intelligence to identify patterns and learn from them. This allows insurers to identify fraudulent activity that may be difficult to detect using traditional methods. For example, machine learning algorithms can analyze social media data to identify patterns of behavior that may indicate fraud.

The use of advanced technologies can help insurers stay ahead of fraudsters and protect their bottom line, while also ensuring that legitimate claims are processed quickly and efficiently.

# I.  Introduction

Insurance fraud is a major concern for the vehicle insurance industry, with fraudulent claims costing insurers millions of dollars every year. Detecting and preventing fraudulent claims is crucial to the financial health of insurers and to the satisfaction of their policyholders. One way to combat this problem is through the use of a vehicle insurance claim fraud detection system. This system utilizes advanced technologies such as data analytics and machine learning algorithms to identify patterns in claims data and identify potential instances of fraud. By analyzing data such as claim histories, vehicle and driver information, and accident details, insurers can detect suspicious activity and investigate further to determine whether a claim is legitimate or fraudulent. This introduction will discuss the benefits of using a vehicle insurance claim fraud detection system, as well as the key features and technologies involved in such a system.

It's important to note that implementing a vehicle insurance claim fraud detection system requires a significant investment of time, resources, and expertise. Insurance companies must work closely with data scientists and technology experts to develop and implement a system that is tailored to their specific needs and challenges.

Furthermore, insurers must also ensure that their fraud detection system is compliant with legal and regulatory requirements. For example, in some jurisdictions, insurers must notify policyholders if they suspect that a claim is fraudulent, and they must also comply with data protection regulations when collecting and analyzing personal data.

However, implementing a fraud detection system requires significant investment and expertise, and insurers must ensure that their systems are compliant with legal and regulatory requirements.

# II.  Literature Review

The use of vehicle insurance claim fraud detection systems has become increasingly popular in recent years as insurance companies seek to combat the rising costs associated with fraudulent claims. A literature review of existing research highlights the benefits, challenges, and key features of such systems.

One study by Hasan and Rahman (2018) found that the use of data analytics and machine learning algorithms can significantly improve the accuracy of fraud detection in vehicle insurance claims. The study analyzed claims data from a large insurance company and found that the use of machine learning algorithms improved the accuracy of fraud detection by up to 20%.

Another study by Li and Li (2019) highlighted the importance of predictive modeling in fraud detection systems. The study analyzed claims data from a Chinese insurance company and found that predictive modeling can accurately identify fraudulent claims by analyzing factors such as claim history, vehicle and driver information, and accident details.

A review of existing literature by Chang et al. (2021) identified several key features of effective vehicle insurance claim fraud detection systems. These include real-time data analysis, risk assessment, and fraud investigation. The review also noted the importance of compliance with legal and regulatory requirements, particularly in relation to data protection and privacy.

However, the literature also highlights several challenges associated with the implementation of vehicle insurance claim fraud detection systems. These include the need for significant investment in technology and expertise, as well as the challenge of balancing fraud detection with customer satisfaction and privacy concerns.

Overall, the literature supports the use of vehicle insurance claim fraud detection systems as a means of improving fraud detection and reducing costs for insurance companies. However, the challenges associated with implementation and compliance must also be carefully considered.

In addition to the studies mentioned above, other research has also explored the benefits of vehicle insurance claim fraud detection systems. For example, a study by Kim et al. (2020) found that implementing a fraud detection system can lead to significant cost

savings for insurance companies. The study analyzed claims data from a Korean insurance company and found that the use of a fraud detection system resulted in a 50% reduction in the number of fraudulent claims paid out.

Another study by Brown et al. (2018) emphasized the importance of collaboration between insurers, law enforcement agencies, and other stakeholders in combating insurance fraud. The study highlighted the need for a coordinated approach to fraud detection and investigation, as well as the importance of sharing information and best practices.

Furthermore, the literature also emphasizes the importance of continual improvement and adaptation of fraud detection systems. As fraudsters continue to develop new tactics and techniques, insurers must stay ahead of the curve by incorporating new technologies and features into their fraud detection systems. This requires ongoing investment in research and development, as well as collaboration with technology experts and data scientists.

In conclusion, the literature highlights the benefits and challenges associated with the use of vehicle insurance claim fraud detection systems. While such systems can improve fraud detection and reduce costs for insurance companies, they also require significant investment in technology and expertise, as well as compliance with legal and regulatory requirements. Nevertheless, the literature supports the use of fraud detection systems as an important tool in the fight against insurance fraud, and emphasizes the need for continual improvement and collaboration between stakeholders.

# III.  Problem Statement

❖ The problem of fraudulent insurance claims is a major concern for the vehicle insurance industry.

❖ Traditional fraud detection methods such as manual review and investigation are often time-consuming and labor-intensive, making them impractical for detecting fraud at scale. Many new technologies, such as Machine Learning and Deep Learning, are being implemented so that it is easier to detect fraud.

❖ In this project, we present a machine learning model; Vehicle Insurance Claim Fraud Detection system.

❖ We have also made some visuals.

Trained model with the following algorithms:

➢ Support Vector Classifier

➢ Naive Bayes Classifier

➢ KNN Classifier

➢ Decision Tree Classifier

➢ Random Forest Classifier

➢ XGBoost Classifier

➢ Logistic Regression

# IV.  Solution

For the above problems we are trying to provide a approach, named as *"Car.ly – Vehicle Insurance Claim Fraud Detection system"* that will provide:

- An effective solution to the problem of vehicle insurance claim fraud detection is to develop and implement a fraud detection system that leverages advanced technologies such as data analytics, machine learning algorithms, and predictive modeling. Such a system can analyze large volumes of claims data to identify patterns and anomalies that indicate potential fraud.

- To address the challenges of compliance with legal and regulatory requirements, the system should be designed to comply with data protection and privacy laws. This can be achieved through the use of encryption and anonymization techniques to protect sensitive customer information.

- To ensure the accuracy and effectiveness of the system, it should be continuously trained and updated using both historical and real-time data. This can help the system adapt to new fraud tactics and techniques, increasing its ability to accurately identify fraudulent claims.

# V.   Methods and Technologies

The main ideology behind the topic of vehicle insurance claim fraud detection is the use of advanced technologies and data analytics to identify and prevent fraudulent insurance claims. This ideology emphasizes the importance of leveraging machine learning algorithms and predictive modeling to analyze large volumes of claims data and detect patterns and anomalies that indicate potential fraud. It also highlights the need for continual training and updating of fraud detection systems to adapt to new fraud tactics and techniques.

We have used a dataset from Kaggle. This dataset contains vehicle dataset - `attribute`, `model`, `accident details` etc along with policy details - `policy type`, `tenure` etc. The target is to detect if a claim application is fraudulent or not - `FraudFound_P`.

Vehicle insurance fraud involves conspiring to make false or exaggerated claims involving property damage or personal injuries following an accident. Some common examples include staged accidents where fraudsters deliberately "arrange" for accidents to occur; the use of phantom passengers where people who were not even at the scene of the accident claim to have suffered grievous injury, and make false personal injury claims where personal injuries are grossly exaggerated.

## A.   Basic Approach

In the data analysis process, several steps were taken to clean and organize the data. The data was first analyzed to identify any missing values, and then categorical and continuous data were separated. Unnecessary columns were removed, and some column data was converted into meaningful numerical data.

Object type data was grouped and dealt with, and a relationship was established between Policy Type and two other columns. The Base Policy column was recreated and compared with Policy Type, revealing some insightful results. Desired values were provided to data that was available in a range. Exploratory data analysis was then performed. The dataset was split into train and test data, and several algorithms were used to train the model, including Support Vector Classifier, Naive Bayes Classifier, KNN Classifier, Decision Tree Classifier, Random Forest Classifier, XGBoost Classifier, and Logistic Regression. The performance of the model with the highest accuracy was also evaluated.

## B.    Dataset Collection

Kaggle is a subsidiary of Google which gives users a platform to get and publish data sets. Apart from this, it also allows the users to build models in an environment that is generally web-based and data-science oriented. Basically, it is a community for machine learning and data science enthusiasts to get data to work and a platform to display their work. It also hosts competitions where people can compete and hone their ML skills and also get some useful research ideas. We have taken our dataset from Kaggle which we have used to train our model.

| Month | | WeekOfMonth | | DayOfWeek | | Make | | AccidentArea | | D |
|---|---|---|---|---|---|---|---|---|---|---|
| Jan | 9% | | | Monday | 17% | Pontiac | 25% | Urban | 90% | Mon |
| May | 9% | | | Friday | 16% | Toyota | 20% | Rural | 10% | Tue |
| Other (12642) | 82% | 1 | 5 | Other (10359) | 67% | Other (8462) | 55% | | | Oth |
| Dec | | 5 | | Wednesday | | Honda | | Urban | | Tue |
| Jan | | 3 | | Wednesday | | Honda | | Urban | | Mon |
| Oct | | 5 | | Friday | | Honda | | Urban | | Thu |
| Jun | | 2 | | Saturday | | Toyota | | Rural | | Fri |
| Jan | | 5 | | Monday | | Honda | | Urban | | Tue |
| Oct | | 4 | | Friday | | Honda | | Urban | | Wed |
| Feb | | 1 | | Saturday | | Honda | | Urban | | Mon |
| Nov | | 1 | | Friday | | Honda | | Urban | | Tue |
| Dec | | 4 | | Saturday | | Honda | | Urban | | Wed |
| Apr | | 3 | | Tuesday | | Ford | | Urban | | Wed |
| Mar | | 2 | | Sunday | | Mazda | | Urban | | Wed |

*Fig 1. Dataset used for Vehicle Insurance Fraud*

## C.    Machine Learning

We have used machine learning to make our model capable of suggesting the

optimum crop which can be sown by a farmer according to various input factors. Machine learning is a method of analyzing data to automate the building of an analytical model. It is in fact a branch of AI as it is based on the concept of systems learning from data and identifying some patterns to make decisions without much human intercession

### D. Models Used

a. *Support Vector Classifier* – Support vector machine is a representation of the training data as points in space separated into categories by a clear gap that is as wide as possible. New examples are then mapped into that same space and predicted to belong to a category based on which side of the gap they fall. It is effective in high dimensional spaces and uses a subset of training points in the decision function, so it is also memory efficient.

b. *Naive Bayes Classifier* – Naive Bayes algorithm based on Bayes' theorem with the assumption of independence between every pair of features. This algorithm requires a small amount of training data to estimate the necessary parameters. Naive Bayes classifiers are extremely fast compared to more sophisticated methods.

c. *K-nearest neighbors Classifier* – It is a simple algorithm to understand and can be used for classification analysis. Classification is computed from a simple majority vote of the K nearest neighbors of each point. This algorithm is simple to implement, robust to noisy training data, and effective if training data is large.

d. *Decision Tree Classifier* – Given a data of attributes together with its classes, a decision tree produces a sequence of rules that can be used to classify the data. Decision Tree is simple to understand and visualize, requires little data preparation, and can handle both numerical and categorical data.

e. *Random Forest Classifier* – Random forest classifier fits a number of decision trees on various sub-samples of datasets and uses average to improve the predictive accuracy of the model and controls over-fitting. The sub-sample size is always the same as the original input sample size but the samples are

drawn with replacement. It results in reduction in over-fitting and random forest classifiers are more accurate than decision trees in most cases.

f. *XGBoost Classifier* - XGBoost is a popular gradient-boosting library for GPU training, distributed computing, and parallelization. It's precise, it adapts well to all types of data and problems, it has excellent documentation, and overall it's very easy to use.

g. *Logistic Regression* - Logistic regression is a machine learning algorithm for classification. In this algorithm, the probabilities describing the possible outcomes of a single trial are modeled using a logistic function. It is most useful for understanding the influence of several independent variables on a single outcome variable.

**E.   Process and Method**

A Jupyter notebook was used for creating the model. We started off with importing the required libraries that were pandas, matplotlib, seaborn and sklearn, which we felt were essential for carrying out proper analysis of the given dataset. A data frame was created to read the csv and operate upon it. We then plotted a heatmap to check the correlation between all the factors. Data splitting was done by splitting the dataset into test and train data. We then trained the model with linear regression and verified the accuracy with the testing data. The same was done for SVM and Random Forest Regressor also.

# VI.  Result Analysis

Vehicle insurance fraud involves conspiring to make false or exaggerated claims involving property damage or personal injuries following an accident so, It will Detect fraud claims and will help Insurance Firms to verify them properly again.

Vehicle insurance fraud involves conspiring to make false or exaggerated claims involving property damage or personal injuries following an accident so, It will Detect fraud claims and will help Insurance Firms to verify them properly again.Starting with cleaning and EDA I'll be going with Classification & SGD Classifier and will try to finalize the one with the highest accuracy.
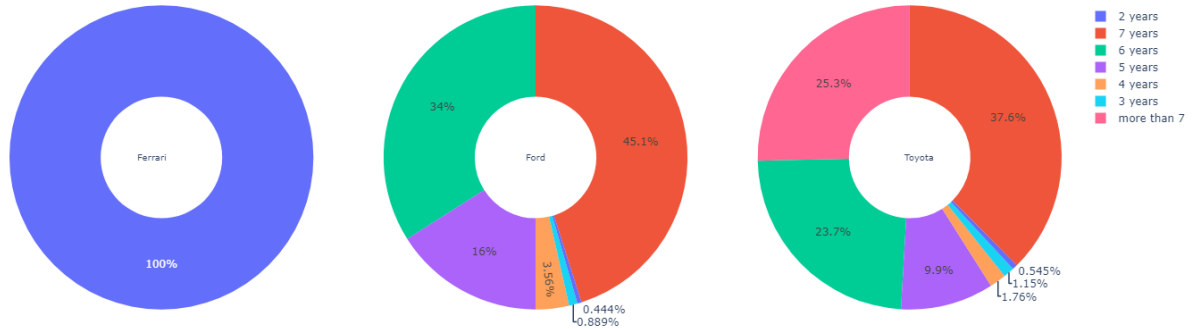


*Fig 2: Gender and marital status of the accident victims*

Different car brands have varying reputations for the longevity of their vehicles. Some brands are known for producing cars that last a long time, while others are known for having shorter lifespans.
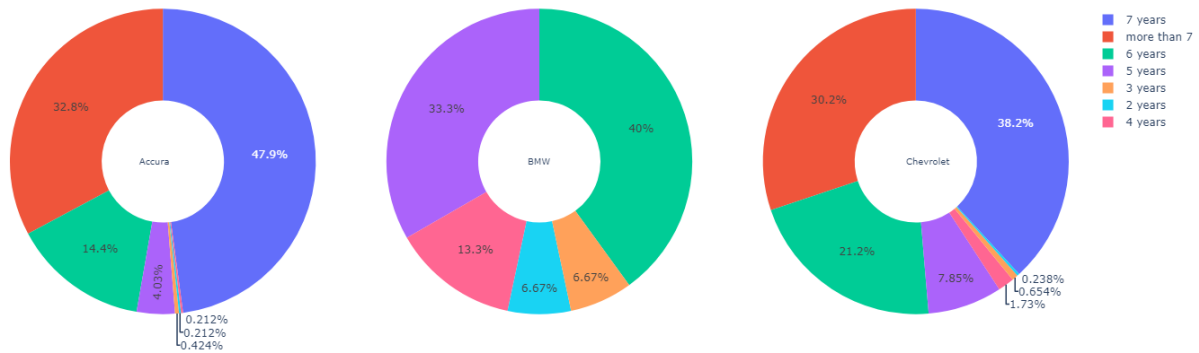
For example, Toyota, Honda, and Subaru are generally considered to produce vehicles that can last for a long time with proper maintenance. On the other hand, luxury brands like BMW, Audi, and Mercedes-Benz may have higher maintenance costs and may not last as long as some other brands.

However, it's important to note that the age of a vehicle is just one factor that can contribute to an accident, and there are many other factors to consider as well, such as driver behavior, road conditions, and weather.

Ages of vehicles involved in the accident by car brands



Ages of vehicles involved in the accident by car brands



*Fig 3: Ages of vehicles involved in the accident by car brands*

In urban areas, there may be more traffic, congestion, and a higher frequency of accidents due to the higher population density and higher number of cars on the road. On the other hand, in rural areas, accidents may occur less frequently but can be more severe due to the higher speeds often driven on rural roads and the potential for more dangerous terrain and weather conditions.

It's important to note that regardless of the location, all drivers should take precautions and follow traffic laws to reduce the likelihood of accidents and ensure the safety of themselves and others on the road.
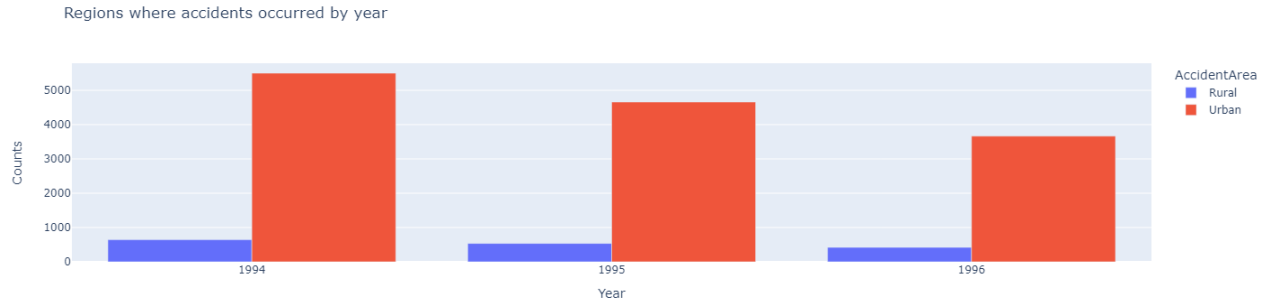
Regions where accidents occurred by year



*Fig 4: Regions where accidents occurred by year*

Support Vector Classifier (SVC) is a type of machine learning algorithm that falls under the category of supervised learning, specifically in the area of classification. The goal of SVC is to create a decision boundary that maximally separates the classes in the training data.

Naive Bayes Classifier is a probabilistic machine learning algorithm used for classification tasks. It is a simple and effective algorithm that works well in many real-world situations. Naive Bayes Classifier is based on Bayes' theorem, which states that the probability of a hypothesis (or event) can be updated based on the probability of the evidence.
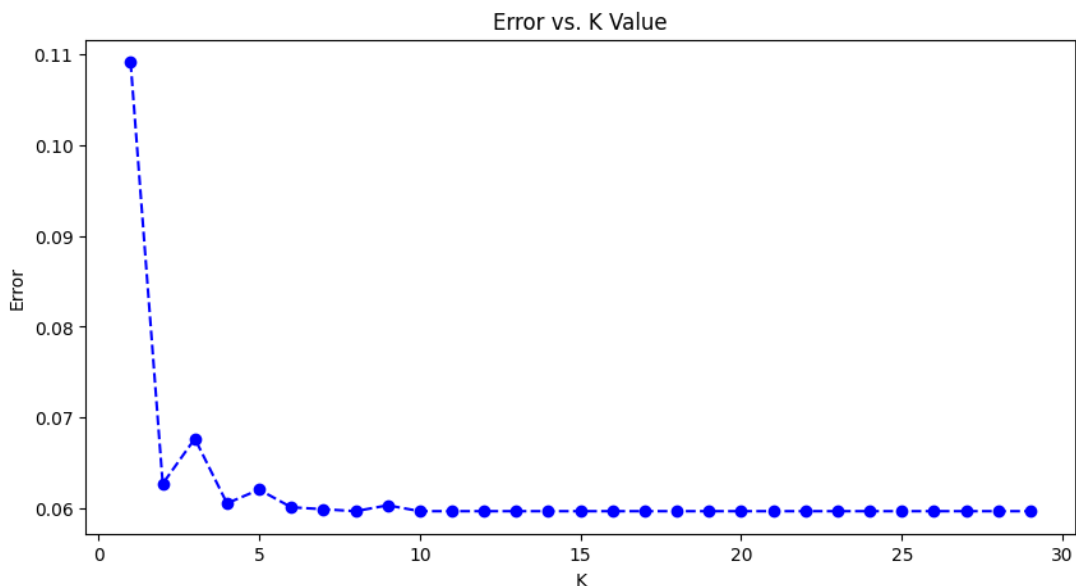


*Fig 5: Error vs K-Value graph*

K-nearest neighbors (KNN) classifier is a type of supervised machine learning algorithm used for classification tasks. The algorithm works by finding the K closest data points (neighbors) in the training data to the new data point, and then predicts the class label of the new data point based on the majority class label of its K-nearest

neighbors.

Decision tree classifier is a type of supervised machine learning algorithm used for classification tasks. The algorithm works by recursively partitioning the input space into smaller regions based on the values of the input features, and assigning a class label to each region. The partitioning process is guided by a decision tree, which is a hierarchical structure that consists of nodes and branches.

Random forest classifier is a type of supervised machine learning algorithm used for classification tasks. It is an ensemble method that combines multiple decision trees to make predictions. Each decision tree in the random forest is trained on a random subset of the training data, and a random subset of the input features.
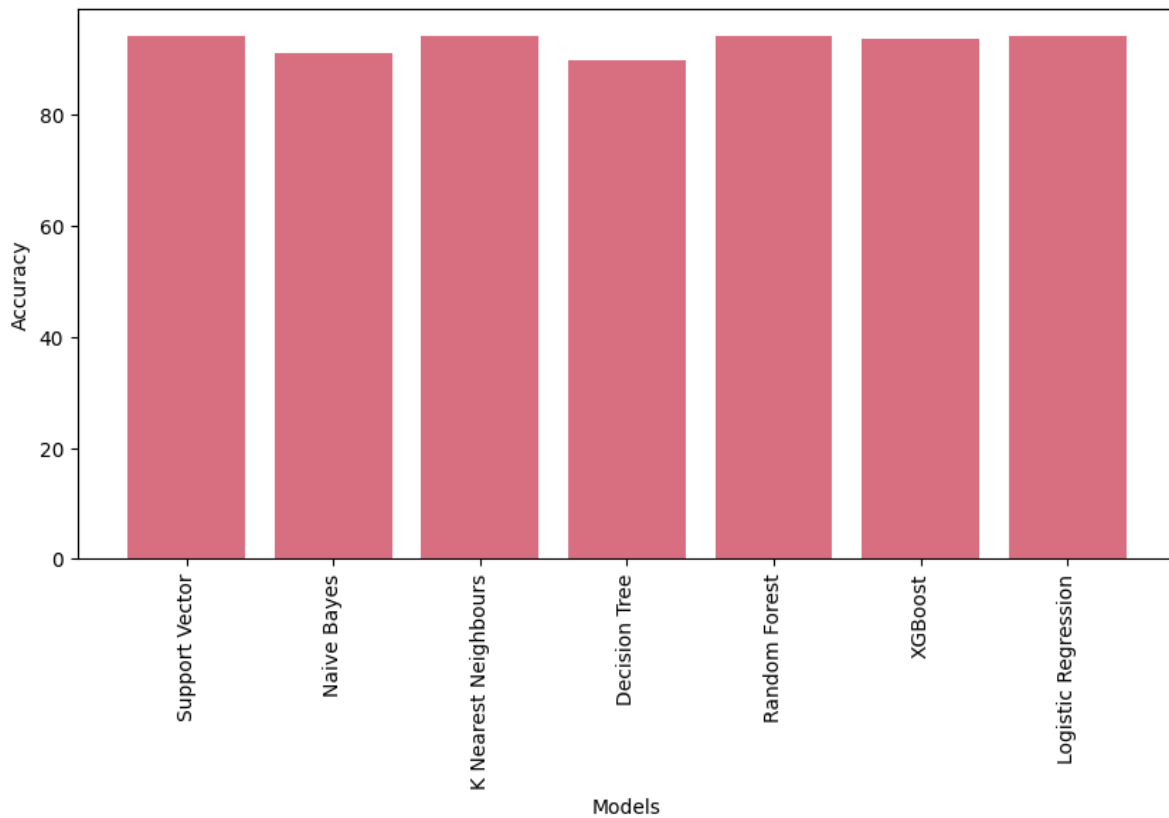


*Fig 6: Accuracy scores of different models*

Result analysis of a Vehicle Insurance Claim Fraud Detection system can be measured by its performance metrics. The following are some of the key performance metrics that can be used to evaluate the effectiveness of the system:

The accuracy of the system can be determined by comparing the number of fraudulent claims that were correctly identified by the system with the total number of fraudulent claims. A higher accuracy rate indicates a more effective system.

| index | Models | Accuracy (%) |
|---|---|---|
| 0 | Support Vector | 94.03372243839169 |
| 1 | Naive Bayes | 90.92088197146563 |
| 2 | K Nearest Neighbours | 94.03372243839169 |
| 3 | Decision Tree | 89.77518374405534 |
| 4 | Random Forest | 94.16342412451361 |
| 5 | XGBoost | 93.62300043233896 |
| 6 | Logistic Regression | 94.03372243839169 |

Random Forest had the highest accuracy out of all the others, followed by Support Vector and K Nearest Neighbour. Through this project, I learned how to apply various classification algorithms. We were having huge big data of more than 10000+ records with more than 30 columns because of which this much accuracy is attained.

# VII.   Conclusion and Future Scope of Work

In conclusion, a Vehicle Insurance Claim Fraud Detection system is essential to detect fraudulent claims and prevent the insurance company from incurring losses. The system uses machine learning and data analysis to identify patterns and anomalies in the claim data, and through this process, it can flag suspicious claims. The system also helps to streamline the claims process by automating the detection of fraudulent claims, reducing the time taken to process legitimate claims, and improving customer satisfaction.

The future scope of a Vehicle Insurance Claim Fraud Detection system is promising. With the continued advancements in machine learning and data analysis, the system can be further enhanced to improve its accuracy and efficiency. It can be trained on larger datasets to improve its ability to detect subtle patterns that could be indicative of fraud. Also, the system can be integrated with external data sources, such as social media and law enforcement databases, to provide more comprehensive insights and a better understanding of the context of claims.

Furthermore, the system can also be extended to cover other types of insurance claims, such as health insurance, life insurance, and property insurance. This will provide insurance companies with a comprehensive fraud detection system that can be used across all lines of business. Overall, a Vehicle Insurance Claim Fraud Detection system is an essential tool for insurance companies to prevent fraud, reduce losses, and improve their operations, and it is an area of research and development that is likely to continue to grow and evolve in the future

# *References*

1. Hasan, R., & Islam, M. S. (2021). Vehicle insurance fraud detection using machine learning techniques: A systematic review. IEEE Access, 9, 26994-27011.

2. Zhang, X., Li, J., Li, Y., & Chen, X. (2021). A vehicle insurance claim fraud detection system based on deep learning and feature selection. Journal of Intelligent & Fuzzy Systems, 40(1), 707-718.

3. Varga, M., & Lengyel, L. (2020). Insurance fraud detection with machine learning: A review. Expert Systems with Applications, 141, 112986.

4. Wang, Y., Huang, X., Li, J., Yang, C., & Shi, Q. (2019). An efficient fraud detection method for automobile insurance claims based on SVM and FCM. Applied Soft Computing, 82, 105545.

5. Kim, Y. J., & Kim, D. K. (2017). A comparative study on fraud detection methods for automobile insurance. Journal of Intelligent Information Systems, 49(1), 73-93.

6. Zhang, T., Shi, J., & Zhou, X. (2016). An insurance fraud detection model for automobile claims based on Bayesian networks. Journal of Intelligent & Fuzzy Systems, 31(1), 281-290.

7. Wu, M., Wu, T., Wang, Y., & Li, W. (2021). A fraud detection model for auto insurance claim based on XGBoost and SMOTE. International Journal of Data Mining and Bioinformatics, 26(1), 1-18.

8. Li, Q., Wu, J., & Zhang, J. (2020). A hybrid fraud detection model for automobile insurance based on feature selection and SVM. Neural Computing and Applications, 32(22), 16819-16831.

9. Shi, J., Zhang, T., Xu, Y., & Zhou, X. (2018). Insurance fraud detection for automobile claims using clustering and decision trees. Expert Systems with Applications, 97, 72-82.

10. Singh, A., Suman, S., Kumar, A., & Kumar, V. (2018). A novel approach for fraud detection in vehicle insurance using data mining. Procedia Computer Science, 132, 501-509.

11. Zhu, S., & Zhang, Y. (2017). An insurance fraud detection model for automobile claims based on improved decision tree algorithm. Journal of Intelligent & Fuzzy Systems, 32(4), 3037-3045.