

# InsureIntel

## Insurance Claim Fraud Detection



Guide - Prof. Asim Tewari

Vishvesh Kodihal      203310013

Amit Meena            190100013

Mudit Rathore        190100078

Dushyant Patil        203170020

Neha Singhal          180070038

# Project Objective

## Aim

The aim of our startup is to provide easy to use and **reliable** solutions to **predict fraud insurance claims** using Machine Learning capabilities and maintaining the privacy of the data of our clients.

The first product predicts the **motor vehicle** fraud insurance claim and later we plan to expand our services to more products.

**Clients** - We provide **B2B** services to **insurance companies** to prevent them from fraud insurance claims in minimum time and affordable prices.

## Brief Description

We wish to develop a **user friendly UI** which uses ML model to predict insurance claim fraud. We train our ML model using the previous insurance claim data provided by our clients.

The model is **trained and updated** regularly for better **accuracy and robustness**. Manual fraud detection requires trained individuals, consumes more time and is not cost effective.

# Problem Definition

<b>Customer Requirement</b>	A <b>cost-effective, accurate</b> , reliable and easy to use product which provides information whether the insurance claim is fraud or not and uses this pattern for future detections. Manual fraud detection requires <b>trained inspectors</b> , consumes more time and cost.
<b>Market Survey</b>	<b>FRISS</b> Fraud Detection at Claims, <b>INSURANALYTICS</b> (Claims AI Cloud), <b>BAE SYSTEMS</b> , Juicy-Score, Bridgei2i, LexisNexis, TransUnion
<b>USP</b>	“ <b>Plug and play</b> ” user-friendly interface, High accuracy (>90%) and specificity. Fast and reliable Insurance claim fraud analysis. Complete data privacy and security.
<b>Protection of USP</b>	User <b>authentication</b> protected UI with password stored in encrypted form. Patenting and copyrighting of our startup solutions.
<b>Barrier to Entry</b>	Patents and copyrights issue from other companies working in this domain. Data Acquisition
<b>Business Case</b>	<ul style="list-style-type: none"><li>● <b>Feasibility:</b> Model will be dynamic and mostly suitable for all the Bike, Car or Four wheeler insurance companies</li><li>● <b>Impacts:</b> It will create a release for task intensive scanning of each Claim (as 80% cases are generally fraud)</li><li>● <b>Benefits:</b> Highly scalable model, continuous support, upgradation, Never ending Demand etc</li></ul>

# Technology Landscape Assessment

<b>Patents</b>	<a href="https://worldwide.espacenet.com/patent/search?q=insurance%20claim%20fraud%20detecti%20on&amp;queryLang=en%3Ade%3Afr">https://worldwide.espacenet.com/patent/search?q=insurance%20claim%20fraud%20detecti%20on&amp;queryLang=en%3Ade%3Afr</a> <a href="https://patents.justia.com/search?q=vehicle+insurance+fraud">https://patents.justia.com/search?q=vehicle+insurance+fraud</a> <a href="https://patents.google.com/patent/US20160117778/en">https://patents.google.com/patent/US20160117778/en</a> <a href="https://worldwide.espacenet.com/patent/search?q=pn%3DCN106600423A">https://worldwide.espacenet.com/patent/search?q=pn%3DCN106600423A</a>
<b>Published literature</b>	<a href="https://www.sciencedirect.com/science/article/pii/S1877050919300079">https://www.sciencedirect.com/science/article/pii/S1877050919300079</a> <a href="https://www.sciencedirect.com/science/article/pii/S1574013721000423">https://www.sciencedirect.com/science/article/pii/S1574013721000423</a> <a href="https://www.sciencedirect.com/science/article/pii/S0167923617302130">https://www.sciencedirect.com/science/article/pii/S0167923617302130</a> <a href="https://ieeexplore.ieee.org/document/8074258">https://ieeexplore.ieee.org/document/8074258</a> <a href="http://www.jcreview.com/fulltext/197-1583405087.pdf">http://www.jcreview.com/fulltext/197-1583405087.pdf</a>
<b>Open libraries</b>	Numpy, matplotlib, keras, Pandas, tensorflow, glob, sklearn
<b>Proprietary libraries</b>	<a href="https://github.com/saritmaitra/Fraud-detection--Insurance">https://github.com/saritmaitra/Fraud-detection--Insurance</a> <a href="https://github.com/sharmaroshan/Fraud-Detection-in-Insurance-Claims">https://github.com/sharmaroshan/Fraud-Detection-in-Insurance-Claims</a> <a href="https://github.com/mehtabhavin10/insurance_fraud_detection">https://github.com/mehtabhavin10/insurance_fraud_detection</a> <a href="https://github.com/ezzaimsoufiane/Auto-Insurance-Claims-Fraud-Detection-with-ML">https://github.com/ezzaimsoufiane/Auto-Insurance-Claims-Fraud-Detection-with-ML</a>

# Project timeline

## 3 week

Team finalised the project topic. Gather Data from various sources

## 1 week

Apply data transformation and visualizations

## 3 week

Models Optimization and Best Model Selection. Dashboard Development

Oct

Nov

## 4 week

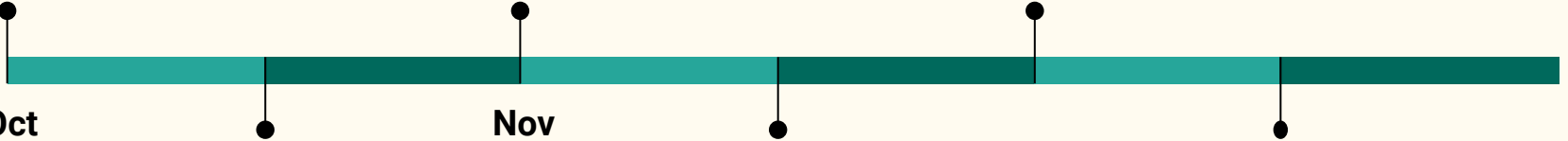
Preparing Data for further analysis

## 2 week

Data Analysis and Features Selection. Training & Evaluating different Models

## 4 week

Post processing and Final Report Submission



# Strategy

Step 1

**Data Collection, Data Cleaning & Filtering, Exploratory Data Analysis (EDA)**

Obtain the dataset from the open sources and clean the data by assigning a valid value to not available values or dropping incomplete rows.

Step 2

**Data Transformation, Data Visualization, Feature selection**

Analyze the relation between features using correlation matrix and the relevance of a feature to the prediction. Drop out unnecessary features (columns).

Step 3

**Modelling, Evaluation & Optimization, Best Model Selection**

Develop ML models such as k-means clustering, logical regression and obtain a model that best fits the data and perform model evaluation using methods as presion, recall, and F1 score.

Step 4

**Dashboard develop-ment, Detailed Reports, Post processing**

Design an authenticated user dashboard where the client can enter the claim data for which prediction is to be done and obtain the predicted results on the dashboard.

# Problem Description

The dataset contains total 37 columns as given below:

1. **months\_as\_customer**: It denotes the number of months for which the customer is associated with the insurance company.
2. **age**: continuous. It denotes the age of the person.
3. **policy\_number**: The policy number.
4. **policy\_bind\_date**: Start date of the policy.
5. **policy\_state**: The state where the policy is registered.
6. **policy\_csl**: How much of the bodily injury will be covered from the total damage.
7. **policy\_deductable**: The amount paid out of pocket by the policy-holder before an insurance provider will pay any expenses.
8. **policy\_annual\_premium**: The yearly premium for the policy.
9. **umbrella\_limit**: An umbrella insurance policy is extra liability insurance coverage that goes beyond the limits of the insured's homeowners, auto or watercraft insurance. It provides an additional layer of security to those who are at risk of being sued for damages to other people's property or injuries caused to others in an accident.
10. **insured\_zip**: The zip code where the policy is registered.

11. **insured\_sex**: It denotes the person's gender.
12. **insured\_education\_level**: The highest educational qualification of the policy-holder.
13. **insured\_occupation**: The occupation of the policy-holder.
14. **insured\_hobbies**: The hobbies of the policy-holder.
15. **insured\_relationship**: Depends on the policy-holder.
16. **capital-gain**: It denotes the monetary gains by the person.
17. **capital-loss**: It denotes the monetary loss by the person.
18. **incident\_date**: The date when the incident happened.
19. **incident\_type**: The type of the incident.
20. **collision\_type**: The type of collision that took place.
21. **incident\_severity**: The severity of the incident.
22. **authorities\_contacted**: Which authority was contacted.
23. **incident\_state**: The state in which the incident took place.
24. **incident\_city**: The city in which the incident took place.
25. **incident\_location**: The street in which the incident took place.
26. **incident\_hour\_of\_the\_day**: The time of the day when the incident took place.
27. **property\_damage**: If any property damage was done.
28. **bodily\_injuries**: Number of bodily injuries.
29. **Witnesses**: Number of witnesses present.



- 30. **police\_report\_available:** Is the police report available.
- 31. **total\_claim\_amount:** Total amount claimed by the customer.
- 32. **injury\_claim:** Amount claimed for injury
- 33. **property\_claim:** Amount claimed for property damage.
- 34. **vehicle\_claim:** Amount claimed for vehicle damage.
- 35. **auto\_make:** The manufacturer of the vehicle
- 36. **auto\_model:** The model of the vehicle.
- 37. **auto\_year:** The year of manufacture of the vehicle.

Using the dataset we predict whether the claim is fraud or valid, thus we use SVM and XGBoost classification model.

**Output :** Fraud\_reported, Y or N (Y: claim fraud, N: claim valid)

Some columns are not required for classification such as 'policy\_number', 'policy\_bind\_date', 'policy\_state', 'insured\_zip', 'incident\_location', 'incident\_date', 'incident\_state', 'incident\_city', 'insured\_hobbies', 'auto\_make', 'auto\_model', 'auto\_year', thus we drop these columns.

In total 26 columns have been used for prediction.