# Case Study- Fraudulent Claim Detection

By: Ilan
Sachin
Disha

# Summary

Plotting data for Fraud in terms of percentage (insured_education_level)



Plotting data for Fraud in terms of percentage (insured_relationship)



Plotting data for Fraud in terms of percentage (incident_severity)

**Objective:**
- Building a model to classify claims as fraudulent or legitimate early in the approval process.
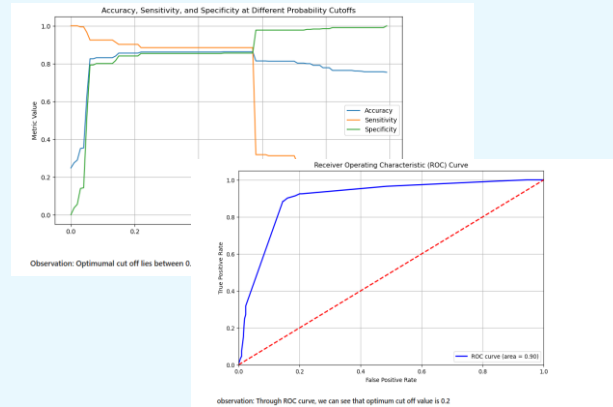
**Approach:**
- Data cleaning has been done by dropping blank rows and columns and imputing junk values.
- Univariant and Bivariate Analysis has been conducted between target variable and categorical variables and Correlation Analysis for numerical feature.
- Class imbalance was addressed by RandomOverSampler method.
- Feature Reduction is done by combining the existing feature
- RFE and Cross Validation methods were used in Logistic Regression for feature selection
- Grid Search for hyperparameter tuning in Random Forest

**Findings:**
Higher proportion if fraud claims observed among:
- Males
- Education Level: College students
- Occupation as: exec-managerial, farming/fishing
- Hobbies- Chess, Crossfit
- Insured relationship: Other relative
- Incident Type: Multi-vehicle and single-vehicle collision
- Incident Severity: Major Damage
- Incident State: Ohio
- Incident City: Arlington
- Number of vehicles involved: 4
- Auto Make: Dodge, Ford, Volkswagen
- Auto Model: A5, Silverado, X6 etc.

# Logistic Regression







**Feature Selection:**

Insured_occupation_exec-managerial, insured_hobbies_camping, insured_hobbies_chess, incident_severity_Trivial Damage, incident_severity_Total Loss, incident_severity_Minor Damage, insured_hobbies_cross-fit were selected as the high contributing features for the fraud by the Recursive Feature Elimination method.

**Optimal Cutoff**

ROC Curve area=0.90,
The cut-off point lies between 0.2 and and 0.3 which is acceptable

**Sensitivity, specificity, precision, recall and F1-score**

Sensitivity: 0.8513513513513513
Specificity: 0.8539823008849557
Precision: 0.65625
Recall: 0.8513513513513513
F1 Score: 0.7411764705882353

# Random Forest







**Feature Selection:**
Selected features with importance score greater than 0.02 and they have been listed below,
Insured_hobbies_cross-fit,
incident_severity_Total Loss, total_claim_ratio,
insured_hobbies_chess,  vehicle_claim,
incident_severity_Minor Damage,
incident_severity_Trivial Damage,
insured_sex_MALE, policy_annual_premium,
age_policy_deductable_interaction,incident_ho
ur_of_the_day, claim_per_vehicle

**Hypertuning:**
max_depth=3, min_samples_leaf=5,
n_estimators=100, max_features='sqrt',
criterion='entropy', class_weight='balanced',

**Sensitivity, specificity, precision, recall and F1-score:**

Sensitivity: 0.8513513513513513
Specificity: 0.8539823008849557
Precision: 0.65625
Recall: 0.8513513513513513
F1 Score: 0.7411764705882353

# Recommendations

- Prioritize Camping and Dancing demography over the others in the insured_hobbies, as they show least fraud rate compared to Chess and Cross-fit which shows highest fraud rate in the feature.
- Investigate Multi-vehicle and Single Vehicle Collisions over Vehicle Theft and Parked Cars in incident_type feature.
- The fraud rate increases with the number of vehicles involved in the accident hence, more investigation is needed with increase in the number of vehicles in an accident.
- College seems to have high fraud rate and other than that, the fraud rate increases with the degrees that needs more number of days to complete.
- Exec-managerial and armed-forces has high fraud rate and has to be investigated further.
- Males tend to have a high fraud rate compared to females, even when the class imbalance is observed in favor of females.

# Business Implications

- Chess and Cross fit has high rate of fraud. Usually these activities doesn't require costly equipment and they'r build quality is strong, so the possibility of a high failure rate is not possible hence, it requires more check for these cases.
-  Multi-vehicle or Single-vehicle collision has high rate of fraud, maybe due to the lenient inquiries or legal implications for these cases than a stolen car.
- With the increase in the number of days required to complete a degree, the fraud rate increases this may be due to high dropout rate.
- Exec-managerial and armed-forces may have a high fraud rate, possibly due to their influence over the higher authorities.

# Questions

- How can we analyse historical claim data to detect patterns that indicate fraudulent claims?
  Yes, it is possible to analyze historical claim data to identify patterns indicative of fraudulent activity. The process begins with data cleaning and imputation to address inconsistencies. This is followed by exploratory data analysis (EDA) to understand data trends and detect imbalances. Feature engineering can then be applied to correct these imbalances and enhance the dataset. Based on the insights gained, an appropriate model can be selected—often, unsupervised learning models are used for such fraud detection tasks.

- Which features are the most predictive of fraudulent behavior?
  The feature **insured_hobbies_cross-fit** is the strongest predictor of fraudulent behavior, with an importance score of **0.134902**.

- Based on past data, can we predict the likelihood of fraud for an incoming claim?
  Yes, it is possible to predict the likelihood of fraud for an incoming claim. After training the model, we can evaluate its performance on test data using metrics such as recall, precision, sensitivity, specificity, and accuracy. These metrics help assess the model's ability to detect potential fraud.

- What insights can be drawn from the model that can help in improving the fraud detection process?

  1. Higher fraud rates were associated with certain demographics, including:

     ○ Males
     ○ Customers with specific educational backgrounds (College students)
     ○ Certain occupations (e.g., exec-managerial, armed-forces)
     ○ Certain hobbies (e.g., chess, cross-fit)

  2. Incident types such as Multi-vehicle and Single Vehicle Collisions showed higher fraud likelihood.

  3. The number of vehicles involved and specific auto makes/models also correlated with fraudulent claims.