

ML:CIA - 2

Churn Analysis On Life Insurance.

- Name: Hardik Shah
- Register No.: 2241131
- Class: 6 BCA - B

Dataset description:

This dataset provides a comprehensive overview of customer churn in the life insurance sector, containing 1000 records with 11 features.

- Features:

- 1.) Customer Name.
- 2.) Customer Address.
- 3.) Company Name.
- 4.) Claim Reason.
- 5.) Data confidentiality (The level of confidentiality of the data).
- 6.) Claim Amount.
- 7.) Category Premium.
- 8.) Premium Ratio (ratio of claim amount to the premium amount).
- 9.) Claim Approved.
- 10.) BMI.
- 11.) Churn (Target variable).

* The target variable "Churn", identifies whether a customer has left the insurance policy or not.

customer name	customer address	company name	claim reason	Data confidentiality	claim amount	category premium	Premium ratio	claim Approved	BMI	churn
TONY	XXX...	WH&P	Travel	low	322	479	0.075	No	21	Yes
Nicole	XXX...	M&P	Medical	Medium	14089	14390	0.100	Yes	24	Yes
Linda	XXX...	HP	Phone	High	1875	1875	1.124	Yes	21	No

* columns like; 'customer Name', 'customer Address', 'company name' & 'BMI' have been dropped as they are irrelevant.

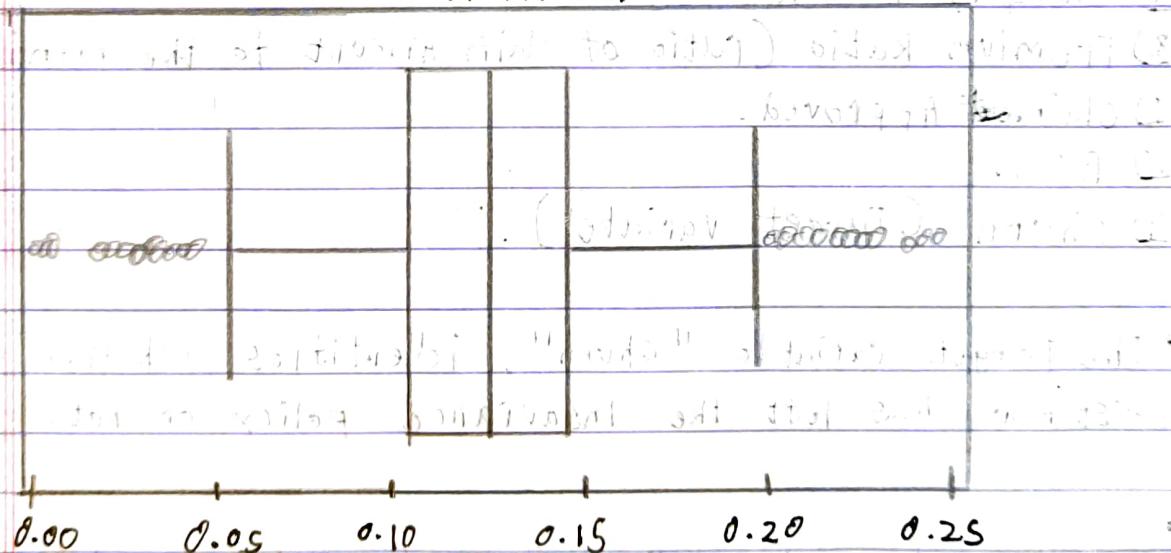
* The column 'BMI' has been dropped as it is highly correlated to the Target variable, which is thorough off the analysis.

• EDA Performed:

- 1.) All rows with Null values have been discarded.
- 2.) Dataset was checked for duplicate rows.
- 3.) 'churn' was renamed to 'Target' for easier Analysis.
- 4.) Outlier Detection ('premium Ratio' was the only column with outliers).

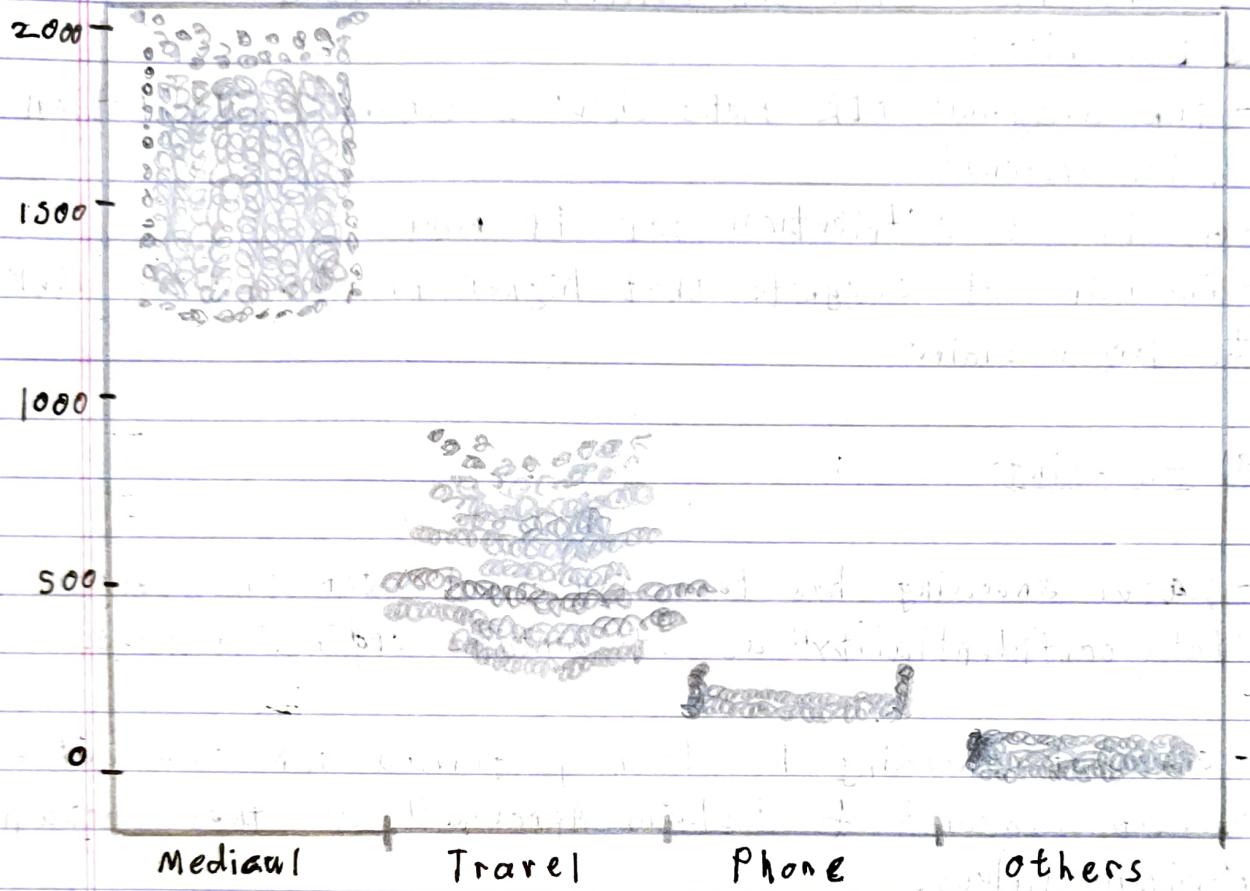
- Box plot:

Premium Ratio



- There are multiple outliers in the column 'Premium Ratio'. Hence we need to get rid of them.
- The IQR Method was used to remove outliers.

5) Swarm plot (For the column 'claim reason').



-Interpretation:

- 1.) Medical claims have the highest variability & the largest claim amount. This might require a higher risk management.
- 2.) Travel claims exhibit a mid-range cost pattern but are more varied than Phone & other claims. This is more moderate.
- 3.) Phone & other claims are associated with consistently low claim amount. They represent minimal risk with low-cost claims.

6) Violin plot. (For the column 'claim reasons').

7) Count plot. (For the column 'claim reasons').

8) Pair Plot:

- The diagonal KDE plots show a binomial distribution for 'claim amount'.
- A skewed distribution for Premiums.
- The pair plot suggests that higher premium are linked to larger claims.

9) Encoding:

- Label encoding has been performed for the column "Data confidentiality" as it contains ordinal values.

- One-Hot encoding has been performed for the columns "claim Reasons" & "claim Approved" as they contain Nominal values.

- Binary Encoding has been performed on the column 'Target' (Churn).

10) Standard scaling: (all values were scaled using the StandardScaler() method (-1 to 1)).

2) Heat map:

Data confidentiality	1	0.83	0.87	0.027	-0.0052
Claim Amount	0.83	1	0.96	0.24	0.0033
Category Premium	0.87	0.96	1	0.021	-0.0039
Premium Ratio	0.027	0.24	0.021	1	-0.02
Target	-0.0052	0.0033	-0.0039	-0.02	1
	Data confidentiality	claim Amount	Category Premium	Premium Ratio	target
			Premium	Ratio	

• Interpretation:

- Category Premium & Claim Amount have a very high correlation (0.96). Indicates that higher premiums are associated with higher claim amount.
- Data Confidentiality is highly correlated with Category Premium & Claim amount, suggesting that this data might be related to high profile clients.
- The strong correlation between Claim amount & Category Premium suggests potential multicollinearity.

ML Analysis:

1) KNN: K nearest neighbour

- Knn is a supervised learning algorithm, used for classification & regression that classifies a data point based on the majority class of its K nearest neighbors.

- Key characteristics:

- i) Lazy Learning: No training phase.
- ii) Instance Based: classifies based on similarity.
- iii) Distance Based: Euclidean Distance.

- Algorithm

- i) choose the number of neighbors 'K'.
- ii) calculate distance between the query point & all data points.
- iii) sort the distance & select the 'k' nearest neighbors.
- iv) perform majority voting (classification) or take the average (Regression).
- v) Return the predicted label / value.

- Error Analysis:

- Using the Error Rate vs. K value plot we select the best 'k' value (minimum Error Rate).
- According to the Error Rate plot the best value for 'k' is 36.

- Results:

- confusion Matrix:

	0	58
0	0	115
1	1	1

$$\therefore \text{Accuracy} = \underline{\underline{66.85\%}}$$

- Interpretation:

i) class 0 (Negative class)

- Precision = 1.00 : The model predicted class 0 perfectly, but it does not account for the fact that it rarely predicted class 0.
- Recall = 0.03 : Only 3% of actual class 0 instances were correctly predicted as class 0. This indicates poor sensitivity for this class.
- F1 score = 0.06 : A low F1 score reflects poor performance of class 0.

Confusion Matrix

- Class 1 (positive class)

- precision = 0.66: when the ~~model~~ predicts class 1, 66% of those predictions were correct.

- Recall = 1.00: The model identified all class 1 instances correctly.

• 2) Decision Tree:

i) Using Gini Index:

- using the Gini criteria the model selected the column 'Claim Amount' as the Root decision node.

- Results: 1) f1 score 0.55, 2) accuracy 66%

- Accuracy: 66%.

- Precision: 0.64

- Recall: 0.66

- F1 score: 0.55

• Confusion Matrix:

0-	3	57
1-	2	113
	0	1

ii) Using Entropy criteria:

- The Entropy model also selected 'claim amount' as the root decision node.

- Results: -

- Accuracy: 66%

- Precision: 0.64

- Recall: 0.66

- F1 Score: 0.55

• Interpretation:

- The decision tree classifiers (Gini & Entropy) achieved moderate accuracy (66.29%) but showed significant bias towards class 1.

- Both models had poor performance for class 0 with low recall & F1 score, while class 1 achieved a high recall (98%) & a stronger F1 score.

- The AUC scores (Gini = 0.54 & Entropy = 0.51) indicates weak discriminatory power.

- The dominant feature in both classes is claim amount, followed by premium ratio.

• Comparative Analysis: (Conclusion).

- Overall Accuracy: Both models achieve similar accuracy, indicating that they are moderately effective overall.

- Class Level Analysis: KNN achieves perfect precision for class 0, but its recall is extremely low. Decision tree has a slightly better recall but low precision. Both models perform well for class 1, with high recall & F1 score.

- Both models have weak discriminatory power, with scores close to 0.5, indicating that neither is reliable at separating classes.

- To improve the performance, particularly for class 0, it is essential to address the class imbalance through data rebalancing techniques.