

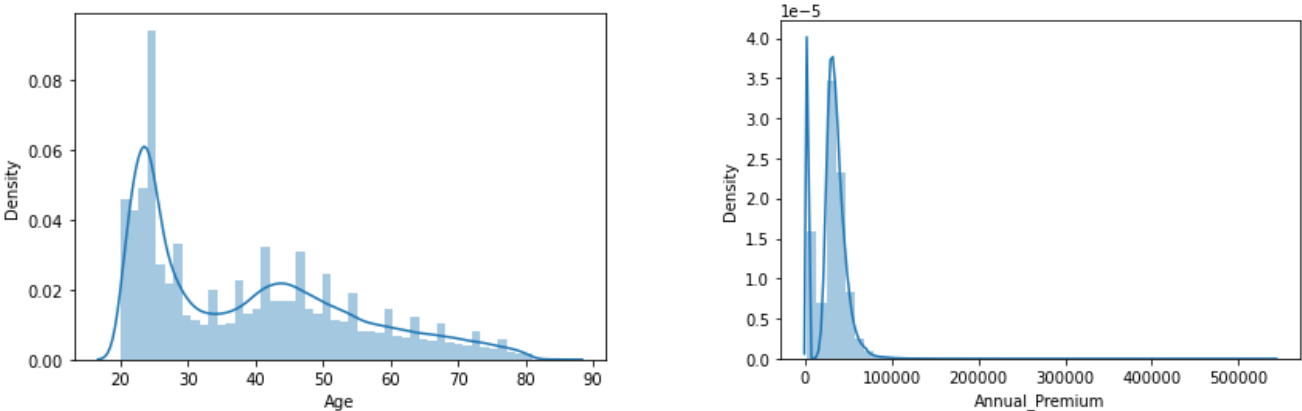
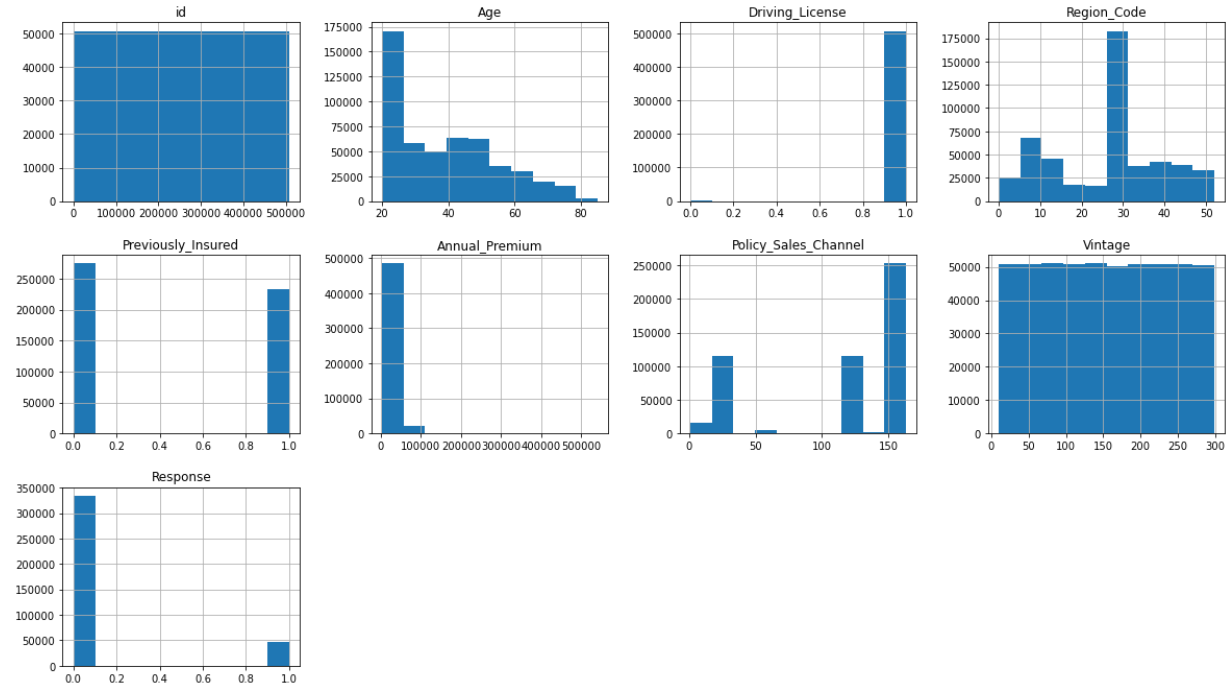
# Cross Prediction Analysis

## Characteristics of dataset:

- No Null Values
- Data types -
  - Numerical:  
'id', 'Age', 'Driving\_License', 'Region\_Code', 'Previously\_Insured', 'Annual\_Premium', 'Policy\_Sales\_Channel', 'Vintage', 'Response'
  - Categorical:  
'Gender', 'Vehicle\_Age', 'Vehicle\_Damage'

## First look insights of data

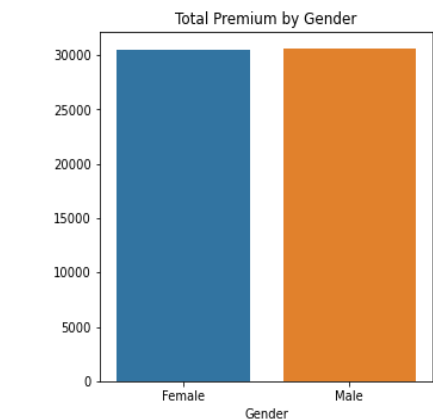
- Age and Annual premium are skewed
- Driving license doesn't add any value, can be dropped
- Gender, Vehicle age, damage, previously insured to be encoded
- Age, Vintage and Annual Premium can be binned
- After binning, vehicle\_Damage', 'Vehicle\_Age', 'Age', 'Gender', 'Previously\_Insured', 'Vintage', 'Annual\_Premium' can be one hot encoded



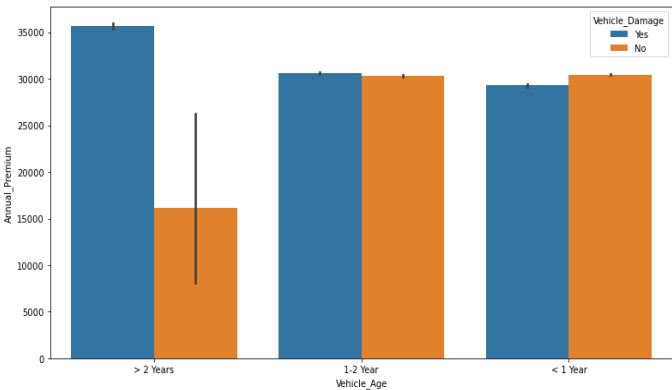
id	1	0.0018	-0.00064	0.0014	0.0013	0.0026	-0.00036	0.0023	0.001	-0.0028	-0.00039	-0.0014
Gender	0.0018	1	-0.15	0.018	-0.00068	0.083	0.11	0.092	-0.0035	0.11	0.0021	-0.052
Age	-0.00064	-0.15	1	-0.08	0.042	-0.26	-0.52	-0.27	0.067	-0.58	-1e-05	0.11
Driving_License	0.0014	0.018	-0.08	1	-0.0002	0.014	0.03	0.016	-0.012	0.043	-0.00085	0.01
Region_Code	0.0013	-0.00068	0.042	-0.0002	1	-0.024	-0.028	-0.027	-0.012	-0.042	-0.0016	0.011
Previously_Insured	0.0026	0.083	-0.26	0.014	-0.024	1	0.18	0.82	0.0048	0.22	0.0017	-0.34
Vehicle_Age	-0.00036	0.11	-0.52	0.03	-0.028	0.18	1	0.17	0.024	0.39	0.0019	-0.1
Vehicle_Damage	0.0023	0.092	-0.27	0.016	-0.027	0.82	0.17	1	-0.0095	0.23	0.0011	-0.35
Annual_Premium	0.001	-0.0035	-0.067	-0.012	-0.012	0.0048	0.024	-0.0095	1	-0.11	-0.00085	0.023
Policy_Sales_Channel	-0.0028	0.11	-0.58	0.043	-0.042	0.22	0.39	0.23	-0.11	1	-0.00081	-0.14
Vintage	-0.00039	0.0021	-1e-05	-0.00085	-0.0016	0.0017	0.0019	0.0011	-0.00085	-0.00081	1	-0.0011
Response	-0.0014	-0.052	0.11	0.01	0.011	-0.34	-0.1	-0.35	0.023	-0.14	-0.0011	1
id	Gender	Age	Driving_License	Region_Code	Previously_Insured	Vehicle_Age	Vehicle_Damage	Annual_Premium	Policy_Sales_Channel	Vintage	Response	

# Story board

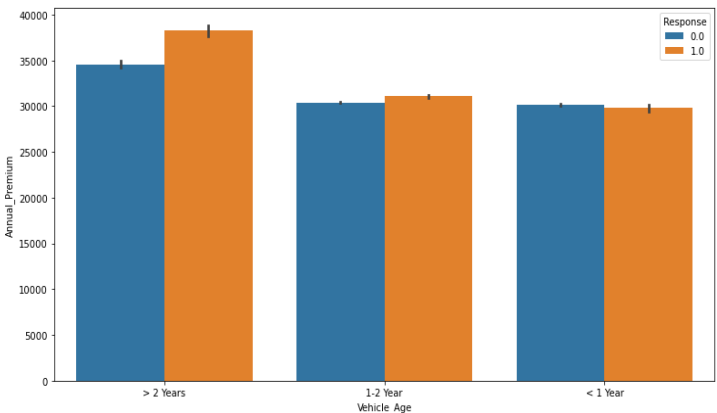
Male population opts for slightly high annual premium



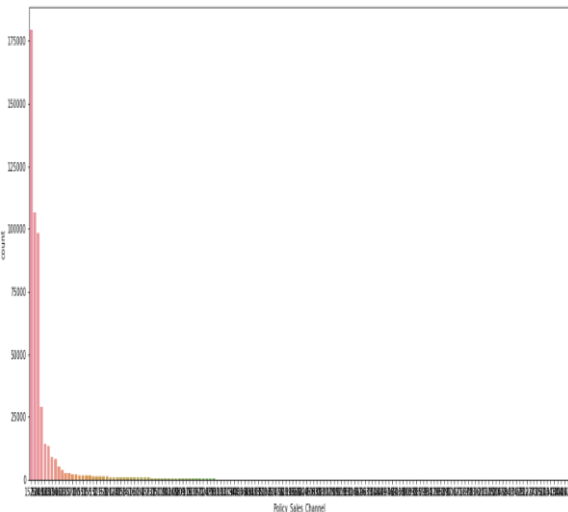
People who own vehicle > 2 years with damage pays for high annual premium



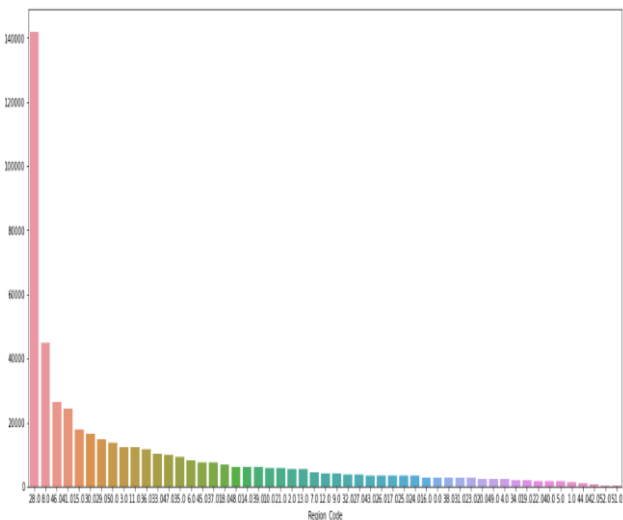
Damaged Vehicle with 1-2 years old and who are insured already are more likely to buy insurance again



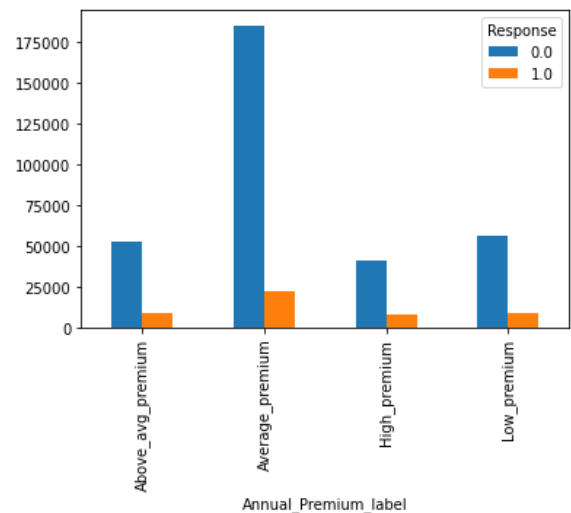
Only few policy channels get more customers



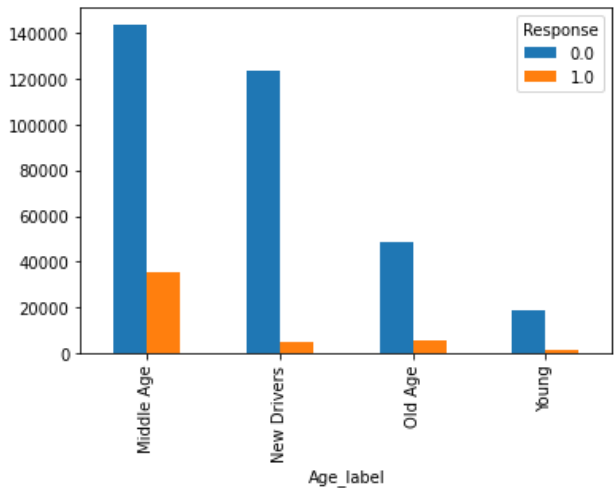
Region code 28 gets the more customers



Generally people do not prefer purchasing high premium



Middle age people are likely to buy premium and own most of the premium



# Approach

## 1. Data Cleaning:

- Convert all the categorical values to numerical values for some features

## 2. Feature Engineering:

- Binning/Discretization –
  1. Pick out the features to perform data binning
    1. Age
    2. Vintage
    3. Annual Premium
  2. Find the best depth of tree to perform binning using decision tree classifier
  3. Check if the bin's relationship with target is linear
  4. Find min and max of bin categories
  5. Label the bins into a new discrete column and assign numerical values to each bins
- Encoding –
  1. Encode the features with less unique values using one hot encoding
  2. Remove the original values

## 3. Model Selection: Try running the classifier models

## 4. Error Metrics: ROC AUC for all models

## 5. Hyper parameter tuning: Use Grid search for choosing the right parameters

## 6. Summary: Load all the Error metrics to a dataframe for better decision making to choose right model

## 7. Best Model: After checking the ROC AUC for all models, it is evident that **CATBOOST** works well for this model, both with and without hypertuning parameters, giving a final result of **85.2** for Test dataset

# Result

	Model	Accuracy Score	F1 Score	Precision Score	Recall Score	ROC AUC
0	Logistic Regression	0.878495	0.000144	1.000000	0.000072	0.837371
1	Naive Bayes CB	0.718734	0.425894	0.283184	0.858562	0.835353
2	Random Forest	0.871087	0.107539	0.338673	0.063917	0.829578
3	Decision Tree	0.871131	0.095074	0.323985	0.055712	0.806867
4	AdaBoost	0.878487	0.000000	0.000000	0.000000	0.846477
5	CatBoost	0.877752	0.053113	0.451613	0.028216	0.852118
6	Logistic Regression - Tuned	0.878495	0.000144	1.000000	0.000072	0.837371
7	Cat Boost - Tuned	0.877752	0.053113	0.451613	0.028216	0.852118