

Customer Conversion Prediction

Problem Statement

Problem statement is to build a ML model that predict if a client will subscribe to insurance or not. We need to use the historical data to reduce the cost of telephonic marketing, still it is effective but incur a lot of cost. It is also important to identify the customers that are most likely to convert beforehand so that they can be specifically targeted via call.

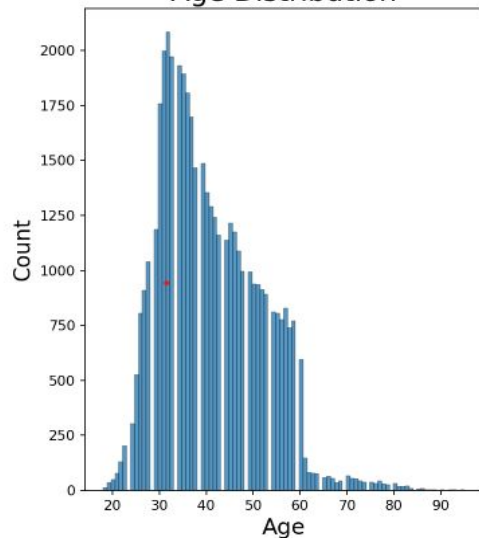
Data Cleaning

1. Checked for null values but found no null values
2. Checked for duplicate values , we got some duplicate values and we removed them
3. Checked for data types but all the data types were found correct
4. Checked outliers by describe function and we got some outliers but we have imbalance data so we are not able to determine this is outliers or not

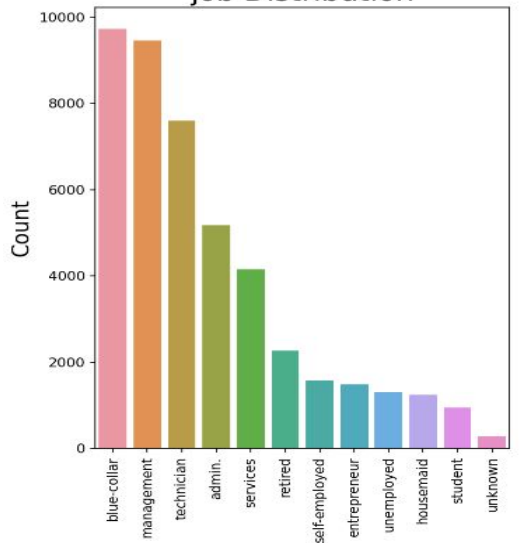
Important insights of EDA

1. Checked linear correlation between the features but we didn't get any relationships between the data
2. Interpretations of feature distribution

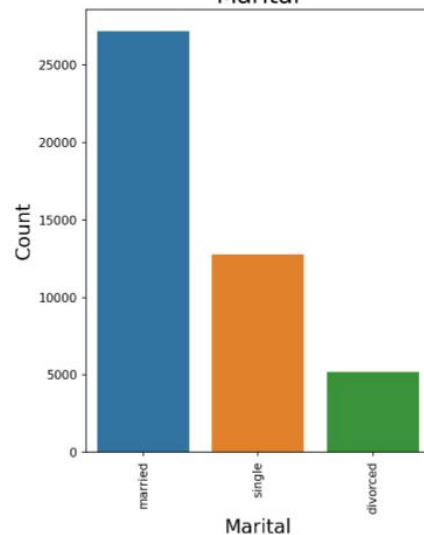
Age Distribution

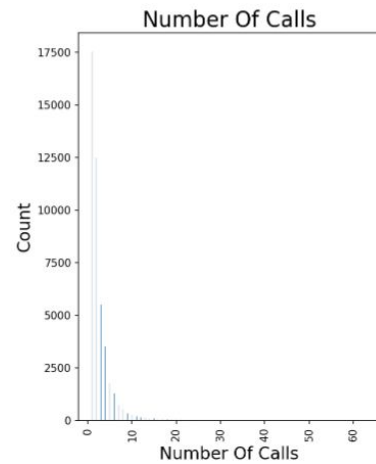
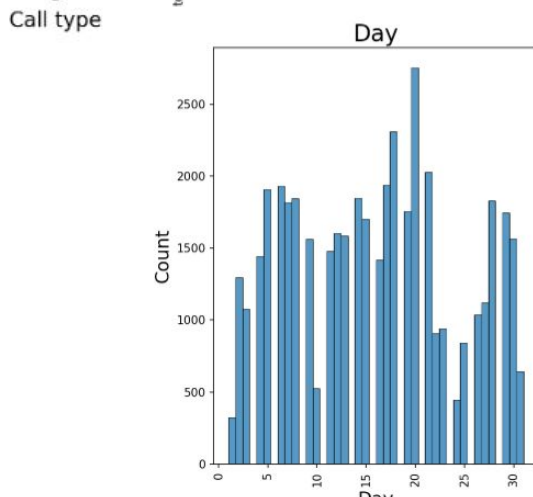
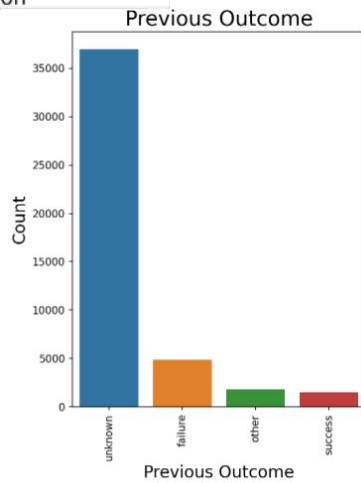
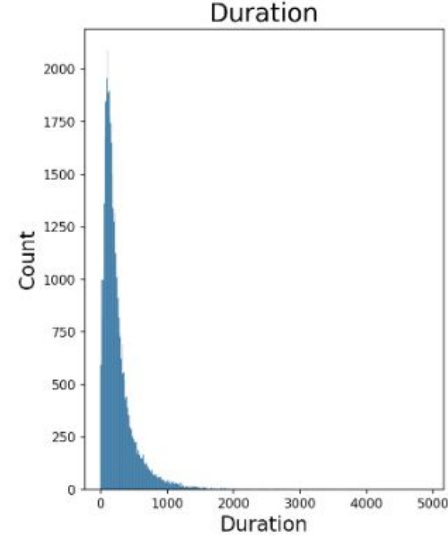
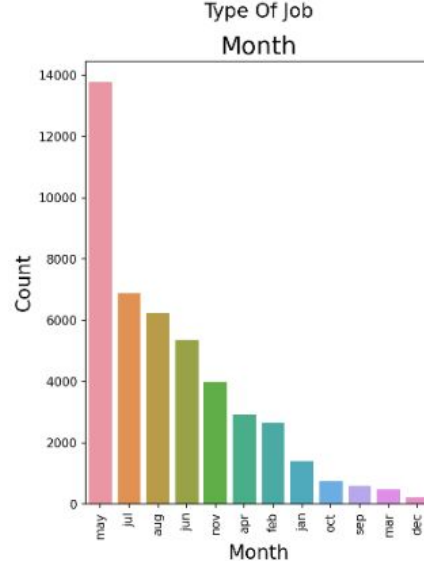
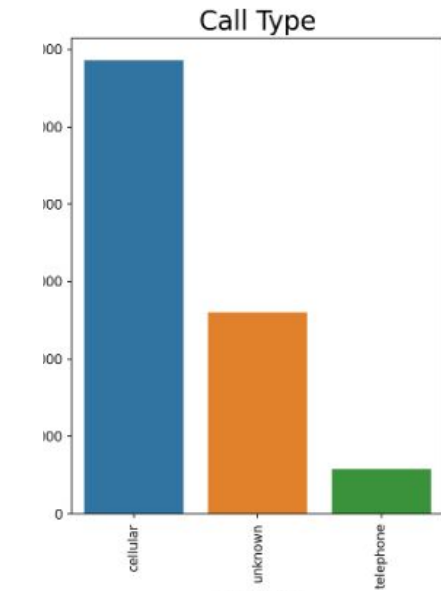
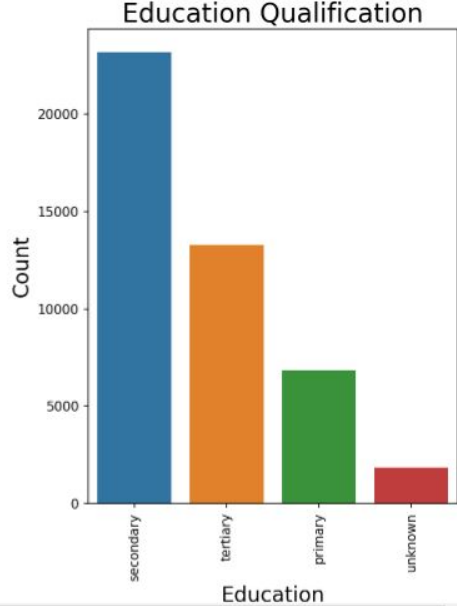


job Distribution



Marital





1. .AGE

- Company targets 30 to 40 ages people to sell their insurance plan
- People aged less than 25 years and more than 60 years are less targeted by the company for insurance

2. JOB

- Blue collar and management people are targeted most by insurance company
- Students are least targeted for insurance plan

3. MARITAL STATUS

- Married people are targeted the most and divorced people are least targeted

4. EDUCATIONAL QUALIFICATION

- Secondary qualified people have been targeted thye most for insurance followed by tertiary and primary

5. CALL_TYPE

- Cellur mode is most used for contacting the people for campaign

6. Month

- May month is when maximum calls are generated for new insurance customers, december is the least no of calls placed by the company

7. Duration

- Maximum number of calls lasted less than 1000 seconds

8. Previous Outcome

- More unknown as status than failure and success

9. Day

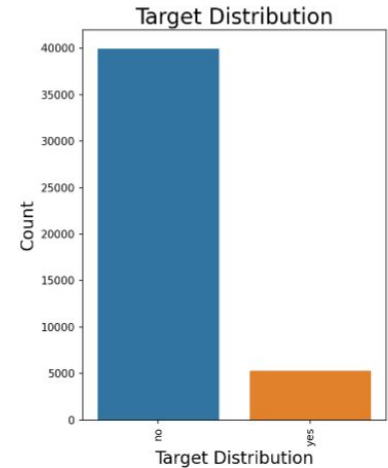
- Most number of calls are between in the mid of month

10. Number of Calls

- Most People are contacted only 1 or 2 times by the company

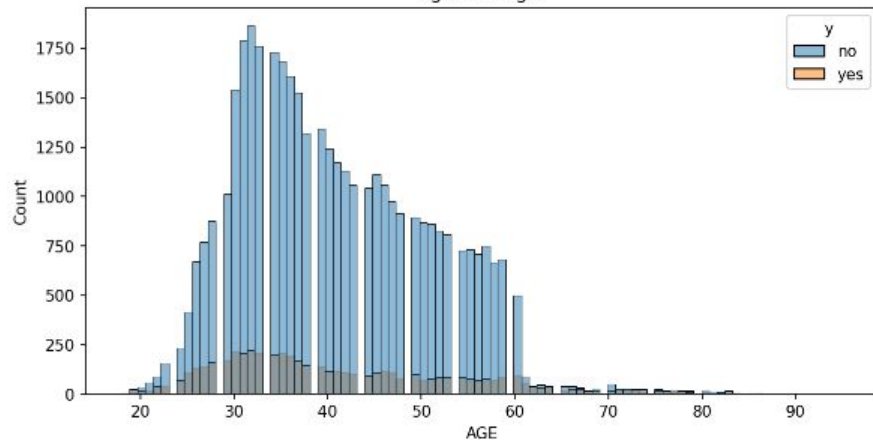
11. Target

- No of people subscribed is very less compared to Unsubscribed people

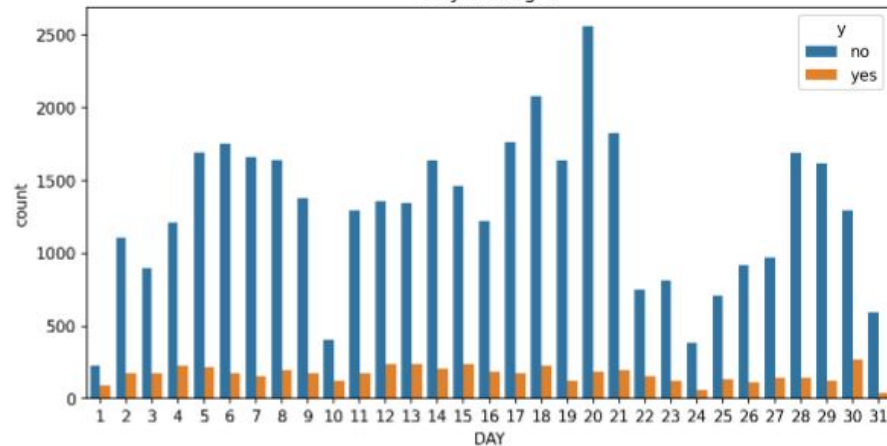


EDA on feature vs Target

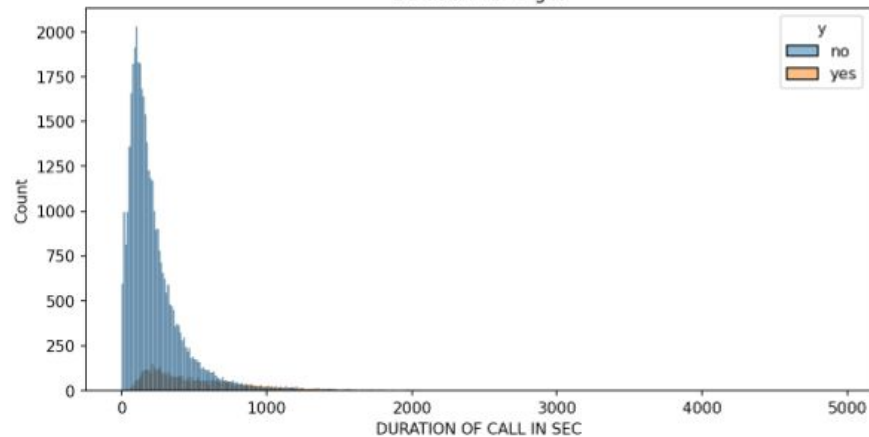
Age Vs Target



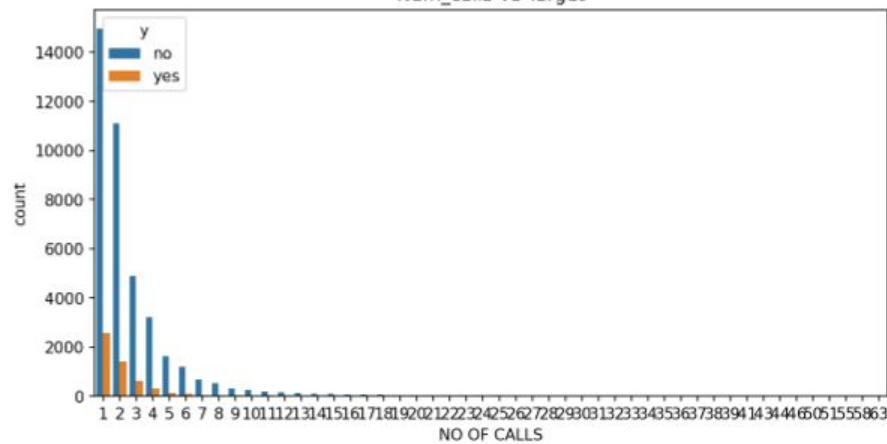
Day Vs Target

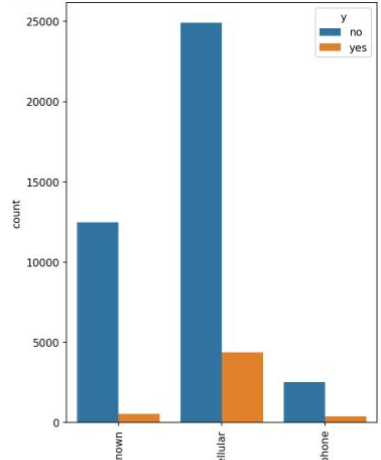
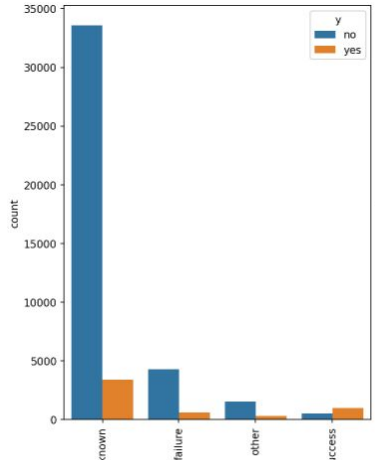
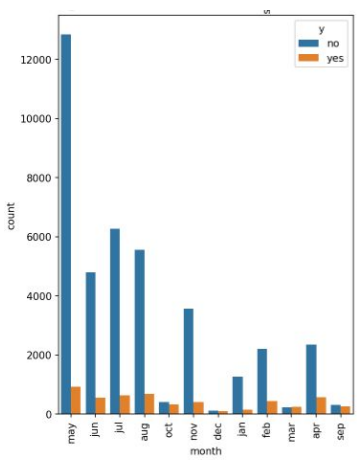
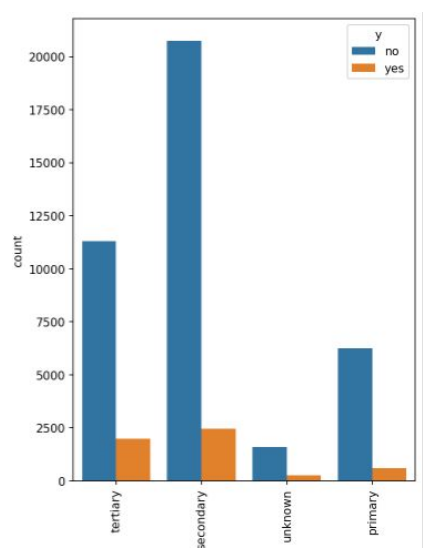
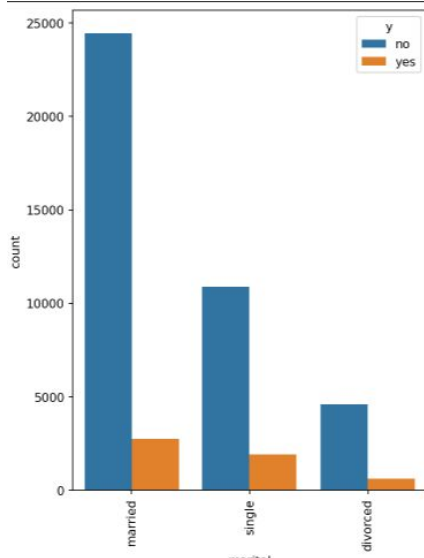
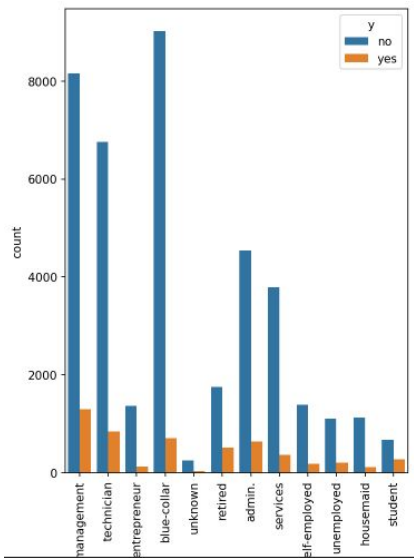


Duration Vs Target



Num_calls Vs Target





Built Models

- We Encode our data by Label and One-hot encoding
- We scaled our data by standard scalar and Split our data
- We imported SMOTETomek module from imblearn.combine to balance our trained data
- We built Logistic regression , KNN, Decision Tree, Random Forest, XG Boost
- With all the features we got AUROC score of 0.77
- From the feature selector we found age,education qualification, duration , day, month to have high feature importance and model was built on these features
- We got a 0.79 AUROC Score, by Random Forest. So we can consider Random Forest is best suitable model