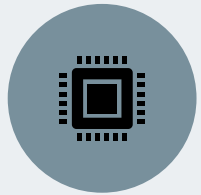# Agenda

Introduction

Domain Knowledge

EDA

Preprocessing

Modeling

Evaluation

Conclusions

# Problem

In the Insurance Industry…

50% of all policies

→

Can be Mispriced by more or Less than 10%

(up to 50%)

# Introduction

Datasets

- Historical dataset of around 400K input samples of Auto Insurance policies
  - 64 input features (ex. vehicle make year, vehicle performance, usage, miles to work etc.)
  - Most of these input features were categorical.
- Testing dataset of 330 policy portfolios
  - Each consisting of at least 1000 policies
  - Included almost all of features above EXCEPT those such as *loss_amount*.

# Introduction

Our Objective

- Predict Missing Loss Amounts in Test Data
- Find *Natural Logarithm (ln_LR)* of the *loss ratio* for each portfolio in testing dataset

$$Total\_Premium = \sum_{i=1}^{N} AnnualPremium(i) \qquad Total\_Losses = \sum_{i=1}^{N} LossAmount(i)$$

**Target: natural log of portfolio loss ratio**

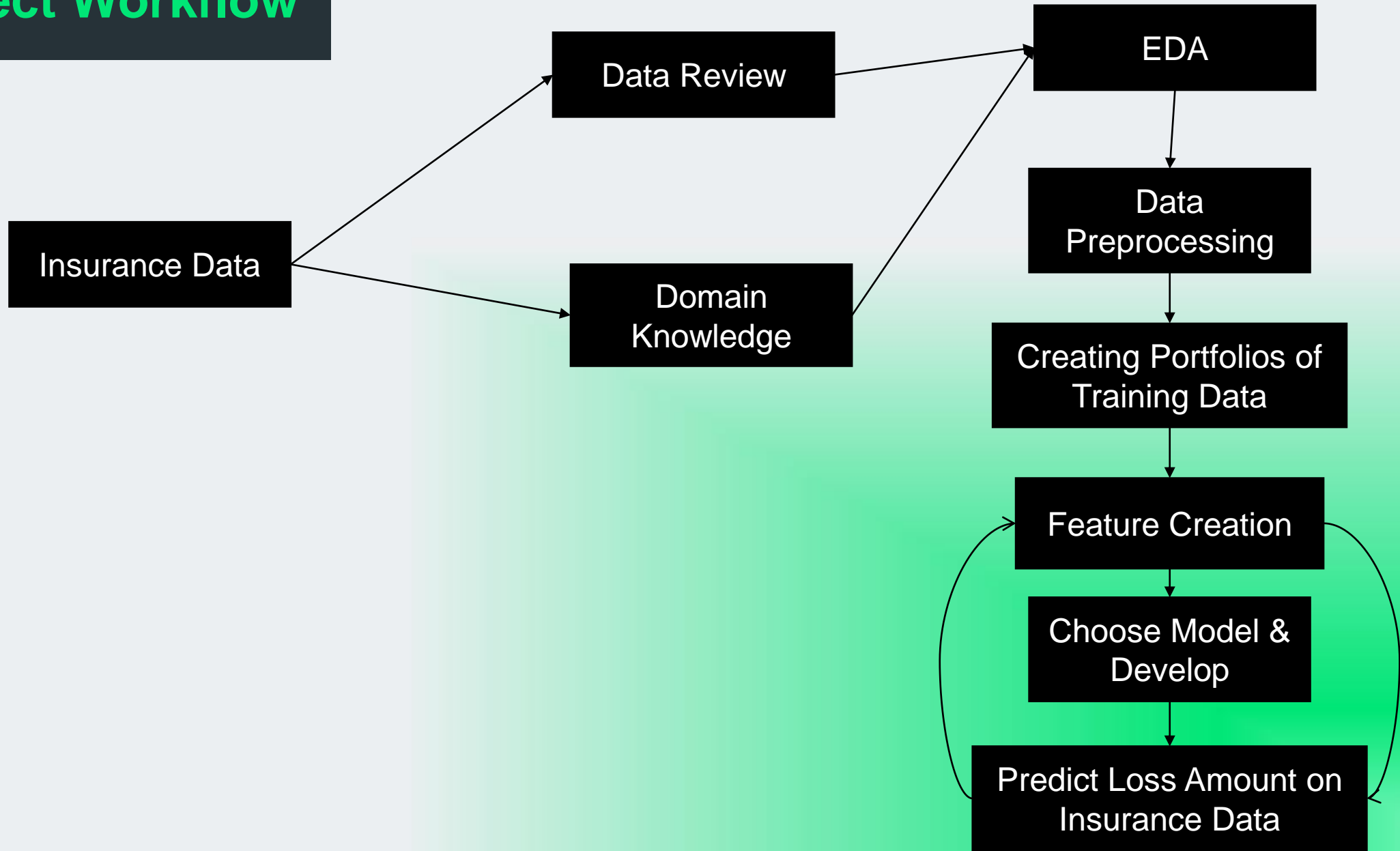$$\ln\_LR = \ln\left(\frac{Total_{Losses}}{Total_{Premium}}\right)$$

# Goal

**Create a competitive advantage for an Insurance Provider by developing a model(s) capable of predicting the loss ratios of policy portfolios, which will enable them to more accurately access the overall risk of a given portfolio as well as make better informed decisions on premium rates.**
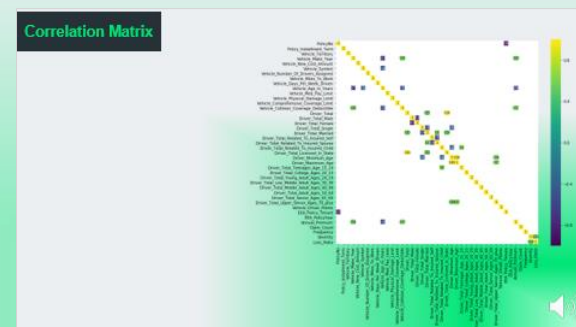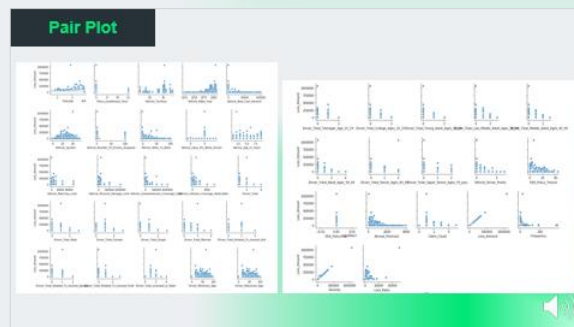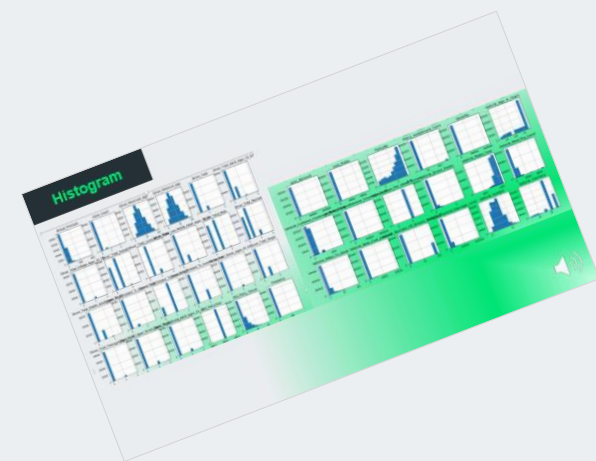
# Domain Knowledge

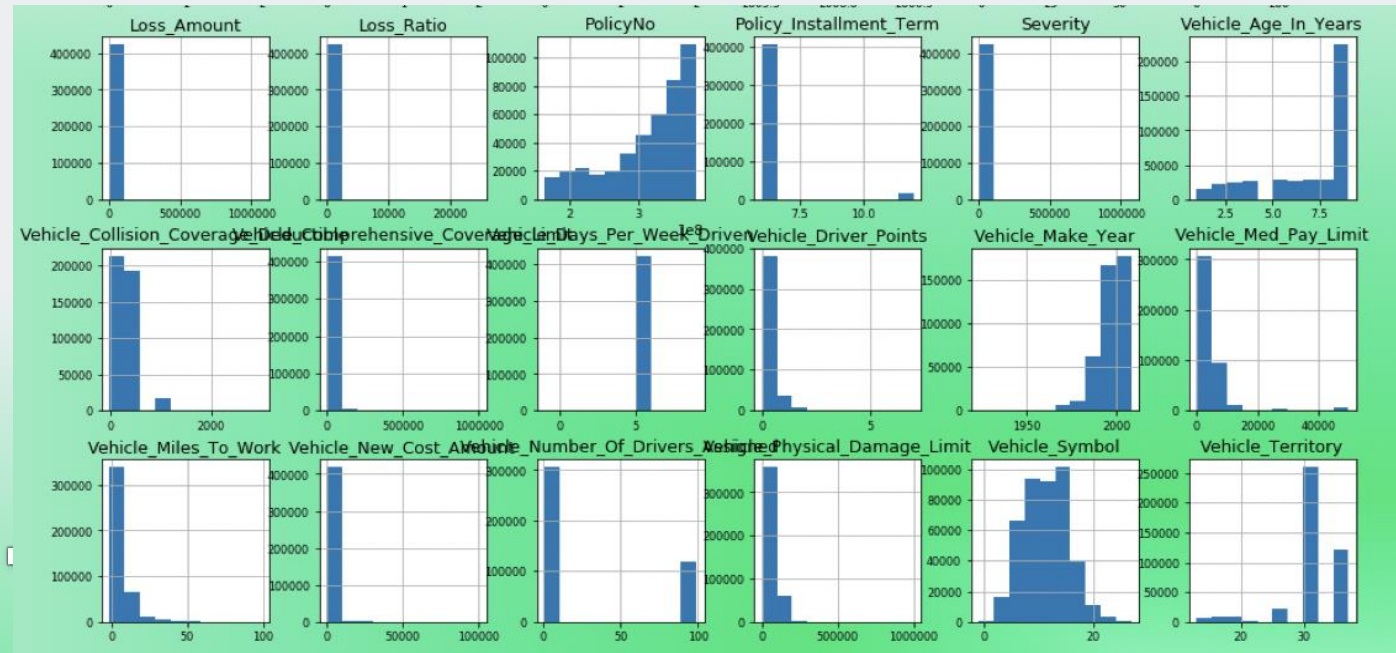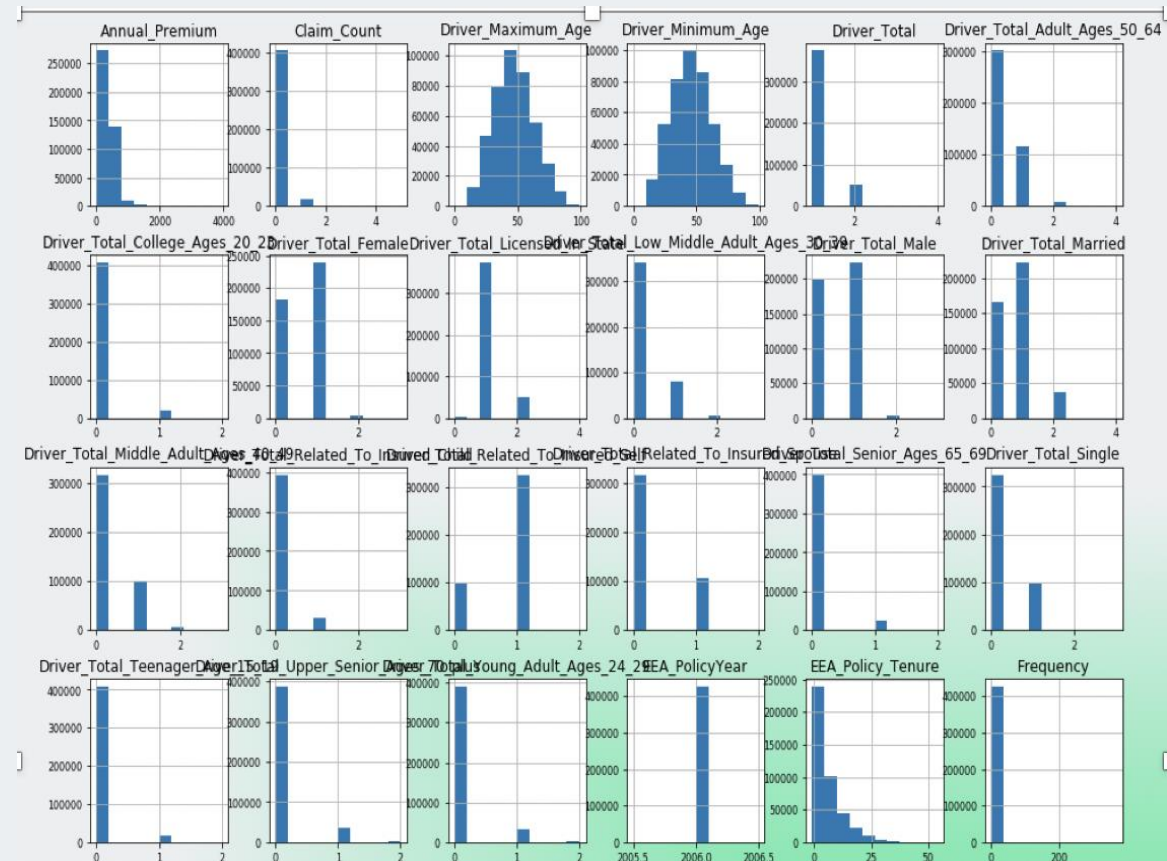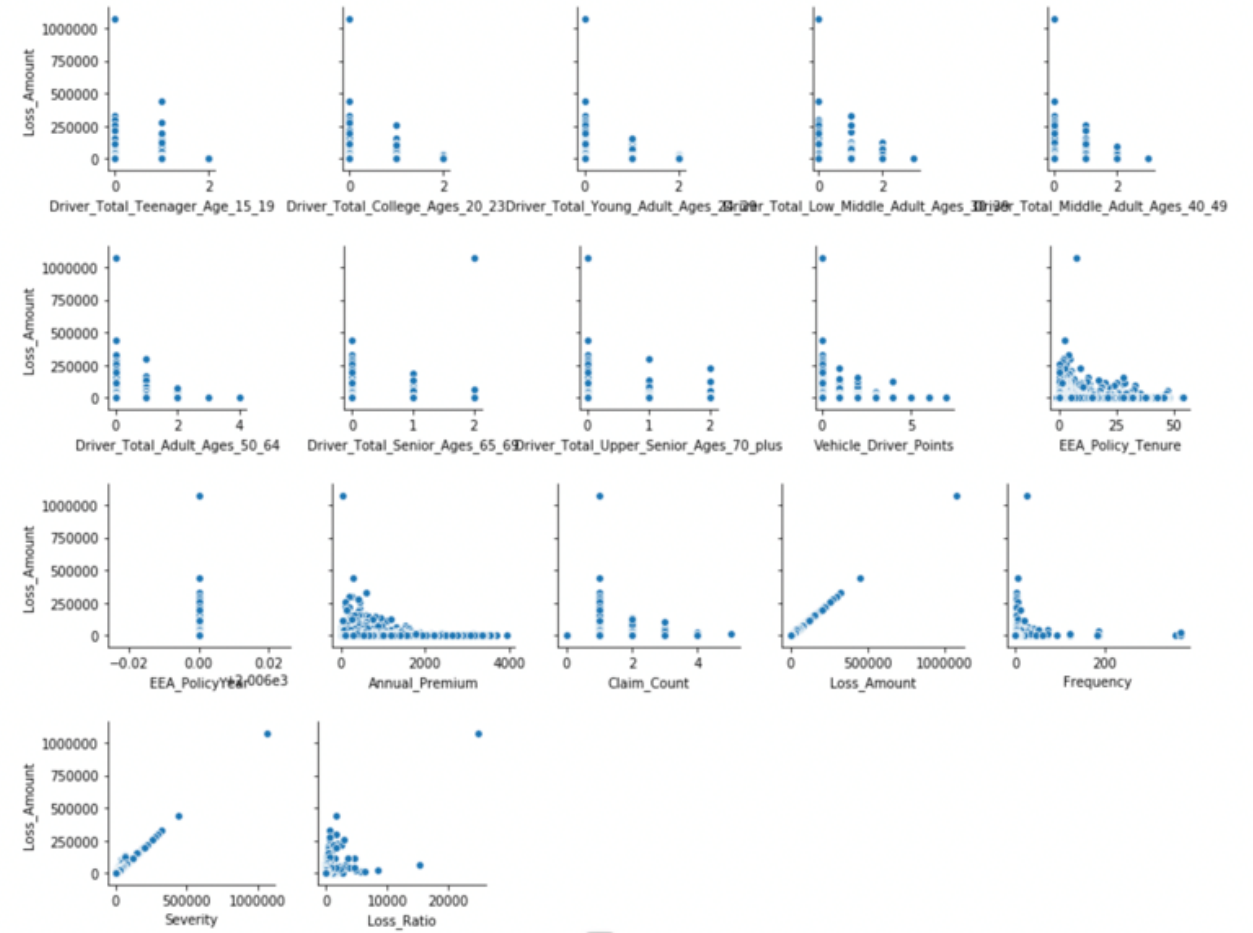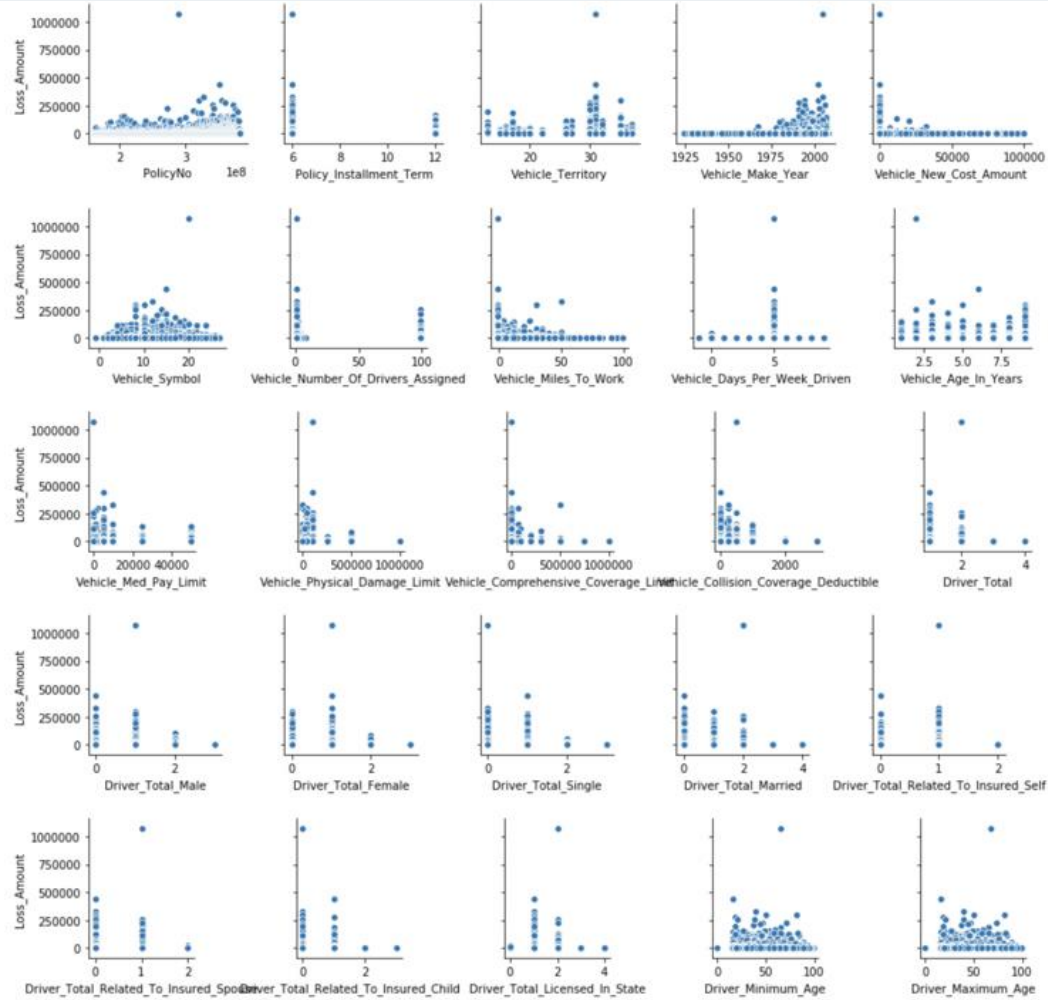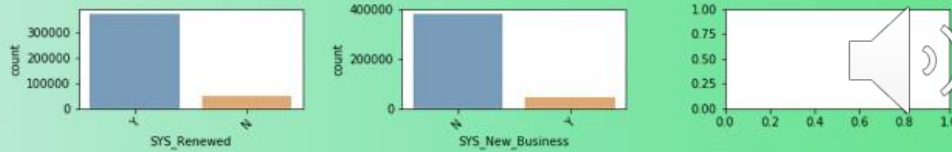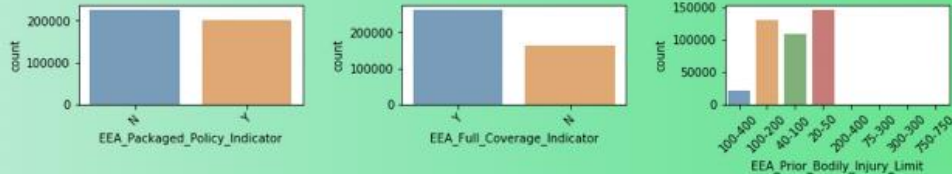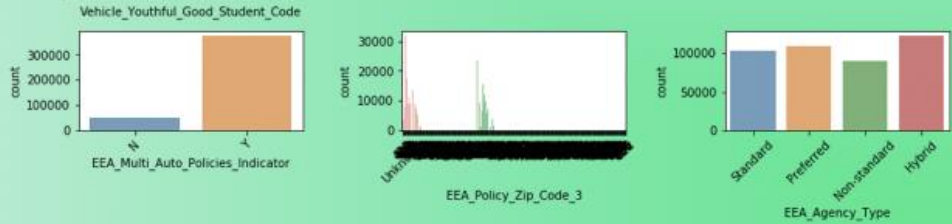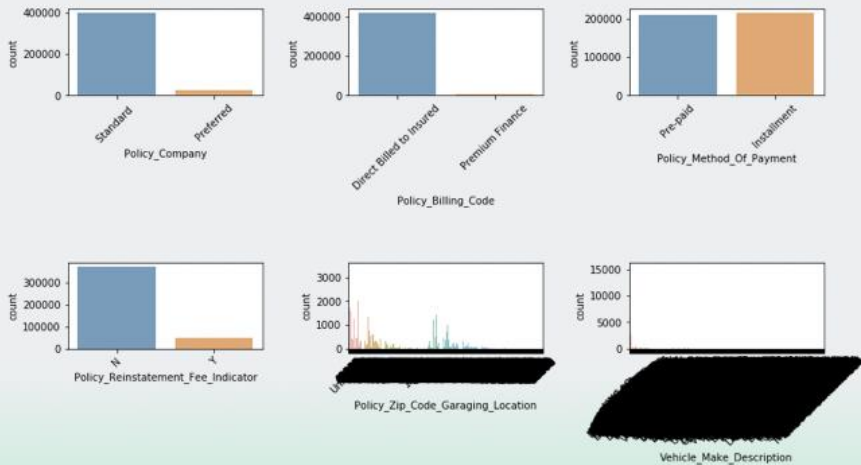| | | Basis for Risk | Data Quality | Description | Predicted Impact | Usage Proposal (Initial - Mike) | Usag Prop (Urm |
|---|---|---|---|---|---|---|---|
| 2 | | | | | | | |
| 3 | 17 | | | | | | 17 |
| 4 | PolicyNo | N/A | | Account Number' of insured person | None, only good as identifier | No | |
| 5 | Policy_Company | N/A | | Provider of Insurance Service | likely not, could be significant depending on how company reports claims | No | |
| 6 | Policy_Installment_Term | MATH | | Duration of term | Could be an issue with comparing policies, would need to evaluate closer the spread of data | Maybe | |
| 7 | Policy_Billing_Code | N/A | | Possibly the method of billing used for customer | Very Likely is NOT significant, values are largely the same from data sampled | No | |
| 8 | Policy_Method_Of_Payment | MATH | | Rather the customer pays the full premium upfront or in installments during the term | This could potentially have a bearing on loss, would be worth investigating more | Maybe | |

# Exploratory Data Analysis

# Histogram

# Pair Plot

# Correlation Matrix

# Bar Chart Comparison

# Data Preprocessing

# Modeling Techniques

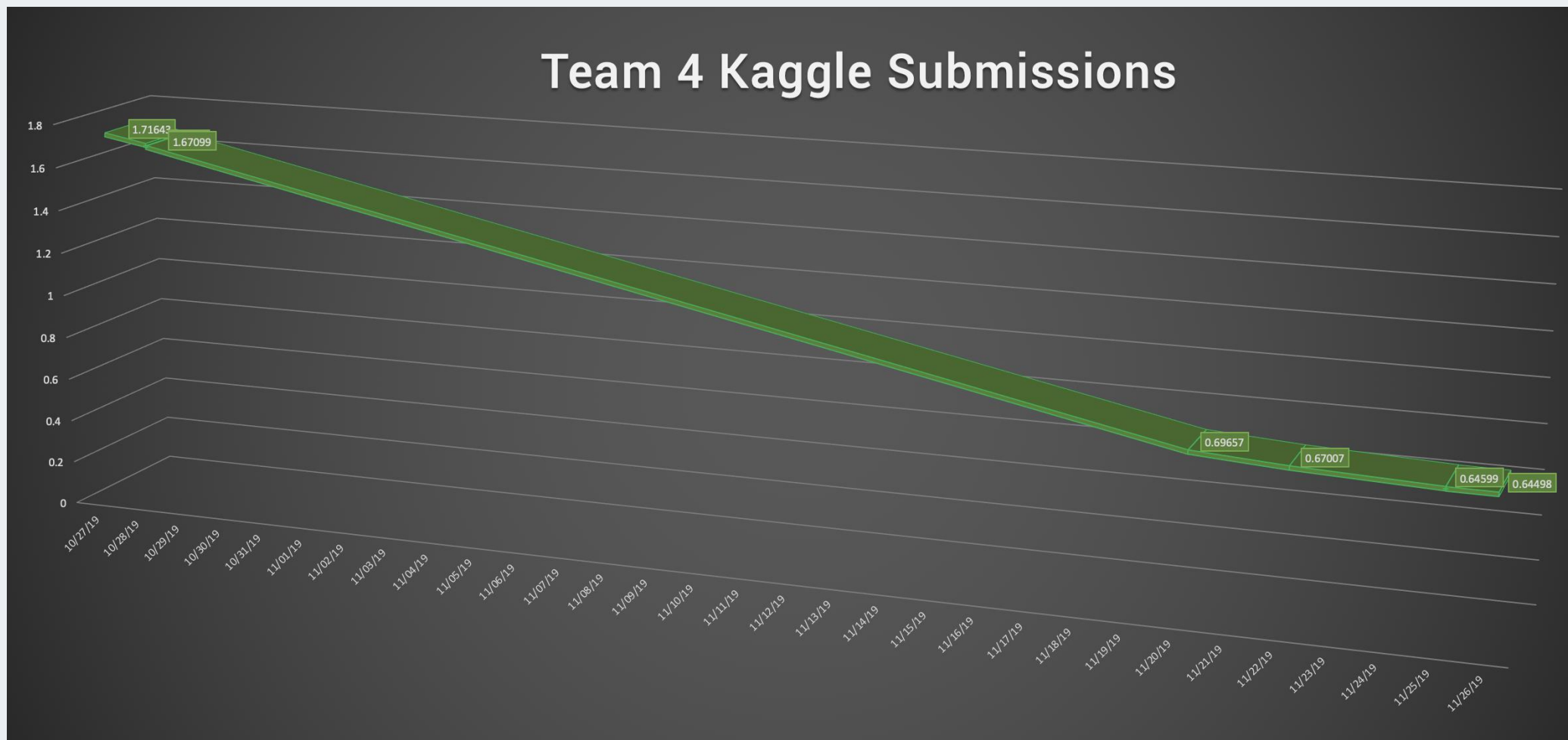Decision trees    Random Forest    Neural Networks    Ada Boost    XGBoost

# Evaluation Methods

- Splitting training portfolios into training and holdout sets

- Calculated the Mean Absolute Error (MAE) on these new splits

- Predicted loss amounts on our test portfolios using our model

- Calculated the loss ratio and corresponding Natural Logarithm on the testing portfolios

- Submitted results on Kaggle which provided a rank based on the MAE of our submission

# Evaluation Methods



Team 4 Kaggle Submissions

# Conclusions

Statistical modeling offers a unique competitive advantage in many industries…

It offers a notable advantage for the Insurance Industry.

# Thank You!