


```
print(train.shape)
print(train_out.shape)

(381109, 12)
(370769, 12)
```

We found roughly 10,000 rows of outliers (less than 3% of the data), so we remove those rows so that our dataset contains no outliers. We use this new dataframe, train_out, for further analysis.

Data Preparation for Modeling

Define predictor/outcome variables and split data frame into predictor/outcome sets.

```
In [5]: #define array of predictors and outcome
predictors = ['Gender', 'Age', 'Driving_License', 'Region_Code',
              'Previously_Insured', 'Vehicle_Age', 'Vehicle_Damage', 'Annual_Premium',
              'Policy_Sales_Channel', 'Vintage']
outcome = 'Response'

In [9]: #create a split of data for training and testing
training, valid = train_test_split(train_out, train_size = 0.6, random_state=1)
print(training.shape, valid.shape)

(222473, 12) (148316, 12)

In [10]: #split train_out_df into two sets
train_X = training[predictors]
train_y = training[outcome]

valid_X = valid[predictors]
valid_y = valid[outcome]
```

Random Oversampling

As we saw during the EDA step, our response variable is heavily imbalanced. If we leave the data imbalanced, our models will likely just predict the majority class which would not be good for future data. Therefore, we want to balance this out so we have the same number of yes's as no's.

```
In [11]: #oversample the minority class
trainOver = RandomOverSampler(sampling_strategy='minority')

In [12]: #apply the oversample
X, y = trainOver.fit_resample(train_X, train_y)
print(X.shape, y.shape)

(390584, 10) (390584,)
```

```
In [13]: #display the count of response variables
y.value_counts()
```

```
Out[13]: 1    195292
         0    195292
         Name: Response, dtype: int64
```

We now have a balanced response variable. We are now ready to begin modeling. As mentioned earlier, we will use the split training set for modeling and model validation. The final model will be used to predict the values of the test set provided by the client.

Modeling

Our approach will be to consider multiple models. We will run a baseline model of the various models. The model which outputs the best baseline performance will then be tuned for optimal performance and that model will be submitted as the final model. The models we will consider are Logistic Regression, Neural Networks, Linear Discriminant Analysis, Random Forests and Decision Trees.

Logistic Regression

```
In [16]: #Run the Logistic Regression model
logReg = LogisticRegressionCV(penalty='l2', solver='lbfgs', cv=10, max_iter=500, n_jobs=-1)
logReg.fit(X, y)

Out[16]: LogisticRegressionCV(cv=10, max_iter=500, n_jobs=-1)

In [17]: #Analyze the Logistic Regression model
classificationSummary(valid_y, logReg.predict(valid_X))

Confusion Matrix (Accuracy 0.6425)

      Prediction
Actual      0      1
0      77528 52814
1      508 17466
```

This model yields terrible results. 64% is not what we are shooting for. Logistic regression is eliminated because of this.

Neural Network

```
In [31]: #Run the Neural Network model
nn = MLPClassifier(hidden_layer_sizes=(11), activation='logistic', solver='lbfgs', random_state=1)
nn.fit(X, y)

Out[31]: MLPClassifier(activation='logistic', hidden_layer_sizes=11, random_state=1,
                    solver='lbfgs')

In [34]: #Analyze the Neural Network model
classificationSummary(valid_y, nn.predict(valid_X))

Confusion Matrix (Accuracy 0.8788)

      Prediction
Actual      0      1
0      130342    0
1      17974    1
```

Neural network yields much better accuracy than Logistic Regression, but as you can see, it doesn't predict any positive values. For the insurance company, this would be useless.

Linear Discriminant Analysis

```
In [35]: #Run the LDA model
lda = LinearDiscriminantAnalysis()
lda.fit(X, y)

Out[35]: LinearDiscriminantAnalysis()
```

```
In [36]: #Analyze the LDA model
classificationSummary(valid_y, lda.predict(valid_X))

Confusion Matrix (Accuracy 0.6410)

      Prediction
Actual      0      1
0      77528 52814
1      430 17544
```

Just like Linear Regression, LDA does not yield great accuracy and therefore will not be used.

Decision Tree

```
In [37]: #Run the Classification Tree model
classTree = DecisionTreeClassifier(random_state=1)
classTree.fit(X, y)

Out[37]: DecisionTreeClassifier(random_state=1)

In [38]: #Analyze the Classification Tree model
classificationSummary(valid_y, classTree.predict(valid_X))

Confusion Matrix (Accuracy 0.8284)

      Prediction
Actual      0      1
0    117914 12428
1    13021 4953
```

Our decision tree model performed relatively well. It yielded an accuracy of 83% while predicting both positive and negative values.

Random Forest

```
In [18]: #Run the Random Forest model
rf = RandomForestClassifier(n_estimators=500, random_state=1, n_jobs=-1)
rf.fit(X, y)

Out[18]: RandomForestClassifier(n_estimators=500, n_jobs=-1, random_state=1)

In [19]: #Analyze the Random Forest model
classificationSummary(valid_y, rf.predict(valid_X))

Confusion Matrix (Accuracy 0.8482)

      Prediction
Actual      0      1
0    121147 9195
1    13115 4659
```

Random Forest turned out to be our best performing model at 85%. Therefore, this will be our chosen model to fine tune and try to improve.

Model and Hyperparameter Tuning

Grid Searching

```
In [20]: #create an initial guess for random forest parameters
param_grid = {
    'max_depth': (None, 10, 15, 20, 25),
    'min_samples_split': (10, 20, 30, 40, 0),
    'min_impurity_decrease': (0, 0.005, 0.001, 0.005, 0.01),
}

gridSearch = GridSearchCV(rf, param_grid, cv=5, n_jobs=-1)
gridSearch.fit(X, y)

print('Initial score: ', gridSearch.best_score_)
print('Initial parameters: ', gridSearch.best_params_)

Initial score: 0.8868053955651168
Initial parameters: {'max_depth': 25, 'min_impurity_decrease': 0, 'min_samples_split': 10}
```

```
In [15]: #adapt initial grid search
param_grid = {
    'max_depth': list(range(23, 27)),
    'min_samples_split': list(range(8, 12)),
    'min_impurity_decrease': (0, 0.001, 0.0001),
}

gridSearch = GridSearchCV(rf, param_grid, cv=5, n_jobs=-1)
gridSearch.fit(X, y)

print('Improved score: ', gridSearch.best_score_)
print('Improved parameters: ', gridSearch.best_params_)

Improved score: 0.8938000506943446
Improved parameters: {'max_depth': 25, 'min_impurity_decrease': 0, 'min_samples_split': 8}
```

```
In [21]: #use the new parameters to rerun the random forest model
rf = RandomForestClassifier(n_estimators=500, random_state=1, max_depth=26, min_impurity_decrease=0, min_samples_split=8,
                           n_jobs=-1)
rf.fit(X, y)
```

```
Out[21]: RandomForestClassifier(max_depth=26, min_impurity_decrease=0,
                               min_samples_split=8, n_estimators=500, n_jobs=-1,
                               random_state=1)
```

```
In [22]: #Analyze the Random Forest model
classificationSummary(valid_y, rf.predict(valid_X))

Confusion Matrix (Accuracy 0.7891)

      Prediction
Actual      0      1
0    104969 25373
1      5900 12074
```

Results and Final Model

We ran a variety of models including Logistic Regression, Neural Network, LDA, Decision Trees and Random Forests. Of these models, Random Forests performed the best with a baseline accuracy of 85%. After tuning the model and finding the best parameters, we were able to obtain a model that will be useful. While the accuracy decreased to 79%, we were able to drastically reduce the number of false negatives (predicted negative but actually positive). This means our model will minimize the number of customers who do not receive the advertisement for car insurance but would actually be interested. In the next section, we identify the customers who we will send advertisements to.

Discussions and Conclusions

Predict the customers who will be interested in car insurance

```
In [27]: #import the test set and view the first few rows
test = pd.read_csv('C:/Users/Leonard/Desktop/Insurance_Test.csv')
test.head()
```

```
Out[27]:
```

	id	Gender	Age	Driving_License	Region_Code	Previously_Insured	Vehicle_Age	Vehicle_Damage	Annual_Premium	Policy_Sales_Channel	
	0	381110	Male	25	1	11.0	1	< 1 Year	No	35786.0	1
	1	381111	Male	40	1	28.0	0	1-2 Year	Yes	33762.0	0
	2	381112	Male	47	1	28.0	0	1-2 Year	Yes	40050.0	1
	3	381113	Male	24	1	27.0	1	< 1 Year	Yes	37356.0	1
	4	381114	Male	27	1	28.0	1	< 1 Year	No	59097.0	1

```
In [28]: #map gender
test['Gender'] = test['Gender'].replace({'Male': 0, 'Female': 1})
test.head()
```

```
Out[28]:
```

	id	Gender	Age	Driving_License	Region_Code	Previously_Insured	Vehicle_Age	Vehicle_Damage	Annual_Premium	Policy_Sales_Channel	
	0	381110	0	25	1	11.0	1	< 1 Year	No	35786.0	1
	1	381111	0	40	1	28.0	0	1-2 Year	Yes	33762.0	0
	2	381112	0	47	1	28.0	0	1-2 Year	Yes	40050.0	1
	3	381113	0	24	1	27.0	1	< 1 Year	Yes	37356.0	1
	4	381114	0	27	1	28.0	1	< 1 Year	No	59097.0	1

```
In [29]: #map vehicle age
test['Vehicle_Age'] = test['Vehicle_Age'].replace({'< 1 Year': 0, '1-2 Year': 1, '> 2 Years': 2})
test.head()
```

```
Out[29]:
```

	id	Gender	Age	Driving_License	Region_Code	Previously_Insured	Vehicle_Age	Vehicle_Damage	Annual_Premium	Policy_Sales_Channel	
	0	381110	0	25	1	11.0	1	0	No	35786.0	1
	1	381111	0	40	1	28.0	0	1	Yes	33762.0	0
	2	381112	0	47	1	28.0	0	1	Yes	40050.0	1
	3	381113	0	24	1	27.0	1	0	Yes	37356.0	1
	4	381114	0	27	1	28.0	1	0	No	59097.0	1

```
In [31]: #map vehicle damage
test['Vehicle_Damage'] = test['Vehicle_Damage'].replace({'No': 0, 'Yes': 1})
test.head()
```

```
Out[31]:
```

	id	Gender	Age	Driving_License	Region_Code	Previously_Insured	Vehicle_Age	Vehicle_Damage	Annual_Premium	Policy_Sales_Channel	
	0	381110	0	25	1	11.0	1	0	0	35786.0	1
	1	381111	0	40	1	28.0	0	1	1	33762.0	0
	2	381112	0	47	1	28.0	0	1	1	40050.0	1
	3	381113	0	24	1	27.0	1	0	1	37356.0	1
	4	381114	0	27	1	28.0	1	0	0	59097.0	1

```
In [33]: #use the test data to predict which customers will be interested in car insurance.
test['Response'] = rf.predict(test[predictors])
```

```
In [34]: #display the number of customers who will be interested vs not
test['Response'].value_counts()
```

```
Out[34]: 0    95314
         1    31523
         Name: Response, dtype: int64
```

After applying the model to the unseen data, we identified 31,523 customers who we should send mailers to. We expect that half of those should be interested. We also expect that only a small subset of the 95,514 customers not sent mailers would have been interested. We therefore maximize the number of customers who will want car insurance while minimizing the cost of mailing advertisements to the customers.

```
In [ ]:
```