BOOZ ALLEN HAMILTON INTERNAL

BAH Machine Learning Presentation
Team 3
Analytics Avengers
Members: Ahmed Sodeinde
Erick Orellana Morales
Rami Tello

April 2022



TEAM 3: Project Car Insurance Claim Prediction

Our Official Project for the BAH
 Machine Learning Program

The Team



AHMED SODEINDE Machine Learning Engineer



ERICK ORELLANA MORALES Machine Learning Engineer



RAMI TELLOMachine Learning
Engineer

Project

Predicting if a car insurance claim will be filed within 6 months based on data related to the policy holder

- 1. Introductions
- 2. Project Objectives
- 3. <u>Project Overview</u>
- 4. <u>Architecture/Design diagram</u>
- 5. Tools/Software Environment used
- 6. <u>Challenges/Issues/Problems/Your</u> <u>stories</u>
- 7. What have you learned/learned?
- 8. What else can we do on the project? future plan- any design and implementation ideas?
- 9. Demo
- 10. Questions and answers

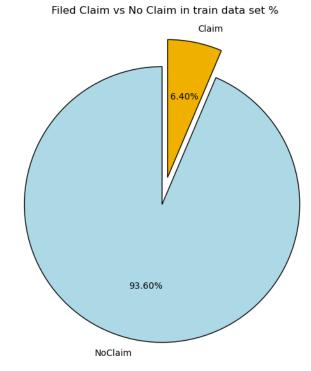


Project Goals & Business Objectives

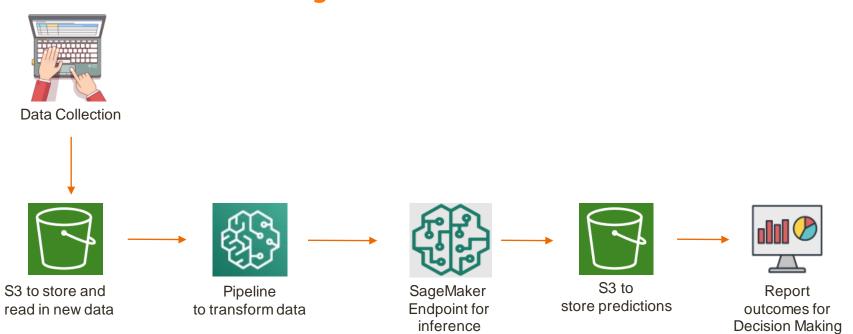
- Project Summary:
 - Building a model that could be used to predict if a car insurance policy holder is likely to submit a claim over the next 6 months primarily based on information about the vehicle
- Business Use Case:
 - Having a prediction for the claims that may be filed can help make predictions about accepting new policy holders/renewing policies, pricing for policies, and making decisions about resources needed to handle the volume of claims
- Out of scope:
 - The type of claim (collision with another vehicle, collision with a stationary object etc)
 - Predicted cost of the claim
 - The type of driver
- Result:
 - It is difficult to predict whether a car insurance claim will be filed based on attributes in the available data.
 These attributes are largely about the car and does not have any information about the driver.

The Dataset

- CSV files
- Approximately 52,500 observations
- 41 attributes
- No duplicates or null values
- Imbalanced Dataset



Architecture and Design



Technologies and Environment Used

- Python based
 - Libraries used:
 - pandas
 - numpy
 - sklearnimblearn
 - seaborn
 - SageMaker
 - boto3
 - XGBoost
- AWS
 - O IAM
 - S3
 - SageMaker
 - Notebooks
 - Studio
 - Registry
 - Endpoints





Amazon SageMaker

The Baseline Models (Dummy and Logistic)

- Baseline Model (Dummy Classifier)
 - Benchmark

Summary of performance:						
The model's accuracy is: The model's F1 score is:						

- Logistic Regression
 - Recursive Feature Elimination

	precision	recall	f1-score	support
0	0.50	1.00	0.66	11446
1	0.00	0.00	0.00	11586
accuracy			0.50	23032
macro avg	0.25	0.50	0.33	23032
weighted avg	0.25	0.50	0.33	23032

The Models (XGBoost)

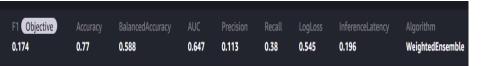
Pre-Hyperparameter Tuning

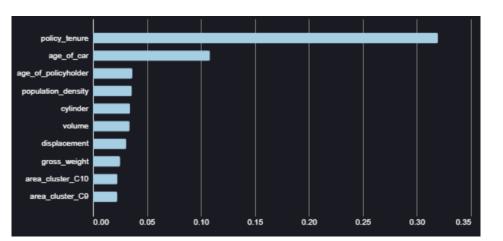
• With Hyperparameter Pruning	g
-------------------------------	---

	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.94	0.98	0.96	10963	0	0.94	0.99	0.96	10963
1	0.16	0.04	0.07	756	1	0.16	0.04	0.06	756
accuracy			0.92	11719	accuracy			0.93	11719
macro avg	0.55	0.51	0.51	11719	macro avg	0.55	0.51	0.51	11719
weighted avg	0.89	0.92	0.90	11719	weighted avg	0.89	0.93	0.90	11719

The Models (AutoML - Autopilot)

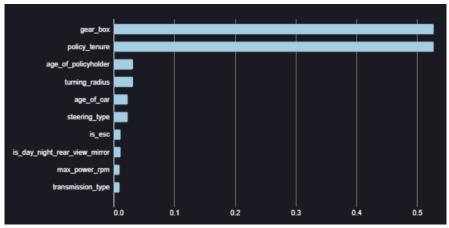
Ensemble





• Hyperparameter Optimization





Demo

Booz | Allen | Hamilton

Experiment Conclusion

H₀: Car features cannot be used to predict claims being filed H_a: Car features can be used to predict claims being filed

- Failed to reject the null Hypothesis
 - The data set did not have the signal to prove that filing claims can be predicted based on car features

Challenges and Learning

Challenges

Data Quality and Lineage

- The data used from Kaggle doesn't have a verifiable source to ensure the quality of the data
- ML Model performance is limited to the quality of the data
- Uncertainty regarding data completeness.

Lesson Learned

- Importance of balancing data
- Importance of data quality and lineage

Opportunity for Enhancement:

Gathering more data for model building and tuning:

- Increasing the number of observations of drivers that filed a claim in the dataset
- Increasing the number of features
- Further feature engineering

Q&A – Ask us anything!



Appendix: Demo Backup

```
Send test data to best model built by AutoML AutoPilot
[146]: # xgb_predictor = tuner_linear.deploy(initial_instance_count=I, instance_type='ml.m4.xlarge')
        xgb_tuned_predictor = sagemaker.predictor.Predictor(
           endpoint name="SageMakerEndpoint-84132023".
           sagemaker_session=sagemaker.Session(),
           serializer=sagemaker.serializers.CSVSerializer()
 [149]: s3_input_validate_df = pd.read_csv('s3://techexcellence.ml.project.team3/CarInsuranceClaim/data/test.csv')
                                                                                                                                                                                                                                  四个少占早ま
       s3_input_validate_df.head()
        policy_id policy_tenure age_of_car age_of_policyholder area_cluster population_density make segment model fuel_type
                                                                                                                                             max_power engine_type airbags is_esc is_adjustable_steering is_tpms is_parking_sensors is_parking_camera rear_brai
                                                                                                                           max_torque
       0 ID58593
                      0.341732
                                     0.00
                                                   0.586538
                                                                    C3
                                                                                                                  CNG 60Nm@3500rpm 40.36bhp@6000rpm
                                                                                                                                                             Engine
                                                                                                                                                            K Series
                       0.307241
                                                   0.442308
                                                                                                                 Petrol 113Nm@4400rpm 88.50bhp@6000rpm
        1 ID58594
                                     0.13
                                                                                                                                                            Dual jet
       2 ID58595
                      0.327924
                                     0.12
                                                   0.451923
                                                                    C8
                                                                                                                        91Nm@4250rpm 67.06bhp@6500rpm
                                                                                                                                                            1.0 SCe
                                                                                                                                                                                                                                            Yes
       3 ID58596
                      0.782654
                                     0.01
                                                   0.461538
                                                                    C5
                                                                                   34738
                                                                                                                  CNG 60Nm@3500rpm 40.36bhp@6000rpm
                                                                                                                                                                                                                           Yes
                                                                                                                                                             Engine
                                                                                                                                                          FSD Petrol
       4 ID58597
                      1.233404
                                    0.02
                                                   0.634615
                                                                    C5
                                                                                  34738
                                                                                                                  CNG 60Nm@3500rpm 40.36bhp@6000rpm
                                                                                                                                                                                                                           Yes
                                                                                                                                                             Engine
-[147_ predictions_list = []
        for i in chunker(s3_input_validate_df, 2500):
           test_data_array = i.values
           xgb_tuned_predictor.serializer = csv_serializer
           predictions = xgb_tuned_predictor.predict(test_data_array).decode('utf-8')
           predictions_array = np.fromstring(predictions[1:], sep='\n')
           predictions_list.append(predictions_array)
           print(predictions_array.shape)
       y_hat_2 = np.concatenate(predictions_list, axis=Mone)
       y_hat_2
 [148]: output_df = pd.DataFrame(data=(s3_input_validate_df('policy_id'), y_hat_2)).T
       output_df.columns = ['policy_id', 'is_claim_hat']
       output_df.head(10)
          policy_id is_claim_hat
       0 ID58593
                          0.0
       1 ID58594
                          0.0
       2 ID58595
                          0.0
       3 ID58596
                          0.0
       4 ID58597
                          0.0
                          0.0
       5 ID68598
       6 ID58599
                           1.0
       7 ID58600
       B ID58601
                          0.0
       9 ID58602
```