

Prediction of Successful Transaction Billing

Svetlana Bogdanova
December 22, 2024

Table of Contents

- Definition3**
 - Domain Background.....3**
 - Problem Statement.....3**
 - Metrics3**
 - Test Dataset3
 - Performance Metrics4
 - ROC Curve4
 - Models Comparison.....4
 - Justification of Metrics.....4
- Analysis and Preparation4**
 - Data Exploration4**
 - Exploratory Visualization5**
 - Transactions status by categories.....5
 - Class Imbalance6
 - Outlier Detection6
 - Data Preprocessing6**
 - Repurposed and Synthetic Features for Transaction Billing:.....6
 - Removal of Sensitive Data7
 - Outlier Detection7
 - Drop identifier column7
 - Duplicate Records.....7
 - Algorithms and Techniques9**
 - Logistic Regression.....9
 - AutoML Solution9
 - XGBoost (Extreme Gradient Boosting)9
 - Model Tuning.....10
 - Frameworks10**
 - Amazon SageMaker10
 - Local Docker Container.....10
- Benchmark Model10**
 - Benchmark Model Development.....11**

Data Preparation.....	11
Model Building.....	11
Model Training.....	11
Model Registration	12
Batch Predictions	12
Model Evaluation.....	12
Comparison of Logistic Binary Classifier and Linear Regression (Threshold = 0.5).....	13
<i>Implementation.....</i>	<i>14</i>
Sequence of steps to train and evaluate models.....	14
List of models trained.....	15
Linear Learner framework in SageMaker	15
AutoGluon in local docker container	15
XGBoost models in local container	16
Model Performance Comparison.....	16
Summary of the performance metrics results.....	16
ROC curves Interpretation.....	17
Discussion	18
<i>Transaction recovery solution architecture</i>	<i>18</i>
<i>Conclusions.....</i>	<i>20</i>
<i>Acknowledgement.....</i>	<i>21</i>
<i>Resources</i>	<i>21</i>

Definition

Domain Background

Subscription-based businesses have become increasingly popular across industries, including entertainment, software, fitness, and e-commerce. However, one major challenge these businesses face is churn due to billing failures. Billing failures can occur for various reasons, such as expired payment methods, insufficient funds, or technical errors, and often lead to unintended subscription cancellations. According to industry reports, involuntary churn due to billing issues accounts for up to 20-40% of total subscription cancellations, representing a significant revenue loss for businesses. [1]

Academic research supports the application of ML in addressing billing-related issues. Studies in predictive analytics and churn modeling demonstrate how machine learning techniques, such as classification algorithms and time-series analysis, can enhance customer retention strategies [2].

This project implements ML predictions to address the problem of stopped subscriptions due to billing failures

Using a dataset originally meant for analyzing customer churn in banks, we adapt it to simulate a subscription business scenario. We added synthetic data fields such as engagement levels and billing retry counts to reflect the subscription model. The ultimate goal is to train machine learning model that can identify customers worth targeting with billing retries and to build API to access this model, helping businesses reduce churn and improve their revenue.

Problem Statement

Billing failures in subscription-based models lead to significant revenue loss. Identifying which transactions are likely to recover after an initial billing failure is crucial for minimizing churn. The challenge is to develop a machine learning model that accurately predicts recovery likelihood, optimizing billing retries and reducing unnecessary operational costs.

Metrics

This project is a **binary classification** problem aimed at predicting the likelihood of a transaction recovery following a failed billing attempt. The following techniques were used to evaluate model performance

Test Dataset

The trained model was applied to a separate test dataset that was not used during the training phase. This dataset serves as an unseen data set to evaluate the model's performance.

Performance Metrics

Various performance metrics were calculated to assess the model's effectiveness. These metrics included:

1. **Accuracy.**
2. **Precision**
3. **Recall**
4. **F1-Score**
5. **AUC-ROC**
6. **Confusion Matrix.**

ROC Curve

The Receiver Operating Characteristic (ROC) curve was plotted to visualize the model's performance across different thresholds, highlighting the trade-off between true positive rate and false positive rate.

Models Comparison

The results were compared with the performance of previous models to identify any improvements or areas needing further tuning.

Justification of Metrics

Precision and recall are prioritized due to the financial implications of retrying transactions. F1 Score and ROC-AUC provide a balanced and comprehensive evaluation of the model's performance in the context of transaction recovery. This thorough evaluation process ensured that the model's performance was accurately assessed

Analysis and Preparation

This step is described and executed in [prepare_data.ipynb](#) notebook

Data Exploration

The primary dataset for this project is **the Bank Customer Churn Dataset** [3], which provides information about customer demographics, account behavior, and churn status. While the dataset is designed to predict customer churn in banking, it will be adapted to represent a subscription context by modifying the target variable to indicate recovery outcomes after billing failure.

First rows of dataset:

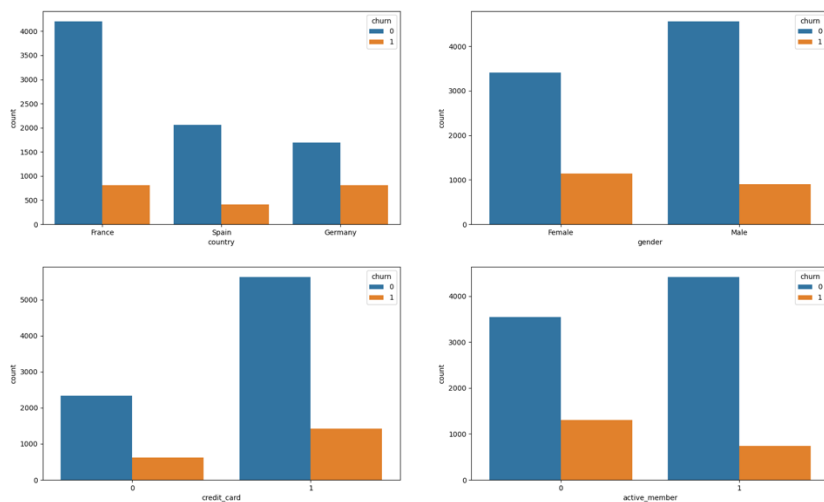
	customer_id	credit_score	country	gender	age	tenure	balance	products_number	credit_card	active_member	estimated_salary	churn
0	15634602	619	France	Female	42	2	0.00	1	1	1	101348.88	1
1	15647311	608	Spain	Female	41	1	83807.86	1	0	1	112542.58	0
2	15619304	502	France	Female	42	8	159660.80	3	1	0	113931.57	1
3	15701354	699	France	Female	39	1	0.00	2	0	0	93826.63	0
4	15737888	850	Spain	Female	43	2	125510.82	1	1	1	79084.10	0

In preparing the dataset for this project, we identified the abnormalities and characteristics to ensure the data's integrity and suitability for the business goals and model training. There was no missing data, we found a few data fields that had customer sensitive data which is normally not available in training public models and not supplied as a part of billing process

We found that the dataset was missing a few data fields that would be important to predict the probability of successful billing such as number of previously failed attempts, billing amount, type of subscription and overall activity score of the customer such as presence of other subscriptions or transaction frequency. We made our best efforts to synthesize the values for these fields based on already existing data.

Exploratory Visualization

Transactions status by categories.



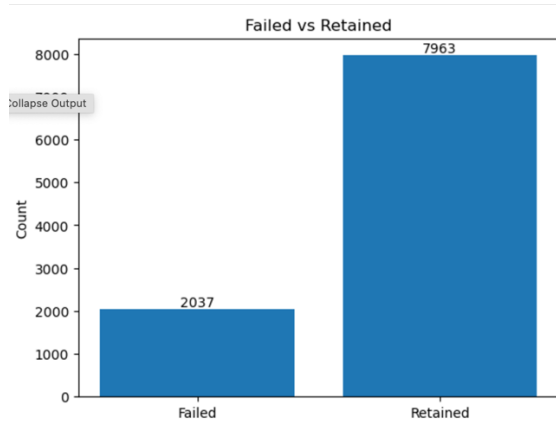
These diagrams suggested potential relationships between customer churn and features such as country, gender, credit card ownership, and active membership.

1. **By country** (France, Spain, Germany). France has the highest number of customers, with fewer churned customers in all countries compared to non-churned.
2. **By gender**, split into churn and non-churn categories. There are slightly more male non-churned customers than female non-churned ones, with a similar trend for churned customers.

3. **By credit card ownership** (0 = no, 1 = yes). A significant majority of customers own a credit card (1), and most of them do not churn.
4. **By active member status** (0 = inactive, 1 = active). Active members (1) are more common, and the number of churned customers is smaller in both active and inactive groups.

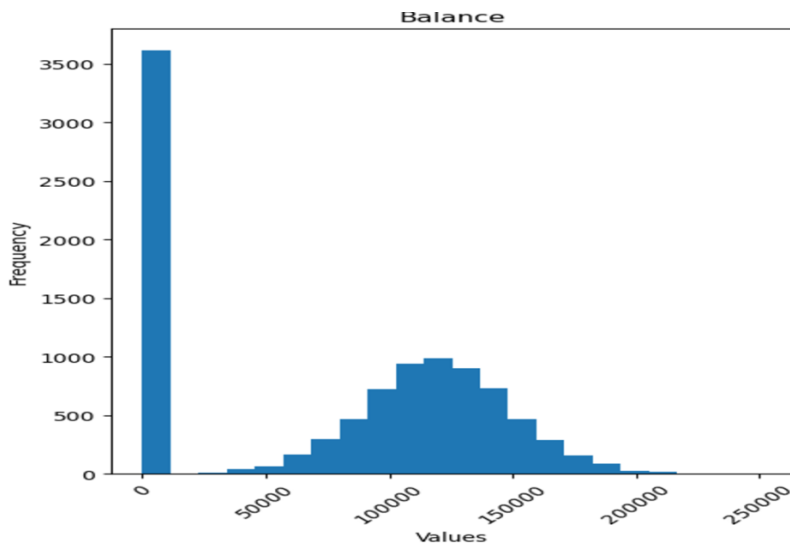
Class Imbalance

The dataset exhibited an imbalance in the target variable (**failed**). This imbalance was addressed by using evaluation metrics such as F1-score and AUC-ROC that are robust to imbalanced datasets.



Outlier Detection

Significant outliers were detected in **balance** (repurposed as **amount**) which were extrapolated to resemble real-life data



Data Preprocessing

Repurposed and Synthetic Features for Transaction Billing:

- **engagement score.** We assume that the was a function of customer being active, number of products used by customer and how long the customer was subscribed:

```
df['engagement_score'] = (
```

```
0.4 * df['tenure'] / df['tenure'].max() +
0.4 * df['products_number'] / df['products_number'].max() +
0.2 * df['active_member'] ) * 100
```

- **subscription_tier.** Categorized into levels such as Basic, Standard, and Premium. Generated based on salary percentiles

```
salary_bins = np.percentile(df['estimated_salary'], [33, 66])
df['subscription_tier'] = pd.cut(
    df['estimated_salary'],
    bins=[-np.inf, salary_bins[0], salary_bins[1], np.inf],
    labels=['Basic', 'Standard', 'Premium']
)
```

- **billing_failure_count.** Represents the number of failed billing attempts for a given transaction. We assume that transactions with higher amount and belonging to customers with lower credit score have more chances to fail. We used Poisson distribution to emulate the value:

```
normalized_credit_score = (df['credit_score'] - df['credit_score'].min()) /
(df['credit_score'].max() - df['credit_score'].min())

df['billing_failure_count'] = np.random.poisson(
    lam=3 * (1 - normalized_credit_score) + 0.5 * (df['amount_due'] >
df['amount_due'].median()),
    size=len(df)
)
```

Removal of Sensitive Data

We dropped unused or sensitive data column and renamed re-purposed columns:

```
df = df.drop(columns=['balance', 'credit_score', 'balance_condition',
'estimated_salary', 'customer_id'])
df = df.rename(columns={'amount_due': 'amount', 'churn': 'failed'})
cleaned_df = df.drop(indices_to_remove)
```

Outlier Detection

Significant outliers were detected in **balance** (repurposed as **amount**) These outliers were capped at \$1 to reduce their influence on model performance while retaining meaningful variability in the data and reflecting real-world scenarios. **balance** field was repurposed to represent the transaction **amount** by dividing by 1000

Drop identifier column. The `customer_id` is a label that uniquely identifies a customer. Including such unique identifiers in a model can cause overfitting

Duplicate Records

A small number of duplicate entries were identified in the dataset that had different churn values

Table 1 summarizes the input features, dropped fields, synthesized fields and actions taken to prepare the dataset to align with real-world constraints

<i>Category</i>	<i>Field Name</i>	<i>Description</i>	<i>Action</i>
Input Features	country	Customer's geographical region (categorical).	Retained
Input Features	gender	Customer's gender (categorical).	Retained
Input Features	age	Customer's age (numeric).	Retained
Input Features	tenure	Duration of the customer relationship (numeric).	Retained
Input Features	products_number	Number of products held by the customer (numeric).	Retained
Input Features	active_member	Whether the customer is actively engaged (binary).	Retained
Input Features	amount_due	Amount owed during billing failures (numeric, repurposed from balance).	Repurposed
Input Features	billing_failure_count	Number of recent billing failures (numeric, synthesized).	Synthesized
Input Features	engagement_score	Overall customer engagement score, derived from tenure, products_number, and active_member.	Synthesized
Input Features	subscription_tier	Subscription level (e.g., Basic, Standard, Premium), based on estimated_salary (before dropping).	Synthesized
Dropped Fields	customer_id	Customer identifier, not useful for modeling.	Dropped
Dropped Fields	credit_score	Credit score, sensitive data unavailable in transactional systems.	Dropped
Dropped Fields	estimated_salary	Estimated salary, sensitive data unavailable in transactional systems.	Dropped

Here is the final data file

	country	gender	age	tenure	products_number	credit_card	active_member	amount	engagement_score	subscription_tier	billing_failure_count	failed
0	France	Female	42	2	1	1	1	1.00	38.0	Standard	2	1
1	Spain	Female	41	1	1	0	1	84.81	34.0	Standard	0	0
2	France	Female	42	8	3	1	0	160.66	62.0	Standard	2	1
3	France	Female	39	1	2	0	0	1.00	24.0	Standard	1	0
4	Spain	Female	43	2	1	1	1	126.51	38.0	Standard	0	0

The dataset was split into training, validation, and test sets for model development and evaluation (60 /20/20)

Here are the line counts for all files (including header record):

```
The file train.csv has 5990 lines.  
The file validate.csv has 1998 lines.  
The file test.csv has 1998 lines.  
The file cleaned_file.csv has 9984 lines.
```

There were 3 categorical columns: ['country', 'gender', 'subscription_tier'] and we used **OneHotEncoder** to encode them as numeric value columns.

We had 3 country values in the dataset, 2 gender values, 3 subscription tiers, we ended up with 8 numeric columns replacing 5 categorical ones. The resulting set was 17 columns that included one target column (**failed**) that we put in the first place for SageMaker to be able to use standard libraries.

Algorithms and Techniques

The objective of this project is to predict the likelihood of successful transaction recovery after billing failure in a subscription-based business. To solve this binary classification problem, a variety of machine learning algorithms and techniques are considered to evaluate their effectiveness in predicting the recovery of failed transactions. The selected techniques balance interpretability, predictive performance, and adaptability to the dataset.

Logistic Regression

It serves as a benchmark model due to its simplicity and interpretability. Despite its basic nature, it can perform well if the relationship between features and the target is approximately linear [4]

AutoML Solution

The AutoML system [5] (e.g., **AutoGluon** in SageMaker) will generate multiple machine learning models, such as:

- Random Forest.
- Gradient Boosting (XGBoost, LightGBM).
- Neural Networks (if dataset size permits).

XGBoost (Extreme Gradient Boosting)

XGBoost (Extreme Gradient Boosting) [6] is a highly efficient and scalable gradient boosting framework used to improve prediction accuracy through optimized decision trees. **Benefits** include high accuracy, efficient handling of missing data, and ability to work with large datasets with many features.

By leveraging these algorithms and techniques, the project aims to develop a robust model for predicting the recovery of failed transactions in a subscription-based business. Each method offers unique strengths, allowing for a comprehensive evaluation of their suitability for this specific binary classification problem.

Model Tuning

Hyperparameter Tuning: We used grid search for the hyperparameters tuning of the XGBoost model to further improve performance [7]

Frameworks

Amazon SageMaker

SageMaker provides pre-built algorithms and easy integration with other AWS services like S3, Lambda, and CloudWatch, simplifying data processing, model monitoring, and inference. Its scalability is a key advantage, as it automatically adjusts to the required compute power for training and deploying models. SageMaker also offers AutoML capabilities.

Local Docker Container

Using Docker containers based on a SageMaker image allows for the portability and isolation of machine learning environments while leveraging the SageMaker ecosystem. A Docker container replicates the SageMaker environment locally, providing the ability to train and test models on any machine that supports Docker. This setup allows for full customization of the environment, including the installation of specific dependencies, while ensuring that the code runs in the same environment as it would in SageMaker. The containerized approach is cost-effective, as it avoids the fees associated with cloud services, and provides flexibility for smaller-scale experiments or when managing resources manually. However, scalability and automated monitoring are more complex compared to a managed service like SageMaker. Here are docker images we used to run Jupiter notebooks and train some models locally

```
Svetlanas-MacBook-Air:figures sab$ docker images
```

REPOSITORY	TAG	IMAGE ID	CREATED	SIZE
sagemaker-custom-image	latest	c3eb3c0265c6	2 months ago	15.2GB
public.ecr.aws/sagemaker/sagemaker-distribution	1.11-gpu	5c13a56bf735	3 months ago	15.2GB

Benchmark Model

Details of implementation of benchmark model is in Jupiter notebook [transaction recovery baseline 2 models notebook.ipynb](#)

The benchmark model for this project is a **Logistic Regression classifier**. Logistic regression is a widely used, interpretable algorithm for binary classification problems and serves as a standard approach in domains such as customer churn prediction. The model predicts the likelihood of recovery after a billing failure based on the adapted features from the Bank Customer Churn Dataset.

We used the same test set to evaluate and compare performance of benchmark model and other models.

Benchmark Model Development

Data Preparation

- Data was prepared according to the steps described in the [Data Preprocessing section](#).
- Training, validation, and test datasets were uploaded to the S3 bucket `saba-transaction-data`.
- Categorical columns were converted to numerical representations.
- The target column, `failed`, was set as the first column in the dataset.
- The header was removed from all CSV files.

Model Building

- I utilized SageMaker's Linear Learner framework to construct the estimator.
- Parameters used for the estimator:
 - `instance_count = 1`
 - `instance_type = 'ml.m5.large'`
 - `feature_dim = 16`
 - `mini_batch_size = 100`
 - `epochs = 10`
 - `predictor_type = 'binary_classifier'`

Model Training

- The model was trained using the training and validation datasets.
- Training completion time was approximately 4 minutes.

Job name	SageMaker metrics time series
linear-learner-2024-12-12-02-51-50-696	Enabled
ARN	Training time (seconds)
arn:aws:sagemaker:us-east-1:190647331505:training-job/linear-learner-2024-12-12-02-51-50-696	144
Status	Billable time (seconds)
✔ Completed	144
View history	Managed spot training savings
	0%
Creation time	Tuning job source/parent
Dec 12, 2024 02:51 UTC	-
Last modified time	IAM role ARN
Dec 12, 2024 02:55 UTC	arn:aws:iam::190647331505:role/saba-sagemaker ↗
Algorithm	
Algorithm ARN	Maximum wait time for managed spot training(s)
-	-
Training image	Managed spot training

Model Registration

- The trained model was registered for future deployments.

Amazon SageMaker AI > Models

Models

< 1 >

	Name	ARN	Creation time
<input type="radio"/>	linear-learner-2024-12-12-02-51-50-696	arn:aws:sagemaker:us-east-1:190647331505:model/linear-learner-2024-12-12-02-51-50-696	12/12/2024, 7:37:06 AM

Batch Predictions

- Batch inference was performed using SageMaker's batch transform feature.
- Parameters for batch transform:
 - `instance_type = "ml.m5.xlarge"`
- The batch transform job processed 1997 records with a runtime of approximately 6–8 minutes.

Amazon SageMaker AI > Batch transform jobs

Batch transform jobs

< 1 >

	Name	Status	Duration	Creation time
<input type="radio"/>	linear-learner-2024-12-12-15-39-17-630	InProgress	a few seconds	12/12/2024, 7:39:17 AM

Model Evaluation

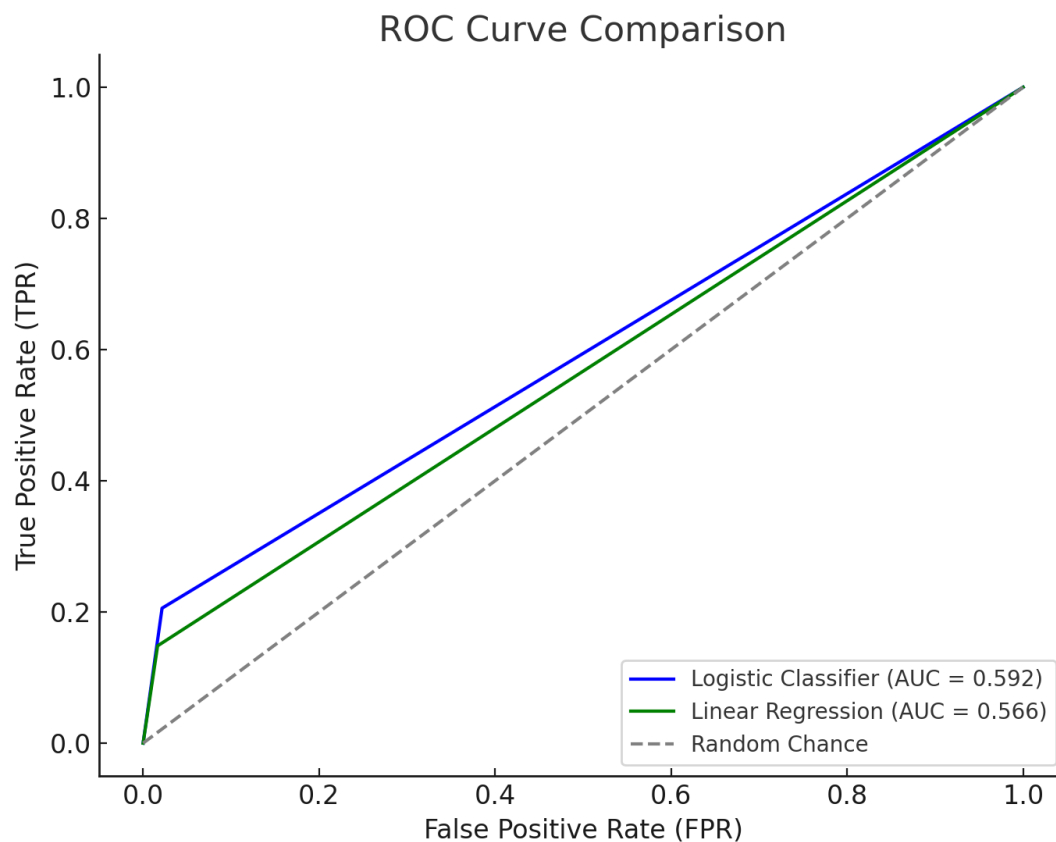
- The model's predictions were evaluated against the ground truth values stored in `test.csv`.
- Performance metrics such as accuracy, precision, and recall were used to assess the model's effectiveness.
- The same steps were repeated twice: first for the linear regression model and then for the logistic binary classifier.
- For the regression model, different thresholds were tested in the range of 0.3 to 0.7, with 0.5 identified as the optimal threshold for our use case.

By following these implementation steps, we ensured a structured and comprehensive approach to developing and evaluating our benchmark model.

Comparison of Logistic Binary Classifier and Linear Regression (Threshold = 0.5)

Performance metrics table

Metric	Logistic Binary Classifier	Linear Regression
Accuracy	0.817 (81.7%)	0.809 (80.9%)
Precision	0.717	0.705
Recall	0.206	0.148
F1-Score	0.320	0.245
Confusion Matrix	[[1545, 34], [332, 86]]	[[1553, 26], [356, 62]]
AUC	0.592	0.566



Observations

1. **Accuracy:**
 - The **logistic binary classifier** achieves slightly higher accuracy (**81.7%**) compared to the linear regression model (**80.9%**).
2. **Precision:**

- Both models have similar precision, with logistic binary classifier at **0.717** and linear regression at **0.705**, indicating that when predicting positive results, the models are fairly accurate.
 - 3. **Recall:**
 - Logistic binary classifier significantly outperforms linear regression in recall (**0.206** vs **0.148**), meaning it is better at identifying actual positives (fewer false negatives).
 - 4. **F1-Score:**
 - The F1-score, which balances precision and recall, is higher for the logistic binary classifier (**0.320**) compared to the linear regression (**0.245**).
 - 5. **AUC Logistic Classifier** outperforms Linear Regression for this problem, as evidenced by its higher AUC
 - 6. **Confusion Matrix:**
 - Both models perform similarly for true negatives (large counts in the top-left corner). However, the **logistic binary classifier** captures more true positives (**86**) compared to the linear regression (**62**). This aligns with its higher recall.
-

The **logistic binary classifier** outperforms the **linear regression** model overall for this binary classification problem. It achieves better recall and F1-score while maintaining slightly higher accuracy and precision. Therefore, it is the preferred model for this task. **We select this model to serve as a benchmark for our project**

Implementation

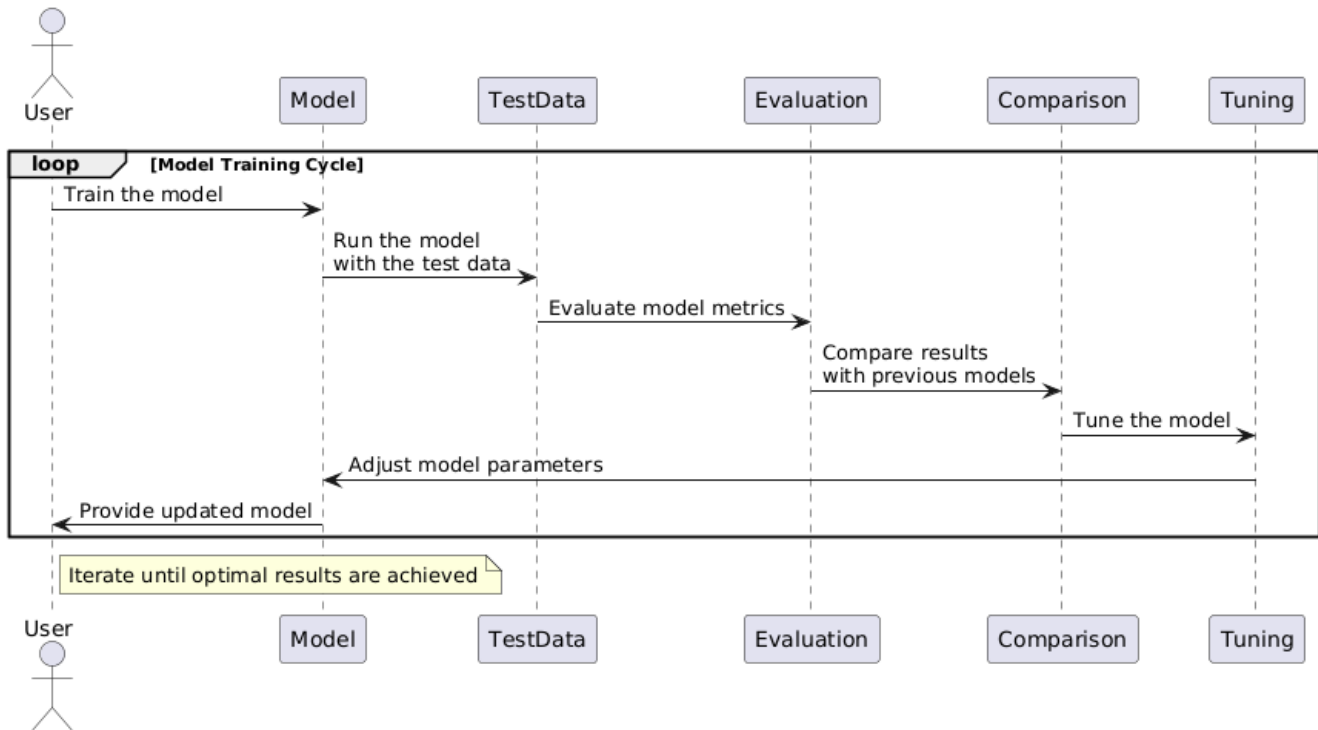
We used the same preprocessed files that were used in training and evaluation of benchmark model. The details of data preprocessing are described in [Data Preprocessing](#) section. The process was shared between all nine trained models

Implementation code is in Jupiter notebook [transaction recovery models.ipynb](#)

Sequence of steps to train and evaluate models

1. **Train the Model:** Begin by training the model with the train and validation data
2. **Run with Test Data:** Apply the trained model to test dataset through Batch Transform.
3. **Evaluate Performance:** Assess the model's performance using appropriate metrics.
4. **Compare Results:** Benchmark the results against those of previous models.
5. **Tune the Model:** Make necessary adjustments and fine-tune the model parameters.
6. **Iterate:** Continue this process until no further improvements are observed.
7. **Select a New Model:** Move on to experimenting with a different model.

We used exactly the same approach when implemented benchmark models. It ensures a thorough and iterative process for model development and optimization.



List of models trained

The following models were trained starting with benchmark model and finishing with the model with the best performance

Linear Learner framework in SageMaker

Model 1 – Linear Regression

Model 2 – Logistic Binary Classifier **selected as benchmark** and is described earlier

AutoGluon in local docker container

Model 3 – AutoGluon with default hyperparameters in local docker container

```

hyperparameters={
    'GBM': {}, # LightGBM
    'CAT': {}, # CatBoost
    'NN_TORCH': {}, # Neural network
    'XGB': {}, # XGBoost
    'RF': {}, # Random Forest
    'KNN': {}, # K-Nearest Neighbors
},

```

Model4 – AutoGluon with extended parameters and 3 models

```

hyperparameters = {
    'GBM': [
        {'num_leaves': 31, 'learning_rate': 0.01, 'feature_fraction': 0.8,
         'bagging_fraction': 0.8, 'bagging_freq': 5},
        {'num_leaves': 50, 'learning_rate': 0.1, 'feature_fraction': 0.9,
         'bagging_fraction': 0.9, 'bagging_freq': 7},
        {'num_leaves': 70, 'learning_rate': 0.05, 'feature_fraction': 0.85,
         'bagging_fraction': 0.85, 'bagging_freq': 10}
    ],
    'CAT': [

```

```

        {'iterations': 500, 'depth': 6, 'learning_rate': 0.01, 'l2_leaf_reg': 3},
        {'iterations': 1000, 'depth': 10, 'learning_rate': 0.1, 'l2_leaf_reg': 5},
        {'iterations': 1500, 'depth': 8, 'learning_rate': 0.05, 'l2_leaf_reg': 4}
    ],
    'XGB': [
        {'n_estimators':100, 'max_depth':3, 'eta': 0.01, 'subsample':0.8, 'colsample_bytree':0.8},
        {'n_estimators':500, 'max_depth':6, 'eta': 0.1, 'subsample':1.0, 'colsample_bytree':1.0},
        {'n_estimators':300, 'max_depth':4, 'eta':0.05, 'subsample':0.9, 'colsample_bytree':0.9}]
    ]
}

```

Model5 – AutoGluon with default parameter and eval_metric='roc_auc'

Model6 – AutoGluon with 3 models

```

hyperparameters= {
    'GBM': {},
    'XGB': {},
    'RF': {},
},

```

XGBoost models in local container

Model7 XGBoost with default hyperparameters

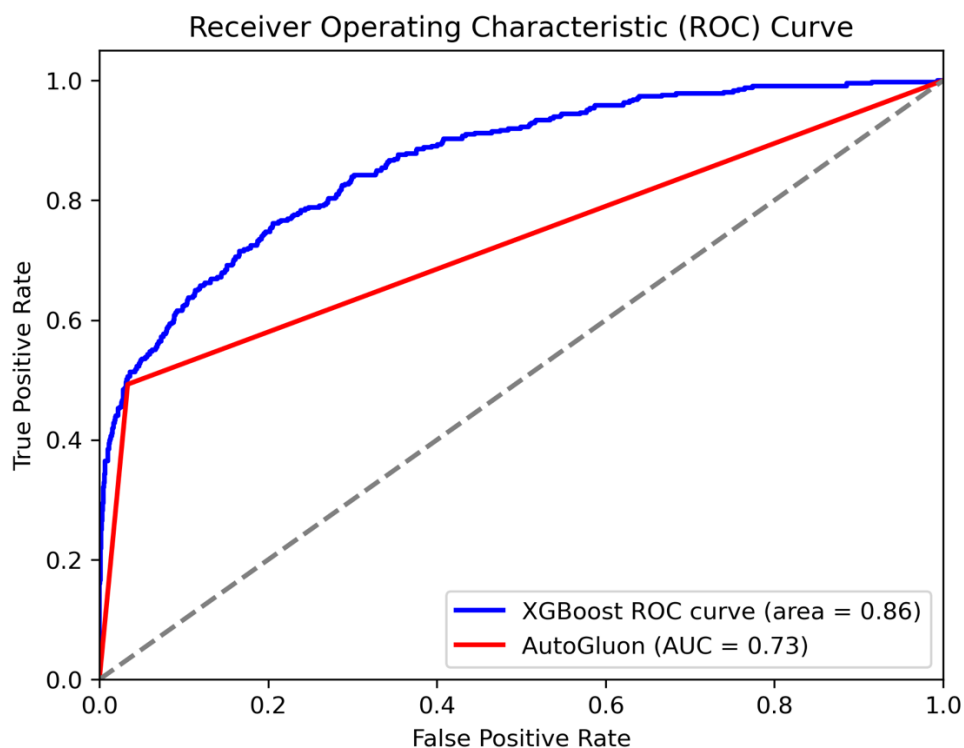
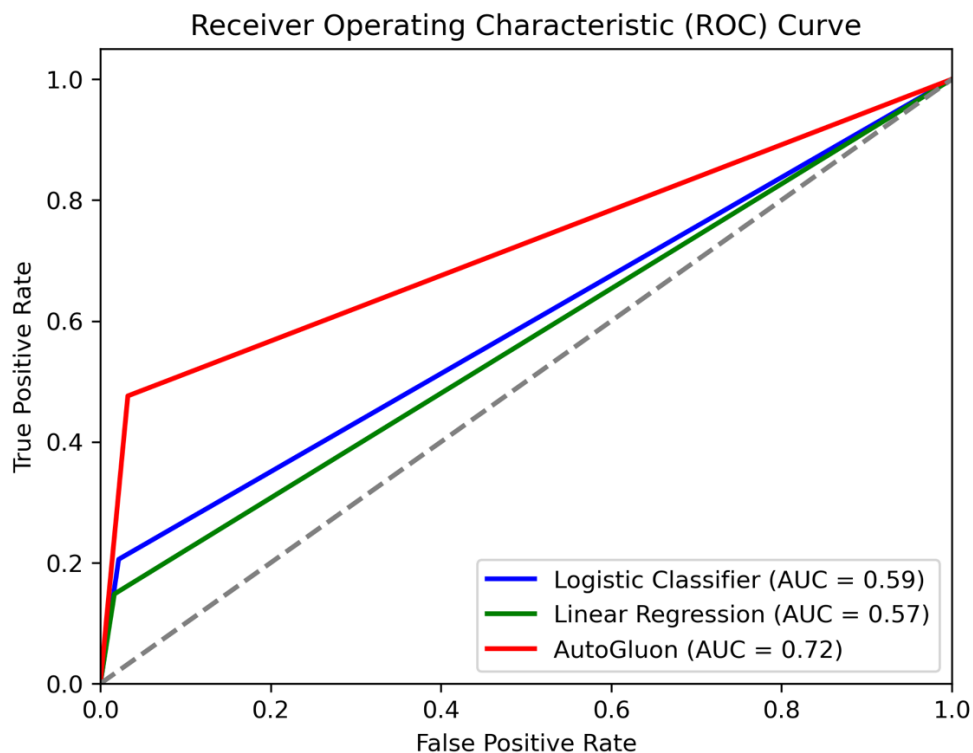
Model8 XGBoost with hyperparameters grid search optimization

Model Performance Comparison

MODEL	1	2	3	4	5	6	7	8
ACCURACY	81.67%	80.87%	86.48%	86.48%	85.88%	86.73%	86.68%	87.18%
Precision	0.717	0.705	0.796	0.796	0.746	0.795	0.86	0.87
Recall	0.206	0.148	0.476	0.476	0.448	0.493	0.87	0.87
F1-Score	0.320	0.245	0.596	0.596	0.559	0.609	0.85	0.86
MCC	0.177	0.229	0.546	0.546	0.504	0.556	-	-
Balanced Accuracy	58.20%	61.90%	72.19%	72.19%	70.47%	72.96%	70.47%	-
AUC-ROC	0.592	0.566	0.720	0.720	0.857	-	0.857	0.864

Summary of the performance metrics results

- **Model 1 (Logistic Classifier):** Decent precision but lower recall and F1-score.
- **Model 2 (Linear Regression):** Similar to Model 1, with slightly lower performance metrics.
- **Model 3 & Model 4:** Strong performance with high accuracy and precision. AUC-ROC at 0.720.
- **Model 5:** Slightly lower performance in accuracy but high AUC-ROC.
- **Model 6:** Best overall performance with the highest accuracy and balanced metrics.
- **Model 7:** Comparable to Model 6 with strong AUC-ROC and weighted metrics.
- **Model 8** shows the highest accuracy among all models at **87.18%**. It maintains very strong precision, recall, and F1-score metrics, indicating balanced performance across both classes. The **AUC-ROC of 0.864** is also excellent, slightly higher than that of Model 7, indicating that **Model 8 has strong discriminatory power.**



ROC curves Interpretation

- The XGBoost model (model 8) outperforms the AutoGluon model (model4), and AutoGluon models (model 4) outperform benchmark models (model 2). This is indicated by higher AUC value and the ROC curve being closer to the top-left corner.

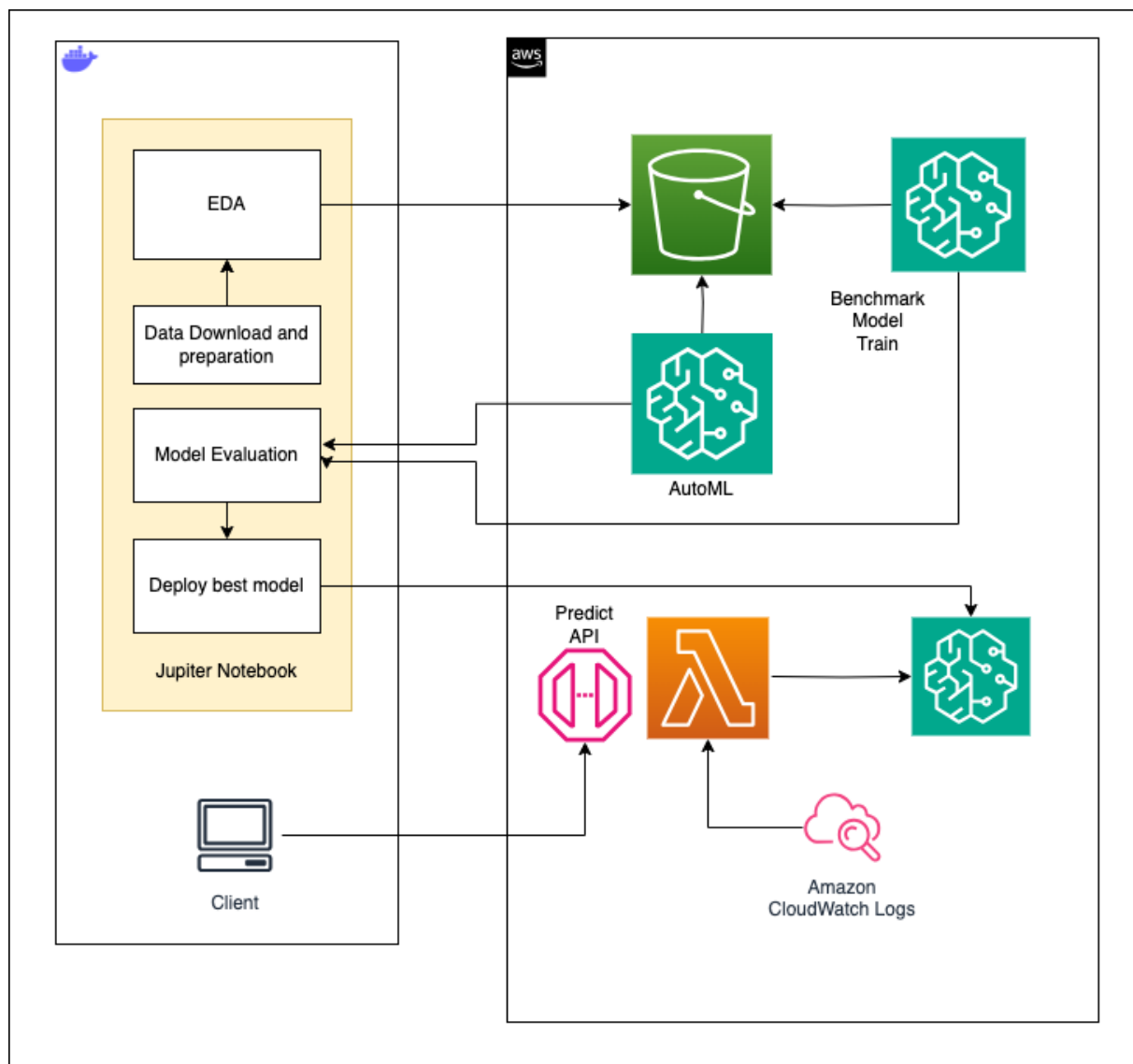
Discussion

Model 8 stands out as the top performer based on accuracy, precision, recall, F1-score, and AUC-ROC, suggesting it is the most reliable model for our task. This model not only achieves the highest accuracy but also maintains a good balance between precision and recall, making it a highly effective choice. This model was trained in the local docker container which was an equivalent of Amazon instance ml.m5.large. Next step was to do wider grid search optimization in Amazon SageMaker, where we could use more grid points and did training jobs in parallel.

The ultimate goal was to deploy the best model to SageMaker endpoint and build an API to predict transactions success rate in real-time.

Transaction recovery solution architecture

We implemented the following architecture to enable customer integration with developed model



Once the best model was trained (model 8 XGBoost with tuned hyperparameters), we reproduced it as Model 9 in SageMaker with the same parameters. We further tuned hyperparameters in SageMaker using grid search.

```
hyperparameter_ranges = {
    'colsample_bytree': ContinuousParameter (0.6, 1.0),
    'eta': ContinuousParameter (0.01, 0.2),
    'max_depth': IntegerParameter(3, 7),
    'subsample': ContinuousParameter (0.6, 1.0)
}
```

SageMaker used different container image than local container, so it explains small discrepancy in performance metrics. These differences did not produce any significant impact considering small amount of data used in the project.

Final values of hyperparameters were found to be:

colsample_bytree	0.800401
eta	0.112352
max_depth	4.0
num_round	429.0
subsample	0.787507
TrainingJobName	XGBoost-tuning-182-5d808c29
TrainingJobStatus	Completed
FinalObjectiveValue	0.85524
TrainingStartTime	2024-12-20 02:08:08+00:00
TrainingEndTime	2024-12-20 02:08:37+00:00
TrainingElapsedTimeSeconds	29.0
TrainingJobDefinitionName	XGBoost

Model 9 (**xgboost-tuning-182-5d808c29**) was deployed and endpoint created:

Amazon SageMaker AI > Endpoints

Endpoints

Update endpoint

Actions ▾

Create endpoint

Q Search endpoints

< 1 >


	Name ▾	ARN ▾	Creation time ▾	Status ▾	Last updated ▾
<div></div>	transaction-predict-1	arn:aws:sagemaker:us-east-1:190647331505:endpoint/transaction-predict-1	12/21/2024, 7:14:37 AM	<div>InService</div>	12/21/2024, 7:18:05 AM


This endpoint was invoked by lambda function (code is [here](#)).
Finally, this lambda function was used as a resource for ***predict-transaction*** API


predict-transaction Throttle Copy ARN Actions

▼ **Function overview** [Info](#)

Diagram Template

 **predict-transaction**


 Layers (0)

 **API Gateway**

+ Add trigger + Add destination

Description
-

Last modified
58 minutes ago

Function ARN
 arn:aws:lambda:us-east-1:190647331505:function:predict-transaction

Function URL [Info](#)
-

Export to Infrastructure Composer Download

[API Gateway](#) > [APIs](#) > Resources - transaction-predict (isx8ecrvxb)


APIs
Custom domain names [Updated](#)
Domain name access associations [New](#)
VPC links


▼ **API: transaction-predict**

[Resources](#)
Stages
Authorizers
Gateway responses
Models

Resources API actions Deploy API

Create resource

 /

 /transaction-predict
POST

Resource details Update documentation Enable CORS

Path
/

Resource ID
yby6each84

Methods (0) Delete Create method

We validated end-to-end process by running API calls through `curl` command:

```
Svetlanas-MacBook-Air:figures sab$ curl -X POST https://isx8ecrvxb.execute-api.us-east-1.amazonaws.com/demo/transaction-predict -H "Content-Type: application/json" -d '{"input_data": "26,1,2,0,0,1.0,24.000000000000004,2,1.0,0.0,0.0,0.0,1.0,0.0,0.0,1.0"}' {"statusCode": 200, "body": "{\"prediction\": \"0.015935370698571205\\n\\n\"}"}Svetlanas-MacBook-Air:figures sab$
```

We ran a few functional tests to validate that the process is actually returning correct values by running transactions from the test data file and comparing with expected results that proved 100% reliability of the implementation

Conclusions

The outcomes of this project are as follows:

1. **Dataset Preparation:** We prepared a dataset that accurately emulates real-world billing transactions.
2. **Benchmark Model Creation:** We established a benchmark model using logistic regression for binary classification.

3. **Model Training:** We trained a sequence of models and identified XGBoost as the best-performing model, achieving an exceptional accuracy of 87%
4. **Model Deployment:** The top-performing model was deployed in the SageMaker environment.
5. **API Development:** An API was created to facilitate inference calls to the deployed model.

The current architecture allows customers to utilize the model to predict billing success rates, thereby avoiding the costs associated with retrying transactions that have a low probability of success. The resulting API is ready for integration into the customer's code base.

Acknowledgement

I thank all the teachers who helped me by providing interesting assignments and learning materials. Kudos to chatGPT (and all developers whose code is used by the tool), I can imagine how much harder this project would be before.

Resources

- [1] [Understanding Involuntary Churn and How to Address It](#)
- [2] [Classifying variety of customer's online engagement for churn prediction with mixed-penalty logistic regression](#)
- [3] [Bank Customer Churn Dataset](#)
- [4] Hosmer, D. W., & Lemeshow, S. (2000). *Applied Logistic Regression*. Wiley.
- [5] Erickson, N., Mueller, J., McWilliams, B., et al. (2020). AutoGluon-Tabular: Robust and Accurate AutoML for Structured Data. *arXiv preprint arXiv:2003.06505*.
- [6] Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 785-794)
- [7] Bergstra, J., & Bengio, Y. (2012). Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13(Feb), 281-305.