

Fraud Analytics (CS6890)

Assignment: 3

Title : Example-dependent cost-sensitive regression

	Name	Roll Number
	<i>Shreesh Gupta</i>	<i>CS23MTECH12009</i>
	<i>Hrishikesh Hemke</i>	<i>CS23MTECH14003</i>
Team Details :		
	<i>Manan Patel</i>	<i>CS23MTECH14006</i>
	<i>Yug Patel</i>	<i>CS23MTECH14019</i>
	<i>Bhargav Patel</i>	<i>CS23MTECH11026</i>

```
In [ ]: 1 # Required Libraries
          2 !pip install pandas numpy matplotlib scikit-learn torch torch-geometric ne
```

Cost-Sensitive Logistic Regression

Overview

This notebook implements various approaches to logistic regression, specifically focusing on cost-sensitive logistic regression. Cost-sensitive learning is a critical area in machine learning where the costs associated with misclassification vary between classes. This notebook provides a comprehensive look into implementing cost-sensitive logistic regression models using both gradient descent and genetic algorithms.

Libraries and Data Loading

We start by importing necessary libraries and loading the dataset:

In [1]:

```

1 import numpy as np
2 from scipy.special import expit as sigmoid # Sigmoid function
3 from sklearn.base import BaseEstimator, ClassifierMixin
4 from scipy.optimize import minimize
5 from sklearn.metrics import accuracy_score, log_loss, confusion_matrix, +
6 import pandas as pd
7 from sklearn.linear_model import LogisticRegression
8 import random
9 import warnings
10 from sklearn.neighbors import kneighbors_graph
11 import torch
12 from torch_geometric.data import Data
13 warnings.filterwarnings('ignore', category=RuntimeWarning)
14 import numpy as np
15 from sklearn.model_selection import train_test_split
16 import plotly.express as px
17 import pandas as pd
18 from sklearn.model_selection import train_test_split
19 import numpy as np
20 from scipy.optimize import minimize
21 from scipy.special import expit as sigmoid # Stable sigmoid function
22 from scipy.special import gamma, gammaln, gammalncinv
23 import numpy as np
24 from scipy.special import expit as sigmoid
25 import scipy
26 import plotly.graph_objects as go
27 from sklearn.metrics import f1_score
28 import plotly.figure_factory as ff
29 import numpy as np

```

Load the data from a CSV file:

In [2]:

```

1 data_path = 'costsensitiveregression.csv'
2 data = pd.read_csv(data_path)
3
4 # Display the first few rows and columns to confirm its structure
5 data.head()

```

Out[2]:

	NotCount	YesCount	ATPM	PFD	PFG	SFD	SFG	WP	WS	AH	AN	Status	FNC
0	2	21	0.0	0.000	0.0	0.0	0.0	0.000000	0.0	0.0	0.0	0	0.0
1	23	0	0.0	0.044	0.0	0.0	0.0	0.306179	0.0	0.0	0.0	1	0.0
2	1	22	0.0	0.000	0.0	0.0	0.0	0.000000	0.0	0.0	0.0	0	0.0
3	5	18	0.0	0.000	0.0	0.0	0.0	0.000000	0.0	0.0	0.0	1	0.0
4	1	22	0.0	0.000	0.0	0.0	0.0	0.000000	0.0	0.0	0.0	0	0.0

Visualization

Visualize the distribution of labels using Plotly:

```
In [3]: 1 # Let's assume the DataFrame 'data' has a 'status' column with binary labels
2 # Let's replace dummy data with your specific column name if different
3 # data['Status'] = data['Status'].replace({1: 'label-1', 0: 'label-0'}) #
4
5 # Calculate the frequency of each label
6 label_counts = data['Status'].value_counts()
7
8 # Create a pie chart
9 fig = px.pie(
10     label_counts,
11     values=label_counts.values,
12     names=label_counts.index,
13     title='Distribution of Labels in Status Column',
14     color_discrete_sequence=px.colors.sequential.RdBu
15 )
16
17 # Show the plot
18 fig.show()
19
```

Data Preprocessing

We preprocess the data by encoding labels and splitting the dataset into training and testing sets:

```
In [4]: 1 X = data.iloc[:, :-3] # Assuming the last three columns are 'Status', 'FNC'
2 y = data['Status']
3 costs = data['FNC']
```

Model Implementation

Cost Matrix Definition

Define the cost matrix for cost-sensitive learning:

```
In [5]: 1 FP = np.full(len(y), 6) # Replace len(y) with the actual length of your labels
2 TP = np.full(len(y), 6) # Same here
3 TN = np.zeros(len(y)) # And here
4 cost_matrix = np.column_stack((TP, data['FNC'], FP, TN)).astype(float)
5
6 # Split the data into train and test sets
7 X_train, X_test, y_train, y_test, cost_matrix_train, cost_matrix_test, costs
8     X, y, cost_matrix, costs, test_size=0.20, random_state=42
9 )
10
11
```

Cost-Sensitive Logistic Regression Using Gradient Descent (Bahnsen approach)

Introduction

This implementation details a cost-sensitive logistic regression model that optimizes its parameters using gradient descent. Unlike traditional logistic regression, this model considers the varying costs of different types of classification errors, which can be crucial in domains like healthcare, finance, or any field where the consequences of errors are asymmetric.

Model Description

Sigmoid Function

The logistic regression model uses the sigmoid function to map the net input (z) (a linear combination of input features (X) and weights (w)) to probabilities between 0 and 1. The sigmoid function is defined as:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Cost-Sensitive Loss Function

The cost-sensitive loss function incorporates costs associated with four types of outcomes: true positives, true negatives, false positives, and false negatives. It is defined as follows:

$$L(y, \hat{y}, C) = \sum_{i=1}^N [y_i \times ((1 - \hat{y}_i) \times C_{FN} + \hat{y}_i \times C_{TP}) + (1 - y_i) \times (\hat{y}_i \times C_{FP} + (1 - \hat{y}_i) \times C_{TN})]$$

Here:

- (y_i) is the actual label of the i th sample.
- (\hat{y}_i) is the predicted probability for the i th sample.
- (C_{TP} , C_{FN} , C_{FP} , C_{TN}) are the costs associated with true positives, false negatives, false positives, and true negatives, respectively.

Gradient Descent Optimization

The model parameters (weights (w) and bias (b)) are updated using gradient descent to minimize the cost-sensitive loss. The gradients for the weights and bias are computed as:

$$\frac{\partial L}{\partial w} = -\frac{1}{N} \sum_{i=1}^N X^T \times ((y - \hat{y}) \times (C_{TP} - C_{FP}))$$

$$\frac{\partial L}{\partial b} = -\frac{1}{N} \sum_{i=1}^N ((y - \hat{y}) \times (C_{TP} - C_{FP}))$$

The weights and bias are then updated iteratively:

$$w = w - \alpha \times \frac{\partial L}{\partial w}$$

$$b = b - \alpha \times \frac{\partial L}{\partial b}$$

where (α) is the learning rate.

Implementation

```
class CostSensitiveLogisticRegressionL:
    def __init__(self, alpha=0.01, epochs=100):
        self.alpha = alpha
        self.epochs = epochs
        self.weights = None
        self.bias = None

    def fit(self, X, y, cost_matrix):
        self.weights, self.bias = gradient_descent(X, y, cost_matrix,
                                                self.alpha, self.epochs)

    def predict_proba(self, X):
        probabilities = sigmoid(np.dot(X, self.weights) + self.bias)
        return probabilities

    def predict(self, X):
        probabilities = self.predict_proba(X)
        return (probabilities >= 0.5).astype(int)
```

In [6]:

```

1 # For Bahnsen et al.-inspired model
2 # Define the sigmoid function
3 def sigmoid(z):
4     return 1 / (1 + np.exp(-z))
5
6 # Define the cost-sensitive loss function
7 def calculate_loss(X, y, weights, bias, cost_matrix):
8     # Predict probabilities
9     probabilities = sigmoid(np.dot(X, weights) + bias)
10
11    # Calculate cost-sensitive loss
12    loss = y * ((1 - probabilities) * cost_matrix[:, 1] + probabilities *
13                + (1 - y) * (probabilities * cost_matrix[:, 0] + (1 - probabilities) * cost_matrix[:, 2]))
14
15    return np.sum(loss)
16
17 # Gradient descent function to minimize the cost-sensitive loss
18 def gradient_descent(X, y, cost_matrix, alpha, epochs):
19     n_samples, n_features = X.shape
20     weights = np.zeros(n_features)
21     bias = 0
22
23     for _ in range(epochs):
24         probabilities = sigmoid(np.dot(X, weights) + bias)
25
26         # Gradient calculation
27         dw = -np.dot(X.T, (y - probabilities) * (cost_matrix[:, 2] - cost_matrix[:, 0] * probabilities - cost_matrix[:, 1] * (1 - probabilities)))
28         db = -np.sum((y - probabilities) * (cost_matrix[:, 2] - cost_matrix[:, 0] * probabilities - cost_matrix[:, 1] * (1 - probabilities)))
29
30         # Update weights and bias
31         weights -= alpha * dw
32         bias -= alpha * db
33
34     return weights, bias
35
36 # Define the cost-sensitive logistic regression model
37 class CostSensitiveLogisticRegression:
38     def __init__(self, alpha=0.01, epochs=100):
39         self.alpha = alpha
40         self.epochs = epochs
41         self.weights = None
42         self.bias = None
43
44     def fit(self, X, y, cost_matrix):
45         self.weights, self.bias = gradient_descent(X, y, cost_matrix, self.alpha, self.epochs)
46
47     def predict_proba(self, X):
48         return sigmoid(np.dot(X, self.weights) + self.bias)
49
50     def predict(self, X):
51         probabilities = self.predict_proba(X)
52         return (probabilities >= 0.5).astype(int)
53

```

Model Training and Evaluation Summary

- **Algorithm:** Cost-sensitive logistic regression optimized via genetic algorithm.
- **Generations:** Trained over 50 generations.
- **Evaluation:**
 - **Accuracy:** 29.9% on the test set, indicating low performance.
 - **Log Loss:** 6.009, suggesting poor prediction calibration.

In [7]:

```

1 # Initialize and fit the model
2 cost_sensitive_lr_bun = CostSensitiveLogisticRegressionL()
3 cost_sensitive_lr_bun.fit(X_train, y_train, cost_matrix_train)
4
5 # Predictions and evaluations
6 predictions_cslr_bun = cost_sensitive_lr_bun.predict(X_test)
7 accuracy_cslr_bun = accuracy_score(y_test, predictions_cslr_bun)
8 loss_cslr_bun = log_loss(y_test, cost_sensitive_lr_bun.predict_proba(X_te
9
10 print("Accuracy:", accuracy_cslr_bun)
11 print("Log Loss:", loss_cslr_bun)
12

```

Accuracy: 0.2987334055811433

Log Loss: 0.6931471805599452

Nikou Gunnemann's approach

Customized Logistic Loss Function Explanation

The `logistic_loss` function in Python computes a variant of the logistic regression loss, incorporating costs that vary based on different model variants. This adaptation allows the model to differently weigh the importance of misclassifications according to their associated costs, crucial in scenarios where different types of errors have significantly different implications.

General Setup for Logistic Regression

Logistic regression predicts the probability (p) that a given instance (x) belongs to the positive class (i.e., ($y = 1$)):

$$p = \sigma(\beta^T x)$$

where (σ) is the sigmoid activation function defined by:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

The conventional logistic loss function, used for training logistic regression models, penalizes predictions based on the log likelihood of the true classes:

$$L(y, p) = -[y \log(p) + (1 - y) \log(1 - p)]$$

This function is summed over all instances, where being wrong and confident leads to higher penalties.

Variant-Specific Cost-Sensitive Loss Functions

Variant A: Linear Weighting

- **Weights ((a_i)):** Directly proportional to the costs associated with each instance, ($a_i = c_i$).
- **Exponents ((b_i)):** Set to 1, maintaining the standard logistic loss form but scaled by cost.

$$L_A = \sum c_i [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

Variant B: Exponential Weighting

- **Weights ((a_i)):** All set to 1.
- **Exponents ((b_i)):** Determined by the inverse gamma function of the costs, modifying the sensitivity of the loss to errors, ($b_i = \text{Gamma}^{-1}(c_i) - 1$).

$$L_B = \sum [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]^{b_i}$$

Variant C: Ratio Control

- **Weights ((a_i)) and Exponents ((b_i)):** Both parameters are tailored to manage the ratio of costs between different types of misclassifications. Calculations involve solving equations to balance the ratios appropriately:
 - (b_i) is calculated such that it satisfies a specific cost-benefit ratio.
 - ($a_i = 1 / \text{Gamma}(b_i + 1)$).

$$L_C = \sum \frac{1}{\Gamma(b_i + 1)} [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]^{b_i}$$

Variant D: Maintaining Constant Correct Classification Cost

- **Weights ((a_i)) and Exponents ((b_i)):** Designed to keep the cost of correct classifications constant while varying the penalties for misclassifications according to their costs.
 - (a_i) and (b_i) are adjusted to ensure that the losses for correct classifications remain constant across different costs.
 - ($a_i = \frac{0.5}{\text{Gamma}(b_i + 1) - \text{Gamma}(b_i + 1, 0.6931)}$).

$$L_D = \sum \frac{0.5}{\Gamma(b_i + 1) - \Gamma(b_i + 1, 0.6931)} [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]^{b_i}$$

These variants modify the standard logistic loss to integrate the cost directly into the model training, making the logistic regression sensitive to the cost associated with incorrect

```
In [8]: 1 def logistic_loss(beta, X, y, costs, variant):
2     predictions = sigmoid(np.dot(X, beta))
3     m = len(y) # Number of instances
4
5     if variant == 'A':
6         # Variant A: Linear weights as per cost
7         ai = costs
8         bi = np.ones_like(costs)
9     elif variant == 'B':
10        # Variant B: Exponential weights determined by the inverse gamma function
11        ai = np.ones_like(costs)
12        bi = gammaln(inv(costs, 0.5) - 1 # This needs validation for correctness
13    elif variant == 'C':
14        # Variant C: Control the ratio with adjusted ai and bi
15        bi = np.array([gammaln(1, (1 + 0.5 * ci) / (ci + 1)) for ci in costs])
16        ai = 1 / gamma(bi + 1)
17    elif variant == 'D':
18        # Variant D: Similar to C but keeps the correct classification loss
19        bi = np.array([gammaln(1, (1 + 0.5 * ci) / (ci + 1)) for ci in costs])
20        ai = 0.5 / (gamma(bi + 1) - gammaln(bi + 1, 0.6931))
21
22    # Calculate the customized Loss
23    loss = np.sum(ai * (y * (-np.log(predictions + 1e-15))**bi + (1 - y) * (1 - predictions)**bi))
24
25    return loss
```

Fitting Cost-Sensitive Logistic Regression Models

The `fit_model` function is designed to fit logistic regression models by optimizing a custom logistic loss function that takes into account different cost structures depending on the specified variant. The function employs the `minimize` method from `scipy.optimize` to find the best model coefficients (`beta`) that minimize the loss.

Function Definition

The `fit_model` function takes the following parameters:

- **X (array-like)**: Feature matrix for the training data.
- **y (array-like)**: Target vector for the training data.
- **costs (array-like)**: Vector of costs associated with each instance in the training data.
- **variant (str)**: Specifies the variant of the cost-sensitive logistic regression to be used.

It initializes the model coefficients to zeros, then uses the BFGS algorithm (a quasi-Newton method) to optimize the coefficients by minimizing the custom logistic loss:

In [9]:

```

1 def fit_model(X, y, costs, variant):
2     beta_init = np.zeros(X.shape[1])
3     result = minimize(logistic_loss, beta_init, args=(X, y, costs, variant))
4     return result.x
5
6
7
8 # Fit models for each variant
9 variants = ['A', 'B', 'C', 'D']
10 models_nik = {}
11 for variant in variants:
12     models_nik[variant] = fit_model(X_train, y_train, costs_train, variant)
13     print(f"Model coefficients for variant {variant}: {models_nik[variant]}")

```

```

Model coefficients for variant A: [ 0.1190579 -0.15401323  0.05821664  0.219
04436 -0.42776506 -0.01477515
0.00892466 -0.05023754 -0.11274034 -0.06083723]
Model coefficients for variant B: [-0.00328367  0.00659971 -0.28634056  0.021
09848  0.02959005 -0.02277006
-0.03964703 -0.00038196 -0.02422468  0.01262785]
Model coefficients for variant C: [ 1.67472449e-01 -1.46788075e-01 -1.4522074
5e+00  1.76324331e-02
1.49062543e-01 4.66728940e-02 -4.56879000e+01 -1.76797565e-02
-9.67787998e-02 3.04078471e-01]
Model coefficients for variant D: [ 1.69746606e-01 -1.53645542e-01 -1.3599763
0e+00  2.23429085e-02
1.46982817e-01 3.17544436e-02 -4.29806934e+01 -1.77911989e-02
-8.88658171e-02 2.86332200e-01]

```

Logistic Regression Overview and Performance Evaluation

Logistic Regression Explanation

Logistic Regression is a statistical method used for binary classification that predicts the probability of the target variable's categories. It models the probability of the default class (usually "1") using the logistic function.

Equation of the Logistic Regression Model

The logistic regression model calculates probabilities using the logistic function:

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1x_1 + \beta_2x_2 + \dots + \beta_nx_n$$

Where:

- (p) is the probability of belonging to the positive class.
- ($\beta_0, \beta_1, \dots, \beta_n$) are the coefficients of the model.
- (x_1, x_2, \dots, x_n) are the predictor variables.

Probability Prediction Using the Sigmoid Function

The logistic function or the sigmoid function is used to convert the linear regression output to a probability:

$$[p = \sigma(z) = \frac{1}{1 + e^{-z}}]$$

where (z) is the linear combination of the features and coefficients, ($z = \beta_0 + \beta_1x_1 + \dots + \beta_nx_n$).

Model Training and Performance

Training Process

A standard logistic regression model was trained using the `LogisticRegression` class from `sklearn`, configured with a maximum of 500 iterations to ensure convergence.

```
standard_lr = LogisticRegression(max_iter=500)
standard_lr.fit(X_train, y_train)
```

Results

- **Accuracy:** The accuracy achieved was 86.58%, indicating the percentage of correctly predicted instances in the test data.
- **Log Loss:** The log loss recorded was 0.3056, reflecting the model's confidence in its predictions. Lower log loss values are better as they indicate higher confidence and accuracy in the predictions.

```
In [10]: 1 standard_lr = LogisticRegression(max_iter=500)
2 standard_lr.fit(X_train, y_train)
3 predictions_lr = standard_lr.predict(X_test)
4
5 accuracy_lr = accuracy_score(y_test, predictions_lr)
6 loss_lr = log_loss(y_test, standard_lr.predict_proba(X_test))
7
8 print("Accuracy:", accuracy_lr)
9 print("Log Loss:", loss_lr)
```

Accuracy: 0.8661609319967488

Log Loss: 0.30547488662730277

Mathematical Explanation and Results Discussion for Cost-Sensitive Logistic Regression

Mathematical Foundation

- **Objective:** Minimize both predictive error and the financial/practical impact of errors, quantified by a cost matrix.
 - **Cost Matrix:** Assigns costs to different types of classification errors:
 - **True Positives (TP)**
 - **False Positives (FP)**
 - **True Negatives (TN)**
 - **False Negatives (FN)**
 - **Cost-Sensitive Loss Function:**

The loss function for cost-sensitive logistic regression, accounting for different types of errors, is defined as:

$$L(y, \hat{y}, C) = \sum_{i=1}^N \left[y_i \times ((1 - \hat{y}_i) \times C_{FN} + \hat{y}_i \times C_{TP}) + (1 - y_i) \times (\hat{y}_i \times C_{FP} + (1 - \hat{y}_i) \times C_{TN}) \right]$$

This function:

```
In [11]: 1 def cost_sensitive_logistic_loss(y_true, y_pred, cost_matrix):
2     total_loss = y_true * ((1 - y_pred) * cost_matrix[:, 1] + y_pred * cos_
3     return np.sum(total_loss)
```

```
In [12]: 1 def evaluate_loss(models, X, y, costs, variant):
2     beta = models[variant]
3     predictions = sigmoid(X.dot(beta))
4     # Calculate the Logistic regression loss
5     loss = logistic_loss(beta, X, y, costs, variant)
6     return loss
7
8 def gamma_func(x):
9     return scipy.special.gamma(x)
```

```
In [13]: 1 cost_cslr_bun = cost_sensitive_logistic_loss(y_test, predictions_cslr_bun)
2 print("Bahnsen approach Cost of cost-sensitive logistic regression:", cost_cslr_bun)
3
4 losses = {}
5 for variant in variants:
6     losses[variant] = evaluate_loss(models_nik, X_test, y_test, costs_test)
7     print(f"Nikou Gunnemann's approach Cost-sensitive Logistic Loss for variant {variant}: {losses[variant]}")
8
9 cost_lr = cost_sensitive_logistic_loss(y_test, predictions_lr, cost_matrix)
10 print("Cost of Standard logistic regression:", cost_lr)
```

Bahnsen approach Cost of cost-sensitive logistic regression: 177168.0
 Nikou Gunnemann's approach Cost-sensitive Logistic Loss for variant A: 486973
 5.186412062
 Nikou Gunnemann's approach Cost-sensitive Logistic Loss for variant B: 8226.0
 9071459818
 Nikou Gunnemann's approach Cost-sensitive Logistic Loss for variant C: 7998.8
 377675579395
 Nikou Gunnemann's approach Cost-sensitive Logistic Loss for variant D: 5092.2
 63001465815
 Cost of Standard logistic regression: 814303.592301

Cost Comparison between Standard and Cost-Sensitive Logistic Regression Models

The bar chart illustrates a stark contrast in the total misclassification costs incurred by two logistic regression models: the standard LR and the cost-sensitive LR.

- **Standard Logistic Regression (Blue Bar):**
 - The cost incurred by the standard logistic regression model is substantially higher, with a total of 812,473.36. This suggests that while the model might predict accurately, it does not account for the varying costs of different types of classification errors, leading to a less economically optimal outcome.
- **Cost-Sensitive Logistic Regression (Green Bar):**
 - In contrast, the cost-sensitive logistic regression model shows a significantly reduced cost of 177,162.0. This model incorporates a cost matrix that assigns specific costs to different types of errors, such as false positives and false negatives. By doing so, it is able to minimize the more economically impactful errors, thus reducing the overall cost.

In [14]:

```

1 # Data for plotting
2 costs = [cost_lr, cost_cslr_bun] + [losses[v] for v in variants]
3 labels = ['Standard LR', 'Bahnsen approach Cost-Sensitive LR'] + [f"Nikou"
4
5 # Create a bar chart
6 fig = go.Figure(go.Bar(
7     x=labels,
8     y=costs,
9     text=costs,
10    textposition='auto',
11    marker_color=['red', 'green'] + ['blue', 'purple', 'yellow', 'orange']
12 ))
13
14 # Add titles and labels
15 fig.update_layout(
16     title='Comparison of Logistic Regression Costs Across Different Approa
17     xaxis_title='Model Type',
18     yaxis_title='Cost',
19     template='plotly_white'
20 )
21
22 # Improve resolution by setting the width and height
23 fig.update_layout(width=1000, height=600) # Adjusted for better visualiza
24
25 # Show the plot
26 fig.show()
27

```

F1 Score Comparison Between Standard and Cost-Sensitive Logistic Regression

F1 Score is the harmonic mean of precision and recall, offering a balance between the two by taking both false positives and false negatives into account. It is particularly useful when the class distribution is imbalanced.

The F1 Score can be mathematically represented as:

$$F1 = 2 \cdot \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

where:

- **Precision** is the ratio of true positive predictions to the total predicted positives:

$$\text{Precision} = \frac{TP}{TP + FP}$$

- **Recall** (also known as sensitivity) is the ratio of true positive predictions to the actual positives:

$$\text{Recall} = \frac{TP}{TP + FN}$$

Results

- **Standard Logistic Regression:**
 - **F1 Score:** 0.763, indicating a relatively high balance between precision and recall. This model may predict positive classes more reliably than the cost-sensitive counterpart.
- **Cost-Sensitive Logistic Regression:**
 - **F1 Score:** 0.460, significantly lower than the standard model, which implies a trade-off between precision and recall that favors minimizing costlier errors over achieving a balance.

Discussion

- The higher F1 Score of the standard LR model indicates better performance in terms of balanced precision and recall. However, it doesn't consider the different costs of misclassification.
- The cost-sensitive LR model's lower F1 Score suggests that while it may be optimizing for cost, it does so at the expense of either precision, recall, or both.
- It's important to note that a lower F1 Score in the cost-sensitive model does not necessarily indicate poorer overall performance, as it may be strategically accepting lower precision and/or recall to reduce more costly errors, aligning with its optimization objectives.

The choice between these models should be guided by the specific context and objectives:

```
In [15]: 1 # Assuming predictions for each variant are available as predictions_nik[variant]
2 predictions_nik = {} # Dictionary to store predictions
3 variants = ['A', 'B', 'C', 'D'] # Define your variants
4 for variant in variants:
5     # Generate predictions for variant; replace this with your actual predictions
6     predictions_nik[variant] = np.round(sigmoid(np.dot(X_test, models_nik[variant])))
7
8 # Calculate F1 scores
9 f1_lr = f1_score(y_test, predictions_lr)
10 f1_cslr_bun = f1_score(y_test, predictions_cslr_bun)
11 f1_scores_nik = {variant: f1_score(y_test, predictions_nik[variant]) for variant in variants}
12
13 # Print F1 scores
14 print("Standard LR F1 Score:", f1_lr)
15 print("Bahnsen approach Cost-Sensitive LR F1 Score:", f1_cslr_bun)
16 for variant in variants:
17     print(f"Nikou Gunnemann's approach Variant {variant} F1 Score:", f1_scores_nik[variant])
18
```

Standard LR F1 Score: 0.7636363636363638
 Bahnsen approach Cost-Sensitive LR F1 Score: 0.46003807139690733
 Nikou Gunnemann's approach Variant A F1 Score: 0.756384222509384
 Nikou Gunnemann's approach Variant B F1 Score: 0.1850216662578694
 Nikou Gunnemann's approach Variant C F1 Score: 0.7663705359769372
 Nikou Gunnemann's approach Variant D F1 Score: 0.7667671216687348

In [16]:

```

1 # Compute the cost according to some method outlined in the PDF
2 def compute_cost(predictions, true_labels, costs):
3     cm = confusion_matrix(true_labels, predictions)
4     tp_cost = costs[:, 0] * cm[1, 1] # True Positives
5     fn_cost = costs[:, 1] * cm[1, 0] # False Negatives
6     fp_cost = costs[:, 2] * cm[0, 1] # False Positives
7     tn_cost = costs[:, 3] * cm[0, 0] # True Negatives
8     total_cost = tp_cost.sum() + fn_cost.sum() + fp_cost.sum() + tn_cost.sum()
9     return total_cost
10
11 # Compute costs for standard Logistic regression and Bahnsen's approach
12 cost_lr = compute_cost(predictions_lr, y_test, cost_matrix_test)
13 cost_cslr_bun = compute_cost(predictions_cslr_bun, y_test, cost_matrix_test)
14
15 # Dictionary to store costs for Nikou Gunnemann's approach variants
16 costs_nik = {}
17 variants = ['A', 'B', 'C', 'D']
18 for variant in variants:
19     predictions_variant = np.round(sigmoid(np.dot(X_test, models_nik[variant])))
20     costs_nik[variant] = compute_cost(predictions_variant, y_test, cost_matrix_test)
21
22 # Print results
23 print("Standard LR Cost:", cost_lr)
24 print("Bahnsen approach Cost-Sensitive LR Cost:", cost_cslr_bun)
25 for variant in variants:
26     print(f"Nikou Gunnemann's approach Variant {variant} Cost:", costs_nik[variant])
27

```

Standard LR Cost: 39207831424.13624

Bahnsen approach Cost-Sensitive LR Cost: 5231416704.0

Nikou Gunnemann's approach Variant A Cost: 42817146185.82045

Nikou Gunnemann's approach Variant B Cost: 104513950986.5759

Nikou Gunnemann's approach Variant C Cost: 37255560550.78639

Nikou Gunnemann's approach Variant D Cost: 37630307950.05574

```
In [17]: 1 def plot_confusion_matrix(cm, class_names, title='Confusion Matrix'):
2     """
3         Plots a confusion matrix using Plotly's Figure Factory.
4
5         Parameters:
6             cm (array-like): A confusion matrix of shape (n_classes, n_classes).
7             class_names (list of strings): List of class names, in the order they
8             title (str): Title of the heatmap. Default is 'Confusion Matrix'.
9             """
10        # Convert the confusion matrix to a DataFrame
11        data = np.array(cm)
12
13        # Create the heatmap
14        fig = ff.create_annotated_heatmap(
15            z=data,
16            x=class_names,
17            y=class_names,
18            annotation_text=data.astype(str),
19            showscale=True,
20            colorscale='Viridis'
21        )
22
23        # Add titles and labels
24        fig.update_layout(
25            title=title,
26            xaxis=dict(title='Predicted Label'),
27            yaxis=dict(title='True Label')
28        )
29
30        # Improve the xaxis and yaxis visibility
31        fig.update_xaxes(side="bottom")
32
33        # Add colorbar
34        fig['data'][0]['showscale'] = True
35
36        # Show the figure
37        fig.show()
38
39
40
41
```

```
In [18]: 1 class_names = ['Negative', 'Positive']
2
```

Confusion Matrix Comparison

```
In [19]: 1 cm_cslr_bun = confusion_matrix(y_test, predictions_cslr_bun)
2 plot_confusion_matrix(cm_cslr_bun, class_names, title='Bahnsen approach Co
```

```
In [20]: 1 variants = ['A', 'B', 'C', 'D'] # Define your variants
2
3 for variant in variants:
4     # Generate predictions for each variant; ensure you have models for each
5     predicted_labels_variant = np.round(sigmoid(np.dot(X_test, models_nik)))
6
7     # Calculate confusion matrix for the variant
8     cm_variant = confusion_matrix(y_test, predicted_labels_variant)
9
10    # Plot the confusion matrix
11    plot_confusion_matrix(cm_variant, class_names, title=f"Variant {variant} Confusion Matrix")
12
```

```
In [21]: 1 cm_lr = confusion_matrix(y_test, predictions_lr)
2 plot_confusion_matrix(cm_lr, class_names, title='Standard Logistic Regression Confusion Matrix')
```

Cost Matrix Analysis for Logistic Regression Models

In [22]:

```
1 def plot_cost_matrix(cost_matrix, title='Cost Matrix'):
2     """
3         Plots a cost matrix using Plotly's Figure Factory.
4
5         Parameters:
6             cost_matrix (DataFrame): A DataFrame representing the cost matrix.
7             title (str): Title of the heatmap.
8         """
9         # Convert the DataFrame to a numpy array for plotting
10        data = cost_matrix.values
11        class_names = cost_matrix.columns.tolist() # Assumes columns are names
12
13        # Create the heatmap
14        fig = ff.create_annotated_heatmap(
15            z=data,
16            x=class_names,
17            y=cost_matrix.index.tolist(),
18            annotation_text=data.astype(str),
19            showscale=True,
20            colorscale='Blues'
21        )
22
23        # Add titles and labels
24        fig.update_layout(
25            title=title,
26            xaxis=dict(title='Predicted Condition'),
27            yaxis=dict(title='Actual Condition')
28        )
29
30        # Improve the xaxis and yaxis visibility
31        fig.update_xaxes(side="bottom")
32
33        # Add colorbar
34        fig['data'][0]['showscale'] = True
35
36        # Show the figure
37        fig.show()
38
```

```

In [23]: 1 # Assume costs for FP and FN are given or calculated elsewhere in the anal
2 cost_fp = 6 # Example cost for False Positives
3 cost_fn = 6 # Example cost for False Negatives
4 cost_tp = 0 # Usually, the cost of True Positives is zero
5 cost_tn = 0 # Usually, the cost of True Negatives is zero
6
7 def calculate_cost_matrix(predictions, y_test):
8     cm = confusion_matrix(y_test, predictions, labels=[1, 0])
9
10    # Calculating costs based on the confusion matrix
11    tp_cost = cm[0, 0] * cost_tp
12    fn_cost = cm[0, 1] * cost_fn
13    fp_cost = cm[1, 0] * cost_fp
14    tn_cost = cm[1, 1] * cost_tn
15
16    cost_matrix = pd.DataFrame({
17        'Predicted Positive': [tp_cost, fp_cost],
18        'Predicted Negative': [fn_cost, tn_cost]
19    }, index=['Actual Positive', 'Actual Negative'])
20
21    return cost_matrix
22
23 # Assume `models_nik` and `X_test` are defined
24 variants = ['A', 'B', 'C', 'D'] # Define your variants
25 cost_matrices_nik = {} # Dictionary to store cost matrices for each variant
26
27 for variant in variants:
28     predicted_labels_variant = np.round(sigmoid(np.dot(X_test, models_nik[variant].T)))
29     cost_matrices_nik[variant] = calculate_cost_matrix(predicted_labels_variant, y_test)
30
31 # Calculate cost matrices for both models
32 cost_matrix_lr = calculate_cost_matrix(predictions_lr, y_test)
33 cost_matrix_cslr_bun = calculate_cost_matrix(predictions_cslr_bun, y_test)
34
35 print("Cost Matrix for Standard Logistic Regression:")
36 plot_cost_matrix(cost_matrix_lr, title="Standard Logistic Regression Cost Matrix")
37
38 print("Bahnsen approach Cost Matrix for Cost-Sensitive Logistic Regression:")
39 plot_cost_matrix(cost_matrix_cslr_bun, title="Bahnsen Approach Cost-Sensitive Logistic Regression Cost Matrix")
40
41 for variant in variants:
42     print(f"Nikou Gunnemann's approach Cost Matrix for Variant {variant}:")
43     plot_cost_matrix(cost_matrices_nik[variant], title=f"Variant {variant} Cost Matrix")
44

```

Cost Matrix for Standard Logistic Regression:

Bahnsen approach Cost Matrix for Cost-Sensitive Logistic Regression:

Nikou Gunnemann's approach Cost Matrix for Variant A:

Nikou Gunnemann's approach Cost Matrix for Variant B:

Nikou Gunnemann's approach Cost Matrix for Variant C:

Nikou Gunnemann's approach Cost Matrix for Variant D:

Savings Calculation Explanation

The `calculate_savings` function compares the cost efficiency between standard logistic regression and another classification technique.

Function Details

- **Purpose:** To quantify the economic benefit of using an alternative technique over standard logistic regression in terms of cost savings.
- **Parameters:**
 - `cost_lr` : The total misclassification cost from using standard logistic regression.
 - `cost_x` : The total misclassification cost from using the alternative technique.

Savings Calculation

The savings are calculated using the formula:

$$\text{Savings} = \frac{\text{Cost}_{\text{LR}} - \text{Cost}_x}{\text{Cost}_{\text{LR}}}$$

- **Interpretation:** The savings score represents the proportion of costs saved by using the new technique. A higher score indicates greater savings.

In [24]:

```
1 def calculate_savings(cost_lr, cost_x):
2     """
3         Calculate the savings of a technique compared to logistic regression.
4
5         Parameters:
6             cost_lr (float): The misclassification cost using logistic regression.
7             cost_x (float): The misclassification cost using the new technique.
8
9         Returns:
10            float: The savings score.
11            """
12
13     savings = (cost_lr - cost_x) / cost_lr
14     return savings
```

```
In [25]: 1 cost_savings_bun = calculate_savings(cost_lr, cost_cslr_bun)
2 print(f"Cost savings for Bahnsen approach: {cost_savings_bun}")
3
4 # Dictionary to store cost savings for Nikou Gunnemann's approach variants
5 cost_savings_nik = {}
6
7 for variant in variants:
8     # Assume cost_matrices_nik contains the cost matrices for each variant
9     cost_savings_nik[variant] = calculate_savings(cost_lr, losses[variant])
10    print(f"Cost savings for Nikou Gunnemann's approach Variant {variant}:")
11
```

Cost savings for Bahnsen approach: 0.8665721486248904
 Cost savings for Nikou Gunnemann's approach Variant A: 0.999875796875024
 Cost savings for Nikou Gunnemann's approach Variant B: 0.9999997901926627
 Cost savings for Nikou Gunnemann's approach Variant C: 0.9999997959887738
 Cost savings for Nikou Gunnemann's approach Variant D: 0.999999870121279

```
In [26]: 1 # Prepare the labels and values for plotting
2 labels = ['Bahnsen approach'] + [f'Variant {v}' for v in variants]
3 savings = [cost_savings_bun] + [cost_savings_nik[v] for v in variants]
4
5 import plotly.graph_objects as go
6
7 # Create a bar chart
8 fig = go.Figure(go.Bar(
9     x=labels,
10    y=savings,
11    text=[f"{s:.2f}%" for s in savings], # Formatting the text on the bars
12    textposition='auto',
13    marker_color='blue' # Color of the bars
14))
15
16 # Add titles and labels
17 fig.update_layout(
18     title='Cost Savings Compared to Standard Logistic Regression',
19     xaxis_title='Model Variant',
20     yaxis_title='Savings (%)',
21     template='plotly_white'
22)
23
24 # Show the plot
25 fig.show()
26
```