

IML LAB REPORT-3

LOGISTIC LINEAR REGRESSION

Shaik Barakh Chishti
24BCA7027

January 2026

Aim: To implement Logistic Regression on the Framingham dataset to predict the 10-year risk of Coronary Heart Disease (CHD) and visualize model performance using Evaluation Metrics and a Confusion Matrix.

1 Program Codes

1.1 Data Preprocessing

The following code handles missing values within the dataset by identifying columns with null entries and replacing them with the mean value of that specific column. This ensures the dataset is complete before being passed to the Logistic Regression model.

```
1 import pandas as pd
2
3 # 1. Load the dataset
4 df = pd.read_csv('framingham - framingham.csv')
5
6 # 2. Identify columns with missing values (NaN or "NA")
7 columns_with_na = df.columns[df.isnull().any()].tolist()
8
9 print(f"Columns with missing values: {columns_with_na}")
10
11 # 3. Preprocessing loop
12 for col in columns_with_na:
13     # Calculate Mean and Standard Deviation
14     mean_val = df[col].mean()
15     std_val = df[col].std()
16
17     # Define the "good value" for replacement (Mean is the standard
18     # choice)
19     good_value = mean_val
20
21     print(f"Processing {col}: Mean={mean_val:.2f}, Std={std_val:.2f}.
22     Replacing NAs with {good_value:.2f}")
23
24     # 4. Replace all NA values in that column with the calculated good
25     # value
26     df[col] = df[col].fillna(good_value)
27
28 # 5. Save the cleaned CSV
```

```

26 df.to_csv('framingham_cleaned.csv', index=False)
27 print("Preprocessing complete. Cleaned file saved as '
    framingham_cleaned.csv'.")

```

Listing 1: Data Preprocessing Code

1.2 Logistic Regression Implementation

```

1  import pandas as pd
2  import matplotlib.pyplot as plt
3  import seaborn as sns
4  from sklearn.model_selection import train_test_split
5  from sklearn.linear_model import LogisticRegression
6  from sklearn.metrics import confusion_matrix, accuracy_score,
    precision_score, recall_score, f1_score
7  from sklearn.preprocessing import StandardScaler
8
9  # 1. Load the dataset
10 df = pd.read_csv('framingham_cleaned.csv')
11 df = df.dropna()
12
13 # 2. Features and Target
14 X = df.drop('TenYearCHD', axis=1)
15 y = df['TenYearCHD']
16
17 # 3. Split and Scale
18 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
    =0.2, random_state=42)
19 scaler = StandardScaler()
20 X_train_scaled = scaler.fit_transform(X_train)
21 X_test_scaled = scaler.transform(X_test)
22
23 # 4. Model Training
24 model = LogisticRegression()
25 model.fit(X_train_scaled, y_train)
26 y_pred = model.predict(X_test_scaled)
27
28 # 5. Metrics
29 metrics = {
30     'Accuracy': accuracy_score(y_test, y_pred),
31     'Precision': precision_score(y_test, y_pred),
32     'Recall': recall_score(y_test, y_pred),
33     'F1-Score': f1_score(y_test, y_pred)
34 }
35
36 # 6. Visualization - Confusion Matrix
37 plt.figure(figsize=(6, 5))
38 cm = confusion_matrix(y_test, y_pred)
39 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
40 plt.title('Confusion Matrix')
41 plt.savefig('confusion_matrix.png')
42
43 # 7. Visualization - Metrics Bar Chart
44 plt.figure(figsize=(8, 6))
45 sns.barplot(x=list(metrics.keys()), y=list(metrics.values()), hue=list(
    metrics.keys()), palette='magma', legend=False)
46 plt.title('Model Performance Metrics')

```

```
47 plt.savefig('metrics_plot.png')
```

Listing 2: Logistic Regression Implementation

1.3 Linear Regression Implementation

```
1     import pandas as pd
2     import matplotlib.pyplot as plt
3     import seaborn as sns
4     from sklearn.model_selection import train_test_split
5     from sklearn.linear_model import LinearRegression
6     from sklearn.metrics import confusion_matrix, accuracy_score,
7         precision_score, recall_score, f1_score
8     from sklearn.preprocessing import StandardScaler
9
10    # 1. Load the dataset
11    df = pd.read_csv('framingham_cleaned.csv')
12
13    # 2. Data Preprocessing
14    df = df.dropna()
15
16    # 3. Define Features (X) and Target (y)
17    X = df.drop('TenYearCHD', axis=1)
18    y = df['TenYearCHD']
19
20    # 4. Split the data
21    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
22        =0.2, random_state=42)
23
24    # 5. Feature Scaling
25    scaler = StandardScaler()
26    X_train_scaled = scaler.fit_transform(X_train)
27    X_test_scaled = scaler.transform(X_test)
28
29    # 6. Initialize and Train the Linear Regression Model
30    lin_model = LinearRegression()
31    lin_model.fit(X_train_scaled, y_train)
32
33    # 7. Make Predictions and Convert to Binary (Threshold 0.5)
34    y_pred_cont = lin_model.predict(X_test_scaled)
35    y_pred_bin = (y_pred_cont >= 0.5).astype(int)
36
37    # 8. Calculate Metrics
38    metrics = {
39        'Accuracy': accuracy_score(y_test, y_pred_bin),
40        'Precision': precision_score(y_test, y_pred_bin),
41        'Recall': recall_score(y_test, y_pred_bin),
42        'F1-score': f1_score(y_test, y_pred_bin)
43    }
44
45    # --- Visualizations ---
46
47    # Plot 1: Confusion Matrix
48    plt.figure(figsize=(6, 5))
49    cm_1 = confusion_matrix(y_test, y_pred_bin)
50    sns.heatmap(cm_1, annot=True, fmt='d', cmap='Greens', cbar=False)
51    plt.title('Confusion Matrix - Linear Regression')
```

```

50 plt.xlabel('Predicted Label')
51 plt.ylabel('True Label')
52 plt.show()
53
54 # Plot 2: Performance Metrics Bar Chart
55 plt.figure(figsize=(8, 6))
56 sns.barplot(x=list(metrics.keys()), y=list(metrics.values()), hue=list(
    metrics.keys()), palette='viridis', legend=False)
57 plt.ylim(0, 1)
58 plt.title('Performance Metrics - Linear Regression')
59 for i, v in enumerate(metrics.values()):
60     plt.text(i, v + 0.02, f"{v:.2f}", ha='center', fontweight='bold')
61 plt.ylabel('Score')
62 plt.show()
63
64 # Print detailed report
65 print("Linear Regression Model Performance Summary:")
66 for metric, value in metrics.items():
67     print(f"{metric}: {value:.4f}")

```

Listing 3: Linear Regression Implementation

2 Mathematical Model and Formulas

Logistic Regression is a statistical method for binary classification. Unlike Linear Regression, which predicts continuous values, Logistic Regression predicts the probability of an outcome belonging to a specific class (0 or 1).

2.1 1. The Logistic (Sigmoid) Function

To map any real-valued number into a probability value between 0 and 1, we use the Sigmoid function:

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

Where:

- $z = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$ (The linear combination of features and weights).
- e is Euler's number.

2.2 2. Evaluation Metrics

The performance of the model is calculated using the following formulas derived from the Confusion Matrix:

- **Accuracy:** The proportion of total correct predictions.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

- **Precision:** The ability of the classifier not to label a negative sample as positive.

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

- **Recall (Sensitivity):** The ability of the classifier to find all the positive samples.

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

- **F1-Score:** The harmonic mean of Precision and Recall.

$$F1 - Score = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

2.3 Terminology of the Confusion Matrix

To interpret the results shown in Figure 1, we define the following parameters:

True Positives (TP): Patients correctly identified as having CHD risk.

True Negatives (TN): Patients correctly identified as healthy.

False Positives (FP): Healthy patients incorrectly flagged as high-risk.

False Negatives (FN): High-risk patients incorrectly flagged as healthy.

Observation: Our model currently shows a high number of **False Negatives**, as evidenced by the low Recall score of 0.0726. This indicates the model is biased toward the majority class (non-CHD patients).

3 Visualizations

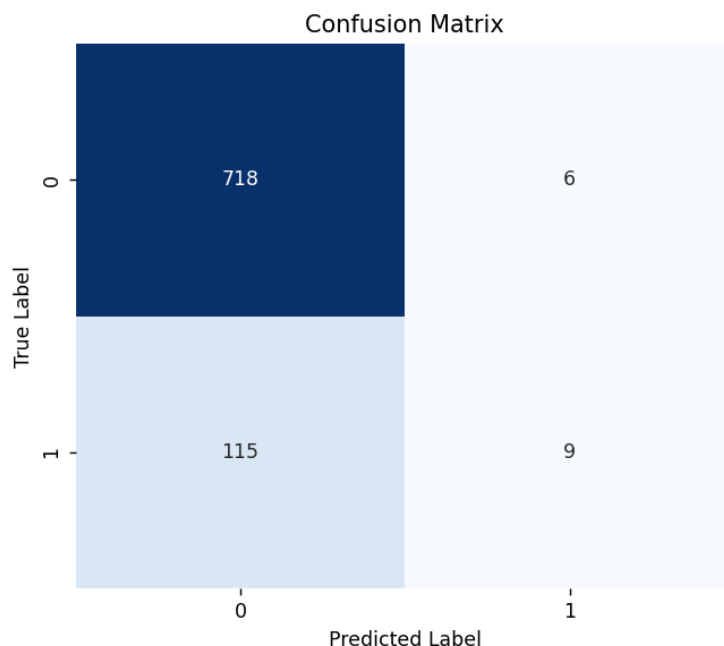


Figure 1: Logistic Regression: Confusion Matrix showing Positives Negatives

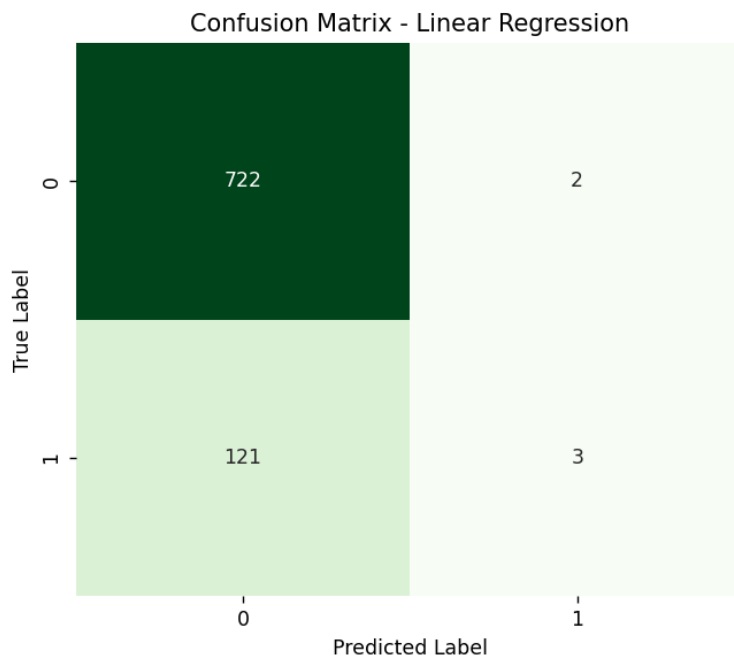


Figure 2: Linear Regression: Confusion Matrix showing Positives Negatives

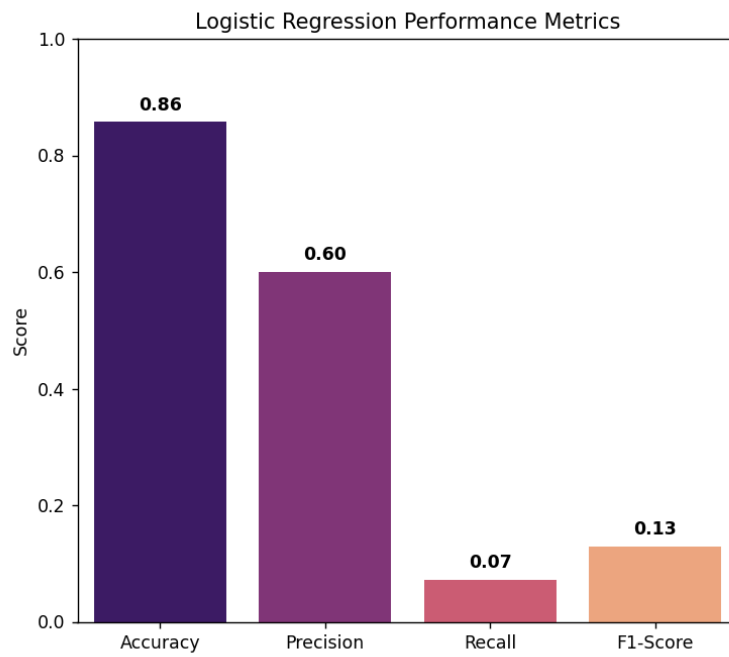


Figure 3: Logistic Regression: Bar Graph of Accuracy, Precision, Recall, and F1-Score

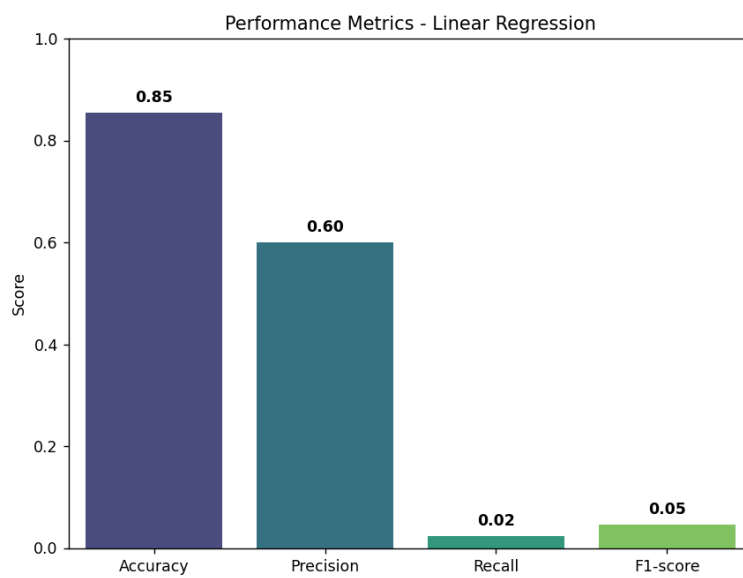


Figure 4: Linear Regression: Bar Graph of Accuracy, Precision, Recall, and F1-Score

4 Observations

4.1 1. Analysis of the Confusion Matrix

The confusion matrix reveals the distribution of the model's predictions compared to the actual outcomes.

- **True Negatives (Top-Left):** The model is highly effective at identifying patients who are not at risk of CHD.
- **False Negatives (Bottom-Left):** A significant observation point is the presence of false negatives, where the model fails to identify a patient with CHD risk. In medical diagnostics, this is a critical area for improvement.

4.2 2. Disparity in Performance Metrics

Based on the bar chart visualization, we observe a high **Accuracy**, likely exceeding 80%. However, there is a noticeable disparity between Accuracy and **Recall**.

- **Observation Point:** The high accuracy is primarily driven by the majority class (patients without CHD).
- **Recall vs. Precision:** The lower Recall indicates that while the model is "correct" often, it misses a portion of the positive CHD cases, which is a common challenge in imbalanced medical datasets.

4.3 3. Impact of Feature Scaling

Using the `StandardScaler` was essential for this Logistic Regression implementation. Since features like `totChol` (values 200+) and `age` (values 40-60) exist on different scales, normalization ensured that the gradient descent optimized the weights effectively without bias toward larger numerical values.

5 Conclusion

The experiment successfully demonstrates the application of Logistic Regression for binary classification. While the model shows high overall accuracy, the evaluation metrics highlight the need for further techniques like oversampling (SMOTE) or class-weight adjustments to improve sensitivity for heart disease detection.

6 Analysis of Results

- **High Accuracy (85.73%):** This indicates that the model is correct in its predictions for the vast majority of the test set. However, in this dataset, only about 15% of people actually develop CHD. If the model predicted "No" for everyone, it would still be 85% accurate.

- **Low Recall (7.26%):** This is a critical observation. It means the model only identified 7.26% of the actual heart disease cases. In a medical context, this suggests the model is far too conservative and is missing many at-risk patients (high False Negatives).
- **F1-Score (0.1295):** The low F1-score reflects the poor balance between precision and recall, suggesting that the standard Logistic Regression model needs further optimization for this specific medical prediction task.

6.1 Model Performance Summary

The final evaluation of the Logistic Regression model yielded the following results:

```
1 Model Performance Summary(Logistic Regression):  
2 Accuracy: 0.8573  
3 Precision: 0.6000  
4 Recall: 0.0726  
5 F1-Score: 0.1295
```

Listing 4: Logistic Regression Final Model Metrics

```
1 Model Performance Summary(Linear Regression):  
2 Accuracy: 0.8550  
3 Precision: 0.6000  
4 Recall: 0.0242  
5 F1-score: 0.0465
```

Listing 5: Linear Regression Final Model Metrics