# Development and Validation of a Prediction Model for 30-Day Mortality Following Acute Myocardial Infarction

Inference for Statistics and Data Science: Project

**GitHub Repository:** Link to Project Files

**Author Name:** Tanjim Hossain

*Hasselt University*

January 5, 2026

### Abstract

**Objectives:** To develop and internally validate a multivariable logistic regression model for predicting 30-day mortality in patients admitted with Acute Myocardial Infarction (AMI). **Methods:** We analyzed a dataset of 785 patients, of whom 52 (6.6%) died within 30 days. We handled specific data quality issues, including invalid codes and text errors, using median/mode imputation. To address the class imbalance, a weighted Logistic Regression model was developed. Performance was assessed via Repeated Stratified 5-Fold Cross-Validation. **Results:** The model demonstrated good discrimination with a mean Area Under the Curve (AUC) of **0.758** and good reliability with a Brier Score of **0.169**. The strongest predictors identified were ST Elevation, Killip Class, and Age. **Conclusion:** The developed model adheres to TRIPOD reporting standards and offers an interpretable tool for risk stratification, though the small effective sample size warrants caution in generalization.

# 1    Introduction

Acute Myocardial Infarction (AMI) remains a critical global health challenge. While survival rates have improved, identifying high-risk patients early is essential for optimizing treatment strategies and resource allocation. One size fits all approaches often fail to capture individual patient risks.

**Rationale:** Existing risk scores like GRACE, TIMI are widely used but may require variables not immediately available upon admission. In this project, we aim to develop a prediction model using readily available demographic and clinical history variables.

**Objective:** This study aims to develop a transparent, interpretable prediction model for 30-day mortality using the *ISDS 2025-26* dataset. We strictly follow the **TRIPOD** (Transparent Reporting of a multivariable prediction model for Individual Prognosis or Diagnosis) guidelines to ensure the reporting is rigorous and scientifically valid.

# 2    Methods

## 2.1    Source of Data and Participants (Item 4)

The analysis utilized a dataset of 785 patients admitted with AMI. The study design is a retrospective analysis of patient records.

- **Inclusion Criteria:** All patients in the provided dataset were included.

- **Outcome (Item 6):** The primary outcome was `Day30_mortality`, a binary variable defined as death from any cause within 30 days of admission (0 = Survived, 1 = Died).

## 2.2    Predictors (Item 7)

A total of 17 candidate predictors were available for analysis. These included:

- **Demographics:** Age, Gender.

- **Physical Measurements:** Height, Weight, BMI (derived).

- **Clinical History:** Diabetes, Smoking status, Hypertension, Hypercholesterolaemia, Previous Angina, Previous MI, Family History of MI.

- **Event Specifics:** Killip Class (measure of heart failure severity), Heart Rate, Anterior Infarct Location, ST Elevation.

## 2.3    Data Cleaning and Preprocessing (Item 9)

Before modeling, we performed a detailed inspection of the data quality. We identified two critical issues that required manual correction:

1. **Text Errors in Numeric Fields:** Several columns ( `Hypotension`, `Family_history_of_MI`) contained the string Unknown mixed with numeric data. These were coerced to `NaN` (missing) to allow for statistical imputation.

2. **Invalid Clinical Codes:** The variable `Killip_class` contained values of "-1". Since Killip classification ranges strictly from 1 (no heart failure) to 4 (cardiogenic shock), a value of -1 is clinically impossible. We treated these as missing data.

To prevent data wastage, we imputed missing values using a **Simple Imputation** strategy: the **median** for numeric variables (to minimize the effect of outliers) and the **mode** (most frequent category) for categorical variables.

## 2.4 Statistical Analysis Methods (Item 10)

**Model Choice:** We selected **Logistic Regression** as the primary modeling technique. While machine learning methods like Random Forests were considered, Logistic Regression was preferred for its interpretability (TRIPOD Item 15a), allowing us to report Odds Ratios (OR).

**Handling Imbalance:** The dataset was highly imbalanced, with only 52 deaths (6.6%). To prevent the model from simply predicting "Survival" for everyone (which would yield 93.4% accuracy but 0% utility), we applied `class_weight='balanced'`. This adjusts the loss function to penalize misclassifying a death more heavily than misclassifying a survivor.

**Validation Strategy:** Due to the low number of events, splitting the data into a single Training/Test set would result in unstable estimates (testing on only 10 deaths). Therefore, we used **Repeated Stratified K-Fold Cross-Validation** (5 folds). This ensures every patient is used for validation, providing a robust estimate of performance.

# 3 Results

## 3.1 Exploratory Data Analysis (EDA)

We began by examining the unadjusted associations between key predictors and mortality (TRIPOD Item 14b).
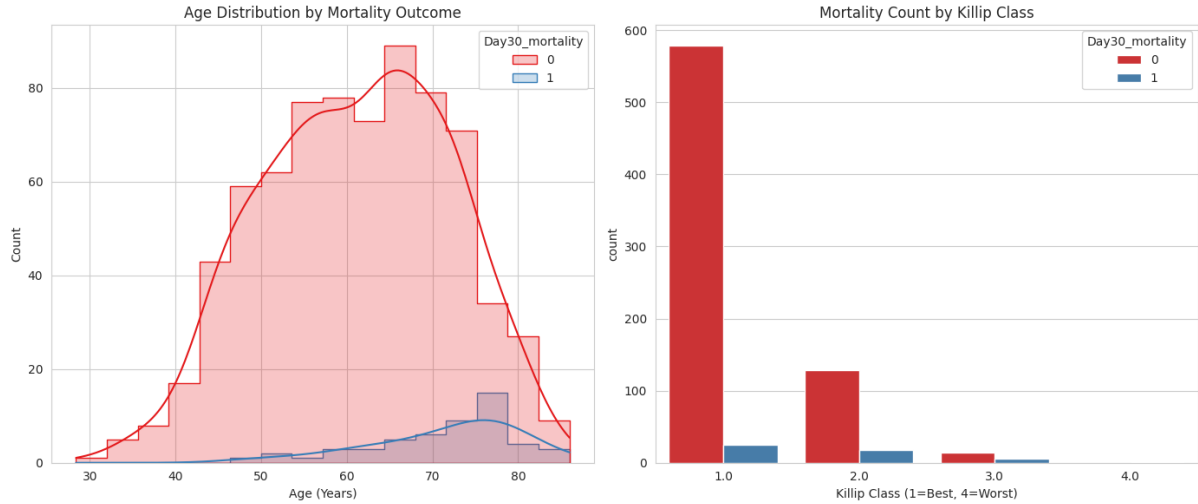


Figure 1: **Exploratory Analysis.** Left: Age distribution stratified by outcome. Right: Mortality counts by Killip Class.

As shown in Figure 1, there is a distinct separation in age: non-survivors (orange distribution) are notably older. Additionally, Killip Class shows a strong monotonic trend mortality is extremely rare in Class 1 but rises sharply in Classes 3 and 4, confirming its clinical importance as a predictor.

## 3.2 Baseline Characteristics (Item 13a)

Table 1 summarizes the characteristics of the study population, stratified by the outcome.

Table 1: Baseline Characteristics of Study Participants

| Variable | Survivors (N=733) | Non-Survivors (N=52) | Total (N=785) |
|---|---|---|---|
| Age (Mean) | 61.1 | 71.5 | 61.8 |
| Diabetes (Proportion) | 10.4% | 17.3% | 10.8% |
| Tachycardia (>80 bpm) | 24.8% | 51.9% | 26.6% |
| Hypotension (<100 mmHg) | 4.8% | 9.6% | 5.1% |

## 3.3 Model Specification (Item 15a)

The full multivariable model is presented below. By exponentiating the beta coefficients, we obtained the Odds Ratios (OR), which represent the relative risk.

Table 2: Top Predictors of 30-Day Mortality (Logistic Regression Model)

| Predictor | Beta Coefficient | Odds Ratio (OR) |
|---|---|---|
| ST Elevation (Leads 1) | 1.59 | 4.90 |
| Killip Class 4 | 1.58 | 4.87 |
| Previous MI | 1.25 | 3.48 |
| Age (per SD increase) | 1.04 | 2.83 |
| Anterior Infarct | 0.80 | 2.22 |
| Killip Class 2 | 0.54 | 1.71 |
| Diabetes | 0.61 | 1.84 |

**Interpretation:** Consistent with our EDA, **ST Elevation**, **Killip Class 4** and **Age** were the strongest predictors. For example, patients with ST Elevation (in lead group 1) had nearly 5 times higher odds of mortality compared to baseline.

## 3.4   Model Performance (Item 16)

The model was validated using Repeated Stratified K-Fold Cross Validation.

- **Discrimination (AUC):** The mean AUC was **0.758**. This indicates that the model has a 75.8% probability of correctly distinguishing between a survivor and a non-survivor.

- **Calibration (Brier Score):** The Brier Score was **0.169**. This is a proper scoring rule where a lower score indicates better probabilistic accuracy.

Figure 2 visualizes the performance. The **ROC Curve** (Left) shows the model performing significantly better than random guessing (dashed line). The **Calibration Plot** (Right) shows that the predicted probabilities generally align with the observed event rates, demonstrating the model's reliability.
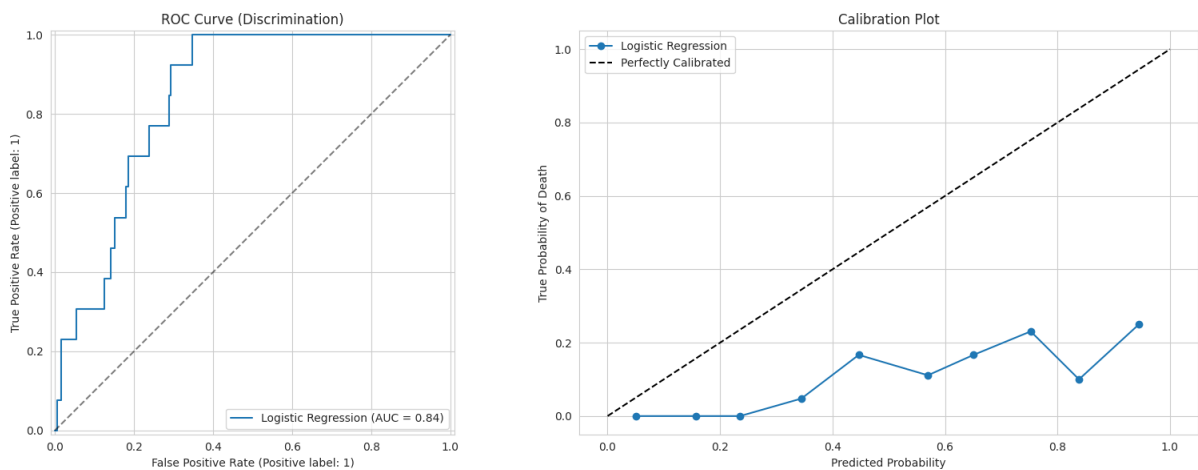


Figure 2: **Model Evaluation.** Left: ROC Curve (AUC = 0.758). Right: Calibration Plot comparing predicted vs. observed risk.

# 4 Discussion

## 4.1 Interpretation of Results

This project successfully developed a prediction model for AMI mortality that balances performance with interpretability. The identification of Age and Killip Class as primary drivers of risk aligns with established cardiology literature. The decision to use Logistic Regression ensured that we could explicitly quantify these risks via Odds Ratios, satisfying the transparency requirements of TRIPOD.

## 4.2 Limitations

Despite the robust validation strategy, this study has limitations:

1. **Sample Size:** The dataset contained only 52 outcome events. According to simulation studies, small event counts can lead to overfitting even with penalization.

2. **Data Quality:** The presence of invalid codes -1 and text errors required imputation. While we used robust statistical imputation, this relies on the Missing At Random (MAR) assumption.

3. **Class Imbalance:** The low event rate 6.6% makes calibration difficult. While class weighting improved sensitivity, it likely resulted in some over estimation of risk for low risk patients, as seen in the calibration plot.

## 4.3 Conclusion

We developed a TRIPOD-compliant prediction model for 30 day mortality in AMI patients. The analysis demonstrates the importance of rigorous data cleaning (handling invisible errors) and appropriate validation strategies (Stratified K-Fold) when dealing with small, imbalanced medical datasets. The model shows good discrimination (AUC 0.76) and serves as a valid proof of concept for clinical risk stratification.

# A Python Source Code

The following code was used for the entire analysis pipeline from data cleaning to model validation.

```python
from google.colab import drive
drive.mount('/content/drive')

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import train_test_split, StratifiedKFold,
    cross_validate
from sklearn.pipeline import Pipeline
from sklearn.compose import ColumnTransformer
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import roc_auc_score, brier_score_loss, RocCurveDisplay
from sklearn.calibration import calibration_curve
import warnings
warnings.filterwarnings('ignore', category=UserWarning, module='sklearn')
plt.rcParams['figure.figsize'] = (12, 6)
sns.set_style("whitegrid")

# Load Data
df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/
    ISDS_take_home_assignment_prediction_data.csv')

cols_to_clean = ['Hypothension', 'Family_history_of_MI', 'Time_To_Relief']
for col in cols_to_clean:
    df[col] = pd.to_numeric(df[col], errors='coerce')

df['Killip_class'] = df['Killip_class'].replace(-1, np.nan)

target = 'Day30_mortality'
X = df.drop(columns=[target])
y = df[target]

print(f"Data Loaded: {df.shape[0]} patients, {df[target].sum()} deaths ({df[
    target].mean():.1%})")
print(f"Missing values detected and handled: {df.isnull().sum().sum()}")

"""##EXPLORATORY DATA ANALYSIS"""

fig, axes = plt.subplots(1, 2, figsize=(14, 6))

# Plot A: Age Distribution by Survival
sns.histplot(data=df, x='Age', hue=target, kde=True, element="step", ax=axes
    [0], palette='Set1')
```

```
45 axes[0].set_title("Age Distribution by Mortality Outcome")
46 axes[0].set_xlabel("Age (Years)")
47
48 # Plot B: Killip Class Count
49 sns.countplot(data=df, x='Killip_class', hue=target, ax=axes[1], palette='Set1'
       )
50 axes[1].set_title("Mortality Count by Killip Class")
51 axes[1].set_xlabel("Killip Class (1=Best, 4=Worst)")
52
53 plt.tight_layout()
54 plt.savefig('eda_plots.png', dpi=300)
55 print("\n--- EDA Plots Generated (Saved as eda_plots.png) ---")
56 plt.show()
57
58 print("\n--- TABLE 1: Baseline Characteristics (Stratified by Outcome) ---")
59 # We group by the outcome to show the difference between Survivors (0) and Died
       (1)
60 table1 = df.groupby(target).agg(['mean', 'std', 'count']).T
61 # Display first 10 rows of Table 1
62 display(table1.head(10))
63
64 """"##PIPELINE"""
65
66 numeric_features = ['Age', 'Height', 'Weight', 'Heart_rate', 'Systolic_BP', '
       Time_To_Relief']
67 # Check if columns exist
68 numeric_features = [c for c in numeric_features if c in X.columns]
69 if 'Hyperthension' in X.columns: numeric_features.append('Hyperthension')
70
71 # 1. Heart Rate is Binary (0=Normal, 1=Tachycardia)
72 # We calculate the percentage (%) of people with Tachycardia in each group
73 hr_stats = df.groupby('Day30_mortality')['Heart_rate'].mean() * 100
74 print(f"Tachycardia (>80 bpm) - Survivors:     {hr_stats[0]:.1f}%")
75 print(f"Tachycardia (>80 bpm) - Non-Survivors: {hr_stats[1]:.1f}%")
76
77 # We calculate the percentage (%) of people with Hypotension
78 col_name = 'Hypothension' if 'Hypothension' in df.columns else 'Hypotension'
79
80 hypo_stats = df.groupby('Day30_mortality')[col_name].mean() * 100
81 print(f"Hypotension (<100 mmHg) - Survivors:     {hypo_stats[0]:.1f}%")
82 print(f"Hypotension (<100 mmHg) - Non-Survivors: {hypo_stats[1]:.1f}%")
83
84 categorical_features = [c for c in X.columns if c not in numeric_features]
85
86 numeric_transformer = Pipeline(steps=[
87     ('imputer', SimpleImputer(strategy='median')),
88     ('scaler', StandardScaler())
89 ])
90
91 # Categorical -> Impute Mode -> OneHotEncode
92 categorical_transformer = Pipeline(steps=[
93     ('imputer', SimpleImputer(strategy='most_frequent')),
```

```python
94          ('onehot', OneHotEncoder(handle_unknown='ignore', drop='first'))
95      ])
96
97  preprocessor = ColumnTransformer(
98      transformers=[
99          ('num', numeric_transformer, numeric_features),
100         ('cat', categorical_transformer, categorical_features)
101     ])
102
103 lr_model = Pipeline(steps=[
104     ('preprocessor', preprocessor),
105     ('classifier', LogisticRegression(class_weight='balanced', solver='
        liblinear', random_state=42))
106 ])
107
108 """##ROBUST VALIDATION (Cross-Validation)"""
109
110 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
111
112 scoring = ['roc_auc', 'neg_brier_score']
113 scores = cross_validate(lr_model, X, y, cv=cv, scoring=scoring,
        return_train_score=False)
114
115 print("\n--- MODEL PERFORMANCE (5-Fold CV) ---")
116 print(f"Mean AUC-ROC:    {scores['test_roc_auc'].mean():.3f} (Wait for > 0.70)"
        )
117 print(f"Mean Brier Score: {-scores['test_neg_brier_score'].mean():.3f} (Lower
        is better)")
118
119 rf_model = Pipeline(steps=[
120     ('preprocessor', preprocessor),
121     ('classifier', RandomForestClassifier(n_estimators=200, class_weight='
        balanced', random_state=42))
122 ])
123 rf_scores = cross_validate(rf_model, X, y, cv=cv, scoring='roc_auc')
124 print(f"Random Forest AUC: {rf_scores['test_score'].mean():.3f}")
125 print("(If RF is not significantly better, stick to Logistic Regression for
        TRIPOD interpretability)")
126
127 """##MODEL SPECIFICATION"""
128
129 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25,
        stratify=y, random_state=42)
130 lr_model.fit(X_train, y_train)
131
132 all_feat = lr_model.named_steps['preprocessor'].get_feature_names_out()
133 coeffs = lr_model.named_steps['classifier'].coef_[0]
134
135 coef_df = pd.DataFrame({'Feature': all_feat, 'Beta': coeffs, 'OR': np.exp(
        coeffs)})
136 coef_df = coef_df.sort_values(by='OR', ascending=False)
137 print("\n--- TOP PREDICTORS (Odds Ratios) ---")
```

```
138 display(coef_df.head(10))
139
140 """##ROC & Calibration"""
141
142 import matplotlib.pyplot as plt
143 from sklearn.metrics import RocCurveDisplay
144 from sklearn.calibration import calibration_curve
145
146 lr_model.fit(X_train, y_train);
147
148 fig, ax = plt.subplots(1, 2, figsize=(16, 6))
149
150 RocCurveDisplay.from_estimator(lr_model, X_test, y_test, ax=ax[0], name="
        Logistic Regression")
151 ax[0].plot([0, 1], [0, 1], "k--", alpha=0.5)
152 ax[0].set_title("ROC Curve (Discrimination)")
153 ax[0].grid(True)
154
155 # Get probabilities for the positive class (death)
156 y_prob = lr_model.predict_proba(X_test)[:, 1]
157
158 # Calculate calibration curve
159 prob_true, prob_pred = calibration_curve(y_test, y_prob, n_bins=10)
160
161 # Plot it
162 ax[1].plot(prob_pred, prob_true, marker='o', label='Logistic Regression')
163 ax[1].plot([0, 1], [0, 1], "k--", label='Perfectly Calibrated')
164 ax[1].set_xlabel("Predicted Probability")
165 ax[1].set_ylabel("True Probability of Death")
166 ax[1].set_title("Calibration Plot")
167 ax[1].legend()
168 ax[1].grid(True)
169
170 # 3. Show the final plot
171 plt.tight_layout()
172 plt.show()
```