

Predicting Pulmonary Edema for Myocardial Infarction Patients

Roshan Salil Parikh



Data Science Institute, Brown University

Providence, Rhode Island, United States of America

December 15, 2024

Contents

1	Introduction	1
1.1	Purpose	1
1.2	Data	1
1.3	Previous Work	1
2	Exploratory Data Analysis	2
3	Methods	5
3.1	Feature Selection	5
3.2	Splitting and Preprocessing	5
3.3	Evaluation Metric	6
3.4	Models Trained	6
4	Results	7
4.1	Model Performance	7
4.2	Global Feature Importance	8
4.3	Local Feature Importance	10
5	Outlook	11
6	GitHub	12

1 Introduction

1.1 Purpose

Myocardial infarctions (MI), also known as heart attacks, affect over 1 million Americans every year. A serious complication of MI is pulmonary edema, where fluid fills the lungs, making breathing more difficult. Pulmonary edema can develop suddenly after MI, and acute pulmonary edema requires immediate treatment. Machine learning prediction methods may allow hospital staff to be better prepared for potential pulmonary edema for hospitalized MI patients [1].

1.2 Data

Data came from the UCI Myocardial Infarction Complications Dataset [2]. The data contains clinical information collected during each patient’s hospital stay after suffering a MI and was collected over a 3-day period by Krasnoyarsk Interdistrict Clinical Hospital (Russia) between 1992-1995. The original dataset has 1700 examples and 111 features. There are 12 target variables available in the dataset, of which ‘OTEK_LANC’, which stood for pulmonary edema, was used, where a value of 1 indicated a pulmonary edema diagnosis and a value of 0 indicated no pulmonary edema diagnosis. Every example had missing values, but none of the target values did.

1.3 Previous Work

Previous work on this dataset conducted by researchers at Urmia University in Iran showed strong performance when predicting fatality from MI complications [3]. However, while f_1 scores were high, accuracy was close to or below baseline.

Newaz, et al., accounted for this dataset’s imbalance through synthetic minority over-

sampling (SMOTE) and cost-based learning [4]. They predicted each of the target variables with different machine learning methods, obtaining accuracies close to baseline but relatively high precision.

2 Exploratory Data Analysis

Exploratory Data Analysis (EDA) was performed to better understand the dataset with respect to the target variable. The data was imbalanced, with 9.40% of patients developing pulmonary edema.

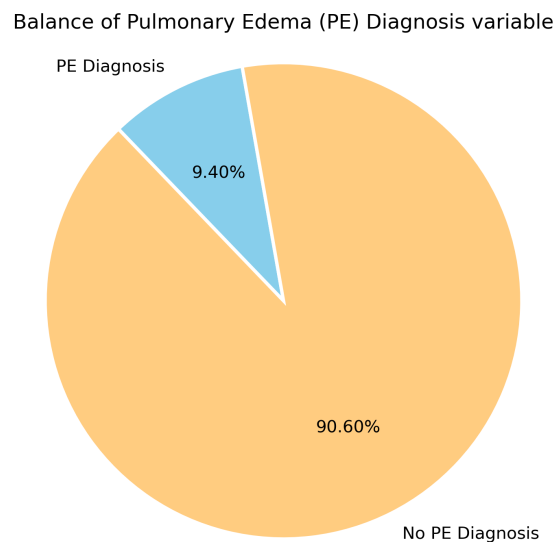


Figure 1: *Pie chart showing share of pulmonary edema diagnoses.*

We first look at metrics available from bloodwork. Aside from MI, causes of pulmonary edema can include infection. White blood cell count and α -1 antitrypsin levels can increase during infection. However, EDA did not indicate a correlation between either and the target.

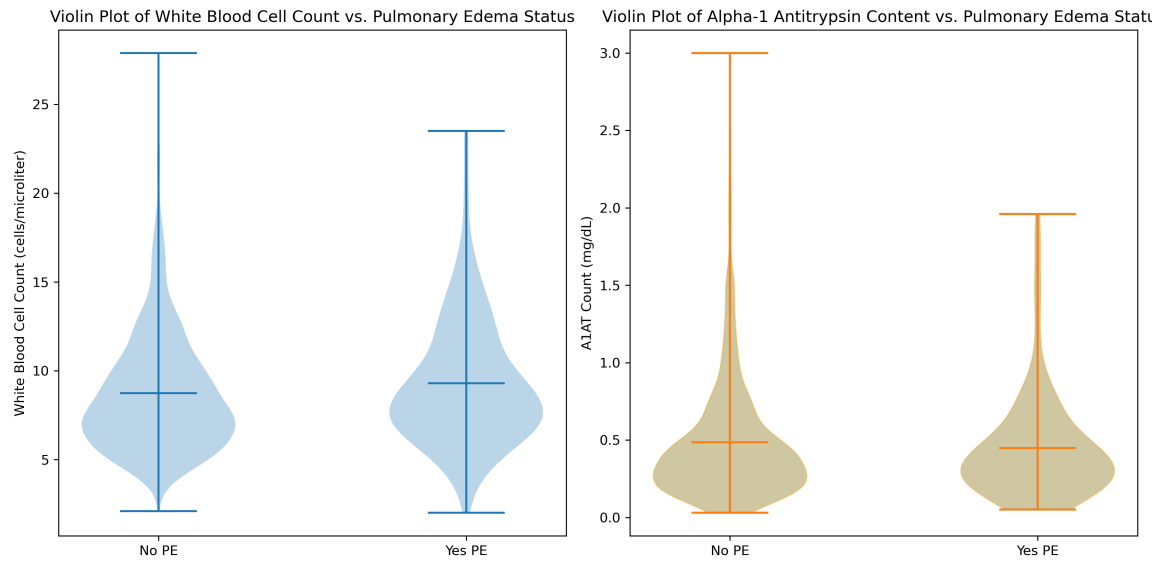


Figure 2: Violin plots showing relationships between white blood cell count or α -1 antitrypsin levels and pulmonary edema.

Serum ion concentrations did not seem to differ by mean between target variable classes, but the spread of concentrations appeared to have lower variance for patients with pulmonary edema.

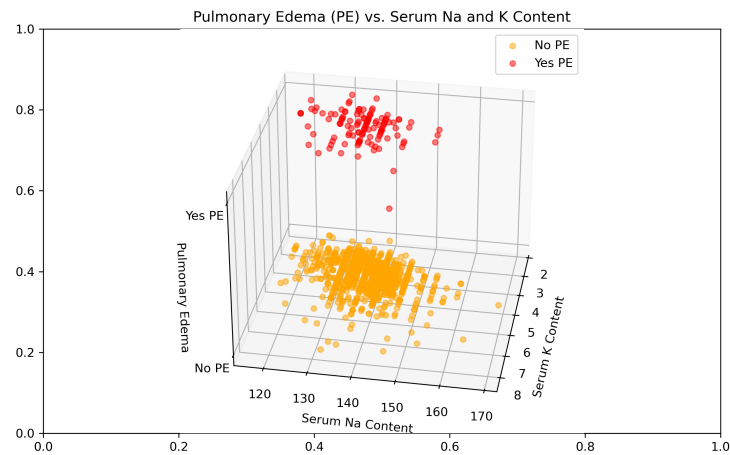


Figure 3: Serum concentrations of cations grouped by pulmonary edema diagnosis.

We also observe at the likelihood of pulmonary edema based on EKG readings.

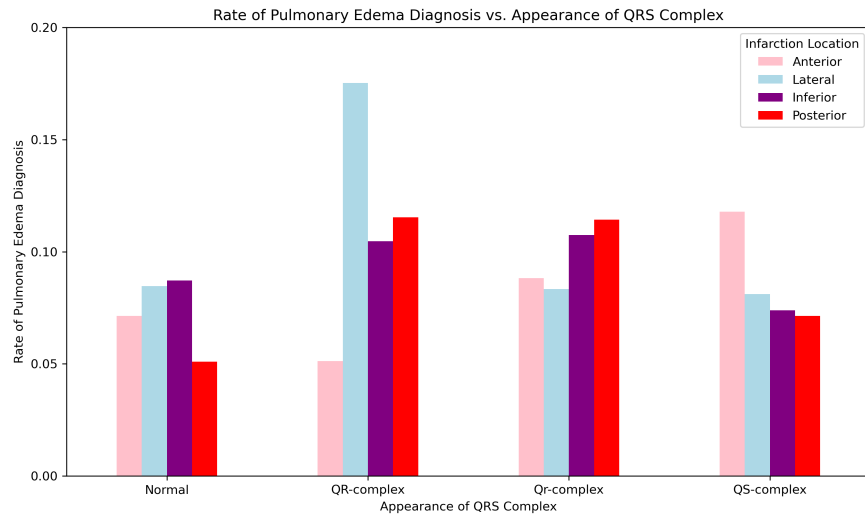


Figure 4: *The prevalence of pulmonary edema based on appearance of the QRS complex combined with the location(s) of the MI. A patient in this dataset could have MI in multiple locations, in which case their data would be included in multiple bars in this plot.*

However, EDA reveals that a patient is more likely to develop pulmonary edema if they are given liquid nitrates in the ICU. While included in the analysis, this variable may be

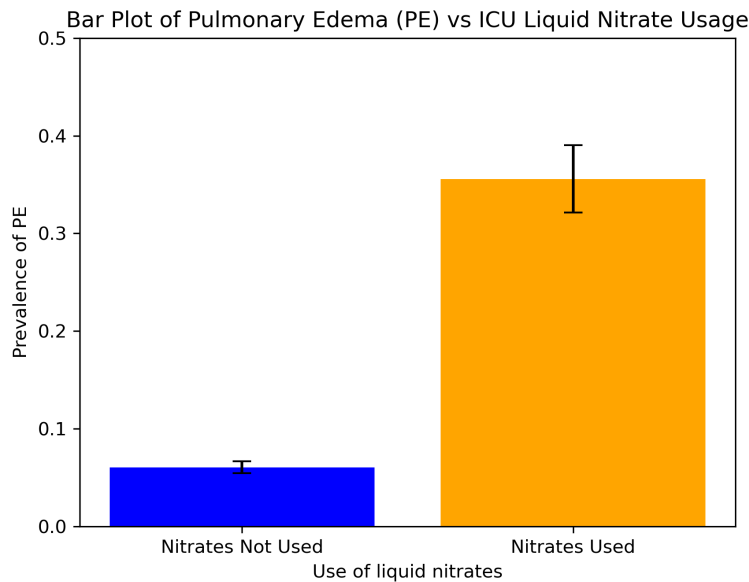


Figure 5: *Bar plots showing.*

misleading, since liquid nitrates can be given as a treatment for MI *and* as a treatment for pulmonary edema. The dataset did not specify when liquid nitrates were administered in

relation to health events.

3 Methods

3.1 Feature Selection

Features that had high Pearson correlations seemed to be clinically relevant, so no features were dropped due to this metric. Of the 111 features, only the 'SEX' feature did not have missing values. Examples that were missing age values were removed from the data, leaving 1692 examples. The 4 features (serum CPK content, heredity on coronary heart disease, and systolic and diastolic blood pressure in the ED) with more than 50% missing values were removed as well.

3.2 Splitting and Preprocessing

We first set aside 20% of the data to serve as a test set for our model. Because of the dataset's imbalance, we use scikit-learn's *StratifiedKFold* method with $n_splits = 3$ to split the remaining data into training and validation sets, each with 902 examples. *StratifiedKFold* ensures an approximately equal distribution of classes between each of the training and validation sets.

Missing values in categorical and ordinal features were labeled as a new category, 'missing,' while iterative imputing was done on continuous features with missing features ¹. With scikit-learn's *ColumnTransformer*, blood pressure and age features are encoded with the *MinMaxScaler*, ordinal features with the *OrdinalEncoder*, and categorical features with the *OneHotEncoder*. All features were then standard scaled to have a mean of 0 and variance 1. The preprocessed data had 276 features.

¹Iterative imputing was performed before implementing every model except for XGBoost, which can handle missing values in continuous features.

3.3 Evaluation Metric

f_1 -score was used as evaluation metric during cross validation. Accuracy (ℓ_{0-1} loss) was also calculated and compared. Baseline scores were calculated by predicting the majority class (class 0) for accuracy and the minority class (class 1) for f_1 -score over 5 tests sets, each with a different random seed, giving a baseline accuracy of $91.150\% \pm 0.009$ and baseline f_1 -score of 0.162 ± 0.015 .

3.4 Models Trained

Multiple types ($n = 5$) of classification models were tuned, trained, and tested. Hyperparameter tuning was performed through cross-validation with the previously-mentioned Stratified K-Folds method and with scikit-learn’s *GridSearchCV*. For each model, hyperparameters were tuned across different random seeds ($n = 5$) before the model with the highest f_1 -score for each random seed was evaluated on the testing set. This process produced 25 trained models.

Method	Hyperparameters tuned
Logistic Regression (Ridge Regularization)	$C = \{10^n n \in \{-8, \dots, 3\}\}$ class weight $\in \{\text{balanced, none}\}$
K-Nearest Neighbors	$n \text{ neighbors} \in \{3, 5, 7, 10, 15, 30, 50, 70, 100\}$ weights $\in \{\text{uniform, distance}\}$ $p \in \{1, 2\}$
Support Vector Classifier with Linear Kernel	$C = \{10^n n \in \{-5, \dots, 3\}\}$ class weight $\in \{\text{balanced, none}\}$
Support Vector Classifier with RBF Kernel	$C = \{10^n n \in \{-5, \dots, 3\}\}$ class weight $\in \{\text{balanced, none}\}$
XGBoost	max depth $\in \{1, 3, 10, 30, 100\}$ column sample by tree $\in \{0.1, 0.25, 0.5, 0.75, 1\}$ positive class weighting $\in \{0.025, 0.05, 0.1, 0.25, 0.5, 1, 5, 10\}$

Table 1: Classification methods and hyperparameters tuned through *GridSearchCV*.

4 Results

4.1 Model Performance

On average, while models had high accuracies, they tended to be slightly lower than baseline (by around 0.07%). f_1 -scores were also generally low, with no single model's f_1 score surpassing 0.36. However, aside from K-Nearest Neighbors, f_1 -scores were higher on-average than the baseline of 0.162 ± 0.015 .

Method	Avg. Accuracy - Baseline (%)	Avg. f_1 -Score
Logistic Regression	-0.104 ± 0.034	0.282 ± 0.064
K-Nearest Neighbors	-0.0136 ± 0.0092	0.0745 ± 0.063
Support Vector Classifier (Linear)	-0.0932 ± 0.036	0.241 ± 0.095
Support Vector Classifier (RBF)	-0.109 ± 0.086	0.265 ± 0.077
XGBoost	-0.0634 ± 0.0099	0.303 ± 0.017

Table 2: Average accuracy minus baseline accuracy (0.911504) and average f_1 score for each method. Models were trained and tested with 5 different random seeds.

Because of the dataset's imbalance, we use the f_1 -score as a measure of performance. The single best model was the SVC model with the RBF kernel ($f_1 = 0.358$), with parameters $C = 1$, class weight = balanced. On average, XGBoost performed the best ($f_1 = 0.303 \pm 0.017$). The single XGBoost model with the highest score ($f_1 = 0.319$) had the following parameters: column sample by tree = 0.25, max depth = 1, and positive class weighting = 5.

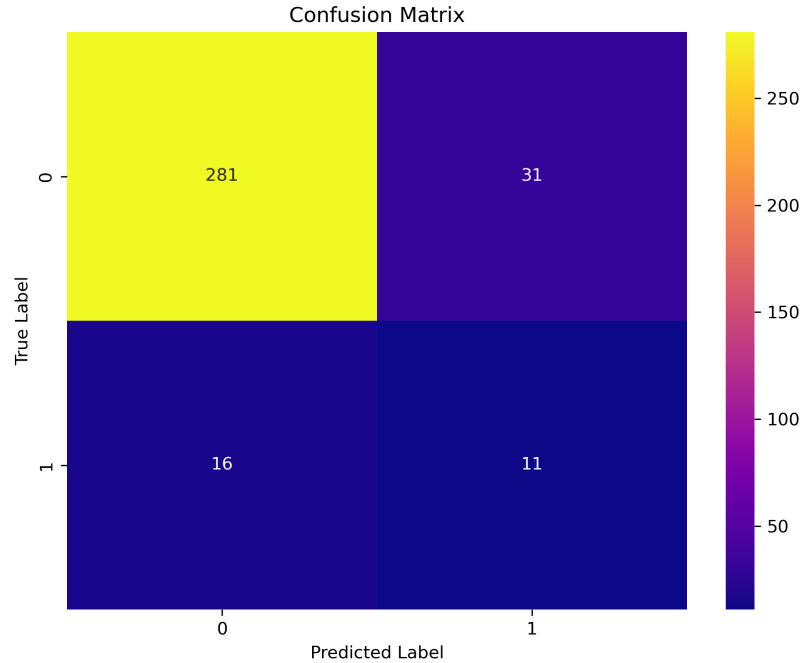


Figure 6: *Confusion matrix for XGBoost model with highest f_1 score.*

Feature importance was assessed with the best performing XGBoost model, since XGBoost had the highest average f_1 -score.

4.2 Global Feature Importance

Global feature importance was measured through various metrics, such as XGBoost's *gain* function, which measures a feature's contribution to reducing the log loss function when splitting a decision tree. Permutation importance was also calculated for each feature. Finally, SHAP values were calculated.

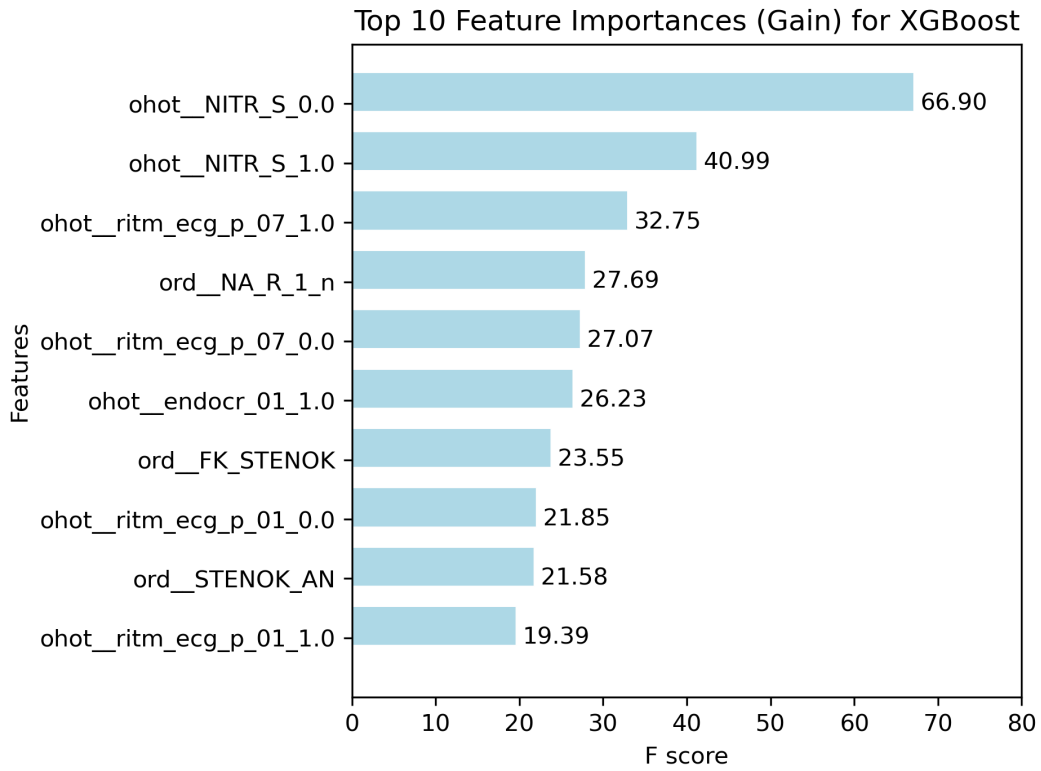


Figure 7: Gain metric for features from XGBoost model with best f_1 score

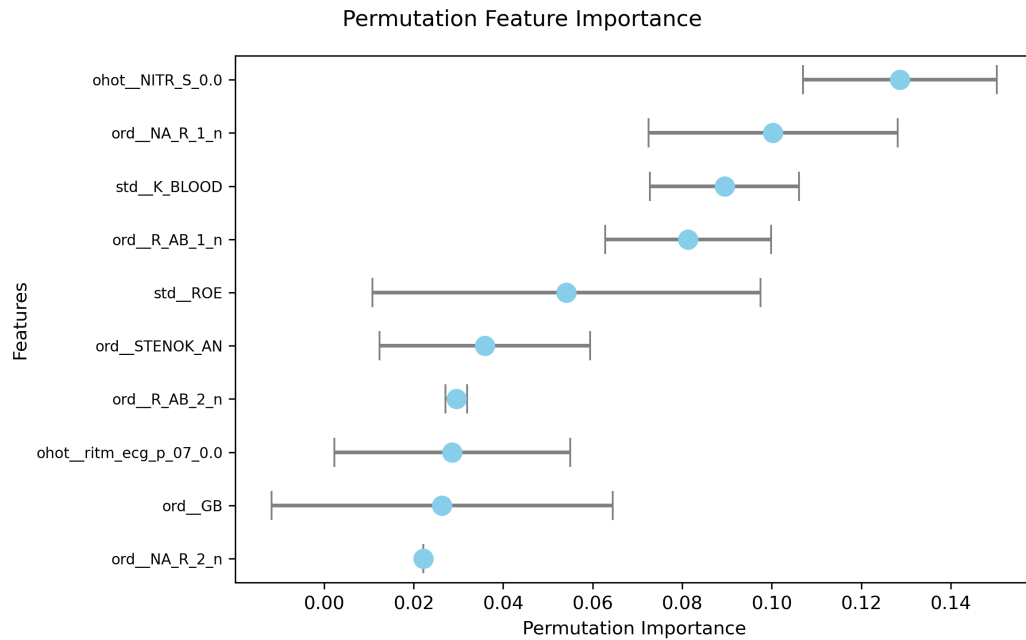


Figure 8: Permutation importance values for features from XGBoost model with best f_1 score. f_1 score was used as the evaluation metric.

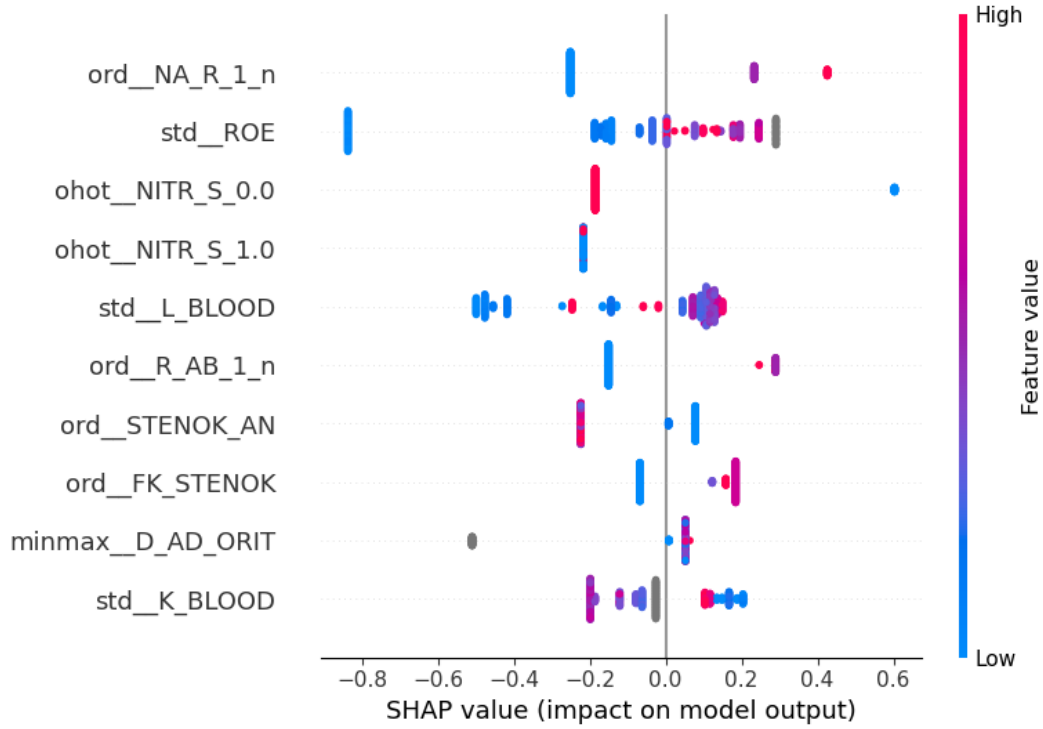


Figure 9: SHAP values for features from XGBoost model with best f_1 score. f_1 score was used as the evaluation metric.

Global feature importance analyses demonstrated that the administration of liquid nitrates ('ohot_NITR_S_0.0', 'ohot_NITR_S_1.0') was among the most important features for the model's predictive power, which followed patterns observed from EDA. Features corresponding to serum potassium levels ('std_K_BLOOD'), hospital opiod administration ('ord_NA_R_1_n', 'ord_NA_R_2_n'), EKG readings ('ohot_ritm_ecg_p_07_1.0', 'ohot_ritm_ecg_p_07_0.0'), and angina pectoris ('ord_FK_STENOK', 'ord_STENOK_AN') were also among the most important features.

4.3 Local Feature Importance

SHAP values were calculated for specific data points to observe local feature importance. Two points were arbitrarily chosen, of which one belonged to class 1.

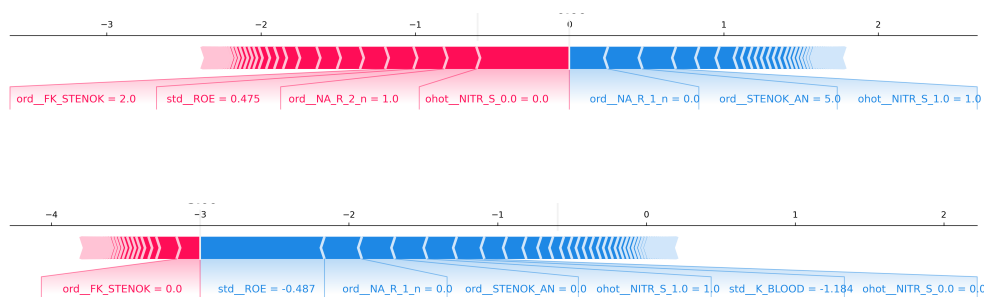


Figure 10: *SHAP force plots for points 2 (top, class 0) and 146 (bottom, class 1)*

One surprising observation is that erythrocyte sedimentation rate (ESR) (‘ROE’) was the most important factor for point 146. While ESR is not known to be directly correlated with pulmonary edema, high ESR combined with preexisting cardiac conditions can increase the risk of developing pulmonary edema. Nitrate usage was also important for both data points, as were angina pectoris features.

Unsurprisingly, features associated with cardiac measurements and cardiac health history tended to be among the top 10 most important features.

5 Outlook

Many factors may have led to the low performance of this model. First, the data contained clinical data collected over 3 days. Many of the features had metadata indicating when the data was collected, but many of them also did not or had imprecise time metadata. Nitrate administration appeared to be the most important factor by some metrics, but the reason for liquid nitrate usage was not provided, as they can be given as a response to either MI (therefore being a potential cause for pulmonary edema) or as a response to pulmonary edema. This complicates practical interpretation of the model.

Improving this model could require reducing the number of features and including the

time when clinical measurements occurred. This process may result in non-i.i.d. data, where each patient has multiple data points, but it may also result in a more powerful model. Data collected more recently and in the United States would also be more practical to eventual deployment of a predictive model, as cardiac care today in American hospitals may be significantly different from cardiac care of the 1990s in Russian hospitals. Supervised machine learning methods such as polynomial support vector classifiers may also be useful. Further, additional evaluation metrics, such as false positive rate, f_β score ($\beta \neq 1$), accuracy, and AUC, could provide more insights. Previous work indicates that fatality due to complications can be successfully predicted. The data may be more separable based on severity (or lethality) of MI complications as opposed to the presence of complications.

6 GitHub

Data, Jupyter Notebooks, and figures can be found at:

<https://github.com/roshanparikh/HeartAttackComplications.git>.

References

- [1] M. Clinic. “Pulmonary edema.” (2022), [Online]. Available: <https://www.mayoclinic.org/diseases-conditions/pulmonary-edema/symptoms-causes/syc-20377009>. (accessed: 12.02.2024).
- [2] S. Golovenkin, V. Shulman, D. Rossiev, *et al.* “Myocardial infarction complications.” DOI: <https://doi.org/10.24432/C53P5M>. (2020), [Online]. Available: <https://archive.ics.uci.edu/dataset/579/myocardial+infarction+complications>.
- [3] R. Ghafari, A. S. Azar, A. Ghafari, *et al.*, “Prediction of the fatal acute complications of myocardial infarction via machine learning algorithms,” *The Journal of Tehran Heart Center*, vol. 18, no. 4, pp. 278–287, 2023. DOI: 10.18502/jthc.v18i4.14827. [Online]. Available: <https://doi.org/10.18502/jthc.v18i4.14827>.
- [4] A. Newaz, M. S. Mohosheu, and M. A. Al Noman, “Predicting complications of myocardial infarction within several hours of hospitalization using data mining techniques,” *Informatics in Medicine Unlocked*, vol. 42, p. 101 361, 2023, ISSN: 2352-9148. DOI: <https://doi.org/10.1016/j.imu.2023.101361>. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2352914823002071>.