**Predicting Readmission for Patients with Diabetes Mellitus**

Vivian Do, Bethany Wang

University of San Diego

Shiley-Marcos School of Engineering

Master of Science, Applied Data Science

June 26, 2023

**Problem Statement**

Rehospitalizations impose a significant financial burden on the healthcare system. According to the Center for Health Information and Analysis, hospital readmissions cost Medicare $26 billion annually, with $17 billion considered avoidable (Reardon, 2015). These avoidable readmissions contribute to the rising healthcare costs and strain resources for both patients and healthcare providers. To incentivize hospitals to reduce their readmission rates, the Centers for Medicare and Medicaid Services (CMS) established the Hospital Readmission Reduction Program (HRRP). Under the program, hospitals are evaluated based on their ability to manage excess readmission and can face payment deductions as a consequence for poor performance (CMS.gov, 2023).

This data science project serves to contribute to the ongoing efforts to reduce hospital readmissions. We focus specifically on diabetes mellitus encounters and hope to gain valuable insights into the risk factors associated with readmission in diabetic patients. This targeted approach serves as a stepping stone towards a broader goal of predicting readmissions for other chronic diseases. Accurately predicting hospital readmission for diabetic patients has two major benefits. First, healthcare providers can identify patients who are most at risk and intervene in a timely manner. Secondly, there are financial incentives for hospitals to reduce admissions as this would keep costs down and prevent financial penalties. These advantages emphasize the importance of developing an accurate predictive model for diabetes with the potential for wider application in other chronic diseases.

**Data Source**

The original data comes from the Health Facts database and represents 10 years (1999-2008) of clinical patient records at 130 hospitals and integrated delivery systems across the United States. The dataset used for this project contains a subset of this database as extracted

by Strack et al. (2013) in their journal article "Impact of HbA1c Measurement on Hospital Readmission Rates: Analysis of 70,000 Clinical Database Patient Records". The extracted dataset, accessed through the UCI Machine Learning Repository, contains 101766 patient encounters and 50 features that satisfy the following criteria:

1) It is an inpatient encounter (a hospital admission)

2) It is a "diabetic" encounter (i.e "diabetes" was entered into the system as either a primary/secondary diagnosis)

3) The length of stay was between 1-14 days

4) Laboratory tests were performed during the encounter

5) Medications were administered during the encounter (Clore et al., 2014).

The 50 features can be divided into three subcategories: (1) patient demographics including race, gender, weight, payer code, (2) admission information including discharge type, admission source, admission type, whether the patient was readmitted, and (3) laboratory and medical interventions during the encounter including number of medications administered, procedures performed, blood glucose test (A1c) result, and the change in dosage for a number of drugs (Clore et al., 2014).

**Data Wrangling**

**Feature Deletion**

The initial data preparation includes the removal of biased observations and irrelevant or missing features. For example, patient encounters who were discharged to hospice or resulted in death were removed to limit bias during modeling. Other categorical variables, such as 'admission_source_id' (25 IDs), 'admission_type_id' (8 IDs), and 'discharge_disposition_id' (30 IDs) contained IDs associated with null values (e.g "Not Available", "NULL", and "Not Mapped"). These levels in their original form would be misleading for modeling and were

encoded as NA to reflect the absence of information. Features representing a patient's weight ('weight') and the medical speciality of the admitting physician ('medical_specialty') were removed due to a high percentage of null values, with 96.85% and 48.94% missing, respectively. Other attributes had less than 2% of null values, and the observations containing these nulls were removed. Degenerated variables and other irrelevant features (e.g payer code, patient ID, encounter ID) were also removed prior to modeling.

**Feature Extraction**

Primary and secondary diagnoses are encoded using the first three digits of ICD-9-CM, the official system for encoding diagnoses and procedures in hospitals in the United States (National Center for Health Statistics, 2021). The dataset contains three attributes for the primary diagnosis ('diag_1') and secondary diagnoses ('diag_2', 'diag_3), with over 800 levels for each attribute. These codes were grouped according to the 2016 edition of ICD-9-CM chapters as follows: 140-239 "Neoplasms", 250 "Diabetes", 320-459 "Circulatory", 460-519 "Respiratory", 520-579 "Digestive", 580-629 "Genitourinary System", 710-739 "Musculoskeletal", 760-779 "Perinatal", and 800-999 " Injury and Poisoning" (ICD.Codes, n.d.). All other chapters with less than 2500 instances in the primary diagnosis were grouped into "Other". These same levels were used to define the secondary diagnoses.

**Target variable**

The target variable, 'readmitted', defines the number of days until the patient is readmitted. It contains three values: "<30" if the patient was readmitted within 30 days, ">30" if the patient was readmitted after 30 days, and "NO" if the patient was not readmitted. For the purpose of this project, we consider all admissions as true regardless of the length of time between encounters. We redefined the target variable as "YES" if the patient was readmitted at

any time and "NO" if the patient was not readmitted. After, there are 41094 (47%) readmissions and 45472 (53%) 1-chance encounters (or there is no record of readmission).

**Final Dataset**

The prepared dataset contains 18 predictors and 86566 patient encounters. The target variable is binary with levels of "YES" and "NO".

**Modeling Preparation**

The dataset was split with stratification of the target variable into a training set containing 70% (59529) and a test set with 30% (25512) observations. Other preprocessing steps included standardization using centering and scaling, removal of outliers, and conversion of categorical variables into *n-1* dummy variables.

## Exploratory Data Analysis

The distributions of the number of medications administered during the encounter ('num_medications') and number of lab procedures ('num_lab_procedures') are approximately normal and centered around 16.06 and 43.42, respectively. All other numeric variables are heavily right skewed, with the majority of values falling close to 0. There is moderate correlation between the time spent in the hospital and the number of procedures performed ($r=0.33$), time spent in the hospital and the number of medications administered ($r=0.46$), as well as the number of medications administered and lab procedures performed ($r=0.38$).

**Figure 1**

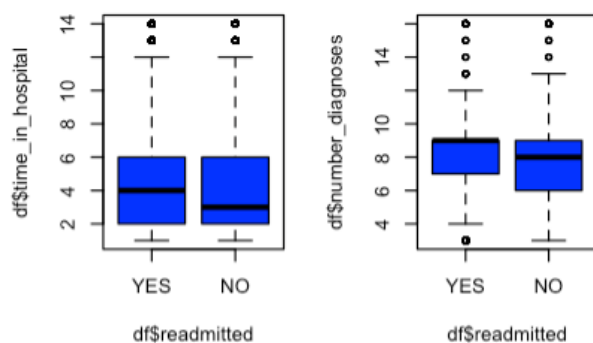*Boxplots for Readmitted vs. Time in Hospital and Number of Diagnoses*

Figure 1 shows boxplots for readmitted vs. time in hospital and number of diagnoses entered into a patient's charts. In general, patients who are readmitted spend more time in the hospital and have more diagnoses entered into the system.

**Figure 2**

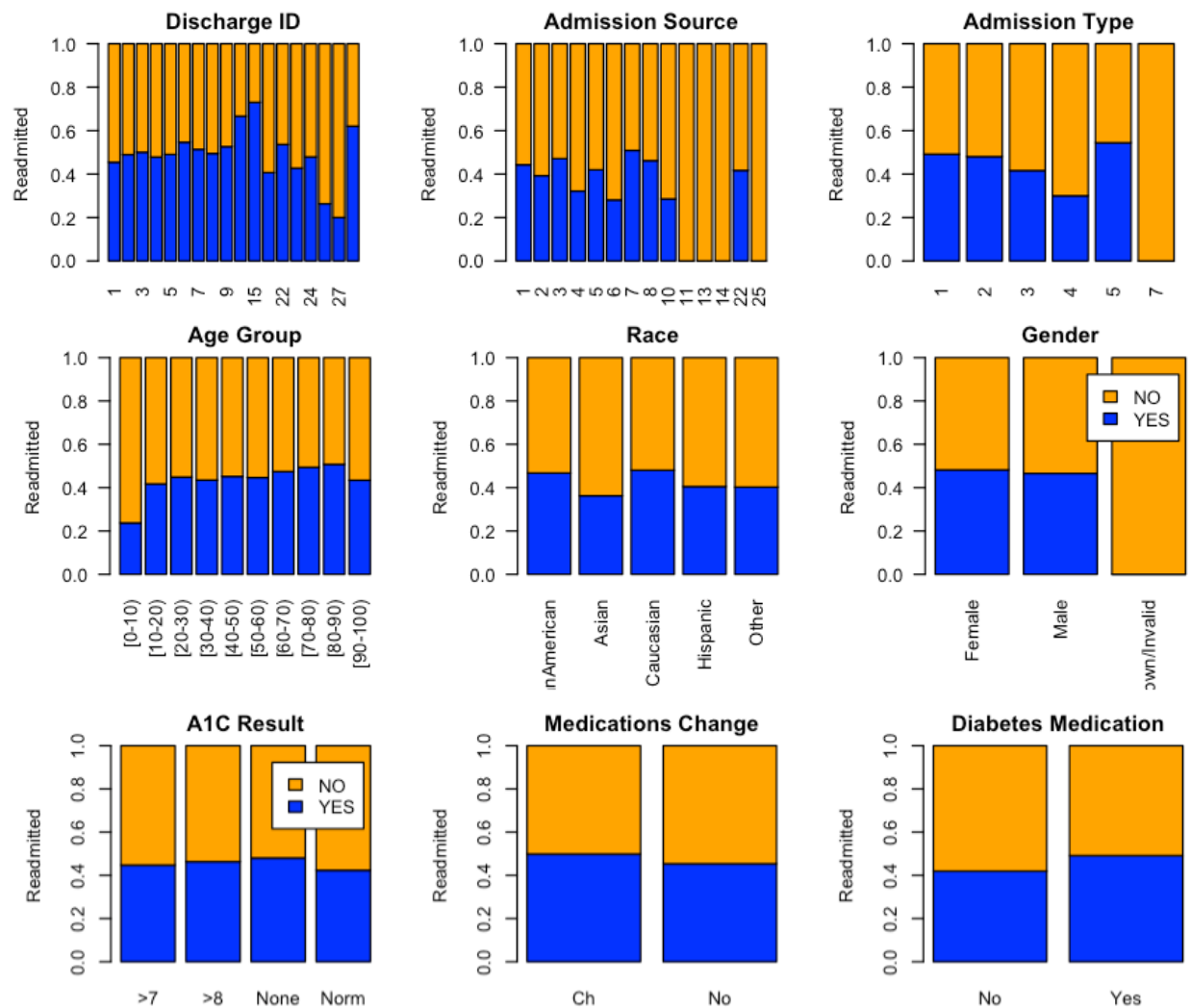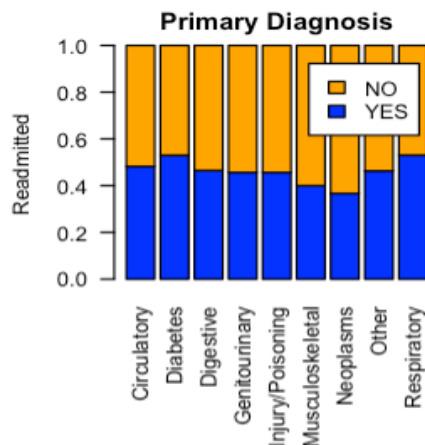*Stacked Bar Charts for Categorical Variables vs. Readmitted*



Figure 2 shows the normalized proportions of all categorical variables and the target variable. Patients who were discharged or transferred to a federal health facility (discharge disposition ID 27) were least likely to be readmitted (IDs 25 and 18 represent null values) while discharges to a Medicare approved swing bed (ID 15), outpatient services (ID 12) and psychiatric unit (ID 28) were most likely to be readmitted. Patients admitted through normal delivery

(admission source ID 11), a sick baby (ID 13), extramural birth (ID 14), and transferred from ambulatory surgery center (ID 25) had no readmissions. Trauma center admission type patients (admission type ID 7) also had no readmissions. Looking at patient demographics, the risk of readmission was seen to increase with age and if the patient was African-American or Caucasian.

The most common primary diagnoses were circulatory (27207, 31.43%), other (18242, 21.07%), respiratory (8618, 9.96%), and digestive (7873, 9.09%). 7185 (8.30%) patients had a primary diagnosis of diabetes. This suggests that most patients are hospitalized for other root causes, and later developed diabetes or are dealing with chronic diabetes concurrently. Patients with a primary diagnosis of diabetes and respiratory diseases were more likely to be readmitted compared to other diseases (Figure 3).

**Figure 3**

*Stacked Barchart for Primary Diagnosis vs Readmitted*



**Modeling**

Using the preprocessed data, we next built different classification models to predict the target variable "readmitted". We used the train dataset for training the models and used the test dataset for testing the models. A 10-fold cross-validation method is applied to optimize the training process.

When selecting the models, we considered if the model is suitable for solving binary classification problems, if it is capable of computing large dimensional data, and if it works with both numerical and categorical features. We explored linear classification models such as logistic regression, penalized logistic regression, and nearest shrunken centroids model. We also explored non-linear models including bagged tree, gradient boosted tree, random forest tree, and K-nearest neighbor models.

**Logistic Regression**

Logistic regression is a linear machine learning algorithm used for binary classification. It requires that each observation is independent of the other and assumes a linear relationship between the independent variables and the log-odds of the dependent variable. It predicts the probability of an instance belonging to a class which is then mapped to a binary value of 0 and 1 corresponding to the binary labels. The cross-validation resampling produced an accuracy of 0.59 and a sensitivity of 0.52.

**Penalized Logistic Regression Model**

Penalized logistic regression models apply regularization methods to the regular logistic regression model. We implemented the glmnet penalized model that adopts ridge and lasso penalties simultaneously, which means a squared penalty and an absolute value penalty are added on the binomial likelihood function. A tuning parameter alpha controls the mixing proportion between pure lasso penalty (alpha =1) and a pure ridge penalty (alpha=0). Another parameter lambda controls the amount of penalization.

We tuned the model with 4 alpha values (0, 0, 0.1, 0.2, 0.4) and 5 lambda values between 0.01 and 0.1. The optimal model is reached when alpha is 0 and lambda is 0.01. The resampling produced an accuracy of 0.59 and a sensitivity of 0.52.

**Nearest Shrunken Centroids Model**

The nearest shrunken centroids model is a linear classification model that assumes the centroids in the feature space are different for each target label. By taking the average of each feature in the training set, the algorithm first summarizes a set of centroids for each class. When predicting on new data, the distance between a given data instance and each centroid is calculated and the closest centroid is used to assign a class label to the query (Brownlee, 2020).

To extend the method for optimal classification, the centroids of each input variable are shrunk towards the centroid of the entire training dataset as means to select features. The amount of shrinkage is a tuning hyperparameter. In training our model, we tuned 20 shrinkage ranging between 0 and 15. The optimal model is reached when the shrinkage threshold is 0. The resampling produced an accuracy of 0.57 and a sensitivity of 0.50.

**Bagged Trees**

Bagged tree is an ensemble machine learning method used to strengthen the single decision tree model. The main purpose is to group a set of weak learners and form a strong learner. Weak learners learn from each other independently in parallel and then are combined to produce a powerful model. Bagging decreases variance, and solves over-fitting issues in a decision tree model (Vadapalli, 2022). With a number of bags parameter 30, the model produced an accuracy of 0.57, and a sensitivity of 0.54.

**Gradient Boosted Trees**

Boosted tree is another ensemble method to enhance the basic decision tree model. It starts with a base/weak learner tree, then constructs multiple tree models in sequence. Each of the trees corrects the previous one's errors. The final tree model works as a strong learner that shows the weighted mean of all the tree models. Boosted tree algorithm decreases the bias error.

In tuning the model, we created a tuning grid with the following perimeters: interaction depth (5, 7), number of trees (500), and shrinkage (0.01, 0.1). The optimal model is achieved when shrinkage is 0.1 and depth is 5. The optimal model produced a resampling performance of accuracy of 0.60 and sensitivity, 0.56.

**Random Forest**

Random forest builds a series of decision trees with a random sample of the training dataset and combines them to decide the final classification. For each decision tree, a subset of features are selected and trained on a different set of sample data. All these decision trees then form a forest. When predicting a new data record, random forest takes the prediction from each tree and decides the final output based on the majority votes.

The tuning parameters include the number of trees and the number of randomly selected predictors. We tuned the model with 100 trees and the number of predictors of 1, 3, 5, and 7. The optimal model is reached with the number of predictors 3. The optimal resampling performance gave an accuracy of 0.59 and sensitivity of 0.53.

**K-Nearest Neighbor (KNN)**

The KNN model is called a lazy learner algorithm because it does not learn from the training set immediately. It stores the available data without training them. It classifies a new data sample into a category based on its similarity to its K nearest neighbors. KNN has one tuning parameter, K, the number of neighbors. When K is too low, the model will form a very complex decision boundary resulting in poor predictions. Therefore, it's important to find the right value of K for the optimal classification result (Aggarwal, 2020). We tuned K with values of 5, 7, and 9 with 9 as the optimal value for K. The training performance produced an accuracy of 0.56 and a sensitivity of 0.53.

## Results

For this data analysis problem, a model's ability to predict the positive (readmitted-Yes) accurately is most important. Therefore, we choose the all positive model as the baseline model. Analyzing the ratio of readmission yes and no in the training dataset, 48% of cases have readmitted label 'Yes'. It indicates that if we assign all predictions as positive, the accuracy for this baseline model to predict the positive is 0.48.

The resampling performance and testing performance are evaluated using various metrics based on the confusion matrix. They include accuracy, sensitivity, specificity, precision, recall, and F score. The comparison of the models on cross-validation performance is seen in Table 1. The comparison on testing performance is seen in Table 2.

**Table 1**

*Metrics for Cross-Validation*

| Metric.Train | LR | GLMN | NSC | GBM | TRBAG | RF | KNN |
|---|---|---|---|---|---|---|---|
| 1 Accuracy | 0.585 | 0.585 | 0.57 | 0.596 | 0.572 | 0.586 | 0.555 |
| 2 Sensitivity | 0.524 | 0.521 | 0.504 | 0.561 | 0.541 | 0.532 | 0.527 |
| 3 Specificity | 0.64 | 0.642 | 0.629 | 0.628 | 0.601 | 0.635 | 0.581 |
| 4 Precision | 0.57 | 0.57 | 0.553 | 0.578 | 0.552 | 0.57 | 0.534 |
| 5 Recall | 0.524 | 0.521 | 0.504 | 0.561 | 0.541 | 0.532 | 0.527 |
| 6 F-Measure | 0.546 | 0.544 | 0.527 | 0.569 | 0.546 | 0.55 | 0.53 |
| 7 ROC | 0.617 | 0.617 | 0.6 | 0.627 | 0.603 | 0.622 | 0.576 |
| 8 AUC | 0.503 | 0.51 | 0.502 | 0.502 | 0.499 | 0.503 | 0.495 |

**Table 2**

*Metrics for Testing*

| Metric.Test | LR | GLMN | NSC | GBM | TRBAG | RF | KNN |
|---|---|---|---|---|---|---|---|
| 1 Accuracy | 0.586 | 0.586 | 0.575 | 0.601 | 0.569 | 0.587 | 0.557 |
| 2 Sensitivity | 0.524 | 0.520 | 0.507 | 0.565 | 0.535 | 0.503 | 0.528 |
| 3 Specificity | 0.643 | 0.646 | 0.637 | 0.633 | 0.601 | 0.663 | 0.583 |
| 4 Precision | 0.572 | 0.572 | 0.559 | 0.583 | 0.549 | 0.576 | 0.535 |
| 5 Recall | 0.524 | 0.520 | 0.507 | 0.565 | 0.535 | 0.503 | 0.528 |
| 6 F-Measure | 0.547 | 0.545 | 0.532 | 0.574 | 0.542 | 0.537 | 0.531 |

From Tables 1 and 2, we see all models produced similar performance metrics. The accuracy ranges from 0.56 to 0.6. The sensitivity ranges from 0.52 to 0.56. The ROC ranges from 0.58 to 0.63. And it shows consistent performance output between cross-validation and testing. Gradient boosted tree model is selected to be the final model with the best performance on all metrics except specificity. It produced an accuracy of 0.60, a sensitivity of 0.57, a specificity of 0.63, a precision of 0.58, and a F1 of 0.57 from testing.

To analyze the importance of features, we checked the first 15 important features from logistic regression, random forest, and KNN model. Though each model displays the important features in different order, some features appeared to be important for all three models, including: number_diagnoses, num_procedures, num_lab_procedures, num_medication, time_in_hospital. etc.

**Conclusion & Discussion**

We created this predictive data analysis project with a goal to help hospitals reduce readmissions for diabetic patients and understand risk factors associated with readmissions. A dataset consisting of 101766 data records with 50 features of patients with diabetic encounters was analyzed to predict patients' possible future readmissions. After wrangling the data, 18 more relevant features were extracted for creating binary classification machine learning models that predict the readmission status of yes or no. Gradient boosted tree model was selected as the final model with its computational efficiency and the best performance in both training and testing. It can improve the baseline model's ability to predict the positive (readmission yes) by 10%.

The advantage of the project is that we have an ample amount of data so that we could remove the missing values and outliers without approximating the data. Another strength is that we were able to build and compare multiple linear and non-linear classification models. All the

seven models we built produced very consistent training and testing outcomes indicating that we have reached the limits that the features might be able to predict the target of readmissions.

The weakness that refrains the models' performance lies in the fact that the dataset we used is not the original data, but extracted. Quite some useful information was lost during the extraction. For example, age, originally a numerical variable was binned into categorical, making it less useful. Some other variables lost their predictive ability likewise. Another challenge we have encountered is the lack of computing power to run more complicated and resource demanding classification models and to tune models with a larger range of tuning parameters. It turned out that when the data is highly dimensional, to train non-linear models with a bigger tuning grid requires a significant amount of running time that we could not afford. As a result, we were not able to examine more non-linear models or tune the built models with larger scale of parameters to discover the ideal model or optimal model results for this project.

To improve the project's outcome further, we should use the original dataset instead of the extracted one to retain informative features. Also, to conquer the computing power limitation, we should try to run the models on a cloud computing system, like AWS.

**References**

Aggarwal, S. (2020, June 8). *K-Nearest Neighbors*. Towards Data Science. Retrieve June 26,

2023 from https://towardsdatascience.com/k-nearest-neighbors-94395f445221

Brownlee, J. (2020, October). *Nearest Shrunken Centroids With Python.* Machine Learning

Mastery. Retrieved June 26, 2023 from Nearest Shrunken Centroids With Python -

MachineLearningMastery.com

Clore, J., Cios, K., DeShazo, J., and Strack, B. (2014). Diabetes 130-US hospitals for years

1999-2008 [Data set]. *UCI Machine Learning Repository.*

https://doi.org/10.24432/C5230J

CMS.gov. (2023). Hospital Readmissions Reduction Program (HRRP). Retrieved June 26, 2023

from Hospital Readmissions Reduction Program (HRRP) | CMS

ICD.Codes. (n.d). ICD9CM Chapters. Retrieved June 26, 2023 from https://icd.codes/icd9cm

National Center for Health Statistics. (2021). International Classification of Diseases, Ninth

Revision, Clinical Modification (ICD-9-CM). *Center for Disease Control and*

*Prevention.* Retrieved June 26, 2023 from ICD-9-CM - International Classification of

Diseases, Ninth Revision, Clinical Modification.

Reardon, S. (2015). Preventable Readmissions Cost CMS $17 Billion. *Revcycle Intelligence.*

Retrieved June 26, 2023 from  Preventable Readmissions Cost CMS $17 Billion

Strack, B., DeShazo, J.P., Gennings, C., Olmo, J.L., Ventura, S., Cios, K.J., & Clore, J.N. (2013).

Impact of HbA1c Measurement on Hospital Readmission Rates: Analysis of 70,000

Clinical Database Patient Records. *BioMed Research International,* 2014, Article

781670.  https://doi.org/10.1155/2014/781670

Vadapalli, P. (2022, October). Difference Between Boosting and Bagging. *upGrad.* Retrieved

June 26, 2023 from https://www.upgrad.com/blog/bagging-vs-boosting/