

Clinical Risk Modeling For Re Hospitalization In Diabetes: Insights From Electronic Health Records

Ritesh Sekar^{AVV 1}, Vikas^{B²}, Abhinandhan^{PS³}, Dr. S Nithya Roopa⁴

^{1, 2, 3, 4} Dept of Information Technology

^{1, 2, 3, 4} Kumaraguru College of Technology, Coimbatore, India

Abstract- Hospital readmissions are a major concern in healthcare, impacting patient well-being and increasing medical costs. This study focuses on predicting hospital readmission rates for diabetic patients using machine learning techniques. By analyzing patient demographics, medical history, and hospitalization details, we aim to identify key risk factors contributing to readmission. The study employs various classification models, including Logistic Regression, Decision Trees, Random Forests, and XGBoost, to determine the most effective predictive approach. Our findings indicate that certain patient attributes, such as time spent in the hospital, number of inpatient visits, and medication changes, play a significant role in readmission likelihood.

I. INTRODUCTION

During past decades, hospital readmissions have been the subject of retrospective surveys and prospective trials with a view to their prevention. A hospital readmission is when a discharged patient gets re-admitted to a hospital within a certain period of time. The need for hospital readmission for certain conditions indicates the hospital quality. In other words, it shows inadequate care was provided to the patient at the time of first admission, and thus poses threat to patient life. Additionally, high hospital readmission rates affect the cost of care adversely. Particularly, early hospital readmissions, usually using the benchmark as less than 30 days, have been recognized as a common and costly occurrence among elderly and high-risk patients. It reported that 20% of beneficiaries are readmitted within 30 days at a cost of more than \$26 billion per year. For diabetes patients particularly, American hospitals spent over \$41 billion on diabetic patients in 2011 who got readmitted within 30 days of discharge. Therefore, to improve the quality of care and a reduction in unnecessary health expense, United States congress passed the Hospital Readmission Reduction Program (HRRP). Starting in October 2012, the Centers for Medicare and Medicaid Services (CMS) began financially penalizing hospitals that perform worse than the national average on risk standardized readmission rates for Medicare patients, using 30-day as a threshold. This lends the critical value for understanding the readmission rate within 30 days, and among our interest group, diabetes patients. Therefore, identifying patients at high risk early in

hospitalization will help to reduce the readmission rate, in that hospitals can focus on preparing readmission for patients at high risk to shorten the length of readmission.

Addressing this critical issue involves the intensive data analysis throughout the research process. This study is a secondary analysis using machine learning methods. Our goal of the analysis is to find the determining factors that lead to higher readmission and correspondingly being able to predict which patients will get readmitted. Therefore, we proposed two research questions:

1. What methods can we use to best predict hospital readmission in this dataset?
2. What are the strongest predictors of hospital readmission in diabetic patients?

The rest of the paper is organized as follows: Section 2 will present a brief overview of the past work and identify the research gaps. Section 3 will specify the methodology we used in the study, including the dataset description and the analysis process. Such process includes data processing, explorative analysis, feature engineering, and modeling and evaluation. Results and discussion are presented in Section 3 with respect to each research questions, followed by conclusion and future work in section 5.

II. RELATED WORK

Numerous previous studies have analyzed the risk factors that predict readmission rates using different type of disease. For example, performed analysis to predict of hospital readmission in general without targeting any specific disease.

For the diabetic patients specifically, focused on subsets of diabetic populations and smaller scale datasets. In predicting the readmission rates, some studies focused on demographic and socioeconomic factors that influence readmission rates. For example, focused on age as a factor and found acute and chronic glycemic control influenced readmission risk for patients over the age of 65 using 29,000 patient's data. used the measurement of HbA1c to study the

relationship between the probability of readmission and the primary diagnosis.

Among the recent studies, predicted diabetes with high risk of readmission through modeling multivariate patient medical records using machine learning classifiers such as Naïve Bayes, Bayesian Networks, Random Forest, Adaboost and Neural Networks. To contribute to the implementation of work into the real world, a cost analysis is used to determine the effective cost. Similarly, [Mingle] addressed the previous research gap that no typical performance metrics of machine learning classifiers is documented.

Our study made the contribution to the field by the following folds:

1. Continue to identify and validate the risk factors in predicting the readmission rates. As the previous literature indicate, knowledge of such factors is likely to be useful in developing protocols for better inpatient care.
2. Explore the unidentified machine learning algorithms to improve the accuracy of predicting performance.

III. METHODOLOGY

In this section, we will provide a description of the dataset, the exploratory data analysis, feature engineering, modeling and evaluation.

3.1 Data Set

To explore this problem, we used a secondary dataset from UCI machine learning repository dataset. The dataset includes 101,766 instances, representing 10 years (1999-2008) of clinical care at 130 US hospitals and integrated delivery networks across the Midwest (18 hospitals), Northeast (58), South (28), and West (16). Most of the hospitals (78) have bed size between 100 and 499, 38 hospitals have bed size less than 100, and bed size of 14 hospitals is greater than 500. The features collected in the dataset are related to patient's demographic information such as race, gender, age, weight; the information related to their hospital diagnosis and treatment, such as num_lab_procedures, num_medications, num_outpatient, diagnosis and medication prescription. The dataset is just an extracted subset set of Health fact dataset. Given this is an open dataset that include the longitudinal and cross-sectional data, and with relatively complete attributes (55 attributes), and released in the recent year (2014), we chose the dataset for exploring the questions.

3.2 Exploratory Analysis

Prior to performing any analysis, we conducted exploratory analysis to preview the data type, attributes, and overall patterns of the data. We are interested in the class label "Readmitted" (See Fig 1), so we checked the distribution of the readmitted, and several categorical variables. For the numerical variables, we used scatter plot to demonstrate the relationship between numerical variables and their distributions (See Fig 2).



Figure 1: Bar plot for the class label.



Figure 2: Scatter plot for numeric features.

3.3 Data Pre-Processing

After the exploratory analysis, we found several challenges lies in the original dataset, and thus some data wrangling tasks such as data cleaning, dealing with missing values, creating new variables, and data transformation needs to be addressed before modeling. The tools used here are Python Packages including Numpy, Pandas, Matplotlib and Seaborn. Specifically, we conducted several pre-data processing processes:

3.3.1 Dealing with Missing Data

We discovered a large number of missing values coded as "?" across nominal variables. As Table 1 shows, this

dataset has 8 variables which contain missing values. Since weight, medical_specialty, and payer_code contains over 35% values, and also because of the irrelevancy toward our study, we decide to drop all of them. Race only includes 2.23% missing values, so we only drop the missing values and keep the rest. Primary (diag_1), secondary (diag_2) and additional (diag_3) diagnoses each has less than 2% missing values, but compared to the total number of instances, we still need to clean them. Technically, our goal is to maintain the most information of the dataset, especially the diagnosis is an important variable related to the diabetes patients. Therefore, we adopted a strategy to drop the missing values when all three diagnosis were missing. We then only drop 3 unknown and invalid instances in our dataset.

Table1: Variables with missing values.

column	count_missing	percent_missing
weight	98569	96.86
Medical_specialty	49949	49.08
Payer_code	40256	39.56
race	2273	2.23
Diag_3	1423	1.4
Diag_2	358	0.35
Diag_1	21	0.02
gender	3	.00003

3.3.2 Dropping Attributes

After a quick view of the current dataset, we found some patients died during the hospital admission who do not have any probability of being readmitted, so we removed those tuples, as the discharge_disposition_id=11. We also drop two variables (drugs named citoglipton and examide) in which all records have the exactly same value. By noticing that two variables called encounter_id and patient_nbr has no relevance with the class label readmission, so we also drop those two variables.

3.3.3 Creation of New Features

1. patient_service: We created a new feature called patient_service, which measures the total number of hospital/clinician services a patient used in the past year. This feature is the sum of original variables for number of inpatient visits, emergency room visits, and outpatient visits. We did not apply weighting for these three variables. The reason for the creation is to lower the dimension of our data and try to make the dataset simpler.

2. med_change: The dataset contains 23 medications of the medicine use for a patient during the stay in hospital. Each of the tuple records when a change was made in this medication or not during the current stay as No-for no medication, Up- for increasing the dose, Down-for decreasing the dose and Steady-for keeping the current dose. Instead of counting changes for each medication, we decide to combine them and count changes for all of them. We define No and Steady as no change, while up and down for change. Doing this step will simplify the model and we can try to find out if the readmission is related with medication changes.

3. num_med: Not only medication changes can be related with readmission, the total number of medications used can also be a key feature, due to the fact that the number of the medicine reflected the severity of certain disease. And thus we created a variable called num_med in order to store the total number of medications a patient used during the stay of hospital.

Recoding Existing Variables

1. Recode diagnoses: in the dataset, we have three features called diag_1, diag_2, and diag_3. All of them were coded according to ICD-9 codes, namely, International Statistical Classification of Diseases and Related Health Problems. The ICD is originally designed as a health care classification system, providing a system of diagnostic code for classifying diseases, including nuanced classifications of a wide variety of signs, symptoms, abnormal findings, complaints, social circumstances, and external causes of injury or disease. This system is designed to map health conditions to corresponding generic categories together with specific variations, assigning for these a designated code, up to six characters long. Thus, major categories are designed to include a set of similar diseases [15]. First, we replaced the unknown value “?” into 1. We then recode the diagnoses into Circulatory-1, Respiratory-2, Digestive-3, Diabetes-4, Injury-5, Musculoskeletal-6, Genitourinary-7, Neoplasms-8, and Others-0. If ICD code is between 390 and 460, or it equals to 785, it belongs to category 1 (circulatory). If ICD code is between 460 and 520 or it equals to 786, it belongs to category 2 (respiratory). If ICD code is between 520 and 580 or it equals to 787, it belongs to category 3 (digestive). If ICD code equals to 250, it belongs to category 4 (diabetes). If ICD code is between 800 and 1000, it belongs to category 5 (injury). If ICD code is between 710 and 740, it belongs to category 6 (musculoskeletal). If ICD code is between 580 and 630 or it equals to 788, it belongs to category 7 (genitourinary). If ICD code is between 140 and 240, it belongs to category 8 (neoplasms). Others belong to category 0

(others). Appendix A shows the details of the recoding process.

2. Recode age: Since we intend to examine how age is related with readmission, we record age with 10 categories into the numerical variables by taking the mean of each age category. For example, if the patient's age category is 10-20 years old, then we use the 15 years old to represent the whole range.
3. Recode readmission: For clinician/hospital and the current interest of the study, we only focus on those whose readmission is less than 30 days. We recode the readmission into two categories, for those who are <30 days, we recode them into 1. For those who >30 days and no need of readmission, we recode them into 0.
4. Recode other variables: For three variables which related with admission type, discharge disposition and admission source, we decided to encode the dummy variables for these categories. For variable "change", we recoded change into 1 and no change into 0. For gender, we recoded male into 1 and female into 0. For diabetes_Med, we recoded yes into 1 and no into 0. For race, we recoded the categorical variables into dummy variables: Caucasian-1, African American-2, Hispanic-3, Asian-4, and others-0. For A1Cresult, we recoded >7 and >8 into 1, Norm into 0, and None into 99. For max_glu_serum, we used the similar method, namely, we recoded >200 and >300 into 1, Norm into 0, and None into 99.

3.4 Feature Engineering

3.4.1 Data Type Conversion

For nominal features, we converted them into object type, in order for the later numerical variables processing.

3.4.2 Log Transformation, Standardization, and Correlation

For numeric features, the scatter plot of the distributions as Figure 1 indicated most of numerical are highly skewed and had high kurtosis. Using the threshold=+/-1 as skewness for normal distribution, if skewness is less than -1 or greater than 1, the distribution is highly skewed. If skewness is between -1 and -0.5 or between 0.5 and 1, the distribution is moderately skewed. If skewness is between -0.5 and 0.5, the distribution is approximately symmetric. As the standard for kurtosis, the threshold = 3 is for normal distribution. Therefore, we used log transformation to normalize the numerical variables to make sure numeric variables had a Gaussian-like or normal distribution.

Since the numerical variable are not using the same scale so we rescale our data using the standardization methods with the following formula:

$$\text{New value} = \frac{\text{Value} - \text{Mean}(\text{Values})}{\text{Standard Deviation}(\text{Values})}$$

After all data are standardized, we checked the correlation between the variables using a heat map to find top 15 correlated variables as Fig. 3 shows. There is not too much correlation between the variables and the correlation listed are self-explainable.

3.4.3 Outliers.

For detecting and processing the outliers, we used the coverage rule for normal distribution to deal with outliers. As Fig. 4 shows, the remaining 0.3% of the data are treated as outliers for this project. And thus, we removed the outliers.

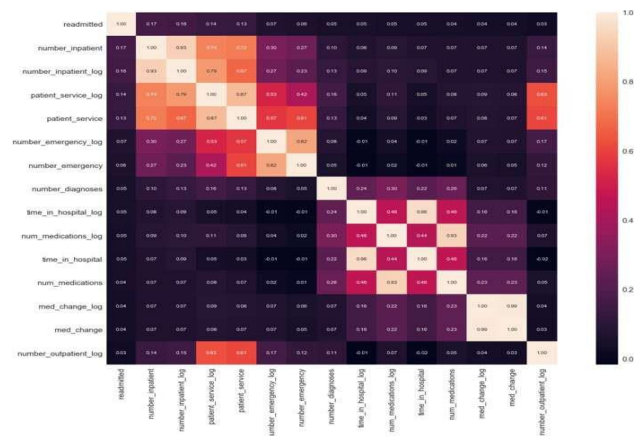


Figure 3: Heat map of top 15 correlated variables.

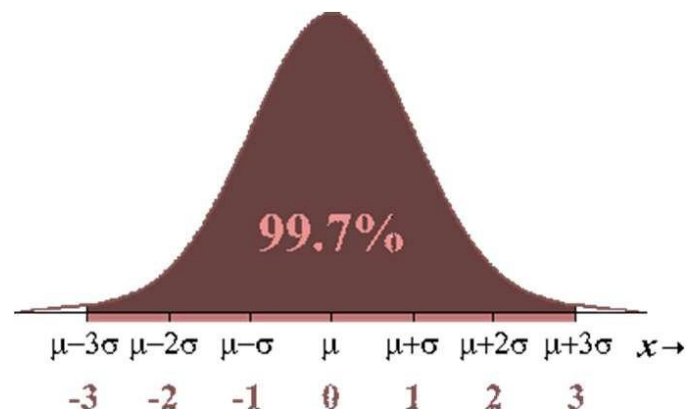


Figure 4: 99.7% of the observations fall within 3 standard deviations of the mean. [8]

3.4.4 Class Imbalance

Before modeling, we checked class label balance after running analysis to check the class balance. As a result, we see there are 79512 tuples belong to class 0, which is >30-day readmission and no admission needed. Only 9607 patients belong to <30-day readmission type. The proportion is above 8:1, with the proportion threshold 10-20%, our data is highly imbalanced which will lead to better accuracy after modeling. We used the confusion matrix to ensure the class label is highly balanced as Fig. 5 shows. Our benchmark model logistic regression showed that 89% accuracy with precision and recall rate equals to zero. We used an over-sampling technique (SMOTE) to our data is balanced by oversampling our underrepresented class of readmissions. Fig. 6 explains how over-sampling and under-sampling works. After using SMOTE, we will get 79512 patients with category 0, and 79512 patients with category 1 too. Fig. 5 also shows the confusion matrix before and after data balancing.

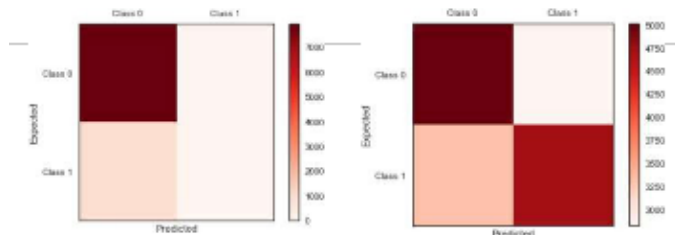


Figure 5: Confusion matrix before data balancing (left) and confusion matrix after data balancing (right).

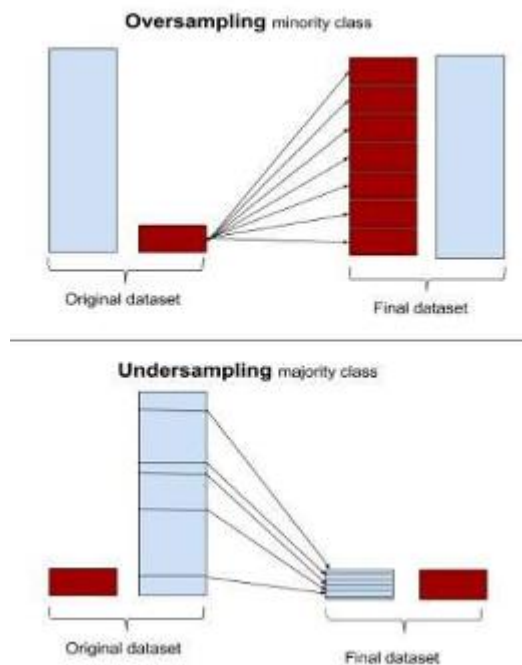


Figure 6: Explanation of over-sampling and under-sampling.

IV. EXPERIMENT

Our aim for experiment of modeling is to identify the factors for high-risk diabetic patients, which was posed as the problem of classifying of whether a patient would be readmitted within 30 days of being discharge or after 30 days of being discharged or never readmitted. Therefore, different classification algorithms are adopted for experimentation and find the best method that achieve the highest accuracy. We adopted and compared the four classification algorithms. Prior to training the classification algorithms, we randomly split our dataset into two distinct sets - the training and the test set. The training and test set consisted of 90% and 10% of the data accordingly. The parameters of each algorithm were chosen based on the classification performance evaluated by 10-fold cross-validation on the training set. The performance of all algorithms was evaluated on the test set. The methods we implemented include:

4.1 Logistic Regression

Logistic regression is used as a benchmark model for our analysis. Since we assume that our data can be modeled as a log likelihood of outcome for the binary class label readmission, logistic regression can help us to understand the relative impact and significance of each attributes. We test this model by using 90% training and 10% testing data and 10-fold cross-validation. We achieved a cross validation score: 61.29% and test set score 61.35%. By looking into the confusion matrix, we can calculate several measures of accuracy:

Accuracy is 0.61
Precision is 0.63
Recall is 0.59
AUC is 0.61

4.2 Decision Trees

Decision trees is a popular tree-based model that is easily to interpret the logic for splitting. Decision trees classify the data by sorting them down the tree from the root to some leaf node, with the leaf node providing the classification to the data. Due to the interactions between variables inherently, we removed the interaction variables from the feature set we did for logistic regression. Similarly, we did 10- fold cross validation score for decision trees too. The score equals to 88.97% and the dev set score is 89.43%, so decision trees look good for this dataset. After checking the score, we analyzed

the confusion matrix for decision trees for both entropy and gini methods. As a result, both of them yielded the same results of measurements:

The result turned out that decision trees performed better than logistic regression based on its accuracy. The following graph showed the splitting process of the tree node. We visualized the trees in first two levels (as Fig 7).

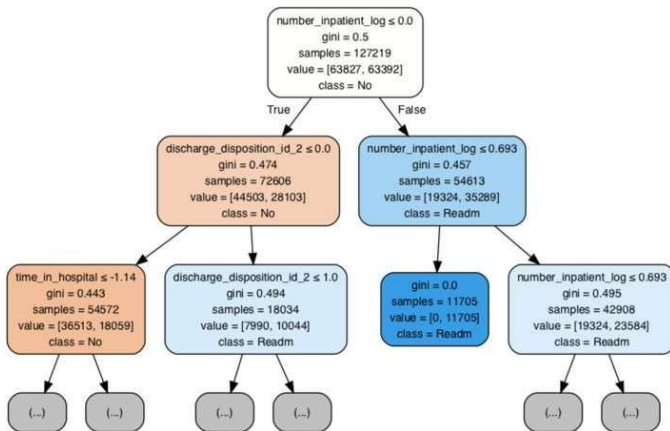


Figure 7: Decision trees for Gini index.

From the graph, it indicated inpatient visits is the first feature this decision tree used in deciding whether a patient will get readmitted.

4.3 Random Forests

Random Forest is composed of a set of decision trees. Each decision tree acts as a weak classifier and pooling the responses from multiple decision trees leads to a strong classifier. Each decision tree is trained independently and determines the class of an input by evaluating a series of greedily learned binary questions. The random forest consisting of 10 trees, with the max_depth of as 25 nodes was used, as it was found to be optimal from the experiment with varying number of trees and depth in the forest. After implementing Random Forest, we achieved similar results of measurements for using gini and entropy methods. Random Forest showed better results than decision tree as regards to prediction accuracy.

Accuracy is 0.92
Precision is 0.98
Recall is 0.87
AUC is 0.92

4.4 Model Improvement

After running the random forest, we decided to use a boosting method by the relatively new algorithm XGboost for model improvement. Boosting is an ensemble method that create a strong classifier based on weak classifiers, according to how correlated are the learners to the actual target variable. The errors of the previous model are corrected by the next predictor, by adding models on top of each other iteratively until the training data is accurately predicted or a maximum number of models are added.

EXtreme Gradient Boosting (XGBoost) is an ensemble machine learning method, which has been very popular since its introduction in 2014. XGBoost is a scalable and accurate implementation of gradient boosting machines and it has proven to push the limits of computing power for boosted trees algorithms as it was built and developed for the sole purpose of model performance and computational speed. XGBoost works for generic loss functions, as shown below and it has more customizable parameters.

$$Obj = \sum_{i=1}^n l(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k)$$

Training loss
Complexity of the Trees

We applied and tuned the algorithm for better performance. We tuned the following three parameters

1. eta: learning rate to prevents overfitting (eta=0.01, 0.02,0.05)
2. max_depth: the max depth of the tree (max_depth=3,4,5,6,7,8,9)
3. cols_sample: the percentage of features can be chosen (cols_sample=0.6,0.7,0.8,0.9,1.0).

We tuned the three parameters one by one and iterate the values to find the least test error and highest accuracy. The best iteration we found is with accuracy 0.94, precision 1.0, recall 0.88 and AUC is 0.94.

4.5 Evaluation

In this section, we will discuss the evaluation of classifier performance and answer the second question of identifying the most import factors.

4.5.1 Classifier Comparison

Each algorithm was evaluated using a 10-fold stratified cross-validation. Cross-validation is an evaluation technique where the dataset is randomly but evenly distributed into a number of folds. Stratified cross-validation attempts to

preserve the class distribution between folds so that each fold is representative of the full dataset. The learning algorithm is trained on 9-fold and tested on 1-fold of dataset. Repeating the cross-validation process ensure that particular random initialization does not bias the overall result.

All of our algorithms are evaluated using the area-under-the-curve (AUC), which is equivalent to the c-statistic in the binary classification scenario. AUC - ROC curve is a performance measurement for classification problem at various thresholds settings. ROC is a probability curve and AUC represent degree or measure of separability. It tells how much model is capable of distinguishing between classes. In our study, the AUC is the probability that a positive instance “<30 as 1” ranks higher than a negative one “0”. The higher the AUC, the better the model is at predicting 0s as 0s and 1s as 1s. Previous research in readmission has achieved AUCs between 0.5 and 0.7.

In comparing of four models, Table 2 shows the XGBoost works the best for predicting the admission rate, which achieves the highest accuracy as 0.94, with 0.61 on AUC, and the second best model is random forest, which achieves 0.92 accuracy and 0.94 on AUC. Also, the model comparison is shown in the Fig 8 for overall performance.

Table2: Comparison between different algorithms

Classifier	Accuracy	Precision	Recall	AUC
Logistic_Regression	0.61	0.63	0.59	0.61
Decision Tree	0.89	0.92	0.87	0.89
Random Forest	0.92	0.97	0.87	0.92
XGBoost	0.94	1.0	0.88	0.94

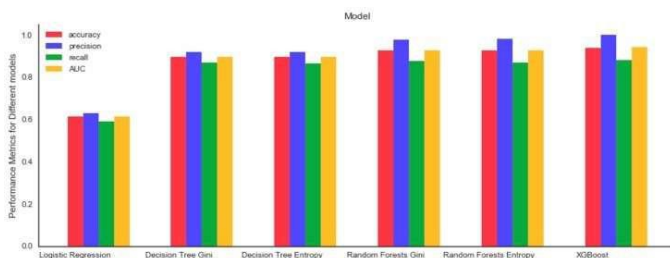


Figure 8: Comparison between models.

4.5.2 Most Important Predictors

For the second question what the strong predictors are contributing to predicting readmission, different algorithms provided different results. Specifically, Fig. 9 illustrated showed the most important variables after the classification for decision tree. We plotted those features whose importance is bigger than 0.01. The most important variables are number_inpatient and time_in hospital, and the

discharge_disposition_id_2, number_procedures, and num_medications are among the top 5 strongest predictors.

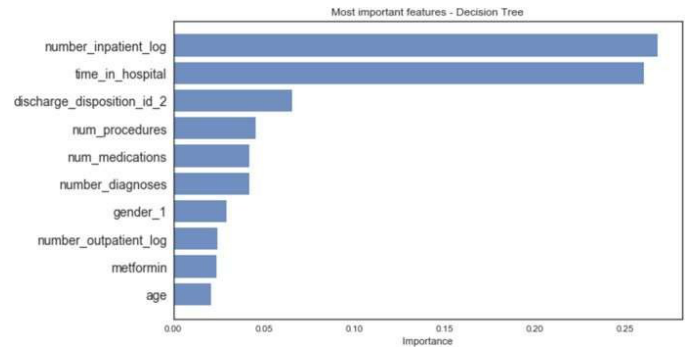


Figure 9: Most important features for decision tree model.

Fig. 10 showed the important features for random forests, which are different from the decision trees, with number_inpatient, time_in_hospital, number_diagnosis, discharge_id_2 and metformin are among the top 5 important predictors.

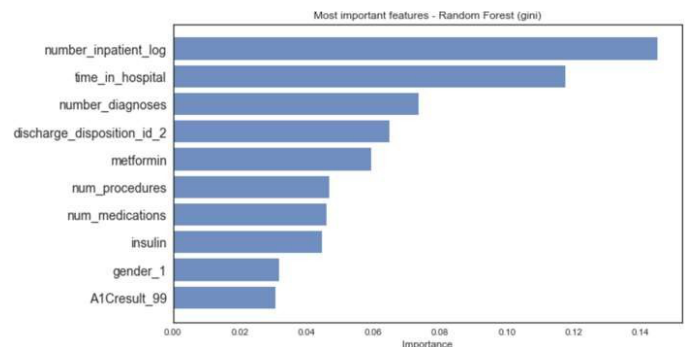


Figure 10: Most important features for random forests.

Fig. 11 indicated the important features for XGBoost which are slightly different than previous with number_medications, time_in_hospital, age, number_procedures, num_diagnosis are among the top 5 important predictors. The results are quite interesting.

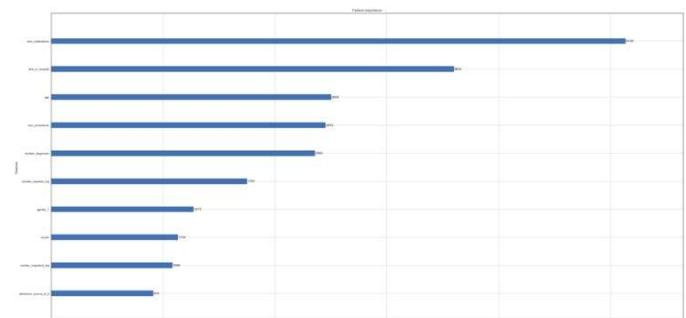


Figure 11: Most important features for XGBoost.

V. CONCLUSION

In this work we adopted machine learning methods to identify high risk patients and evaluated different machine learning algorithms. Compared to the previous analysis, our study achieved high accuracy due to the sophisticated pre-processing procedure. The XGBoost method is reported to be the best method for prediction of the readmission rate for diabetes patients.

We identified the most important factors as the time_in_hospital and number of inpatient, number of diagnosis, which appears to associate with the severity of the disease. Further studies could conduct more exploration when analyzing these factors individually.

REFERENCES

- [1] Benbassat, J. Taragin, M. 2000. Hospital readmissions as a measure of quality of health care advantages and limitations. *Arch Intern Med.* 160(8):1074–1081.
- [2] Leppin, A.L., Gionfriddo, M.R., Kessler, M., Brito, J.P., Mair, F.S., Gall K., Wang, Z., Erwin, P.J., Sylvester, T., Boehmer, K. and Ting, H.H., 2014. Preventing 30-day hospital readmissions: a systematic review and meta-analysis of randomized trials. *JAMA internal medicine*, 174 (7), 1095-1107.
- [3] Hines, A.L., Barrett, M.L., Jiang, H.J. and Steiner, C.A., 2006. Conditions with the largest number of adult hospital readmissions by payer, *Statistical Brief.* 172 (2011).
- [4] <https://medium.com/berkeleyischool/how-to-use-machine-learning-to-predict-hospital-readmissions-part-1-bd137cbdba07>
- [5] Salerno, A.M., Horwitz, L.I., Kwon, J.Y., Herrin, J., Grady, J.N., Lin, J.S. and Bernheim, S.M., 2017. Trends in readmission rates for safety net hospitals and non-safety net hospitals in the era of the US Hospital Readmission Reduction Program: a retrospective time series analysis using Medicare administrative claims data from 2008 to 2015. *BMJ open*, 7(7).
- [6] Dungan, K. M. The effect of diabetes on hospital readmissions.,2012. *Journal of diabetes science and technology*, 6(5), 1045–1052.
- [7] Eby, E., Hardwick, C., Yu, M., Gelwicks, S., Deschamps, K., Xie, J. and George, T., 2015. Predictors of 30-day hospital readmission in patients with type 2 diabetes: a retrospective, case-control, database study. *Current medical research and opinion*, 31(1), 107-114.
- [8] Howell, S., Coory, M., Martin, J. and Duckett, S., 2009. Using routine inpatient data to identify patients at risk of hospital readmission. *BMC Health Services Research*, 9(1), 96.
- [9] Jiang, H.J., Stryer, D., Friedman, B. and Andrews, R., 2003. Multiple hospitalizations for patients with diabetes. *Diabetes*