

# PYTHON FOR DATA ANALYSIS

Diabetes 130-US hospitals for years 1999-2008

---

Oscar Pastural - Richard Goudelin - Capucine Boudin



# DEFINITION OF THE DATASET

The dataset provided consists of hospital records from 130 US hospitals and integrated delivery networks, covering a period of ten years from 1999 to 2008. The records pertain to patients who were diagnosed with diabetes and include information on laboratory tests, medications, and hospital stays of up to 14 days..



## Content

Ten years (1999-2008) of clinical care at 130 US hospital

Patients diagnosed with diabetes, who underwent laboratory, medications, and stayed up to 14 days.



## Shape

- Feature type : Multivariate
- Instances : 101766
- Features : 47

# **SUMMARY**

**I. Problem Definition**

**II. Data Preprocessing**

**1. Data Cleaning**

**2. Mapping the data**

**3. Encoding the data**

**4. Re-Writing the data**

**III. Exploring Plots**

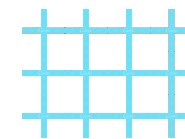
**1. Parameters Overview**

**2. Correlations with the readmission distribution in time**

**3. Correlations with the readmission status**

**IV. Studying Machine Learning models**

# PROBLEM DEFINITION



**How can we predict the readmission status of a patient ?**

How does the problem fit in our context ?

- The total annual cost of diabetes in 2022 is \$412.9 billion
- 9.3% of the population in the United States have diabetes

Predicting the readmission status involves identifying the key parameters that play a significant role and reduce complications associated with diabetes.



# DATA PREPROCESSING

## Data cleaning

- Removing :
  - Features with a high percentage of missing values ('weight', 'medical\_specialty', 'payer\_code', 'max\_glu\_serum', 'A1Cresult')
  - Redundancies (29423 patient\_nbr)
  - Useless columns ('encounter\_id', 'patient\_nbr', 'payer\_code')
  - 3 missing values in the gender column
  - Drugs with zero variance → feature selection based on a variance threshold at 0.05. (11 features are kept, 15 are deleted)

## Encoding data

- Transformed non-numerical columns into a data\_encoded dataframe :
  - ['gender', 'admission\_type\_id', 'discharge\_disposition\_id', 'admission\_source\_id', 'diag\_1', 'diag\_2', 'diag\_3', 'metformin', 'repaglinide', 'glimepiride', 'glipizide', 'glyburide', 'pioglitazone', 'rosiglitazone', 'insulin', 'change', 'diabetesMed']

## Mapping the data

- Variables mapped by following the IDS mapping provided with the dataset :
  - admission\_type\_id
  - admission\_source\_id
  - discharge\_disposition\_id
- Variables mapped by information found along the dataset :
  - diag\_1, diag\_2 and diag\_3

## Re-Writing the data

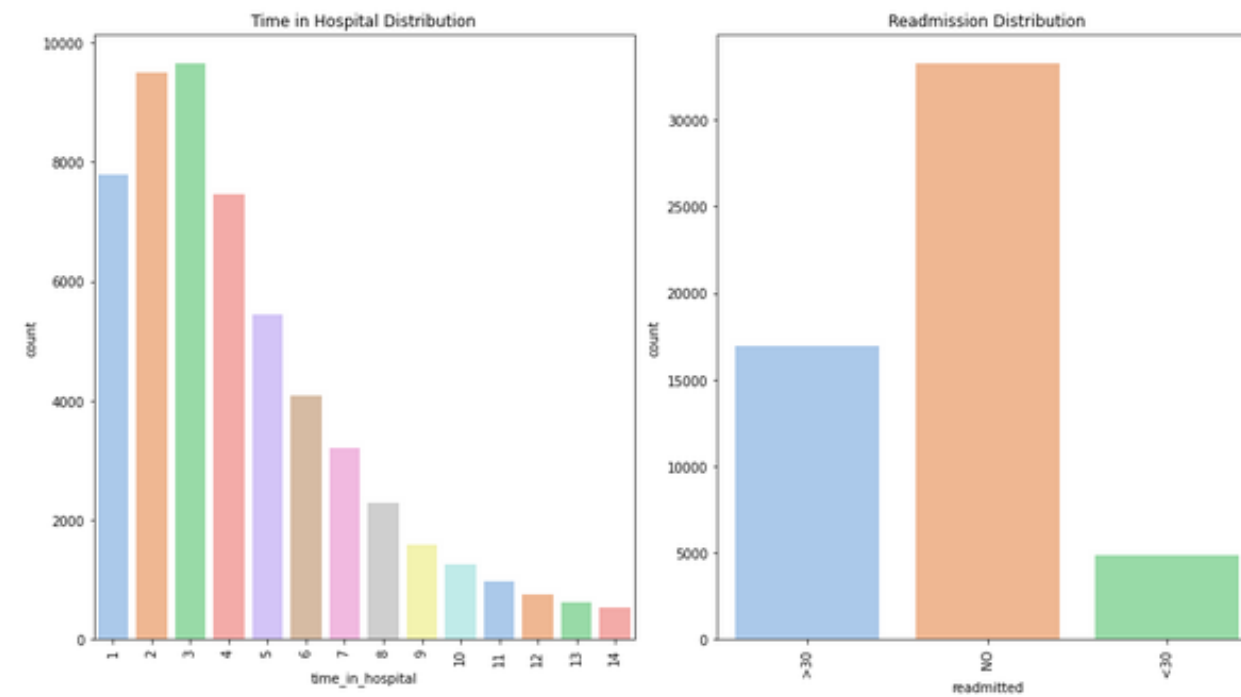
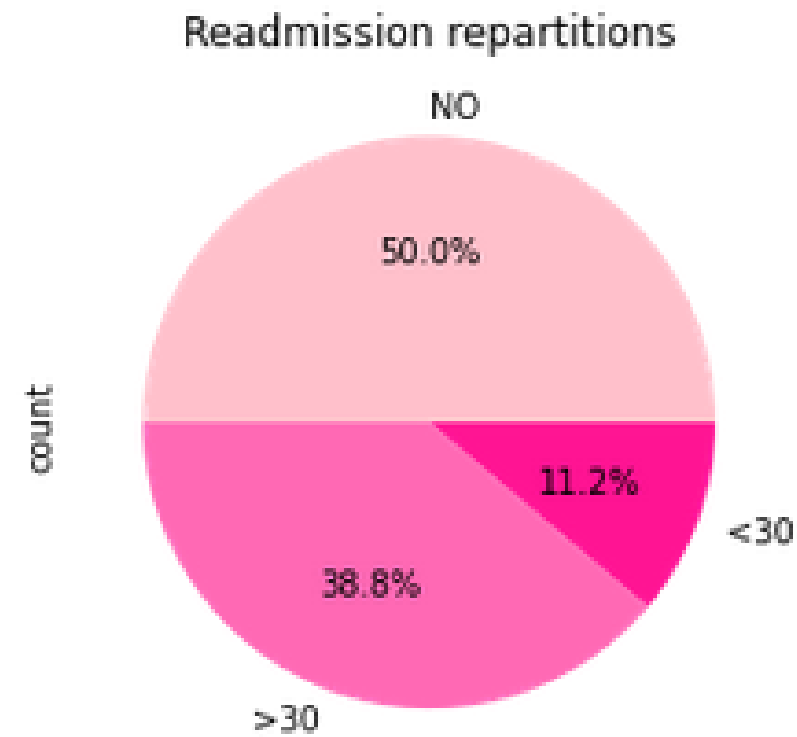
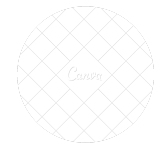
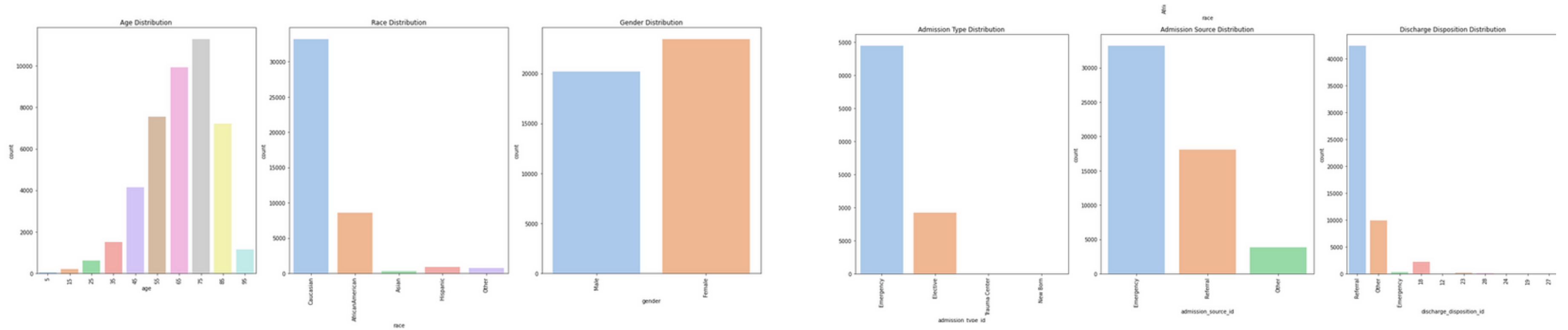
- The age was given as classes such as [0–10), [10–20), [20–30) etc...
- We converted this column to a numeric feature as the following :

Figure 3



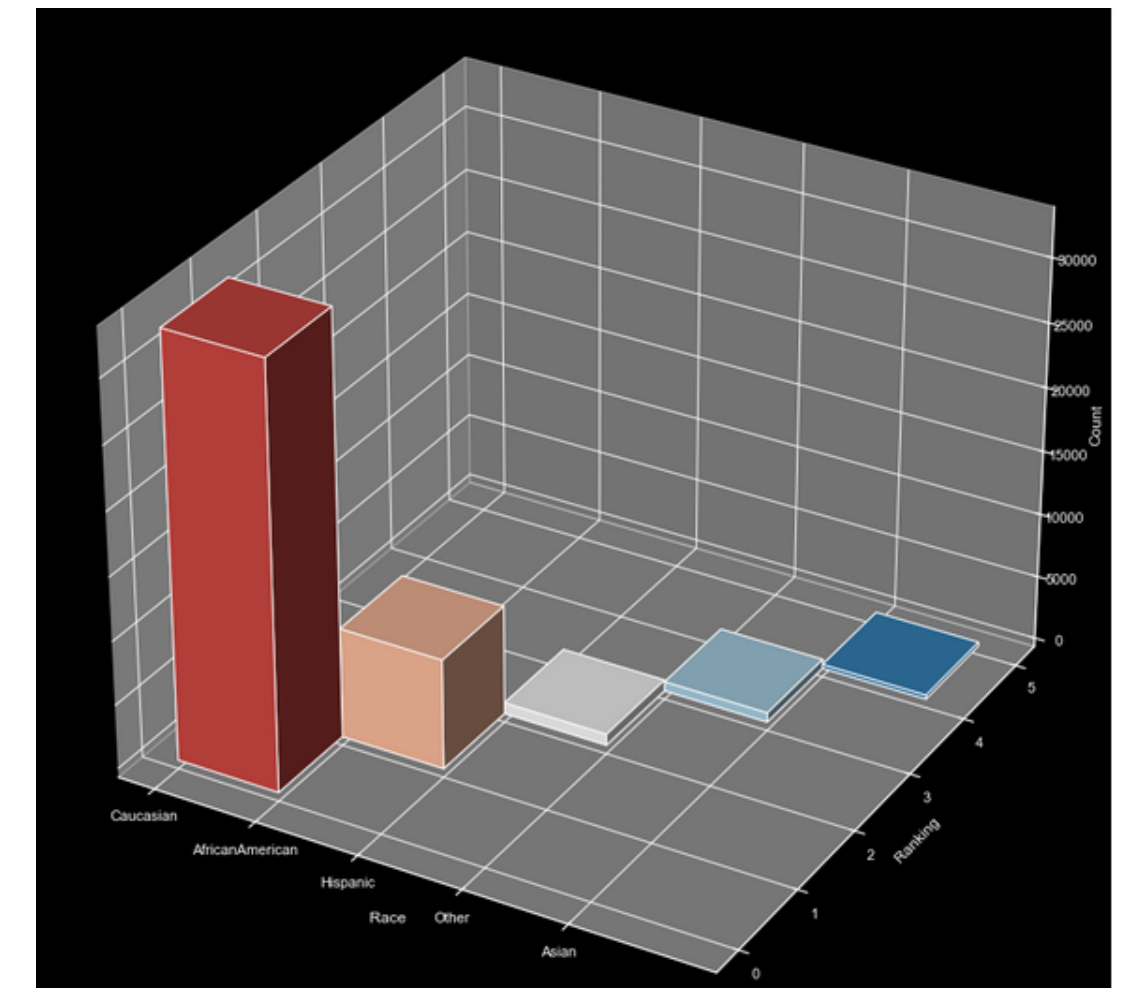
Predicted (adjusted for covariates) readmission rates by the primary diagnosis and HbA1c measurement. Readmission rates were predicted on the reference values of other predictors and the mean value of time in hospital (Table 3). The error bars represent the 95% confidence intervals for the predicted values. The following abbreviations are used for particular icd9 codes: “circulatory” for icd9: 390–459, 785, “digestive” for icd9: 520–579, 787, “genitourinary” for icd9: 580–629, 788, “diabetes” for icd9: 250.xx, “injury” for icd9: 800–999, “musculoskeletal” for icd9: 710–739, “neoplasms” for icd9: 140–239, “respiratory” for icd9: 460–519, 786, and “other” for otherwise.

```
replaceDict = {'[0-10)' : 5,  
               '[10-20)' : 15,  
               '[20-30)' : 25,  
               '[30-40)' : 35,  
               '[40-50)' : 45,  
               '[50-60)' : 55,  
               '[60-70)' : 65,  
               '[70-80)' : 75,  
               '[80-90)' : 85,  
               '[90-100)' : 95}  
  
data['age'] = data['age'].apply(lambda x : replaceDict[x])
```

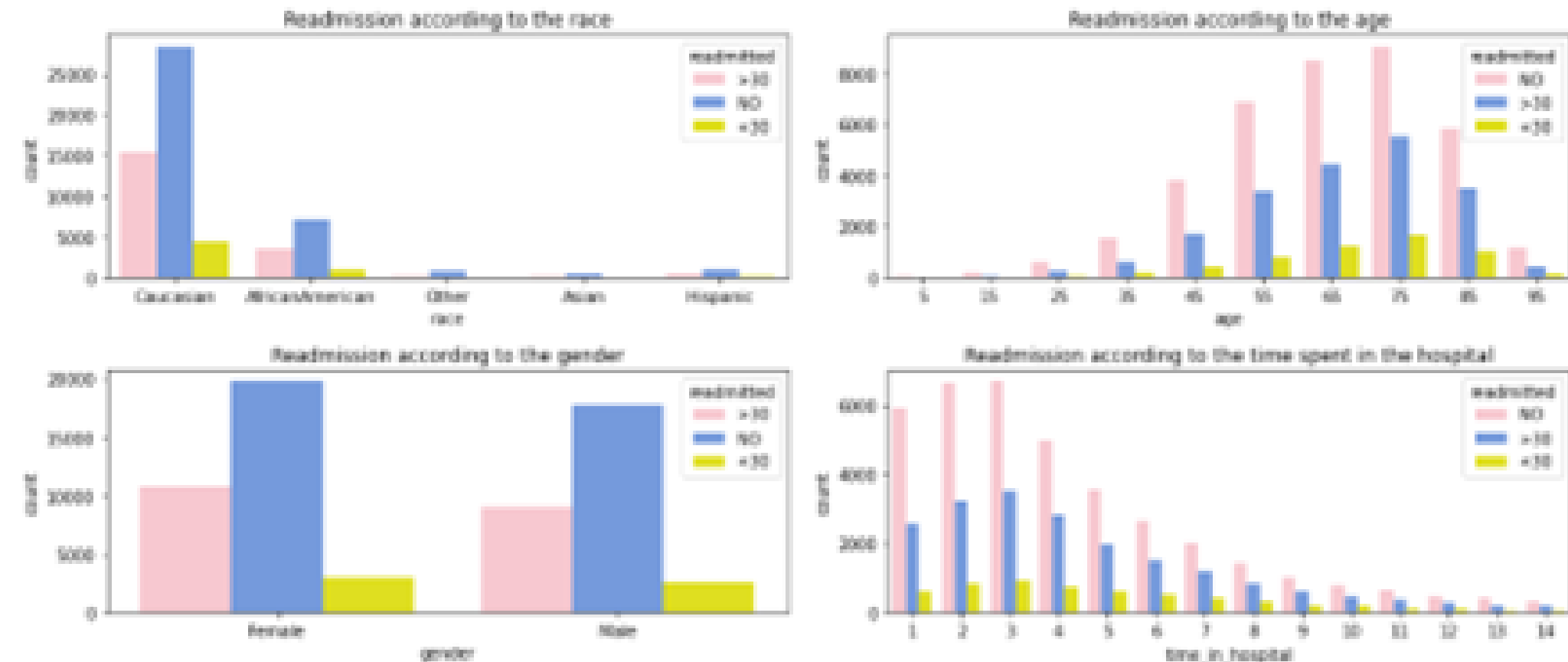
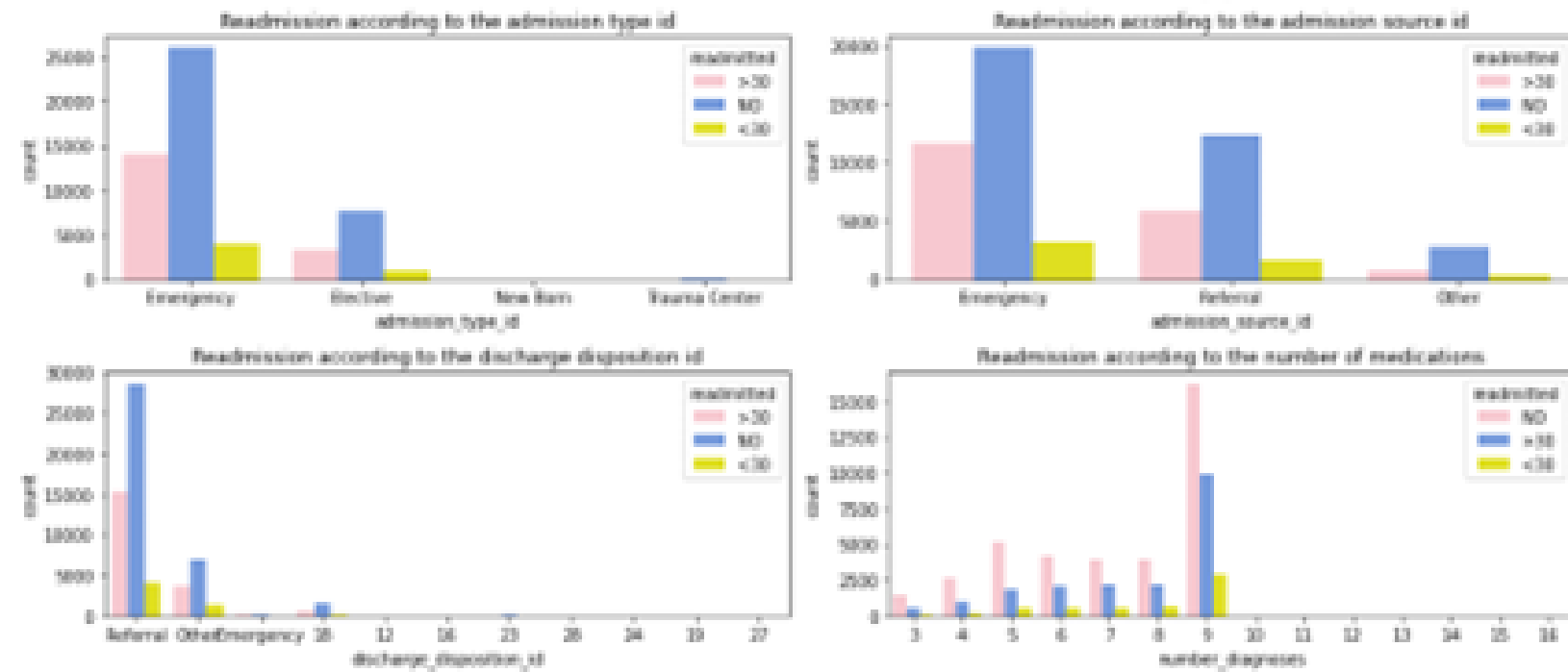


# EXPLORING PLOTS

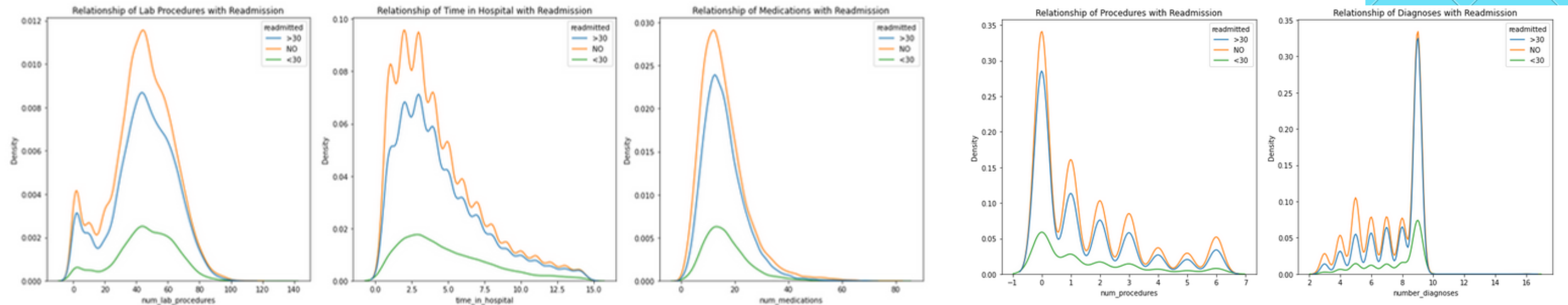
Parameters overview

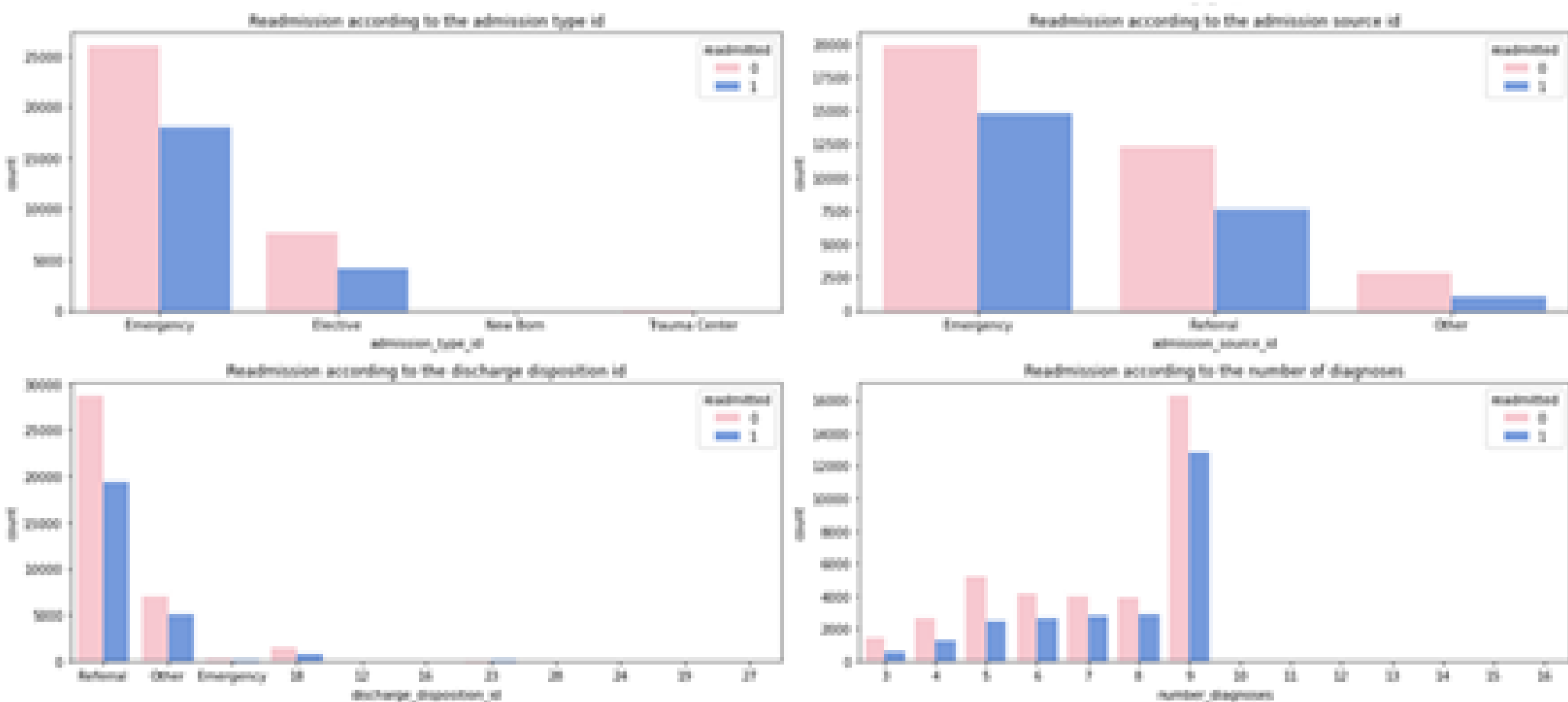


- We now consider the parameters in correlation with the readmission within more or less than 30 days or no readmission at all



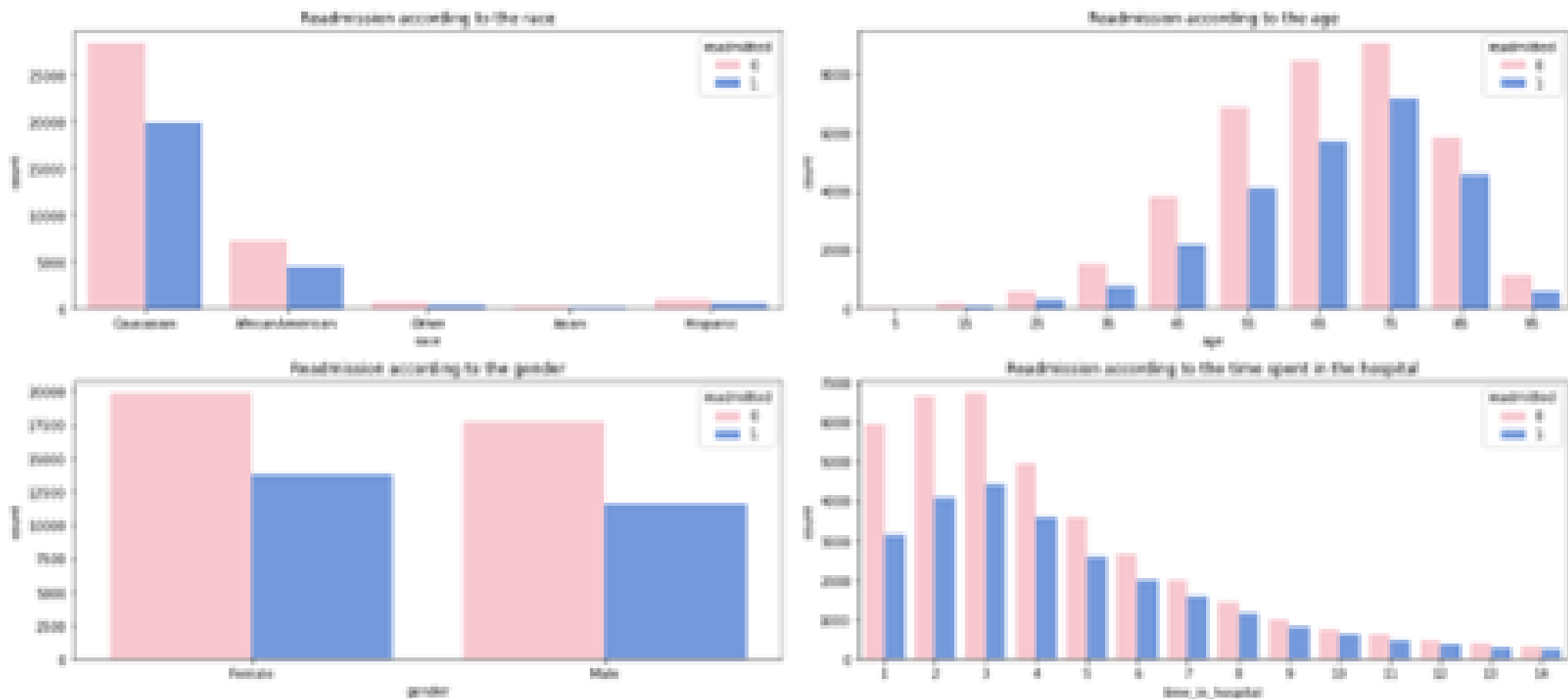
# CORRELATIONS WITH THE READMISSION DISTRIBUTION IN TIME





•We now consider the parameters in correlation with the readmission status. The patient is now either readmitted or not readmitted so that the problem is more balanced

# CORRELATIONS WITH THE READMISSION STATUS







# STUDYING MACHINE LEARNING MODELS

Thus this far, we have fathomed what machine learning models we could apply to our dataset: RandomForest, DecisionTree, Linear Regression and DecisionTree classifiers. We will try to predict our values and evaluate our models.



## Linear Regression

- Simplicity and Interpretability: It models the relationship between a dependent variable and one or more independent variables.
- Baseline Model: Can be used as a starting point for regression problems to establish a baseline performance



## Random Forest

- Non-Linearity: Random Forest can handle non-linear data effectively. It combines multiple decision trees to improve the model's accuracy and prevent overfitting.
- Feature Importance: Provides insights into which features are most important in making predictions.
- Versatility: Works well for both classification and regression tasks and can handle large datasets with higher dimensionality

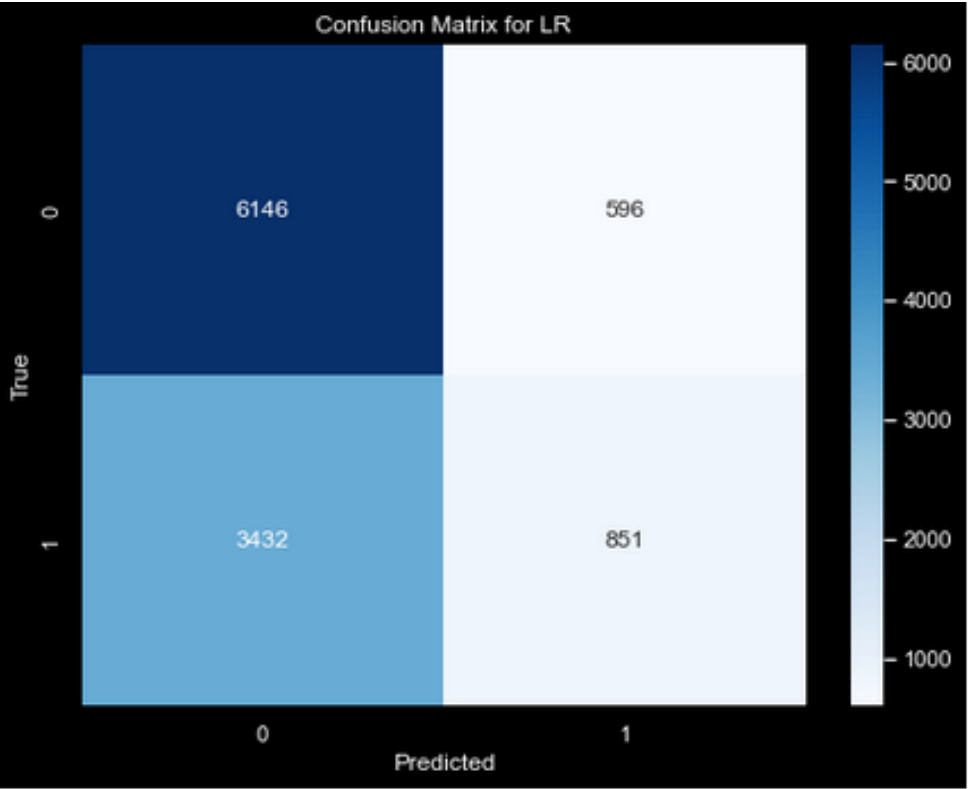


## Decision Tree

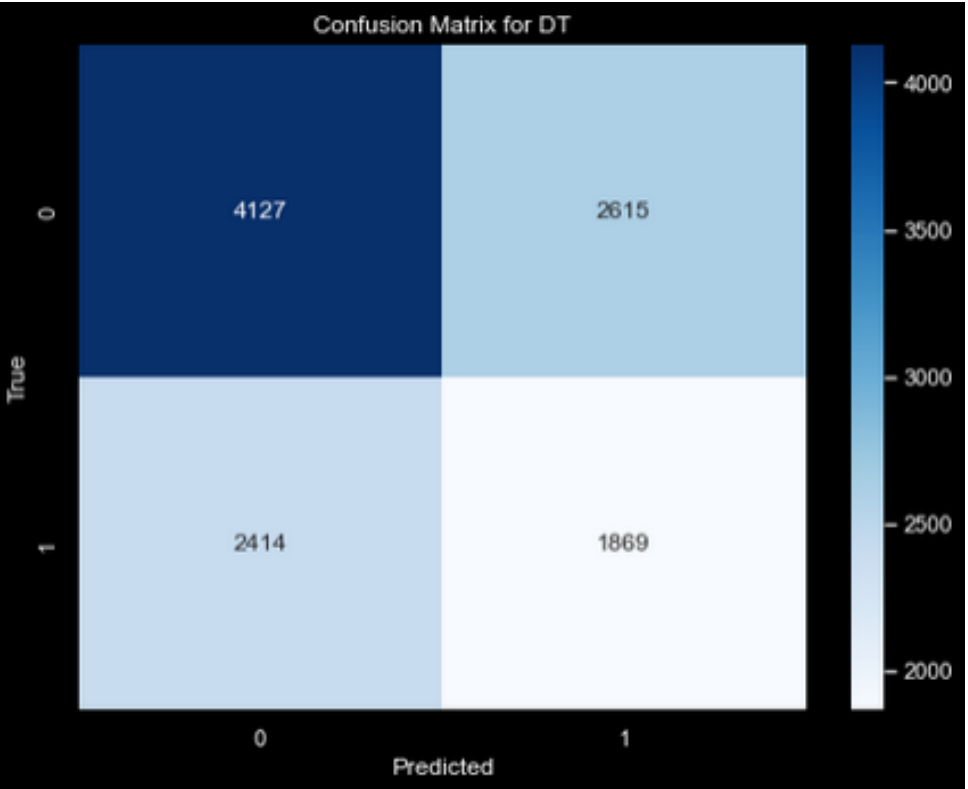
- Easy to Interpret: Decision Trees can be visualized easily, making them highly interpretable.
- Non-Parametric: Does not require much data preprocessing, such as normalization or scaling. They can handle both numerical and categorical data.

# CONFUSION MATRIX

## Linear Regression



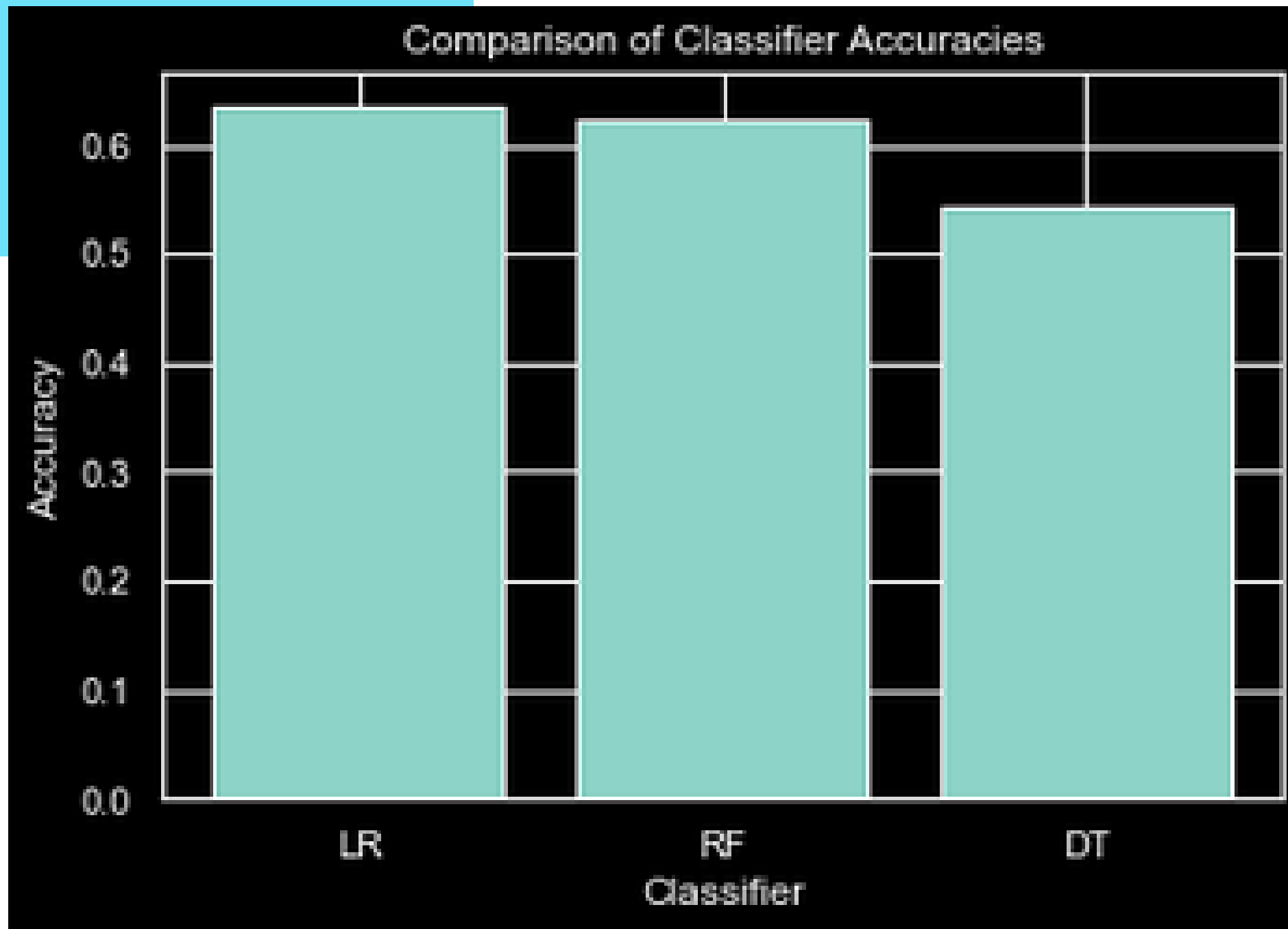
## Doctor appointment



## Random Forest



# COMPARISON OF CLASSIFIER ACCURACIES

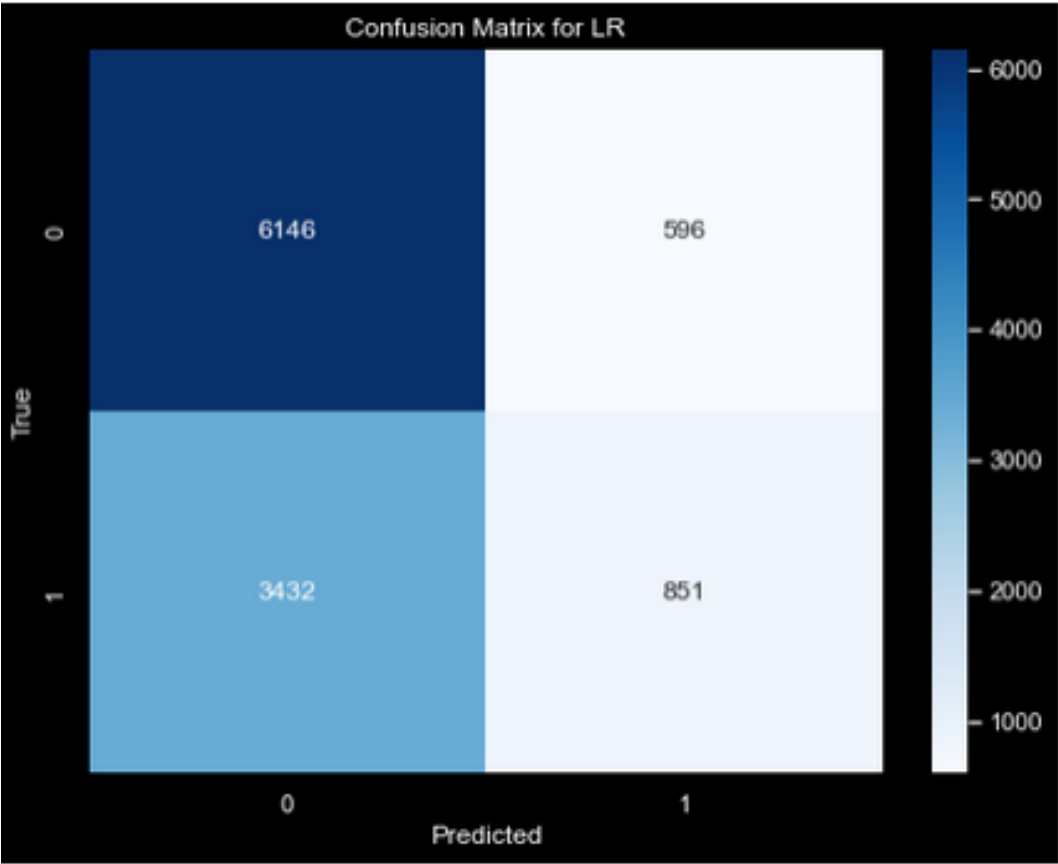


## Observation

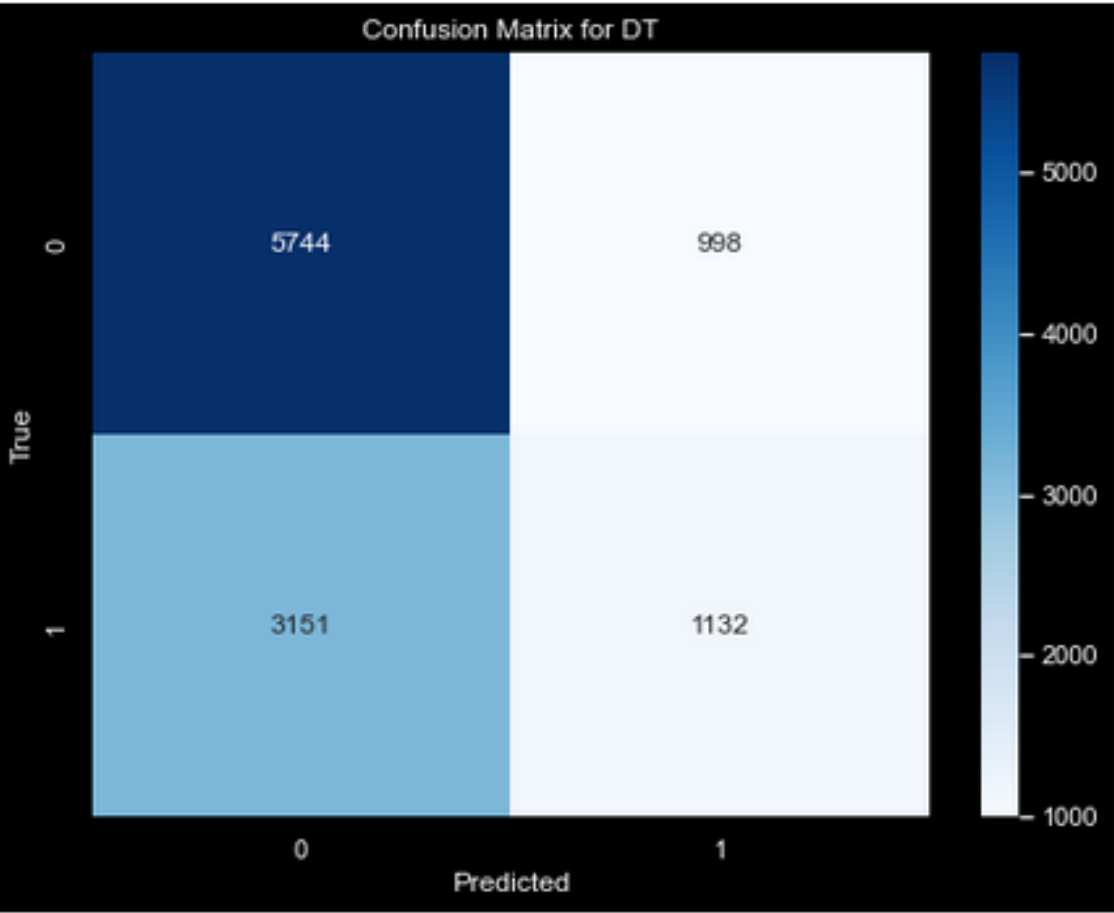
Our first results show us promising results in bout Linear Regression and Random Forest Classifiers, but Decision Tree seems to lack a few parameters to be completely operational.

# RE-EVALUATING THE MODEL WITH THE BEST PARAMETERS

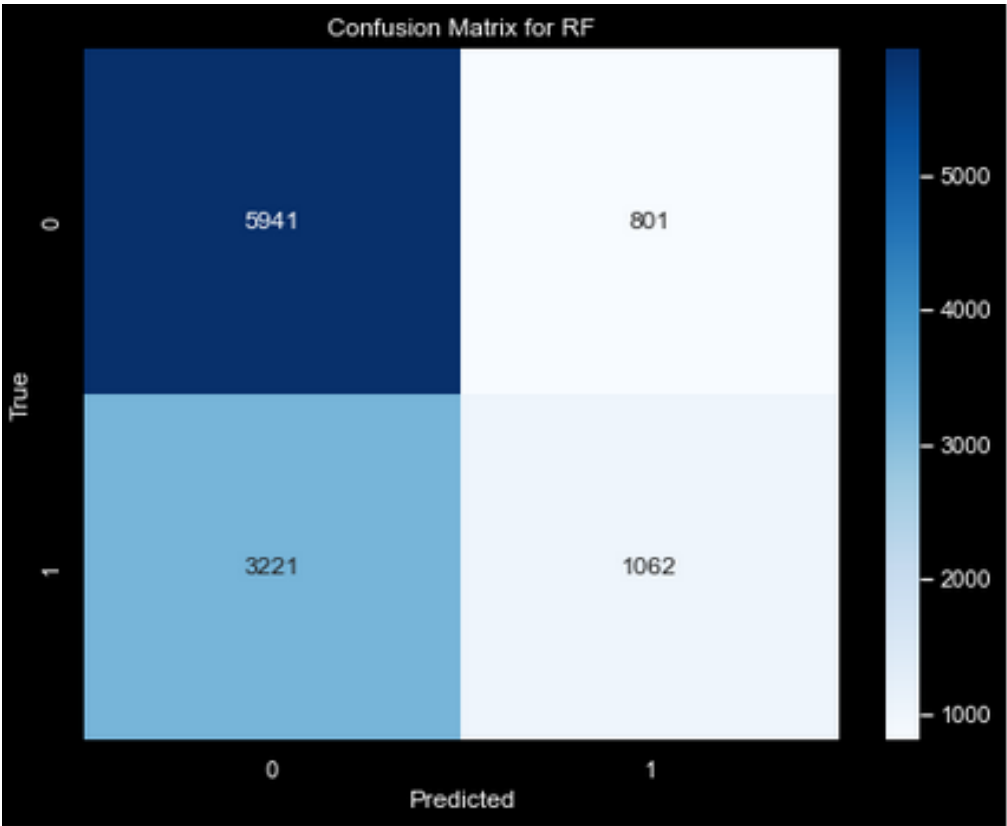
Linear  
Regression



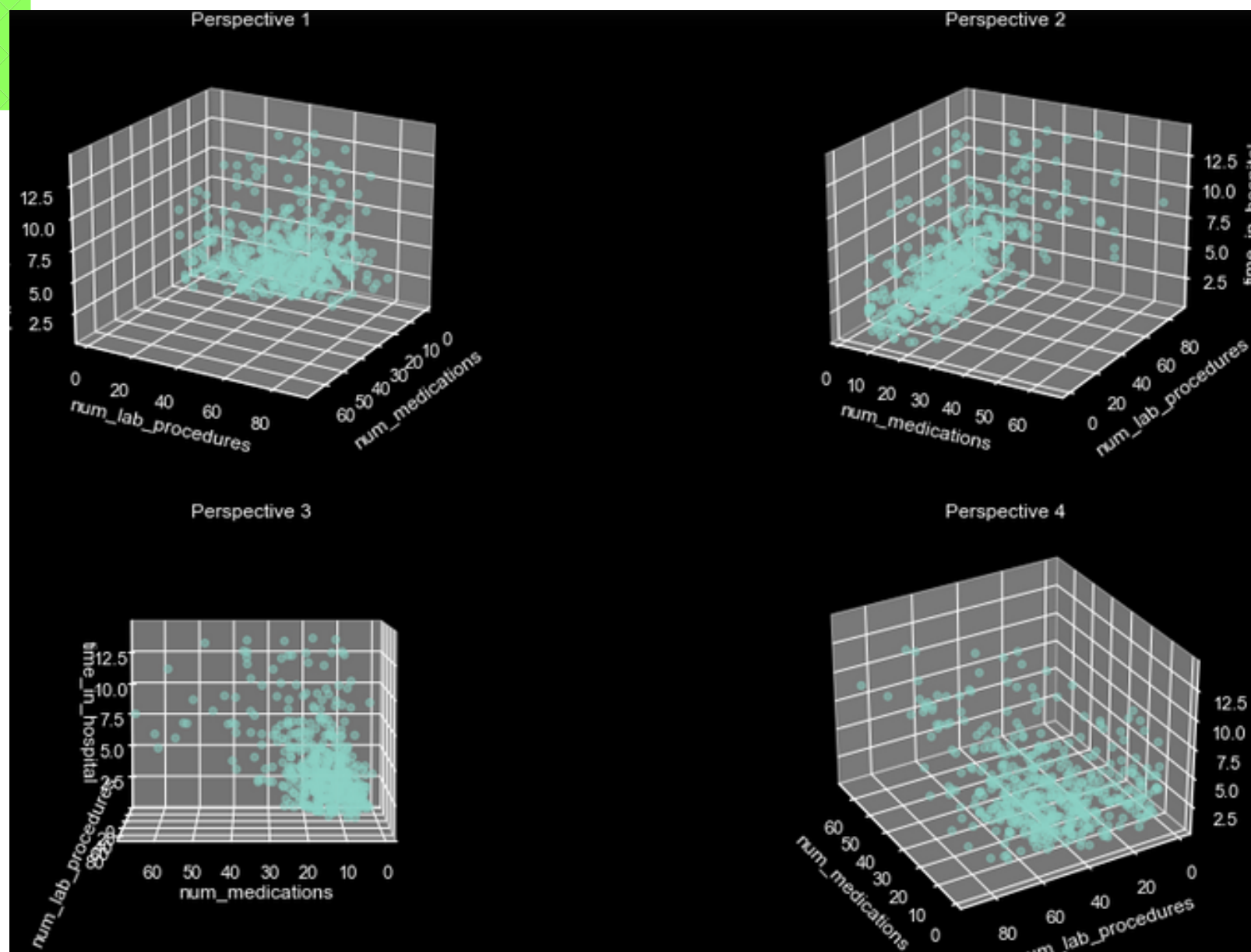
Doctor  
appointment



Random Forest





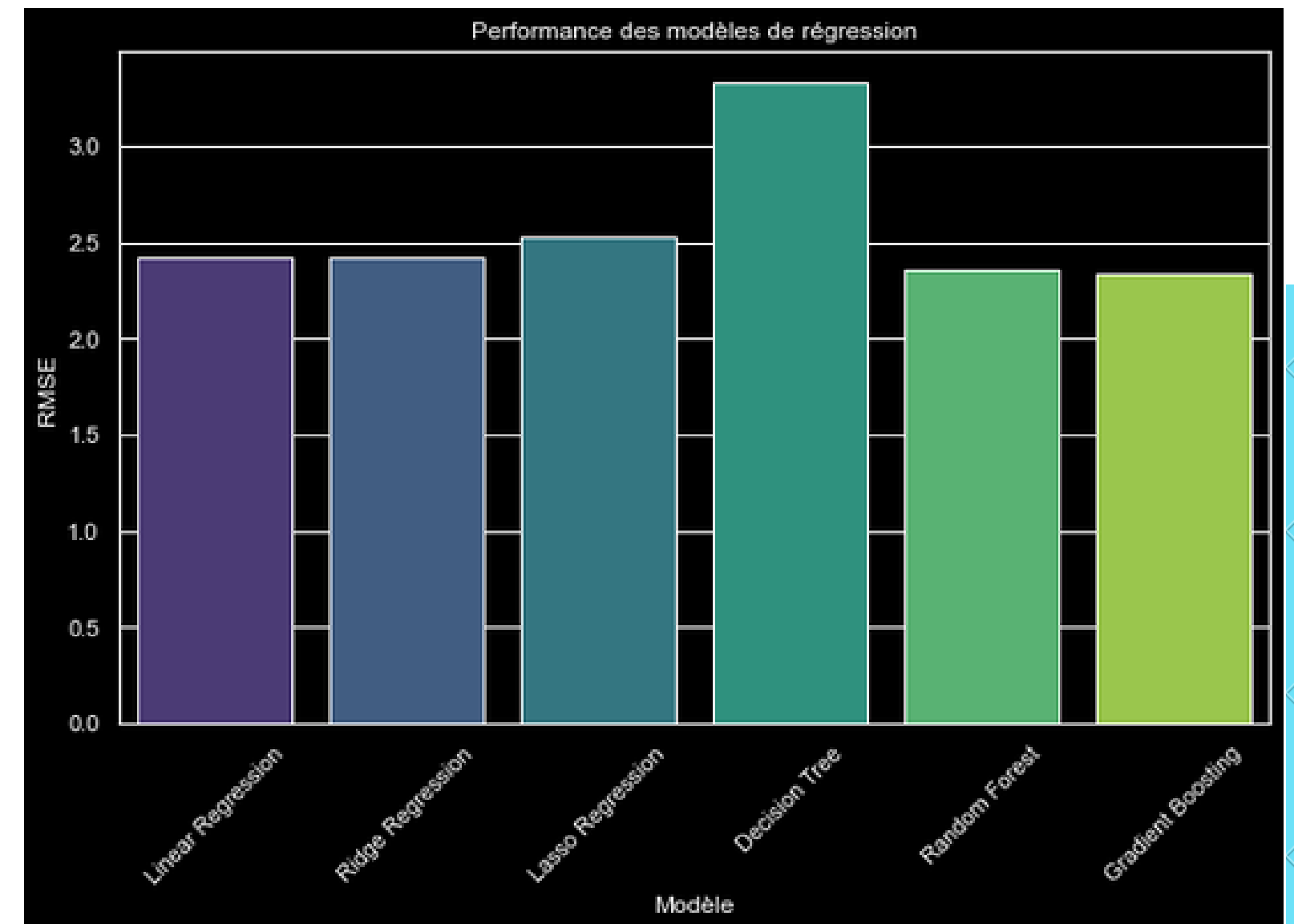


## Data explained

Our approach involved a straightforward process: excluding health parameters. Health, in this context, is an aggregate of components considering the number of days spent in the hospital. The solution to our inquiry was inherent in the question itself. Additionally, we had to convert discrete variables into a binary format using dummy variables, ensuring compatibility with the sklearn library.

## REGRESSION STUDY

we observe that the Gradient Boosting approach appears to be the most effective, although it is closely trailed by the other methods, with the exception of the decision tree, which performs notably worse. Regarding the most influential parameters, we can confidently assert that they are 'num\_medications' and 'num\_lab\_procedures.' This correlation makes perfect sense, as a higher count of lab procedures is associated with a prolonged hospital stay. However, it comes as a surprise to us that the number of medications is not ranked higher; this finding is unexpected





# CONCLUSION

We have now studied this dataset from the Machine Learning Models point of view, and can conclude to good results with the information we have extracted

## Reflexions

With the obtained results, we can already make some previsions that lean more towards preventive measures than pinpoint predictions. This shift from a purely reactive to a more preventive stance in healthcare could significantly change how we approach chronic diseases like diabetes. It emphasizes the importance of early intervention and continuous, personalized care to manage the disease more effectively.

In conclusion, while our study has made considerable progress in predicting patient readmission, it also highlights the need for further research. This research should aim to uncover the less visible variables affecting diabetes, improve measurement methods, and ultimately, guide us towards more effective and personalized healthcare solutions

# THANK YOU

---