



**Data Glacier**

Your Deep Learning Partner

# Drug Persistency Project

**Virtual Internship:** Week 11 Presentation on EDA on the Dataset

**Group Name:** Attack on Data

**Group Members:**

Armin Khayati (United Arab Emirates)

Ezzuldin Zaky (United Arab Emirates)

Orcun Sami Tandogan (Turkey)

**Date:** 10-May-2022

# Problem description

- ❑ One of the challenges for all pharmaceutical companies is to understand the persistency of drug as per the physician prescription. To solve this problem ABC pharma company approached an analytics company to automate this process of identification.
- ❑ Objective : Find important features and prepare them by Feature Engineering and Feature Selection techniques for training with machine learning algorithms.
- ❑ The analysis has been divided into several parts:
  - Data Understanding
  - Data Cleaning
  - Data insights and visualization
  - Recommendations

# Data Exploration

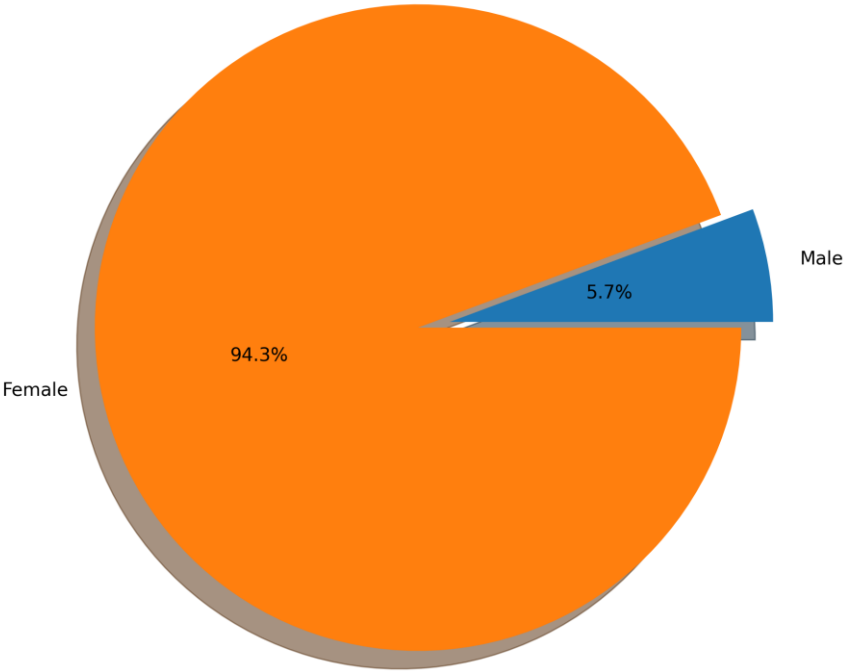
- 68 Features, two numerical and 66 categorical, including :
  - General features such as (Patient general info)
  - Diseases/Drugs Factors
  - Clinical Factors
- Total number of patients : 3424

## **Assumptions:**

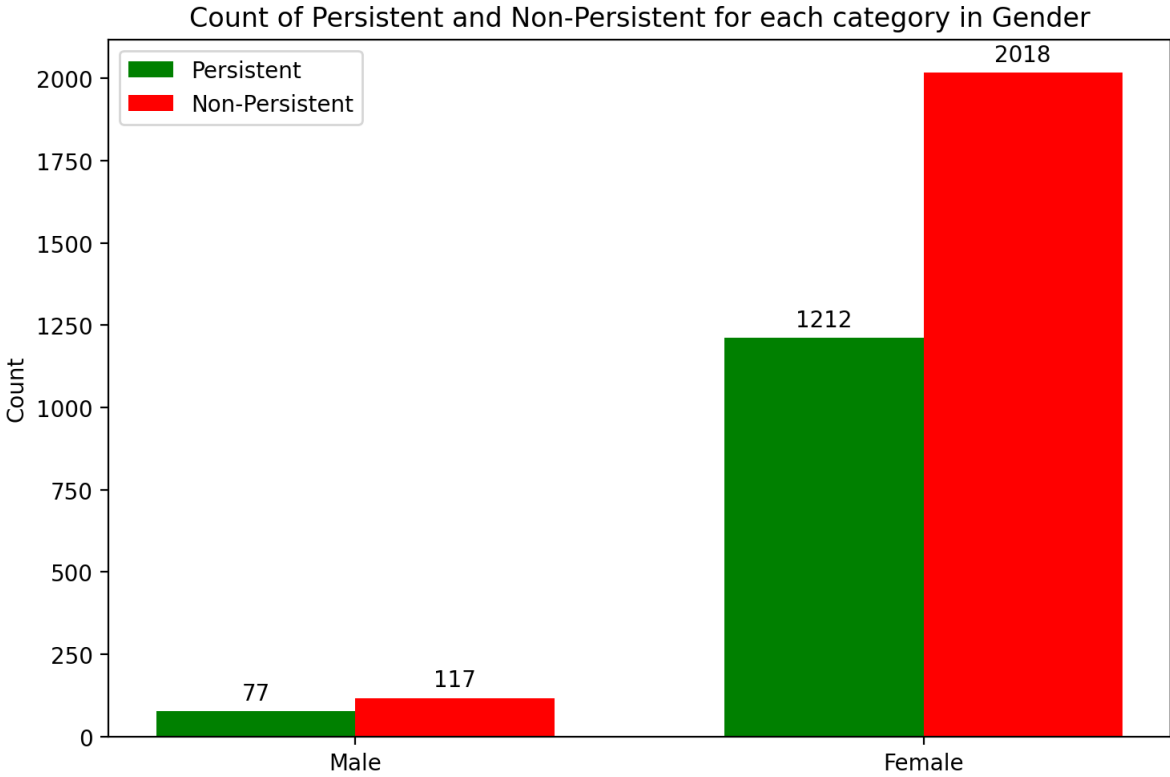
- Dataset is imbalanced.
- It follows a normal distribution.
- Patients' data were gathered accurately without any errors in testing or examination.

# Patients General Info Analysis

Gender

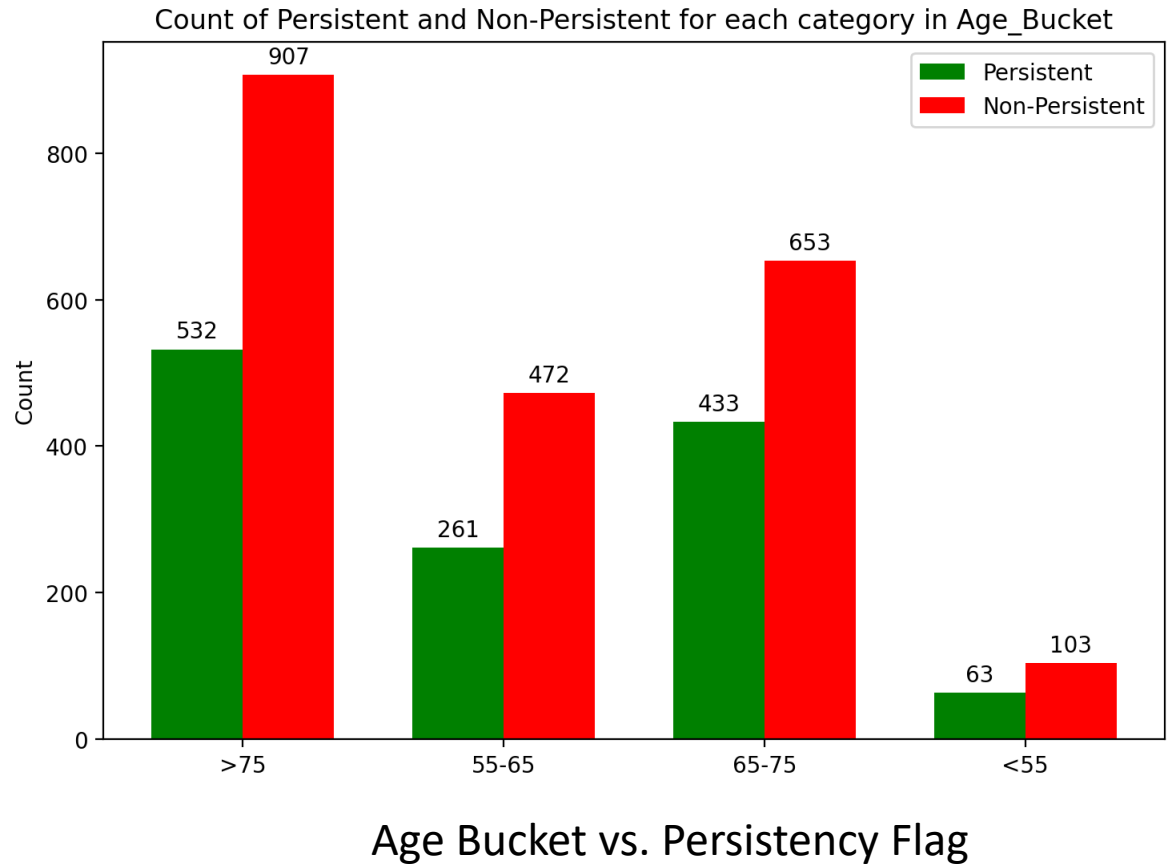
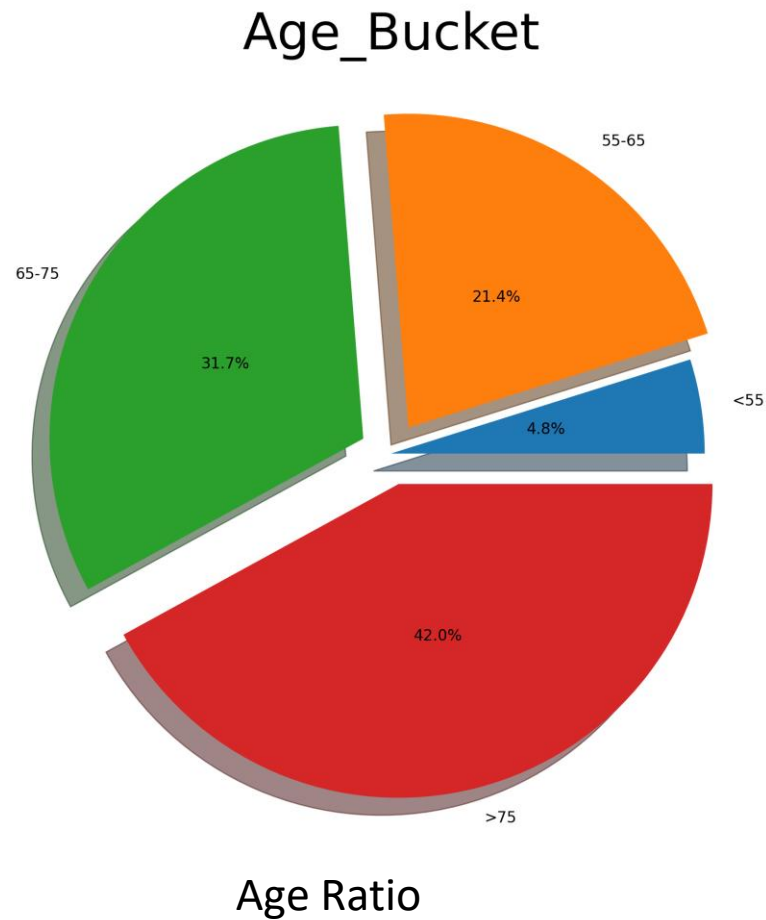


Gender Ratio



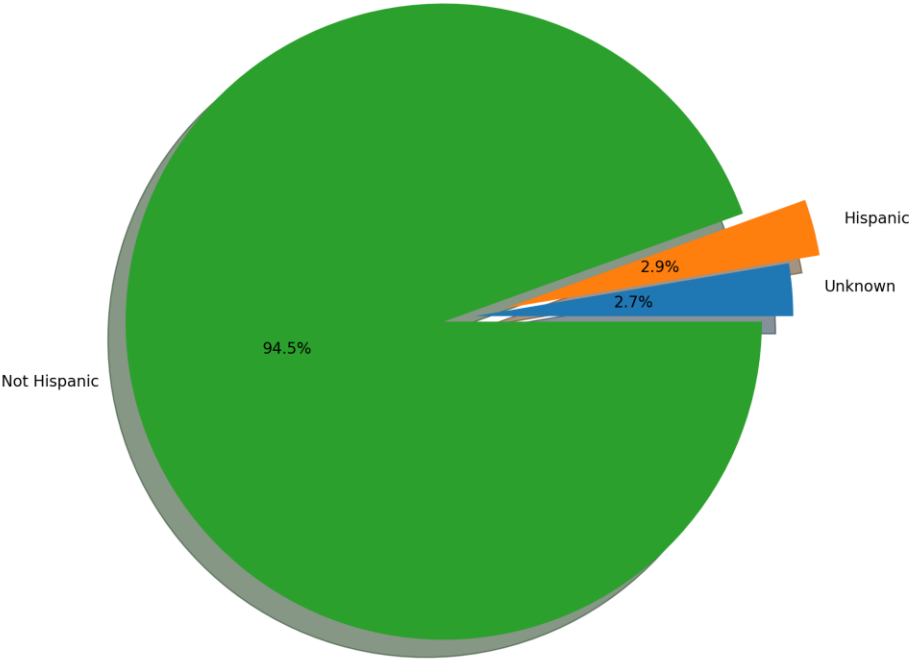
Gender Ratio vs. Persistency Flag

# Patients General Info Analysis



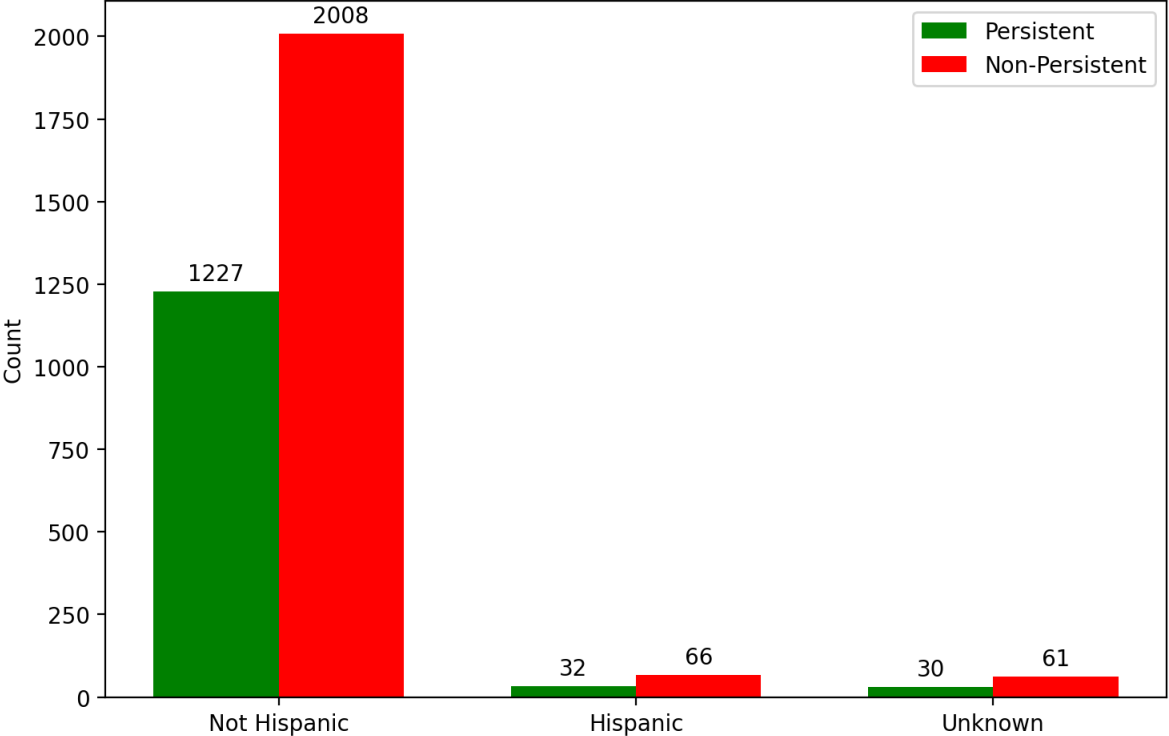
# Patients General Info Analysis

Ethnicity



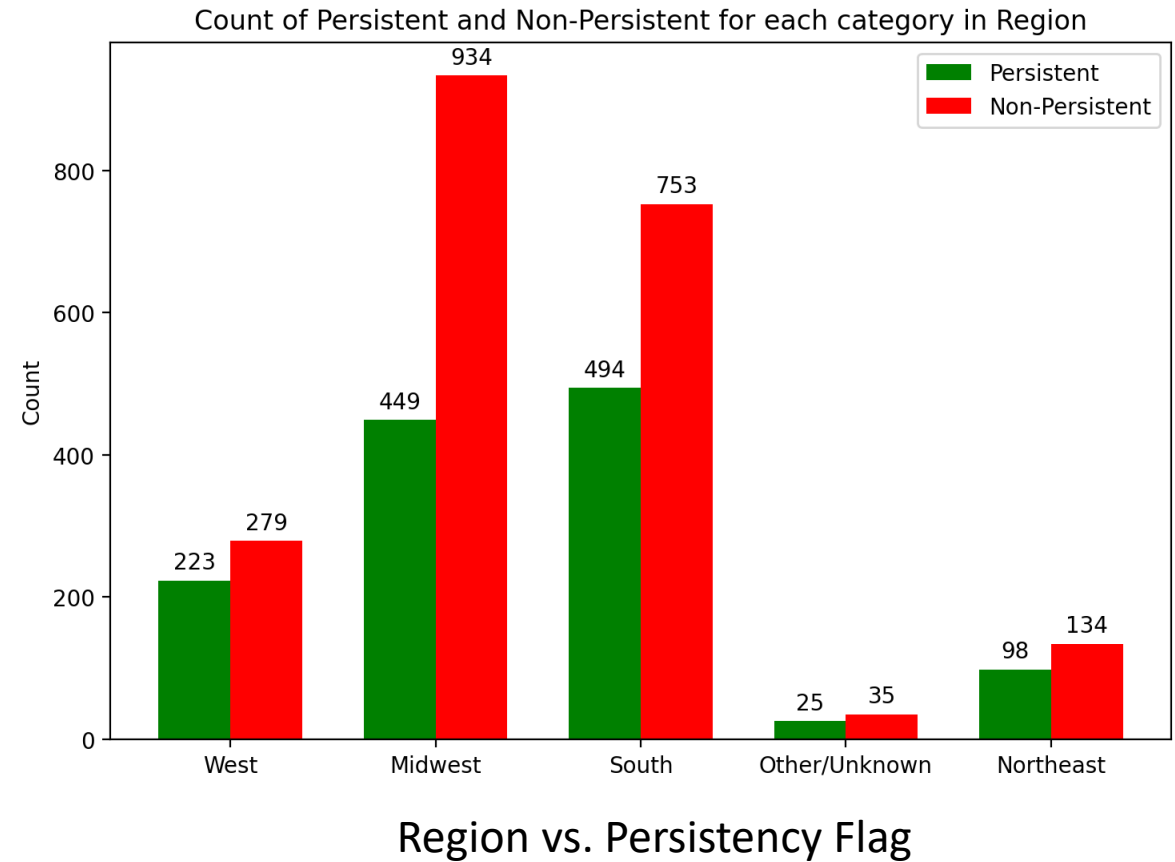
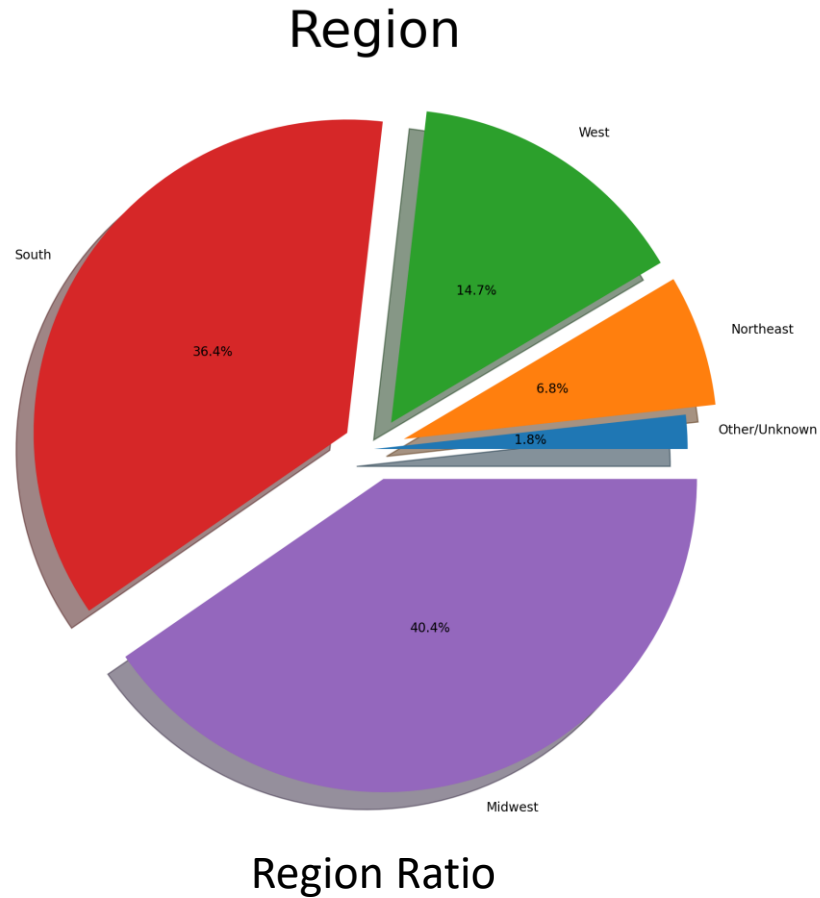
Ethnicity Ratio

Count of Persistent and Non-Persistent for each category in Ethnicity

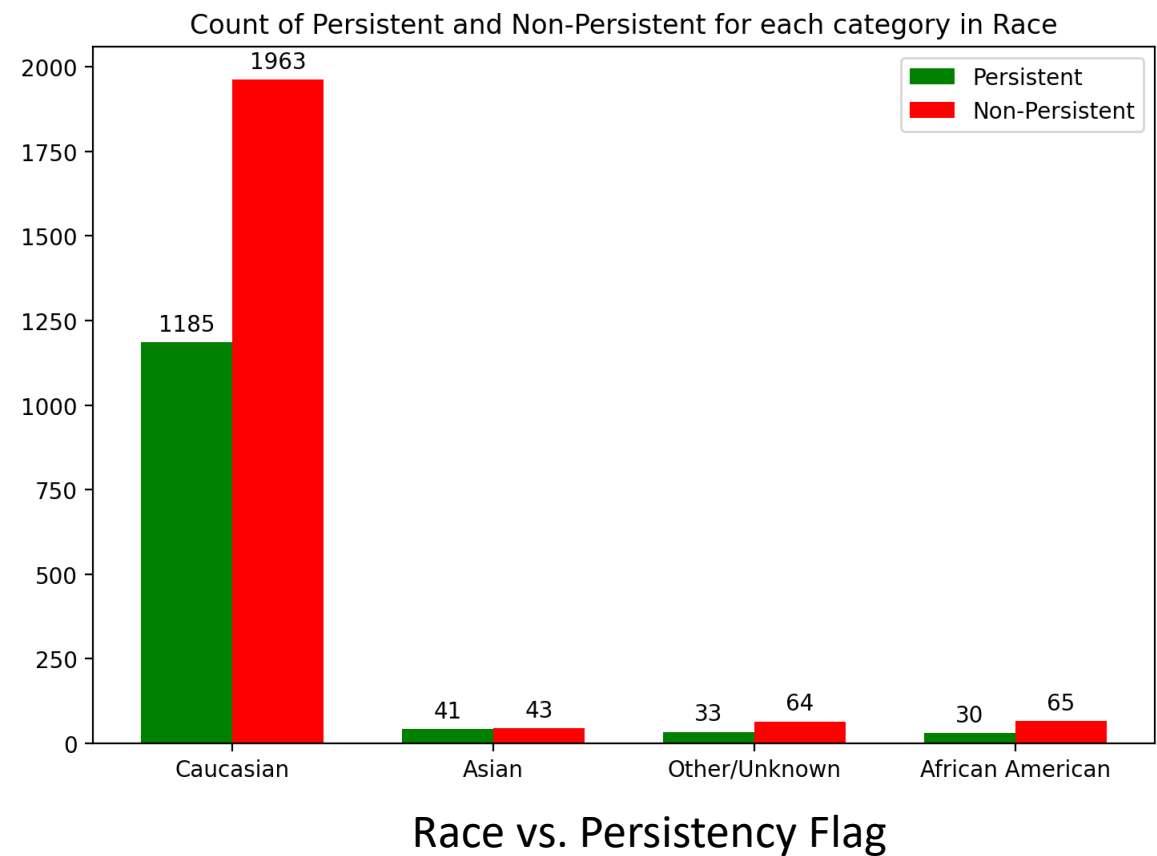
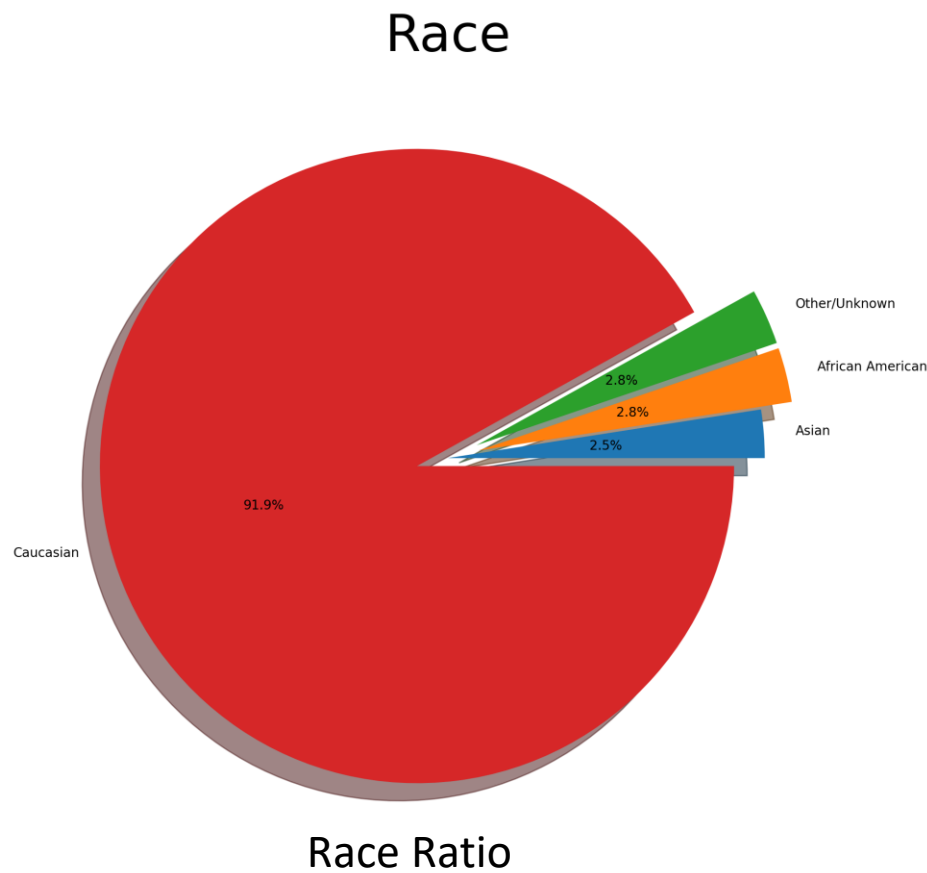


Ethnicity vs. Persistency Flag

# Patients General Info Analysis

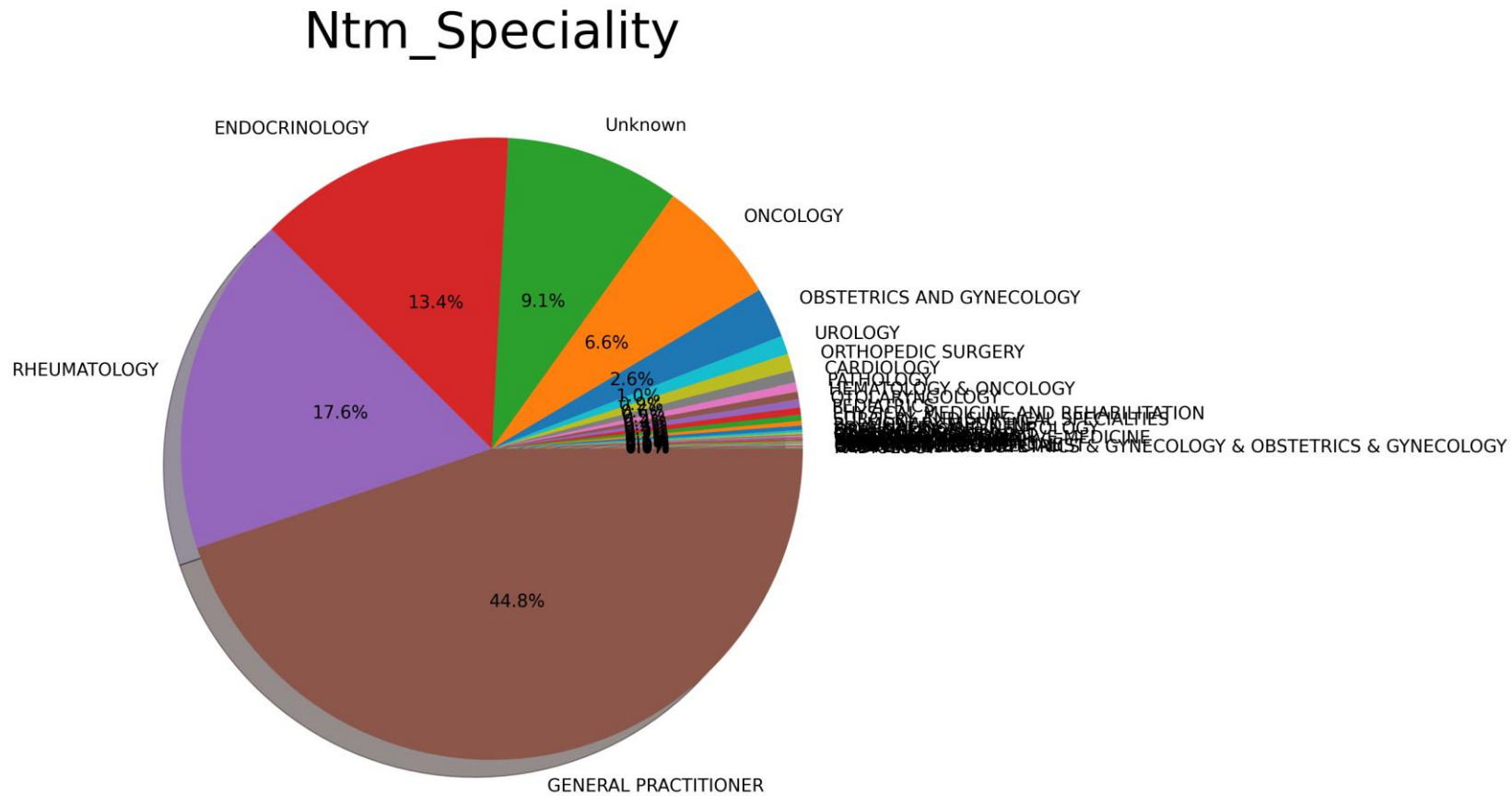


# Patients General Info Analysis





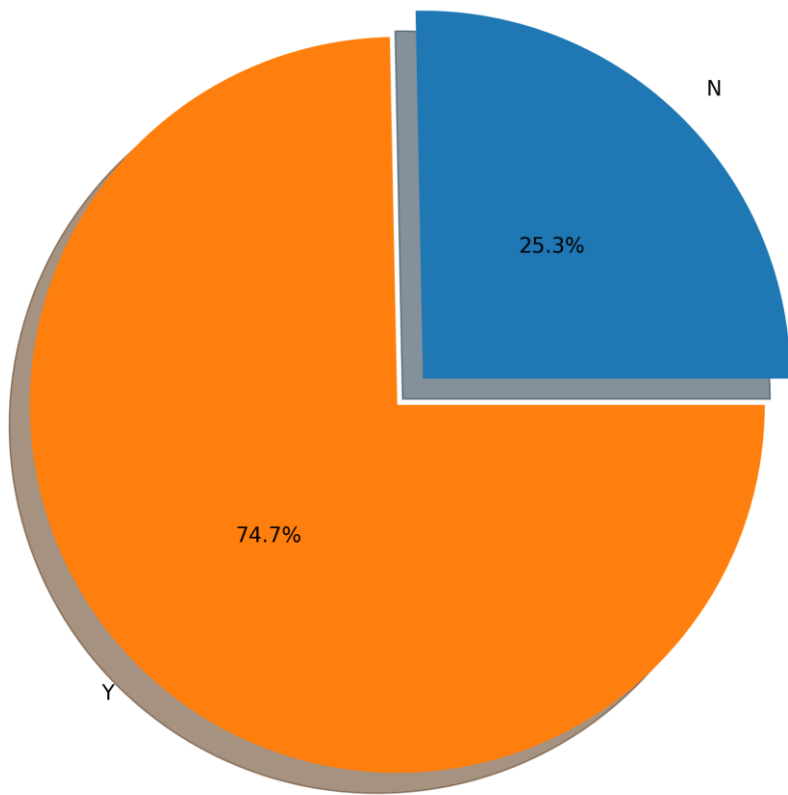
# Patients General Info Analysis



## NTM\_Speciality Ratio

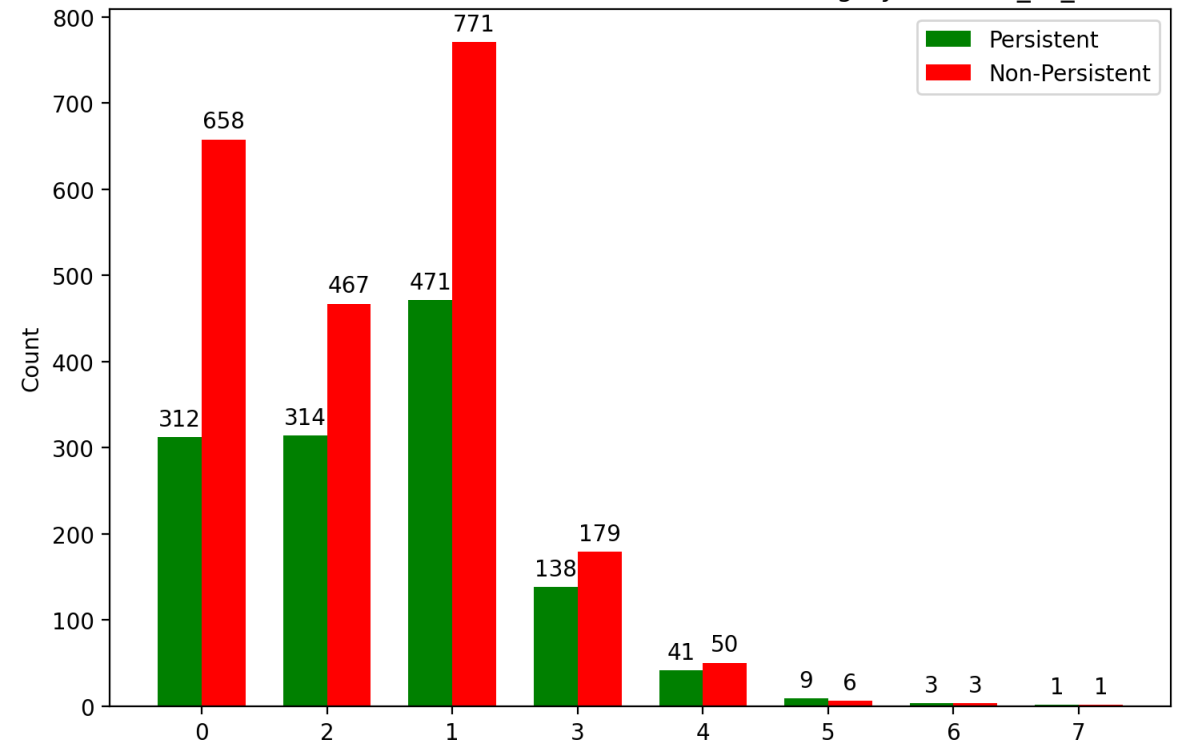
# Patients General Info Analysis

Idn\_Indicator



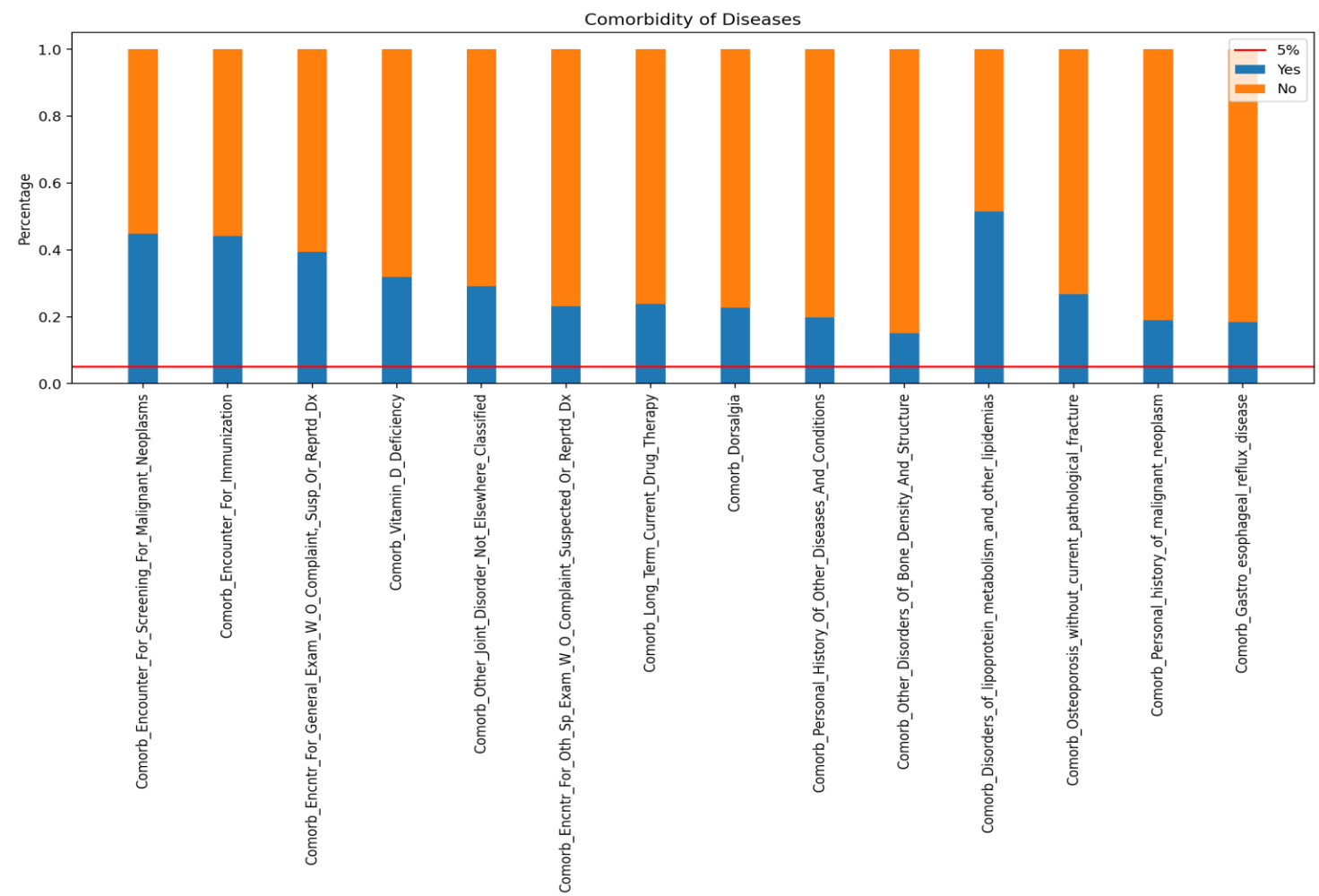
IDN Indicator Ratio

Count of Persistent and Non-Persistent for each category in Count\_Of\_Risks



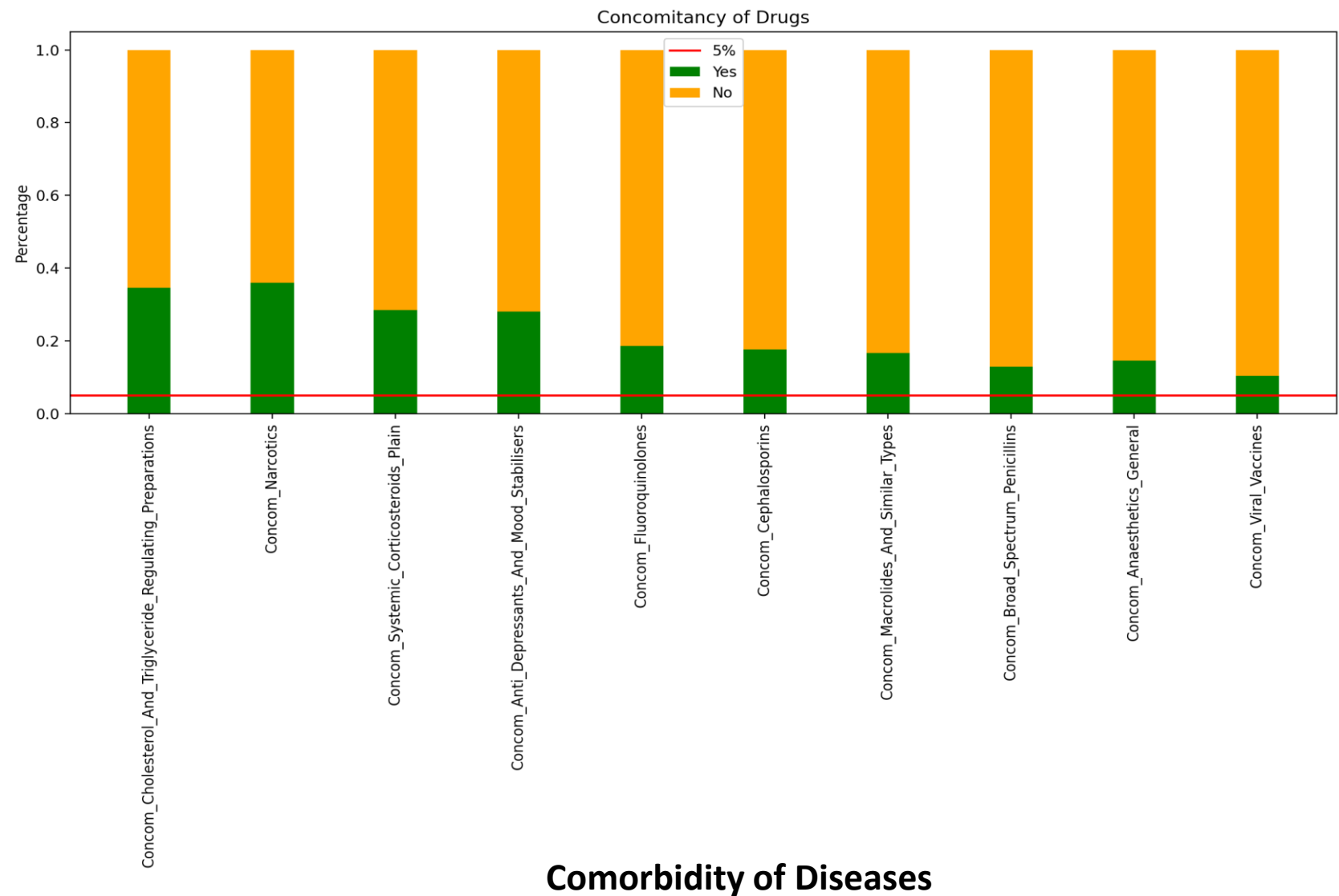
Count of Risk vs . Persistency Flag

# Concomitancy of Drugs

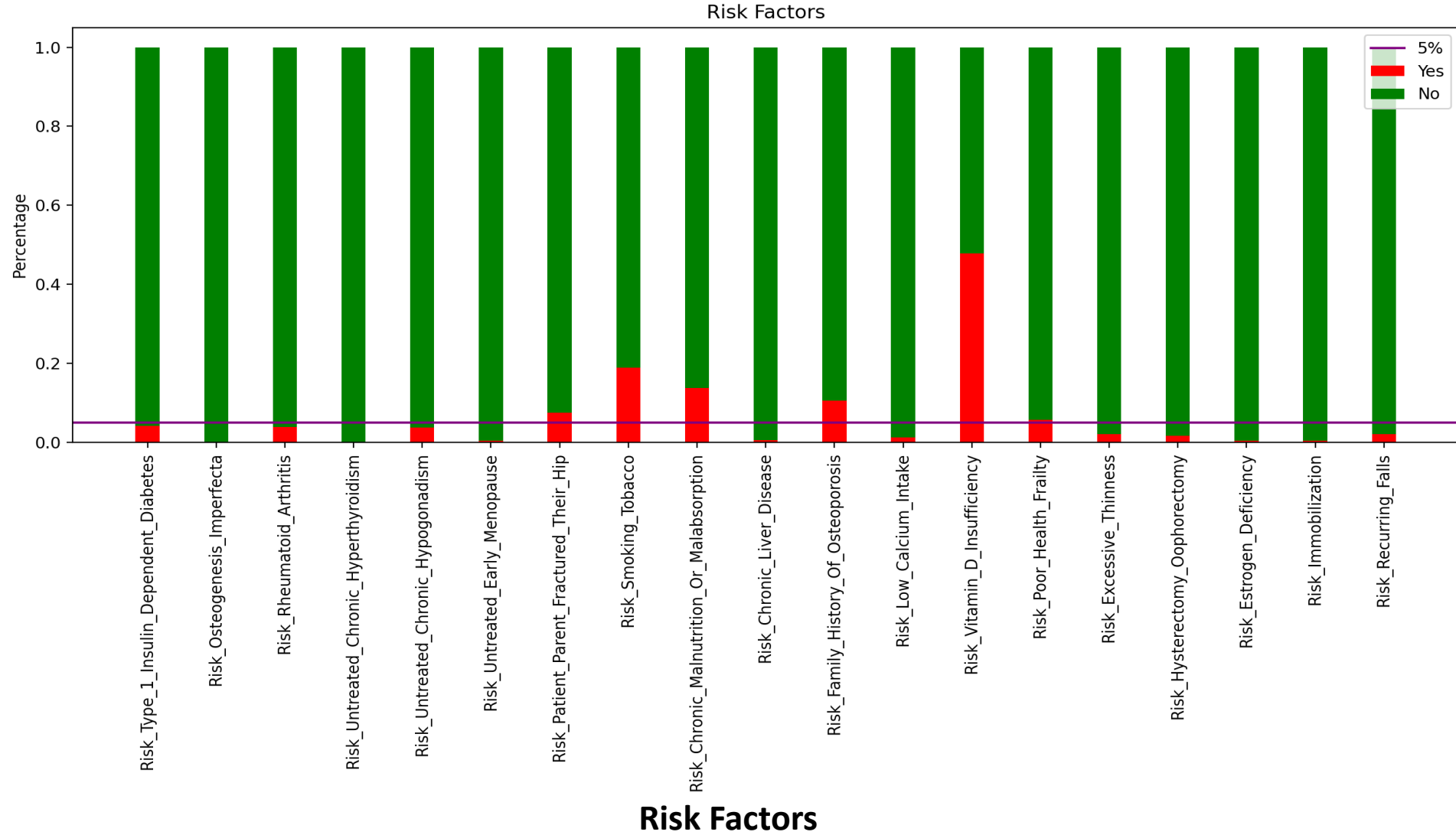


Concomitancy of Drugs

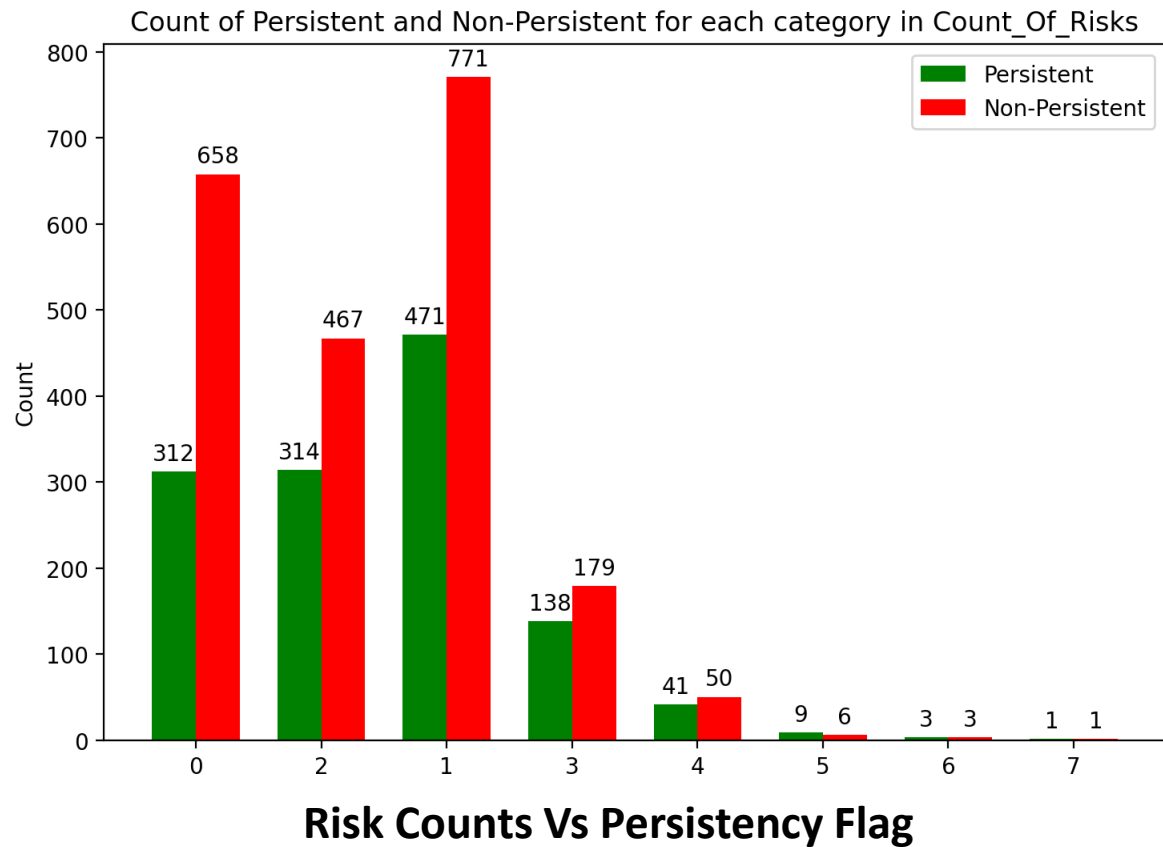
# Comorbidity of Diseases



# Risk Factors

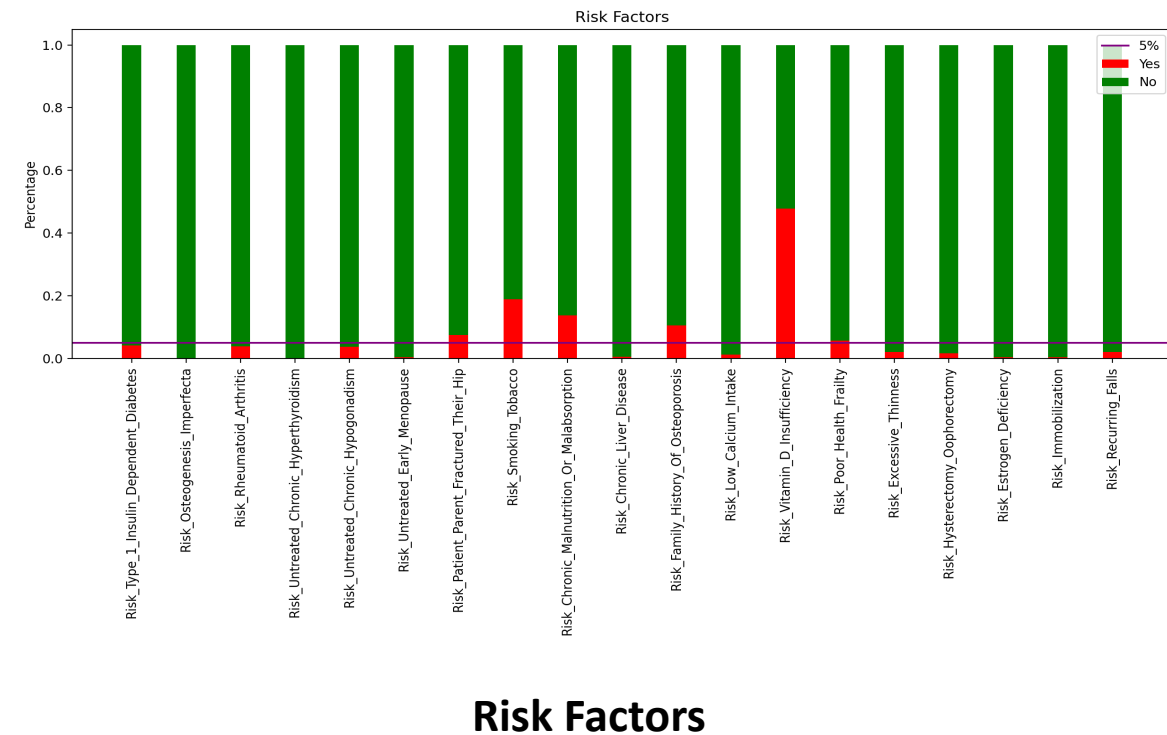


# Risk Factor Analysis



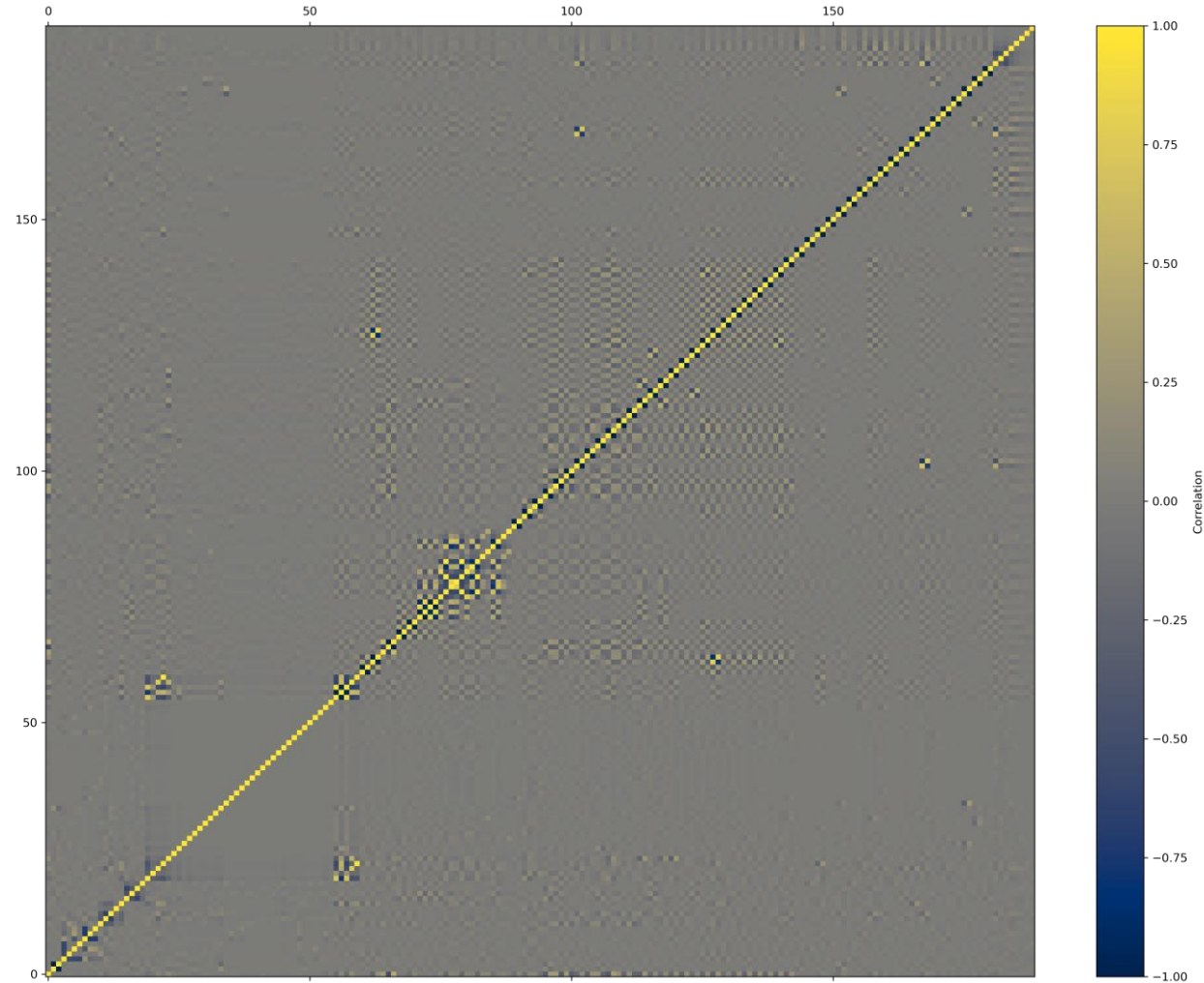
- Number of non-persistence is higher in lower counts of risks.
- Patients with zero count of risk have the highest Non-Persistent Ratio
- Low risk patients were found to be less persistent than the high-risk ones.

# Risk Factors



- Among the risk factors most of them have less than 5% chance to endanger treatment.
- The risk factor with highest chance is Vitamin D Insufficiency and others above 5% are:
  - Poor Health Frailty
  - Family History Of Osteoporosis
  - Chronic Malnutrition Or Malabsorption
  - Smoking Tobacco
  - Patient Parent Fractured Their Hip
- Rest of the factors have less than 5% risk to endanger treatment.

# Features Correlation



**Correlation Matrix**

- We will be Removing variables with more than 98% correlation.



# EDA Recommendations

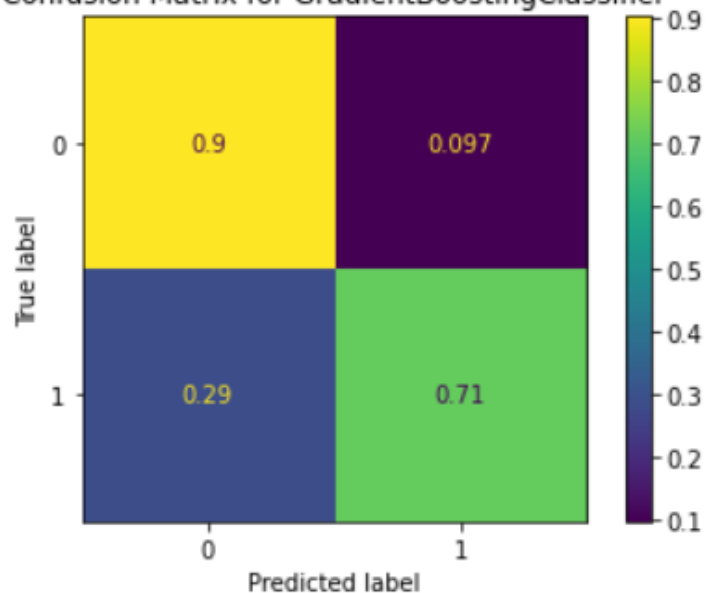
From the Exploratory Data Analysis (EDA) done on the dataset, we will recommend these instructions:

1. Handling Unknown values for Race, Region, and Ethnicity Variables
  - Using mode as an imputer as an imputer on Race and Ethnicity variables.
  - For Region variable, because most of the people with Unknown Region have Not Hispanic Ethnicity, and Most of people with Not Hispanic Ethnicity, have Midwest Region, we will replace Unknown Regions with Midwest.
2. Handling Rare Labels: Finding categories less than 5 percent in each variable, then merging those categories into one or drop them if the variable only has 2 categories (e.g., Y/N) and cardinality of one them is less than 5 percent.
3. Grouping integer values of Count\_Of\_Risks variable into two bins: Bin 1 is [0,1,2,3] and Bin 2 is [4,5,6,7].
4. **One hot encoding** all the variables after doing above tasks
5. Removing variables with more than **98% correlation**.

# Model selection and building

Accuracy GradientBoostingClassifier: 0.83

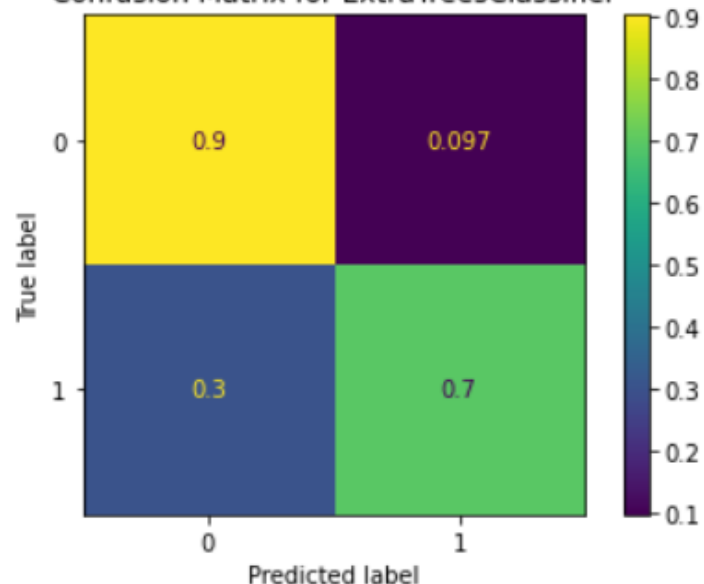
Confusion Matrix for GradientBoostingClassifier



	precision	recall	f1-score	support
0	0.83	0.90	0.86	515
1	0.83	0.71	0.77	341
accuracy			0.83	856
macro avg	0.83	0.81	0.81	856
weighted avg	0.83	0.83	0.82	856

Accuracy ExtraTreesClassifier: 0.82

Confusion Matrix for ExtraTreesClassifier

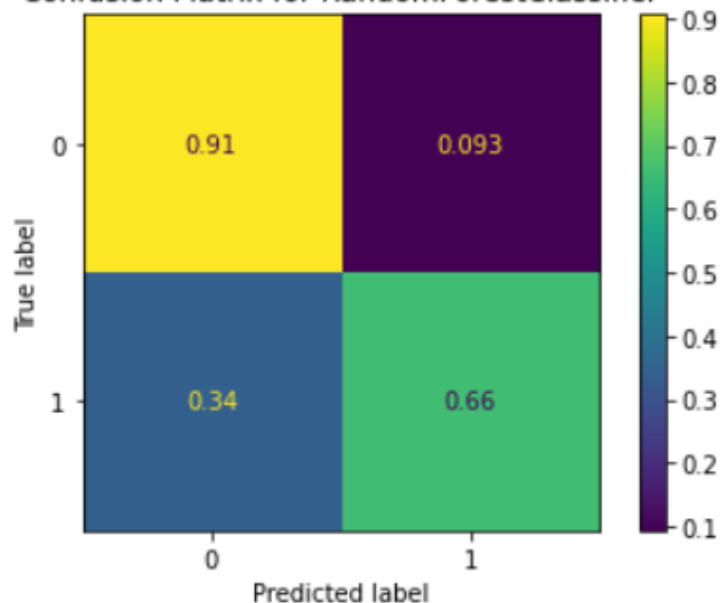


	precision	recall	f1-score	support
0	0.82	0.90	0.86	515
1	0.83	0.70	0.76	341
accuracy			0.82	856
macro avg	0.82	0.80	0.81	856
weighted avg	0.82	0.82	0.82	856

# Model selection and building

Accuracy RandomForestClassifier: 0.81

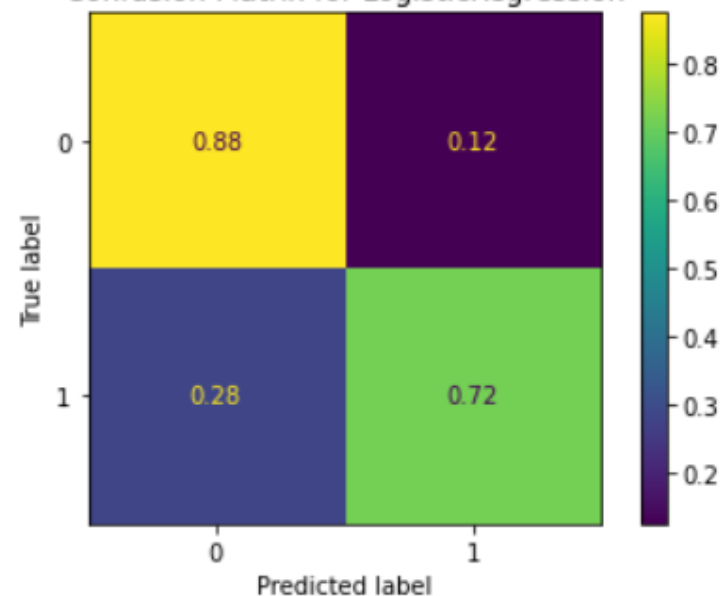
Confusion Matrix for RandomForestClassifier



	precision	recall	f1-score	support
0	0.80	0.91	0.85	515
1	0.82	0.66	0.73	341
accuracy			0.81	856
macro avg	0.81	0.78	0.79	856
weighted avg	0.81	0.81	0.80	856

Accuracy LogisticRegression: 0.81

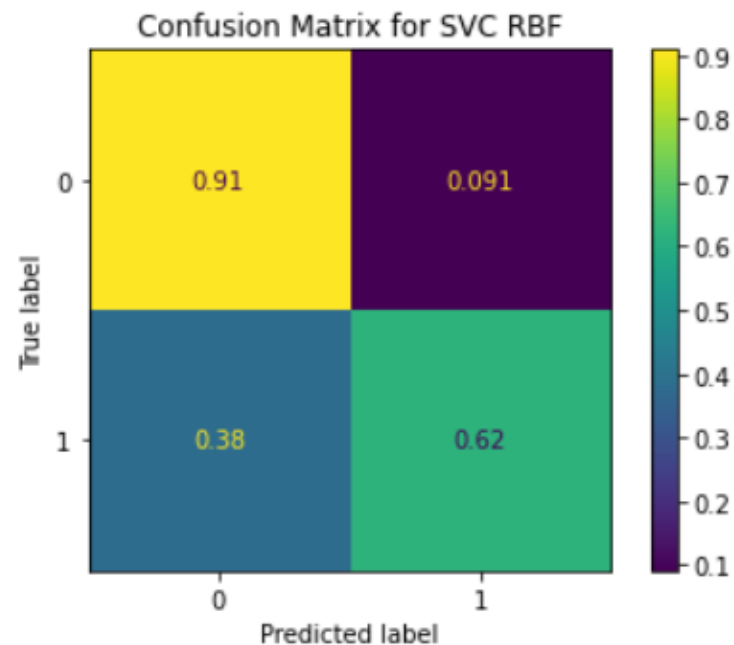
Confusion Matrix for LogisticRegression



	precision	recall	f1-score	support
0	0.82	0.88	0.85	515
1	0.79	0.72	0.75	341
accuracy			0.81	856
macro avg	0.81	0.80	0.80	856
weighted avg	0.81	0.81	0.81	856

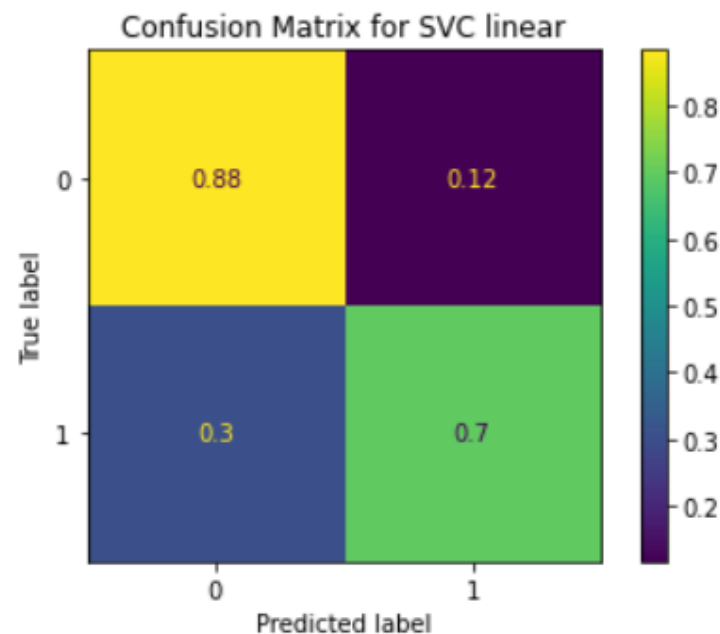
# Model selection and building

Accuracy SVC RBF: 0.79



	precision	recall	f1-score	support
0	0.78	0.91	0.84	515
1	0.82	0.62	0.70	341
accuracy			0.79	856
macro avg	0.80	0.76	0.77	856
weighted avg	0.80	0.79	0.79	856

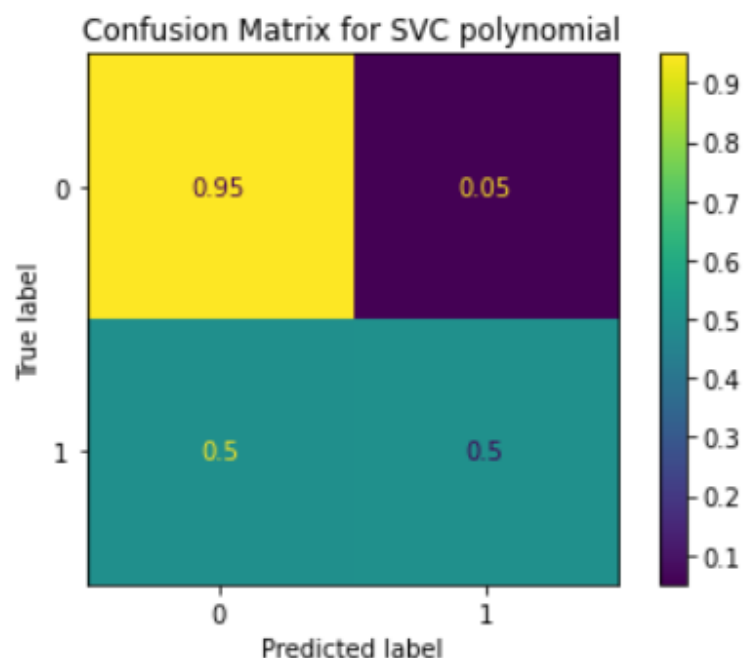
Accuracy SVC linear: 0.81



	precision	recall	f1-score	support
0	0.81	0.88	0.85	515
1	0.80	0.70	0.74	341
accuracy			0.81	856
macro avg	0.81	0.79	0.80	856
weighted avg	0.81	0.81	0.81	856

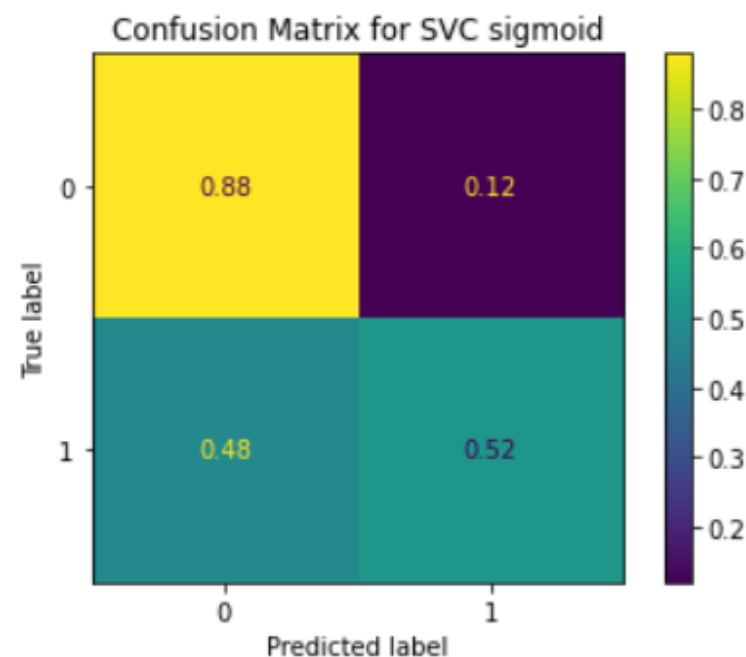
# Model selection and building

Accuracy SVC polynomial: 0.77



	precision	recall	f1-score	support
0	0.74	0.95	0.83	515
1	0.87	0.50	0.64	341
accuracy			0.77	856
macro avg	0.81	0.73	0.73	856
weighted avg	0.79	0.77	0.75	856

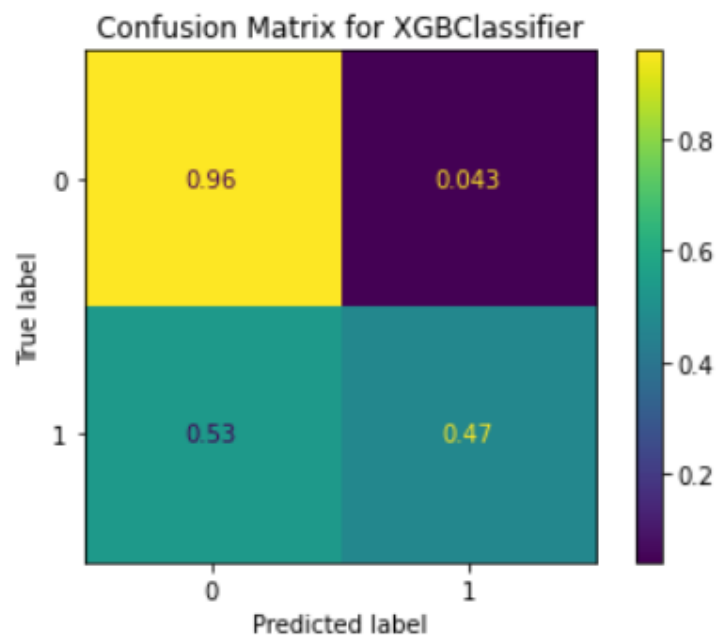
Accuracy SVC sigmoid: 0.74



	precision	recall	f1-score	support
0	0.74	0.88	0.80	515
1	0.74	0.52	0.61	341
accuracy			0.74	856
macro avg	0.74	0.70	0.71	856
weighted avg	0.74	0.74	0.73	856

# Model selection and building

Accuracy XGBClassifier: 0.76



	precision	recall	f1-score	support
0	0.73	0.96	0.83	515
1	0.88	0.47	0.61	341
accuracy			0.76	856
macro avg	0.80	0.71	0.72	856
weighted avg	0.79	0.76	0.74	856

# Final Recommendation

Then Based on previous slides, we recommend these machine learning techniques:

- Gradient Boosting
- Extra trees classifier
- Random forest classifier

Justification: They are the most accurate techniques and best predictors of drug non persistency, which is a safer error to make than the contrary.

[Link of code](#)

# Thank You



**Data Glacier**

Your Deep Learning Partner