



**Data Glacier**

Your Deep Learning Partner

# Exploratory Data Analysis

HEALTH CARE – PERSISTENCY OF A DRUG

**Data Science Project**

LISUM07

**25/04/2021**

# Background

One of the challenge 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.

# Problem Statement

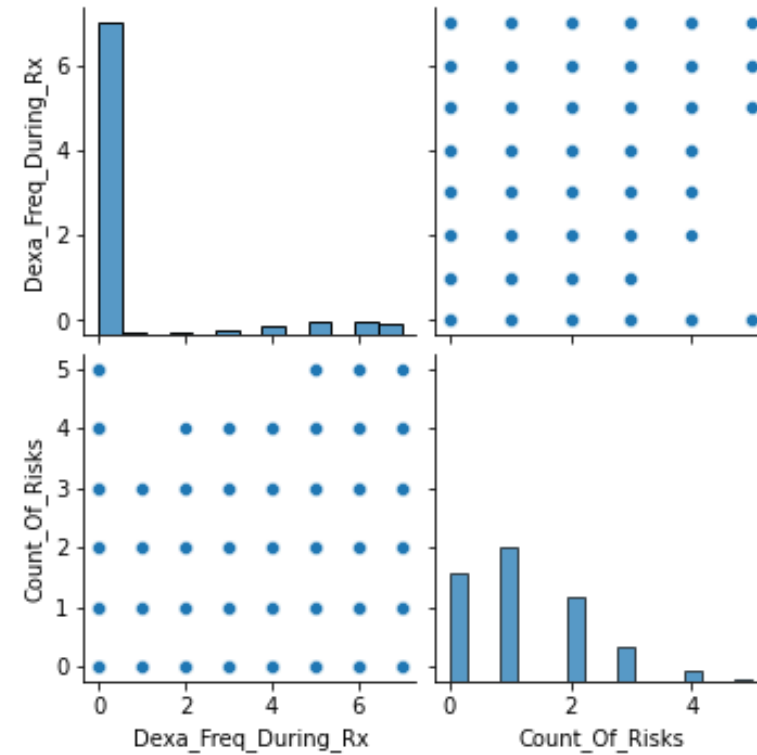
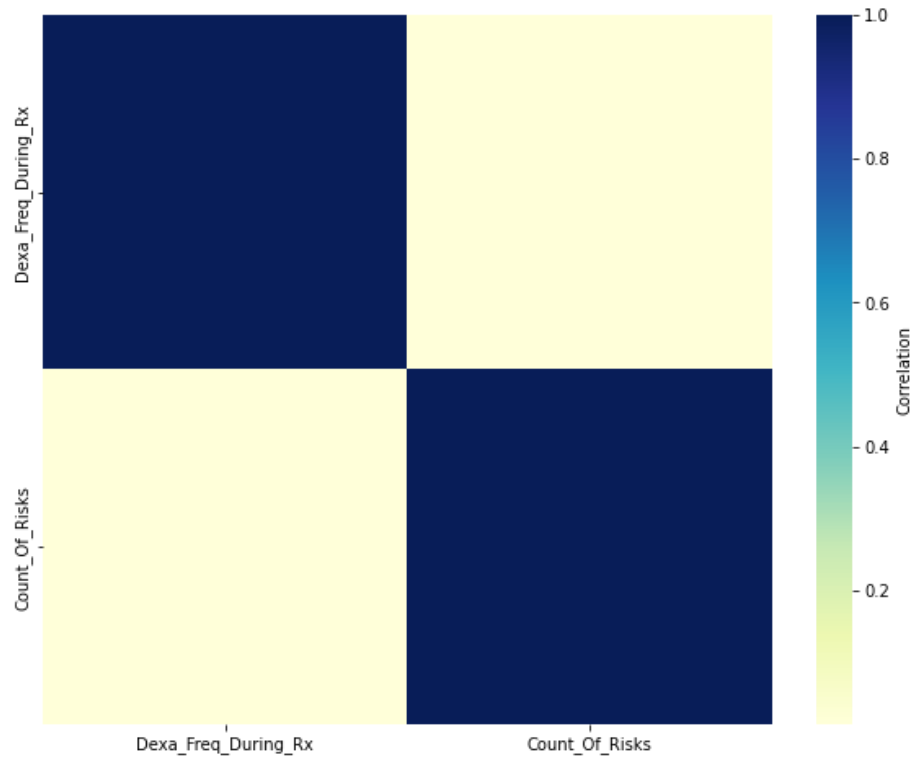
# Problem Statement

- Classification the process of identification of the persistency of the drug as per prescription provided by the physician.
- Understanding the persistency is an issue for pharmaceutical companies and so, ABC pharma company has approached XYZ analytics company to provide insights into the same.
- The role of XYZ analytics company is to undertake this project and provide a detailed understanding regarding the drug persistency.

EDA

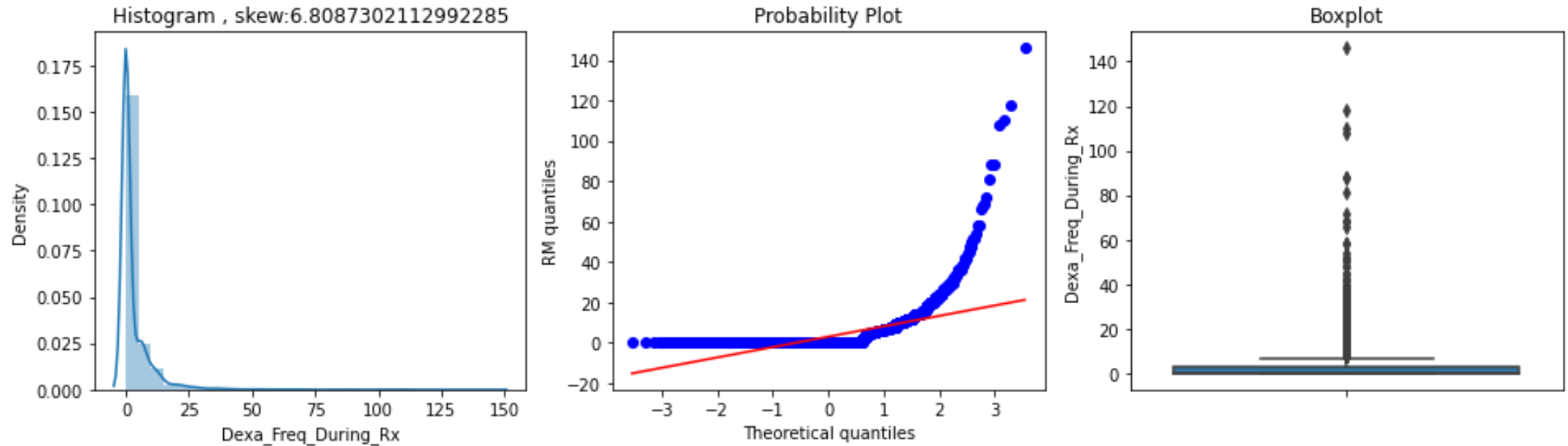
# Numerical Data Analysis

# Correlation



- The correlation values indicate that there is a low correlation between the two variables.
- On the other hand, `Dexa_Freq_During_Rx` has the least distributed values.

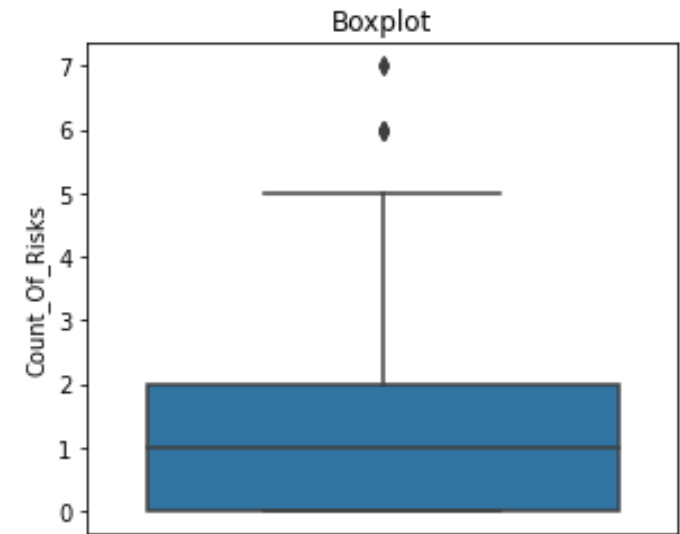
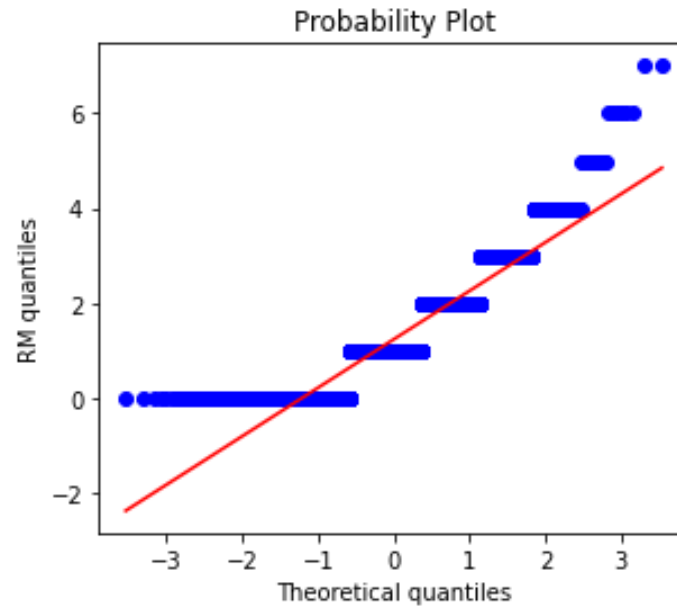
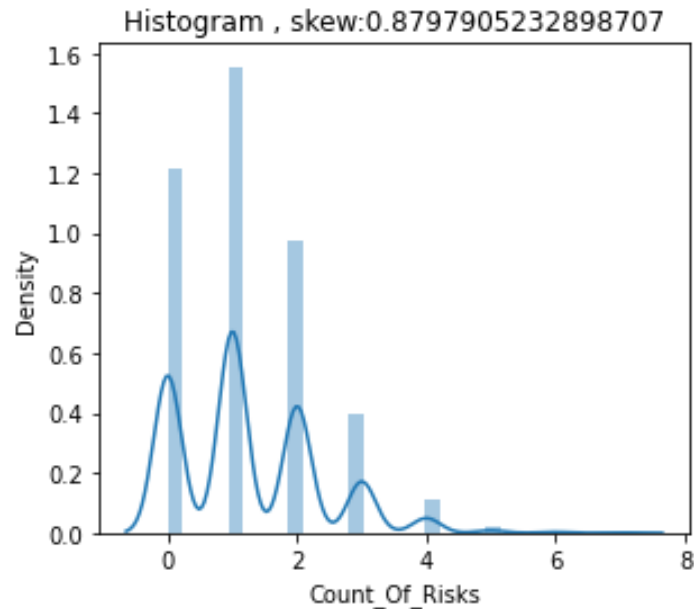
# Dexa\_Freq\_During\_Rx



- From the diagrams above, we can see that most of the frequencies lie between 0 and 20.
- The minimum frequency is 0 and the maximum is around 140.
- The data is highly skewed.
- It has many data outliers, with a mean close to zero.



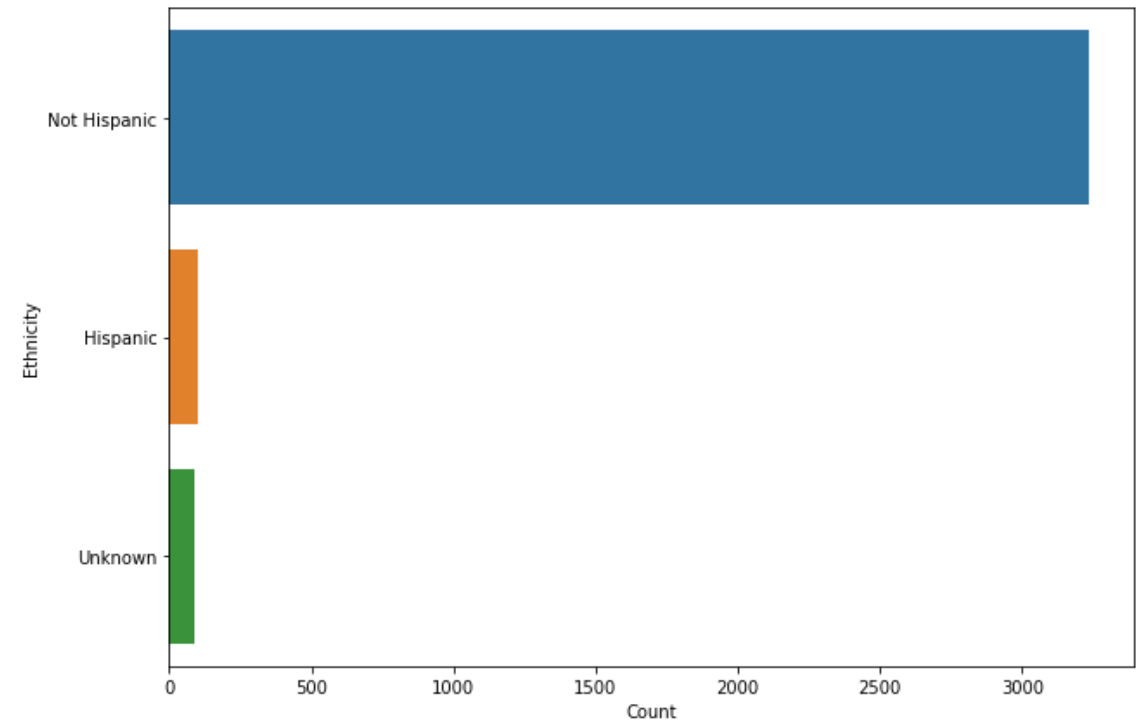
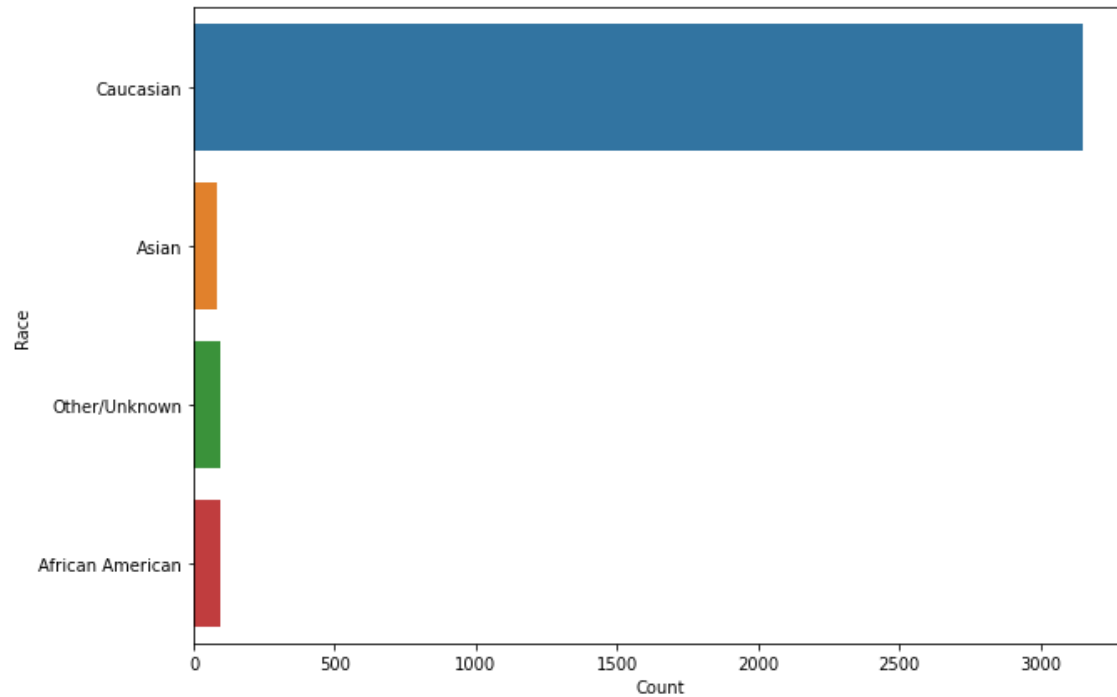
# Count\_of\_Risks



- Most of the count of risks lie between 0 and 1.
- The data is slightly skewed.
- There's a slight difference in the distribution of the count of risks between persistent patient's and non-persistent patients.
- It has a quite data outliers, with a mean close to one.

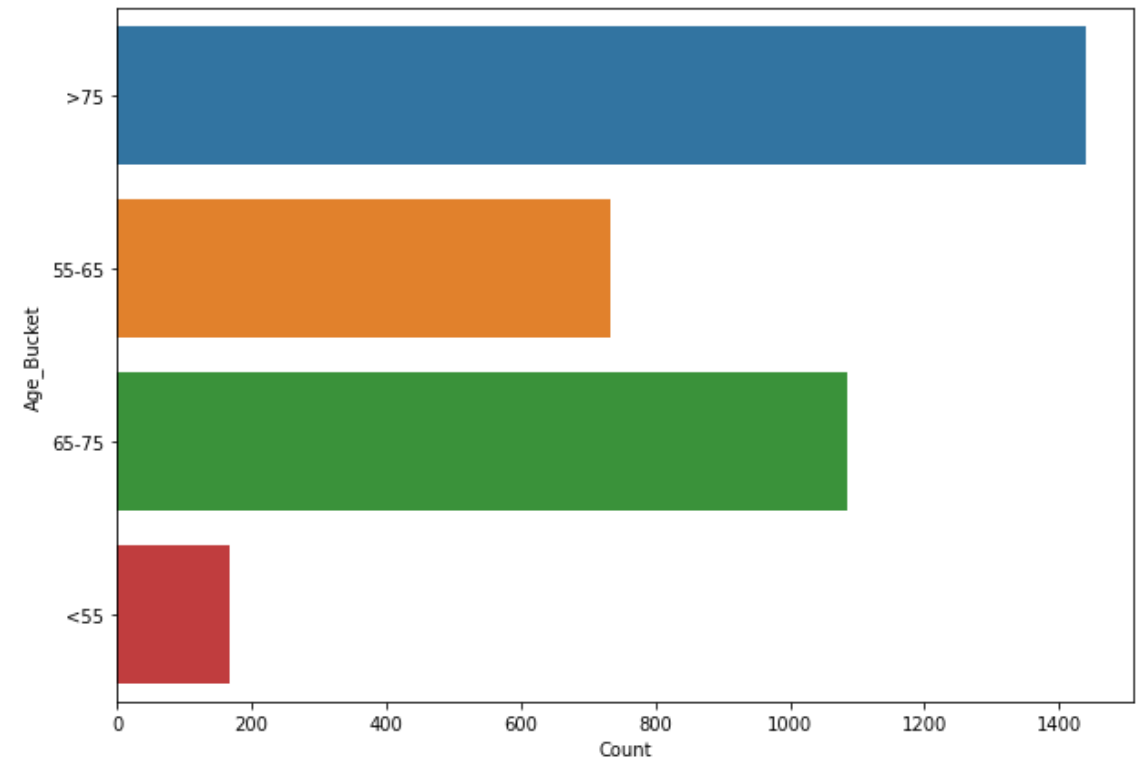
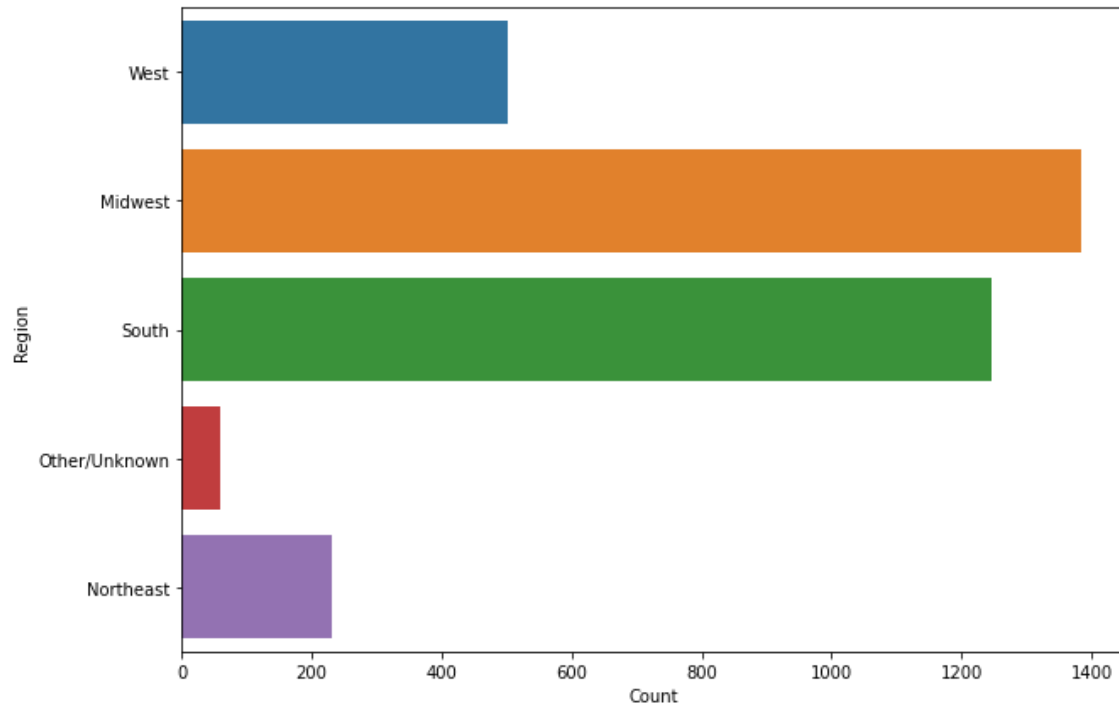
# Categorical Data Analysis

# Race/Ethnicity



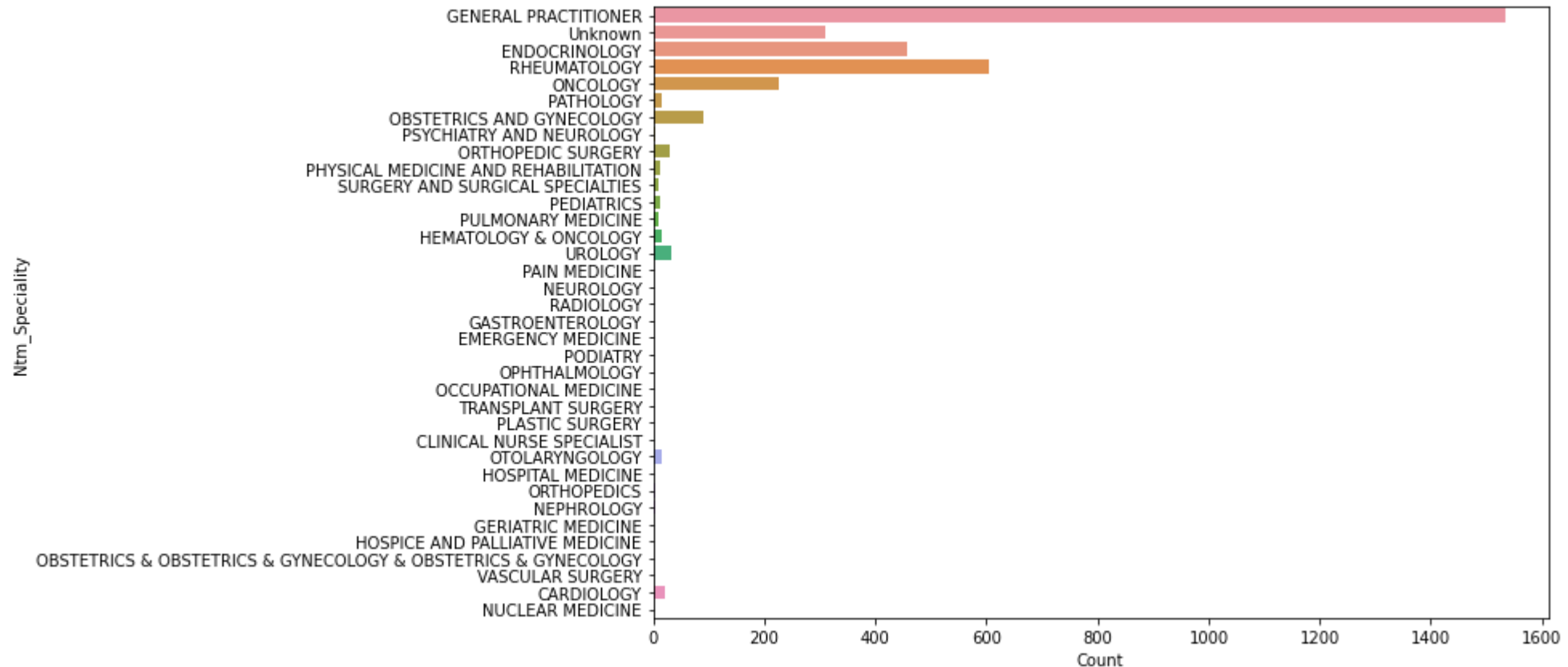
- RACE unbalanced data, greater predominance in Caucasian.
- Unbalanced data ETHNICITY, greater predominance in Not Hispanic.

# Region/Age\_Bucket



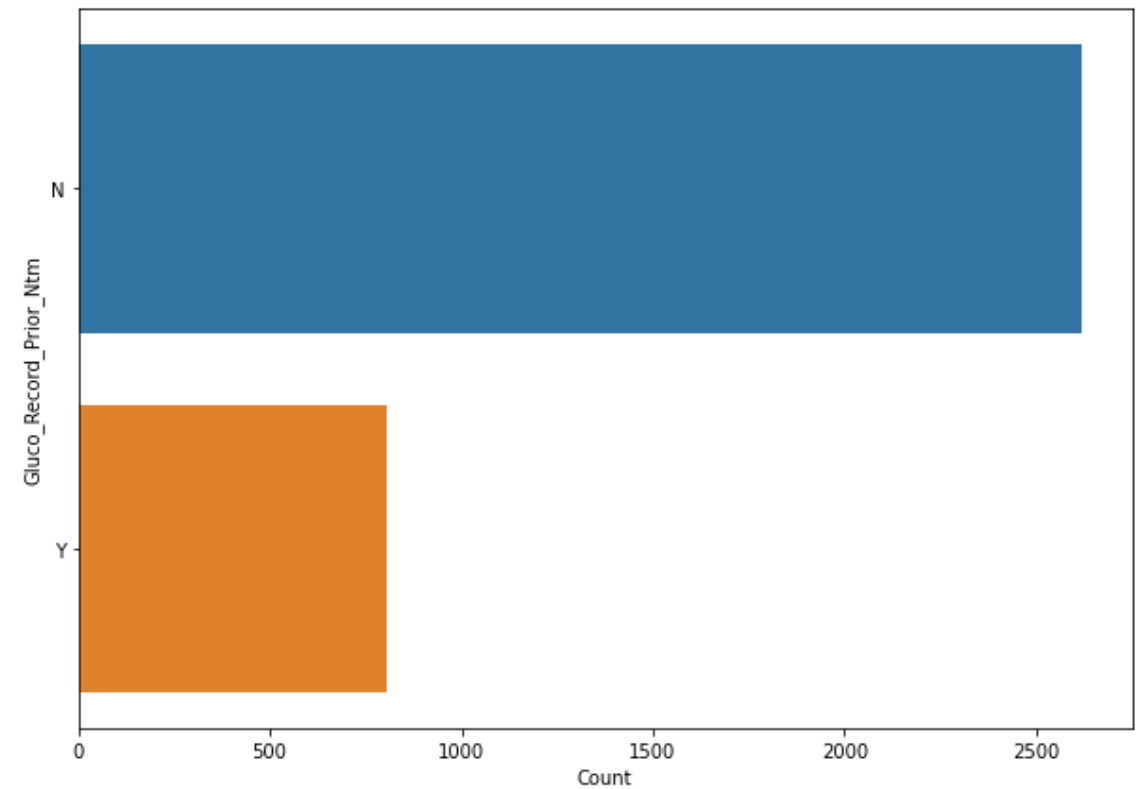
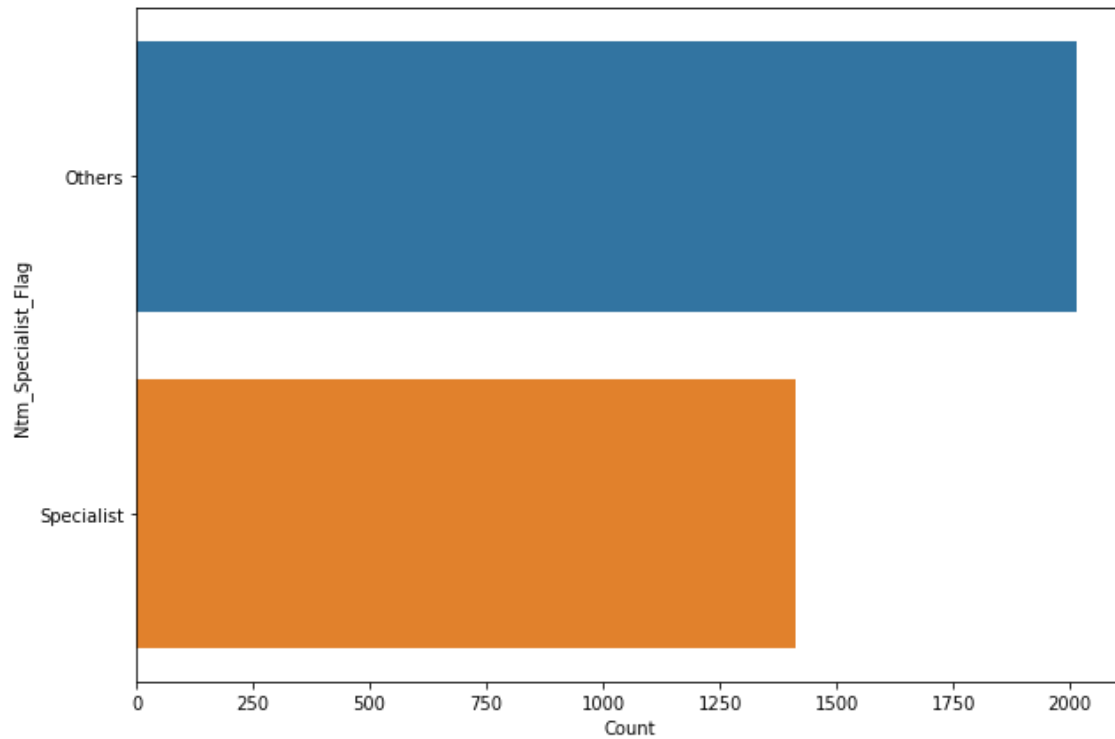
- REGION unbalanced data, greater predominance between Midwest and South.
- Unbalanced data AGE\_BUSCKET, greater predominance between >75 & 64-75.

# Num\_Speciality



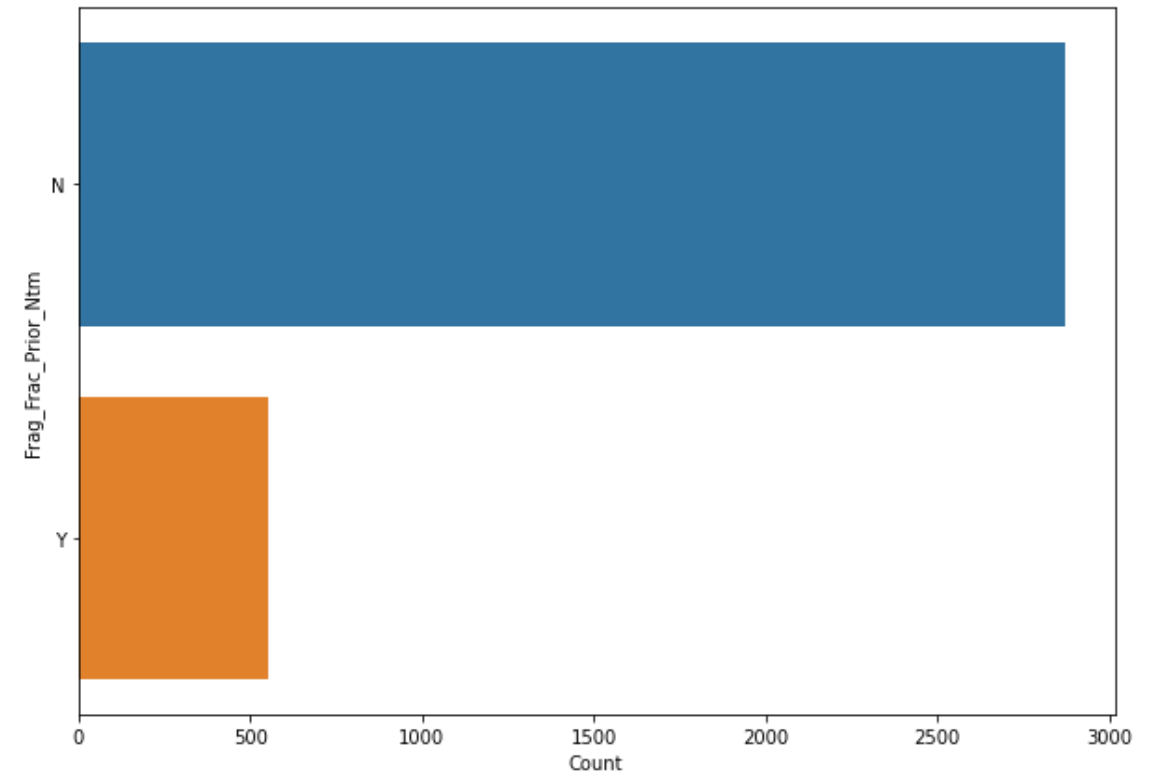
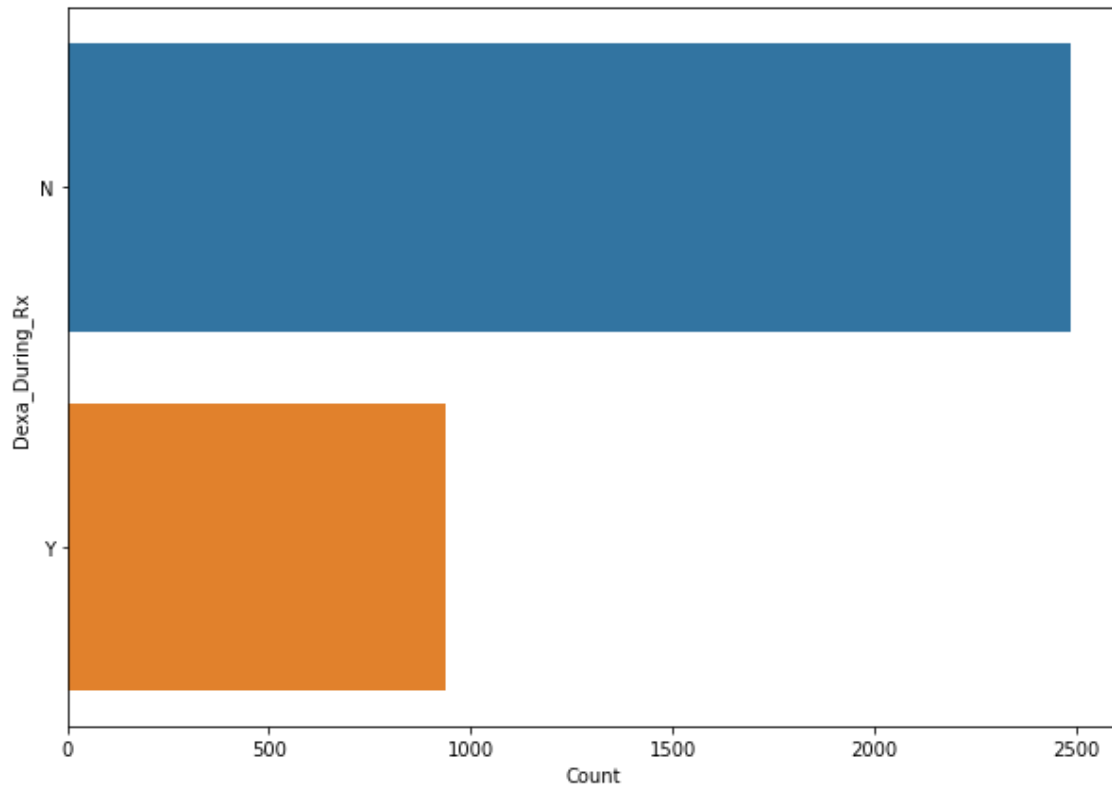
- NUM\_SPECIALITY unbalanced data, greater predominance in GENERAL PRACTITIONER.

# Num\_Speciality\_Flag/Gluco\_Record\_Prior\_Num



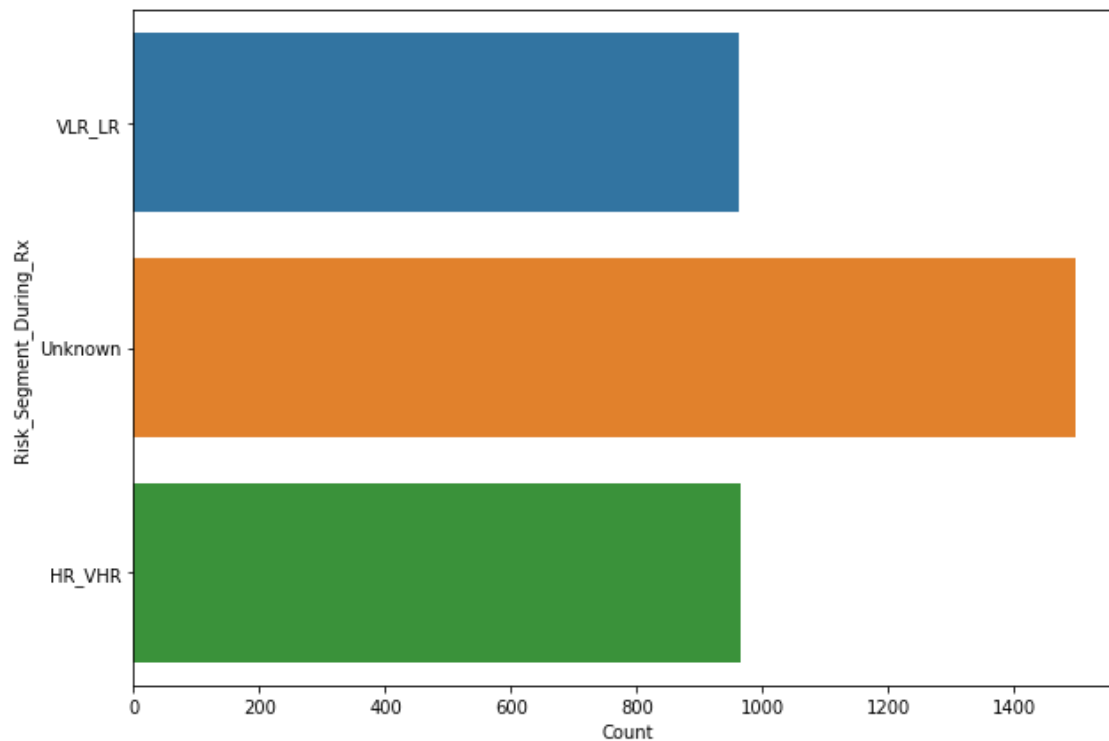
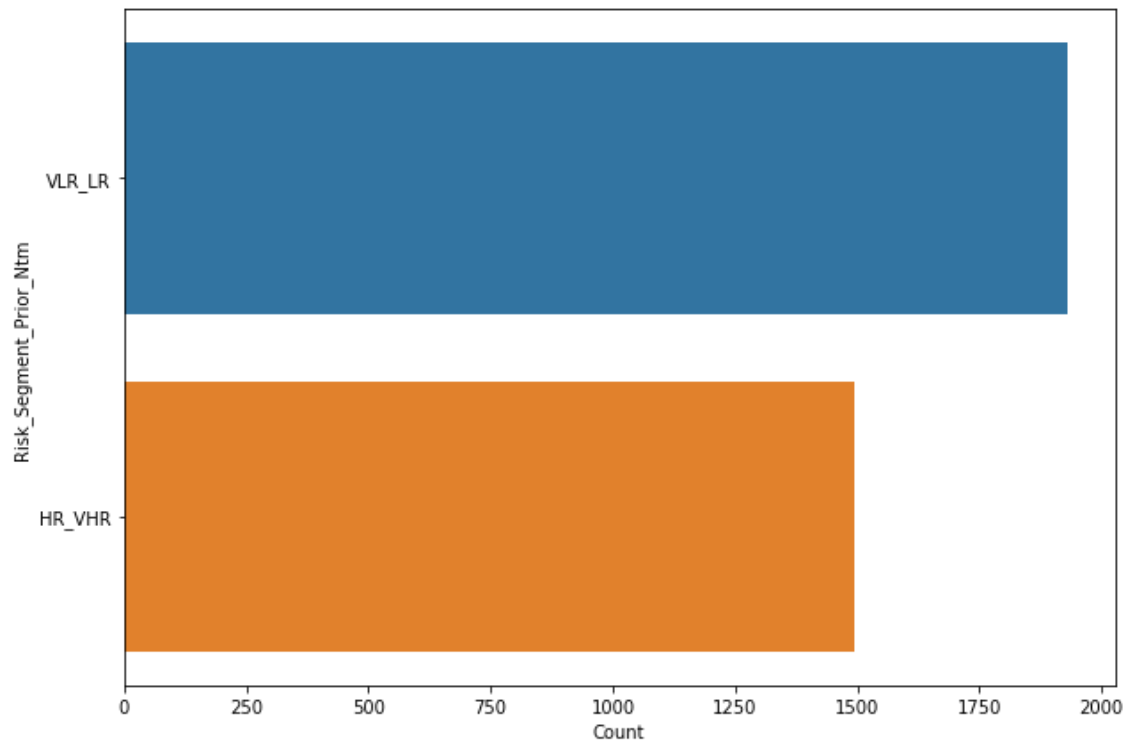
- NUM\_SPECIALIST\_FLAG balanced data.
- Unbalanced data GLUCO\_RECORD\_PROIR\_NUM, greater predominance in N.

# Dexa\_During\_Rx/Frag\_Frac\_Prior\_Num



- DEXA\_DURING\_RX unbalanced data, greater predominance between N.
- Unbalanced data FRAG\_FRAC\_PRIOR\_NUM, greater predominance in N.

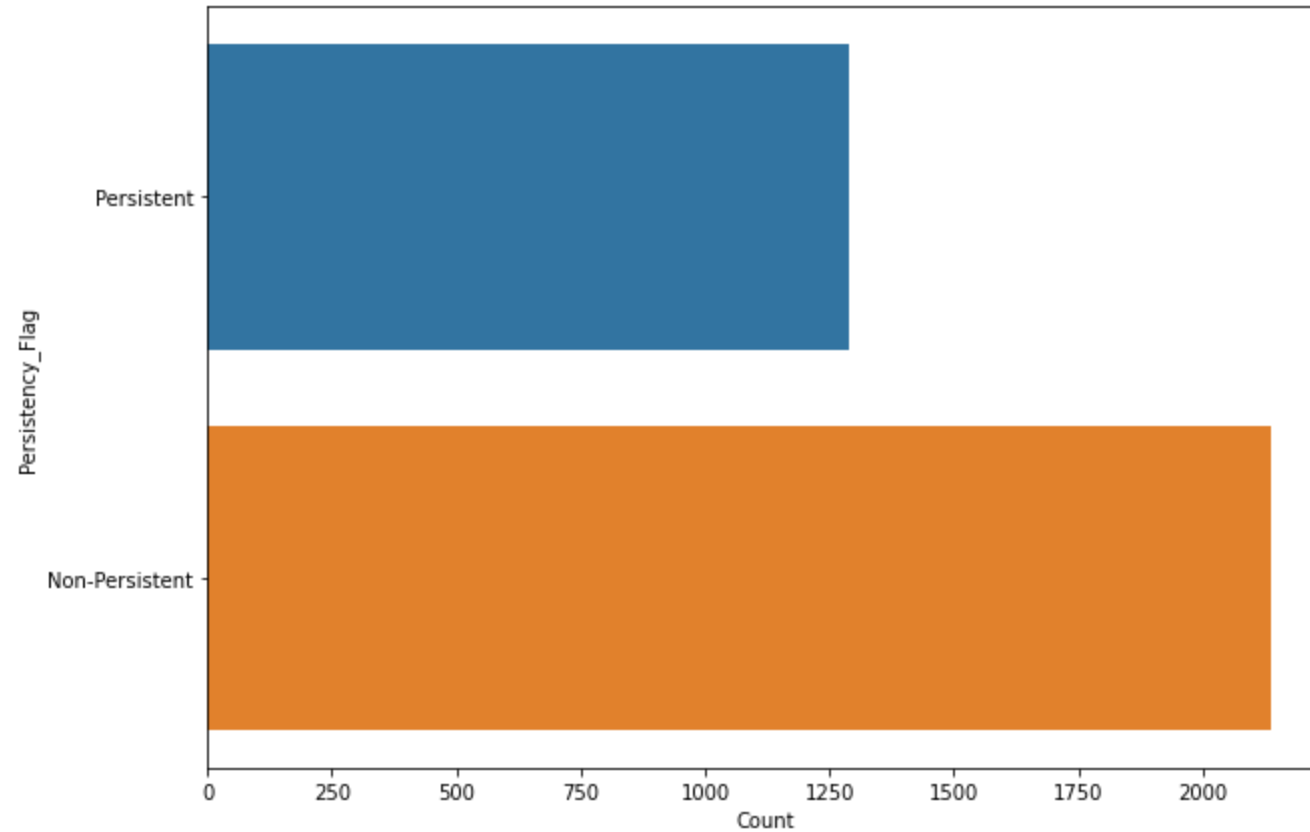
# Risk\_Segment\_Prior\_Num/Risk\_Segment\_During\_Rx



- RISK\_SEGMENT\_PRIOR\_NUM balanced data.
- Balanced data RISK\_SEGMENT\_DURING\_RX.



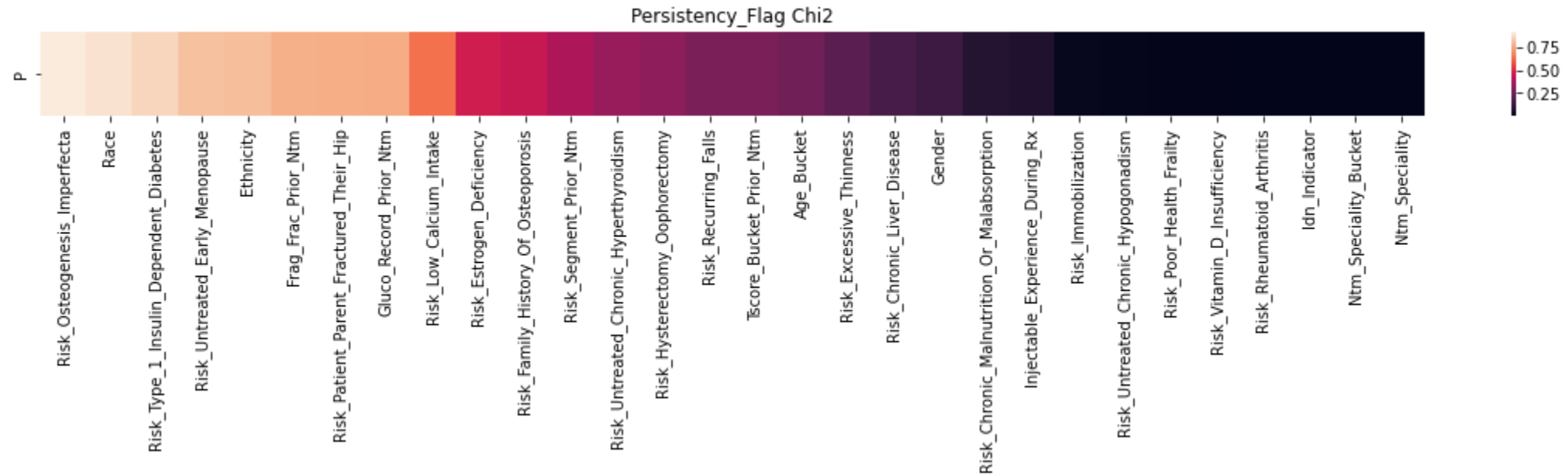
# Persistence\_Flag



- PERSISTENCY\_FLAG unbalanced data, greater predominance in Non-Persistent.

# Categorical Data Analysis Pt2

# Chi Square



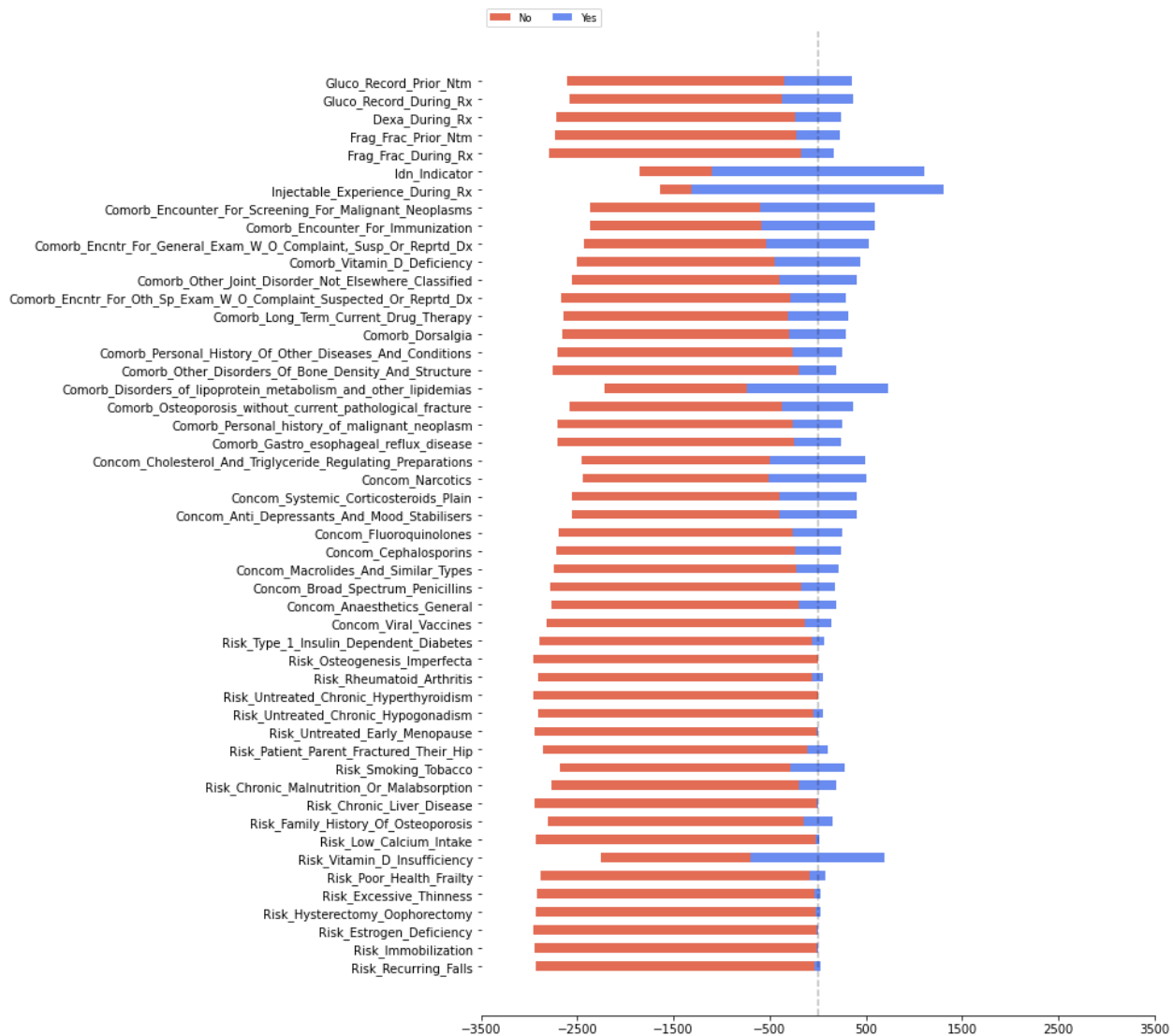
- The Chi-Square Test of Independence determines whether there is an association between categorical variables.

# Correspondence Analysis



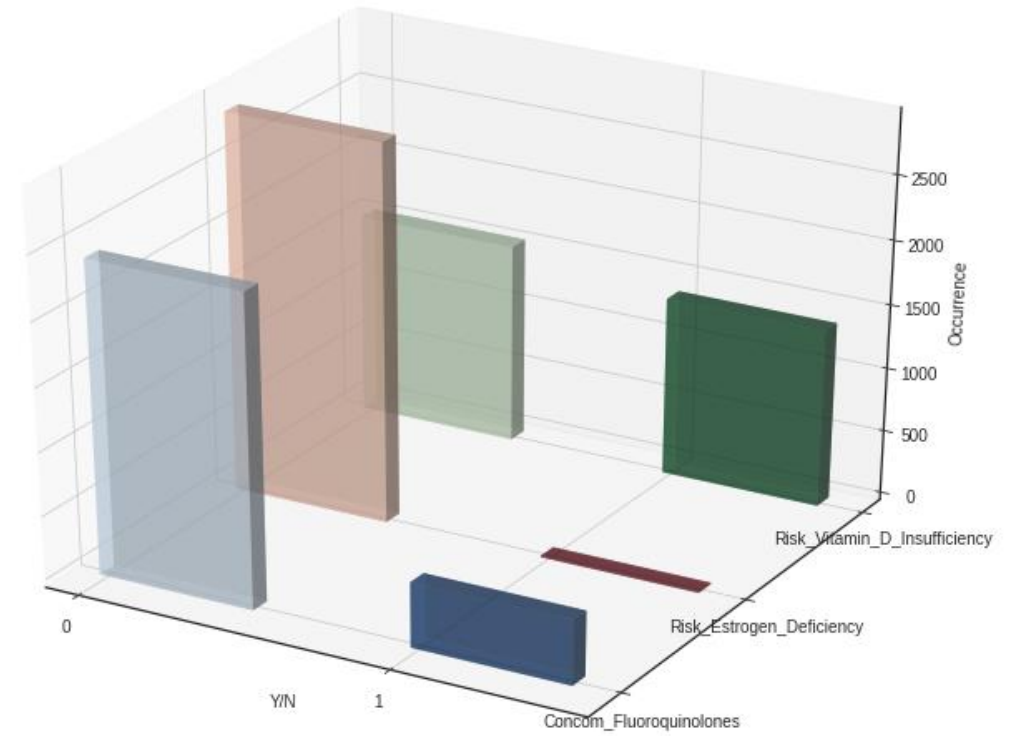
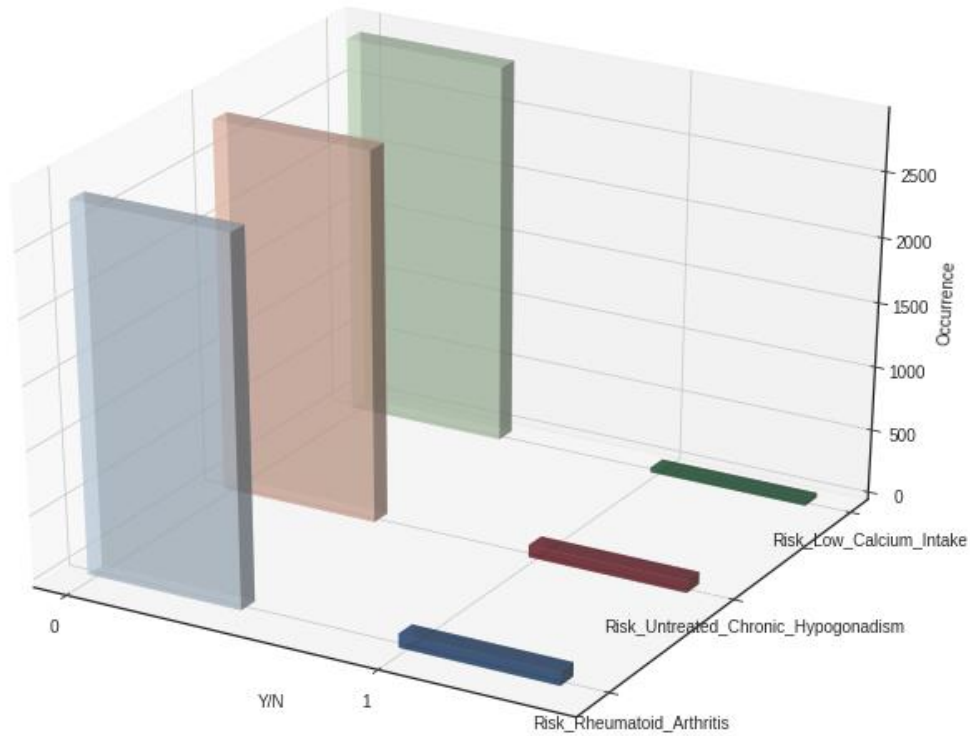
- Correspondence analysis, also called reciprocal averaging, is a useful data science visualization technique for finding out and displaying the relationship between categories.

# Tornado Chart



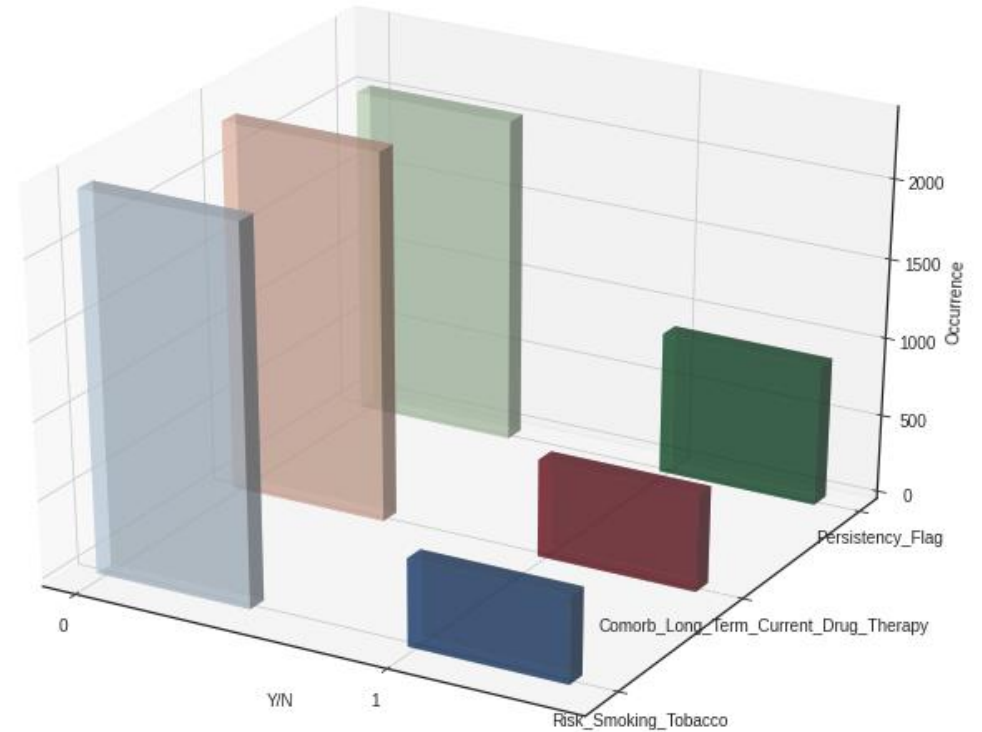
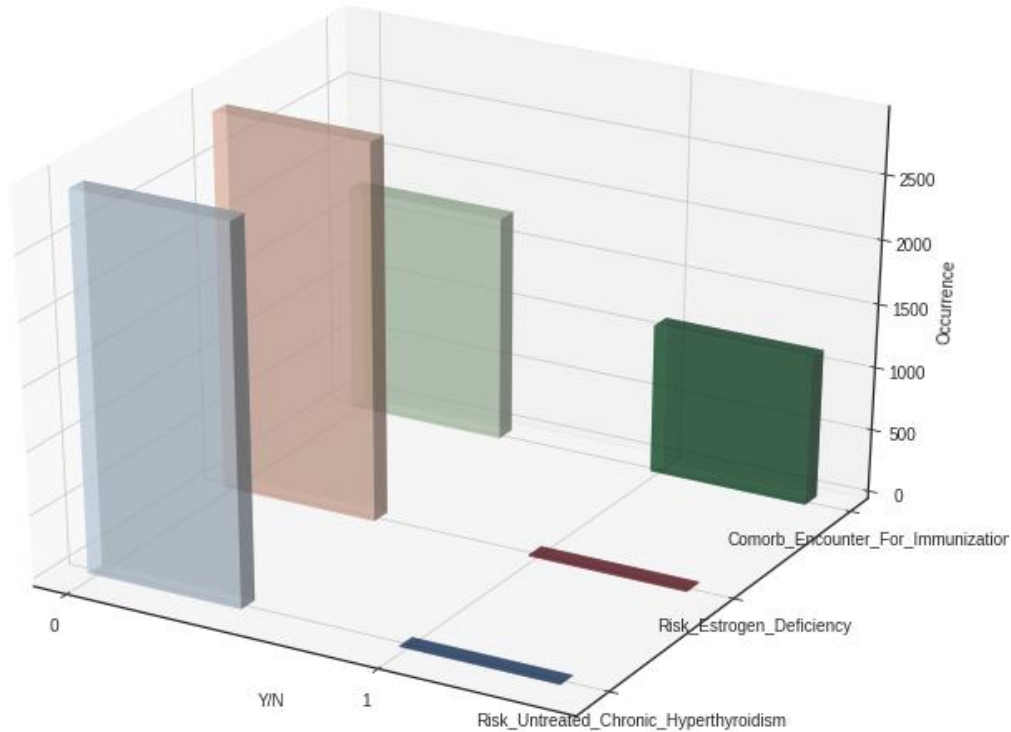
- Butterfly charts, also called Tornado or Divergent Chart, are essentially bar charts comparing two different metrics at a time.

# Crosstabulation



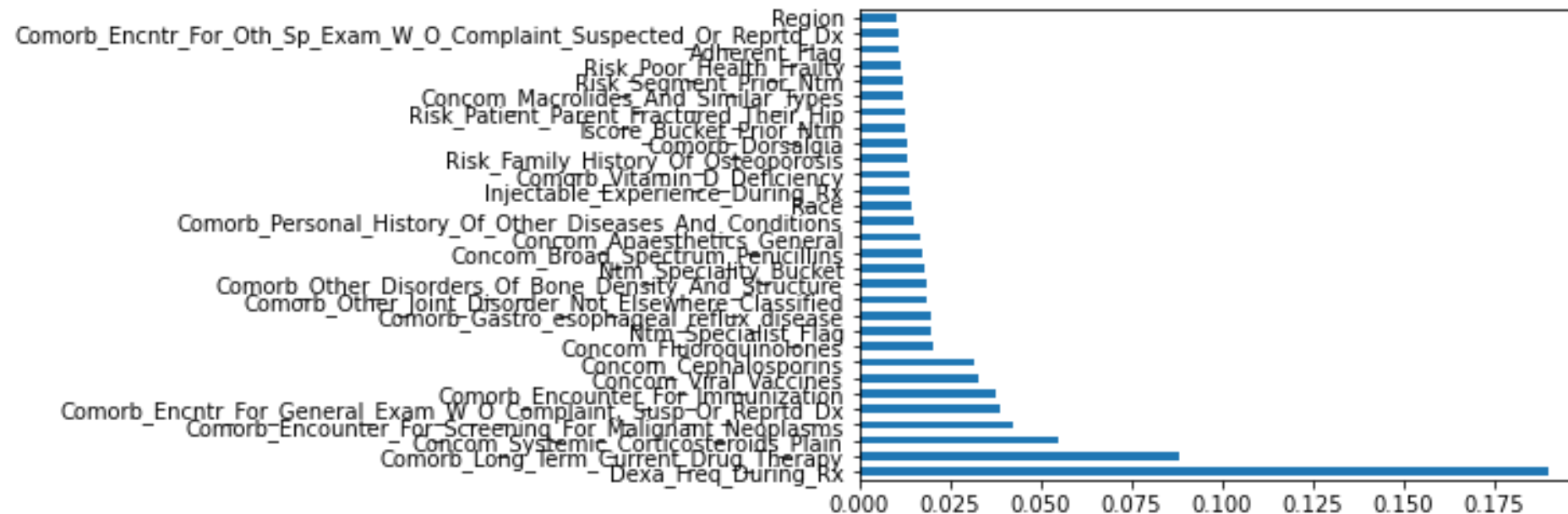
- Takes a dataframe and at least two variables as input, conducts a crosstabulation of the variables.

# Crosstabulation



- PERSISTENCY\_FLAG unbalanced data, greater predominance in Non-Persistent.

# Feature Importance



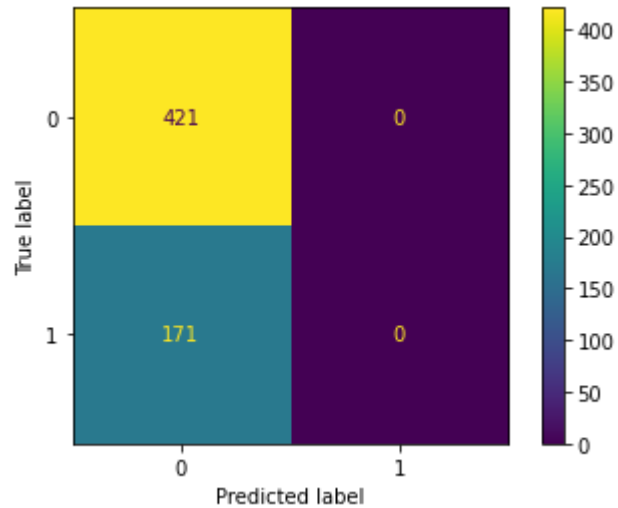
- Using xboosting helps us define the most relevant variables. In our case, 'Dexa\_Freq\_During\_Rx' and the first 20 have very high relevance values with respect to 'Persistency\_Flag'.



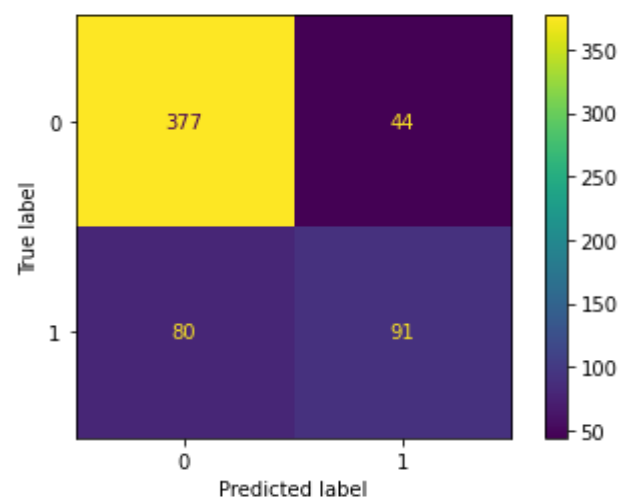
# Model Building

# Models

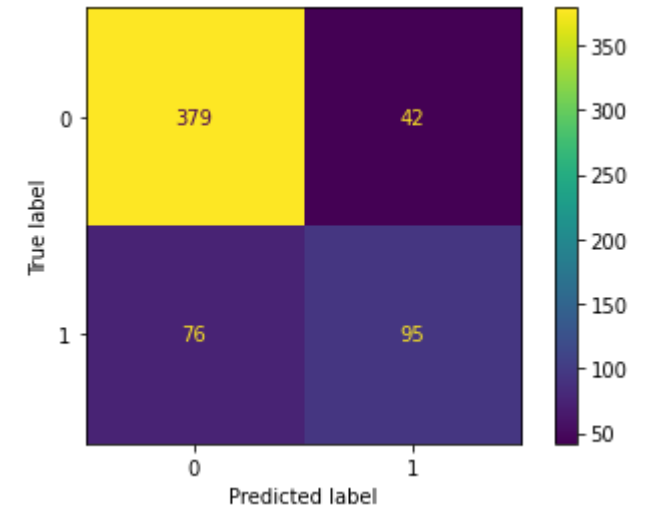
DummyClassifier, acc:71%



LogisticRegression, acc:79%

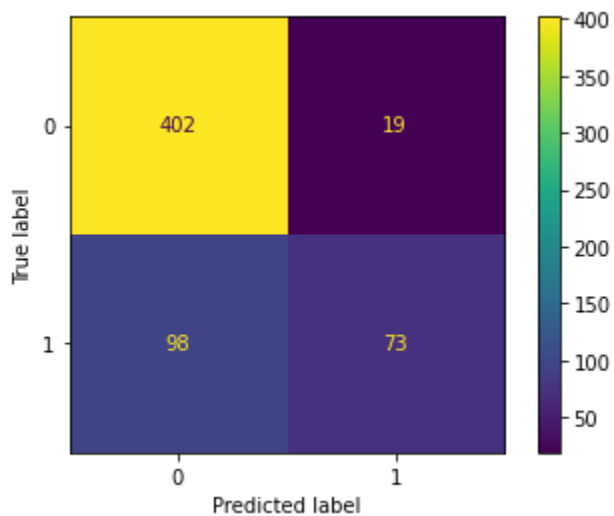


SVC, acc:80%

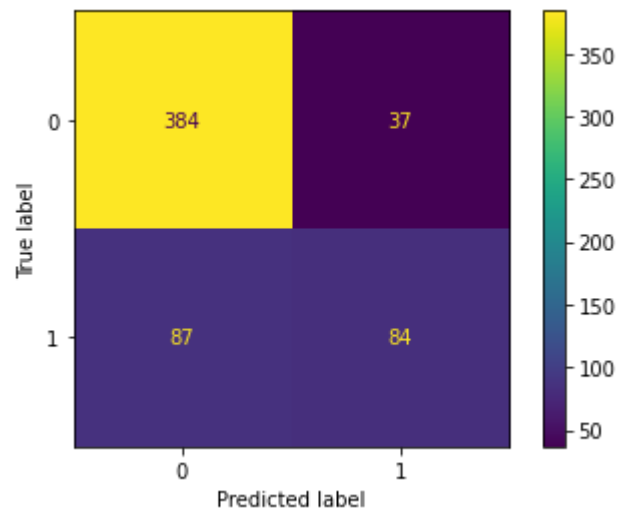


# Models

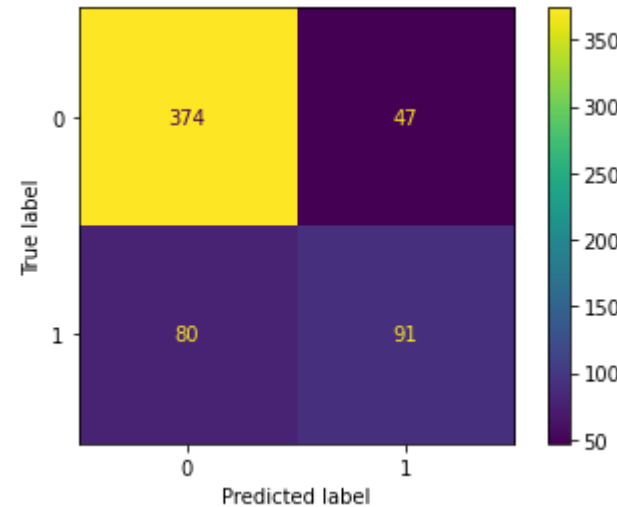
KNN, acc:80%



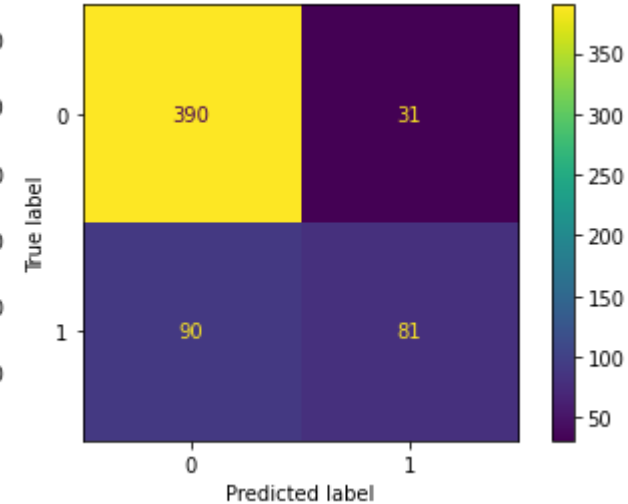
RandomForest, acc:80%



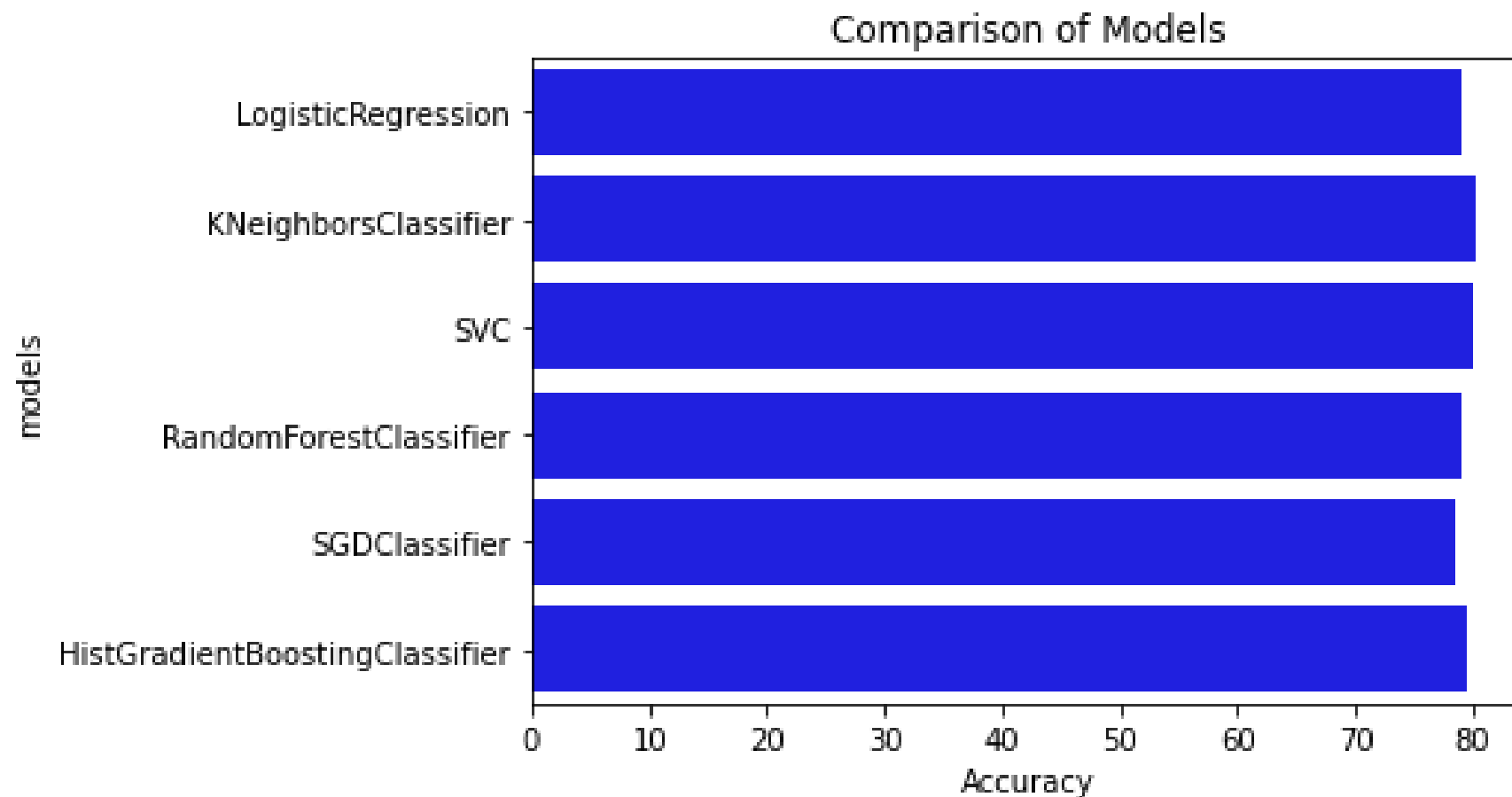
SGDClass, acc:78%



HistGB, acc:80%

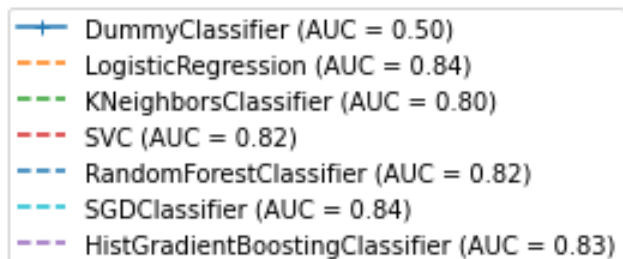
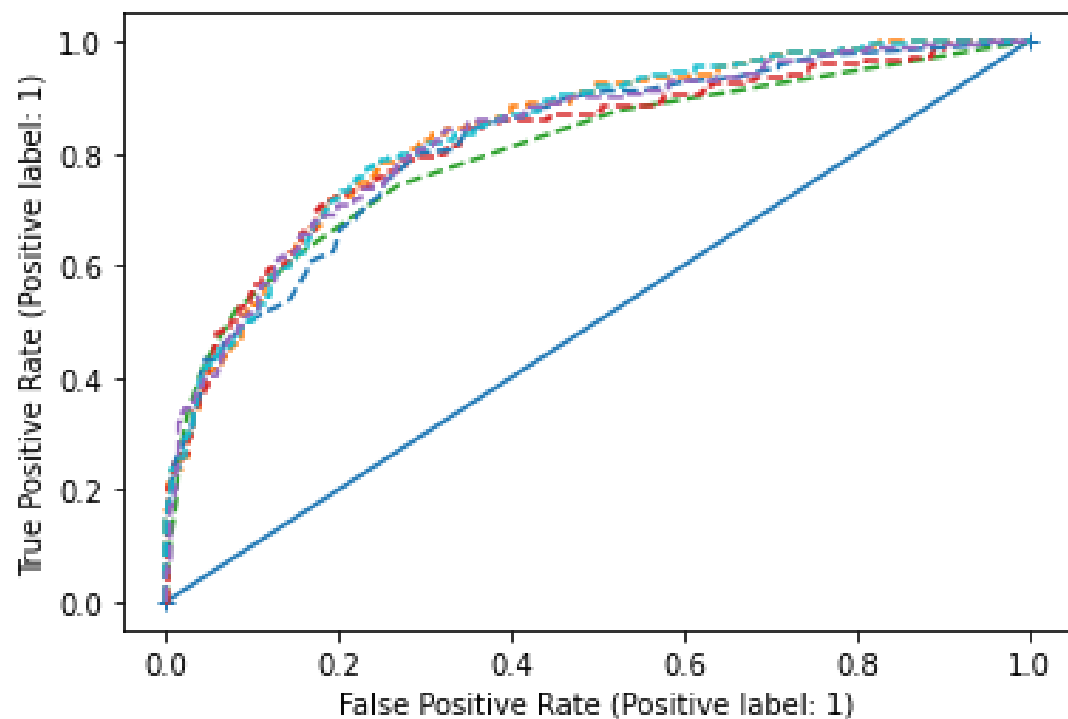


# Models

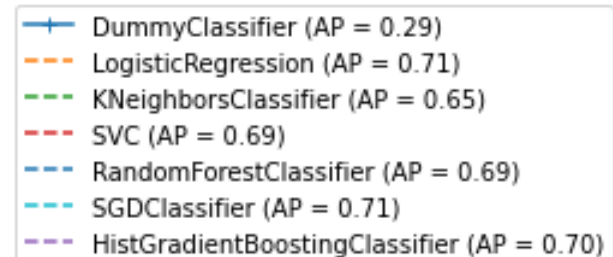
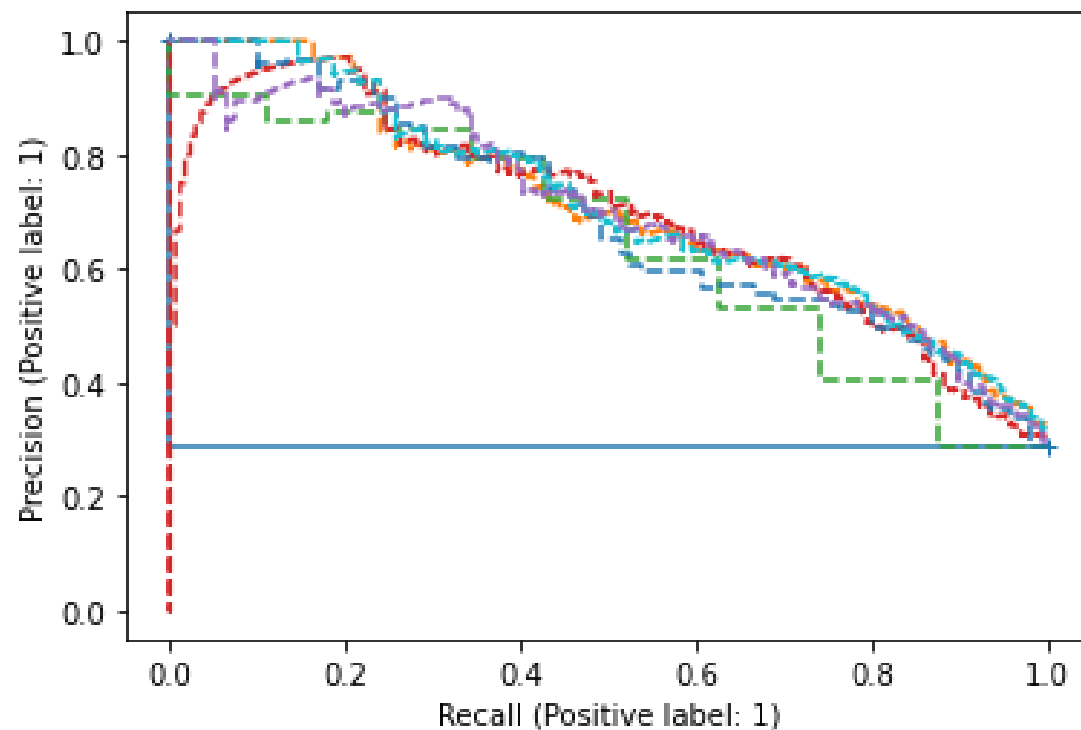


# Metrics

ROC AUC curve



Precision-recall curve



# Final Recommendations

# Final Recommendations

- It is recommended to apply classification algorithms.
- **Precision:** precision is the ratio or percentage of correct classifications of our classifier. The highest precision was obtained with SGDClassifier and LogisticRegression
- **Recall:** the recall or sensitivity of our model is the ratio of positives detected in the dataset by our classifier. The highest recall was found with SGDClassifier and LogisticRegression.
- The models used had good performance, but they can be improved with data processing, since the most relevant ones were taken.

# Thank You



Your Deep Learning Partner