



**Data Glacier**

Your Deep Learning Partner

# Exploratory Data Analysis

HEALTH CARE – DRUG PERSISTENCY

**DECEMBER 2022**

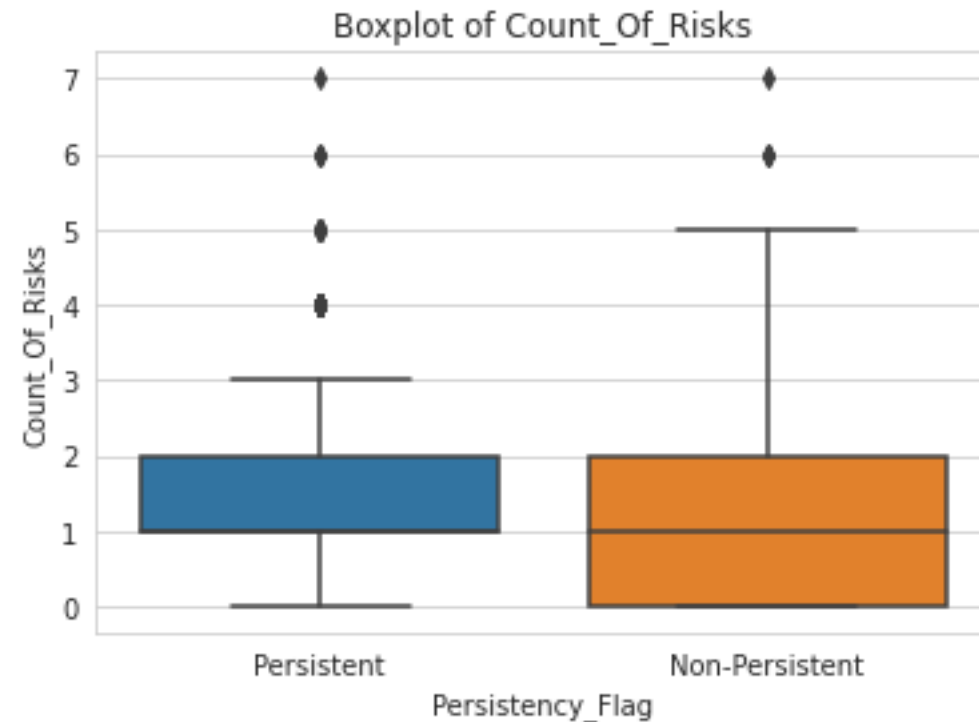
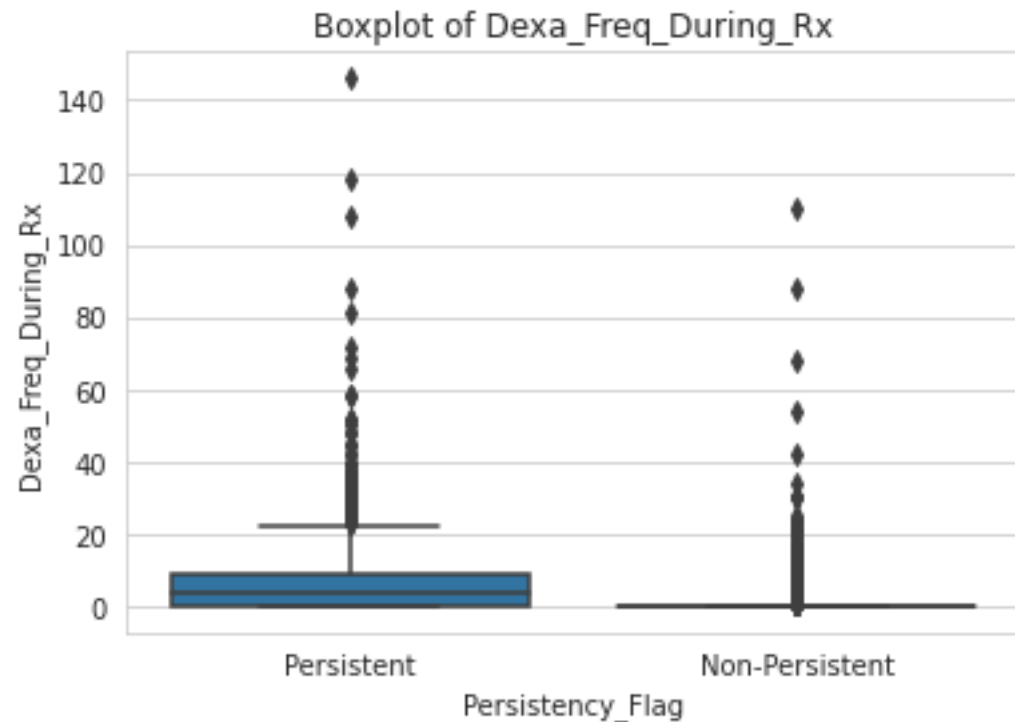
# Data Exploration

- One file used for the dataset
- 3,424 data points
- 75 features/variables (6 derived)

# Analysis of Numerical Features

- Dexa\_Freq\_During\_Rx vs Persistency\_Flag
- Count\_Of\_Risks vs Persistency\_Flag

## Dexa\_Freq\_During\_Rx and Count\_Of\_Risks Analysis



- There's a big difference in the distribution of the dexamethasone scan frequency during prescription between persistent patients and non-persistent patients.
- 50% of persistent patients have count of risks between 1 and 2, while 50% of non-persistent patients have count of risks between 0 and 2.

# Numerical Features Correlation

	Persistency_Flag	Dexa_Freq_During_Rx	Count_Of_Risks	log_Dexa	log_Count_Risks
Persistency_Flag	1.000000	0.517337	0.082431	0.517315	0.083250
Dexa_Freq_During_Rx	0.517337	1.000000	0.063414	0.990813	0.064419
Count_Of_Risks	0.082431	0.063414	1.000000	0.067388	0.966552
log_Dexa	0.517315	0.990813	0.067388	1.000000	0.067405
log_Count_Risks	0.083250	0.064419	0.966552	0.067405	1.000000

# Analysis of Categorical Features

- MUTUAL INFORMATION (MI) SCORE
- Class separation by categorical features
- Class separation by Dexa\_During\_Rx
- Class separation by Comorb\_Long\_Term\_Current\_Drug\_Therapy
- Class separation by Comorb\_Encounter\_For\_Screening\_For\_Malignant\_Neoplasms
- Class separation by Comorb\_Encounter\_For\_Immunization
- Class separation by Comorb\_Encntr\_For\_General\_Exam\_W\_O\_Complaint,\_Susp\_Or\_Reprtd\_Dx
- And go on...

# MUTUAL INFORMATION (MI) SCORE

Mutual information is the measurement of how much information one can obtain about a random variable given the value of another variable.

<b>Persistency_Flag</b>	6.623046e-01
<b>Dexa_During_Rx</b>	1.211142e-01
<b>Comorb_Long_Term_Current_Drug_Therapy</b>	6.101576e-02
<b>Comorb_Encounter_For_Screening_For_Malignant_Neoplasms</b>	5.251710e-02
<b>Comorb_Encounter_For_Immunization</b>	5.002321e-02
...	...
<b>Gluco_Record_Prior_Ntm</b>	1.659650e-05
<b>Risk_Untreated_Early_Menopause</b>	1.416782e-05
<b>Risk_Family_History_Of_Osteoporosis</b>	6.108238e-06
<b>Risk_Osteogenesis_Imperfecta</b>	3.532791e-06
<b>Frag_Frac_Prior_Ntm</b>	5.052294e-08

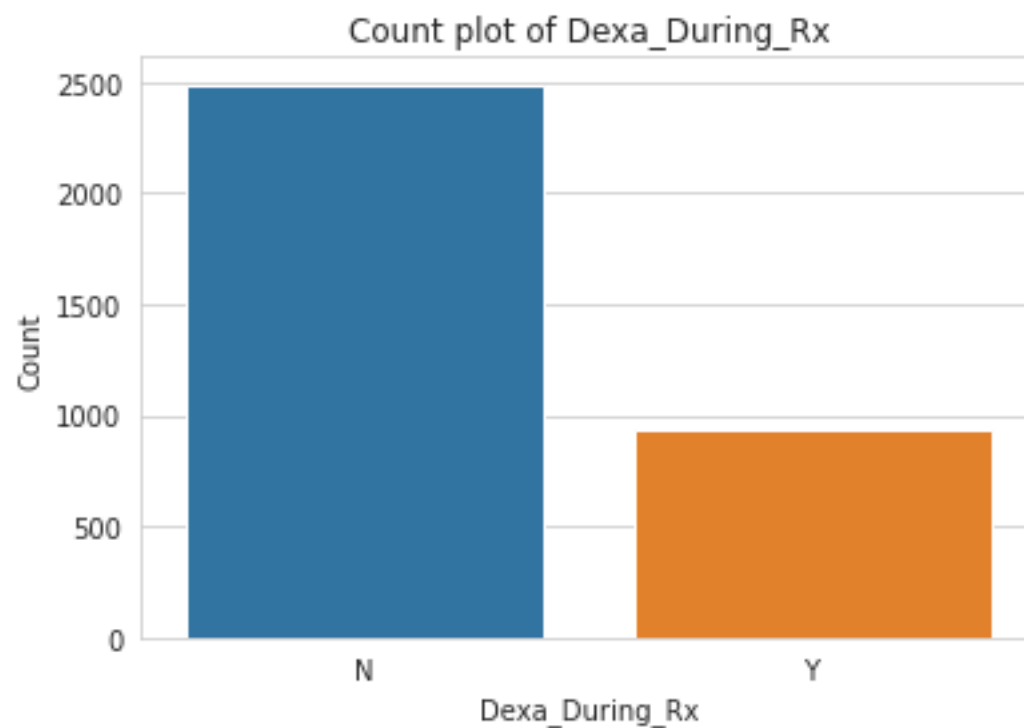
66 rows x 1 columns

# Class separation by categorical features

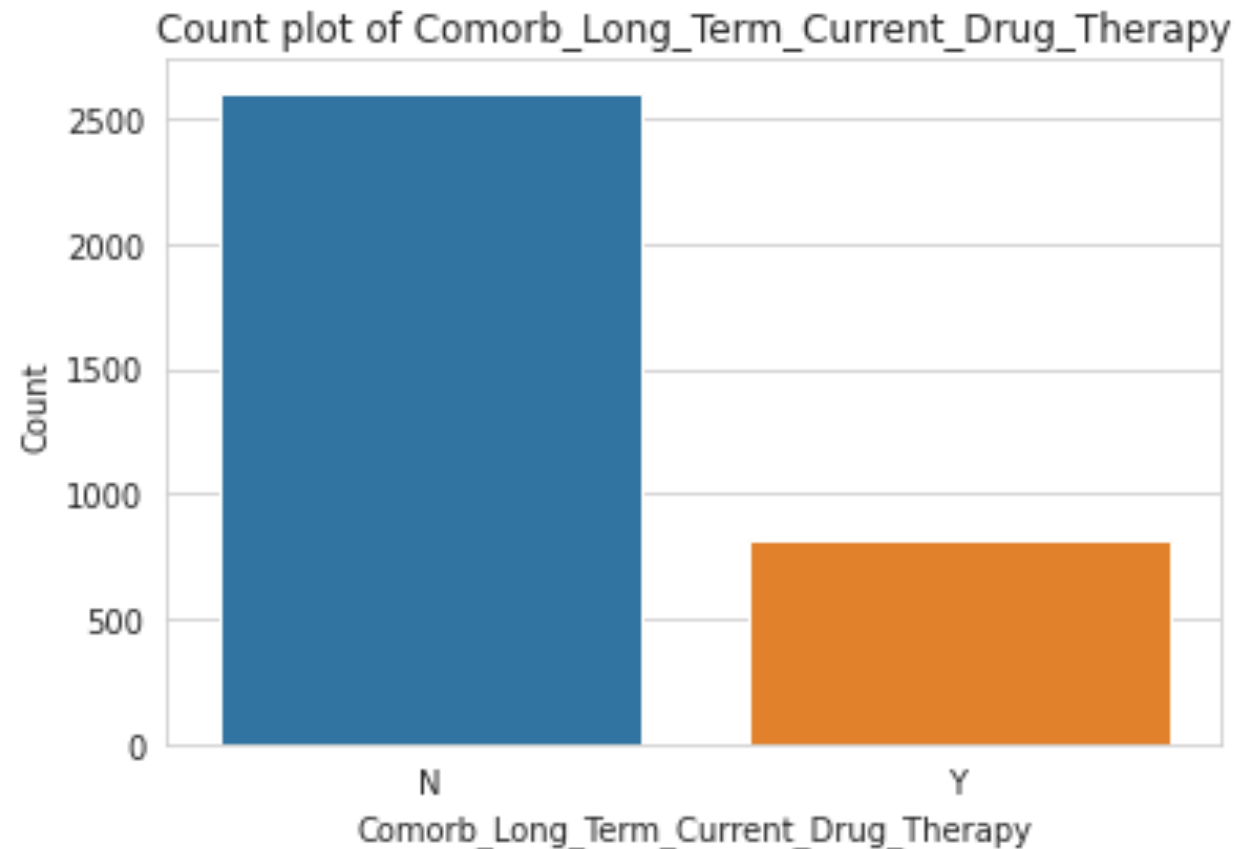
Using bar plots to visualize the separation of the target variable by the various categorical features.



# Class separation by Dexa\_During\_Rx

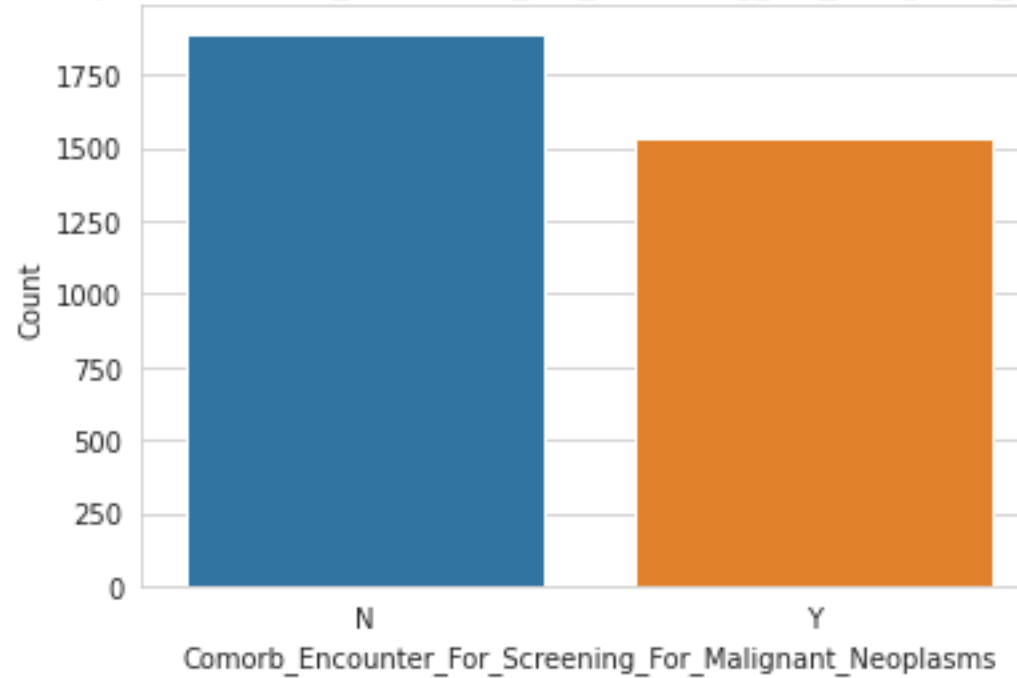


## Class separation by Comorb\_Long\_Term\_Current\_Drug\_Therapy



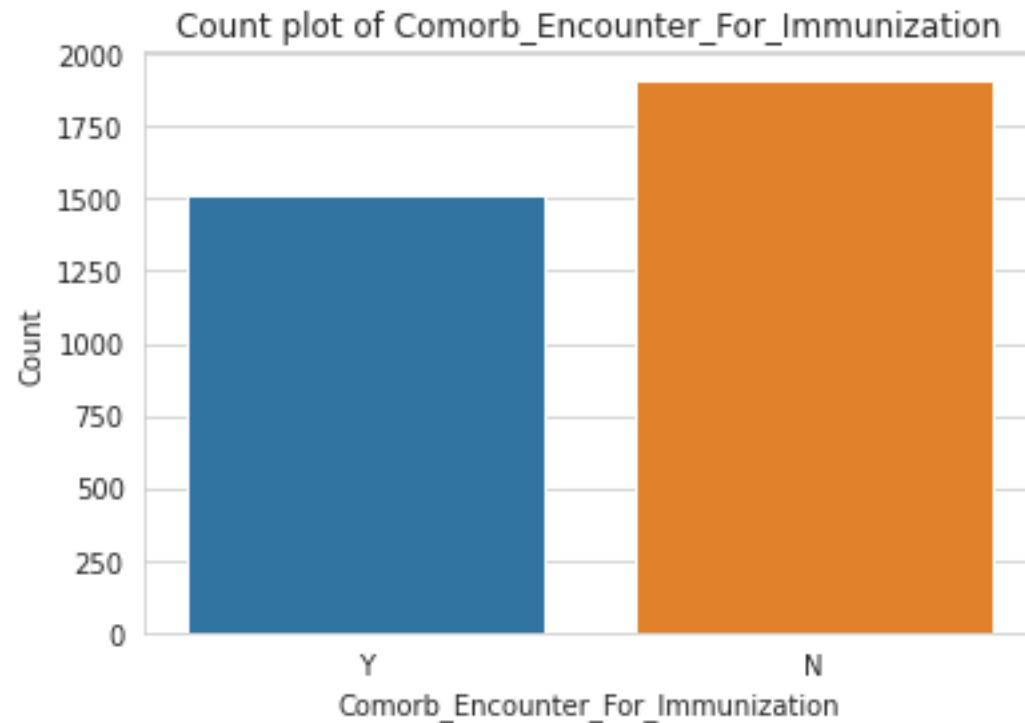
# Class separation by Comorb\_Encounter\_For\_Screening\_For\_Malignant\_Neoplasms

Count plot of Comorb\_Encounter\_For\_Screening\_For\_Malignant\_Neoplasms



- It may be a misleading data compared to other results. The results are close to each other.

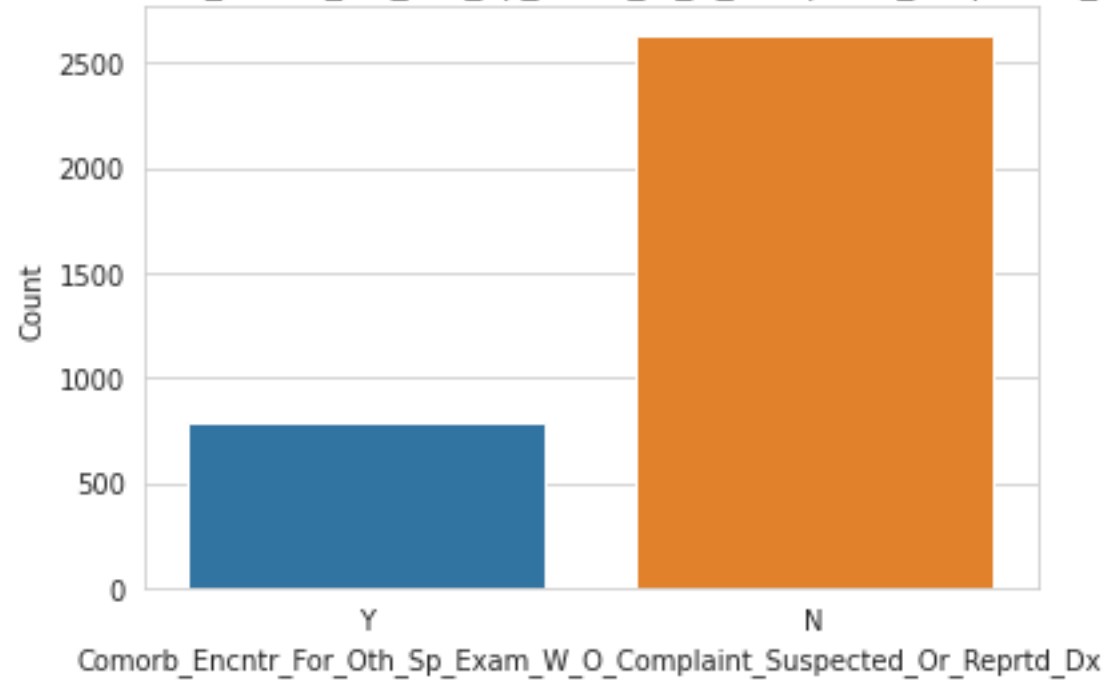
## Class separation by Comorb\_Encounter\_For\_Immunization



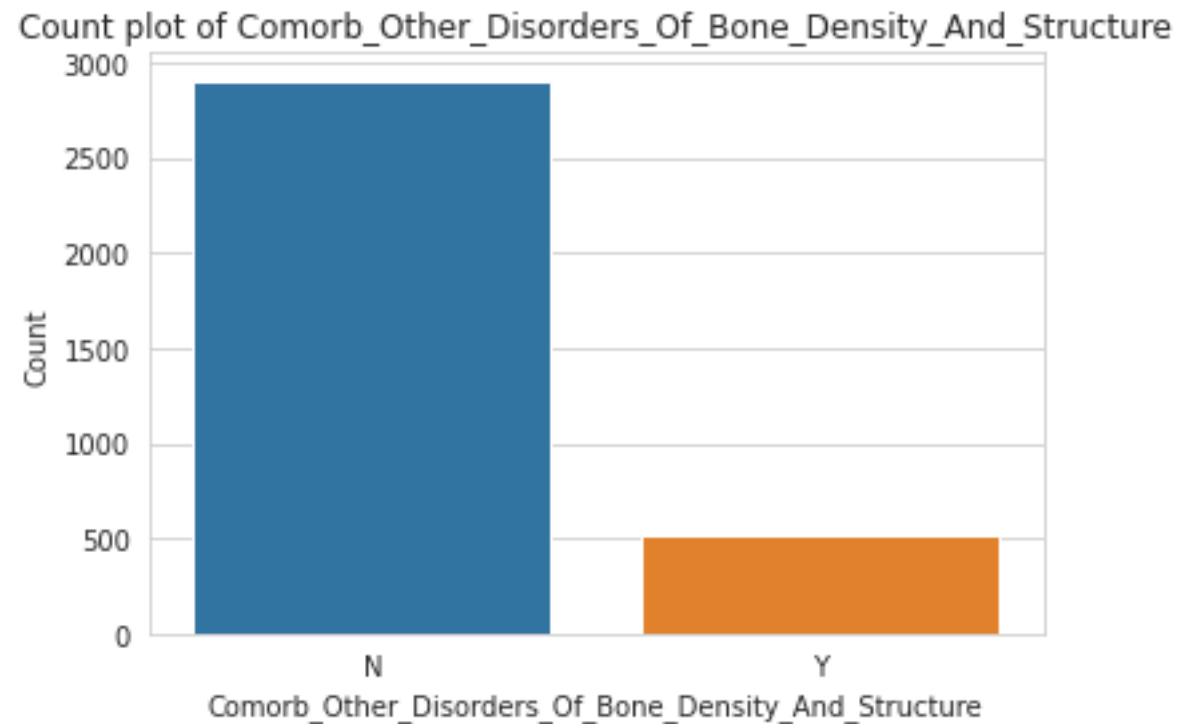
- The data are quite close.

# Class separation by Comorb\_Encntr\_For\_General\_Exam\_W\_O\_Complaint, \_Susp\_Or\_Reprtd\_Dx

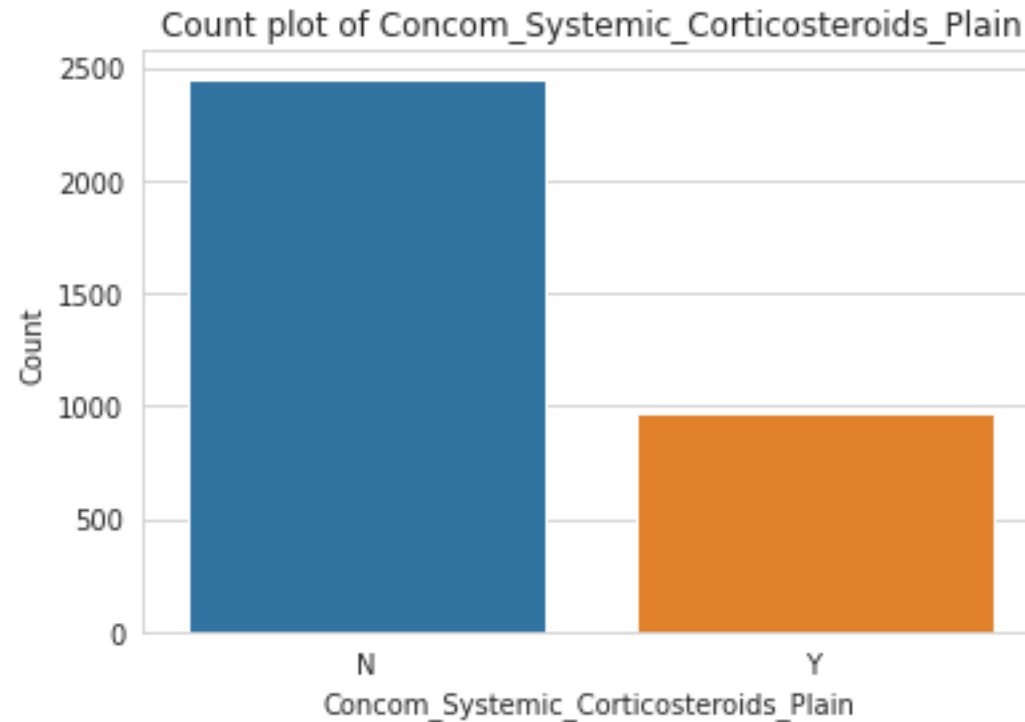
Count plot of Comorb\_Encntr\_For\_Oth\_Sp\_Exam\_W\_O\_Complaint\_Suspected\_Or\_Reprtd\_Dx



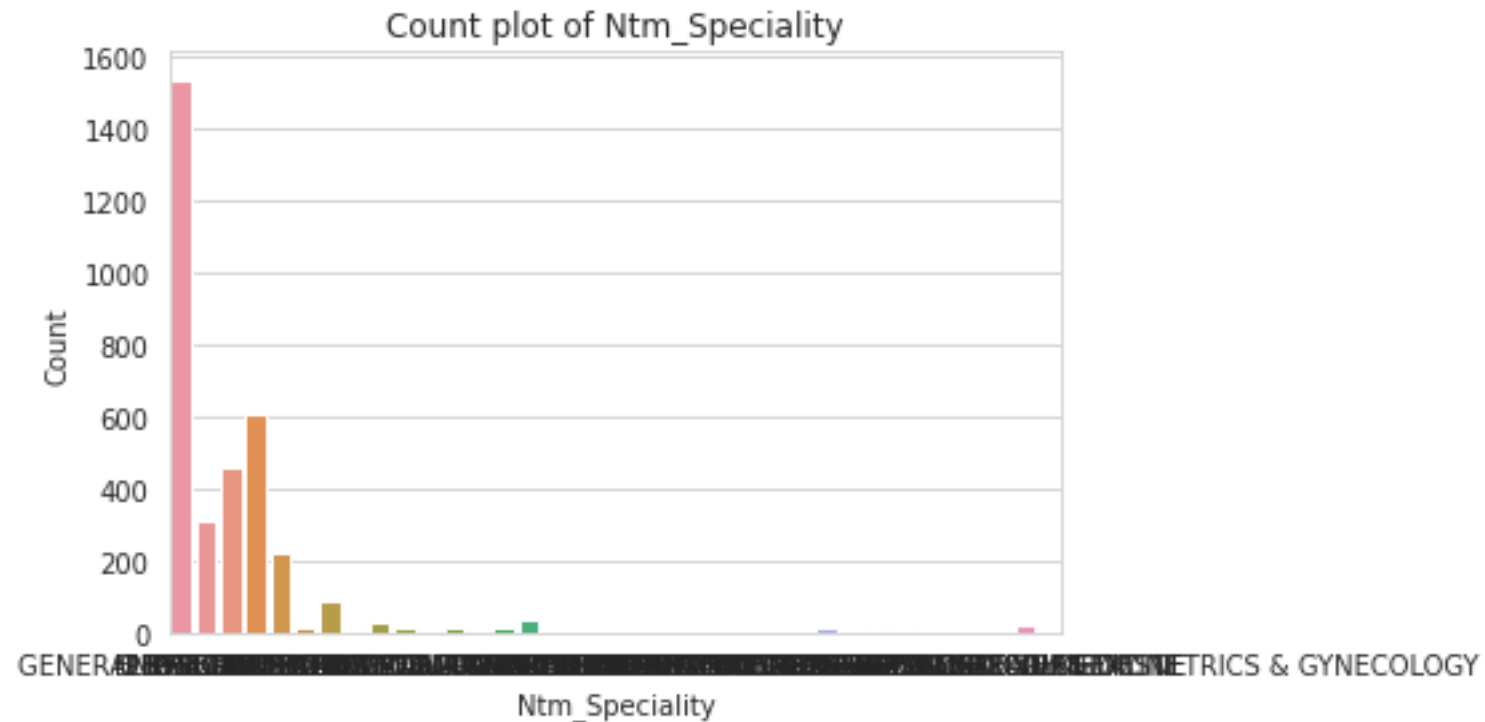
# Class separation by Comorb\_Other\_Disorders\_Of\_Bone\_Densit y\_And\_Structure



# Class separation by Concom\_Systemic\_Corticosteroids\_Plain



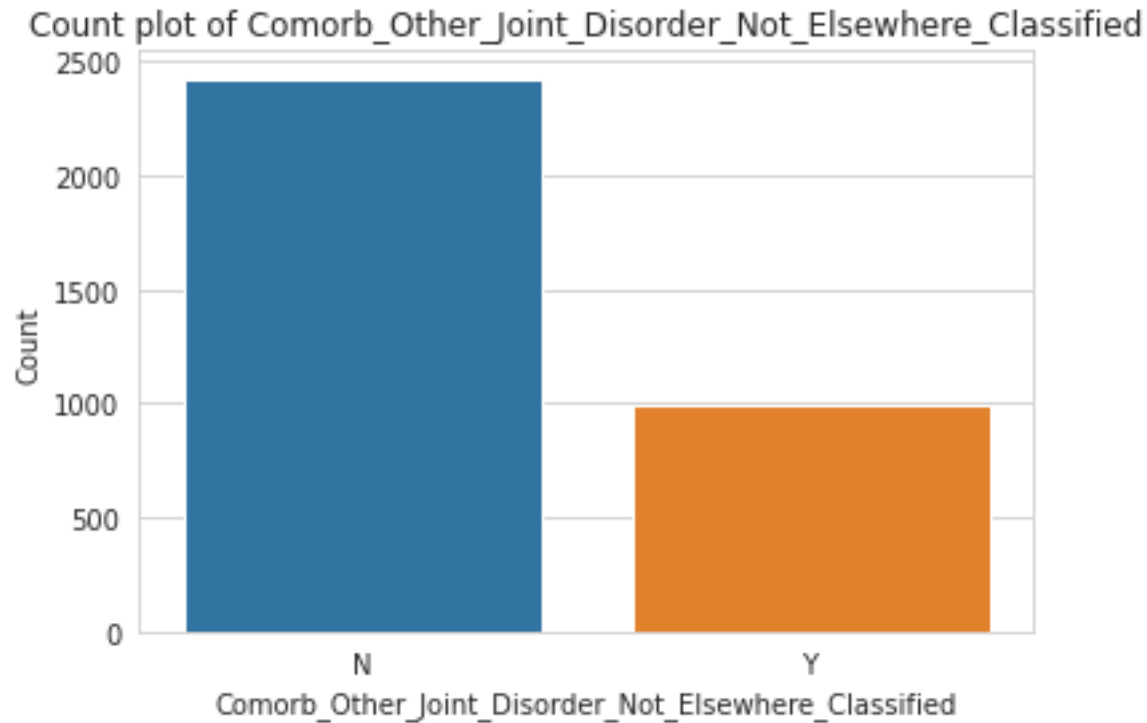
# Class separation by Ntm\_Speciality



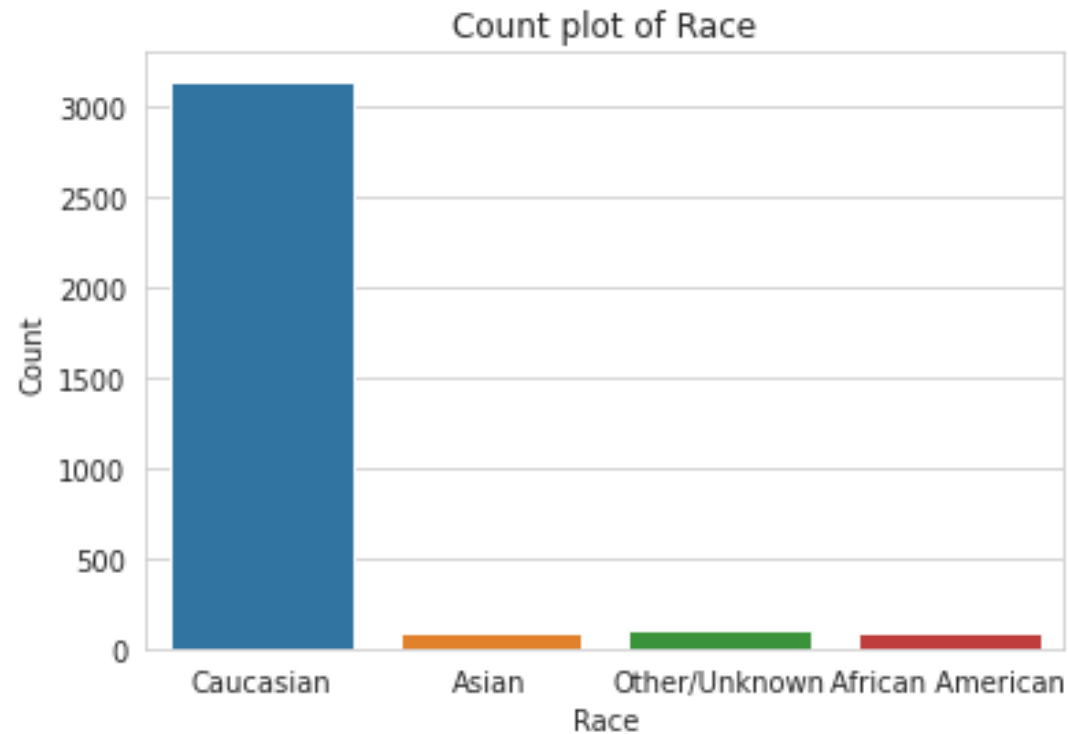
Although some of the data is spread over the whole, the majority is concentrated in one region.



# Class separation by Comorb\_Other\_Joint\_Disorder\_Not\_Elsew here\_Classified

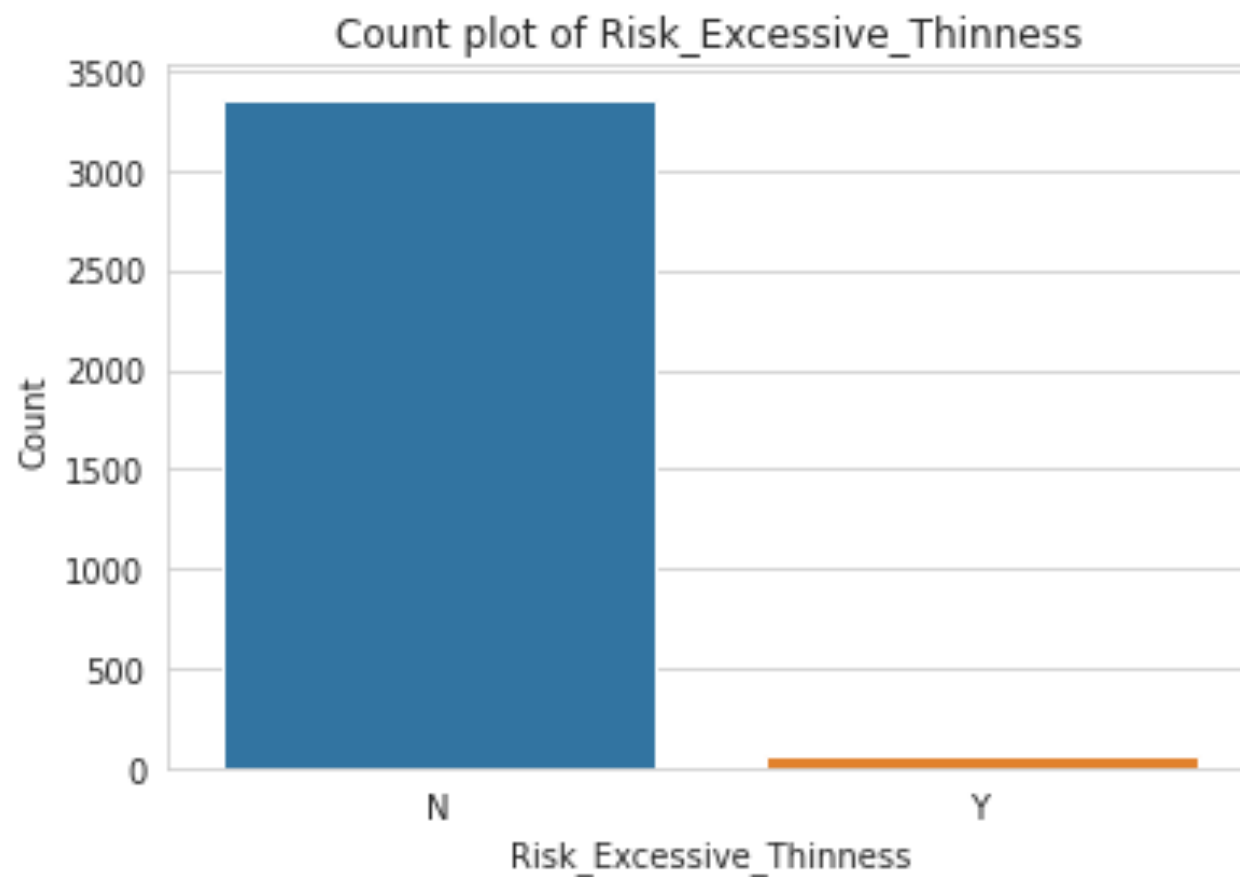


# Class separation by Race

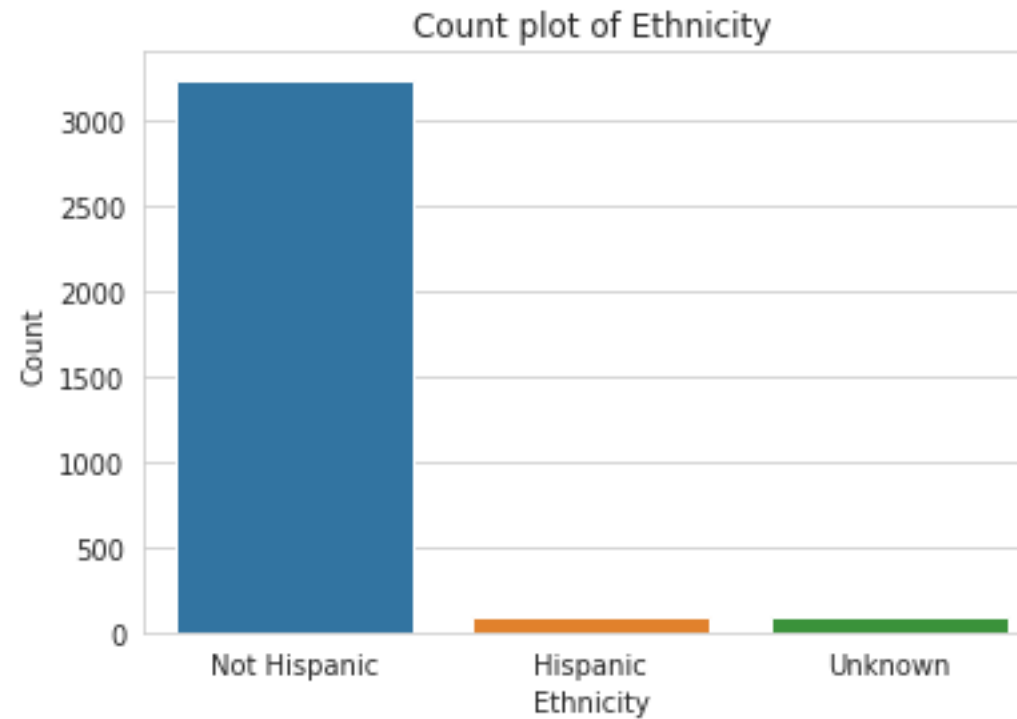


- The data is highly skewed.

# Class separation by Risk\_Excessive\_Thinness



# Class separation by Ethnicity



- The graph changed in 3 different factors.

## EDA Summary

- From the Exploratory Data Analysis done, we are able to find how the different features/variables affects drug persistency.
- The categorical variables that have higher MI scores have a greater effect in the drug persistency as compared to those that have low MI scores.

# Recommendations

For the purpose of automating the process of drug persistency identification, the following machine learning models can be used:

- **Logistic regression** – It is a type of linear model that is used for binary classification. It predicts output which is a categorical dependent variable. Such predictions are like yes or no, A or B, etc.
- **Decision tree**-They are good classifiers which are robust against outliers
- **LightGBM Classifier** – This is high-performance gradient boosting framework based on decision tree that is used for classification.

# Thank You