

# Data Science Persistency of a Drug Project

Exploratory Data Analysis

<b>Name</b>	<b>Email</b>	<b>Country</b>	<b>College</b>	<b>Specialization</b>
<b>Han-Fu Lin</b>	<b>hanfu.lin@mail.utoronto.ca</b>	<b>Canada</b>	<b>University of Totonto</b>	<b>Data science</b>
<b>Aişe Refia Yılmaz</b>	<b>aiserefia.yilmaz@gmail.com</b>	<b>Turkey</b>	<b>Beykent University</b>	<b>Data science</b>
<b>Anıl Ilgın Büyüksaraç</b>	<b>Anililginb@gmail.com</b>	<b>Turkey</b>	<b>Y.t.u</b>	<b>Data science</b>



# Includes

- Data Analysis Approach
- Problem Understanding
- Data Preparation
- Exploratory Data Analysis
- EDA Summary
- Model Building
- Final Recommendation
- Conclusion

# Data Analysis Approach

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.

- Explore and Understand the data
- Prepare and clean the data

Analyze the data

Find the features of drug persistency

Give recommendations



# Problem Understanding

Since the data is very large, we know that ABC pharmaceutical company wants to run the whole process faster and to get precise results.

We have considered and measured many factors that do or do not affect permanence.

# DATA PREPARATION

- We used pandas library and Google Colab Notebook
- Data was cleaned and prepared for analyzation.
- Approaches:
  - From null and missing values.
  - Virtualization of data.

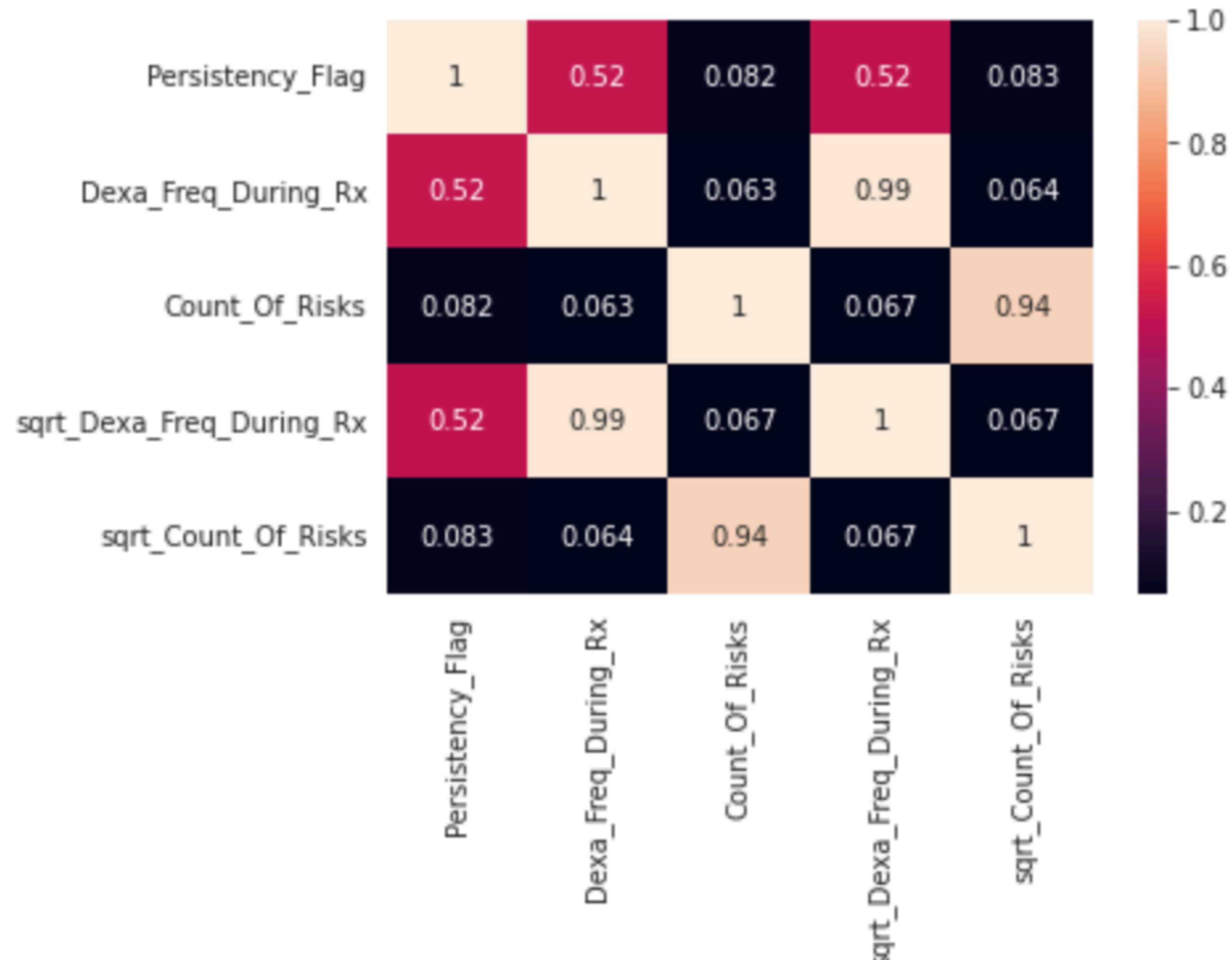


# EXPLORATORY DATA ANALYSIS

- The demographic data characteristics, gender, race, ethnicity, age, region etc are analyzed.
- Many graphics used for visualization.
- Some analyzes achieved very high accuracy rates.

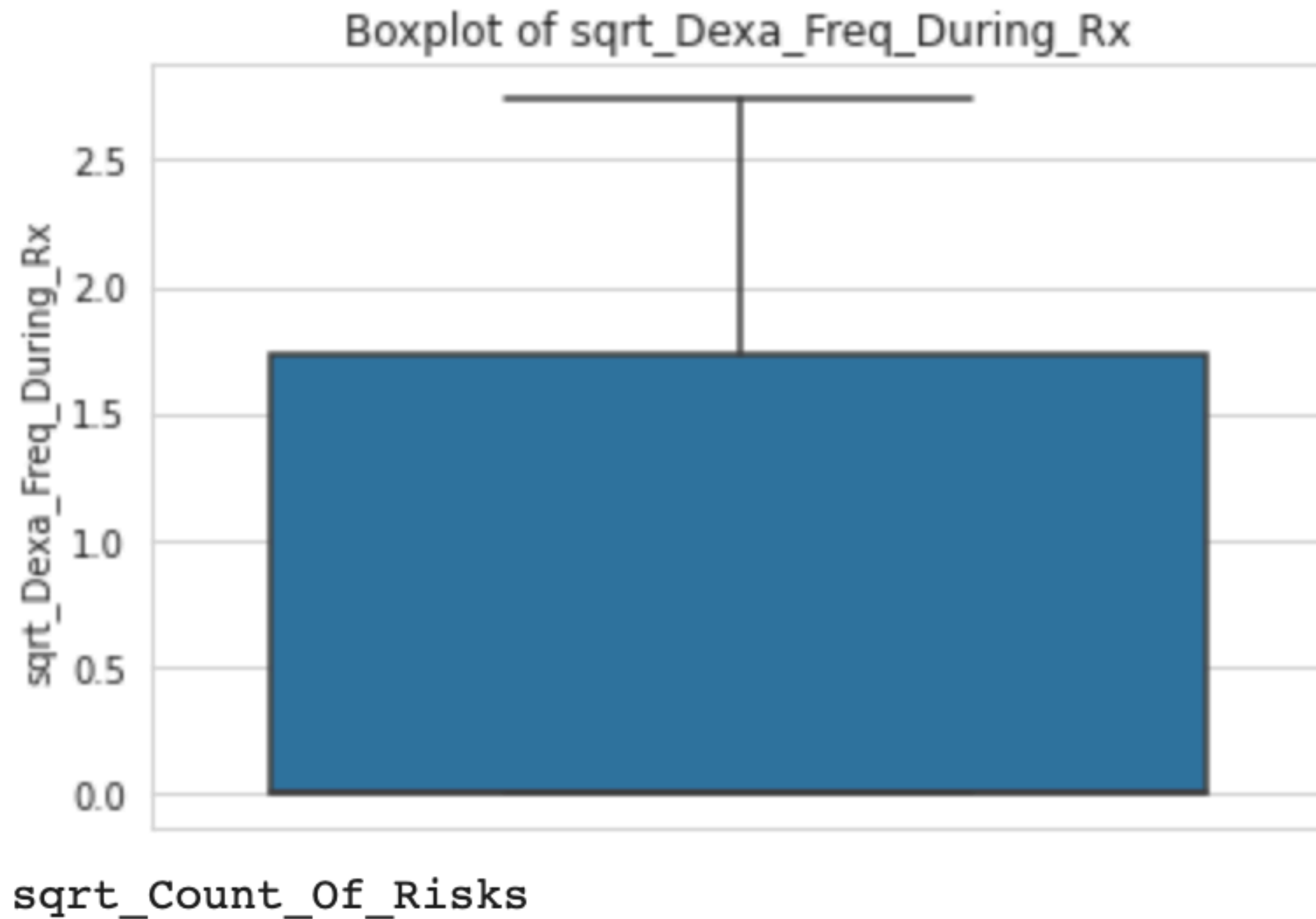
As a result of the analyzes made, we got different outputs every week and saved them.

# Heatmap

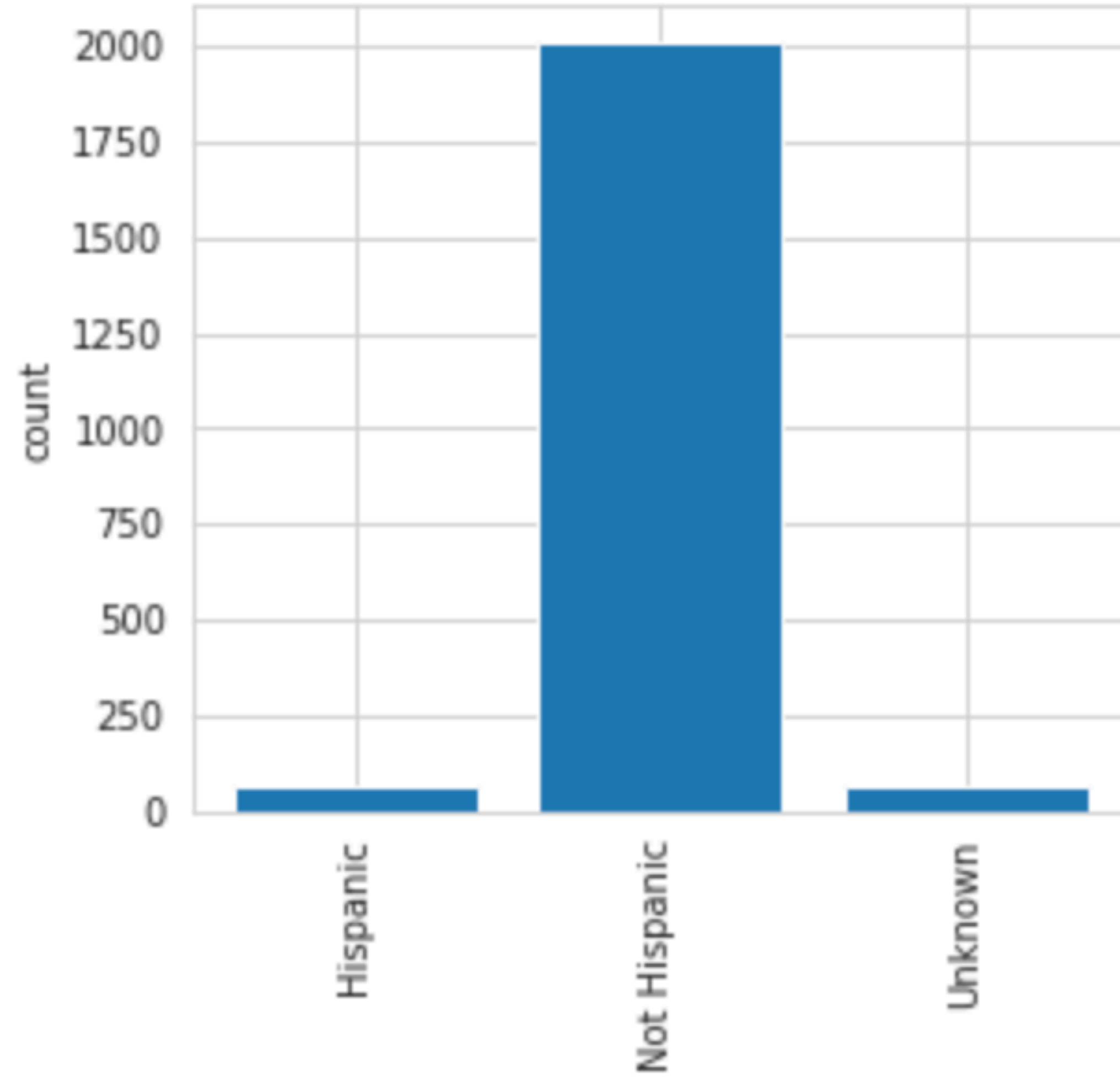




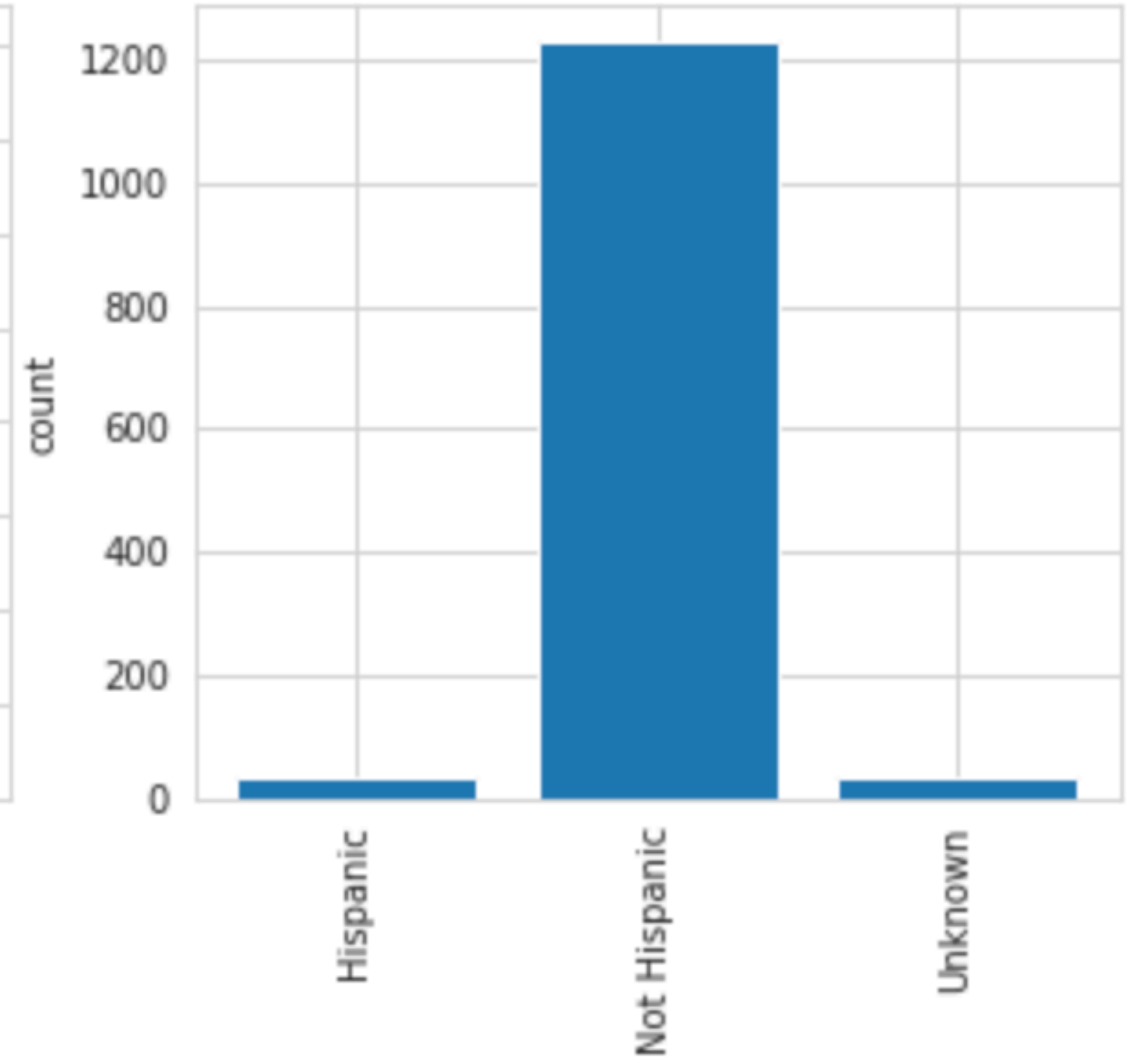
# Outlier Test



Counts for Ethnicity  
not persistent

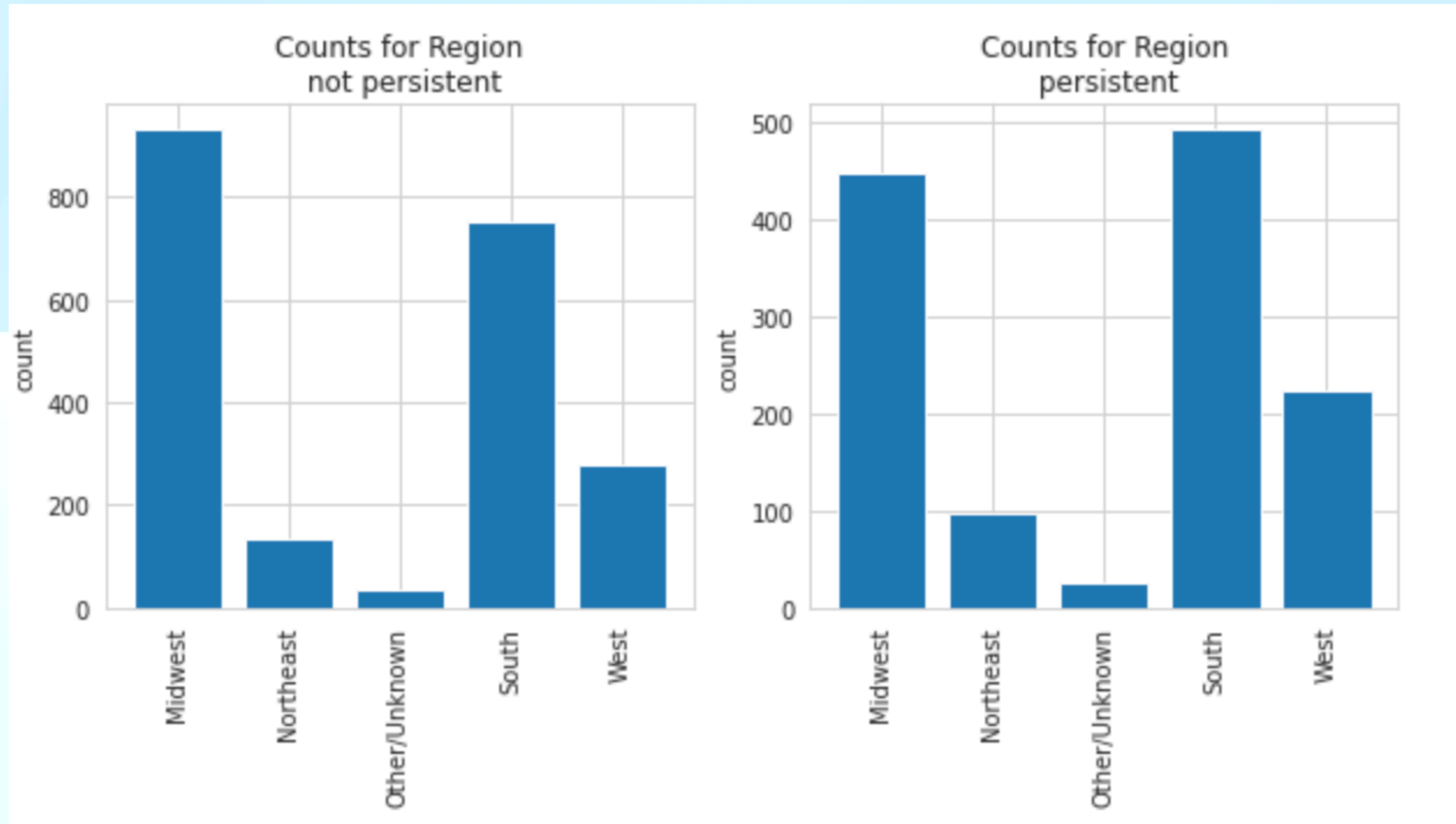


Counts for Ethnicity  
persistent

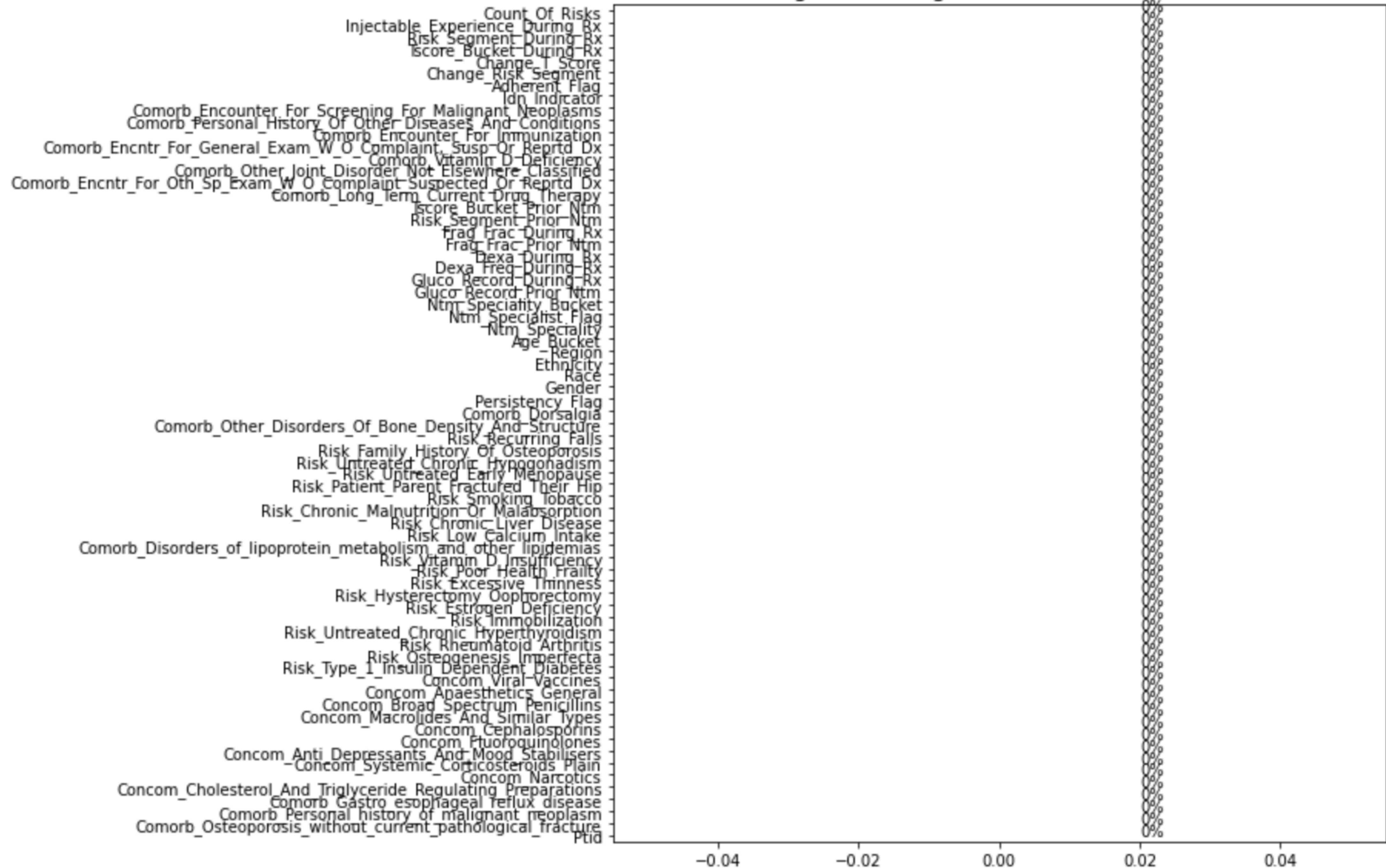




# Persistency and Non-Persistent



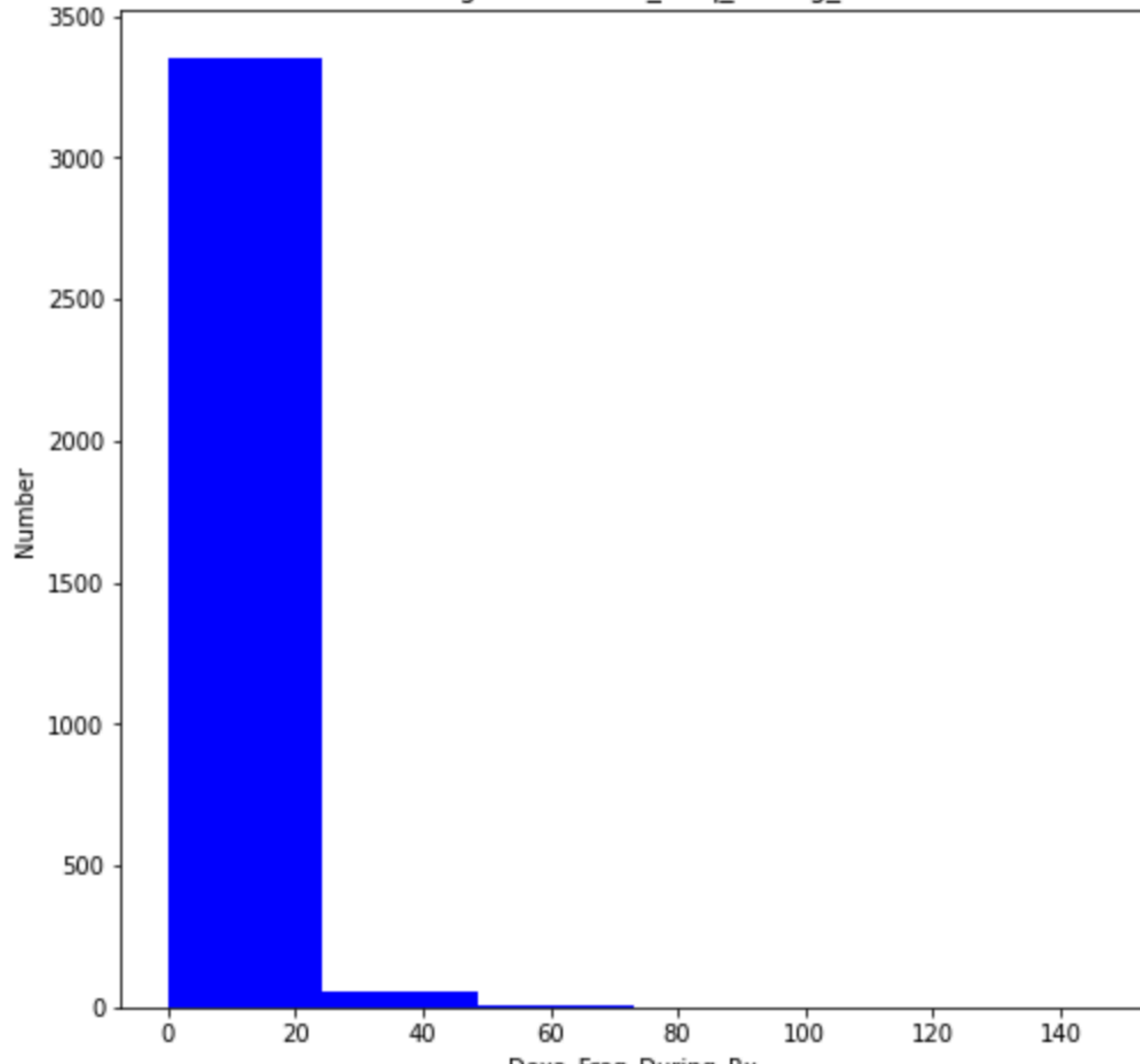
Percentage of Missing Values Per Column in Train Set



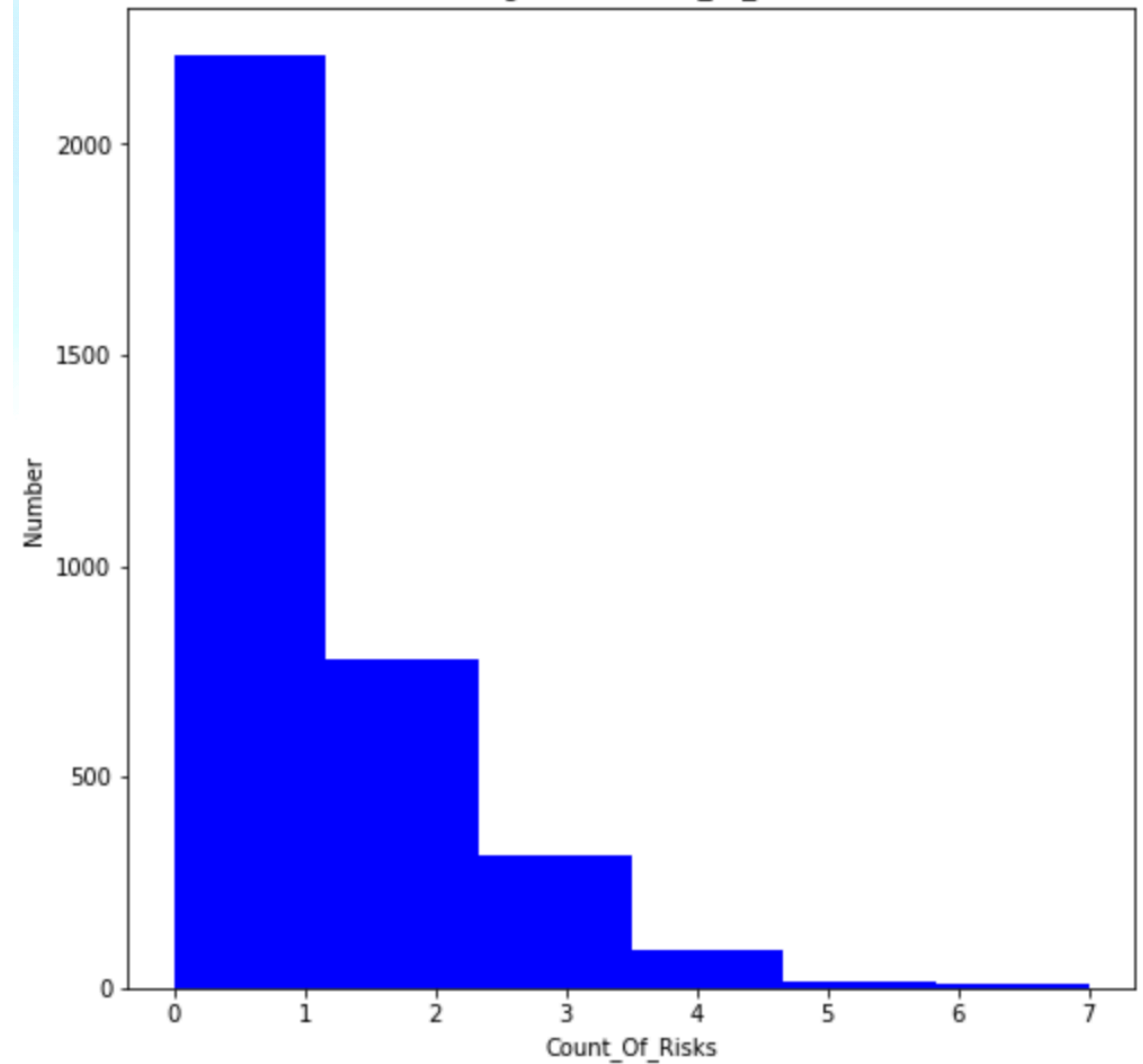


# Skewed Histogram

Histogram of Dexa\_Freq\_During\_Rx



Histogram of Count\_Of\_Risks



# EDA Summary

- We tried many ways to learn the permanence of the drugs on the patient, and as a result of the necessary analyzes, the effect of factors such as race, gender, origin was better understood. Apart from these, although there are many different factors, we achieved good results with a very good margin of precision as a result of the data we have.



# Model Building

- Categorical data were converted into numerical data.
- Tree Classifier, Random Forest Classifier and LightGBM Classifier regression techniques were used.

## Final Recommendation

- The Random Forest Classifier model is the most suitable model because it has the highest accuracy rate.
- It will give the best results in persistent and non-persistent drug effects compared to other models.