

# Data Science Internship at Data Glacier

## Healthcare: Persistency of a drug (Data Science)

### Week 8: Deliverables

**Name:** Chooladeva Piyasiri

**University:** National Institute of Business Management (NIBM)

**Email:** [chooladevapiyasiri@gmail.com](mailto:chooladevapiyasiri@gmail.com)

**Country:** Sri Lanka

**Specialization:** Data Science

**Batch Code:** LISUM18

**Date:** 26 March 2023

**Submitted to:** Data Glacier

# **Table of Contents**

**Problem Statement**

**Project Plan**

**Data Understanding**

**Data Problem**

## Problem Statement

ABC is a pharmaceutical company. Medication persistency following a patient's doctor's prescription is a problem that ABC, a pharmaceutical company, wants to comprehend. Persistence of drugs is the duration of time a patient takes medication, from initiation to discontinuation of therapy. To automate this identifying process, this organization has contacted an analytics firm.

The analytics firm has to create a classification for the given dataset with the aim of gathering insights on the factors that are affecting the persistency.

## Project Plan

Weeks	Date	Plan
Weeks 07	19 March 2023	Problem Statement, Data Collection, Data Intake Report
Weeks 08	26 March 2023	Data Understanding, Data Problems
Weeks 09	02 April 2023	Feature Extraction
Weeks 10	09 April 2023	Data Cleansing and Transformation
Weeks 11	16 April 2023	EDA Presentation and proposed modeling technique
Weeks 12	23 April 2023	Model Selection and Model Building/Dashboard
Weeks 13	30 April 2023	Final Project Report and Code

## Data Understanding

The dataset contains 3423 observations with 26 columns and 68 features. The feature names and their data types are shown below.

Ptid	object
Persistency_Flag	object
Gender	object
Race	object
Ethnicity	object
Region	object
Age_Bucket	object
Ntm_Speciality	object
Ntm_Specialist_Flag	object
Ntm_Speciality_Bucket	object
Gluco_Record_Prior_Ntm	object
Gluco_Record_During_Rx	object
Dexa_Freq_During_Rx	int64
Dexa_During_Rx	object
Frag_Frac_Prior_Ntm	object
Frag_Frac_During_Rx	object
Risk_Segment_Prior_Ntm	object
Tscore_Bucket_Prior_Ntm	object
Risk_Segment_During_Rx	object
Tscore_Bucket_During_Rx	object
Change_T_Score	object
Change_Risk_Segment	object
Adherent_Flag	object
Idn_Indicator	object
Injectable_Experience_During_Rx	object
Comorb_Encounter_For_Screening_For_Malignant_Neoplasms	object
Comorb_Encounter_For_Immunization	object
Comorb_Encntr_For_General_Exam_W_O_Complaint,_Susp_Or_Reprtd_Dx	object
Comorb_Vitamin_D_Deficiency	object
Comorb_Other_Joint_Disorder_Not_Elsewhere_Classified	object
Comorb_Encntr_For_Oth_Sp_Exam_W_O_Complaint_Suspected_Or_Reprtd_Dx	object
Comorb_Long_Term_Current_Drug_Therapy	object
Comorb_Dorsalgia	object
Comorb_Personal_History_Of_Other_Diseases_And_Conditions	object
Comorb_Other_Disorders_Of_Bone_Density_And_Structure	object
Comorb_Disorders_of_lipoprotein_metabolism_and_other_lipidemias	object
Comorb_Osteoporosis_without_current_pathological_fracture	object
Comorb_Personal_history_of_malignant_neoplasm	object
Comorb_Gastro_esophageal_reflux_disease	object
Concom_Cholesterol_And_Triglyceride_Regulating_Preparations	object
Concom_Narcotics	object
Concom_Systemic_Corticosteroids_Plain	object
Concom_Anti_Depressants_And_Mood_Stabilisers	object
Concom_Fluoroquinolones	object
Concom_Cephalosporins	object
Concom_Macrolides_And_Similar_Types	object
Concom_Broad_Spectrum_Penicillins	object
Concom_Anaesthetics_General	object
Concom_Viral_Vaccines	object
Risk_Type_1_Insulin_Dependent_Diabetes	object
Risk_Osteogenesis_Imperfecta	object
Risk_Rheumatoid_Arthritis	object
Risk_Untreated_Chronic_Hyperthyroidism	object
Risk_Untreated_Chronic_Hypogonadism	object
Risk_Untreated_Early_Menopause	object
Risk_Patient_Parent_Fractured_Their_Hip	object
Risk_Smoking_Tobacco	object
Risk_Chronic_Malnutrition_Or_Malabsorption	object
Risk_Chronic_Liver_Disease	object
Risk_Family_History_Of_Osteoporosis	object
Risk_Low_Calcium_Intake	object
Risk_Vitamin_D_Insufficiency	object
Risk_Poor_Health_Frailty	object
Risk_Excessive_Thinness	object
Risk_Hysterectomy_Oophorectomy	object
Risk_Estrogen_Deficiency	object
Risk_Immobilization	object
Risk_Recurring_Falls	object
Count_OF_Risks	int64
dtype: object	

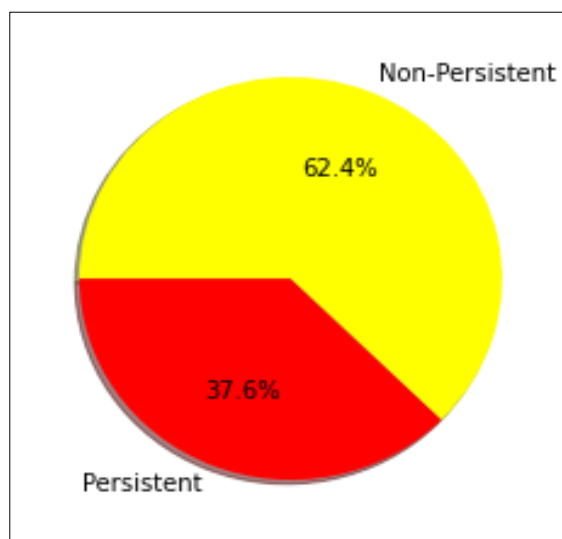
There are two numerical columns in the dataset, per the data type description. Which are:

- **Dexa\_Freq\_During\_Rx**
- **Count\_Of\_Risks**

The Descriptive Statistics of those numerical features of the dataset are shown below.

	Dexa_Freq_During_Rx	Count_Of_Risks
count	3424.000000	3424.000000
mean	3.016063	1.239486
std	8.136545	1.094914
min	0.000000	0.000000
25%	0.000000	0.000000
50%	0.000000	1.000000
75%	3.000000	2.000000
max	146.000000	7.000000

The target variable of the dataset is "Persistency\_Flag," - a flag indicating if a patient was persistent or not. There are 2135 non-persistent patients and 1289 persistent patients in this sample.



## Data Problems

### NA values

The dataset does not contain any NA values.

### Skewness & Kurtosis

#### Count\_Of\_Risks

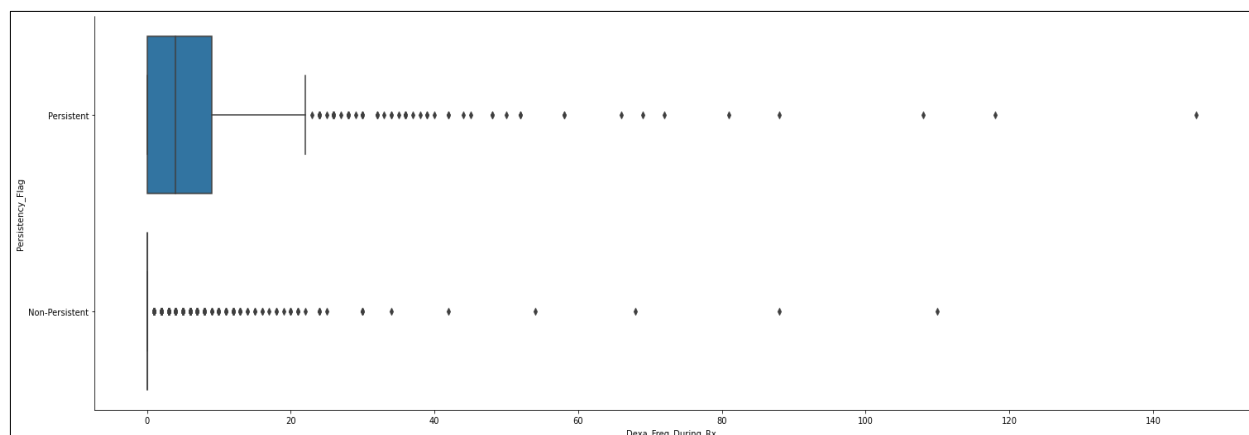
- ✓ The Count\_Of\_Risks distribution is moderately skewed. (0.879)
- ✓ The Count\_Of\_Risks distribution is Platykurtic (kurtosis <3). Compared to a normal distribution, its tails are shorter and thinner, and often its central peak is lower and broader. (0.900)

#### Dexa\_Freq\_During\_Rx

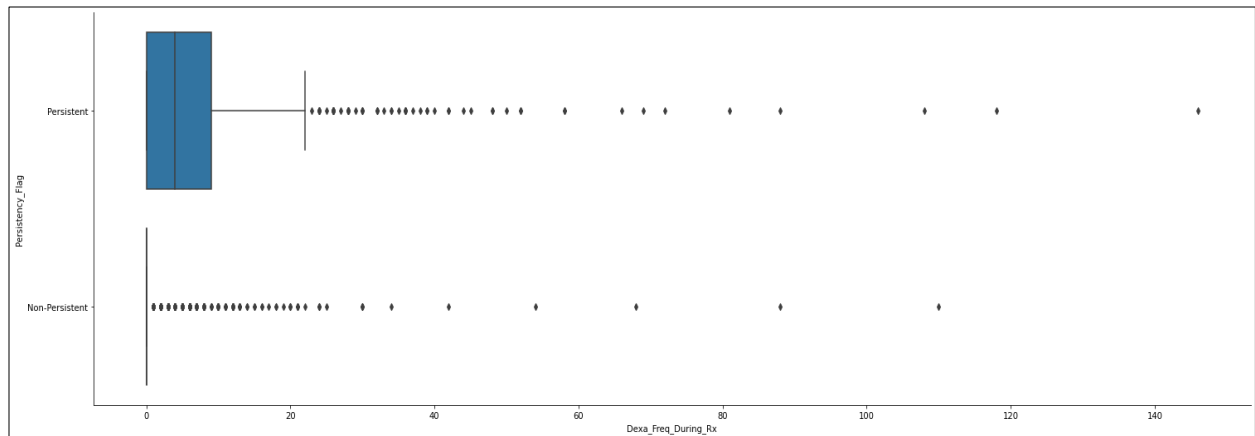
- ✓ The Dexa\_Freq\_During\_Rx distribution is highly skewed. (6.808)
- ✓ The Dexa\_Freq\_During\_Rx distribution is Leptokurtic (kurtosis >3). Compared to a normal distribution, its tails are longer and fatter, and often its central peak is higher and sharper. (74.758)

### Outliers

#### Dexa\_Freq\_During\_Rx



## Count\_Of\_Risks



Both numerical features contain some outliers.