# Data Science Project

**Project:** Healthcare - Persistency of a Drug

**Week 8** Deliverables

**Team Name:** Team Healthy Bones

**Member:** Jeeyeon Shon

**E-mail:** shon.jeeyeon@gmail.com

**College/Company:** Data Glacier

**Batch Code:** LISUM 10

**Country:** United States of America

**Specialization:** Data Science

**Submission Date:** July 25, 2022

**Submitted to:** Data Glacier

**Github Link:** https://github.com/shonjeeyeon/DG_Week_8

# Table of Contents

# Problem Description

A model will be established and deployed to automate identifying persistency of a certain pharmaceutical product.

Data of patients who take the medication will be used for analysis, and correlation between medication persistency and other factors such as patient demographics, provider attributes, clinical factors, and disease/treatment factors will be investigated. Finally, an optimal model to predict persistency based on above features will be selected and developed.

# Data Understanding

The dataset includes 3,424 records of patients on a certain medication. 69 features pertaining the demographics of the patient, attributes of the prescriber, and clinical/disease/treatment factors of the disease progression are present.

The client requested to build a model to predict a patient's drug persistency, so the column 'Persistent' will be the target variable. The prediction will use a classification process since the values of the target column are binary ('Persistent' vs. 'Non-Persistent')

# Summary of Columns and Data Types

| Bucket | Variable | index # | Dtype | Notes |
|---|---|---|---|---|
| Target | Persistency | 1 | Object | Non-Persistent: 62.35% Persistent: 37.65% (→ Imbalanced data) |
| Unique Row ID | Patient ID | 0 | | |
| Demographics | Gender | 2 | | |
| | Race | 3 | | NaN='Other/Unknown' (2.85%) Mode='Caucasian' (91.94%) |
| | Ethnicity | 4 | | NaN='Unknown' (2.66%) Mode='Non-Hispanic' (94.48%) |
| | Region | 5 | | NaN='Other/Unknown' (1.75%) Mode='Midwest' (40.39%) |
| | Age_Bucket | 6 | | . |
| | Ntm_Speciality | 7 | | NaN='Unknown' (9.05%) |

| | | | | |
|---|---|---|---|---|
| Prescriber Attributes | | | | Mode='General Practitioner' (44.83%) |
| | Ntm_Specialist_Flag | 8 | | |
| | Ntm_Speciality_Bucket | 9 | | |
| Clinical Factors | Gluco_Record_Prior_Ntm | 10 | | |
| | Gluco_Record_During_Rx | 11 | | |
| | Dexa_Freq_During_Rx | 12 | int64 | • Outlier issues<br>• The data is **skewed** (6.81)<br><br><table><tr><td>Count</td><td>3,424</td></tr><tr><td>Mean</td><td>3.02</td></tr><tr><td>Std</td><td>8.14</td></tr><tr><td>Min</td><td>0.00</td></tr><tr><td>25%</td><td>0.00</td></tr><tr><td>50%</td><td>0.00</td></tr><tr><td>75%</td><td>3.00</td></tr><tr><td>Max</td><td>146.00</td></tr></table> |
| | Dexa_During_Rx | 13 | Object | |
| | Frag_Frac_Prior_Ntm | 14 | | |
| | Frag_Frac_During_Rx | 15 | | |
| | Risk_Segment_Prior_Ntm | 16 | | |
| | Tscore_Bucket_Prior_Ntm | 17 | | NaN='Unknown' **(43.72%)** |

| | | | | |
|---|---|---|---|---|
| | | | | **The other two categories have very few differences in percentages** |
| | | | | HR_VHR 28.18% |
| | | | | VLR_LR 28.10% |
| | Risk_Segment_During_Rx | 18 | | NaN='Unknown' **(43.72%)** **The other two categories have very few differences in percentages** |
| | | | | <=-2.5 29.70% |
| | | | | >-2.5 26.56% |
| | Tscore_Bucket_During_Rx | 19 | | NaN='Unknown' **(43.72%)** Mode='No Change' (48.48%) |
| | Change_T_Score | 20 | | NaN='Unknown' **(65.01%)** Mode='No Change' (30.72%) |
| | Change_Risk_Segment | 21 | | |
| Disease/ Treatment Factors | Adherent_Flag | 22 | | |
| | Idn_Indicator | 23 | | |
| | Injectable_Experience_During_Rx | 24 | | |
| | Comorbidities columns (Column names start with 'Comorb_') | 25-38 | | |

| | | | | |
|---|---|---|---|---|
| | Concomitant drugs use columns (Column names start with 'Concom_') | 39-48 | | |
| | Risk factors columns | 49-67 | | |
| | Count_of_Risks | 68 | Dtype: int64 | • Outlier issues<br>• The data is **skewed** (0.88) |

| Count | 3,424 |
|---|---|
| Mean | 1.24 |
| Std | 1.09 |
| Min | 0.00 |
| 25% | 0.00 |
| 50% | 1.00 |
| 75% | 2.00 |
| Max | 7.00 |

## Problems and Suggested Actions

| Problem | Column | Index | Details | Categorical/ Quantifiable | Suggested Actions |
|---|---|---|---|---|---|
| Missing Data | Race | 3 | NaN='Other/Unknown' (2.85%) Mode='Caucasian' (91.94%) | Categorical | Impute with mode |
| | Ethnicity | 4 | NaN='Unknown' (2.66%) Mode='Non-Hispanic' (94.48%) | | |
| | Ntm_Speciality | 7 | NaN='Unknown' (9.05%) Mode='General Practitioner' (44.83%) | | |
| >40% Missing Data | Tscore_Bucket_Prior_Ntm | 17 | NaN='Unknown' **(43.72%)** <br> HR_VHR 28.18% <br> VLR_LR 28.10% | Categorical | Delete the columns because the proportions of missing data are |
| | Risk_Segment_During_Rx | 18 | NaN='Unknown' **(43.72%)** | | |

| | | | | | |
|---|---|---|---|---|---|
| | | | **The other two categories have very few differences in percentages** | | too large to impute without contributing to potential biases |
| | | | <=-2.5 : 29.70%<br>>-2.5 : 26.56% | | |
| | Tscore_Bucket_During_Rx | 19 | NaN='Unknown' **(43.72%)**<br>Mode='No Change' (48.48%) | | |
| | Change_T_Score | 20 | NaN='Unknown' **(65.01%)**<br>Mode='No Change' (30.72%) | | |
| Outliers/ Skews | Dexa_Freq_During_Rx | 12 | • Outlier issues<br>• The data is **skewed** (6.81)<br><br>Count : 3,424<br>Mean : 3.02<br>Std : 8.14<br>Min : 0.00<br>25% : 0.00 | Quantifiable | Use Tukey's rule to remove outliers<br><br>Also implement other |

| | | | 50% | 0.00 | | methods to remove skews, such as log transformation or box-cox method |
| | | | 75% | 3.00 | | |
| | | | Max | 146.00 | | |
| | Count_of_Risks | 68 | • Outlier issues<br>• The data is **skewed** (0.88) | | | |
| | | | Count | 3,424 | | |
| | | | Mean | 1.24 | | |
| | | | Std | 1.09 | | |
| | | | Min | 0.00 | | |
| | | | 25% | 0.00 | | |
| | | | 50% | 1.00 | | |
| | | | 75% | 2.00 | | |
| | | | Max | 7.00 | | |
| Imbalanced Target Data | Persistency | 1 | Non-Persistent: 62.35%<br>Persistent: 37.65% | | Categorical | Consider SMOTE |
| Encoding | Applies to every categorical column | | Categorical values are written in alphabet, which ML cannot process | | Categorical | Label, dummy, or one-hot encoding |

| Basic Cleaning | All columns | | Will need to remove upper cases, special characters, or spaces | Categorical/ Quantifiable | Switch the col names and values to all lower cases Remove special characters and spaces |
| --- | --- | --- | --- | --- | --- |

# Link to the Repository

https://github.com/shonjeeyeon/DG_Week_8