

Project: Data Science :: Persistency of a Drug

Week 12: Deliverables

Presented by: Sammy Suliman, Olivia Foster, and Tahsin Azad

Table of Contents:

| | |
|---------------------------|---|
| Team Member Detail..... | 2 |
| Problem Statement..... | 3 |
| Project Deadlines..... | 3 |
| Data Intake Report..... | 4 |
| Data Cleaning..... | 5 |
| Data Preprocessing..... | 6 |
| Exploratory Analysis..... | 6 |
| Model Training..... | 7 |
| GitHub Links..... | 7 |

Team Member Details:

Group Name: Call it Version 1.0

Member 1:

- **Name:** Sammy Suliman
- **Email:** sammysuliman@gmail.com
- **Country:** USA
- **College/Company:** UC Santa Barbara
- **Specialization:** Data Science

Member 2:

- **Name:** Olivia Foster
- **Email:** livia.n.foster@gmail.com
- **Country:** USA
- **College/Company:** NT Logistics
- **Specialization:** Data Science

Member 3:

- **Name:** Tahsin Azad
- **Email:** tahsinazad31@gmail.com
- **Country:** USA
- **College/Company:** University of California, Santa Barbara
- **Specialization:** Data Science

Problem Statement:

One of the challenges faced by Pharmaceutical companies is the persistence of a drug (that is, the extent to which a patient will act in accordance with the prescribed time interval, and dose of a medication) as the physician prescribed it. In this problem, we will automate the process of classifying factors that determine the persistence of a drug through Machine Learning and Python.

Drug persistence is a task of classifying different disorders and a patient's medical history to determine the dose and length of dose. In order to train our model, we will need to classify risk factors, medical histories, and disorders. To do this, we will be using a dataset based on over 3000 patients' records.

Project Deadlines:

| Week: | Due Date: | Plan: | Completed: |
|---------|-----------|---|-------------|
| Week 7 | 1/19/2023 | Problem Statement, Data Collection, Data Report | Yes |
| Week 8 | 1/26/2023 | Data Preprocessing | Yes |
| Week 9 | 2/2/2023 | Feature Extraction | Yes |
| Week 10 | 2/9/2023 | Building Model | Yes |
| Week 11 | 2/16/2023 | Model Result Evaluation | Yes |
| Week 12 | 2/23/2023 | Flask Development and Web Application | In Progress |
| Week 13 | 2/28/2023 | Final Submission (Report + Code + Final Submission) | No |

Data Intake Report:

Name: Final Project -- Data Science:: Healthcare - Persistency of a drug:: Group Project

Report date: 1/17/2023

Internship Batch: LISUM16

Version:<1.0>

Data intake by: Olivia Foster

Data intake reviewer: Tahsin Azad

Data storage location:

https://github.com/LiviaNFoster/DataGlacierFinalProject/blob/main/Healthcare_dataset.xlsx

Tabular data details:

| | |
|-------------------------------------|--------|
| Total number of observations | 3424 |
| Total number of files | 1 |
| Total number of features | 69 |
| Base format of the file | xlsx |
| Size of the data | 922 KB |

Proposed Approach:

- Unique Row identified using with the “Ptid” column, a unique ID for each patient
- No duplicates of Unique ID confirmed
- Assumptions
 - No false positives or errors during any of the testing
 - No dishonest and unbiased data recorded

Data Cleaning:

We began this process last week and continued to move forward with it.

Null Values:

There were no null values.

Yes versus No:

We had over half of the data be defined as either Y or N, “Yes” or “No” respectively. We replaced Y and N with 1 and 0 respectively.

Renaming Columns:

Because many of the columns had long and difficult names or were not great descriptors, we went ahead and renamed many of them to make referencing easier and to add clarity.

Case sensitive:

Made all words lowercase to make referencing easier without having to worry about case sensitivity.

Acknowledgement of outliers:

In exploring our data, we discovered that because our data is mostly binary (such as above in the ‘Yes versus No’ section) that there were no significant outliers at this time.

‘Age Bucket’ addressed:

In the original data, age was not presented in the same manner as in the description. Instead of producing quantifiable numbers for analysis, the values were grouped together as objects. To rectify this, each ‘bucket’ was replaced by a normal random output of values that fit the parameters of each ‘bucket.’

‘Unknown’s:

Later in the exploration of our data, we discovered Nulls that were disguised as ‘Unknown’s. To address this, we used hotkey encoding to deal with the unknown values. Depending on the contents of the column however, we would elect to drop the whole column if we did not find the contents of it relevant.

Data Preprocessing:

Encoding Categorical Data:

The data had certain categorical values. These values were changed into integer format, which will later be used for the model. Target Encoding, Label Encoding, One-Hot Encoding and Smoothing were observed.

Key Errors fixed:

Small key errors fixed, allowing for clean, correct encoding.

Exploratory Analysis:

Correlation:

Before building our model, we prepared several correlation maps to investigate what would have the greatest influence on our model for persistence. We found several higher level influencing factors but that reach at most a correlation 0.35.

Distribution charts:

A distribution chart for the age variable was prepared, displaying the mean and median age for the patients. This helps to visualize the data.

Box plot:

A box plot was created to look at any outliers. Since nearly all of the data is binary, only the boxplot for the number of “Risk_Counts” was observed. The boxplot proved that the “Risk_Counts” column had outliers.

Pie Chart:

Pie charts of Race, Ethnicity, Region and Gender were made to show the percentages and relationships between the different variables.

Bar plots:

Since much of our data was binary, we presented a majority of it as bar plots with percentages. This also was another great way to confirm that all the data was used and accounted for.

Model Training:

Flask using Logistic Regression:

Using Logistic Regression, the model was deployed on Flask and pushed to a server such as Heroku.

Ensemble Models:

Different types of ensemble models were trained and observed. These include Decision Tree Classifier, Random Forest Classifier, Ada Boosting, Bootstrap Aggregation, and Voting Tree Classifier

Stacking Model:

Logistic Regression, Decision Tree Classifier, and Random Forest Classifier were “stacked” to yield a better result.

Github Repo links:

General Repository:

<https://github.com/LiviaNFoster/DataGlacierFinalProject.git>

Week 12 Jupyter Notebook:

Please note that this file includes all of our work, where we each worked on it one at a time.

<https://github.com/LiviaNFoster/DataGlacierFinalProject/blob/main/Week12code>

Presentation:

<https://github.com/LiviaNFoster/DataGlacierFinalProject/blob/main/BusinessPresentation.pdf>