

Week 8: Data Overview

Submitted by: Cyprian Nyamwaro

Email: cypriannyamwaro@gmail.com

Country: Kenya

Submission date: August 12, 2021

Internship Batch: LISUM01

Problem Description:

ABC is a pharmaceutical company that wants to understand the persistency of a drug as per the physician's prescription for a patient. This company has approached an Analytics company to automate this process of identification. This Analytics company has given responsibility to CNN and has asked to come up with a solution to automate the persistency of a drug for the client ABC.

Data Understanding

The healthcare dataset is considerable size data having 69 columns altogether. The label is the Persistency Flag which is a binary data having value as True or False depending on the other features. The first column being the unique id of the patient is of no use for us as this is not going to help us in training the model. Hence the first thing we do is drop that column. We have analyzed the dataset and upon analysis we found that there are only few numerical data column and rest are either binary or string value.

Exploratory Data Analysis

We analyzed and dataset and found that there are many binary data having "Y" and "N" values. We mapped each if the data to 1 and 0 respectively. Also, we find no null values in the dataset and hence there was no need of handling it. We checked the numerical columns and found that one of the feature is having some outliers and we handled it using log transformation.

GitHub Repository:

Project Link: https://github.com/cnyamwaro/Data_Glacier_Cyprian.git