

**Team member's details:****Group Name: Data Dominators**

<a href="mailto:chau.devin031602@gmail.com">chau.devin031602@gmail.com</a>	Devin Chau	United States	San Jose State University	Data Science
--	------------	---------------	---------------------------	--------------

<a href="mailto:rohankhatri0507@gmail.com">rohankhatri0507@gmail.com</a>	Rohan Khatri	United States	San Jose State University	Data Science
--	--------------	---------------	---------------------------	--------------

<a href="mailto:ethan05dy@gmail.com">ethan05dy@gmail.com</a>	Ethan Dy	United States	San Jose State University	Data Science
--	----------	---------------	---------------------------	--------------

**Project: Data Science:: Healthcare - Persistency of a drug****Problem description**

We are building a predictive model that classifies patients into “persistent” or “non-persistent” categories based on factors like their demographics, medical history, physician characteristics, and treatment details. Factors like the patient level such as their age, risk factors, previous test results, or provider type allows for insights into why some patients continue therapy while others drop off. Thus understanding “persistence” levels. By analyzing these data points and finding patterns, the predictive model helps explain patient behavior and supports the creation of targeted interventions to improve adherence.

**Data Cleansing and Transformation**

The data cleansing and transformation process included a combination of techniques to enhance data quality and ensure compatibility with modeling tools:

- **Data Normalization and Duplicate Checks:**
  - MinMaxScaler was used to normalize numerical columns, bringing all variables to a common scale and ensuring equal contribution to the model.

```
[ ] 1 scaler = MinMaxScaler()
    2 encoded_data[numerical_cols] = scaler.fit_transform(encoded_data[numerical_cols])
    3
```

- Duplicate rows were detected based on the 'PatientID' column to prevent skewed analyses, and they were logged for further review.

```
1 duplicates = df.duplicated(subset='PatientID').sum()
```

- Columns were standardized to appropriate data types, reducing potential errors during analysis.
- **Outlier Handling:**
  - Devin Chau removed outliers from the Dexa\_Freq\_During\_Rx column by excluding entries where its value was zero.

```
1 df = df[df['Dexa_Freq_During_Rx'] != 0]
```

- Rohan Khatri applied the IQR method to cap extreme values within acceptable ranges, addressing potential data skew.

```
1 # Using IQR
2 Q1 = df['Dexa_Freq_During_Rx'].quantile(0.25)
3 Q3 = df['Dexa_Freq_During_Rx'].quantile(0.75)
4 IQR = Q3 - Q1
5 lower_bound = Q1 - 1.5 * IQR
6 upper_bound = Q3 + 1.5 * IQR
7
8 df['Dexa_Freq_During_Rx_No_Outliers'] = np.where(
9     df['Dexa_Freq_During_Rx'] > upper_bound, upper_bound,
10     np.where(df['Dexa_Freq_During_Rx'] < lower_bound, lower_bound, df['Dexa_Freq_During_Rx'])
11 )
12 df['Dexa_Freq_During_Rx_No_Outliers'].describe()
13
```

- **Binary Encoding:**
  - Devin Chau converted categorical binary values ('N' and 'Y') into numerical representations (0 and 1), ensuring compatibility with classifier models.

```
1 binary_columns = [col for col in df.columns if set(df[col].dropna().unique()) == {'N', 'Y'}]
2
3 for col in binary_columns:
4     df[col] = df[col].replace({'N': 0, 'Y': 1}).astype(int)
```

- **Risk Level Encoding:**

- Ethan Dy categorized patient risk levels into bins such as 'None', 'Low', 'Moderate', and 'High' based on Count\_Of\_Risks values, enabling targeted risk analysis.

```
1 risk_bins = [-1, 0, 2, 5, float('inf')]  
2 risk_labels = ['None', 'Low', 'Moderate', 'High']  
3 encoded_data['Risk_Level'] = pd.cut(encoded_data['Count_Of_Risks'].astype(int),  
4                                     bins=risk_bins,  
5                                     labels=risk_labels)  
6
```

- **Feature Scaling:**

- Ethan Dy scaled numerical features using MinMaxScaler to normalize their range, ensuring model fairness.

## **Consolidated Results**

### **Observations**

- **Outlier Removal:** Improved data consistency by removing erroneous values from Dexa\_Freq\_During\_Rx.
- **Binary Encoding:** Facilitated the use of binary data in models by converting values into integers.
- **Risk Level Encoding:** Categorized patient risk levels for analysis.
- **Normalization:** Ensured all variables contributed equally to the model.
- **Duplicate Checks:** Identified potential data redundancies.
- **Logistic Regression:** Provided risk predictions with an evaluated accuracy.
- **IQR Outlier Handling:** Capped extreme values to enhance data quality.