

**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****Data Normalization and Duplicate Checks**

- **MinMaxScaler:** The numerical columns were normalized to a common scale using MinMaxScaler to ensure that all variables contributed equally to the model.
- **Duplicate Checks:** Detected duplicate rows based on the 'PatientID' column, which could skew analysis. The duplicates were logged for further review.
- **Standardization of Data Types:** Columns were converted to appropriate data types to ensure consistency across the dataset, minimizing errors during analysis.

**Consolidated Results**

## My Contribution

### Standardizing data types:

```
1 categorical_cols = df.select_dtypes(include=['object']).columns.tolist()
2 numerical_cols = df.select_dtypes(include=['int64', 'float64']).columns.tolist()
3
4 df[categorical_cols] = df[categorical_cols].astype('category')
5
6 encoded_data = df.copy()
7 for col in categorical_cols:
8     encoded_data[col] = encoded_data[col].cat.codes
9
10
```

### Use MinMaxScaler to normalize and check for duplicate rows in the PatientID column:

```
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

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

[ ] 1 duplicates = df.duplicated(subset='PatientID').sum()
2
3    summary_stats = encoded_data.describe(include='all')
4
5    {
6        "categorical_cols": categorical_cols,
7        "numerical_cols": numerical_cols,
8        "duplicates_found": duplicates,
9        "summary_stats": summary_stats
10    }
11
```

- **Summary of Steps:**
  - Encoded risk levels into categorized bins.
  - Scaled numerical columns using MinMaxScaler.
  - Detected duplicates and provided a summary of the dataset.
- **Observations:**

- Duplicate entries detected: duplicates.
- Statistical summary of all variables provided insights into data distribution.

**Github Repo link**

<https://github.com/ethan05d/DataGlacier-Internship/tree/main/Week%209>