Name: Yusuf Aishat
Email: aishatyusuf12@gmail.com
Country: Nigeria
Specialization: Data Science

## PROBLEM DESCRIPTION

ABC Pharma wants to automate the identification of drug persistence as per the physician's prescription. The goal is to build a classification model to predict whether a patient will persist with their treatment and to identify the factors that influence drug persistence.

## DATA UNDERSTANDING

The dataset provided tracks patient persistence with their NTM (Non-Tuberculous Mycobacterial Disease) treatment. The dataset contains 3424 records, each representing patients with various demographics, attributes of their health care provider, clinical factors, disease, and treatment factors.

Total number of rows: 3424
Total number of columns: 69

| Bucket | Feature | Feature Description |
|---|---|---|
| Unique Row Id | Ptid | Data Type: Object<br>Missing Values: none |
| Target Variable | Persistency_Flag | Data Type: Object |

| | | Unique Values (in order of frequency):<br>● 'Non-Persistent': 62.35%<br>● 'Persistent': 37.64%<br>Missing Values: none |
|---|---|---|
| Demographics | Age_Bucket | Data Type: Object<br>Unique Values (in order of frequency):<br>● '>75': 42.02%<br>● '65-75': 31.71%<br>● '55-65': 21.40%<br>● '<55': 4.84%<br>Missing Values: none |
| | Race | Data Type: Object<br>Unique Values (in order of frequency):<br>● 'Caucasian': 91.93%<br>● 'Other/Unknown': 2.83%<br>● 'African American': 2.77%<br>● 'Asian': 2.45%<br>Missing Values: 'Other/Unknown' |
| | Region | Data Type: Object<br>Unique Values (in order of frequency):<br>● 'Midwest': 40.39%<br>● 'South': 36.41% |

| | | |
|---|---|---|
| | | • 'West': 14.66%<br>• 'Northeast': 6.77%<br>• Other/Unknown: 1.75%<br>Missing Values:<br>'Other/Unknown' |
| | Ethnicity | Data Type: Object<br>Unique Values (in order of frequency):<br>• 'Not Hispanic': 94.48%<br>• 'Hispanic': 2.86%<br>• 'Unknown': 2.65%<br>Missing Values: 'Unknown' |
| | Gender | Data Type: Object<br>Unique Values (in order of frequency):<br>• 'Female': 94.33%<br>• 'Male': 5.66%<br>Missing Values: none |
| | Idn_Indicator | Data Type: Object<br>Unique Values (in order of frequency):<br>• 'Y': 74.67%<br>• 'N': 25.32%<br>Missing Values: none |
| Provider Attributes | Ntm_Specialty | Data Type: Object<br>Unique Values:<br>'GENERAL PRACTITIONER',<br>'Unknown', |

| | | |
|---|---|---|
| | | 'ENDOCRINOLOGY',<br>'RHEUMATOLOGY',<br>'ONCOLOGY',<br>'PATHOLOGY',<br>'OBSTETRICS AND GYNECOLOGY',<br>'PSYCHIATRY AND NEUROLOGY',<br>'ORTHOPEDIC SURGERY',<br>'PHYSICAL MEDICINE AND REHABILITATION',<br>'SURGERY AND SURGICAL SPECIALTIES',<br>'PEDIATRICS',<br>'PULMONARY MEDICINE',<br>'HEMATOLOGY & ONCOLOGY',<br>'UROLOGY',<br>'PAIN MEDICINE',<br>'NEUROLOGY',<br>'RADIOLOGY',<br>'GASTROENTEROLOGY',<br>'EMERGENCY MEDICINE',<br>'PODIATRY',<br>'OPHTHALMOLOGY',<br>'OCCUPATIONAL MEDICINE',<br>'TRANSPLANT SURGERY',<br>'PLASTIC SURGERY', |

| | | 'CLINICAL NURSE SPECIALIST', 'OTOLARYNGOLOGY', 'HOSPITAL MEDICINE', 'ORTHOPEDICS', 'NEPHROLOGY', 'GERIATRIC MEDICINE', 'HOSPICE AND PALLIATIVE MEDICINE', 'OBSTETRICS & OBSTETRICS & GYNECOLOGY & OBSTETRICS & GYNECOLOGY', 'VASCULAR SURGERY', 'CARDIOLOGY', 'NUCLEAR MEDICINE' , Missing Values: 'Unknown' |
| | Ntm_Specialist_Flag | Data Type: Object<br>Unique Values (in order of frequency):<br>● 'Others': 58.79%<br>● 'Specialist': 41.21%<br>Missing Values: none |
| | Ntm_Specialist_Bucket | Data Type: Object<br>Unique Values (in order of frequency):<br>● 'OB/GYN/Others/PCP/ Unknown': 61.44% |

| | | |
|---|---|---|
| | | • 'Endo/Onc/Uro': 20.91% <br> • 'Rheum': 17.64 <br> Missing Values: none |
| Clinical Factors | Gluco_Record_Prior_N tm | Data Type: Object <br> Unique Values (in order of frequency): <br> • 'N': 76.48% <br> • 'Y': 23.51% <br> Missing Values: none |
| | Gluco_Record_During _Rx | Data Type: Object <br> Unique Values (in order of frequency): <br> • 'N': 73.65% <br> • 'Y': 26.34% <br> Missing Values: none |
| | Dexa_Freq_During_Rx | Data Type: Integer , Mean: 3.01 <br> Standard deviation: 8.13 <br> Min: 0 <br> 25%: 0 <br> 50%: 0 <br> 75%: 3 <br> Max: 146 <br> Missing Values: none |
| | Dexa_During_Rx | Data Type: Object <br> Unique Values (in order of frequency): <br> • 'N': 72.66% |

| | | |
|---|---|---|
| | | ● 'Y': 27.33%<br>Missing Values: none |
| | Frag_Frac_Prior_Ntm | Data Type: Object<br>Unique Values (in order of frequency):<br>● 'N': 83.87%<br>● 'Y': 16.12%<br>Missing Values: none |
| | Frag_Frac_During_Rx | Data Type: Object<br>Unique Values (in order of frequency):<br>● 'N': 87.82%<br>● 'Y': 12.17%<br>Missing Values: none |
| | Risk_Segment_Prior_N<br>tm | Data Type: Object<br>Unique Values (in order of frequency):<br>● 'VLR_LR': 56.39%<br>● 'HR_VHR': 43.60%<br>Missing Values: none: |
| | Tscore_Bucket_Prior_<br>Ntm | Data Type: Object<br>Unique Values (in order of frequency):<br>● '>-2.5': 56.98%<br>● '<=2.5': 43.01%<br>Missing Values: none |
| | Risk_Segment_During<br>_Rx | Data Type: Object<br>Unique Values (in order of frequency): |

| | | |
|---|---|---|
| | | • 'Unknown': 43.72%<br>• 'HR_VHR': 28.18%<br>• 'VLR_LR': 28.09%<br><br>Missing Values: 'Unknown': |
| | Tscore_Bucket_During_Rx | Data Type: Object<br>Unique Values (in order of frequency):<br>• 'Unknown': 43.72<br>• '<=2.5': 29.70%<br>• '>-2.5': 26.57%<br>Missing Values: none |
| | Change_T_Score | Data Type: Object<br>Unique Values (in order of frequency):<br>• 'No change': 48.48%<br>• 'Unknown': 43.72%<br>• 'Worsened': 5.05%<br>• Improved: 2.74%<br>Missing Values: 'Unknown': |
| | Change_Risk_Segment | Data Type: Object<br>Unique Values (in order of frequency):<br>• 'Unknown': 65.09%<br>• 'No change': 30.72%<br>• 'Worsened': 3.53%<br>• Improved: 0.64%<br>Missing Values: 'Unknown': |

| Disease/Treatment Factor | Injectable_Experience_During_Rx | Data Type: Object<br>Unique Values (in order of frequency):<br>• 'Y': 89.25%<br>• 'N': 10.74%<br>Missing Values: none |
|---|---|---|
| | NTM - Risk Factors (19 risk factor columns) | Data Type: Object<br>Unique Values:'Y', 'N'<br>Missing Values: none |
| | Count_Of_Risks | Data Type: Integer,<br>Mean: 1.23<br>Standard deviation: 1.09<br>Min: 0<br>25%: 0<br>50%: 1<br>75%: 2<br>Max: 7<br>Missing Values: none |
| | NTM - Comorbidity (14 comorbidity columns) | Data Type: Object<br>Unique Values:'Y', 'N'<br>Missing Values: none |
| | NTM - Concomitancy (10 concomitancy columns) | Data Type: Object<br>Unique Values:'Y', 'N'<br>Missing Values: none |
| | Adherent_Flag | Data Type: Object<br>Unique Values (in order of frequency): |

| | | ● 'Adherent': 94.94% |
| | | ● 'Non-Adherent': 5.05% |
| | | Missing Values: none |

## DATA QUALITY ISSUES

- **Race**: ~2.83% of entries are labeled Other/Unknown. This entries will be relabeled as "Other" because there are other races outside of Caucasian, African American and Asian

- **Ethnicity**: ~2.65% were entered as Unknown. These will be replaced by the mode ("Not Hispanic") which is about 94.5%

- **Region**: ~1.75% are labeled as "Other/Unknown". These will also be replaced by the mode ("Midwest")

- **NTM Speciality**: At approximately 9.05%, the Unknown category is quite a sizeable portion and for now it will be left as is.
  An entry labeled as "Obstetrics & Obstetrics & Gynecology & Obstetrics & Gynecology" seems to be a data entry error and it will be added to the "Obstetrics and Gynecology" category.

- **Risk_Segment_During_Rx, Change_Risk_Segment, Tscore_Bucket_During_Rx, Change_T_Score**:
  These attributes have >40% of their entries as Unknown, as they are not adding much information to the data, they will be removed.

GITHUB REPOSITORY:
https://github.com/aishatyusuf/drug_persistence_abc_pharma