

Data Science HealthCare Project Drug Persistence ABC Pharma



By : Sourour Cherif

Mail : [Sourour.cherif@esprit .tn](mailto:Sourour.cherif@esprit.tn)

Data Science Engineering Student Esprit Tunisia

Problem description

ABC Pharma is looking for an automated way better than the traditional debilitating methods currently used to assess persistence of drugs as per the physician prescription, in order to have a deeper understanding on the factors impacting the persistence of their drug. The aim is to know if a patient, based on his/her information, will follow the prescription of the physician and continue taking the drug for all the treatment time. We have been provided with a dataset which contains patients' details.

Business understanding

We will create a classification model as a solution that divides patients into categories depending on their information, to determine if a patient was persistent or not.

Our goal is to create a web application that might be used as an automated solution to this process of identification.

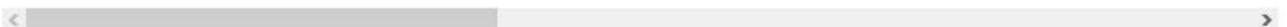
Data understanding

To fit any predictive model on a dataset, we need to understand the complexity of the dataset before deciding which predictive model to use to get optimal performance .

```
persist_pd.head()
```

	Ptid	Persistence_Flag	Gender	Race	Ethnicity	Region	Age_Bucket	Ntm_Speciality	Ntm_Specialist_Flag	Ntm_Speciality_Bucket	...	Risk_F
0	P1	Persistent	Male	Caucasian	Not Hispanic	West	>75	GENERAL PRACTITIONER	Others	OB/GYN/Others/PCP/Unknown	...	
1	P2	Non-Persistent	Male	Asian	Not Hispanic	West	55-65	GENERAL PRACTITIONER	Others	OB/GYN/Others/PCP/Unknown	...	
2	P3	Non-Persistent	Female	Other/Unknown	Hispanic	Midwest	65-75	GENERAL PRACTITIONER	Others	OB/GYN/Others/PCP/Unknown	...	
3	P4	Non-Persistent	Female	Caucasian	Not Hispanic	Midwest	>75	GENERAL PRACTITIONER	Others	OB/GYN/Others/PCP/Unknown	...	
4	P5	Non-Persistent	Female	Caucasian	Not Hispanic	Midwest	>75	GENERAL PRACTITIONER	Others	OB/GYN/Others/PCP/Unknown	...	

5 rows × 69 columns



Type of data

```
persist_pd.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 3424 entries, 0 to 3423
```

```
Data columns (total 69 columns):
```

#	Column	Non-Null Count	Dtype
0	Ptid	3424 non-null	object
1	Persistency_Flag	3424 non-null	object
2	Gender	3424 non-null	object
3	Race	3424 non-null	object
4	Ethnicity	3424 non-null	object
5	Region	3424 non-null	object
6	Age_Bucket	3424 non-null	object
7	Ntm_Speciality	3424 non-null	object
8	Ntm_Specialist_Flag	3424 non-null	object
9	Ntm_Speciality_Bucket	3424 non-null	object
10	Gluko_Record_Prior_Ntm	3424 non-null	object
11	Gluko_Record_During_Rx	3424 non-null	object
12	Dexa_Freq_During_Rx	3424 non-null	int64
13	Dexa_During_Rx	3424 non-null	object
14	Frag_Frac_Prior_Ntm	3424 non-null	object
15	Frag_Frac_During_Rx	3424 non-null	object
16	Risk_Segment_Prior_Ntm	3424 non-null	object
17	Tscore_Bucket_Prior_Ntm	3424 non-null	object
18	Risk_Segment_During_Rx	3424 non-null	object
19	Tscore_Bucket_During_Rx	3424 non-null	object
20	Change_T_Score	3424 non-null	object
21	Change_Risk_Segment	3424 non-null	object
22	Adherent_Flag	3424 non-null	object
23	Idn_Indicator	3424 non-null	object
24	Injectable_Experience_During_Rx	3424 non-null	object
25	Comorb_Encounter_For_Screening_For_Malignant_Neoplasms	3424 non-null	object
26	Comorb_Encounter_For_Immunization	3424 non-null	object
27	Comorb_Encntr_For_General_Exam_W_O_Complaint,_Susp_Or_Reprtd_Dx	3424 non-null	object
28	Comorb_Vitamin_D_Deficiency	3424 non-null	object
29	Comorb_Other_Joint_Disorder_Not_Elsewhere_Classified	3424 non-null	object

29	Comorb_Other_Joint_Disorder_Not_Elsewhere_Classified	3424	non-null	object
30	Comorb_Encntr_For_Oth_Sp_Exam_W_O_Complaint_Suspected_Or_Reprtd_Dx	3424	non-null	object
31	Comorb_Long_Term_Current_Drug_Therapy	3424	non-null	object
32	Comorb_Dorsalgia	3424	non-null	object
33	Comorb_Personal_History_Of_Other_Diseases_And_Conditions	3424	non-null	object
34	Comorb_Other_Disorders_Of_Bone_Density_And_Structure	3424	non-null	object
35	Comorb_Disorders_of_lipoprotein_metabolism_and_other_lipidemias	3424	non-null	object
36	Comorb_Osteoporosis_without_current_pathological_fracture	3424	non-null	object
37	Comorb_Personal_history_of_malignant_neoplasm	3424	non-null	object
38	Comorb_Gastro_esophageal_reflux_disease	3424	non-null	object
39	Concom_Cholesterol_And_Triglyceride_Regulating_Preparations	3424	non-null	object
40	Concom_Narcotics	3424	non-null	object
41	Concom_Systemic_Corticosteroids_Plain	3424	non-null	object
42	Concom_Anti_Depressants_And_Mood_Stabilisers	3424	non-null	object
43	Concom_Fluoroquinolones	3424	non-null	object
44	Concom_Cephalosporins	3424	non-null	object
45	Concom_Macrolides_And_Similar_Types	3424	non-null	object
46	Concom_Broad_Spectrum_Penicillins	3424	non-null	object
47	Concom_Anaesthetics_General	3424	non-null	object
48	Concom_Viral_Vaccines	3424	non-null	object
49	Risk_Type_1_Insulin_Dependent_Diabetes	3424	non-null	object
50	Risk_Osteogenesis_Imperfecta	3424	non-null	object
51	Risk_Rheumatoid_Arthritis	3424	non-null	object
52	Risk_Untreated_Chronic_Hyperthyroidism	3424	non-null	object
53	Risk_Untreated_Chronic_Hypogonadism	3424	non-null	object
54	Risk_Untreated_Early_Menopause	3424	non-null	object
55	Risk_Patient_Parent_Fractured_Their_Hip	3424	non-null	object
56	Risk_Smoking_Tobacco	3424	non-null	object
57	Risk_Chronic_Malnutrition_Or_Malabsorption	3424	non-null	object
58	Risk_Chronic_Liver_Disease	3424	non-null	object
59	Risk_Family_History_Of_Osteoporosis	3424	non-null	object
60	Risk_Low_Calcium_Intake	3424	non-null	object
61	Risk_Vitamin_D_Insufficiency	3424	non-null	object
62	Risk_Poor_Health_Frailty	3424	non-null	object
63	Risk_Excessive_Thinness	3424	non-null	object
64	Risk_Hysterectomy_Oophorectomy	3424	non-null	object
65	Risk_Estrogen_Deficiency	3424	non-null	object
66	Risk_Immobilization	3424	non-null	object
67	Risk_Recurring_Falls	3424	non-null	object
68	Count Of Risks	3424	non-null	int64

```
persist_pd.describe()
```

	Dexa_Freq_During_Rx	Count_Of_Risks
count	3424.000000	3424.000000
mean	3.016063	1.239486
std	8.136546	1.094914
min	0.000000	0.000000
25%	0.000000	0.000000
50%	0.000000	1.000000
75%	3.000000	2.000000
max	146.000000	7.000000

Let's see how we can find out unique elements in a column of the dataset.

```
persist_pd["Race"].unique()
```

```
array(['Caucasian', 'Asian', 'Other/Unknown', 'African American'],  
      dtype=object)
```

```
persist_pd["Age_Bucket"].unique()
```

```
array(['>75', '55-65', '65-75', '<55'], dtype=object)
```

```
persist_pd["Ntm_Speciality"].unique()
```

```
array(['GENERAL PRACTITIONER', 'Unknown', 'ENDOCRINOLOGY', 'RHEUMATOLOGY',  
      'ONCOLOGY', 'PATHOLOGY', 'OBSTETRICS AND GYNECOLOGY',  
      'PSYCHIATRY AND NEUROLOGY', 'ORTHOPEDIC SURGERY',  
      'PHYSICAL MEDICINE AND REHABILITATION',  
      'SURGERY AND SURGICAL SPECIALTIES', 'PEDIATRICS',  
      'PULMONARY MEDICINE', 'HEMATOLOGY & ONCOLOGY', 'UROLOGY',  
      'PAIN MEDICINE', 'NEUROLOGY', 'RADIOLOGY', 'GASTROENTEROLOGY',  
      'EMERGENCY MEDICINE', 'PODIATRY', 'OPHTHALMOLOGY',  
      'OCCUPATIONAL MEDICINE', 'TRANSPLANT SURGERY', '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'], dtype=object)
```

Data Problems

Data problems such as irrelevant columns , Null values , duplicates , skewed data ,outliers and many others may cause bad predictions ...

So we need to check if we have one of them to know then how to overcome it .

- Skewed Data :


```

✓ [497] def measure_skew_kurtosis(cols):
0s     for col in cols:
        print(col)
        result = data[[col]].agg(['skew', 'kurtosis']).transpose()
        print(result)
    measure_skew_kurtosis(numeric_col)

```

```

Dexa_Freq_During_Rx
                                skew  kurtosis
Dexa_Freq_During_Rx  6.80873  74.758378
Count_Of_Risks
                                skew  kurtosis
Count_Of_Risks  0.879791  0.900486

```

```

✓ [498] #skew and kurtosis values
0s     data.agg(['skew', 'kurtosis']).transpose()

```

	skew	kurtosis
Dexa_Freq_During_Rx	6.808730	74.758378
Count_Of_Risks	0.879791	0.900486

- Outliers

```

▶ # creating a box plot of numerical columns against persistency flag to identify outliers
def boxplot(data, cols):
    for col in cols:
        sns.set_style('whitegrid')
        sns.boxplot(x='Persistency_Flag', y=col, data=data)
        plt.title('Boxplot of ' + col)
        plt.ylabel(col) #setting text for y axis
        plt.show()
    boxplot(data, numeric_col)

```

Solutions :

Having duplicates in the dataset is not advisable and it often leads to overfitting.

```
persist_pd.duplicated().sum()
```

0

⇒ There is no duplicates

Missing Values

```
persist_pd.isnull().sum()
```

```
Ptid 0
Persistency_Flag 0
Gender 0
Race 0
Ethnicity 0
Risk_Hysterectomy_Oophorectomy ..
Risk_Estrogen_Deficiency 0
Risk_Immobilization 0
Risk_Recurring_Falls 0
Count_Of_Risks 0
Length: 69, dtype: int64
```

⇒ No missing Values

Eliminating Skewed data

```
✓ [501] #Checking skew after transformation
0s data.agg(['skew', 'kurtosis']).transpose()
```

	skew	kurtosis
Dexa_Freq_During_Rx	6.808730	74.758378
Count_Of_Risks	0.879791	0.900486
log_Dexa	1.405860	0.624570
log_Count_Risks	-0.091583	-1.006414

Example of removing 99% Percentile

```
✓ [505] # To remove the 99th percentile
0s q = data['Dexa_Freq_During_Rx'].quantile(0.99)
data_1 = data[data['Dexa_Freq_During_Rx'] < q]
data_1.describe()
```

	Dexa_Freq_During_Rx	Count_Of_Risks	log_Dexa	log_Count_Risks
count	3389.000000	3389.000000	3389.000000	3389.000000
mean	2.440248	1.240484	0.572915	0.685941
std	5.183446	1.095904	0.997375	0.499826
min	0.000000	0.000000	0.000000	0.000000
25%	0.000000	0.000000	0.000000	0.000000
50%	0.000000	1.000000	0.000000	0.693147
75%	3.000000	2.000000	1.386294	1.098612
max	34.000000	7.000000	3.555348	2.079442

Trying to :

- Remove no significant columns
- Handling with Skewed data
- Detect Outliers using different vizualisation tools