



**Data Glacier**

Your Deep Learning Partner

# Healthcare - Persistency of Drugs

Groupe name: Ehiafarin

Country: Iran

Groupe member: Alireza Ehiaei

Email: [arh.ehiaei@yahoo.com](mailto:arh.ehiaei@yahoo.com)

Specialization: Data Science

Data Cleansing and Transformation

**September-2021**

**TABLE OF CONTENTS**

**01**

**Introduction**

**02**

**Problem Statement**

**03**

**Dealing with missing values, method 1**

**04**

**Dealing with missing values, method 2**

**05**

**Dealing with missing values, method 3**

# Introduction

Every dataset has missing values that need to be treated appropriately to create a robust model. In this practice, I have discussed two ways to handle missing values. There is no unique rule to handle missing values in a specific manner.

One can use various methods on different features depending on how and what the data is about. Indeed, each dataset needs some specific approaches to handle missing values.

# Problem Statement

The dataset in this research is related to drug treatment, and except two features which are integer and have no missing values, there are 67 other features with categorical type which some of them have unknown values.

The aim is predicting drug persistency of a patient. And there is Persistency-Flag column as target column, with two labels of Persistent and Non-Persistent. It is a binary classification problem.

# Problem Statement

There are different ways to handle missing values including:

- Deleting Rows with missing values

- Impute missing values for continuous variable

- Impute missing values for categorical variable

- Using Algorithms that support missing values

- Prediction of missing values

<https://towardsdatascience.com/7-ways-to-handle-missing-values-in-machine-learning-1a6326adf79e>

# Problem Statement

## **Delete columns or rows with missing values:**

Missing values can be handled by deleting the rows or columns having null values. For example, one can choose If 50% or 60% of or more of the rows as null then the entire column can be dropped. Similarly, If a row has e.g. one or more columns values as null can also be dropped.

# Dealing with missing values, method 1:

## Dropping columns or rows including missing values:

For example, we can remove the columns with more than 30%, and drop the rows with 20% missing values. There are 3424 rows, we remove 4 columns which have more than 30% missing values.

```
1 list_column_with_nulls = []
2
3 for column in df.columns:
4     n_missed = len(df[df.loc[:,column]==])
5     percent = n_missed/3424
6     if percent > 0.3:
7         print(column)
8         print("percentage of missing values : ", percent)
9         list_column_with_nulls.append(column)
```

```
Risk_Segment_During_Rx
percentage of missing values : 0.44
Tscore_Bucket_During_Rx
percentage of missing values : 0.44
Change_T_Score
percentage of missing values : 0.44
Change_Risk_Segment
percentage of missing values : 0.65
```

```
1 len(list_column_with_nulls)
```

# Dealing with missing values

New data set has 4 less columns.

```
1 drug_Data.head()
```

	Ptid	Persistency_Flag	Gender	Race	Ethnicity	Region
0	P1	Persistent	Male	Caucasian	Not Hispanic	West
1	P2	Non-Persistent	Male	Asian	Not Hispanic	West
2	P3	Non-Persistent	Female	Other/Unknown	Hispanic	Midwest
3	P4	Non-Persistent	Female	Caucasian	Not Hispanic	Midwest
4	P5	Non-Persistent	Female	Caucasian	Not Hispanic	Midwest

```
1 drug_Data.shape
```

(3424, 69)

```
1 df = drug_Data
```

```
1 df.drop([col for col in list_column
```

```
1 df.head()
```

	Ptid	Persistency_Flag	Gender	Ra
0	P1	Persistent	Male	Caucasi
1	P2	Non-Persistent	Male	Asi
2	P3	Non-Persistent	Female	Other/Unknow
3	P4	Non-Persistent	Female	Caucasi
4	P5	Non-Persistent	Female	Caucasi

```
1 df.shape
```

(3424, 65)



# Dealing with missing values

Similarly, we drop the rows with just 2% missing values.

```
1 rows_with_missing_values = []
2 for i in range(len(df)) :
3     n=0
4     s = []
5     for j in df.columns:
6         if (df.loc[i,j]=='Other/Unknown') | (df.loc[i,j]=='Unknown'):
7             n=n+1
8             s.append(j)
9     if (n/69 > 0.02):
10         print('row :',i, '    number of missing values :', n)
11         print('Columns including missing values :',s)
12         rows_with_missing_values.append(i)
```

```
row : 154    number of missing values : 2
Columns including missing values : ['Race', 'Ethnicity']
row : 218    number of missing values : 2
Columns including missing values : ['Race', 'Ethnicity']
row : 221    number of missing values : 2
Columns including missing values : ['Race', 'Ntm_Speciality']
row : 406    number of missing values : 2
Columns including missing values : ['Race', 'Ethnicity']
```

# Dealing with missing values

Finally, we have a data set with 3379 rows and 65 columns.

	Ptid	Persistency_Flag	Gender	Race
0	P1	Persistent	Male	Caucasian
1	P2	Non-Persistent	Male	Asian
2	P3	Non-Persistent	Female	Other/Unknown
3	P4	Non-Persistent	Female	Caucasian
4	P5	Non-Persistent	Female	Caucasian
<div>&lt; <div></div></div>				
1	df.shape			
(3379, 65)				

# Dealing with missing values, method 2

## **Imputation**

Since the columns including missing values in this research are categorical, their missing values can be replaced with the most frequent category. We could use this method as first step, but if the number of missing values is very large the accuracy of final prediction would be reduced. Now, we removed some rows and columns, and impute the rest of missing values in each column with the mode of that column.

# Dealing with missing values

By assuming as nulls, still there are four columns including missing values with maximum 290 missed values:

Race	64
Ethnicity	58
Region	55
Age_Bucket	0
Ntm_Speciality	290

# Dealing with missing values

After replacing nulls with the mode of each column, no missing value remains.

**Filling nulls of each column with the most frequent (mode) of each column**

```
1 df = df.fillna(df.mode().iloc[0])
```

**Any null in data?**

```
1 df.isnull().values.any()
```

False

# Dealing with missing values, method 3

## Prediction of missed values

Another way to treat with missing values is predicting them. For example, assume that we want to predict all null values of the column 'Risk\_Segment\_During\_Rx' based on other columns. First, object data type have been converted to float type to be used in the LinearDiscriminantAnalysis model. Then, by finding missed values in target column and removing their rows we have data set for training the model.

# Dealing with missing values, method 3

Finally, by using the model the missed values of the last column are predicted and filled.

```
1 from numpy import nan
2 from pandas import read_csv
3 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
4 from sklearn.model_selection import KFold
5 from sklearn.model_selection import cross_val_score
6
7 X = Data_without_nulls.iloc[:, :-1]
8 y = Data_without_nulls.iloc[:, -1]
9 # define the model
10 model = LinearDiscriminantAnalysis()
11 # define the model evaluation procedure
12 cv = KFold(n_splits=3, shuffle=True, random_state=1)
13 # evaluate the model
14 result = cross_val_score(model, X, y, cv=cv, scoring='accuracy')
15 # report the mean performance
16 print('Accuracy: %.3f' % result.mean())
```

Accuracy: 0.926

# Dealing with missing values, method 3

Finally, there is no missed values:

**Predicting the last column using the model**

```
] : 1 data_ = []  
    2 data_ = pd.DataFrame(data_)  
    3 data_ = pd.DataFrame(model.predict(data_without_null_1.iloc[:, :-1]))  
    4
```

**Replacing predicted values just for missed values**

```
] : 1 for i in range(len(data_without_null_1)) :  
    2     if(pd.isnull(data_without_null_1.loc[i, 'Risk_Segment_During_Rx'])):  
    3         data_without_null_1.loc[i, 'Risk_Segment_During_Rx'] = data_.iloc[i, 0]
```



# Github link

Github link of notebook code:

[https://github.com/Alireza-Ehiaei/Data\\_Sciences/tree/main/Drug\\_Persistence1/Cleaning\\_data](https://github.com/Alireza-Ehiaei/Data_Sciences/tree/main/Drug_Persistence1/Cleaning_data)

Thank You