

## 1<sup>st</sup> step – Importing required libraries

## 2<sup>nd</sup> step – Loading data and basic EDA

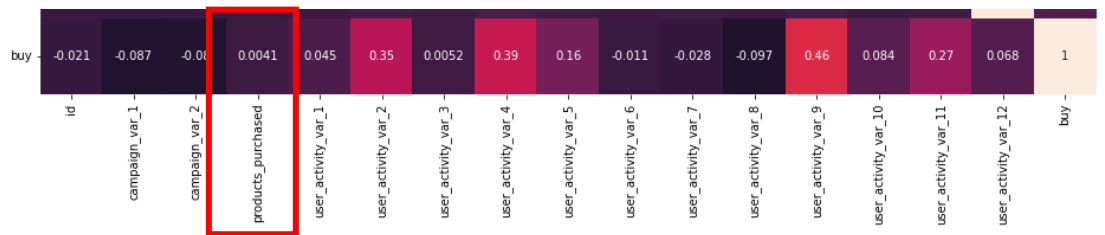
1. Info() and describe functions explain a lot about the dataset and the features.
  - a. Info() shows datatypes of each variable and non-null values
  - b. Describe lists all stats of all variables by looking at mean and median we can tell if there is any **skewness** in the data or not since the difference between mean and median is less for most of the columns, so we conclude that very less skewness is there and it will not affect much. So we can proceed further.
2. By isnull and sum function it's clear that 2 columns have null values present

```
✓ df.isnull().sum() # now printing null values count
```

```
id          0
created_at  0
campaign_var_1  0
campaign_var_2  0
products_purchased  20911
signup_date  15113
user_activity_var_1  0
user_activity_var_2  0
user_activity_var_3  0
user_activity_var_4  0
user_activity_var_5  0
user_activity_var_6  0
user_activity_var_7  0
user_activity_var_8  0
user_activity_var_9  0
user_activity_var_10  0
user_activity_var_11  0
user_activity_var_12  0
buy          0
dtype: int64
```

But instead of directly removing those null values I plotted correlation plot to check if these columns really have a major role in predicting our target value.

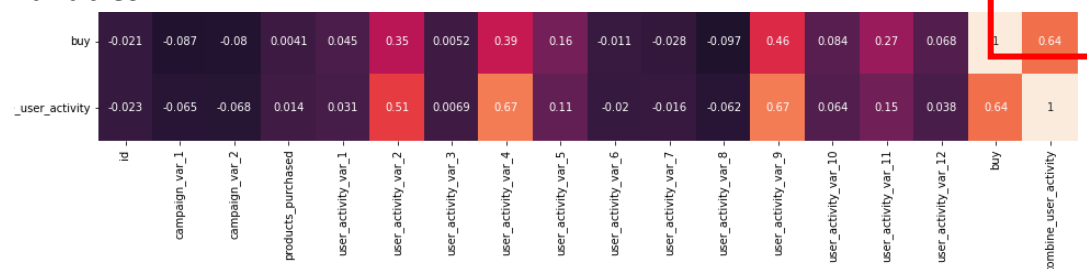
3. In correlation plot we can see the columns which have null values have very less correlation with target variable and we can ignore them



### 3<sup>rd</sup> Step –

**Feature engineering** – Now after seeing the correlation plot i thought to combine few features and make single column so that it will be easy for the model as well and also we can get feature with even higher correlation. I have written python code for the same.

- I've combined user\_activity\_var\_2, user\_activity\_var\_4 and user\_activity\_var\_9 to produce single column name combine\_user\_activity
- I've used simple approach to make column value as 1 if any of these above mention column is 1 else 0
- I've got the correlation value 0.64 which is highest among other variables



d.

### 4<sup>th</sup> Step –

**Feature selection** – Now for model building I've selected features with high correlation to avoid training the model with noise. I've selected the following features

```
[11] # feature selection - selecting only features with high correlation
X= df[['user_activity_var_5','user_activity_var_10','user_activity_var_11','combine_user_activity']]
Y= df[['buy']]
```

### 5<sup>th</sup> Step –

**Data splitting** – I've splitted data in ratio of 20 : 80 for testing and training respectively.

### **6<sup>th</sup> Step –**

**Model selection** – So I've tried different models and at last selected decision tree classifier since its producing better results.

After model selection I've initialized the model and trained using training data

### **7<sup>th</sup> step –**

**Model evaluation** – so it's giving accuracy of **97.5 %** on test data and precision of **93%**.

### **8<sup>th</sup> step –**

#### **Generating solution file**

Loading the testing data given on the site to produce solution file.

So I've performed the same operation of feature engineering and feature selection as we did before model training

Predicted the values on testing data set.

Mapped the id and predicted value as sample submission file and saved in csv file.