Data Analyst Intern at Data Glacier

Week 9: Deliverables

Project: Customer Segmentation

Name: Ece Yavuzyılmaz

University: Duzce University

Email: eceyavuzyilmaz@gmail.com

Country: Turkey

Specialization: Data Analyst

Internship Batch: LISUM21

Date: 02 July 2023

Table of Contents

1.	P	Prob	blem Description	3		
2.	E	Business Understanding				
3.	F	Project Lifecycle3				
4.	Ι	Data	a understanding	4		
	4.1.		Changing Column Name	5		
	4.2.		Data Types	5		
5.	Ι	Data	a Cleaning	6		
	5.1.		Deleting Unnecessary Column	6		
	5.2.		Missing Values	6		
	5.3.		Outlier Detection	7		

1. Problem Description

XYZ bank wants to roll out Christmas offers to their customers. But Bank does not want to roll out same offer to all customers instead they want to roll out personalized offer to particular set of customers. If they manually start understanding the category of customer then this will be not efficient and also they will not be able to uncover the hidden pattern in the data (pattern which group certain kind of customer in one category). Bank approached ABC analytics company to solve their problem. Bank also shared information with ABC analytics that they don't want more than 5 group as this will be inefficient for their campaign.

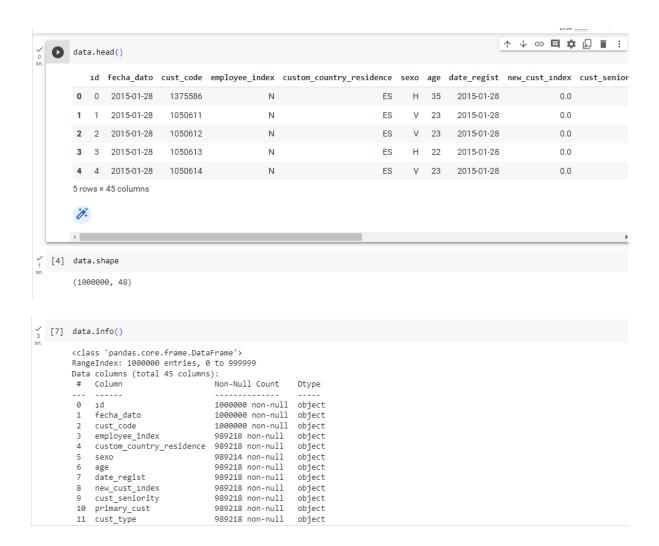
2. Business Understanding

ABC analytics proposed customer segmentation approach to Bank. ABC analytics assigned this talk to their analytics team and instructed their team to come up with the approach and feature which group similar behavior customer in one category and others in different category.

3. Project Lifecycle

Weeks	Date	Plan
Week 7	18 June 2023	Business Understanding
Week 8	26 June 2023	Data Understanding
Week 9	2 July 2023	EDA
Week 10	9 July 2023	Feature Engineering, Model Building
Week 11	16 July 2023	Model Evaluation
Week 12	23 July 2023	Presentation
Week 13	30 July 2023	Document the challenges

4. Data understanding



```
989218 non-null
                                                   object
     12 cust_relation_type
[7]
     13 residence index
                                  989218 non-null
                                                   object
     14 foreigner_index
                                 989218 non-null
                                                   obiect
     15 channel
                                 989139 non-null
                                                   object
     16 deceased_index
                                 989218 non-null
                                                   object
                                989218 non-null
     17 addres_type
                                                   object
     18 cod_prov
                                 982266 non-null
                                                   object
     19 province name
                                 982266 non-null
                                                   object
                                989218 non-null
     20 activity_index
                                                   obiect
     21 gross_income
                                824817 non-null
     22 save account
                                 1000000 non-null
                                                  obiect
                                1000000 non-null object
     23 current_account
     24 derivada_account
25 payroll_account
                               1000000 non-null object
                                 1000000 non-null
                                                  object
     26 junior_account
                                 1000000 non-null object
      27 mas_particu_account 1000000 non-null object
     28 particu_account
                                 1000000 non-null object
     29 particu_plus_account
                                 1000000 non-null object
      30 short term deposit
                                1000000 non-null object
     30 short_term_aeposit
31 medium_term_deposit
                                 1000000 non-null object
     31 medium_cc...__ ,
32 long_term_deposit
                                 1000000 non-null object
     33 e account
                                 1000000 non-null object
                                 1000000 non-null object
     34 funds
     35 mortgage
                                 1000000 non-null object
                                1000000 non-null object
     36 pensions
                                 1000000 non-null object
     37 loans
      38
        taxes
                                  1000000 non-null
                                  1000000 non-null object
     39 credit card
   40 securities
                                 1000000 non-null object
   41 home account
                                  1000000 non-null object
                                994598 non-null object
   42 payroll
   43 a pension
                                994598 non-null object
   44 direct debit
                                1000000 non-null object
  dtypes: object(45)
  memory usage: 343.3+ MB
```

4.1. Changing Column Name

First, some column names have been changed for better understanding and interpretation the data.

```
#changing some column name
data.rename({'incodpers':'cust_code', 'ind_empleado':'employee_index', 'fecha_alta': 'date_regist'}, axis=1, inplace=Tr
data.rename({'Unnamed: 0':'id','pais_residencia':'custom_country_residence', 'ind_nuevo':'new_cust_index', 'antiguedad
data.rename({'indrel':'primary_cust', 'ult_fec_cli_1t':'last_data_primary_cust', 'indrel_1mes': 'cust_type'}, axis=1,
data.rename({'tiprel_1mes':'cust_relation_type', 'indresi':'residence_index', 'indext': 'foreigner_index'}, axis=1, in
data.rename({'conyuemp':'spouse_index', 'canal_entrada':'channel', 'indfall':'deceased_index', 'tipodom':'addres_type'}
data.rename({'ind_cco_fin_ult1':'current_account', 'ind_cder_fin_ult1':'derivada_account', 'ind_cno_fin_ult1':'payroll
data.rename({'ind_cco_fin_ult1':'mas_particu_account', 'ind_ctop_fin_ult1':'medium_term_deposit', 'ind_dela_fin_ult1':
data.rename({'ind_deco_fin_ult1':'funds', 'ind_hip_fin_ult1':'mortgage', 'ind_plan_fin_ult1':'pensions', 'ind_pres_fin_
data.rename({'ind_reca_fin_ult1':'taxes', 'ind_tjcr_fin_ult1':'credit_card', 'ind_recibo_ult1':'securities', 'ind_viv
data.rename({'ind_nomina_ult1':'payroll', 'ind_nom_pens_ult1':'a_pension', 'ind_recibo_ult1':'direct_debit'}, axis=1,
```

4.2. Data Types

All data types are seen as 'objects'. These data types have been changed as datetime64[ns](2), float64(3), int64(29), object(11).

```
data['cust_code']=data['cust_code'].astype("int64")
data['cust_seniority']=data['cust_seniority'].astype("int64")
                                                                                                     ↑ ↓ ⊕ 目 ‡ ♬ 👔
    data['activity_index']=data['activity_index'].astype("int64")
     data['save_account']=data['save_account'].astype("int64")
    data['current account']=data['current account'].astype("int64")
     data['derivada_account']=data['derivada_account'].astype("int64")
    data['payroll_account']=data['payroll_account'].astype("int64")
     data['junior_account']=data['junior_account'].astype("int64")
     data['mas_particu_account']=data['mas_particu_account'].astype("int64")
     data['particu_account']=data['particu_account'].astype("int64")
     data['particu_plus_account']=data['particu_plus_account'].astype("int64")
    data['short_term_deposit']=data['short_term_deposit'].astype("int64")
     data['medium_term_deposit']=data['medium_term_deposit'].astype("int64")
     data['long_term_deposit']=data['long_term_deposit'].astype("int64")
     data['e_account']=data['e_account'].astype("int64")
     data['e_account']=data['e_account'].astype("int64")
    data['funds']=data['funds'].astype("int64")
     data['mortgage']=data['mortgage'].astype("int64")
     data['pensions']=data['pensions'].astype("int64")
     data['loans']=data['loans'].astype("int64")
     data['taxes']=data['taxes'].astype("int64")
     data['credit_card']=data['credit_card'].astype("int64")
     data['securities']=data['securities'].astype("int64")
     data['home_account']=data['home_account'].astype("int64")
     data['direct_debit']=data['direct_debit'].astype("int64")
```

5. Data Cleaning

5.1. Deleting Unnecessary Column

```
#deleting unnecessary column
data.drop('ult_fec_cli_1t',axis=1, inplace=True)
data.drop('conyuemp',axis=1, inplace=True)
data.drop('ind_aval_fin_ult1',axis=1, inplace=True)
```

5.2. Missing Values

Two different methods were used to fill in the missing value. These methods are; fillna() and interpolate().

fillna()

We can use *fillna()* function to fill NaN values. In this here, we can use different methods such 'backfill', 'bfill', 'pad', 'ffill' and mean/median/mode based approach to fill the missing values.

```
↑ ↓ ⊖ 目 🛊 🗓 📋
data['age']=data['age'].fillna(data['age'].mean()).astype("int64")
    data['cust_type']=data['cust_type'].fillna(method ='ffill').astype("int64")
    data['primary_cust']=data['primary_cust'].fillna(method ='ffill')
    data['cust_relation_type']=data['cust_relation_type'].fillna(method ='ffill')
    data['residence_index']=data['residence_index'].fillna(method='ffill')
    data['foreigner_index']=data['foreigner_index'].fillna(method='ffill')
    data['channel']=data['channel'].fillna(method='ffill')
    data['deceased index']=data['deceased index'].fillna(method='ffill')
    data['addres_type']=data['addres_type'].fillna(method='ffill')
    data['province_name']=data['province_name'].fillna(method='ffill')
    data['activity_index']=data['activity_index'].fillna(method='ffill')
    data['payroll']=data['payroll'].fillna(method='ffill').astype("int64")
    data['a_pension']=data['a_pension'].fillna(method='ffill').astype("int64")
    data['employee_index']=data['employee_index'].fillna(method='ffill')
    data['custom_country_residence']=data['custom_country_residence'].fillna(method='ffill')
    data['sexo']=data['sexo'].fillna(method='ffill')
    data['new_cust_index']=data['new_cust_index'].fillna(method='ffill')
    data['cod_prov']=data['cod_prov'].fillna(method='ffill')
```

interpolate()

```
# Interpolate backwardly across the column:
data['gross_income'].interpolate(method ='linear', limit_direction ='backward', inplace=True)

#data['gross_income']=data['gross_income'].fillna(data['gross_income'].mean())
```

5.3. Outlier Detection

