Customer Segmentation Project

Week 9

Done By: Shaimaa Al-Khawlani

1. Problem description

Most banks around the world have variant large customer base with different income levels, ages, characteristics, values and lifestyles.

XYZ bank wants to increase the production and the satisfactions of all customers categories by 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).

2. Data Understanding

The existing data, which was provided by the bank, is the bank's customers data. However, the data contains many columns that will help the analytics team analyze the data and build a customer segmentation approach for the bank.

Since the data does not contain a dependent variable or (Target), We believe that machine learning (clustering) techniques would be appropriate to use for this type of data.

Size: 1000000 records, 48 columns.

• Columns Description:

Column Name	Description
fecha_dato	The table is partitioned for this column
ncodpers	Customer code
ind_empleado	Employee index: A active, B ex employed, F filial, N not employee,
	P pasive
pais_residencia	Customer's Country residence
sexo	Customer's sex
age	Age
fecha_alta	The date in which the customer became as the first holder of a
	contract in the bank
ind_nuevo	New customer Index. 1 if the customer registered in the last 6
	months.
antiguedad	Customer seniority (in months)

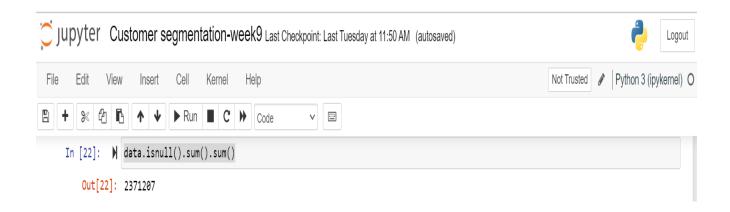
indrel	1 (First/Primary), 99 (Primary customer during the month but not at the end of the month)
ult_fec_cli_1t	Last date as primary customer (if he isn't at the end of the month)
indrel_1mes	Customer type at the beginning of the month ,1 (First/Primary customer), 2 (co-owner),P (Potential),3 (former primary), 4(former co-owner)
tiprel_1mes	Customer relation type at the beginning of the month, A (active), I (inactive), P (former customer),R (Potential)
indresi	Residence index (S (Yes) or N (No) if the residence country is the same than the bank country)
indext	Foreigner index (S (Yes) or N (No) if the customer's birth country is different than the bank country)
conyuemp	Spouse index. 1 if the customer is spouse of an employee
canal_entrada	channel used by the customer to join
indfall	Deceased index. N/S
tipodom	Addres type. 1, primary address
cod_prov	Province code (customer's address)
nomprov	Province name
ind_actividad_cliente	Activity index (1, active customer; 0, inactive customer)
renta	Gross income of the household
ind_ahor_fin_ult1	Saving Account
ind_aval_fin_ult1	Guarantees
ind_cco_fin_ult1	Current Accounts
ind_cder_fin_ult1	Derivada Account
ind_cno_fin_ult1	Payroll Account
ind_ctju_fin_ult1	Junior Account
ind_ctma_fin_ult1	Más particular Account
ind_ctop_fin_ult1	particular Account
ind_ctpp_fin_ult1	particular Plus Account
ind_deco_fin_ult1	Short-term deposits
ind_deme_fin_ult1	Medium-term deposits
ind_dela_fin_ult1	Long-term deposits
ind_ecue_fin_ult1	e-account
ind_fond_fin_ult1	Funds
ind_hip_fin_ult1	Mortgage
ind_plan_fin_ult1	Pensions
ind_pres_fin_ult1	Loans
ind_reca_fin_ult1	Taxes
ind_tjcr_fin_ult1	Credit Card

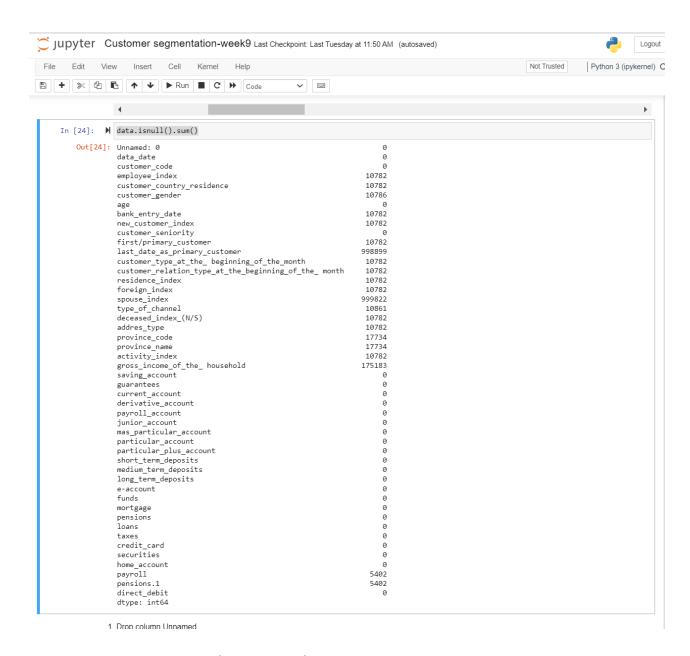
ind_valo_fin_ult1	Securities
ind_viv_fin_ult1	Home Account
ind_nomina_ult1	Payroll
ind_nom_pens_ult1	Pensions
ind_recibo_ult1	Direct Debit

• Data Cleaning:

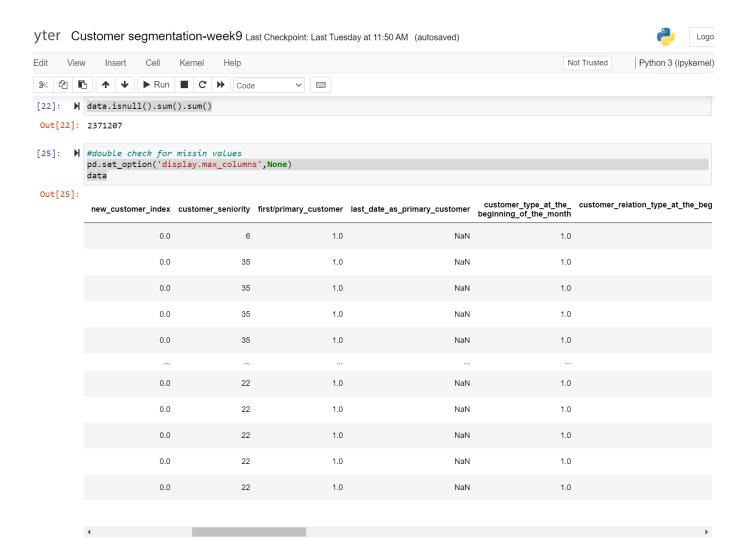
1. Handling Missing Values

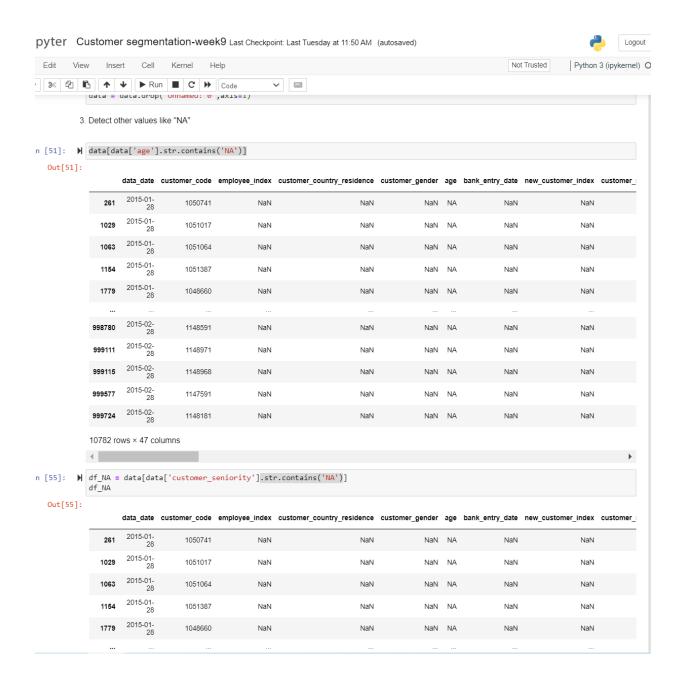
From the below print screen we have 2371207 missing values that must be filled or removed.





The missing values contains (NaN and NA) values as the below print screens





The above print screen shows that there are 10782 row that are considered as null, so it is recommend to remove them. We retrieved these values based on "NA" Columns of "Age" and "Customer Stationary". There is no need for these row as all of them not indicate and benefit sight so they have been removed.

```
In [70]: ► data.isnull().sum()
   Out[70]: data_date
                                                                            0
                                                                            0
             customer_code
             employee_index
                                                                            0
                                                                            0
             customer_country_residence
             customer_gender
                                                                            4
             age
             bank_entry_date
                                                                            0
             new_customer_index
                                                                            0
             customer_seniority
             first/primary_customer
                                                                            0
             last_date_as_primary_customer
                                                                       988117
             customer_type_at_the_ beginning_of_the_month
                                                                            0
             customer_relation_type_at_the_beginning_of_the_ month
                                                                            0
             residence index
                                                                            0
             foreign_index
                                                                            0
             spouse_index
                                                                       989040
             type of channel
             deceased_index_(N/S)
                                                                           0
             addres_type
                                                                           0
             province_code
                                                                         6952
                                                                         6952
             province name
             activity_index
             gross_income_of_the_ household
                                                                       164401
             saving_account
                                                                            0
             guarantees
             current_account
                                                                            а
             derivative_account
                                                                            0
             payroll_account
             junior_account
                                                                            0
             mas_particular_account
                                                                            Θ
             particular_account
             particular_plus_account
             short_term_deposits
                                                                            0
             medium_term_deposits
                                                                            0
             long_term_deposits
                                                                            0
             e-account
                                                                            0
             funds
             mortgage
                                                                            0
             pensions
             loans
                                                                            0
                                                                            0
             taxes
             credit_card
                                                                            0
             securities
                                                                            0
             home account
                                                                            0
             payroll
                                                                          100
             pensions.1
                                                                          100
             direct debit
             dtype: int64
In [71]: ► data.isnull().sum().sum()
   Out[71]: 2155745
```

After Removing the 10782 rows, the null values decrease significantly from 2371207 to 2155745

Then we replaced the rest of missing values with ZEROS to git ride of all missing values as the below print screen.

```
In [76]: M fill_data.isnull().sum()
   Out[76]: data date
              customer_code
employee_index
              customer_country_residence
                                                                             0
              customer_gender
                                                                             0
              bank_entry_date
              new_customer_index
              customer_seniority
              first/primary_customer
              last_date_as_primary_customer
customer_type_at_the_ beginning_of_the_month
              customer_relation_type_at_the_beginning_of_the_ month
              residence_index
              foreign_index
spouse_index
              type_of_channel
              deceased\_index\_(N/S)
              addres_type
province_code
              province_name
              activity_index
              gross_income_of_the_ household
              saving_account
              guarantees
              current_account
              derivative_account
              payroll_account
              junior_account
              mas_particular_account
              particular_account
              particular_plus_account
short_term_deposits
              medium_term_deposits
              long_term_deposits
              e-account
              funds
              mortgage
              pensions
              loans
              taxes
              credit_card
              securities
              home_account
              payroll
                                                                             0
              pensions.1
              direct_debit
              dtype: int64
In [77]: | fill_data.isnull().sum().sum()
   Out[77]: 0
 In [ ]: 📕
```

Then convert the types of some columns in order to get best insights:

1. Convert Dates columns from Object to Date.

2. Convert the below columns from Object to integer.

3. Check unique values to convert that needs to convert to category type.

Check Unuique Values to convert that needs to convert to category type

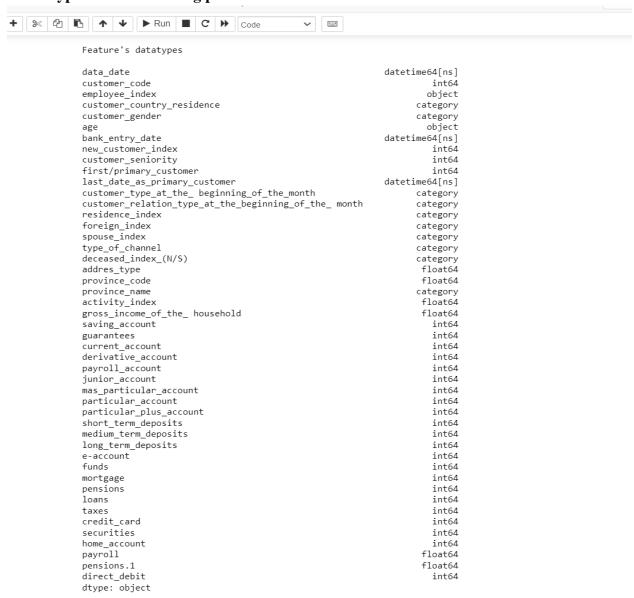
```
Out[88]:
                                                  Column_Name Num_Unique
                                                    addres_type
                                                       data_date
                                                derivative_account
                                                  junior_account
            30
                                               particular_account
            32
                                             short term deposits
                                           medium_term_deposits
            25
                                                 current_account
                                              long_term_deposits
            36
                                                         funds
                                                        mortgage
             38
                                                        pensions
                                                           loans
                                                      taxes
                                                      credit_card
                                                      securities
             43
                                                   home_account
                                                       e-account
                                                  saving_account
direct_debit
                                                    activity_index
                                             ceased_index_(N/S)
                                             new_customer_index
foreign_index
                                                     pensions.1
                                               customer_gender
            12
                   ustomer relation type at the beginning of the...
                    customer_type_at_the_ beginning_of_the_month
                                                employee_index
                                   last_date_as_primary_customer
                                     province_code
                                                                             53
                                                  province_name
                                      customer country residence
```

```
Check unique values of the columns that have object type
           In [76]: | data['customer_country_residence'].unique()
                  Out[76]: array(['ES', 'CA', 'CH', 'CL', 'IE', 'AT', 'NL', 'FR', 'GB', 'DE', 'DO', 'BE', 'AR', 'VE', 'US', 'MX', 'BR', 'II', 'EC', 'PE', 'CO', 'HN', 'FI', 'SE', 'AL', 'PT', 'MZ', 'CN', 'TW', 'PL', 'IN', 'CR', 'NI', 'HK', 'AD', 'CZ', 'AE', 'MA', 'GR', 'FR', 'RO', 'II', 'RU', 'GT', 'GA', 'NO', 'SN', 'MR', 'UA', 'BG', 'PR', 'RO', 'II', 'RU', 'GT', 'SA', 'CI', 'QA', 'LU', 'PA', 'BA', 'BO', 'AU', 'BY', 'KE', 'SG', 'HR', 'MD', 'SK', 'TR', 'AO', 'CU', 'GG', 'EG', 'ZA', 'DK', 'UY', 'GE', 'TH', 'DZ', 'LB', 'JP', 'NG', 'PK', 'TN', 'TG', 'KR', 'GH', 'RS', 'VN', 'PH', 'KN', 'NZ', 'MM', 'KH', 'GI', 'SL', 'GN', 'GM', 'CG', 'LV', 'LT', 'ML', 'MK', 'HU', 'IS', 'LY', 'CF', 'GM', 'KZ', 'CD', 'BZ'], dtype=object)
           In [77]: M data['customer_gender'].unique()
                    Out[77]: array(['H', 'V', 0], dtype=object)
           Out[79]: array([1, 3, 2], dtype=int64)
           In [80]: ) data['customer_relation_type_at_the_beginning_of_the_ month'].unique()
                    Out[80]: array(['A', 'I', 'P'], dtype=object)
           In [81]: M data['residence_index'].unique()
                    Out[81]: array(['S', 'N'], dtype=object)
           In [82]: | data['foreign_index'].unique()
                    Out[82]: array(['N', 'S'], dtype=object)
          In [83]: | data['spouse_index'].unique()
     Out[83]: array([0, 'N', 'S'], dtype=object)
In [84]: | data['type_of_channel'].unique()
                      Out[84]: array(['KHL', 'KHE', 'KHD', 'KFA', 'KFC', 'KAT', 'KAZ', 'RED', 'KHC',
                                                                                                                                                                                           'KBG', 0, 'KGC',
                                                                'KHK', 'KGN', 'KHM', 'KHO',
'KHF', 'KFK', 'KHN', 'KHA',
                                                                                                                                      'KDH', 'KEH',
                                                                                                                                                                         'KAD',
                                                                                                                                                                        'KFD',
                                                                                                                                       'KAF',
                                                                                                                                                        'KGX',
                                                                                                                    'KAH',
                                                                                                                                       'KAR',
                                                                                                                                                                          'KFL',
                                                                                                                                                                                            'KAI',
                                                                                                                                                                                                              'KFU',
                                                                 'KAB', 'KCC', 'KAE',
                                                                                                                                                        'KFJ',
                                                                'KAQ', 'KFS', 'KAA', 'KAP', 'KDE', 'KFV',
                                                                                                                                       'KAJ',
                                                                                                                                                        'KFN',
                                                                                                                                                                          'KGV'.
                                                                                                                                                                                           'KGY',
                                                                                                                     'KFP',
                                                                                                                                                                                                             'KFF'
                                                                                                                    '013',
                                                                                                                                       'K00',
                                                                                                                                                        'KAK',
                                                                                                                                                                          'KCK',
                                                                                                                                                                                            'KCL',
                                                                                                                                                                                                              'KAY',
                                                                                                                                                                                           'KBQ',
                                                                                                                                       'KCG',
                                                                                                                                                                          'KDY',
                                                                                                                                                                                                              'KDA',
                                                                 'KBU', 'KDR', 'KAC',
                                                                                                                    'KDT',
                                                                                                                                                        'KDO',
                                                                                                                    'KBZ',
                                                                                                                                                                         'KAS',
                                                                'KBO', 'KCI', 'KEC', 'KCA', 'KAL', 'KDC',
                                                                                                                                      'KES', 'KDX', 'KAS', 'KCS', 'KCB', 'KDU',
                                                                                                                                                                                           '007',
'KDQ',
                                                                                                                                                                                                             'KEU'
                                                                                                                    'KAW',
                                                                                                                                                                                                              'KCN',
                                                                                                   'KCH',
                                                                                                                    'KCD',
                                                                                                                                       'KCE',
                                                                                                                                                                          'KBL',
                                                                                                                                                                                           'KEA',
                                                                 'KCM', '004',
                                                                                                                                                        'KEV',
                                                                                                                                                                                                              'KBH
                                                                                                                    'KAO',
                                                                                                                                      'KEJ', 'KEO', 'KEI',
                                                                                                                                                                                           'KEW',
                                                                  'KDV', 'KFT', 'KEY',
                                                                                                                                                                                                             'KDZ'
                                                                  'KBV', 'KBR', 'KBF',
                                                                                                                    'KDP',
                                                                                                                                       'KCO',
                                                                                                                                                        'KCF',
                                                                                                                                                                          'KCV',
                                                                                                                                                                                           'KAM',
                                                                                                                                                                                                              'KEZ',
                                                                                                                    'KCT',
                                                                                                                                       'KDD',
                                                                                                                                                                          'KCU',
                                                                                                                                                                                           'KBX',
                                                                 'KBD', 'KAN',
                                                                                                   'KBY',
                                                                                                                                                        'KBW',
                                                                                                                                                                                                              'KDB
                                                                                                                    'KBP',
                                                                                                                                      'KBN', 'KEB', 'KDS',
                                                                 'KBS', 'KBE', 'KCX', 'KDF', 'KEF', 'KCP',
                                                                                                                                                                                          'KEL', 'KDG',
'KFI', 'KBM',
                                                                                                                    'KDM',
                                                                                                                                                        'KDW',
                                                                                                                                                                         'KBJ',
                                                                                                                                      'KBB',
                                                                'KEG', 'KEN', 'KEC', 'KAV', 'KFH', 'KFM', 'KAU', 'KED', 'KFR', 'KKM', 'KAU', 'KED', 'KFR', 'KKM', 'KDN', 'KDN', 'KEE', 'KCR', 'KCQ', 'KEM', 'KCJ'], dtype=object)
              In [85]: | data['province_name'].unique()
                     Out[85]: array(['MALAGA', 'CIUDAD REAL', 'ZARAGOZA', 'TOLEDO', 'LEON', 'GIPUZKOA', 'CACERES', 'GIRONA', 'ZAMORA', 'BARCELONA', 'SALAMANCA', 'BURGOS', 'HUESCA', 'NAVARRA', 'AVILA', 'SEGOVIA', 'LUGO', 'LERIDA', 'MADRID', 'ALICANTE', 'SORIA', 'SEVILLA', 'CANTABRIA', 'BALEARS, ILLES', 'VALLADOLID', 'PONTEVEDRA', 'VALENCIA', 'TERUEL', 'CORUÑA, A', 'OURENSE', 'JAEN', 'CUENCA', 'BIZKAIA', 'CASTELLON', 'RIOJA, LA', 'ALBACETE', 'BADAJOZ', 'MURCIA', 'CADIZ', 'ALMERIA', 'GUADALAJARA', 'PALENCIA', 'PALMAS, LAS', 'CORDOBA', 'HUELVA', 'GRANADA' 'ASTIRITAS', 'CANTA CRIZT DE TENRETEFE', 'METETEFE', 'METETETE', 'MET
```

'GRANADA', 'ASTURIAS', 'SANTA CRUZ DE TENERIFE', 'MELILLA', 'TARRAGONA', 'ALAVA', 0, 'CEUTA'], dtype=object)

4. Convert to Category type

Print Types after converting process



5. Drop column "addres_type" since it has only one value and it will not affect data insights.

```
In [98]: #drop Undeeded Columns
data=data.drop("addres_type",axis=1)
```

Finally, after the above process they have been saved to new excel file.

```
In [102]: M data.to_csv('custSeg_cleaned.csv')
```