# Project: Customer Segmentation

**Team member:**

- MAMADI FOFANA

- mamadi.fofana@edu.dsti.institute

- Republic of Guinea

- Data Science Tech Institute (DSTI)

- Data Science Specialization

# Problem description

**Problem Statement:** 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.

## Business understanding (Customer Segmentation)

We will propose customer segmentation approach to the Bank.

Our goal will be to create models which we use in turn to find a solution

We will adopt three predictive models:

- Clustering models which group similar behavior customer in one category and others in different category.

- Classification or Recommender models to predict group of a customer and which will be used to evaluate quality of the clustering model.

## Project lifecycle along with deadline

1. Business Understanding    ( Week 7)

2. Data Understanding        (week 8)

3. EDA                       (Week 8)

4. Feature Engineering       (Week 9)

4. Model Building            (Week 10)

5. Model Evaluation          (Week 11)

6. Presentation               (week 12)

7. Document the challenges (Week 13)

# Data understanding

## Data source

### Data source used is: cust_seg.csv

**Data source Link:**
https://drive.google.com/drive/folders/1bfCpJIKmp6IHxiLPWvOS2nU1dc24pViB

| 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 |

## Type of data

Data source contains several sources of data:

- Numerical: integer and float
- Categorical:

## Sweetviz (EDA Tool)  overview

```
1000000          ROWS
      0       DUPLICATES
 1.1 GB          RAM
     47        FEATURES
     39       CATEGORICAL
      3        NUMERICAL
      5          TEXT
```

## Pandas overview

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000000 entries, 0 to 999999
Data columns (total 47 columns):
 #   Column                 Non-Null Count     Dtype
---  ------                 --------------     -----
 0   fecha_dato             1000000 non-null   object
 1   ncodpers               1000000 non-null   int64
 2   ind_empleado           989218 non-null    object
 3   pais_residencia        989218 non-null    object
 4   sexo                   989214 non-null    object
 5   age                    1000000 non-null   object
 6   fecha_alta             989218 non-null    object
 7   ind_nuevo              989218 non-null    float64
 8   antiguedad             1000000 non-null   object
 9   indrel                 989218 non-null    float64
 10  ult_fec_cli_1t         1101 non-null      object
 11  indrel_1mes            989218 non-null    float64
 12  tiprel_1mes            989218 non-null    object
 13  indresi                989218 non-null    object
 14  indext                 989218 non-null    object
 15  conyuemp               178 non-null       object
 16  canal_entrada          989139 non-null    object
 17  indfall                989218 non-null    object
 18  tipodom                989218 non-null    float64
 19  cod_prov               982266 non-null    float64
 20  nomprov                982266 non-null    object
 21  ind_actividad_cliente  989218 non-null    float64
 22  renta                  824817 non-null    float64
 23  ind_ahor_fin_ult1      1000000 non-null   int64
```

```
24   ind_aval_fin_ult1        1000000 non-null   int64
25   ind_cco_fin_ult1         1000000 non-null   int64
26   ind_cder_fin_ult1        1000000 non-null   int64
27   ind_cno_fin_ult1         1000000 non-null   int64
28   ind_ctju_fin_ult1        1000000 non-null   int64
29   ind_ctma_fin_ult1        1000000 non-null   int64
30   ind_ctop_fin_ult1        1000000 non-null   int64
31   ind_ctpp_fin_ult1        1000000 non-null   int64
32   ind_deco_fin_ult1        1000000 non-null   int64
33   ind_deme_fin_ult1        1000000 non-null   int64
34   ind_dela_fin_ult1        1000000 non-null   int64
35   ind_ecue_fin_ult1        1000000 non-null   int64
36   ind_fond_fin_ult1        1000000 non-null   int64
37   ind_hip_fin_ult1         1000000 non-null   int64
38   ind_plan_fin_ult1        1000000 non-null   int64
39   ind_pres_fin_ult1        1000000 non-null   int64
40   ind_reca_fin_ult1        1000000 non-null   int64
41   ind_tjcr_fin_ult1        1000000 non-null   int64
42   ind_valo_fin_ult1        1000000 non-null   int64
43   ind_viv_fin_ult1         1000000 non-null   int64
44   ind_nomina_ult1           994598 non-null   float64
45   ind_nom_pens_ult1         994598 non-null   float64
46   ind_recibo_ult1          1000000 non-null   int64
dtypes: float64(9), int64(23), object(15)
memory usage: 358.6+ MB
```

## Problems in the data

After descriptive analysis (univariate analysis), correlation analysis (bivariate analysis), and exploratory data analysis, we notice following:

- **Missing value in some fields**
  For missing value in categorical columns, we replace missing value with a new category
  For missing value in numerical values, we replace missing value with mean value of the columns
  - We have a total of **2 371 207** missing values as per column
    - ind_empleado        10782
    - pais_residencia     10782
    - sexo                10786

- fecha_alta            10782
- ind_nuevo             10782
- indrel                10782
- ult_fec_cli_1t        998899
- indrel_1mes           10782
- tiprel_1mes           10782
- indresi               10782
- indext                10782
- conyuemp              999822
- canal_entrada         10861
- indfall               10782
- tipodom               10782
- cod_prov              17734
- nomprov               17734
- ind_actividad_cliente 10782
- renta                 175183
- ind_nomina_ult1       5402
- ind_nom_pens_ult1     5402

- **Duplicate values**

No duplicated values were found before data preprocessing, but found out duplicated values after removing some columns

Those duplicated value has been removed

- **High percentage of missing values in some columns in features:**

We dropped columns with high percentage of missing value which cannot decrease performance of the model

- ult_fec_cli_1t
- conyuemp

- **Column with unique categorical value**

We remove column with unique categorical value since it brings no insight to the model

- tipodom

- **Low variance in some numeric attributes**

So far, no variance issue has been observed in the numerical column

- **Low entropy in some categorical attributes.**

The identified problem is a very small entropy of the feature, which means that most of the records have the same categorical values. In most cases, it is safe to remove categorical attributes with low entropy. This will not harm the performance of the model, and it can reduce the complexity of the model.

- ind_nuevo
- indrel
- indrel_1mes
- indresi
- indfall
- ind_ahor_fin_ult1
- ind_aval_fin_ult1
- ind_cder_fin_ult1
- ind_deco_fin_ult1
- ind_deme_fin_ult1
- ind_pres_fin_ult1
- ind_viv_fin_ult1
- ind_hip_fin_ult1
- ind_plan_fin_ult1
- ind_valo_fin_ult1
- ind_fond_fin_ult1
- ind_ctma_fin_ult1
- ind_ctju_fin_ult1

- indext
- ind_empleado
- pais_residencia

- **high correlated columns (*cor* with ind_nom_pens_ult1 >0.8 )**

We remove columns highly correlated with others (we choose abs(correlation) > 0.8)

- ind_cno_fin_ult1
- ind_nomina_ult1
- fecha_alta

- **Columns with high cardinality**

The identified problem occurs when the number of unique values is too large for categorical attributes. This high cardinality creates problems for the typical one-hot-encoding process, creating a representation in an extremely high-dimensional space

Removing the feature if the number of occurrences per unique value of the attributes is too low

- ncodpers
- fecha_alta

- **Data in wrong format**

Some data are in wrong format need to be converted in the right format

- Age
- Antiguedad

- **Outliers and errors management**

Field *antiguedad* contains an errors value (-999 999) which need to be dropped

- **Data not scaled for some numerical variable**

For the clustering algorithm to consider all numerical attributes as equal, they must all have the same scale

- Age
- Antiguedad
- Renta

We have to scale those columns using Z-score transformation

- **Data skewness**

Columns *renta* is right skewed and we use log transformation to reduce its skewness.