

**DATA GLACIER VIRTUAL INTERNSHIP**  
**CROSS SELLING RECOMMENDATION-GROUP PROJECT**  
**WEEK 9: DELIVERABLES**  
**GROUP NAME: HEGY**

**Team members:**

**Name:** B. Harika  
**Email:** [harikabreddy444@gmail.com](mailto:harikabreddy444@gmail.com)  
**Country:** India  
**College/ Company:** Data Glacier  
**Specialization:** Data Analyst

**Name:** Yusuf Yuhan  
**Email:** [yusufyuhan98.yy@gmail.com](mailto:yusufyuhan98.yy@gmail.com)  
**Country:** Srilanka  
**College/ Company:** The Open University of Srilanka  
**Specialization:** Data Analyst

**Name:** Ebaghae Imhanlahimi  
**Email:** [imhanlahimiw@gmail.com](mailto:imhanlahimiw@gmail.com)  
**Country:** America  
**College/ Company:** Data Glacier  
**Specialization:** Data Analyst

**Name:** Gladys Kalas  
**Email:** [gladys@kalas.me](mailto:gladys@kalas.me)  
**Country:** USA  
**College/ Company:** Data Glacier  
**Specialization:** Data Analyst

## Contents

|  |    |
|--|----|
| Problem description:.....                          | 3  |
| Data understanding:.....                           | 3  |
| Type of data:.....                                 | 4  |
| Problems and solutions for the Data:.....          | 6  |
| Data Cleaning:.....                                | 8  |
| Cleaning and Transformation Done on the Data:..... | 8  |
| Ebaghae Imhanlahim.....                            | 8  |
| B. Harika.....                                     | 9  |
| Gladys Kalas.....                                  | 9  |
| Yusuf Yuhan.....                                   | 9  |
| GitHub Repo Link:.....                             | 10 |

## **Problem description:**

XYZ Credit Union is a financial institution based in Latin America that offers a variety of banking products to its customers, including credit cards, deposit accounts, retirement accounts, and safe deposit boxes. While the credit union has been successful in selling these products individually, it has not been as successful in cross-selling its products to existing customers. The lack of success in cross-selling suggests that there may be several barriers preventing XYZ Credit Union from selling additional products to its existing customers. To address this problem, XYZ Credit Union has decided to work with ABC Analytics, a data analytics consulting firm, to identify the barriers to cross-selling and develop strategies to overcome them. ABC Analytics will work with the credit union to analyze Customer data and information to identify patterns and trends, and develop targeted marketing strategies that are designed to increase their possibilities and revenues in the credit union's quest to cross sell banking products to the customers.

## **Data understanding:**

The data available for analysis was obtained from the data bank of XYZ credit union. It contains information about XYZ bank customers and the financial product holdings that XYZ offers to its customers. The data for cross-selling recommendation is a large csv file, with file size on disk of 2.13 GB. Upon primary understanding the features appear in Spanish literacy. It comprises of 48 features and 13647309 observations (The feature names are changed to English for better understanding). The dataset contains both numerical and categorical variables.

Data contains demographic characteristics of the customers:

- Age
- Address
- Income
- Etc.
- Gender
- Nationality
- Life status

Financial products offered by the bank:

- ind\_ahor\_fin\_ult1 / Saving Account
- ind\_cco\_fin\_ult1 / Current Accounts
- ind\_ctju\_fin\_ult1 / Junior Account
- ind\_deme\_fin\_ult1 / Medium-term deposits
- ind\_ecue\_fin\_ult1 / e-account
- ind\_hip\_fin\_ult1 / Mortgage
- ind\_pres\_fin\_ult1 / Loans
- ind\_tjcr\_fin\_ult1 / Credit Card
- ind\_viv\_fin\_ult1 / Home Account
- ind\_nom\_pens\_ult1 / Pensions
- ind\_aval\_fin\_ult1 / Guarantees
- ind\_cno\_fin\_ult1 / Payroll Account
- ind\_deco\_fin\_ult1 / Short-term deposits
- ind\_dela\_fin\_ult1 / Long-term deposits
- ind\_fond\_fin\_ult1 / Funds
- ind\_plan\_fin\_ult1 / Pensions
- ind\_reca\_fin\_ult1 / Taxes
- ind\_valo\_fin\_ult1 / Securities
- ind\_nomina\_ult1 / Payroll

**Type of data:**

|                     |           |
|---------------------|-----------|
| File name           | Train.csv |
| No. of observations | 13647309  |
| No. features        | 48        |
| File Size           | 2.13 GB   |
| File type           | CSV       |
| No. of files        | 1         |

## Cross-Selling Recommendation

| Column Names                            | Data Types |
|---|------------|
| fecha dato/ Date                        | Object     |
| ncodpers/ Customer code                 | Int64      |
| ind empleado/ Employee index            | Object     |
| pais residencia/ Country                | Object     |
| Sexo/ Gender                            | Object     |
| age/ Age                                | Object     |
| fecha alta/ Customer join date          | Object     |
| ind nuevo/ Customer index               | Float64    |
| antiguedad/ Customer senoirity          | Object     |
| indrel/ primary customer                | Float64    |
| ult fec cli 1t/ Customer leave date     | Object     |
| indrel 1mes/ Customer type              | Object     |
| tiprel 1mes/ Customer relation          | Object     |
| indresi/ Residence index                | Object     |
| indext/ Foreigner index                 | Object     |
| conyuemp/ Spouse index                  | Object     |
| canal entrada/ Channel                  | Object     |
| indfall/ Deceased index                 | Object     |
| tipodom/ Primary address                | Float64    |
| cod prov/ Customer address              | Float64    |
| nomprov/ province name                  | Object     |
| ind actividad cliente/ Activity index   | Float64    |
| renta/ Gross income                     | Float64    |
| segmento/ Segmentation                  | Object     |
| ind ahor fin ult1/ Saving account       | Int64      |
| ind aval fin ult1/ Guarantees           | Int64      |
| ind cco fin ult1/ Current account       | Int64      |
| ind cder fin ult1/ Derivative account   | Int64      |
| ind cno fin ult1/ Payroll account       | Int64      |
| ind ctju fin ult1/ Junior account       | Int64      |
| ind ctma fin ult1/ More private account | Int64      |
| ind ctop fin ult1/ Private account      | Int64      |
| ind ctpn fin ult1/ Private plus account | Int64      |
| ind deco fin ult1/ Short term deposits  | Int64      |
| ind deme fin ult1/ Medium term deposits | Int64      |
| ind dela fin ult1/ Long term deposits   | Int64      |
| ind ecue fin ult1/ E account            | Int64      |
| ind fond fin ult1/ Funds                | Int64      |
| ind hip fin ult1/ Mortgage              | Int64      |
| ind plan fin ult1/ Pensions             | Int64      |
| ind pres fin ult1/ Loans                | Int64      |
| ind reca fin ult1/ Taxes                | Int64      |

## Cross-Selling Recommendation

|                                 |         |
|---------------------------------|---------|
| ind tjcr fin ult1/ Credit card  | Int64   |
| ind dela fin ult1/ Securities   | Int64   |
| ind dela fin ult1/ Home account | Int64   |
| ind dela fin ult1/ Payroll      | Float64 |
| ind dela fin ult1/ Pensions 2   | Int64   |
| ind dela fin ult1/ Direct debit | Int64   |

## Problems and solutions for the Data:

| S/<br>N | Problem   | Proposed Solution                             | Reason   |
|---------|---|---|--|
| 1       | Column names recorded in Spanish                              | Rename column names in English interpretation | For ease in understanding and analyzing the data.                                |
| 2       | Column data types interpreted wrongly                         | Convert column data types                     | To ensure accuracy in analysis and improve memory efficiency                     |
| 3       | Some Column information recorded in Spanish                   | Replace records with English interpretation   | For ease in understanding and analyzing the data.                                |
| 4       | Employee_index / ind_empleado has 27734 null values           | Values to be deleted.                         | It contains categorical data and requires further information from the company.  |
| 5       | Country / pais_residencia has 27734 null values               | Values to be deleted                          | It contains demographic data and requires more information from the company.     |
| 6       | Gender / sexo has 27804 null values                           | Values to be deleted                          | It contains demographic data and requires accurate information from the company. |
| m<br>7  | Customer_join_date / fecha_alta has 27734 null values         | Value to be imputed.                          | Values can be imputed based on existing records.                                 |
| 8       | Customer_index / ind_nuevo has 27734 null values              | Value to be imputed                           | Values can be imputed based on existing records.                                 |
| 9       | Primary_customer / indrel has 27734 null values               | Value to be imputed                           | Values can be imputed based on existing records.                                 |
| 10      | Customer_leave_date / ult_fec_cli_1t has 13622516 null values | Value to be imputed                           | Values can be imputed based on existing records.                                 |
| 11      | Customer_type / indrel_1mes has 149781 null values            | To be deleted                                 | It contains categorical data and requires accurate information from the company. |
| 12      | Customer_relation / tiprel_1mes has 14781 null values         | Values to be deleted                          | It contains categorical data and requires accurate information from the company. |

## Cross-Selling Recommendation

|    |  |                      |  |
|----|--|----------------------|--|
| 13 | Residence_index / indresi<br>has 27734 null values                 | Values to be deleted | It contains demographic data and requires accurate information from the company. |
| 14 | Foreigner_index / indext<br>has 27734 null values                  | Values to be deleted | It contains demographic data and requires accurate information from the company. |
| 15 | Spouse_index / conyuemp<br>has 13645501 null values                | Values to be imputed | Values can be imputed based on existing records.                                 |
| 16 | Channel / canal_entrada<br>has 186126 null values                  | Values to be deleted | It contains categorical data and requires accurate information from the company. |
| 17 | Deceased_index / indfall<br>has 27734 null values                  | Values to be deleted | It contains demographic data and requires accurate information from the company. |
| 18 | Primary_address / tipodom<br>has 27735 null values                 | Value to be imputed  | It contains demographic data and requires accurate information from the company. |
| 19 | Customer_address /<br>Cod_prov has 93591 null values               | Values to be deleted | It contains demographic data and requires accurate information from the company. |
| 20 | Province_name / nomprov<br>has 93591 null values                   | Values to be deleted | It contains demographic data and requires accurate information from the company. |
| 21 | Activity_index /<br>ind_actividad_cliente has<br>27734 null values | Values to be imputed | Values can be imputed based on existing records.                                 |
| 22 | Gross_income / renta has<br>2794375 null values                    | Values to be imputed | Values can be imputed based on existing records.                                 |
| 23 | Segmentation / segmento<br>has 189368 null values                  | Values to be imputed | Values can be imputed based on existing records.                                 |
| 24 | Payroll / ind_nomina_ult1<br>has 16063 null values                 | Values to be imputed | Values can be imputed based on existing records.                                 |
| 25 | Pensions_2 /<br>ind_nom_pens_ult1 has<br>16063 null values         | Values to be imputed | Values can be imputed based on existing records.                                 |



## Data Cleaning:

- Several missing values have been dropped from the variables.
- Column names are translated and renamed accordingly.
- The mean, mode, median, and zeroes are used to impute null values.
- Columns like gender, residence index, spouse index, customer relations, employee index, etc.; variables are assigned to their respective categories.
- Outliers are detected using different methods and treated accordingly.

## Cleaning and Transformation Done on the Data:

### Ebaghae Imhanlahim

- Renamed column names from Spanish to English interpretation
- Checked for nulls in the data.
- Deleted null values of categorical and demographic columns with sensitive information.
- Filled some columns null values with the mode of the column.
- Filled some categorical column null values with appropriate alternative values according to the column description.
- Replaced the Spanish values in columns (Gender, Customer type, Customer relation, Employee index, Residence index, Foreigner index, Deceased index, etc.) with their English interpretations.
- Replaced categorical columns values (Activity index and segmentation) with their equivalent column values interpretations.
- Checked each column's value to ensure they are correct and correspond with description.
- Changed data types of columns as required to save memory.
- Checked for duplicates in the data.
- Checked for outliers in the data.
- Treated outliers found in the data.

### B. Harika

- Translated the columns names from Spanish to English and renamed the columns.
- Checked the data for duplications.
- Checked for the Null values in data.
- Treated the null values using Median method, using dropna function to remove few column values that were missing important information, and Fillna methods.
- Renaming the Columns values description with proper words.
- Replaced the Spanish values for columns such as: Gender [H : Male, V: Female], Activity Index, Residence Index, Foreigner Index, Spouse Index, Deceased Index these columns with ['N':No, 'S': Yes], Customer\_type, Employee\_index, Customer\_relation with their appropriate names.
- Checked for the outliers detection using the Standard Deviation method.
- Treated the outliers using .clip method.

### Gladys Kalas

- The data received is in its raw format with language in Spanish literacy, erroneous inputs, incompatible data types which need to be addressed to obtain a clean organized data for further analysis.
- Changed the feature names from Spanish literacy to English for easy understanding.
- Checked for null values in the data set Train.csv.
- Imputed some features' missing /null values with the Mean value of the respective column.
- Filled some categorical column null values with highest occurring class in the features [value\_counts().idmax()]
- Dropped columns Customer\_leave\_date and Spouse\_index as they contained 99% null values in their respective features.
- Replaced some features with their equivalent category values interpretations.  
Example- ['H'-Male, 'V'- 'Female']  
['S' to 'Yes' and 'N' to 'No']
- Changed data types of columns Age and Seniority, Customer\_join\_date etc.
- Detection of outliers using univariate visualizations, distplots, box plots.
- Detected and treated Outliers in the data using IQR method and trimming outliers.

Yusuf Yuhan

- Use `df.rename(columns={'old_col_name': 'new_col_name'})` to rename columns.
- Use `df.isnull().sum()` to check for null values.
- Use `df.fillna(method='ffill')`, `df.fillna(method='bfill')`, and `df['col_name'].fillna(df['col_name'].mode()[0])` to replace null values with forward fill, backward fill, and mode, respectively.
- Use `df['col_name'].astype(new_data_type)` to change data types for specific columns.
- Use `sns.boxplot(x=df['col_name'])` from the seaborn library to detect outliers by visualization.
- Use the interquartile range (IQR) and upper and lower caps to treat outliers by capping their values using a lambda function.

## GitHub Repo Link:

[https://github.com/HarikaReddyB/Cross\\_selling\\_recommendation---Group\\_Project/tree/main/Week%2009](https://github.com/HarikaReddyB/Cross_selling_recommendation---Group_Project/tree/main/Week%2009)