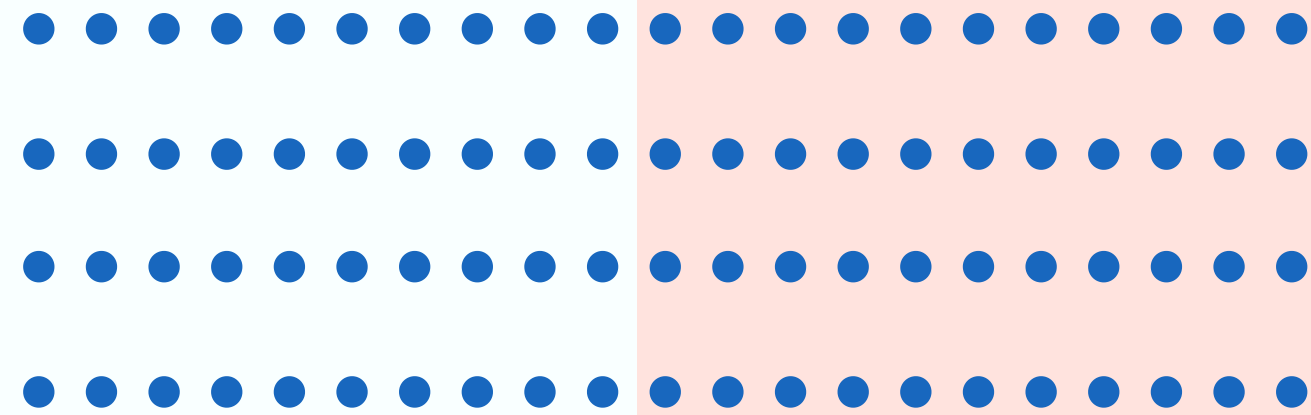


20 DESEMBER 2021



# Final Project Presentation

Derivative team

# HackerEarth\_how not to lose a customer in 10 days

## WHAT HAVE WE DONE

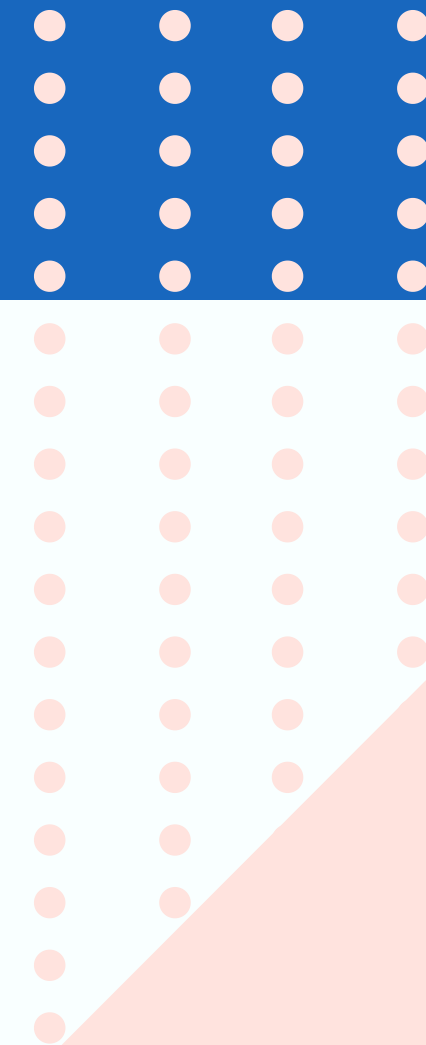
Dataset Understanding

Identify the to do list

Analyzing the Data

Pre-processing

Model development



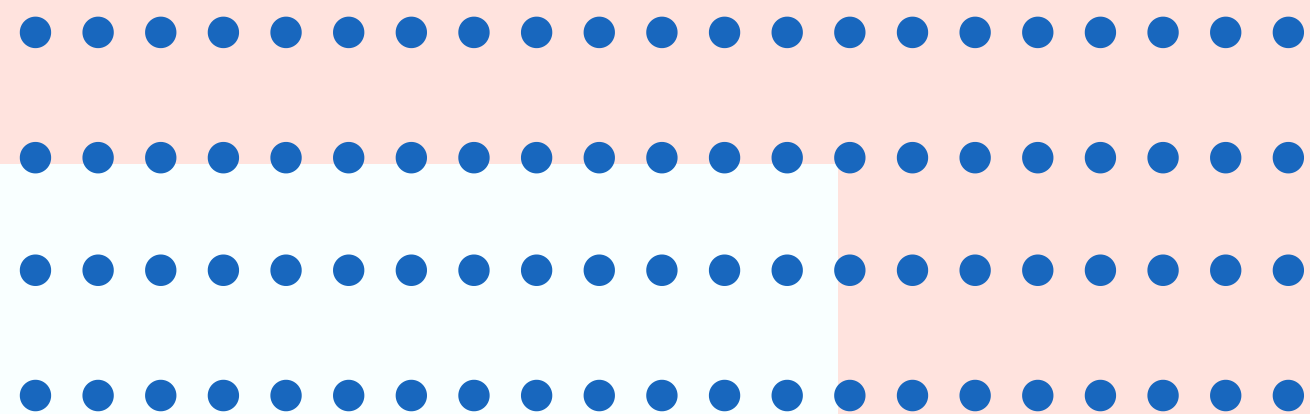
# Preface

Churn rate is a marketing metric that describes the number of customers who leave a business over a specific time period.

Churn rate value may be predicted based on multiple factors such as the user's demographic, their browsing behavior, historical purchase data, etc.

On this Final Project, We try to predict the Churn Rate Value between 1 and 5.

# Dataset Understanding



**HACKEREARTH:** HOW TO NOT  
LOSE A CUSTOMER IN 10 DAYS

THE GIVEN DATA MOSTLY  
CONSISTS OF USERS DATA  
RECORD AND ACTIVITY RECORD

THEIR RESPECTIVE LABELS  
CONSISTS OF THEIR CHURN RISK  
SCORE FROM 1 TO 5 CATEGORIC  
ORDINAL SCALE

**25 Columns:**  
– 24 Features  
– 1 Label  
**36992 Rows**

# Features (Pure Labels)

CUSTOMER\_ID

the unique identification number of a customer

NAME

the name of a customer

SECURITY\_NO

a unique security number that is used to identify a person

REFERRAL\_ID

a referral ID

# Features (Categoric Nominals)

GENDER

the gender of a customer

REGION\_CATEGORY

the region that a customer belongs to

MEMBERSHIP\_CATEGORY

the category of the membership that a customer is using

JOINED\_THROUGH\_REFERRAL

Represents whether a customer joined using any referral code or ID

derivative team

# Features (Categoric Nominals)

PREFERRED\_OFFER\_TYPES

the type of offer that a customer prefers

MEDIUM\_OF\_OPERATION

the medium of operation that a customer  
uses for transactions

INTERNET\_OPTION

the type of internet service a customer uses

USED\_SPECIAL\_DISCOUNT

Represents whether a customer uses special  
discounts offered

derivative team



# Features (Categoric Nominals)

OFFER\_APPLICATION\_PREFERENCE  
Represents whether a customer prefers offers

PAST\_COMPLAINT  
Represents whether a customer has raised  
any complaints

COMPLAINT\_STATUS  
Represents whether the complaints raised by  
a customer was resolved

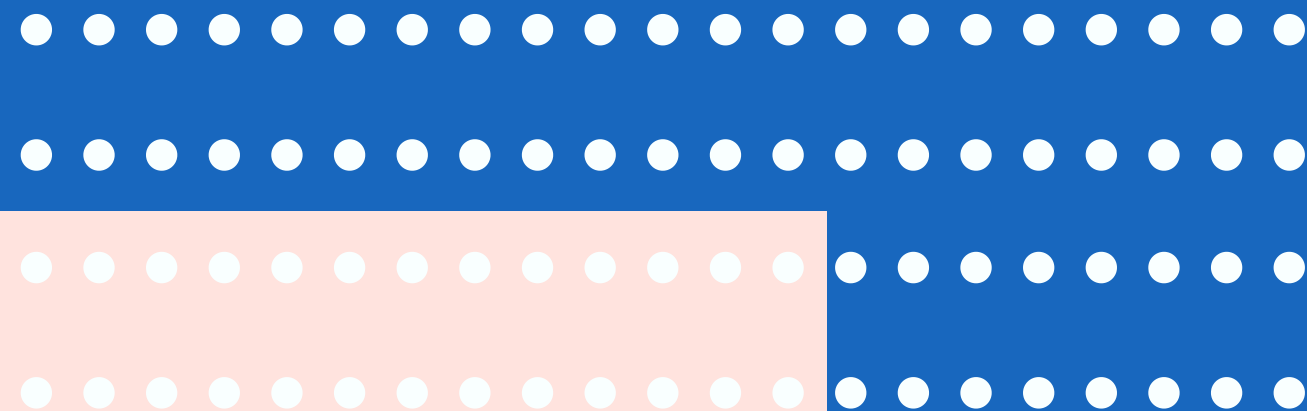
LAST\_VISIT\_TIME  
Represents whether a customer has raised  
any complaints

JOINING\_DATE  
the date when a customer became a member  
@ derivative team

# Features (Categoric Ordinals)

FEEDBACK

the feedback provided by a customer



# Features (Numeric)

AGE

the age of a customer

DAYS\_SINCE\_LAST\_LOGIN

the no. of days since a customer last logged into the website

AVG\_TIME\_SPENT

the average time spent by a customer on the website

AVG\_TRANSACTION\_VALUE

the average transaction value of a customer

AVG\_FREQUENCY\_LOGIN\_DAYS

the no. of times a customer has logged in to the website

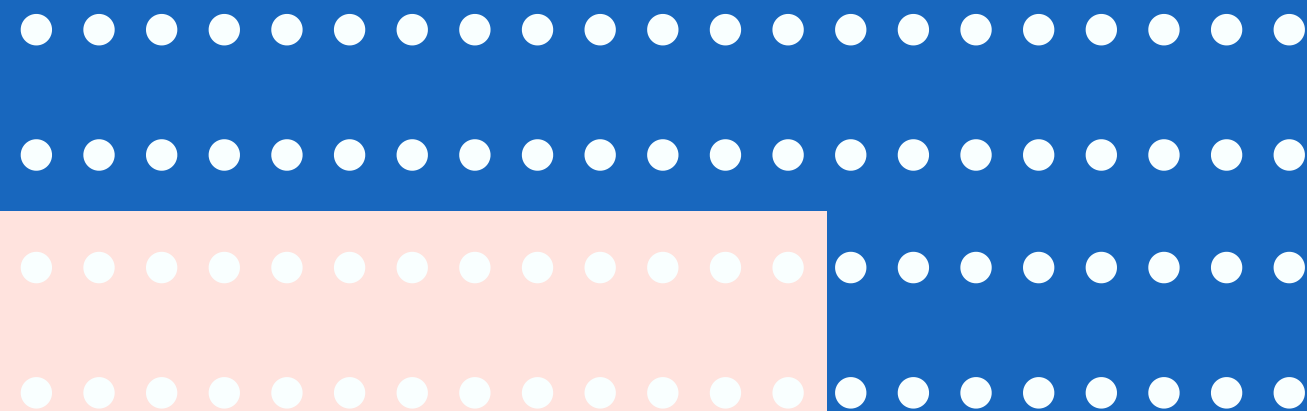
POINTS\_IN\_WALLET

the points awarded to a customer on each transaction

# Label (Categoric Ordinal)

CHURN\_RISK\_SCORE

Represents the churn risk score that ranges  
from 1 to 5



# EXPLORATORY DATA ANALYSIS

Now we are going to breakdown all of the data we have to get the better understanding of it



## CHURN RISK SCORE

There is a negative value in churn risk score column, but in application there shouldn't be any negative prediction. The lowest should be 1

## VISUALIZATION FOR THE DATA

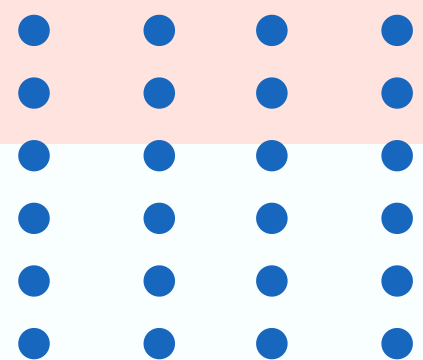
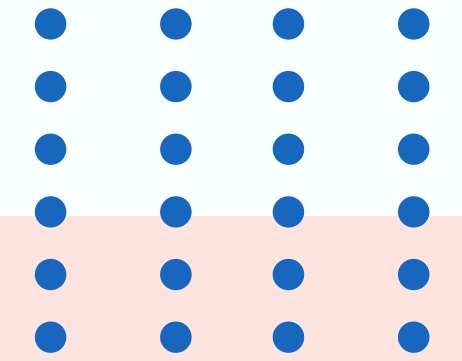
Visualize the data to further understand about the statistical condition of the data

## FEATURE DATA TYPE CLEANSING

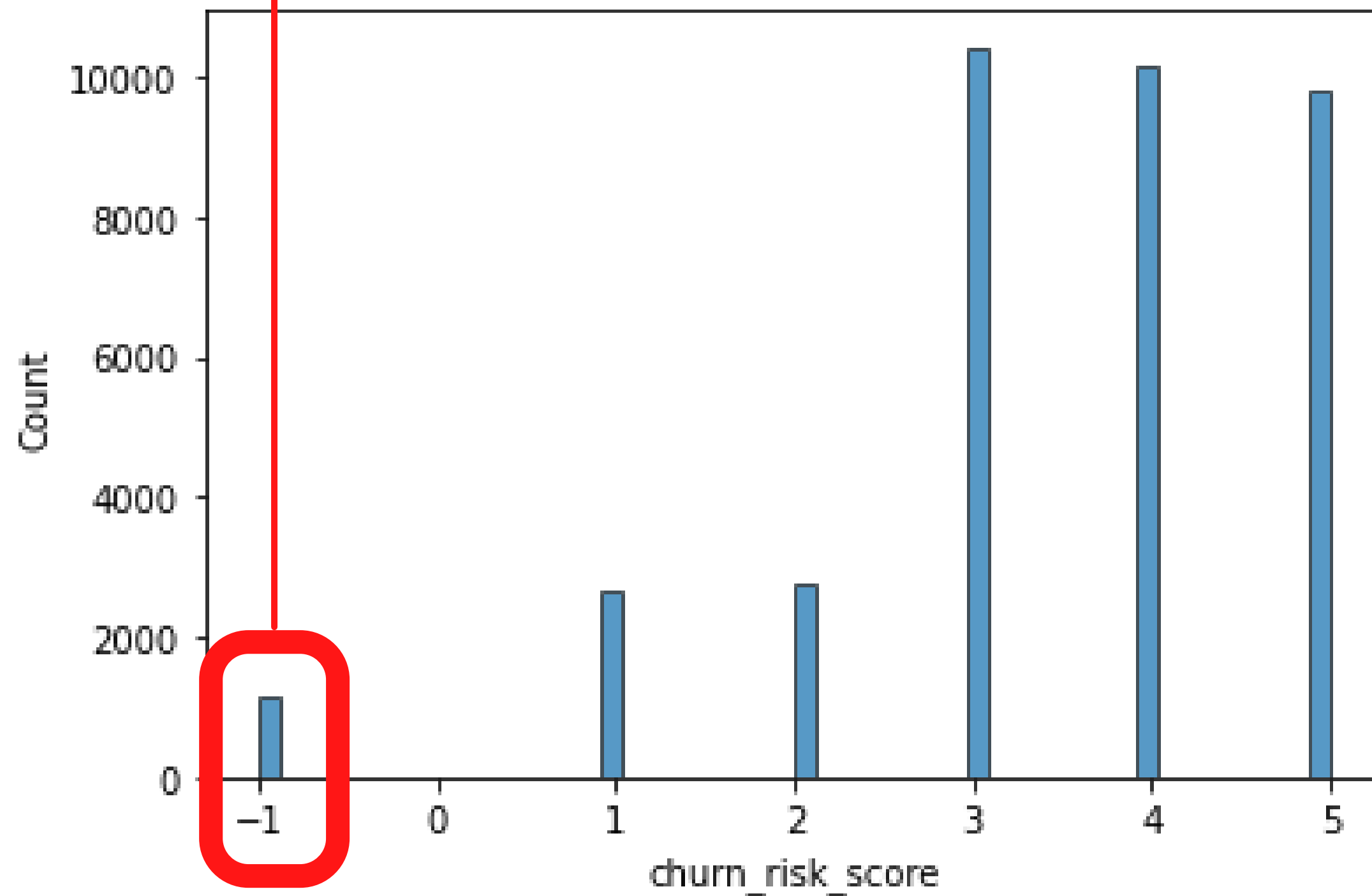
We will make sure every data there were represented by its respective true data-type

## CALCULATE THE CORRELATION BETWEEN FEATURES

Correlation counting to see how every numeric features correlate with one-another

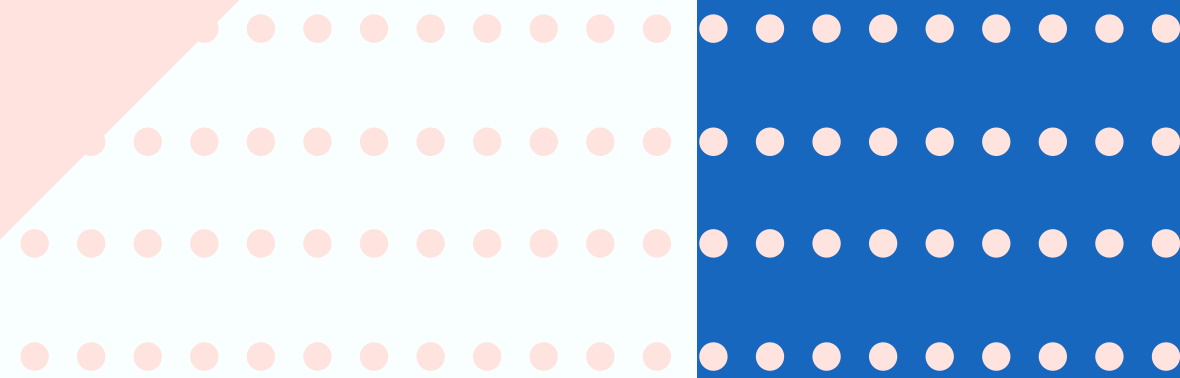


EDA  
Processing  
derivative team



There are invalid  
churn\_risk\_scores.

We will process further  
without these data



## WHAT WE DO?

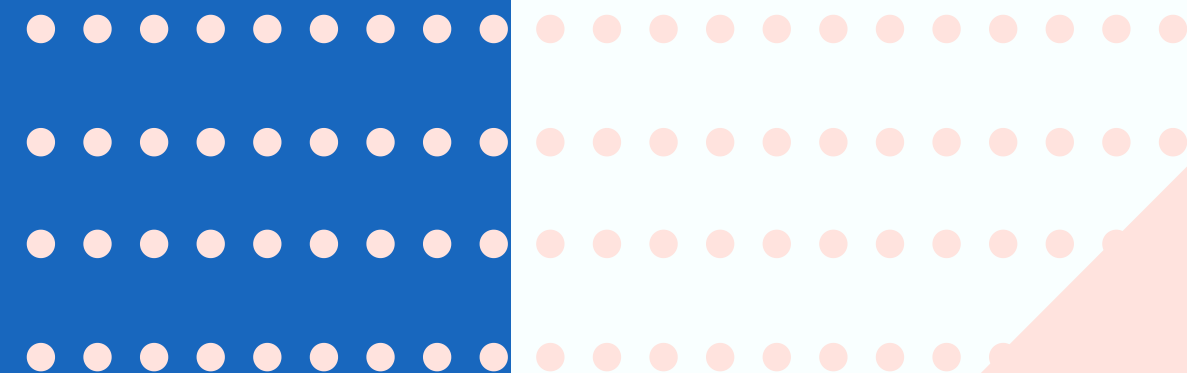
Replace the 'Error' value into NaN for avg\_frequency\_login\_days feature

Change the Data Type into Float for avg\_frequency\_login\_days feature

Change the joining\_date feature data type into datetime

## FEATURE DATA TYPE CLEANSING

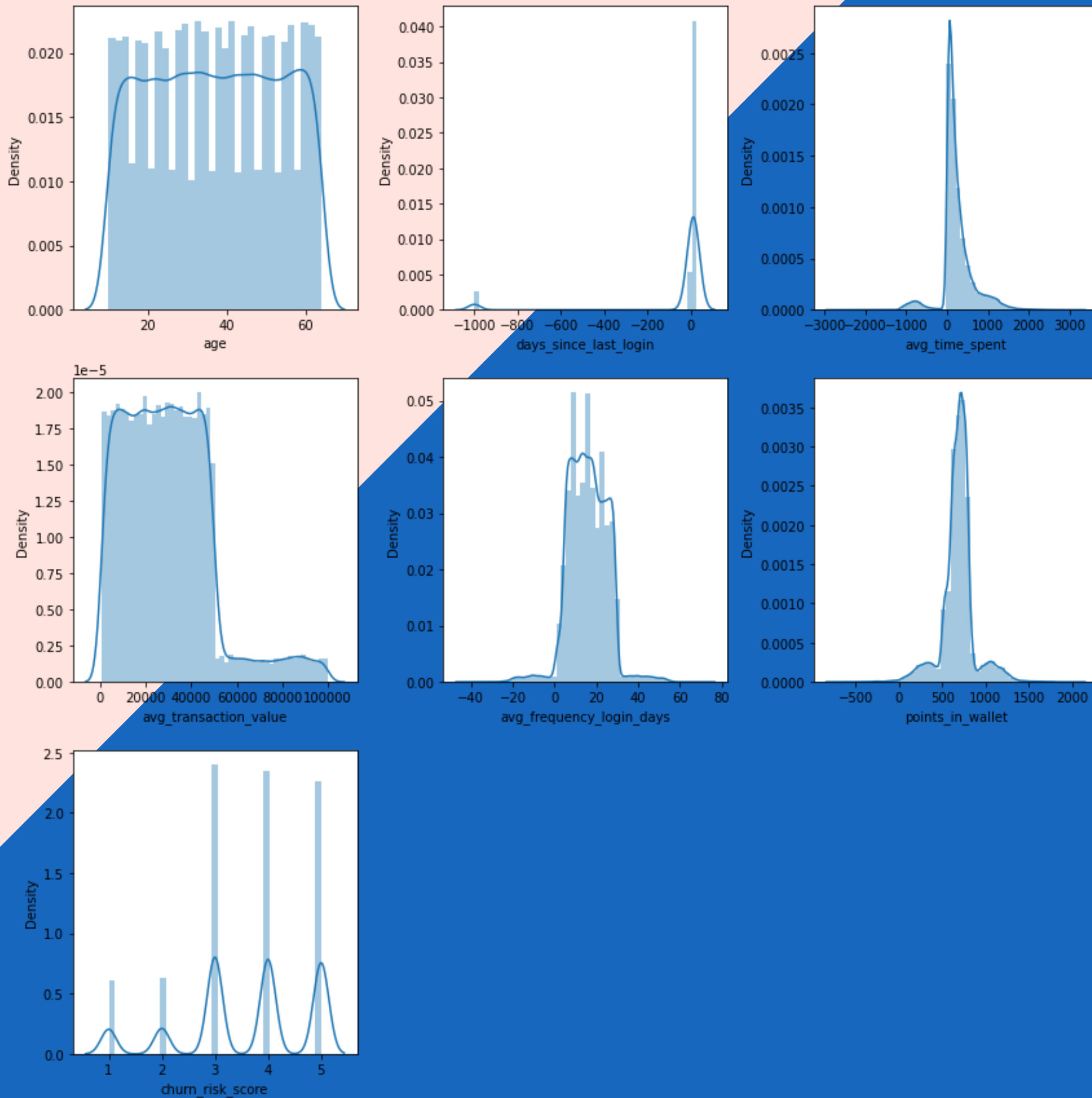
We will make sure every feature represented by their true data type (i.e.: 123 should be either int or float, ['female', 'male'] should be str or object, and '11-12-2008' into datetime64)





# DATA VISUALIZATION

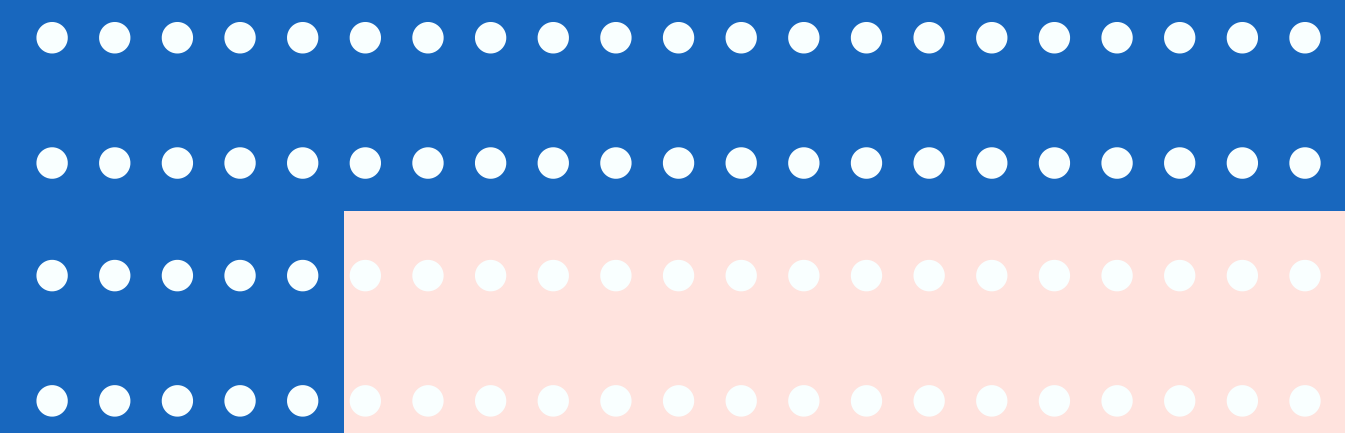
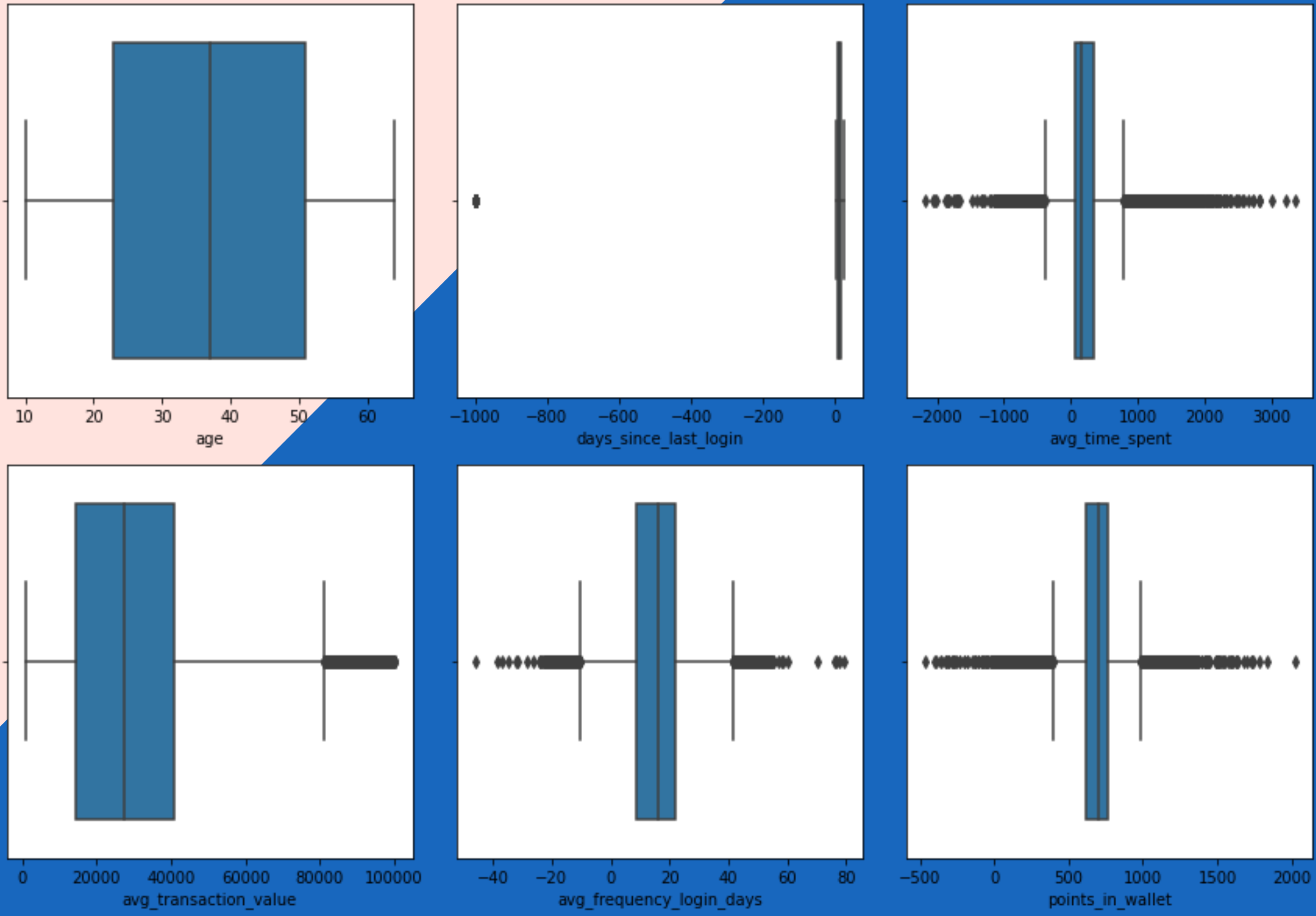
## -DISTRIBUTION PLOT-



derivative team

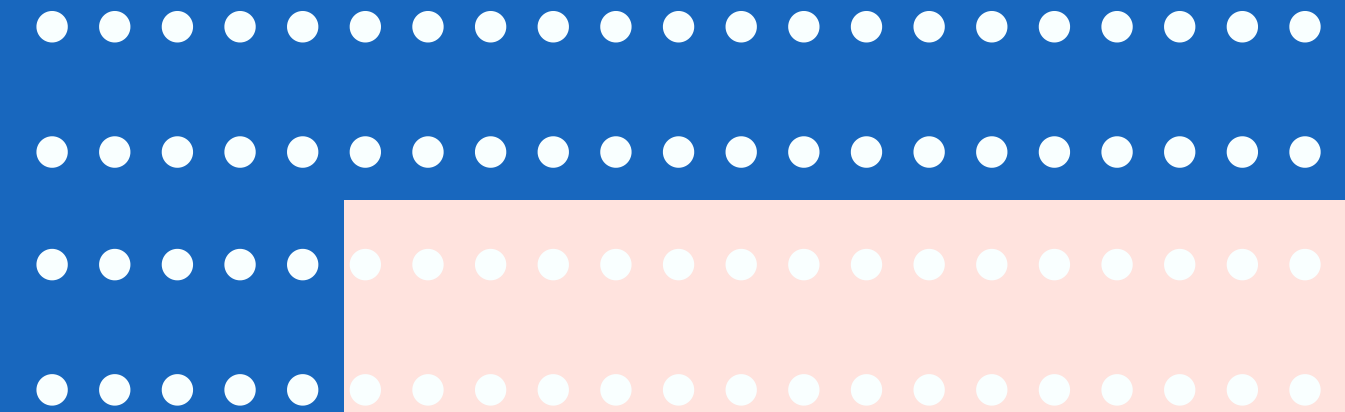
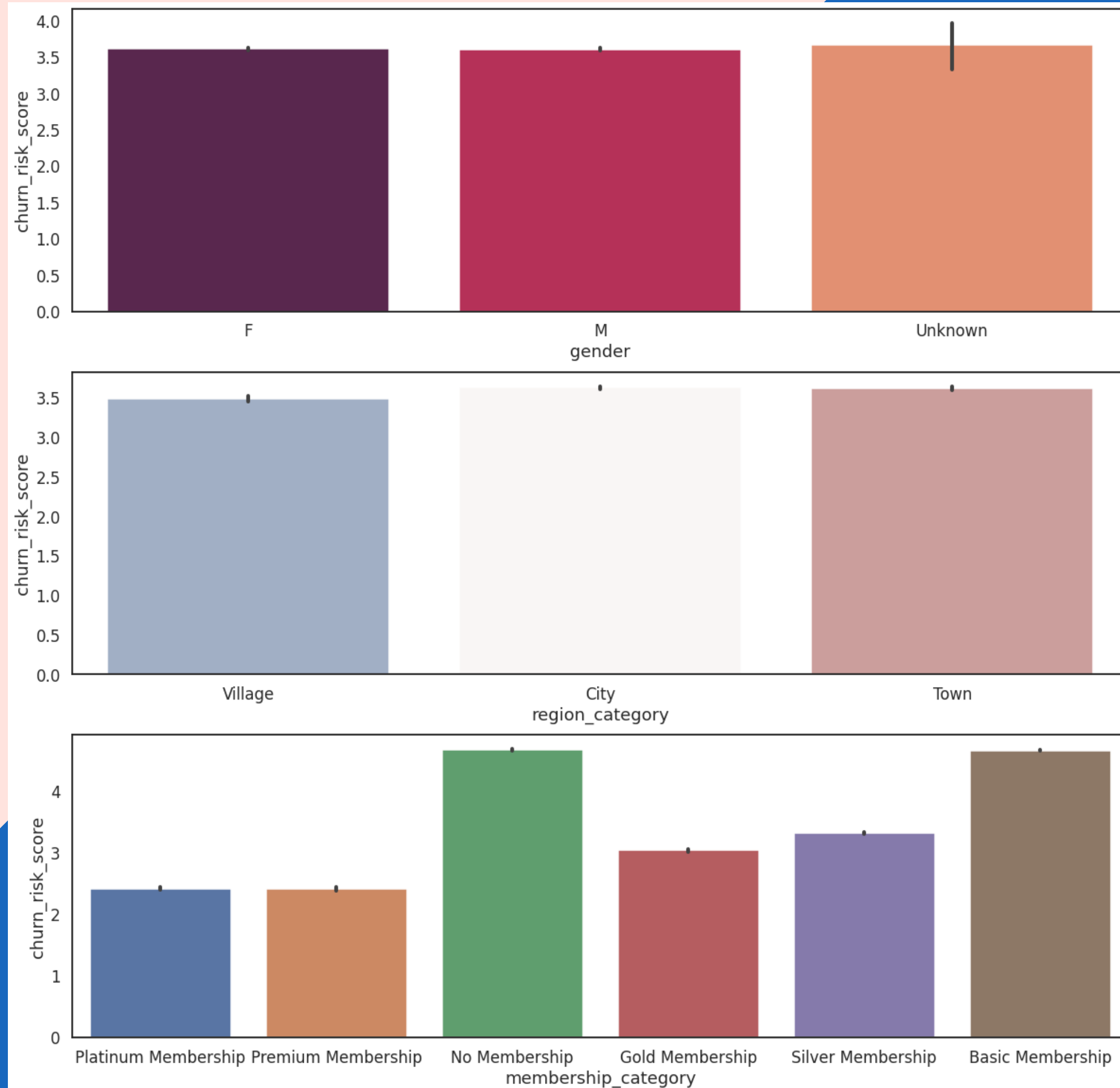
# DATA VISUALIZATION

## -BOX PLOT-

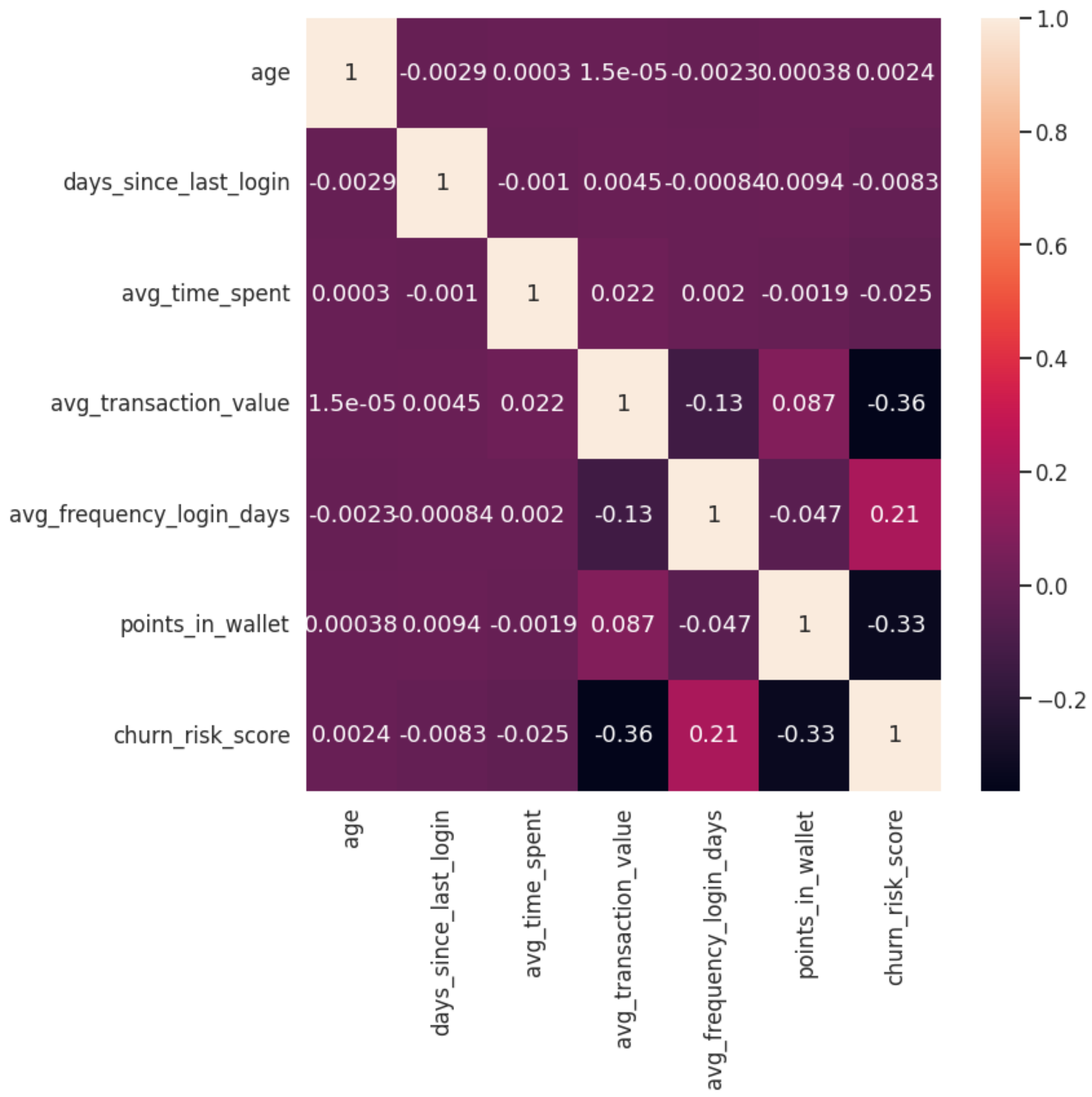


# DATA VISUALIZATION

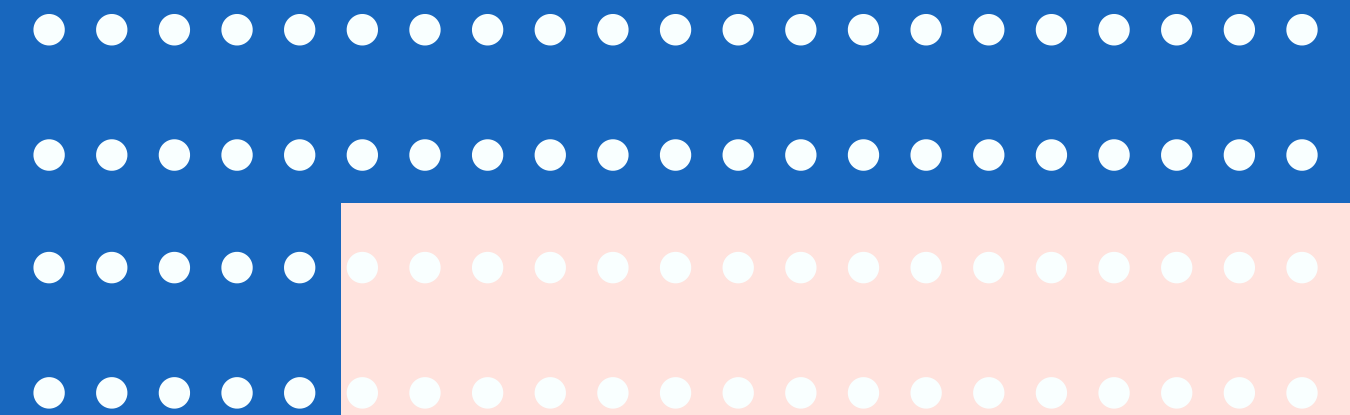
## -BAR CHART-



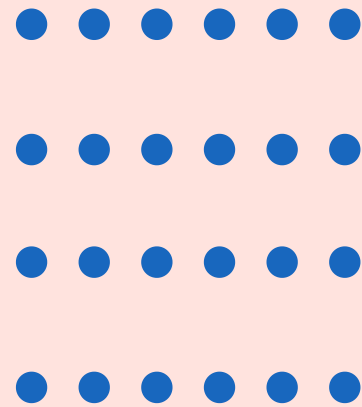
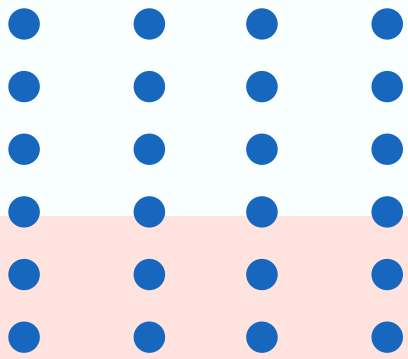
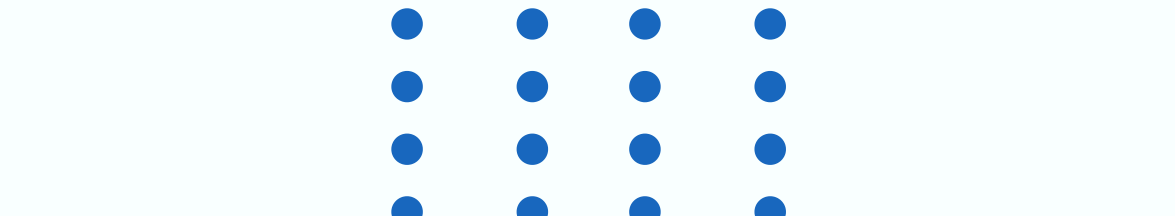
derivative team



# Calculating Correlation Between Features



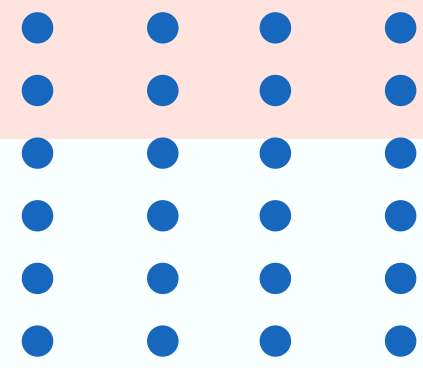
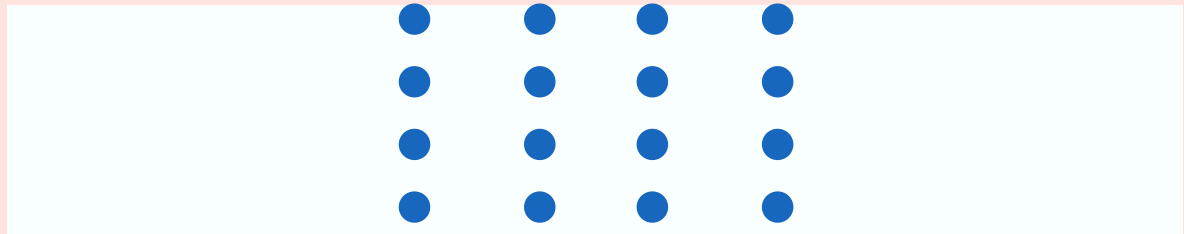
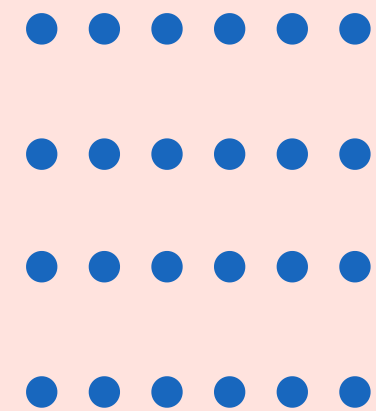
derivative team



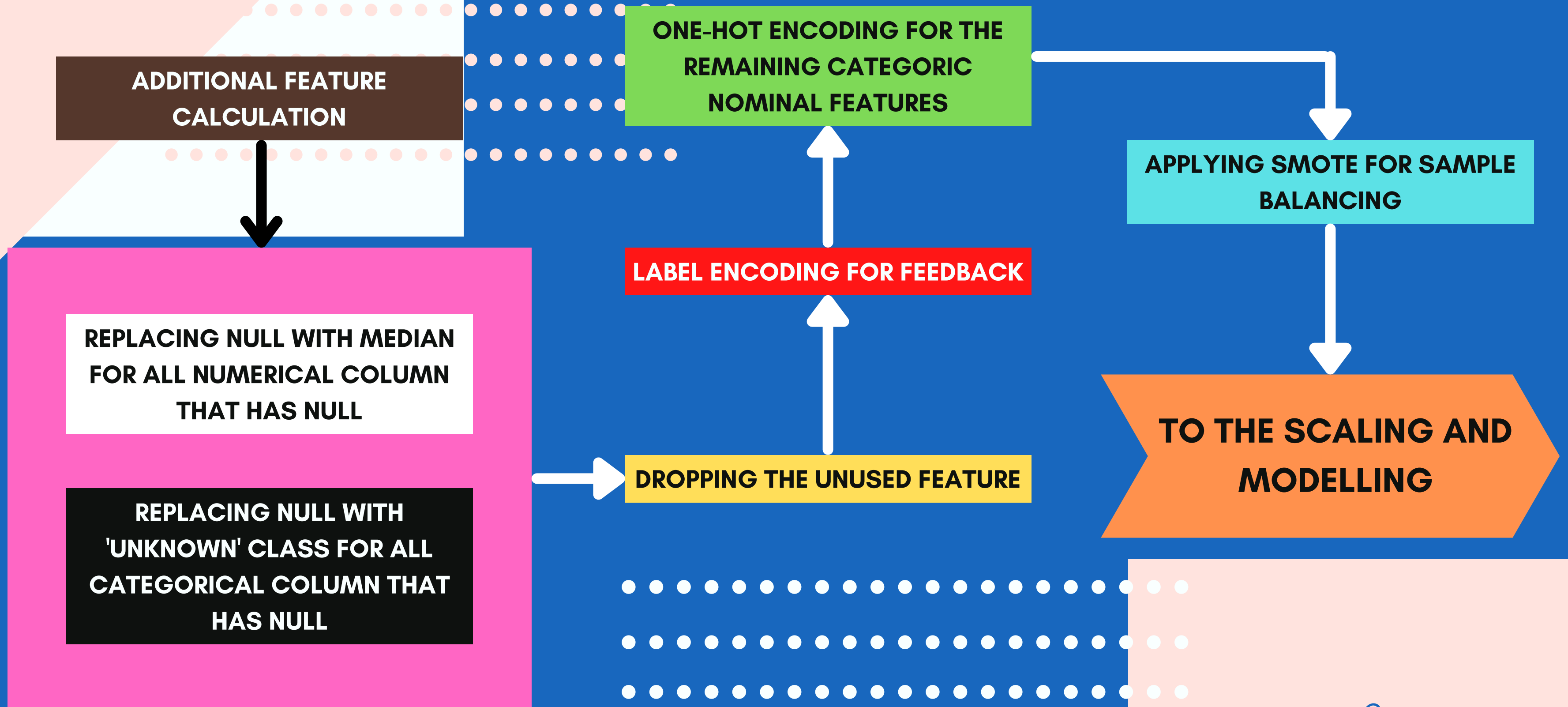
derivative team

# DATA PRE-PROCESSING

SO HOW WAS IT?



# PRE-PROCESSING STEPS



## JOINING DAYS

From **joining\_date** feature, will be extracted *the number of days* a user had been registered and the calculation will be written in a new column of **joining\_days** which has an integer data type

**JOINING\_DAYS = LAST\_JOIN\_DATE - JOINING\_DATE**

ADDITIONAL FEATURE  
CALCULATION

## REPLACE NULL WITH 'UNKNOWN' CLASS

Untuk sisa kolom fitur kategorik  
yang memiliki nilai Null:  
"region\_category",  
"preferred\_offer\_types"

## REPLACE NULL WITH MEDIAN

Untuk sisa kolom fitur numerik yang  
memiliki nilai Null:  
"avg\_frequency\_login\_days",  
"points\_in\_wallet"



# DROPPING THE UNUSED FEATURE

CUSTOMER\_ID

NAME

SECURITY\_NO

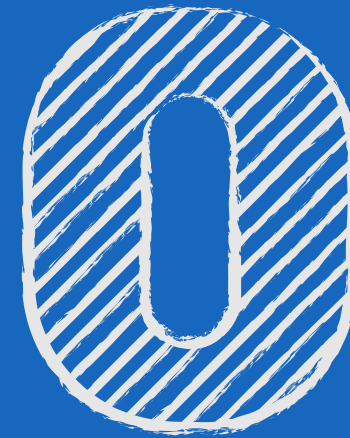
JOINING\_DATE

REFERRAL\_ID

LAST\_VISIT\_TIME

LAST\_JOIN\_DATE

# LABEL ENCODING FOR 'FEEDBACK'



'POOR WEBSITE'

'NO REASON  
SPECIFIED'

'POOR PRODUCT  
QUALITY'

'POOR CUSTOMER  
SERVICE'

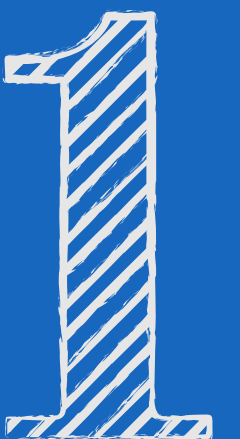
'TOO MANY ADS'

'PRODUCT ALWAYS  
IN STOCK'

"QUALITY  
CUSTOMER CARE'

'USER FRIENDLY  
WEBSITE'

'REASONABLE  
PRICE'



# ONE HOT ENCODING FOR CATEGORIC NOMINAL COLUMN



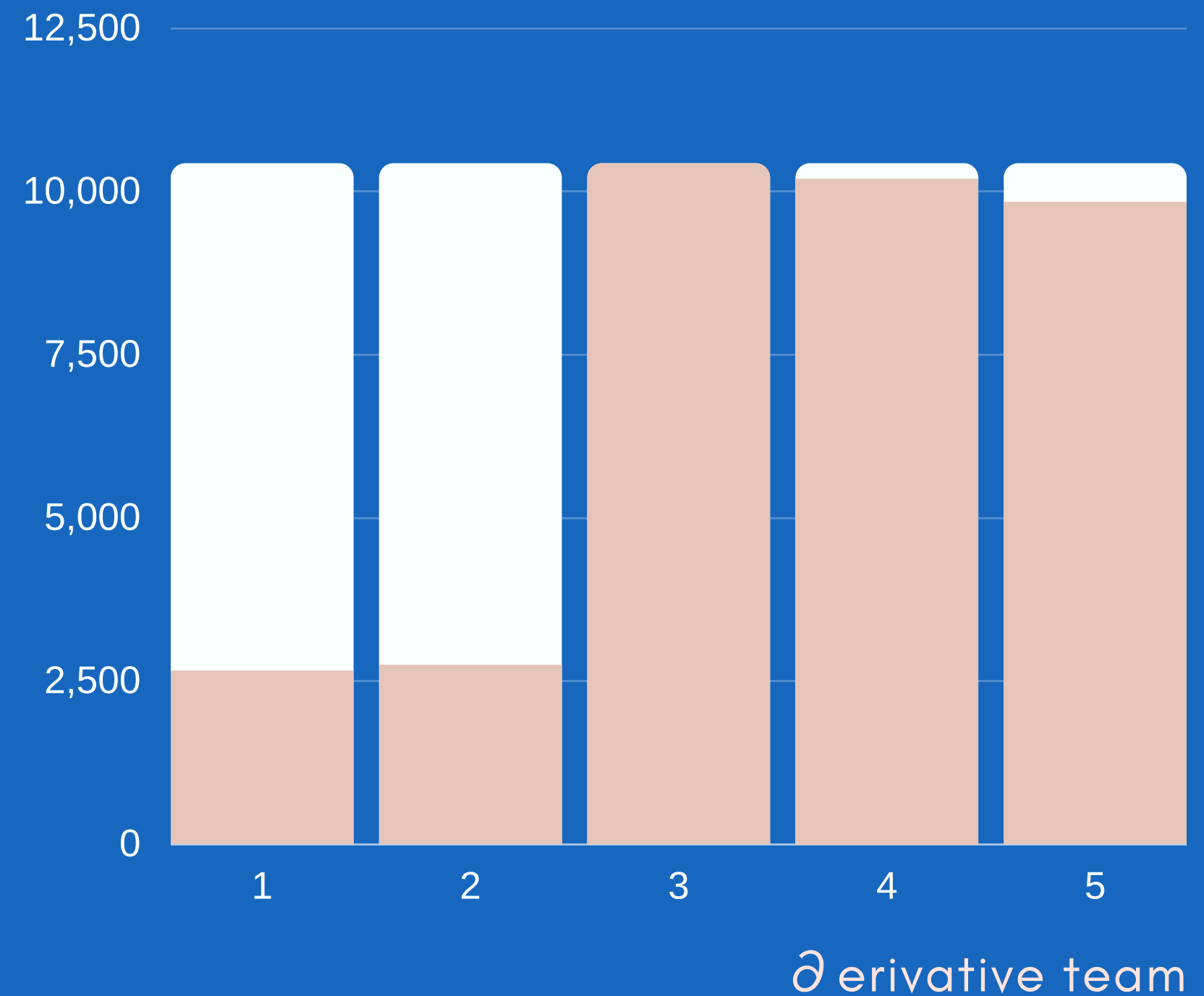
GENDER  
REGION\_CATEGORY  
MEMBERSHIP\_CATEGORY  
JOINED\_THROUGH\_REFERRAL  
PREFERRED\_OFFER\_TYPES  
MEDIUM\_OF\_OPERATION  
INTERNET\_OPTION  
USED\_SPECIAL\_DISCOUNT  
OFFER\_APPLICATION\_PREFERENCE  
PAST\_COMPLAINT  
COMPLAINT\_STATUS

# SMOTE USAGE

TO BALANCE THE SAMPLE BY  
SYNTHETIZING NEW DATA USING K-NN  
METHOD

 BEFORE OVERSAMPLING

 AFTER OVERSAMPLING

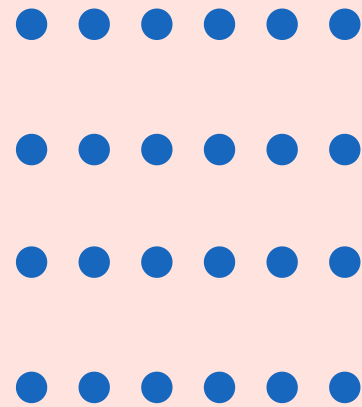
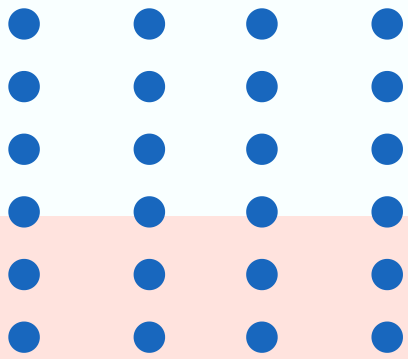


# USING ROBUST SCALING



```
array([[ -0.77042858,  0.625      ,  0.44389089, ...,  0.        ,
        0.        ,  0.        ],
       [ -0.20274436,  0.5        ,  0.46372345, ...,  1.        ,
        0.        ,  0.        ],
       [  0.28384211,  0.25       ,  1.19249174, ...,  0.        ,
        1.        ,  0.        ],
       ...,
       [ -0.33719064,  1.37205026, -0.47785962, ...,  0.        ,
        0.        ,  0.        ],
       [  0.63767368, -0.24357934,  0.92370912, ...,  0.98287825,
        0.01712175,  0.        ],
       [  0.07573499,  0.07765599, -0.37156102, ...,  0.        ,
        0.        ,  0.        ]])
```

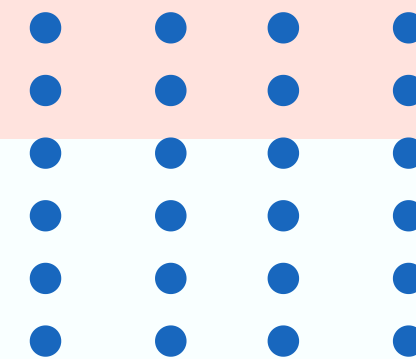
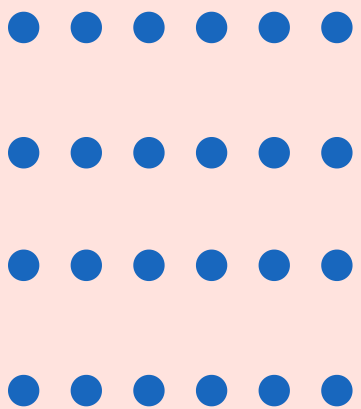
TO MAKE SURE EVERY FEATURE HAS THE  
FAIR RANGE OF VALUE



derivative team

# MODELLING

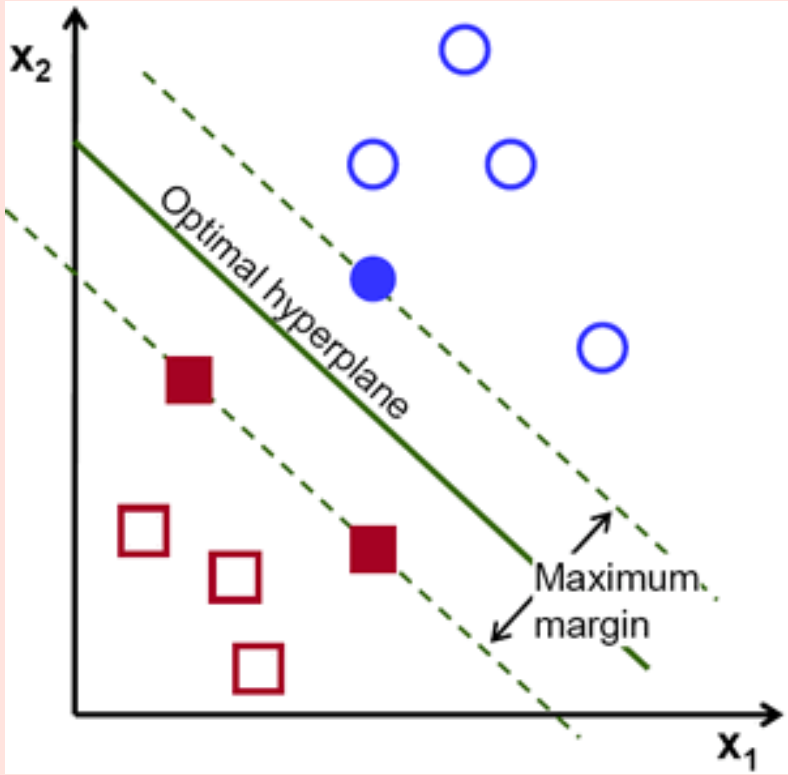
OF COURSE NOT A CATWALK MODEL



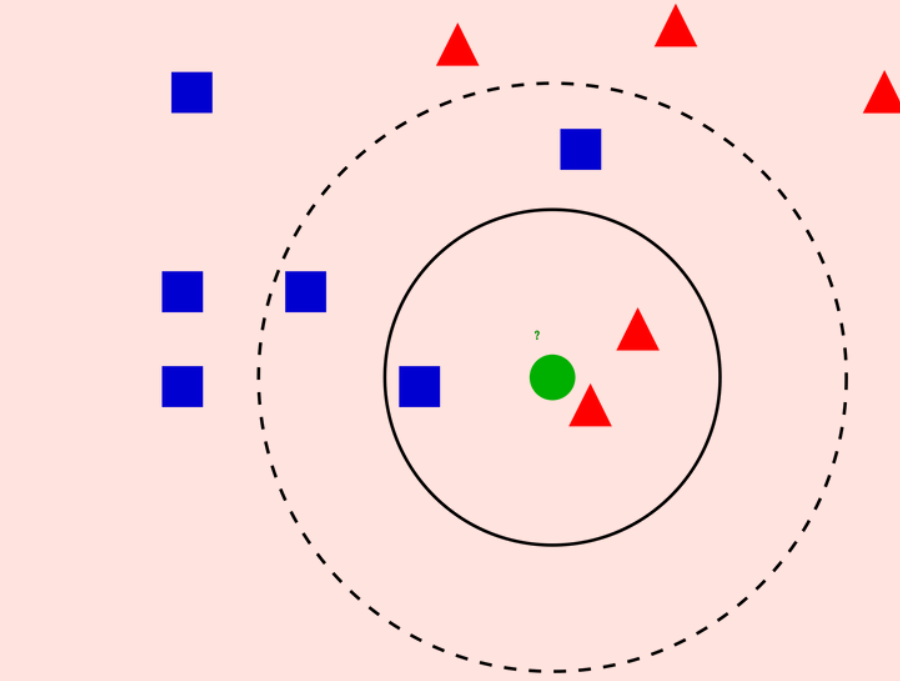
# METHOD USED



RANDOM FOREST



SVM



K-NEAREST NEIGHBOR

*dmlc*  
***XGBoost***

Derivative team

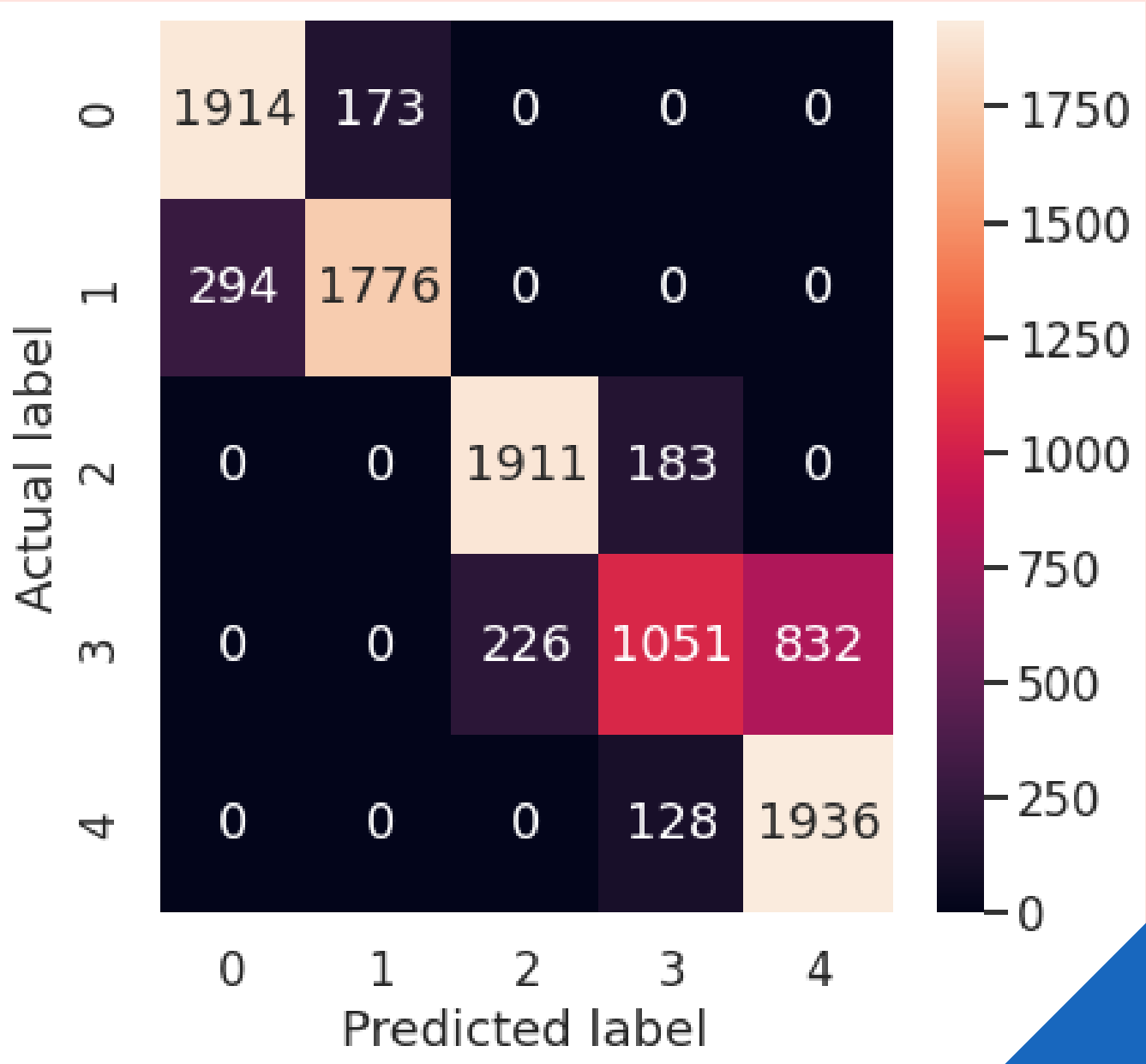
# RANDOM FOREST

f1-score :0.8170634148645475

accuracy : 0.8238679969301612

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 1            | 0.87      | 0.92   | 0.89     | 2087    |
| 2            | 0.91      | 0.86   | 0.88     | 2070    |
| 3            | 0.89      | 0.91   | 0.90     | 2094    |
| 4            | 0.77      | 0.50   | 0.61     | 2109    |
| 5            | 0.70      | 0.94   | 0.80     | 2064    |
| accuracy     |           |        | 0.82     | 10424   |
| macro avg    | 0.83      | 0.82   | 0.82     | 10424   |
| weighted avg | 0.83      | 0.82   | 0.82     | 10424   |

|                         |          |          |
|-------------------------|----------|----------|
| Submission ID: 67376341 | Result   | Score    |
| 21 seconds ago          | Accepted | 75.28612 |



derivative team

Our  
Model



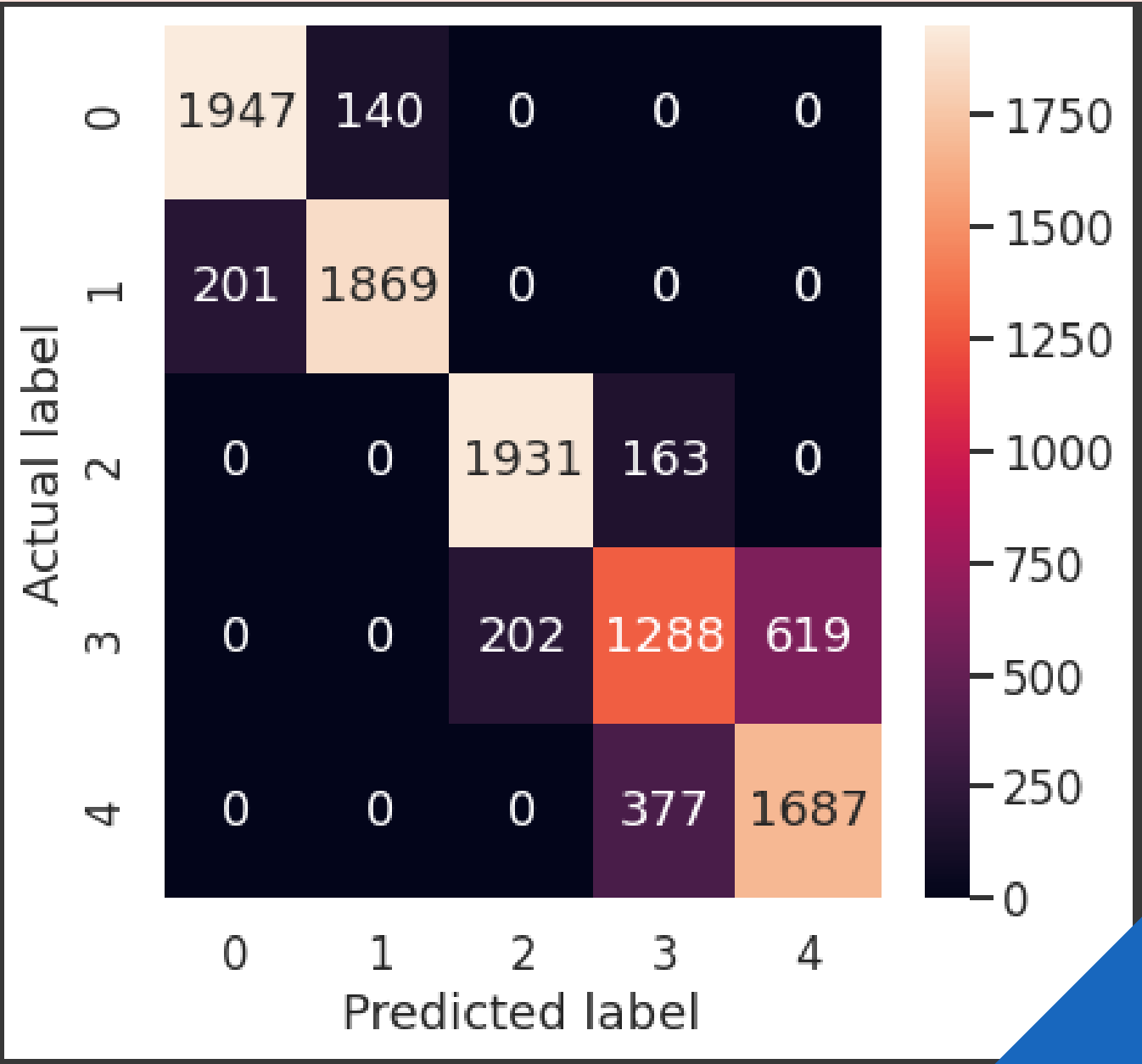
# XGBOOST

f1-score : 0.8351839259364959

accuracy : 0.8367229470452802

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 1            | 0.91      | 0.93   | 0.92     | 2087    |
| 2            | 0.93      | 0.90   | 0.92     | 2070    |
| 3            | 0.91      | 0.92   | 0.91     | 2094    |
| 4            | 0.70      | 0.61   | 0.65     | 2109    |
| 5            | 0.73      | 0.82   | 0.77     | 2064    |
| accuracy     |           |        | 0.84     | 10424   |
| macro avg    | 0.84      | 0.84   | 0.84     | 10424   |
| weighted avg | 0.84      | 0.84   | 0.83     | 10424   |

|                         |          |          |
|-------------------------|----------|----------|
| Submission ID: 67370677 | Result   | Score    |
| 15 seconds ago          | Accepted | 75.63348 |



derivative team

Our  
Model

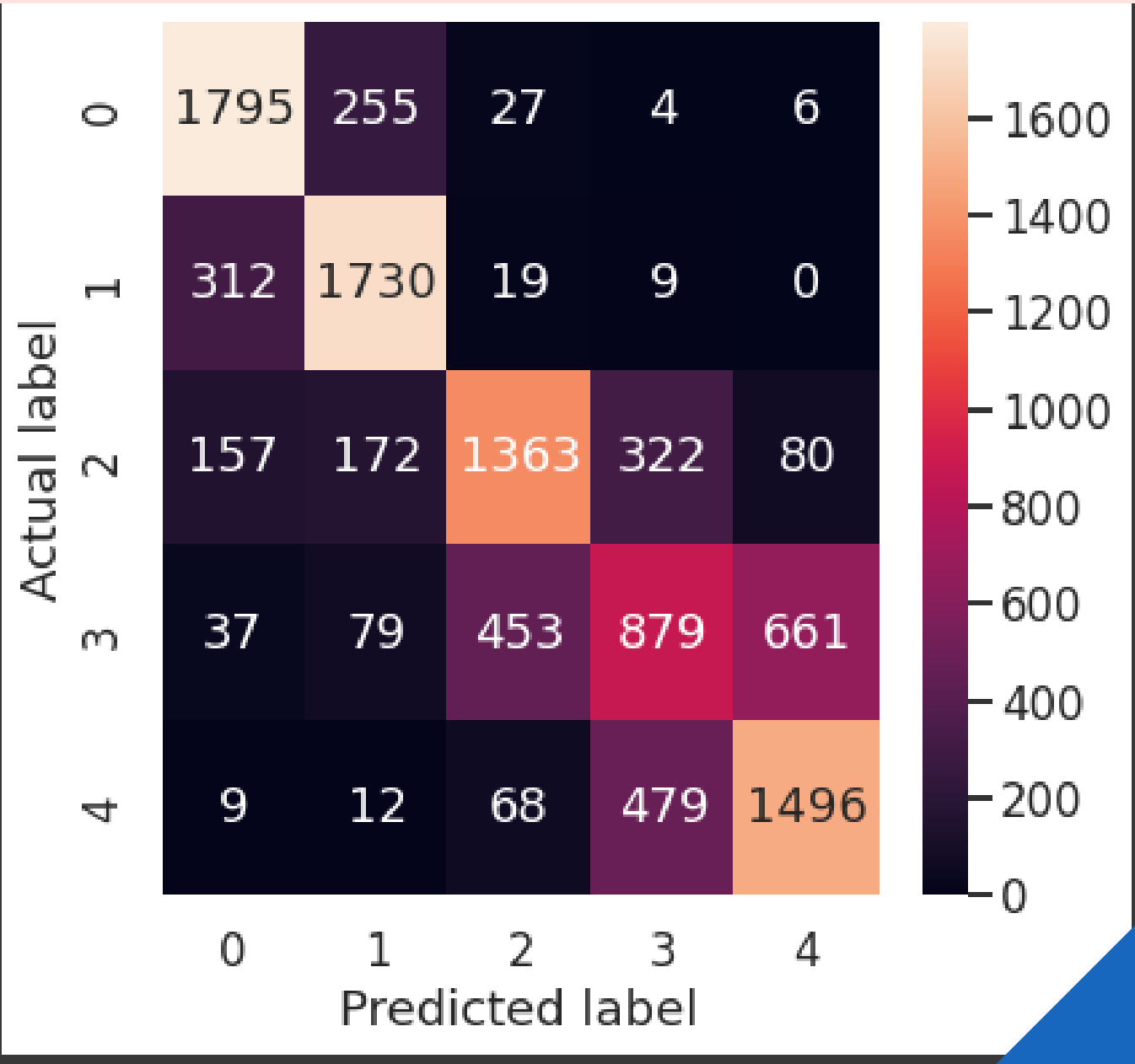
# KNN

f1-score : 0.6904538694592525

accuracy : 0.6967574827321565

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 1            | 0.78      | 0.86   | 0.82     | 2087    |
| 2            | 0.77      | 0.84   | 0.80     | 2070    |
| 3            | 0.71      | 0.65   | 0.68     | 2094    |
| 4            | 0.52      | 0.42   | 0.46     | 2109    |
| 5            | 0.67      | 0.72   | 0.69     | 2064    |
| accuracy     |           |        | 0.70     | 10424   |
| macro avg    | 0.69      | 0.70   | 0.69     | 10424   |
| weighted avg | 0.69      | 0.70   | 0.69     | 10424   |

|                         |          |          |
|-------------------------|----------|----------|
| Submission ID: 67376491 | Result   | Score    |
| 19 seconds ago          | Accepted | 53.34862 |



derivative team

Our  
Model

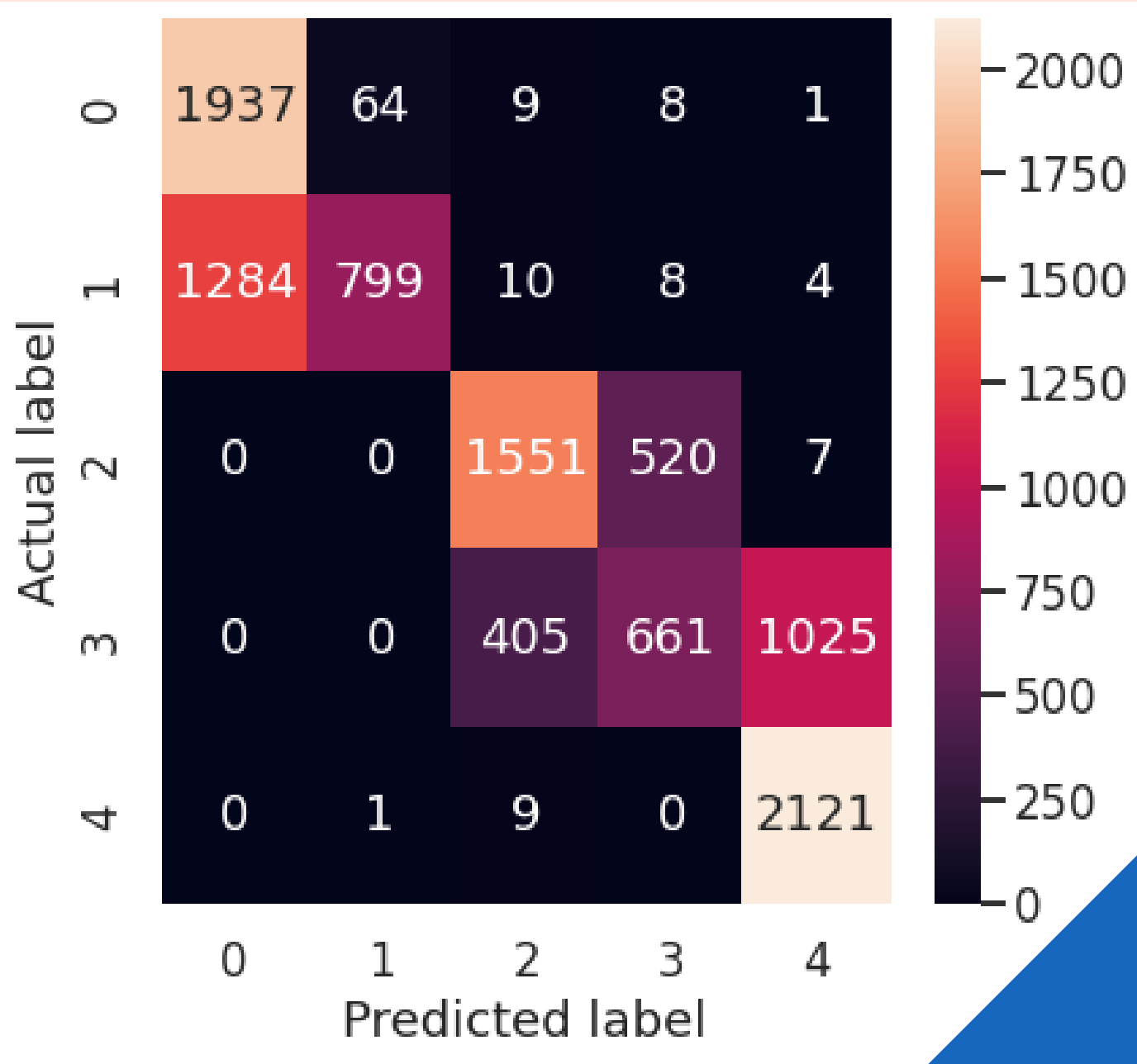
# SVM

f1-score : 0.6490629314667464

accuracy : 0.6781465848042978

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 1            | 0.60      | 0.96   | 0.74     | 2019    |
| 2            | 0.92      | 0.38   | 0.54     | 2105    |
| 3            | 0.78      | 0.75   | 0.76     | 2078    |
| 4            | 0.55      | 0.32   | 0.40     | 2091    |
| 5            | 0.67      | 1.00   | 0.80     | 2131    |
| accuracy     |           |        | 0.68     | 10424   |
| macro avg    | 0.71      | 0.68   | 0.65     | 10424   |
| weighted avg | 0.71      | 0.68   | 0.65     | 10424   |

|                         |          |          |
|-------------------------|----------|----------|
| Submission ID: 67377054 | Result   | Score    |
| 6 seconds ago           | Accepted | 61.44761 |



derivative team

Our  
Model

# Summary

XGBoost and Random Forest have highest accuracy both on test and validation dataset (about 83 % on test data and 75% on validation data)

While KNN and SVM don't performed well on this dataset

(KNN has 53.35 % accuracy and SVM has 61.45 % accuracy on validation dataset )

Derivative team



derivative team

# Temui Tim Kami



BENNITO  
SIANTURI



YUDHI  
NUGRAHA  
RIYANSYAH



MUHAMMAD  
GALANG  
GARDAMUKTI



ALDRICHO A.  
POLLARDO