

# **International Institute of Informational technology**

## **AI511 Machine Learning**

Project: Its Fraud

Team enigma

### **Team members**

Sajal Gupta  
Darshitkumar Pipariya

MT2022096  
MT2022035

# Table of Content

|  |           |
|--|-----------|
| <b>1. Introduction</b>   | <b>3</b>  |
| 1.1 Problem Statement  | 3         |
| 1.2 Approach Overview  | 3         |
| <b>2. Dataset Description</b>  | <b>4</b>  |
| 2.1 Files  | 4         |
| 2.2 Columns  | 4         |
| <b>3. Exploratory Data Analysis</b>  | <b>5</b>  |
| 3.1 Object columns   | 5         |
| 3.1.1 ProductCD  | 5         |
| 3.1.2 card4  | 6         |
| 3.1.3 card6  | 7         |
| 3.1.4 P_emaildomain  | 8         |
| 3.1.5 M6   | 9         |
| 3.2 Numerical columns  | 10        |
| 3.2.1 card1 ,card2 , card3 , card5   | 10        |
| 3.2.2 addr1 and addr2  | 14        |
| 3.2.3 D1 , D4 , D10 , D15  | 16        |
| 3.2.4 C1-C14 columns   | 20        |
| 3.2.5 Vxxx: Vesta engineered rich features, including ranking, counting, and other entity relations. | 21        |
| <b>4. Training Models</b>  | <b>22</b> |
| <b>Conclusion</b>  | <b>23</b> |

# 1. Introduction

## 1.1 Problem Statement

The traditional approaches used for detecting transaction fraud involve manual monitoring through humans which also involves interaction with the card holders. This approach is not efficient as it not only consumes a lot of time but it is also quite expensive in terms of the resources it uses to detect the fraudulence of a transaction. Moreover, in the real world the majority of the transactions being manually monitored turn out to be legit and only a few of them turn out to be fraudulent which wastes a lot of time and energy. And, hence in this competition, we need to develop an automated screening system that requires minimalistic human intervention to detect the legitimacy of transactions using Machine Learning.

We will be predicting whether a transaction is fraudulent or not and to predict this we will be building an ML model using a real-world e-commerce dataset.

## 1.2 Approach Overview

Following are the steps followed in Project

- Here the goal is to create a classification model that can predict whether the transaction is fraudulent or not.
- Understand Dataset and perform the EDA on Data and preprocessing of Data.
- Train the Model that has the best accuracy, F1-score and AUC score for that do the following
  - Use a different classification model
  - Use sampling techniques,
  - Use cross-validation
  - Use Perform hyperparameter tuning for all models

## 2. Dataset Description

### 2.1 Files

- **train.csv**: The training dataset comprising the transaction details, identity data and target label for fraud transactions.
- **test.csv**: same as train dataset but without target label.
- **sample\_submission.csv** - a sample submission file in the correct format

### 2.2 Columns

- **TransactionDT**: timedelta from a given reference DateTime (not an actual timestamp) "TransactionDT "corresponds to the number of seconds in a day.
- **TransactionAMT**: transaction payment amount in USD
- **ProductCD**: product code, the product for each transaction
- **card1 - card6**: payment card information, such as card type, card category, issue bank, country, etc.
- **addr**: address  
"both addresses are for purchaser  
addr1 as billing region  
addr2 as billing country"
- **dist**: distances between (not limited) billing address, mailing address, zip code, IP address, phone area, etc."
- **P\_ and (R\_) emaildomain**: purchaser and recipient email domain  
"certain transactions don't need a recipient, so R\_emaildomain is null."
- **C1-C14**: counting, such as how many addresses are found to be associated with the payment card, etc.
- **D1-D15**: timedelta, such as days between the previous transaction, etc.
- **M1-M9**: match, such as names on card and address, etc.
- **Vxxx**: Vesta engineered rich features, including ranking, counting, and other entity relations.
- **id01 to id11** are numerical features for identity.  
Other columns such as network connection information (IP, ISP, Proxy, etc), and digital signature (UA/browser/os/version, etc) associated with transactions are also present

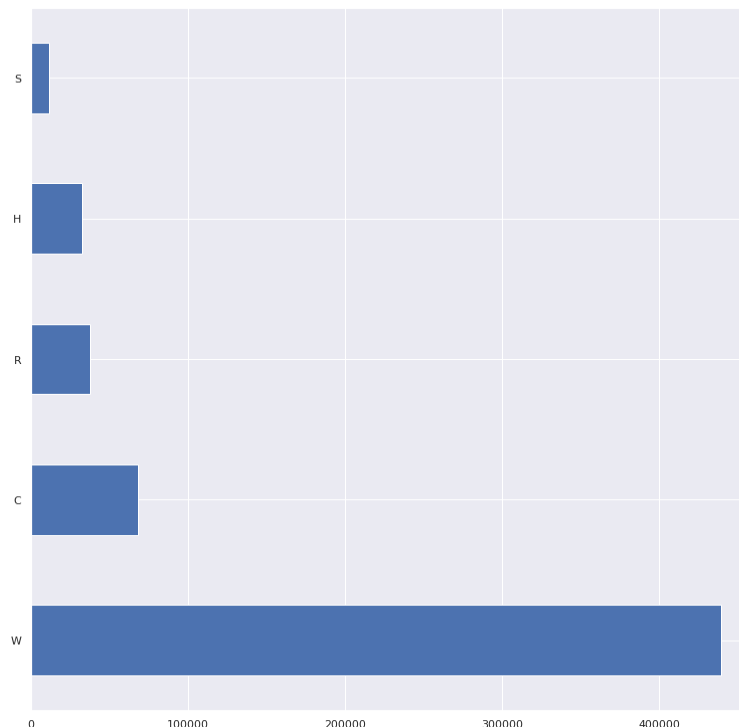
## 3. Exploratory Data Analysis

- In train data we have 434 columns and 442905 rows and In test data we have 433 columns and 147635 rows. We have contacted both train and test data column wise to perform the same data processing.
- We have 400 columns of float64 type, 3 columns of int64 and 31 columns of object type.
- Memory usage is 1.9+ GB and to run a model on it we have to reduce it.
- We have 232 columns that have null value more than 40%.we are dropping them now we have 202 column left

### 3.1 Object columns

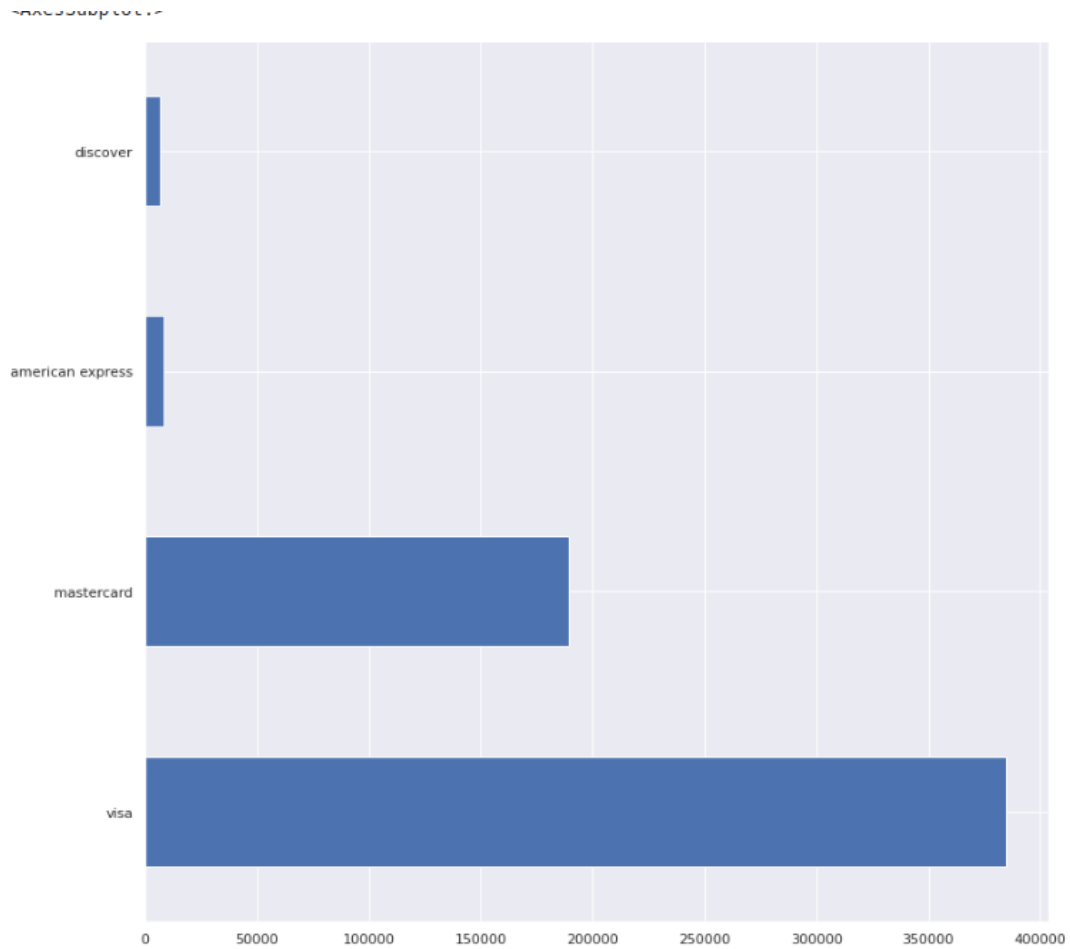
#### 3.1.1 ProductCD

It does not have any null value. It has 5 categories. We have used a one-hot encoder for this column.



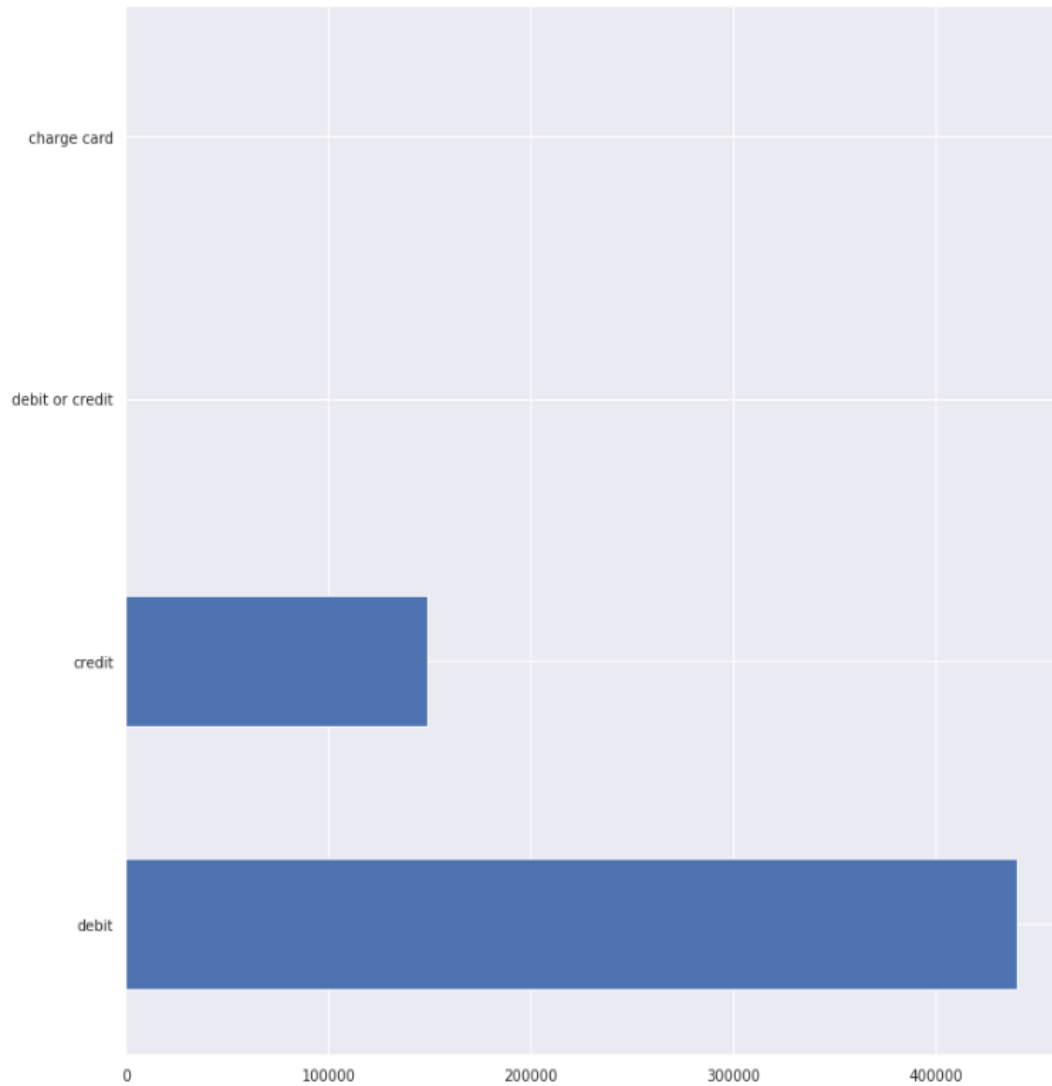
### 3.1.2 card4

It has 1577 null values.to fill null we have used mode of card4.it has 4 categories .we have used one-hot encoder



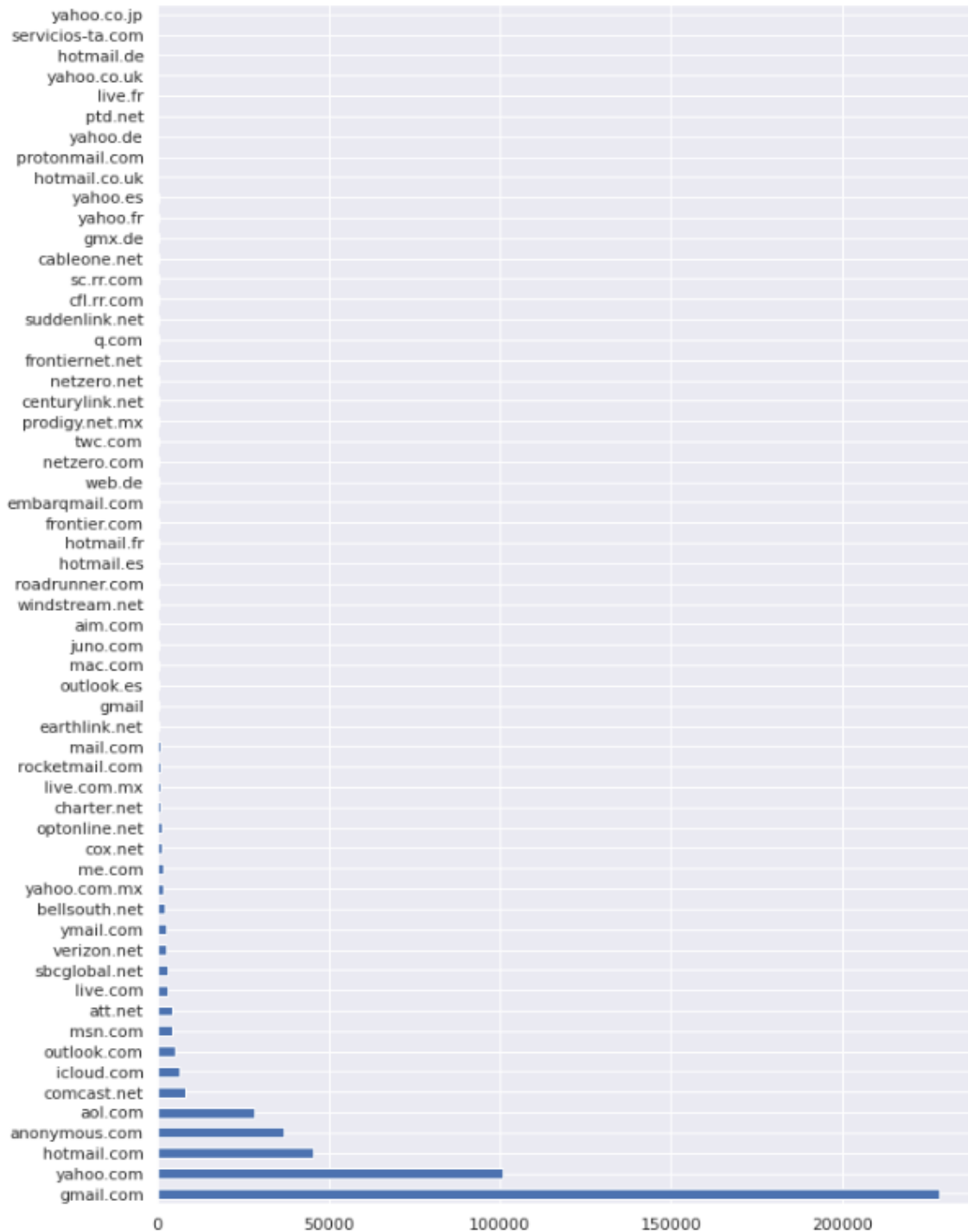
### 3.1.3 card6

It has 1571 null values.to fill null we have used mode of card4.it has 4 categories.we have used one-hot encoder.



### 3.1.4 P\_emaildomain

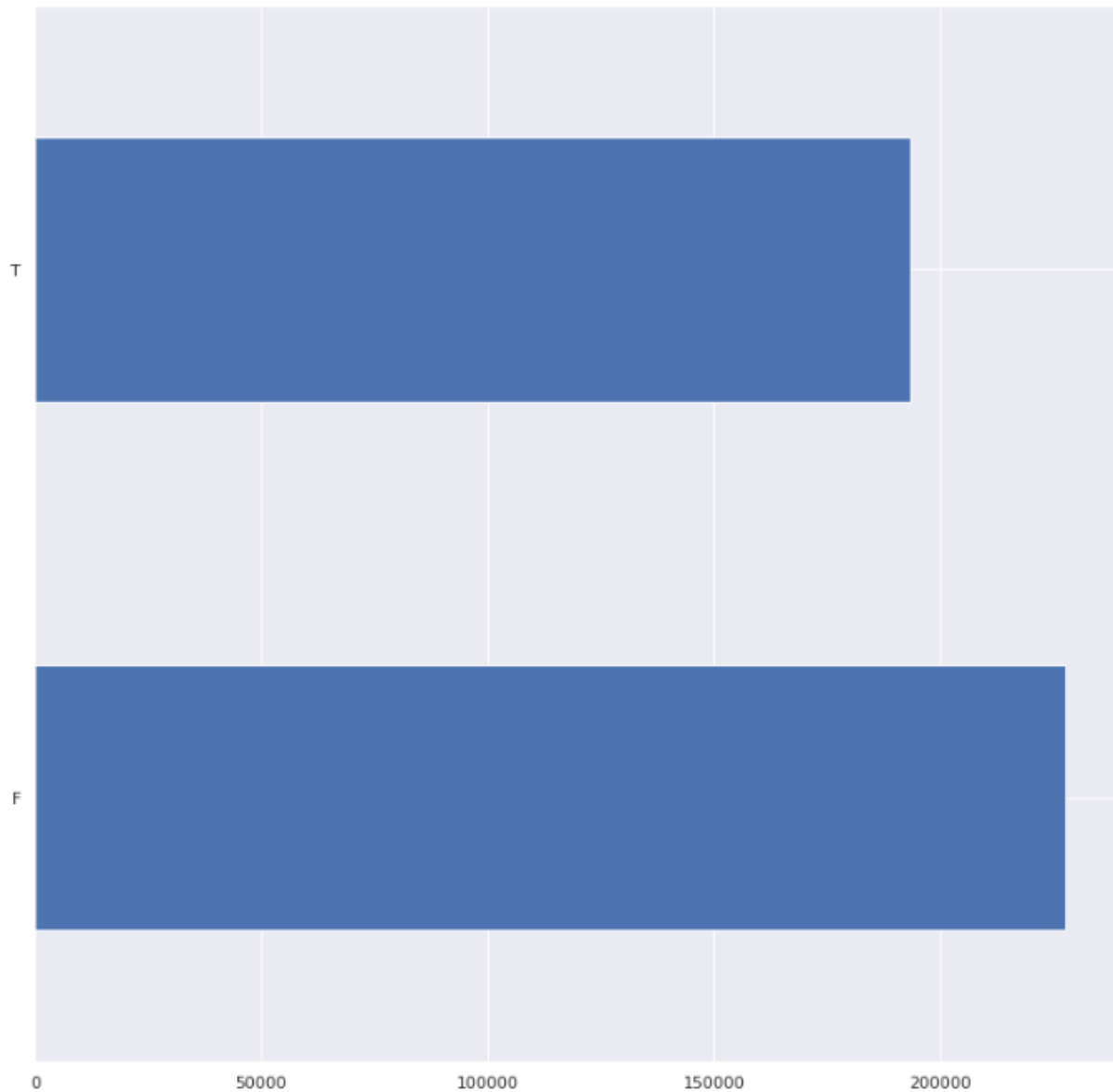
It has 94456 null values and 55 categories. Here we fill null values according to category distribution and we are combining 18 categories which do not have fraud transactions as goodmail. Now we have 48 different categories. to encode it we are using a one-hot encoder.





### 3.1.5 M6

It has 169360 null values and 2 categories. Here we fill null values according to categories distribution. and to encode it we are using label encoder.

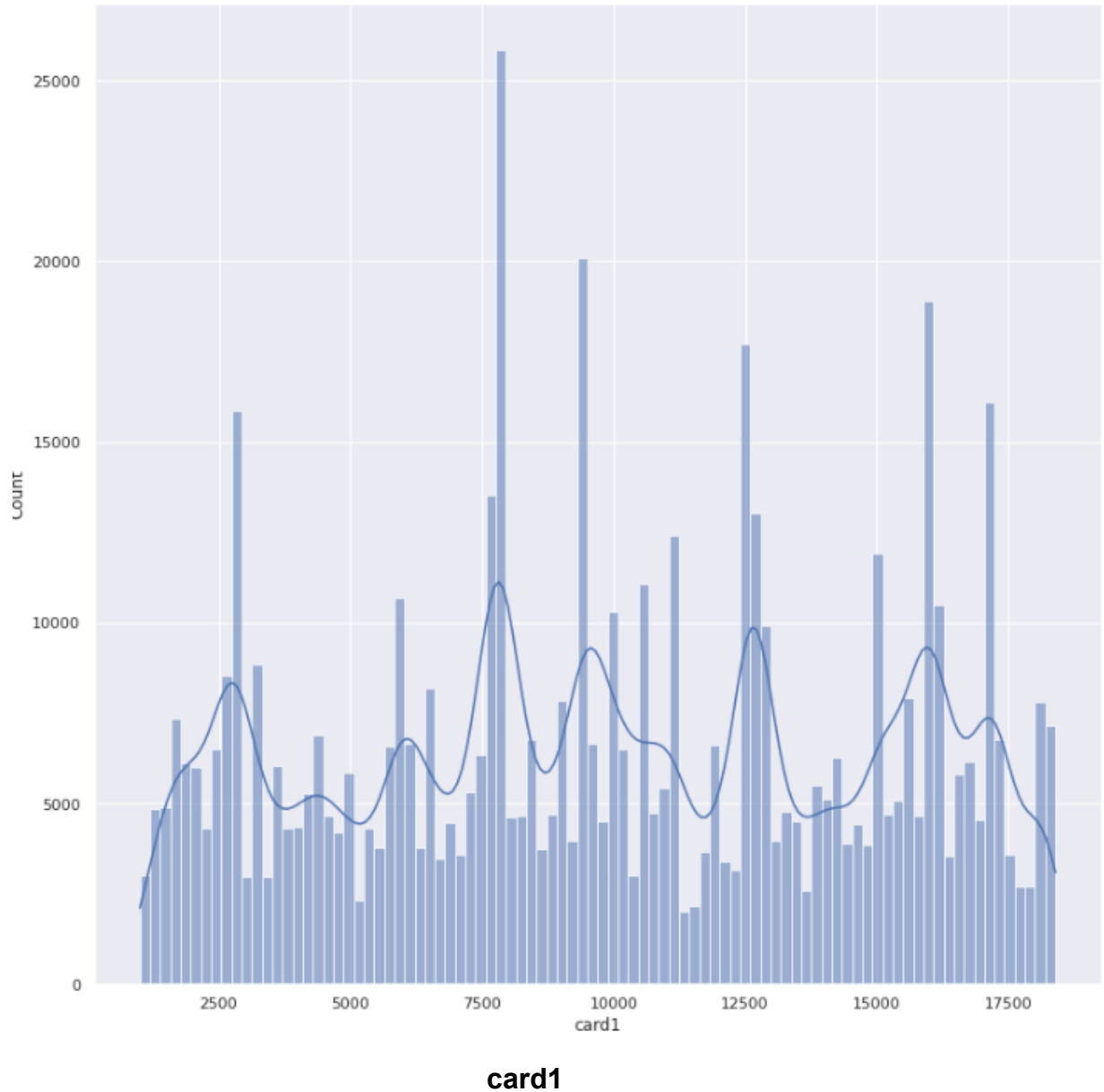


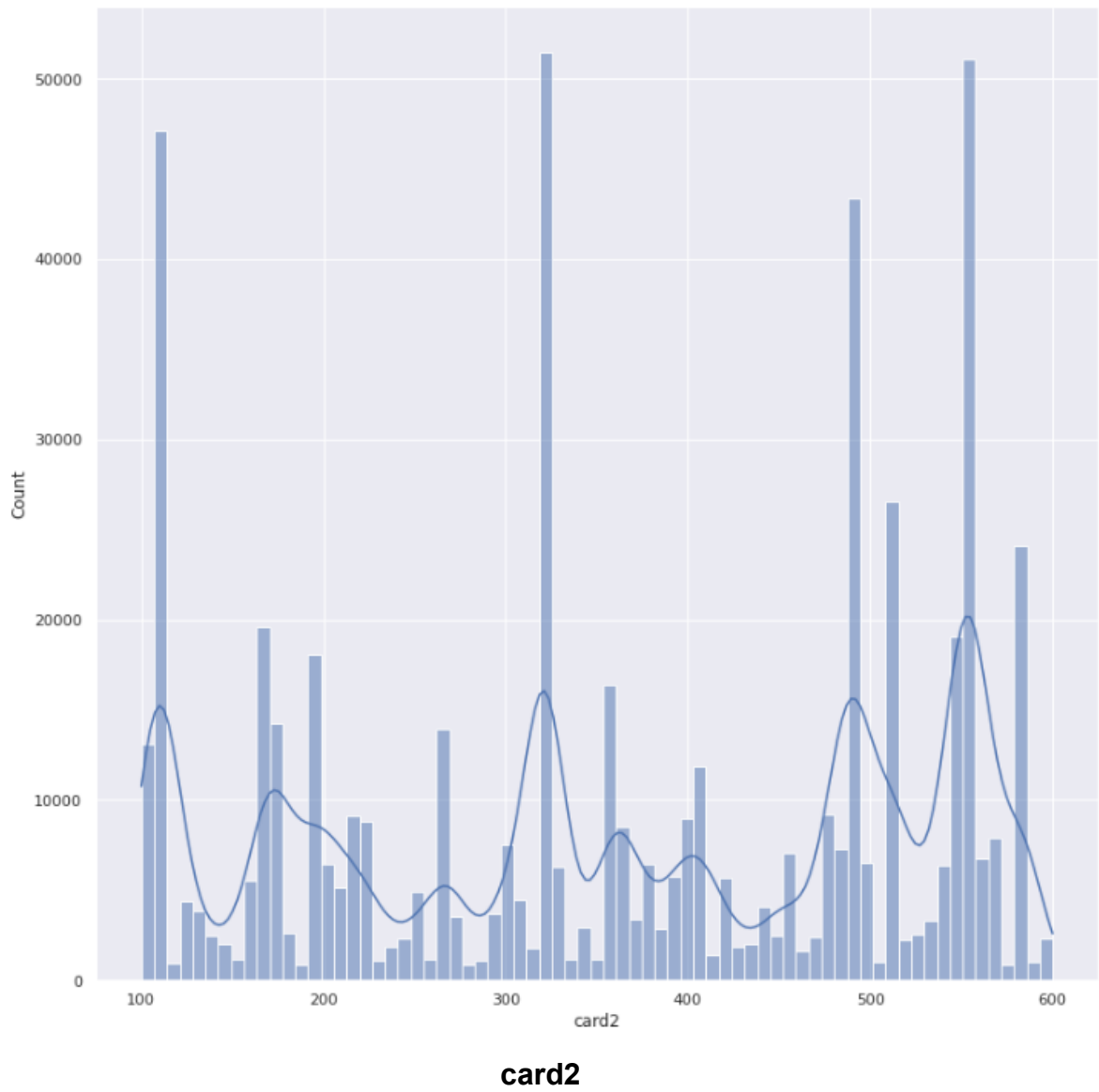
## 3.2 Numerical columns

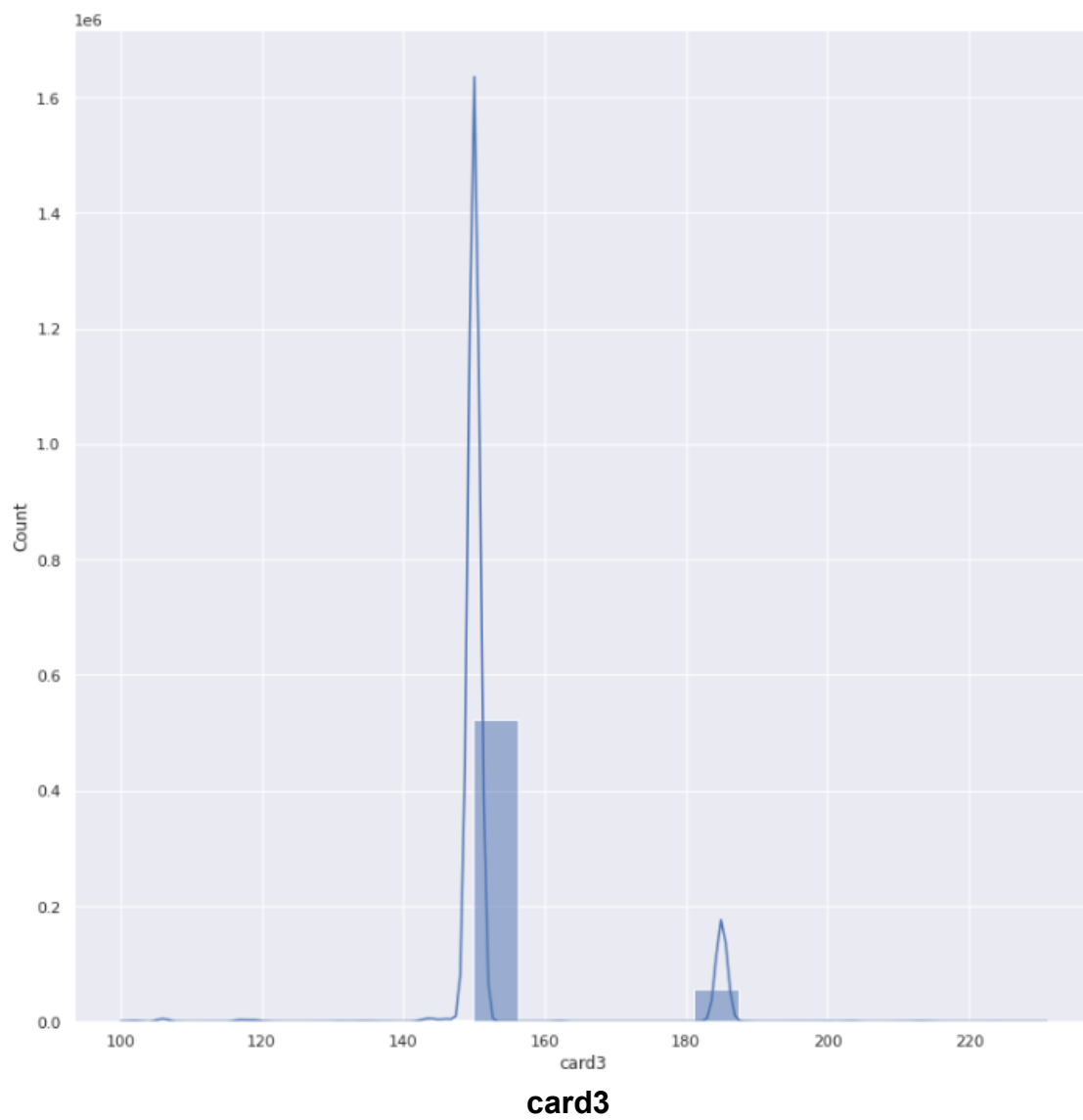
### 3.2.1 card1 ,card2 , card3 , card5

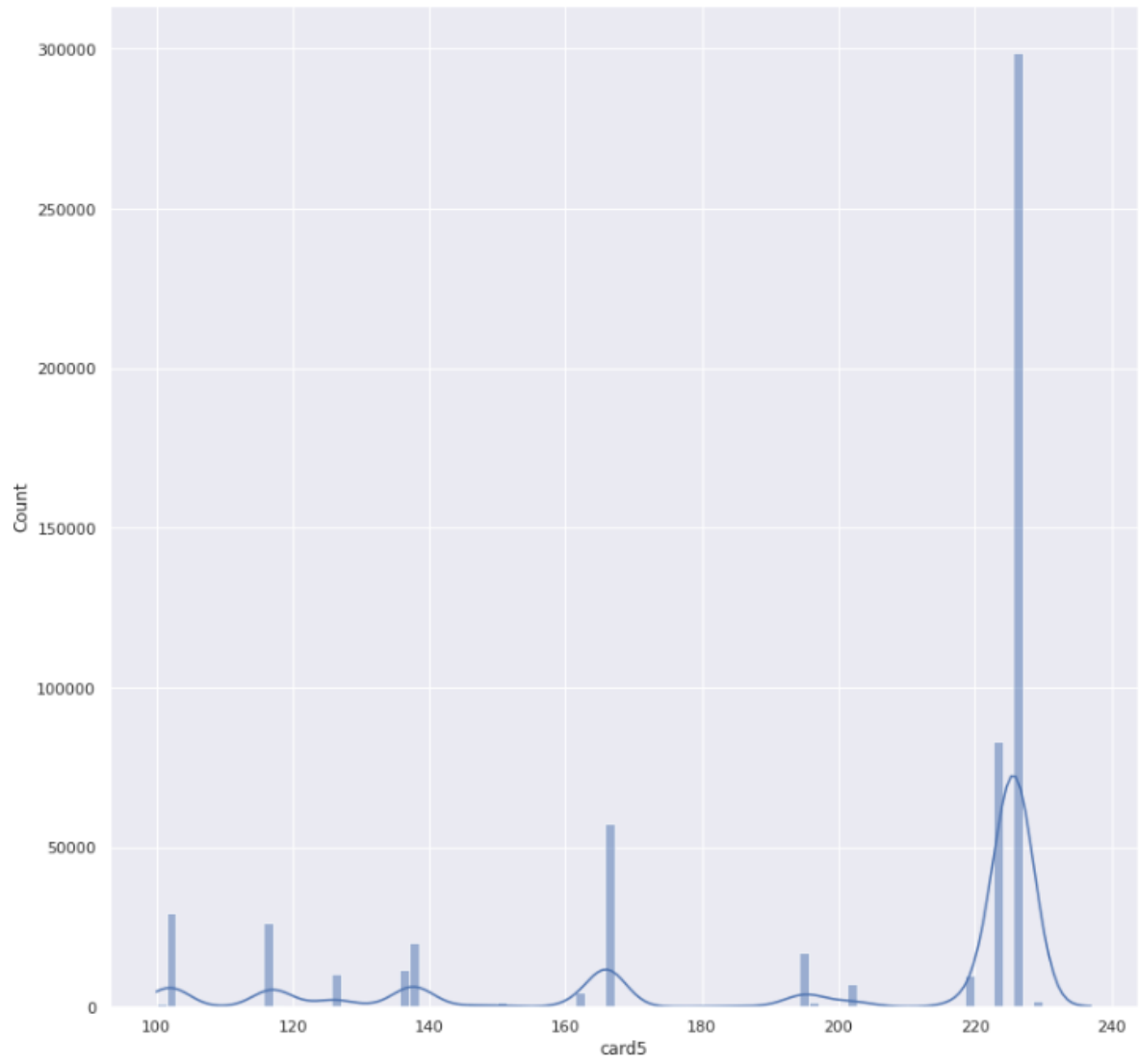
card1 has no null value.

card2 , card3 , card4 columns have 8933,1565,4259 null values respectively.we will fill the null value by the distribution of the numerical values as you can see in below graphs. Here numerical value represents the categories.





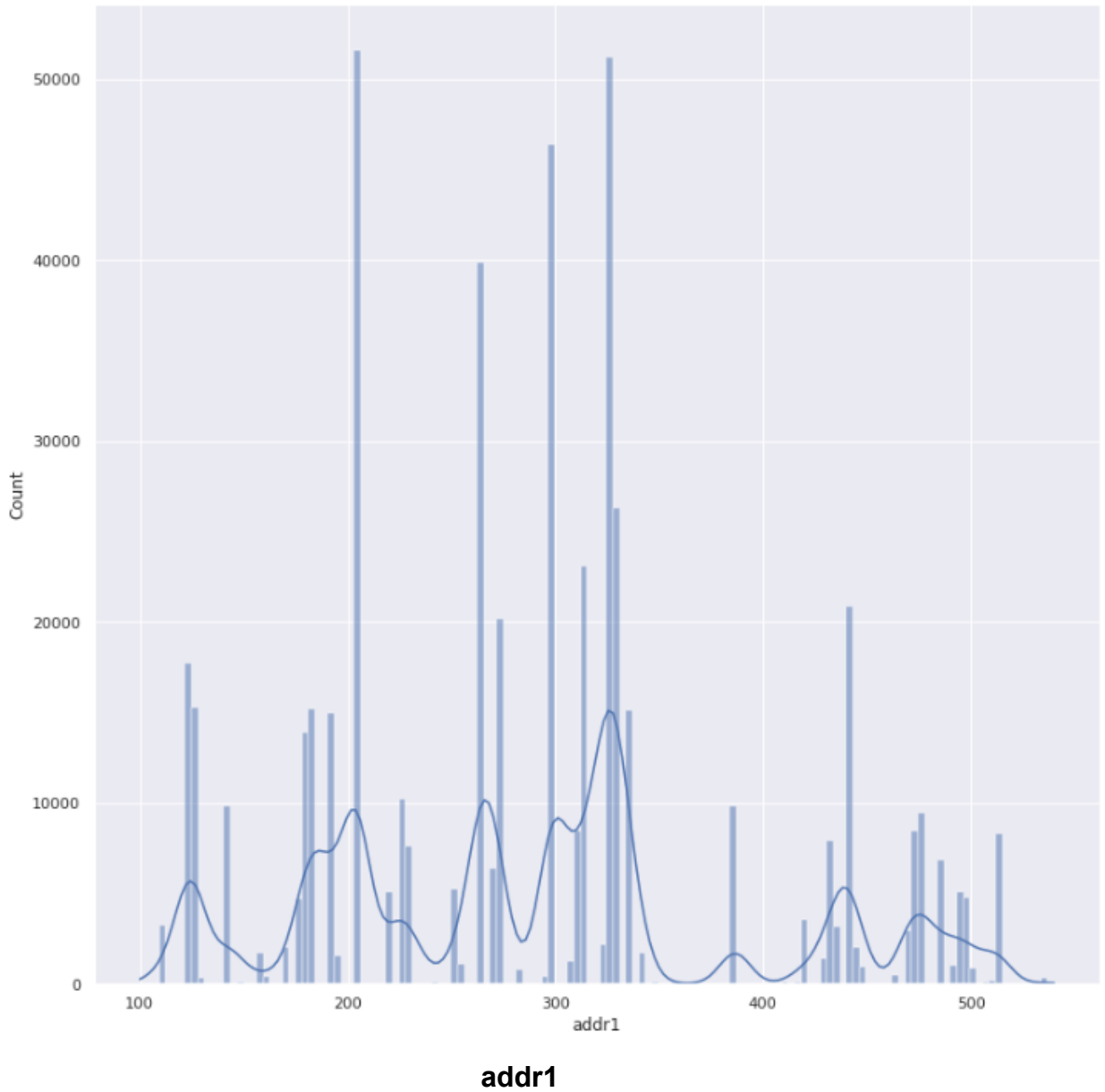


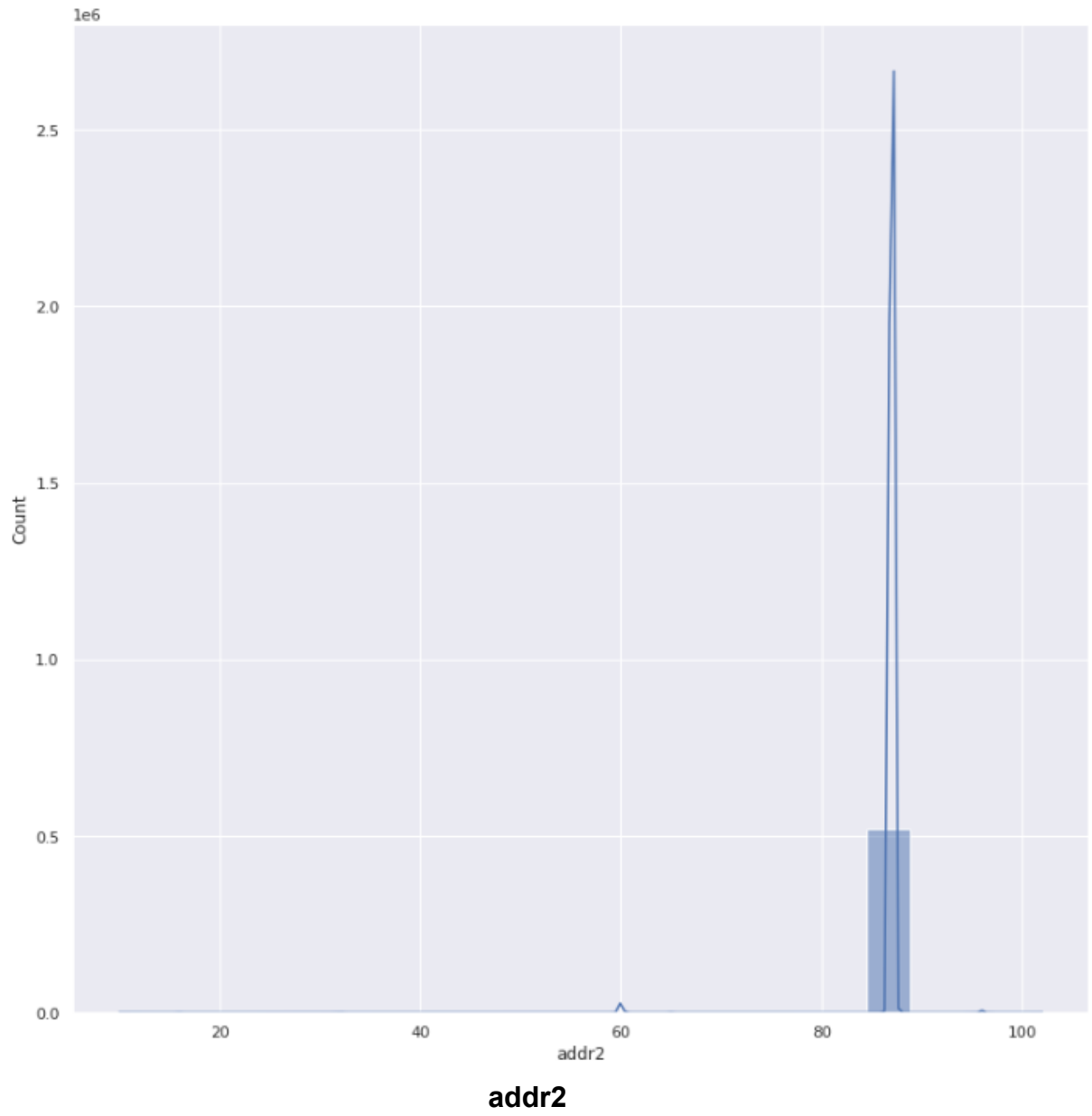


**card5**

### 3.2.2 addr1 and addr2

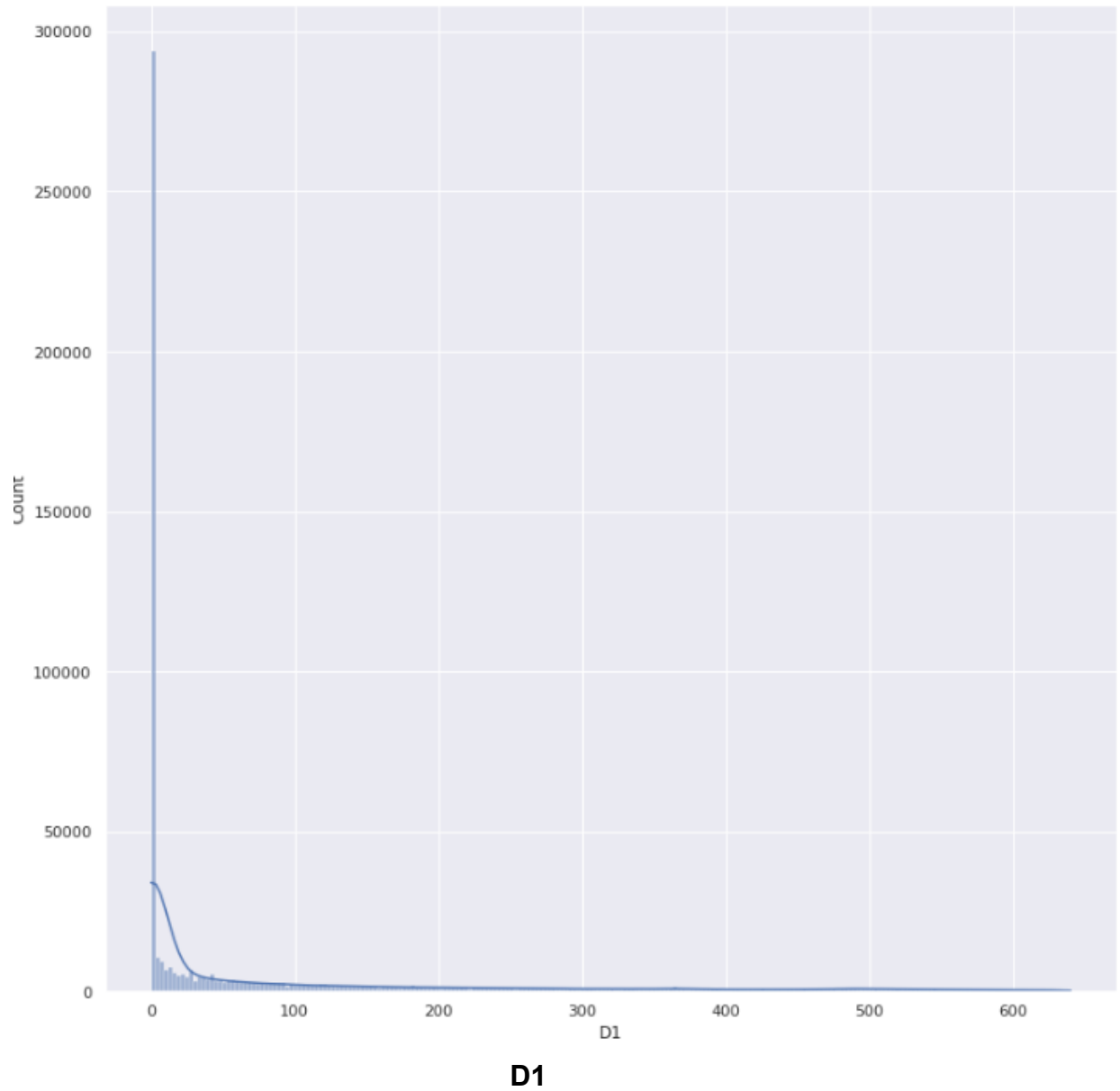
Both addr1 and addr2 columns have 65706 null values. we will fill the null value by the distribution of the numerical values as you can see in below graphs. Here numerical value represents the categories.



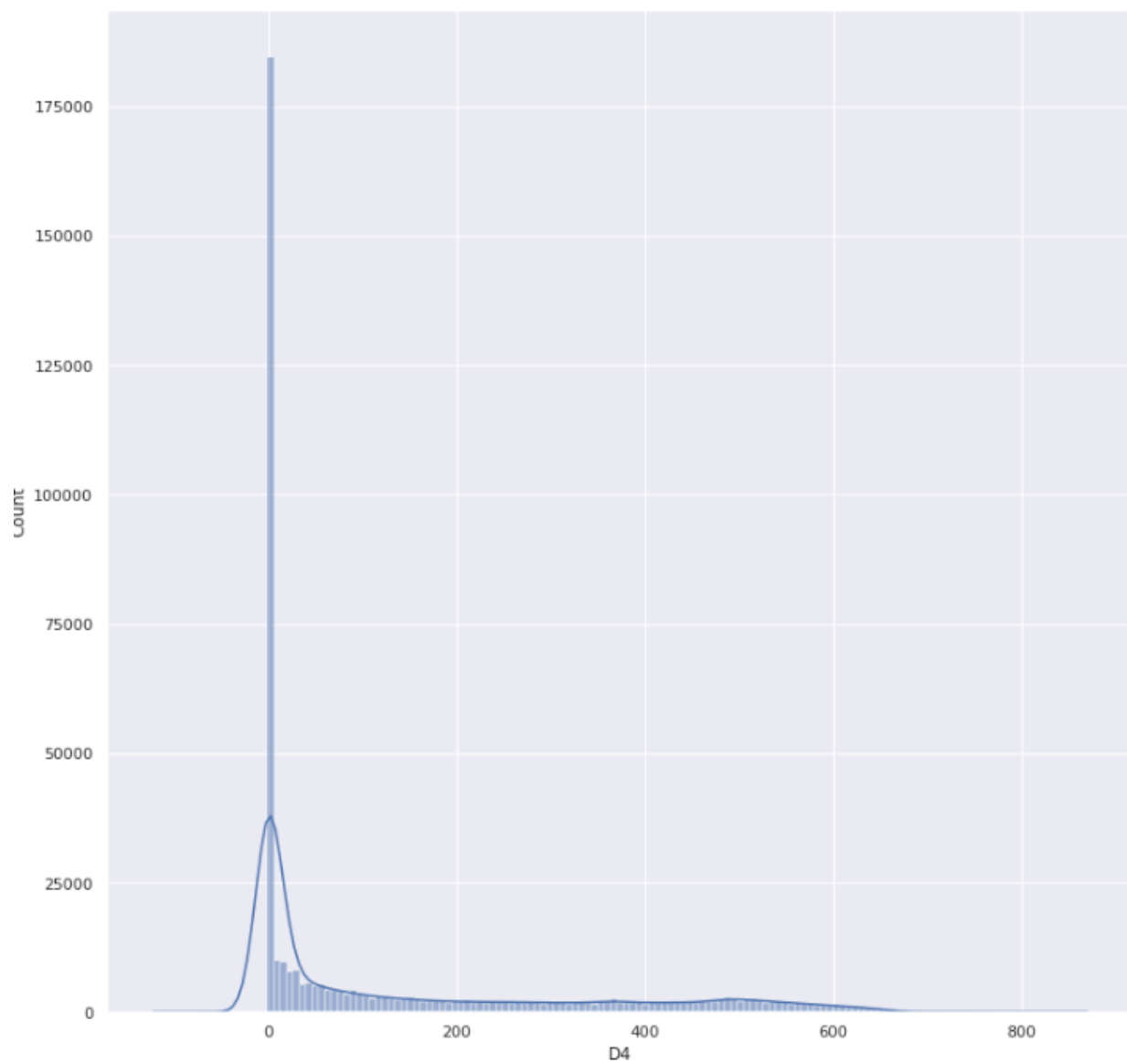


### 3.2.3 D1 , D4 , D10 , D15

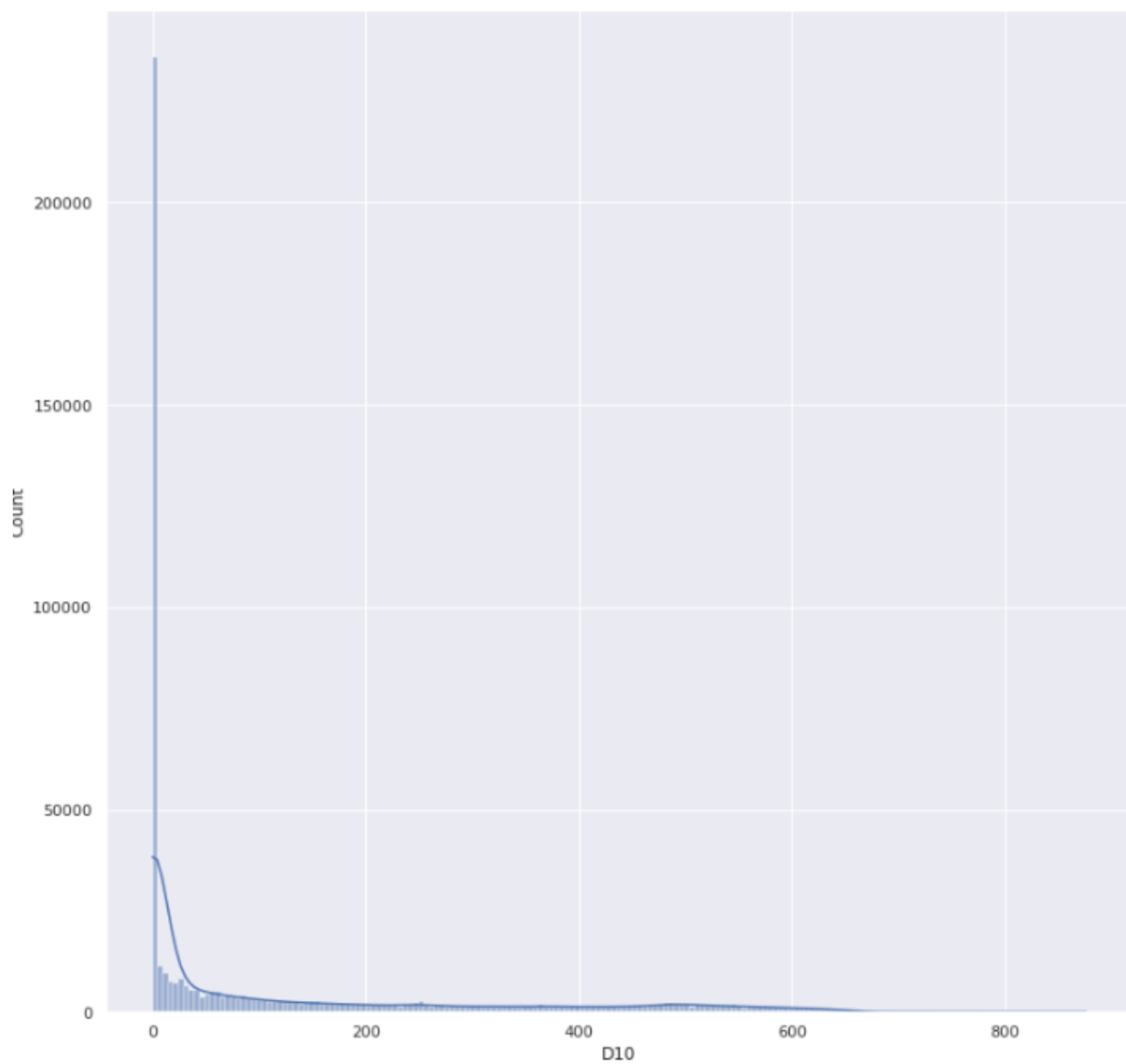
D1 , D4 , D10 and D15 columns have 1269,168922,76022,89133 null values respectively. we will fill the null value by mode of each column.



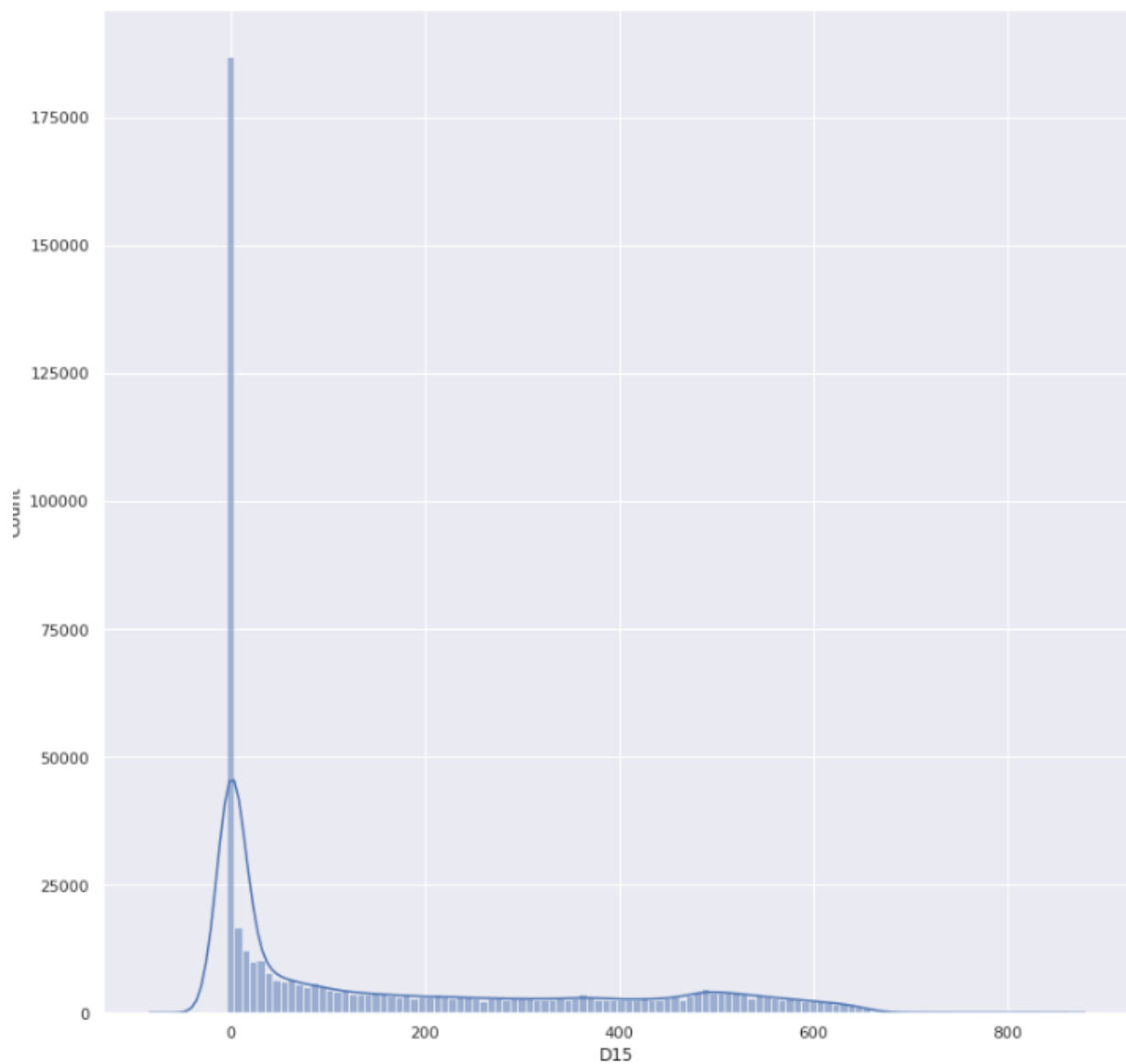




**D4**



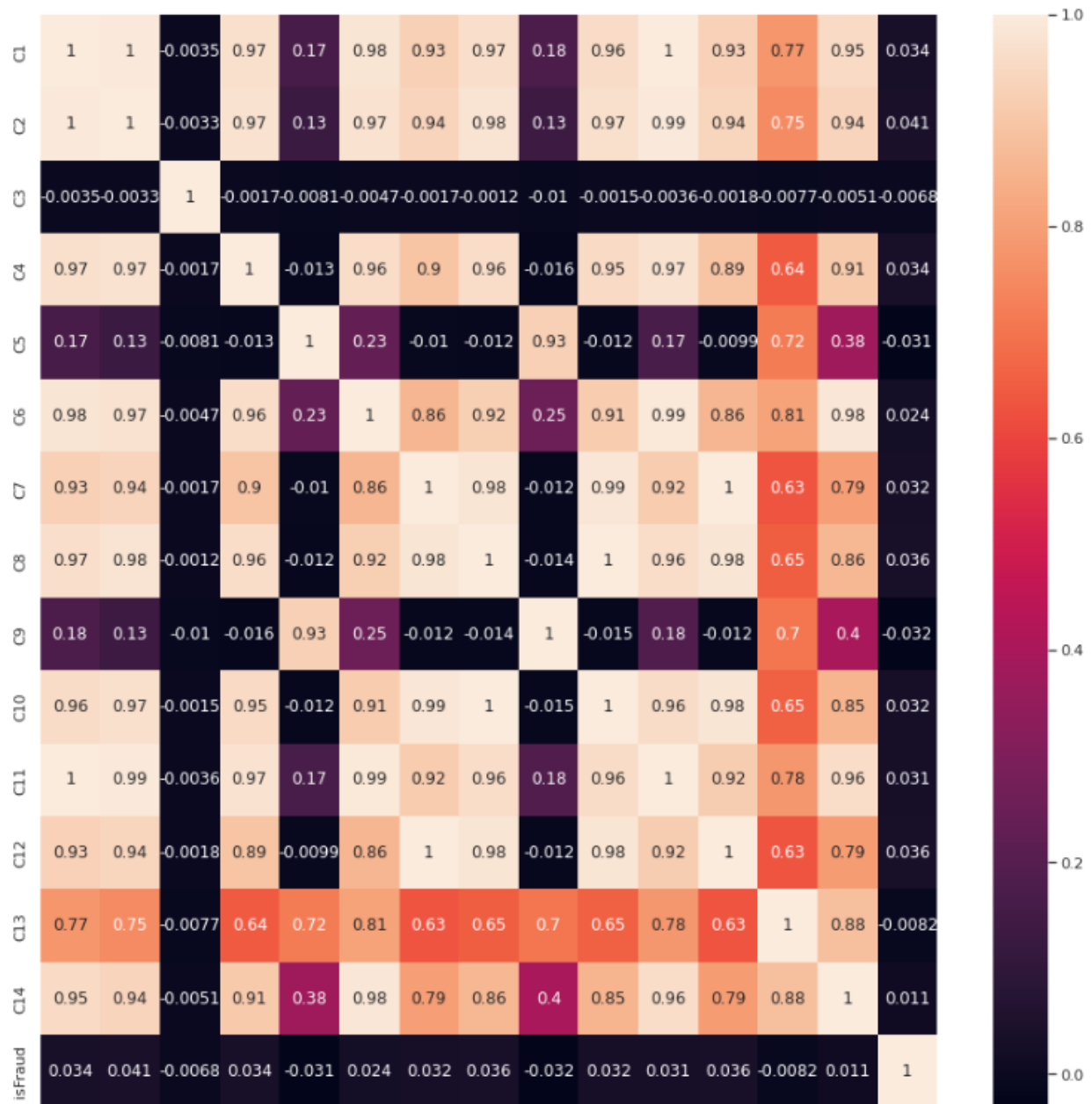
**D10**



**D15**

### 3.2.4 C1-C14 columns

c1-c14 has no null values. Some of the columns have high correlations as we can see in the graph below. we will drop the columns which have correlation greater than 90%. Ex. C2, C4, C6, C7, C8, C9, C10, C11, C12, C14



### 3.2.5 Vxxx: Vesta engineered rich features, including ranking, counting, and other entity relations.

V12-V34 has 76073 null values. V35-V52 has 168969 null values. V53-V74 has 77096 null value. V75-V94 has 89164 null value. V95-V137 has 314 null values. V279, V280, V284-V287, V290-V295, V297-V299, V302-V312, V316-V321 have 12 null values. V281-V283, V288-V289, V296, V300, V301, V313-V315 has 1269 null values.

We are filling null values using the mode of the columns. There are 86 columns which have correlation greater than 90. We will drop them.

Ex. 'V13', 'V16', 'V18', 'V20', 'V21', 'V22', 'V28', 'V30', 'V31', 'V32', 'V33', 'V34', 'V36', 'V40', 'V42', 'V43', 'V45', 'V48', 'V49', 'V50', 'V51', 'V52', 'V54', 'V57', 'V58', 'V59', 'V60', 'V63', 'V64', 'V68', 'V69', 'V70', 'V71', 'V72', 'V73', 'V74', 'V76', 'V79', 'V80', 'V81', 'V84', 'V85', 'V88', 'V89', 'V90', 'V91', 'V92', 'V93', 'V94', 'V96', 'V97', 'V101', 'V102', 'V103', 'V105', 'V106', 'V113', 'V126', 'V127', 'V128', 'V132', 'V133', 'V134', 'V137', 'V279', 'V280', 'V292', 'V293', 'V294', 'V295', 'V296', 'V297', 'V298', 'V299', 'V301', 'V302', 'V304', 'V306', 'V307', 'V308', 'V309', 'V315', 'V316', 'V317', 'V318', 'V321'

## 4. Training Models

Models are evaluated on ROC-AUC Score. For hyper parameter tuning we have used RandomSearchCV except XGboost. In XGboost we have used the hyperopt library.

| Model  | Hyper Parameters  | Value              | Model Score |
|--|-------------------|--------------------|-------------|
| <a href="#">Logistic regression classifier</a> | max_iter          | 1000               | 0.57204     |
|  | solver            | libliner           |             |
|  | penalty           | l2                 |             |
|  | c                 | 0.01               |             |
| <a href="#">Naive Bayes</a>                    | var_smoothing     | 0.0533669923120631 | 0.70151     |
| <a href="#">Ridge Classifier</a>               | alpha             | 1                  | 0.53441     |
| <a href="#">Perceptron</a>                     | penalty           | l1                 | 0.52756     |
|  | max_iter          | 1000               |             |
|  | eta0              | 0.0001             |             |
| <a href="#">Decision Tree</a>                  | min_sample_split  | 5                  | 0.62807     |
|  | min_samples_leaf  | 50                 |             |
|  | max_feature       | sqrt               |             |
|  | max_depth         | 20                 |             |
|  | criterion         | gini               |             |
| <a href="#">Random Forest</a>                  | n_estimators      | 1635               | 0.56897     |
|  | min_samples_split | 5                  |             |
|  | min_samples_leaf  | 100                |             |
|  | max_feature       | sqrt               |             |
|  | max_depth         | 10                 |             |
|  | criterion         | entropy            |             |
|  | bootstrap         | TRUE               |             |

### [Xgboost](#)

|                         |                        |                |
|-------------------------|------------------------|----------------|
| <b>alpha</b>            | <b>2</b>               | <b>0.87369</b> |
| <b>colsample_bytree</b> | <b>0.85</b>            |                |
| <b>eta</b>              | <b>0.125</b>           |                |
| <b>gamma</b>            | <b>0.75</b>            |                |
| <b>lambda</b>           | <b>1.8</b>             |                |
| <b>max_depth</b>        | <b>12</b>              |                |
| <b>n_estimators</b>     | <b>7</b>               |                |
| <b>min_child_weight</b> | <b>500</b>             |                |
| <b>scale_pos_weight</b> | <b>60</b>              |                |
| <b>subsample</b>        | <b>0.7</b>             |                |
| <b>objective</b>        | <b>binary:logistic</b> |                |
| <b>eval_metric</b>      | <b>auc</b>             |                |
| <b>tree_method</b>      | <b>hist</b>            |                |
| <b>bootstrap</b>        | <b>gbtree</b>          |                |

## Conclusion

In this Project we had a large dataset that contains columns with high null values.to handle null values we have dropped some columns ,for others we have filled with either mode or values according to distribution. To handle object type columns we have to use a one-hot encoder and label encoder. We have used many classification models with cross validation and parameter tuning using RandomSearch or hypertune.we got best result in XGboost model and parameter tuning using hypetopt.we got approximately 87% ROC-AUC Score.