



MICRO CREDIT PREDICTION

Submitted by:
KONATALA MOHIT

ACKNOWLEDGMENT

I am grateful to the writers, developers and authority of the website pages machinelearningmastry, geekforgeeks, wikipedia and stackoverflow which have helped me to refer for whenever there was need for guidance. Articles on Indonesian Telecom sector, Bing networks and Economics in finance helped me to get a better understanding of the data to analyse and explore data. I am mostly grateful to DataTrained who have guided me in learning and enhancing my data science skills required to analyse and solve this project.

INTRODUCTION:

In this report I will be discussing about micro credit prediction using few of the machine learning models via python and its libraries. In simple terms micro credit means a loan of small amount. The primary objective of micro credit is to help raise the income and standard of living of the below poverty line category families. Micro credit is only applicable to lower income families and it is provided by micro credit institutes of the country. Communication devices have become primary digital gadget to pursue our normal life in today's life. On an average a person uses his phone for 4 hours of the day i.e., texting, internet surfing, streaming online or calling etc. Telecom industries have also initiated micro credit service to help provide services to lower income class which would help sustain their lives.

The drawback of micro credit services is the uncertainty of the subscriber to payback the amount allotted to them within 5 days of time. To get the better idea of this we will look at the various factors of data provided by one of the clients to showcase if the telecom subscriber can return borrowed loan amount within the given time period.

Review of Literature:

For the better understanding of the data referred to telecom sector of Indonesia as the data is of that particular country and micro credit in finance. Micro credit was originated from Grameen Bank that was developed by economist Muhammad Yunus i.e., established in 1976 in Bangladesh. Telecom micro credit is a value-added service particularly designed to meet consumer needs. It is a service which the telecom can offer to subscribers and resellers both. Through this service, telco can provide ease of access of airtime stock credit to their resellers when they run out of stock balance, without disrupting the normal flow of business. On the subscriber side, they can remain connected without having to recharge through a voucher or digital topup. This loan is disbursed in real-time and is calculated based on intelligent reseller/subscriber profiling.

Undertaken Problem:

Objective of the project is to predict if the loan borrowed by the customer can be cleared by that person with a span of 5 days or not i.e., by analysing various factors that affect the label outcome. Micro Credit prediction is Classification type prediction as the label is either yes or no.

Mathematical/ Analytical Modelling of the Problem:

Label is the target data in the dataset and the variables present the target data are in binary form, therefore Classification Machine Learning Models are used in this project. Classification is for estimating categorical form of data by using established relation

between feature and target variables. Descriptive Analysis is used to study and observe the data.

Data Sources and their formats:

The source of the data is provided by the client. The data is provided in the csv format. Data contains 209593 entries having 37 variables in the dataset.

Data Pre-processing:

Data pre-processing has two main steps i.e., Data Cleaning and Data Transforming.

Data Cleaning:

Data Cleaning is one of the most important steps creating a machine model. If an uncleaned data is fed to a machine learning model, then the model will perform very poorly.

The Dataset does not contain any null values in it. Unnamed: 0 (serial number), msisdn (mobile number) these two columns have unique element for every individual row therefore dropping them. Pcircle column has only one unique element and it is repeated in every row so it won't help us the outcome prediction therefore dropping this variable. Pdate columns can be simplified by splitting it into day, month and year columns by using split method and assigning values to new columns. Removing the pdate columns as it has been divided into 3 new columns and also removing year column as the data is of only of year 2016. Day and month are in object format converting them into numerical format using pd.numeric.

Few of the columns contains negative values which is can't possible so locating those negative values and multiplying them with -1 to make them as positive values. The dataset contains outliers. Age of the cellular maximum value is 999860 days which is inappropriate in general sense. So, limiting the data to 1465 days as considering the client is a new telecom industry launched 2012. In Indonesia the average monthly expenditure per capita in mobile bill in 2016 is 22182 Indonesian Rupiah. According to that per day daily money spent can be around 740 (maximum). Therefore, limiting the maximum value of daily value to 740 for both 30 and 90 days. We know that 22182 is the maximum value and loan is used by the lower income families so the upper limit plan of lower income used is around 7500 so any value around that is rounded to 7500 as monthly recharge and for 90 days it is 10000 Indonesian Rupiah. Last recharge date has been limited to 30 days as most of the package used by lower income class is monthly plan. Recharge frequency has been limited to 1 as most of the people try to recharge only once a month unless it is for additional data. The max loan amount values have only 3 variety 0, 5 and 10 i.e., values been 1-4 as marked as 5 and any other value above is marked as 10. For remaining columns outliers can be handled by using IQR method.

Note: Label column is imbalanced.

Data Transforming:

After Data Cleaning the data must be transformed into numerical form as one can't feed ordinal data to machine learning model and also the data must be normalized as normalizing the data the machine learning model gives equal importance to all the data.

In this project I am using label encoder to transform the ordinal data present in the data set into numerical form. Using Standard Scalar function to normalize the data.

Note: Before normalize the data separated the target variable from feature variables.

Using VIF to check that value is under 5 as having higher number indicates that the dataset the multi collinearity between the independent variables which must be arrested else we cannot achieve optimal machine learning model.

After standardizing the data, we must split the data into train and test sets. Train Test Split method is to split the data into train and test data

Data Inputs- Logic- Output Relationships:

After Separating the Input (feature) and output (target) data we split them into train and test division one part of the data is used to train the ML model and other part of the data to predict the output. After splitting the data, we must balance the training data in order to minimize the ML being biased the majority result. Using SMOTE to balance the training data. When the train data is fed to machine learning model it generates an algorithm or simply put an equation that is applicable to all the data and when test data is fed to it implements the trained data equation to the current input values to predict the outcome. If the input has no outliers and is clean that there is no over or underfitting in the outcome.

Hardware and Software Requirements and Tools Used:

Hardware Required for Jupyter Notebook Software is as follows:

Memory and disk space required per user: 1GB RAM + 1GB of disk + .5 CPU core.

Server overhead: 2-4GB or 10% system overhead (whatever is larger), .5 CPU cores, 10GB disk space.

Port requirements: Port 8000 plus 5 unique, random ports per notebook

Libraries Used:

Pandas library: To frame raw data, visualize and perform task on it via other libraries.

Numpy library: To perform mathematical functions on the framed data numpy is used. In this project used it to location nan values and replace them with desired value also to find mean and standard values.

Matplotlib library: This library is used to visualize the data. Used to visualize univariate and bivariate analysis (pie plot, count plot and scatter plot) also to visualize outliers via box plot.

Seaborn library: heatmap to see co-relation between feature variable to arrest high collinearity.

Sklearn library: Imported this library to normalize the data, split the data into train and test data, various machine learning model and cross validation techniques.

Warnings library: To ignore filter warning shown while compiling block of codes

Pickle: To save the trained machine learning model.

Algorithm Used:

1. XGB Classifier
2. Random Forest Classifier
3. K-Neighbors Classifier
4. Logistic Regression

Model Training and Selection:

XGB Classifier:

XGBoost stands for "Extreme Gradient Boosting" and it is an implementation of gradient boosting trees algorithm. The XGBoost is a popular supervised machine learning model with characteristics like computation speed, parallelization, and performance

```

model3 =XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                      colsample_bynode=1, colsample_bytree=0.4,
                      enable_categorical=False, gamma=0, gpu_id=-1,
                      importance_type=None, interaction_constraints='',
                      learning_rate=0.2, max_delta_step=0, max_depth=16,
                      min_child_weight=4, monotone_constraints='()',
                      n_estimators=350, n_jobs=8, num_parallel_tree=1, predictor='auto',
                      random_state=0, reg_alpha=0, reg_lambda=1, scale_pos_weight=1,
                      subsample=0.6, tree_method='exact', validate_parameters=1,
                      verbosity=None)
model3.fit(X_over,Y_over)
p3=model3.predict(X_test)
print(classification_report(p3, y_test))

```

[12:20:09] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1
 0, the default evaluation metric used with the objective 'binary:logistic' was cl
 t eval_metric if you'd like to restore the old behavior.

	precision	recall	f1-score	support
0	0.51	0.56	0.53	4833
1	0.94	0.93	0.94	37086
accuracy			0.89	41919
macro avg	0.73	0.74	0.73	41919
weighted avg	0.89	0.89	0.89	41919

Random Forest Classifier:

A random forest is a meta estimator that fits a number of classifying decision trees on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting

```

model4 =RandomForestClassifier(max_depth=16, max_features=10, max_leaf_nodes=50,
                               max_samples=0.3, min_samples_leaf=300,
                               min_samples_split=4, n_estimators=300)
model4.fit(X_over,Y_over)
p4=model4.predict(X_test)
print(classification_report(p4, y_test))

```

	precision	recall	f1-score	support
0	0.71	0.39	0.51	9523
1	0.84	0.95	0.89	32396
accuracy			0.83	41919
macro avg	0.78	0.67	0.70	41919
weighted avg	0.81	0.83	0.81	41919

K-Neighbors Classifier:

KNN is a non-parametric method that, in an intuitive manner, approximates the association between independent variables and the continuous outcome by averaging the observations in the same neighbourhood

```

model2 = KNeighborsClassifier(leaf_size=40, n_neighbors=1, weights='distance')
model2.fit(X_over,Y_over)
p2=model2.predict(X_test)
print(classification_report(p2, y_test))

```

	precision	recall	f1-score	support
0	0.50	0.35	0.41	7505
1	0.87	0.92	0.89	34414
accuracy			0.82	41919
macro avg	0.69	0.64	0.65	41919
weighted avg	0.80	0.82	0.81	41919

Logistic Regression:

Logistic regression is a process of modelling the probability of a discrete outcome given an input variable. The most common logistic regression models a binary outcome, something that can take two values such as true/false, yes/no, and so on.

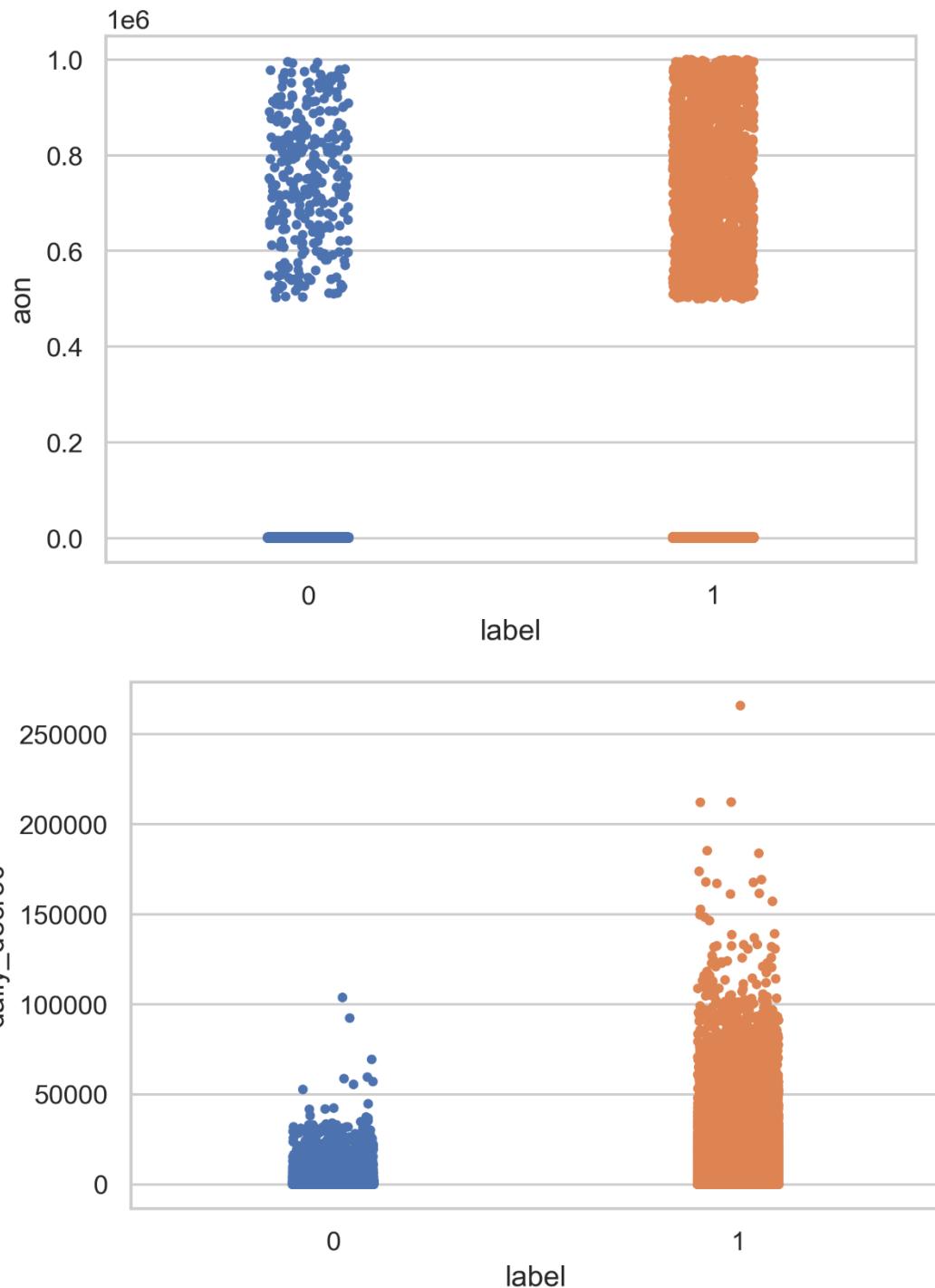
```

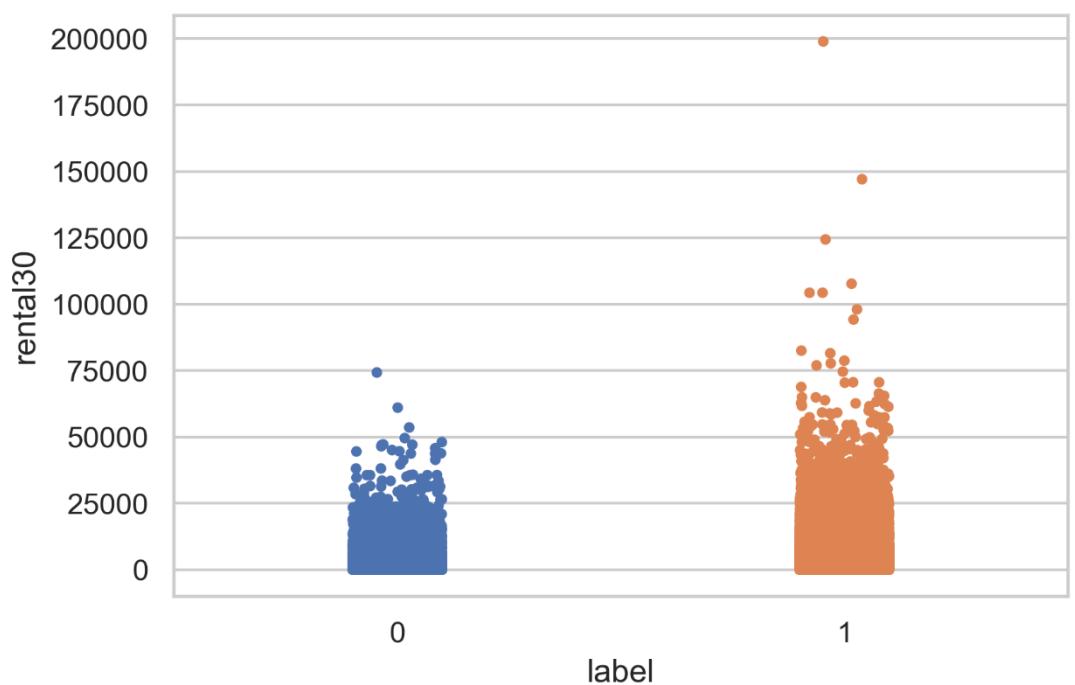
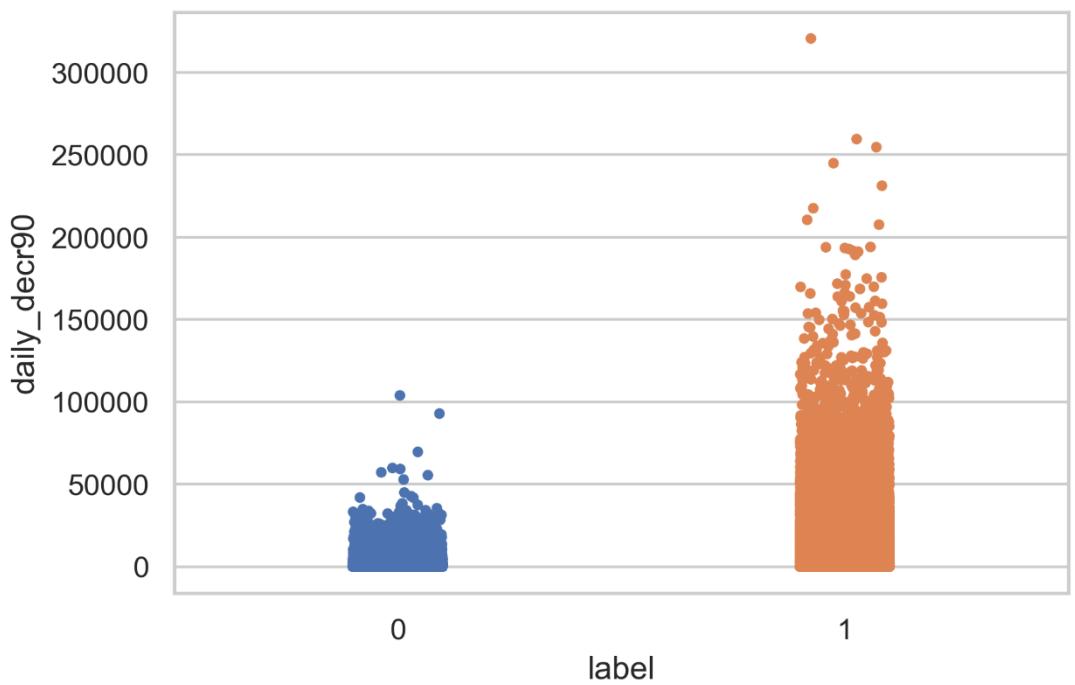
model1 = LogisticRegression(solver='newton-cg')
model1.fit(X_over,Y_over)
p1=model1.predict(X_test)
print(classification_report(p1, y_test))

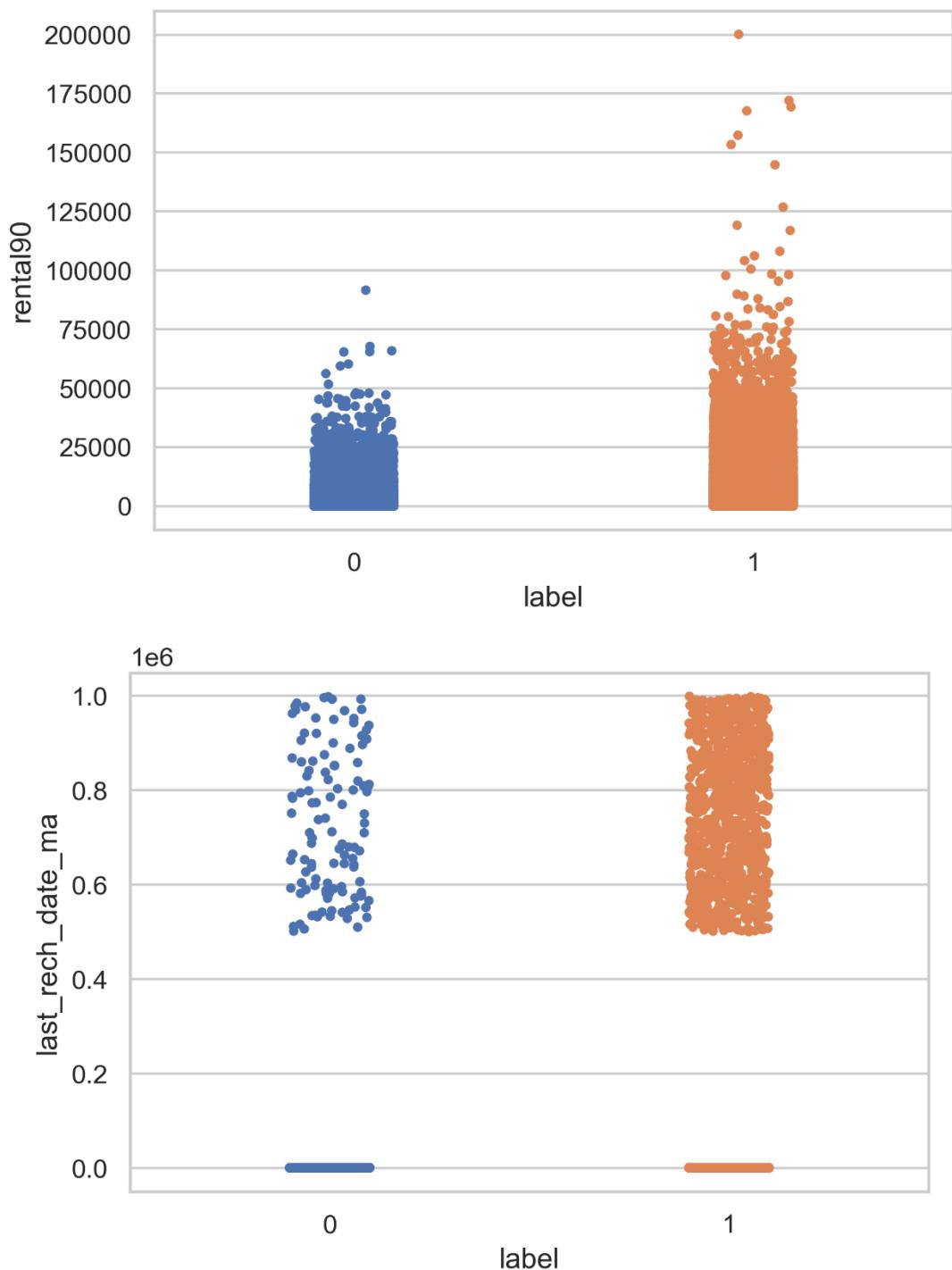
```

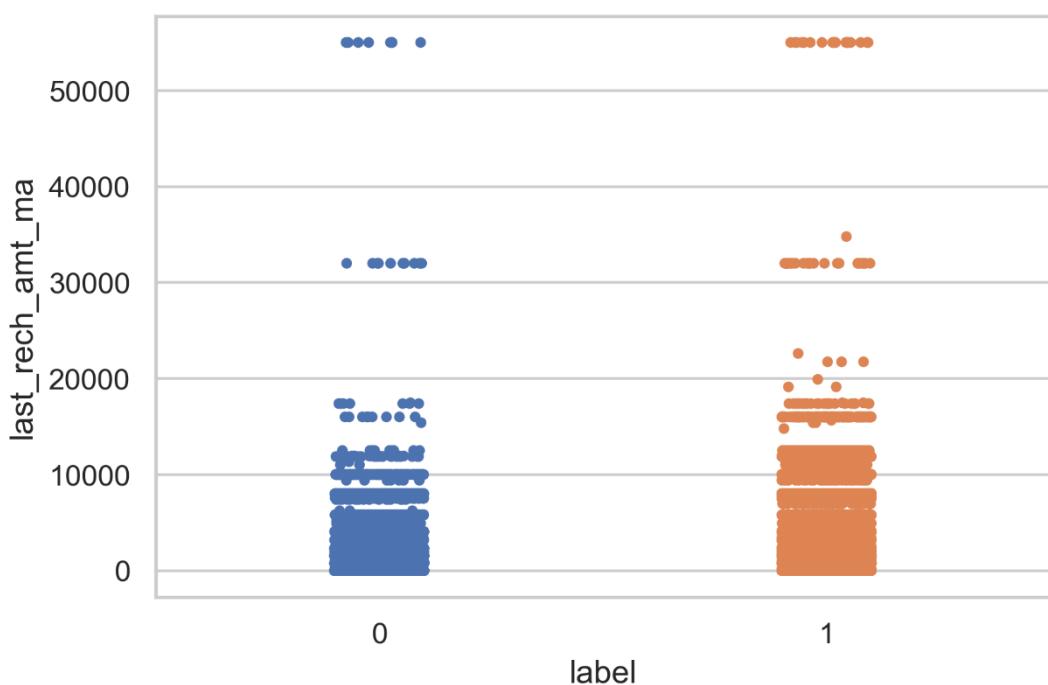
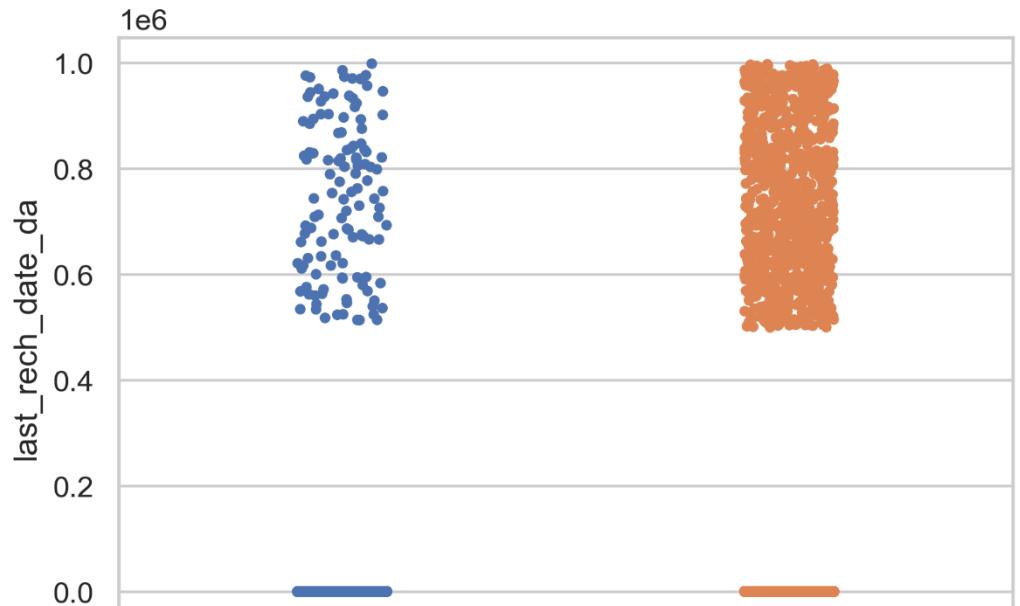
	precision	recall	f1-score	support
0	0.77	0.28	0.42	14191
1	0.72	0.96	0.82	27728
accuracy			0.73	41919
macro avg	0.75	0.62	0.62	41919
weighted avg	0.74	0.73	0.69	41919

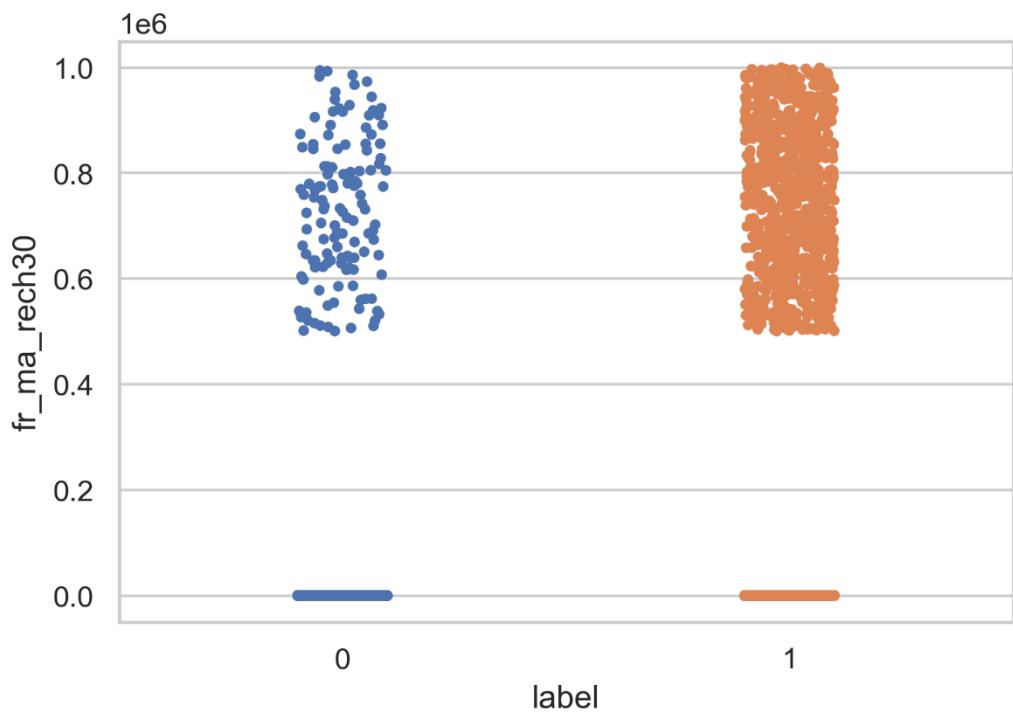
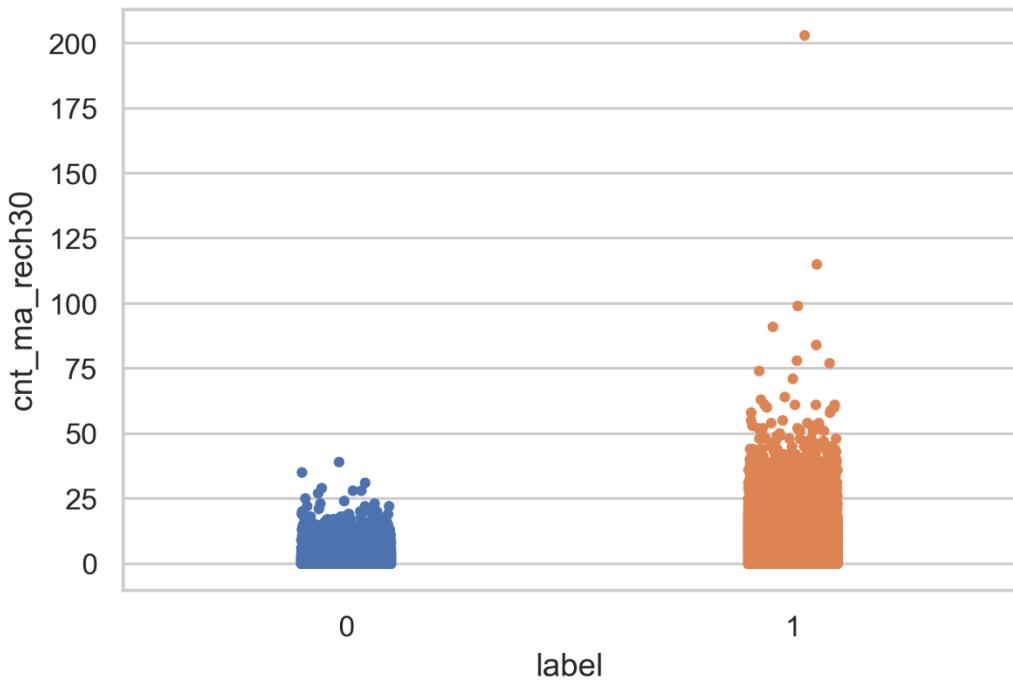
Visualization:

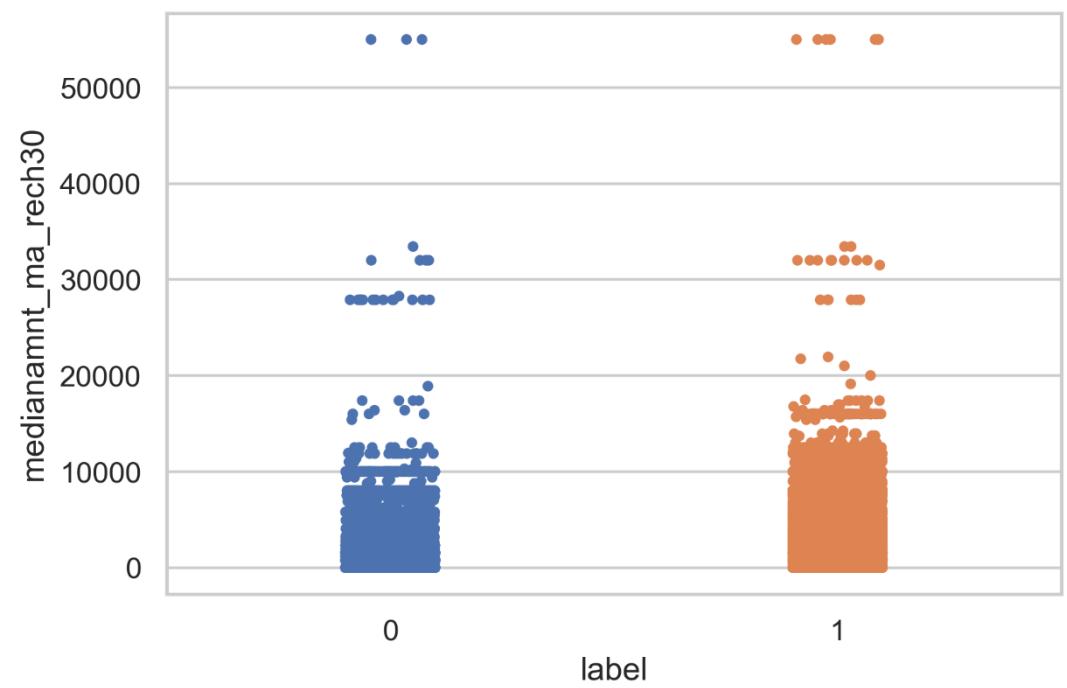
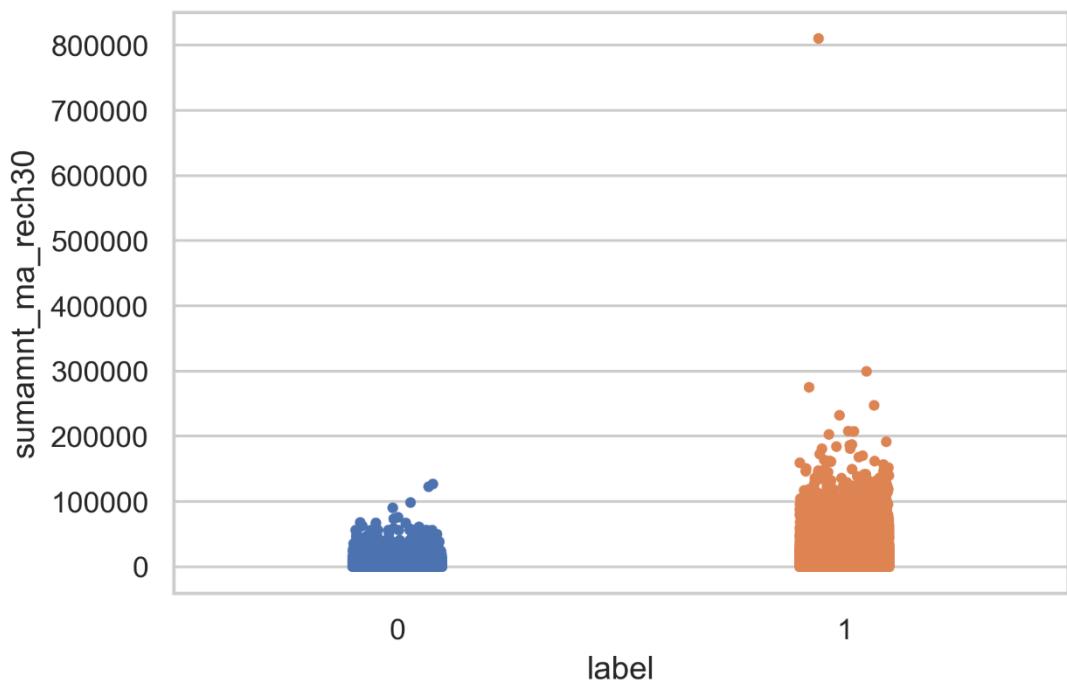


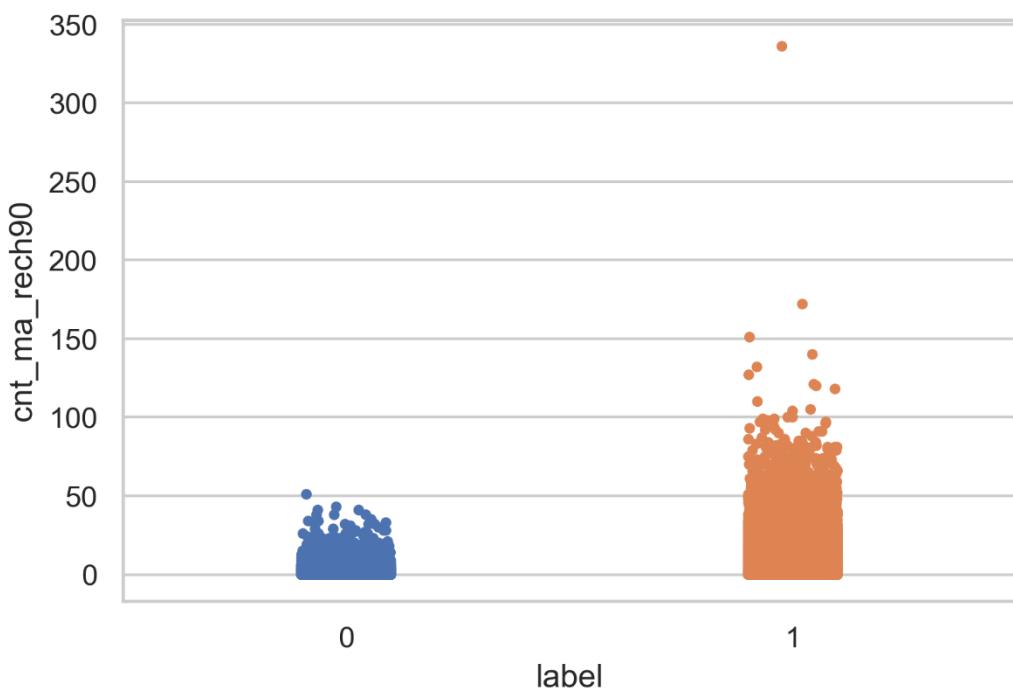
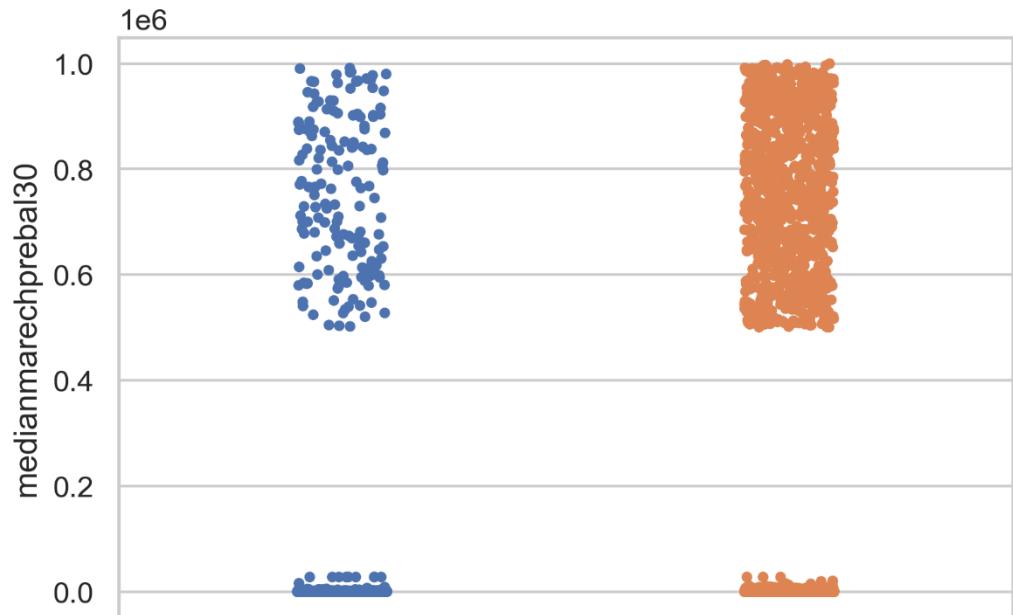


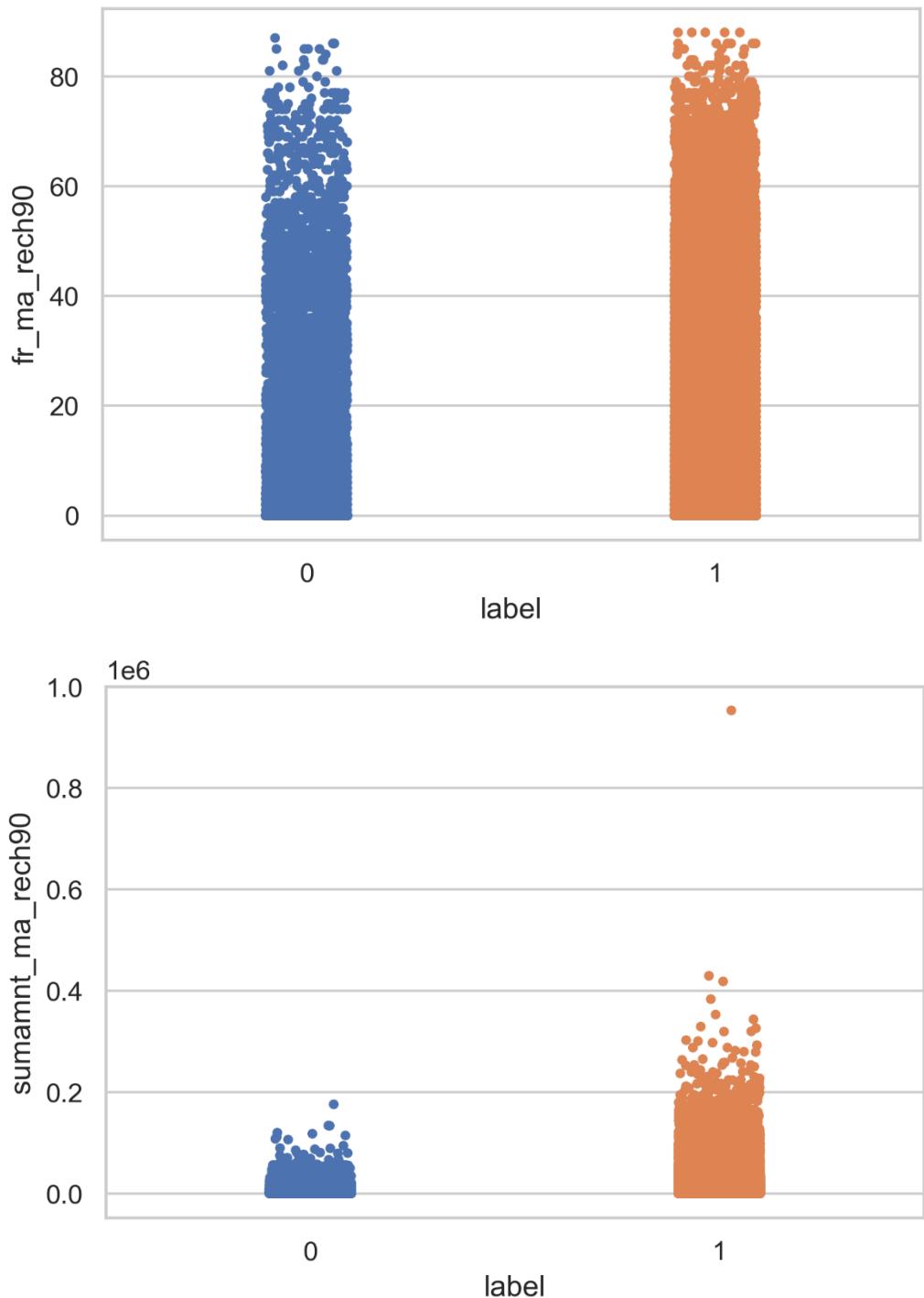


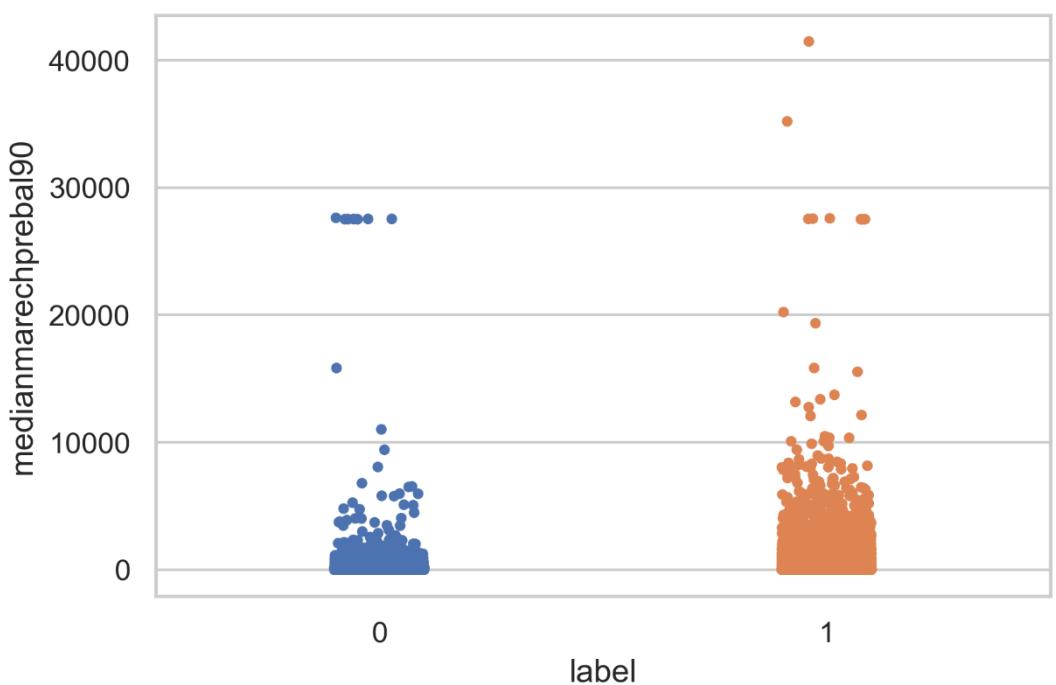
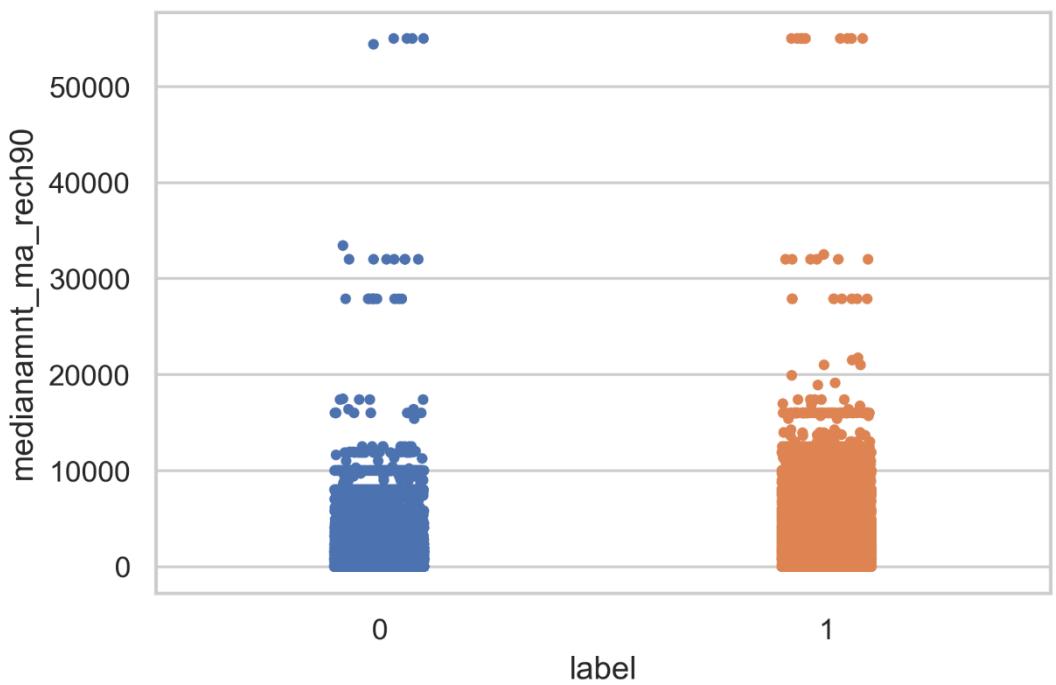


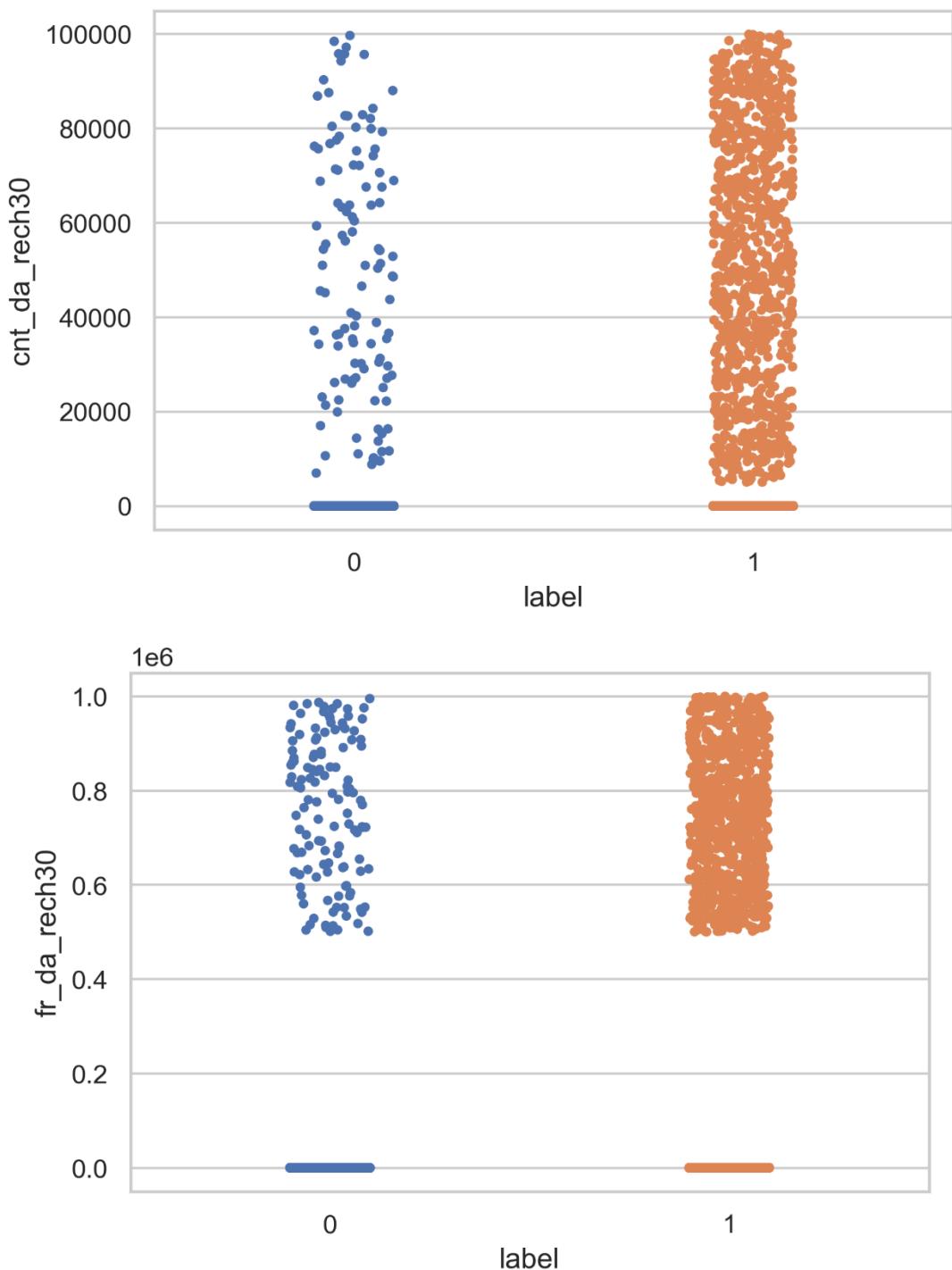


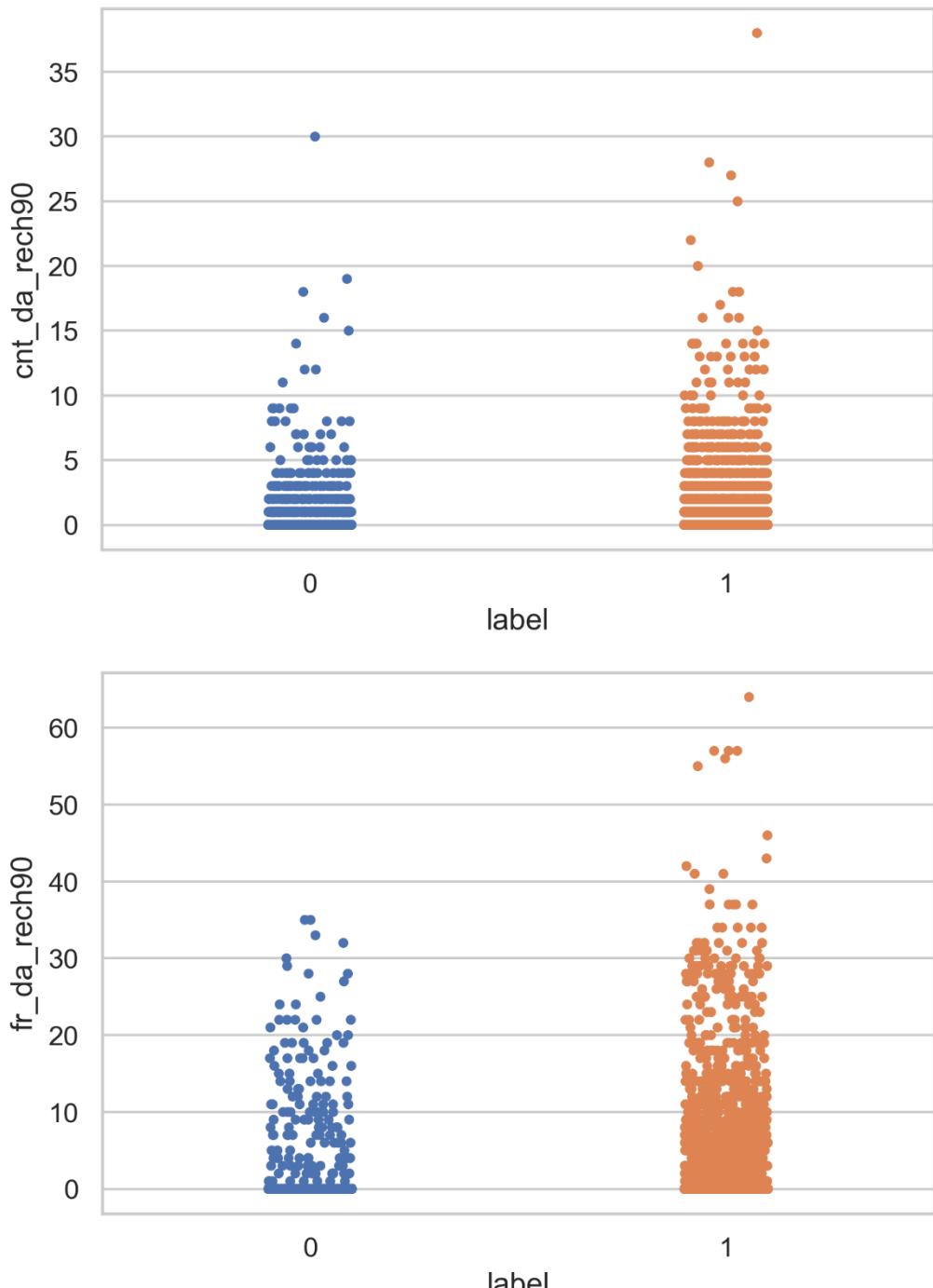


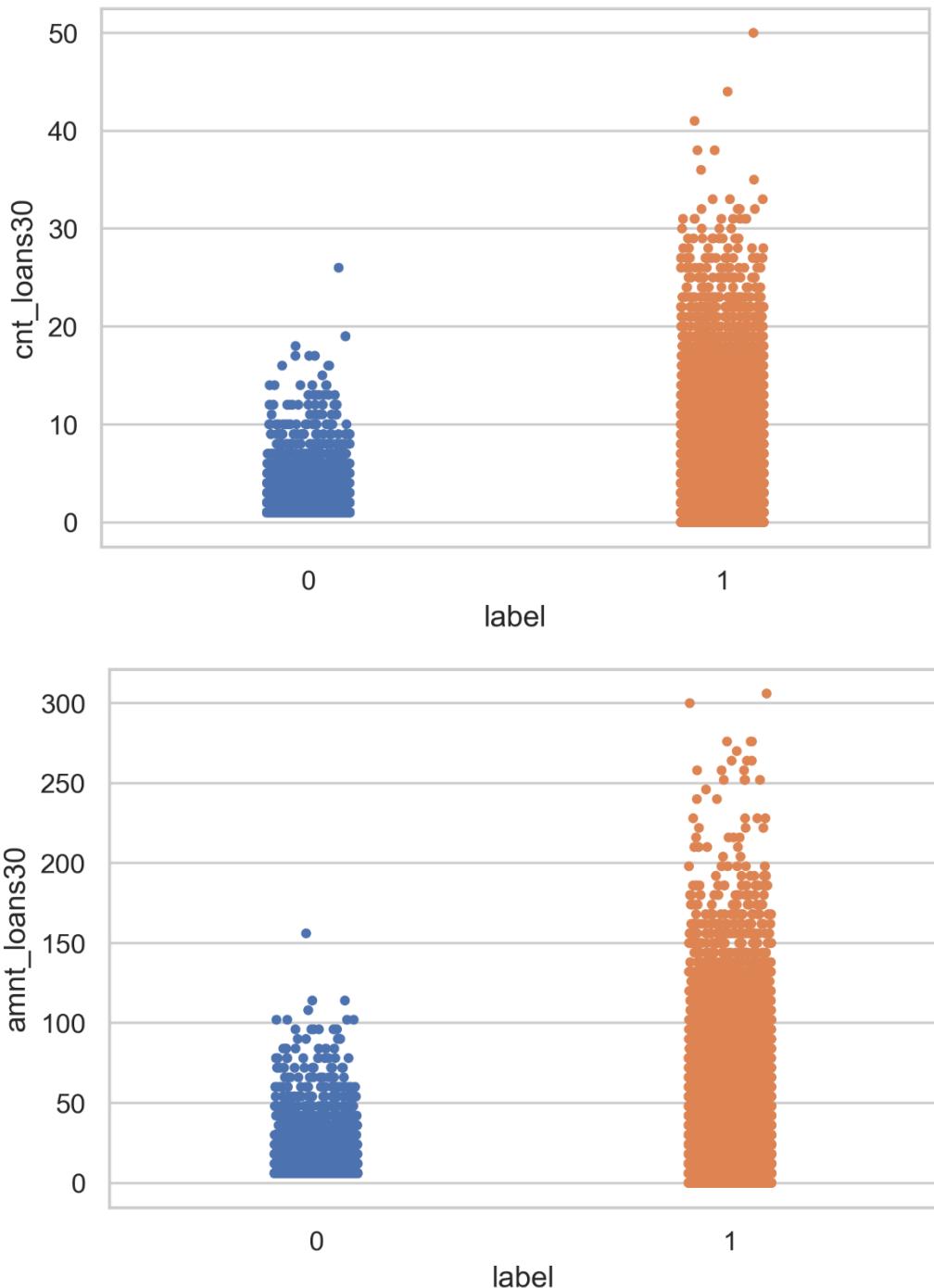


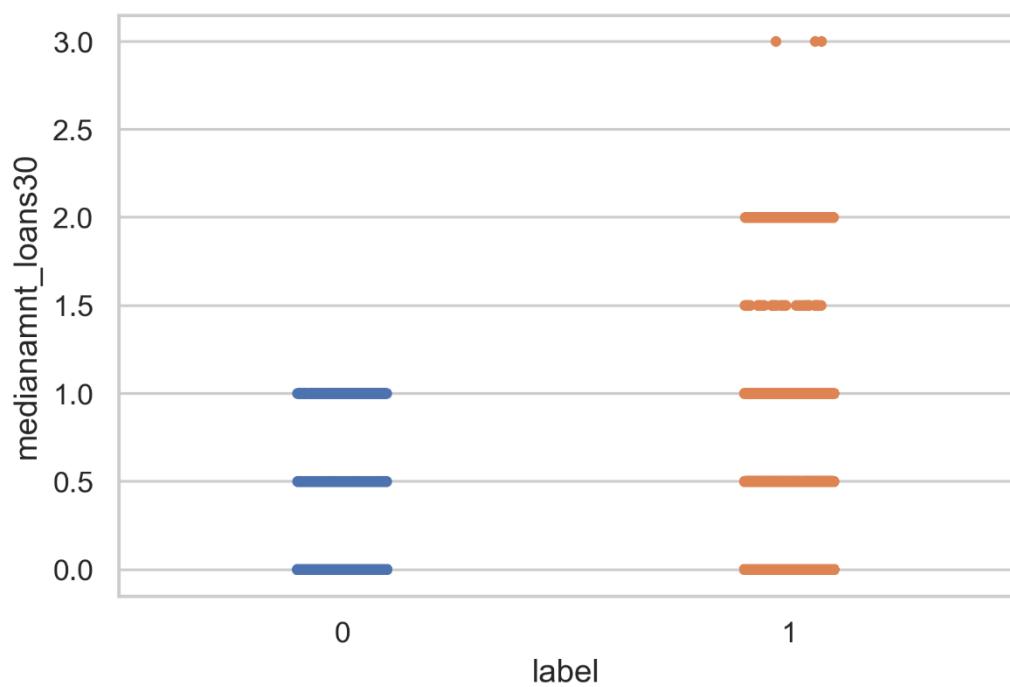
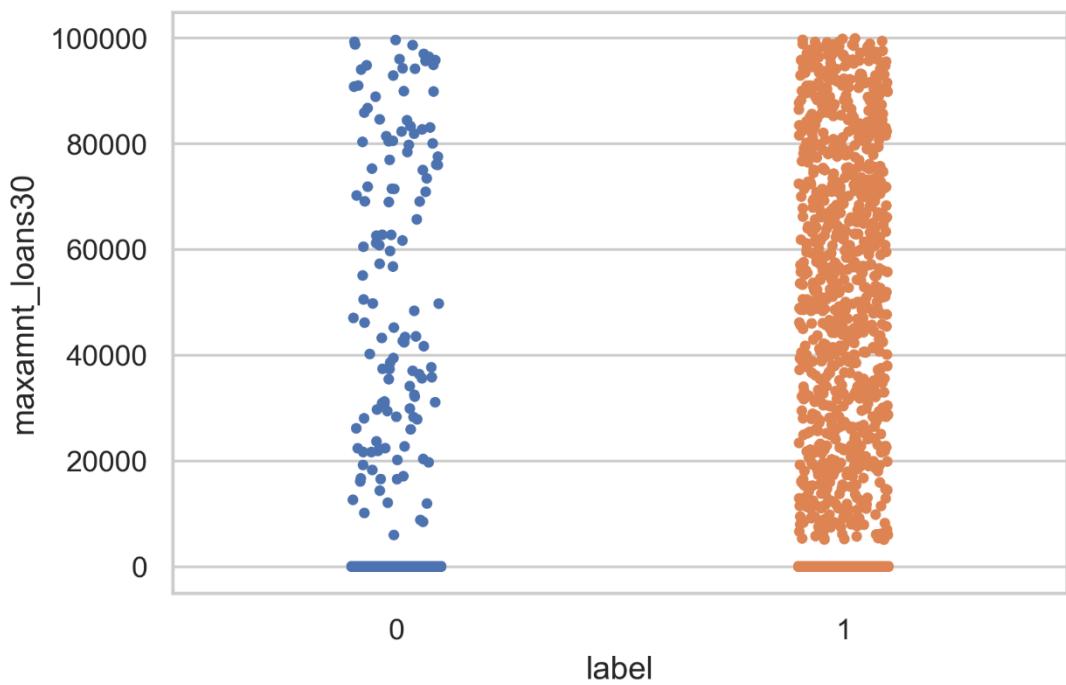


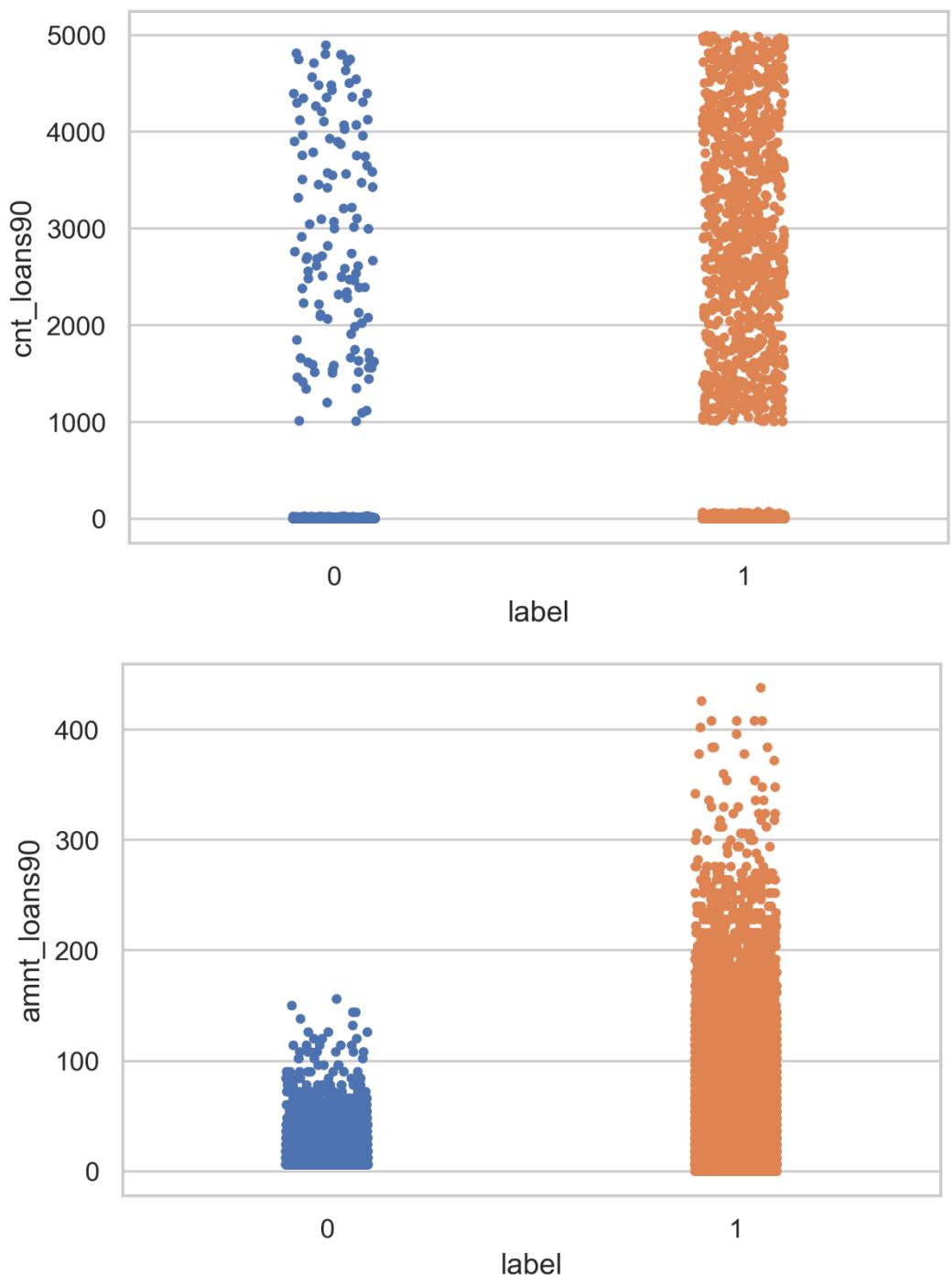


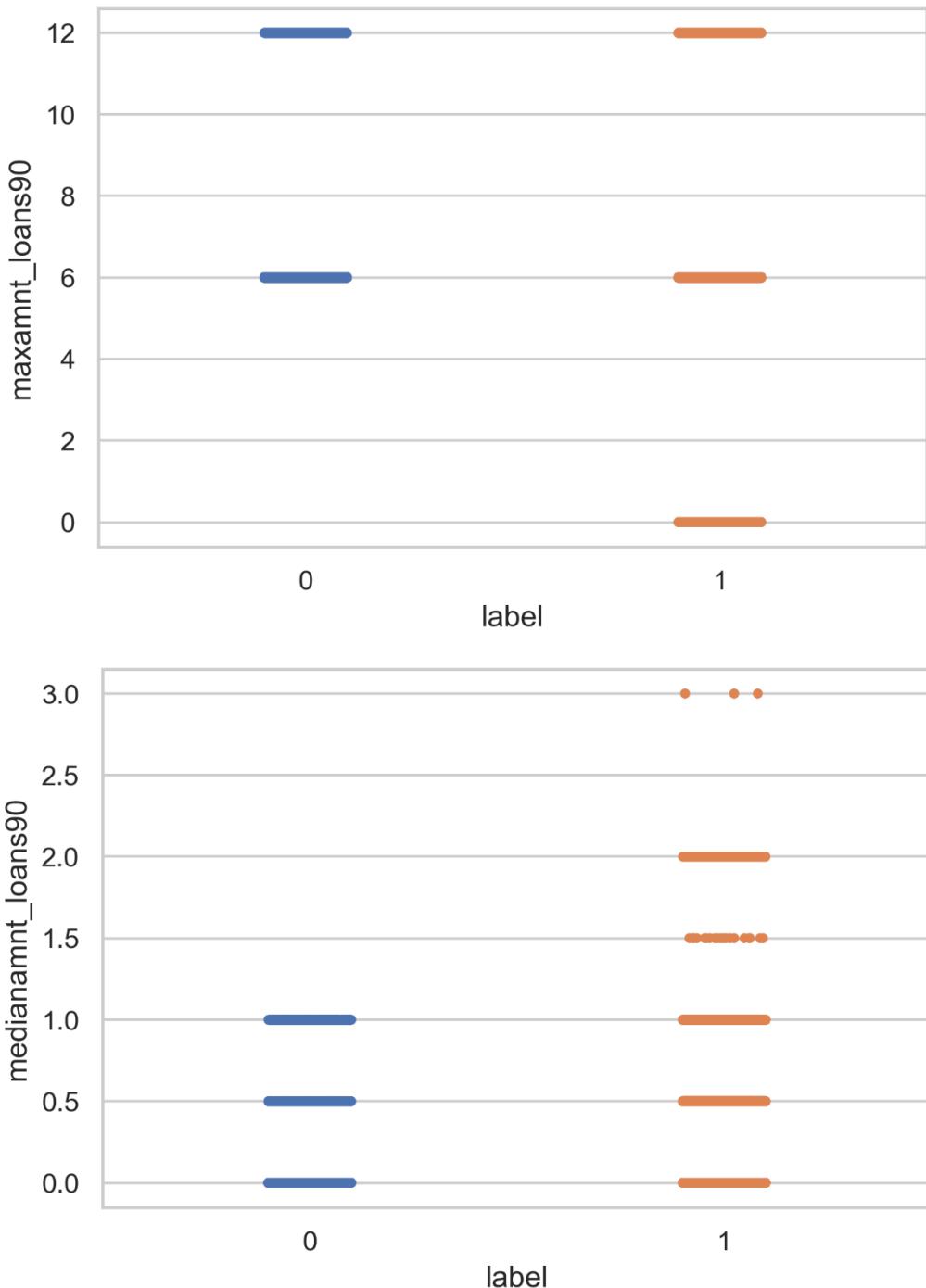


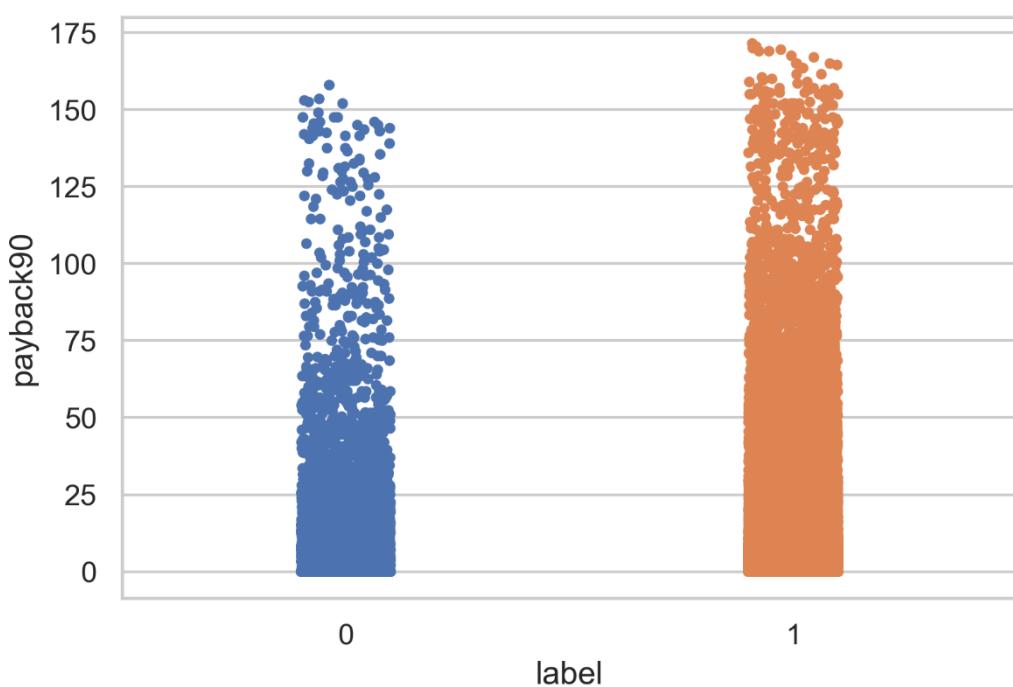
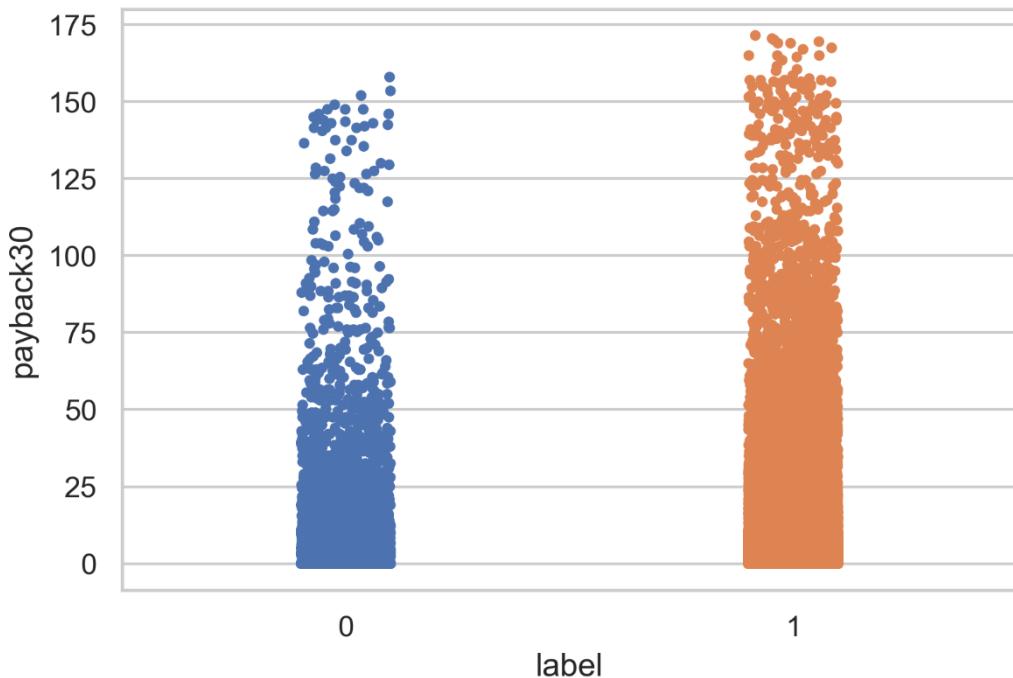


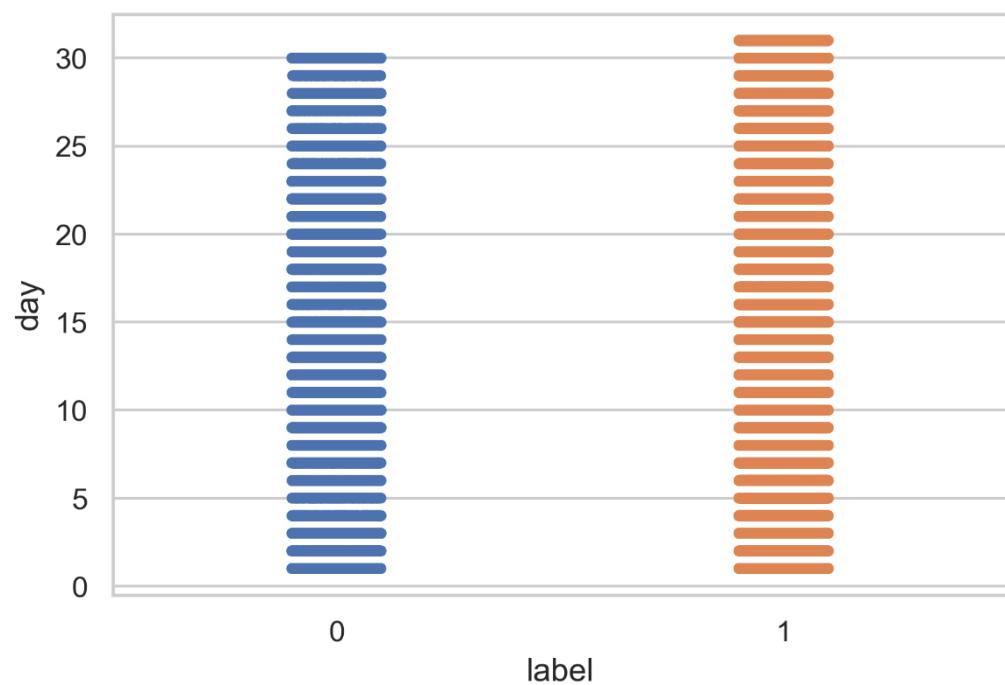
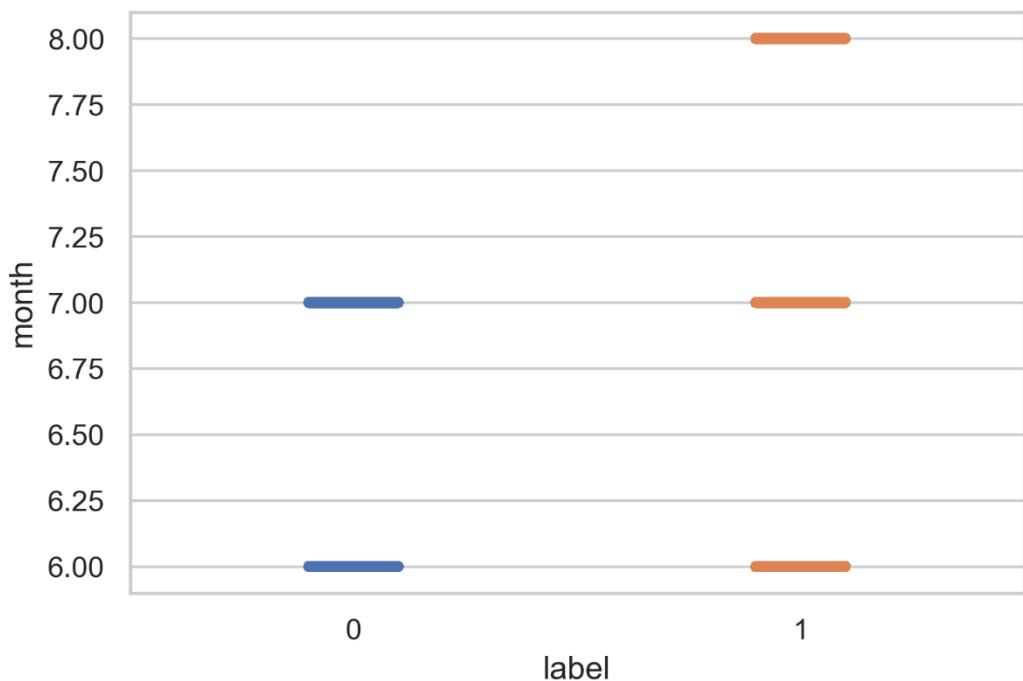


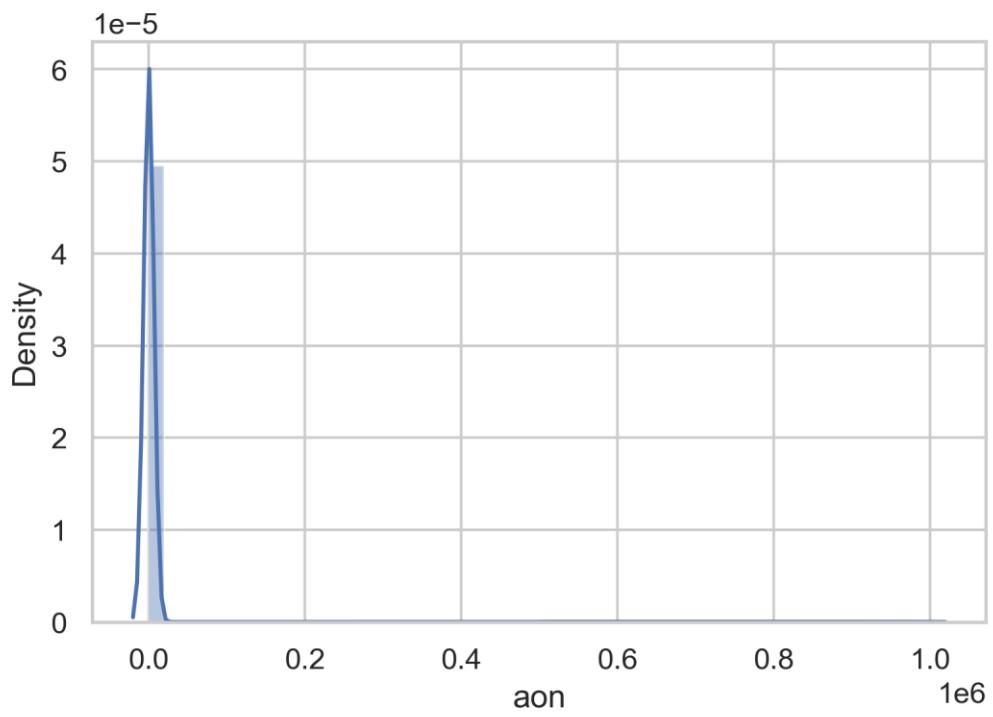
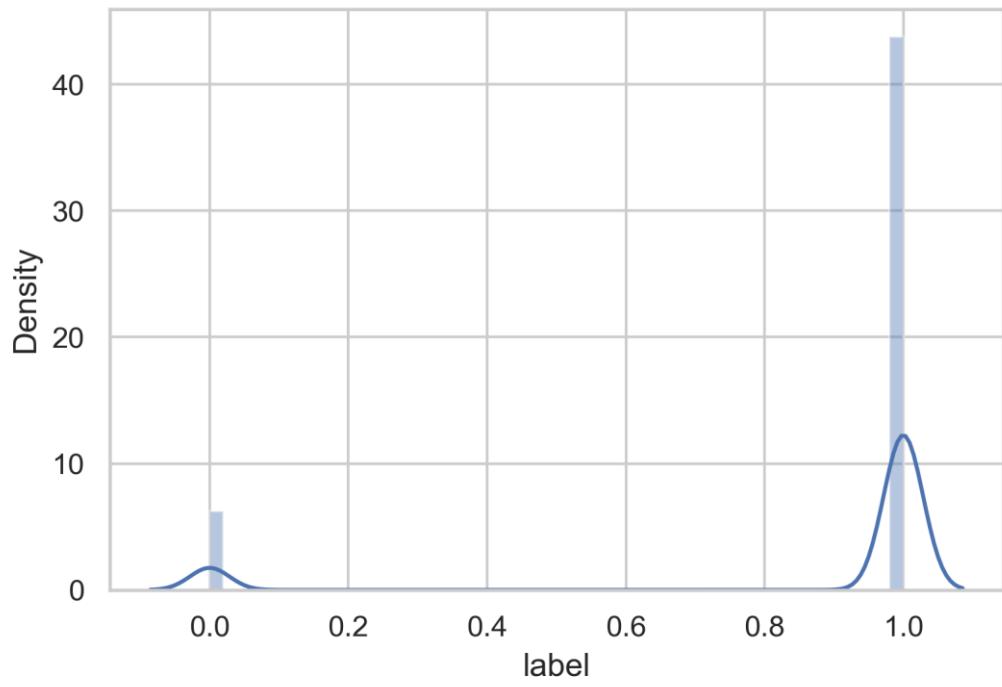


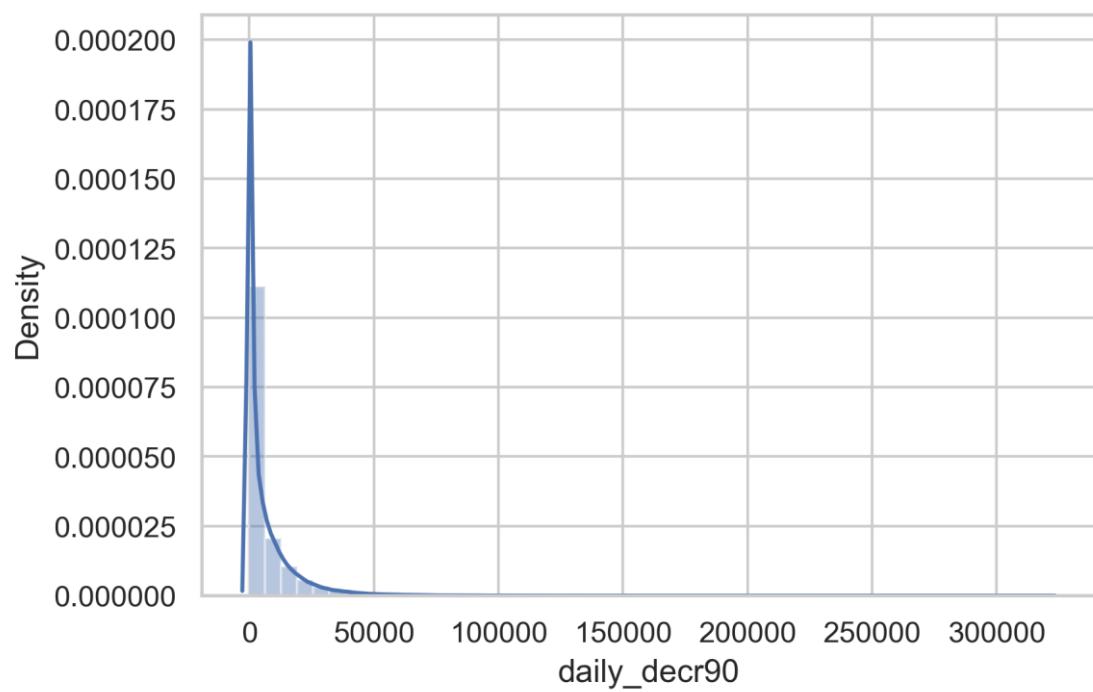
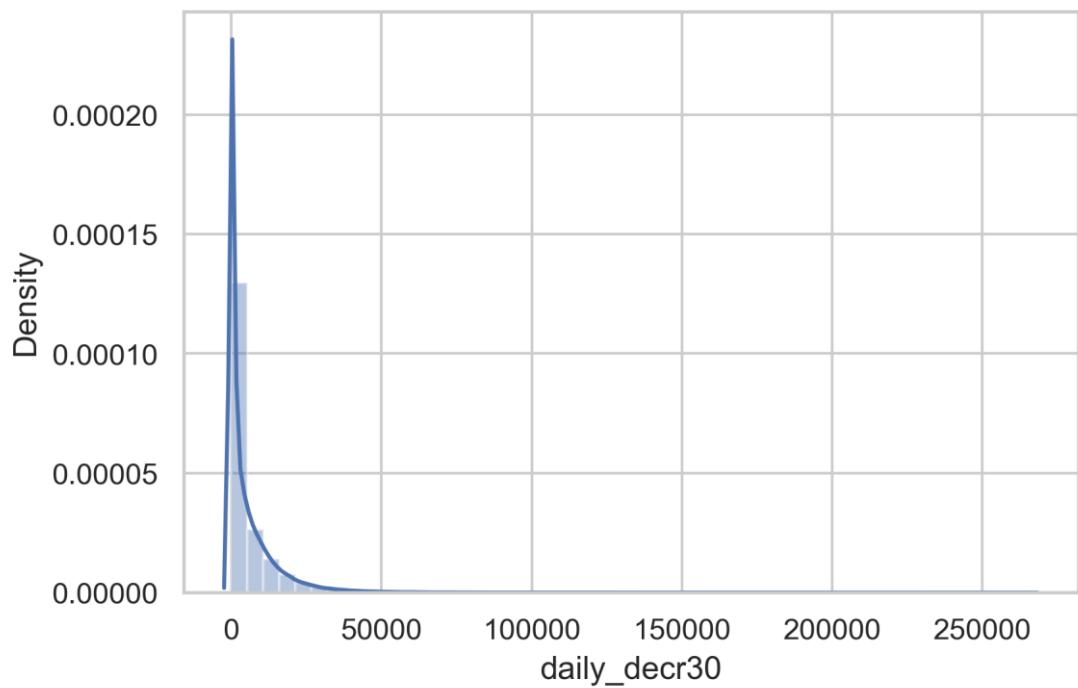


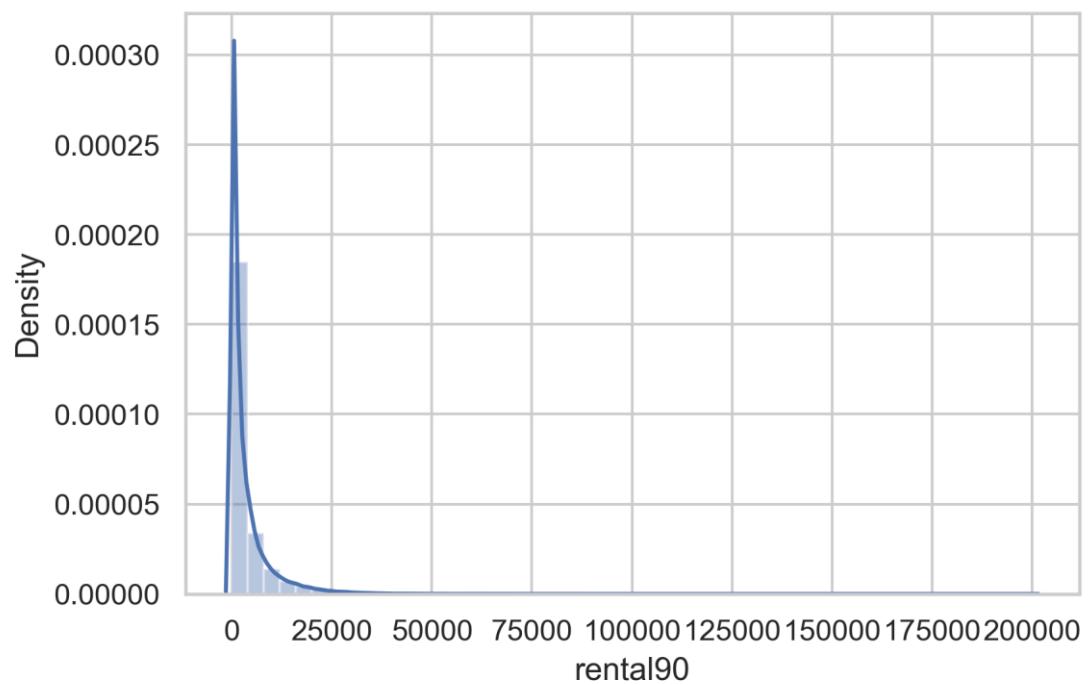
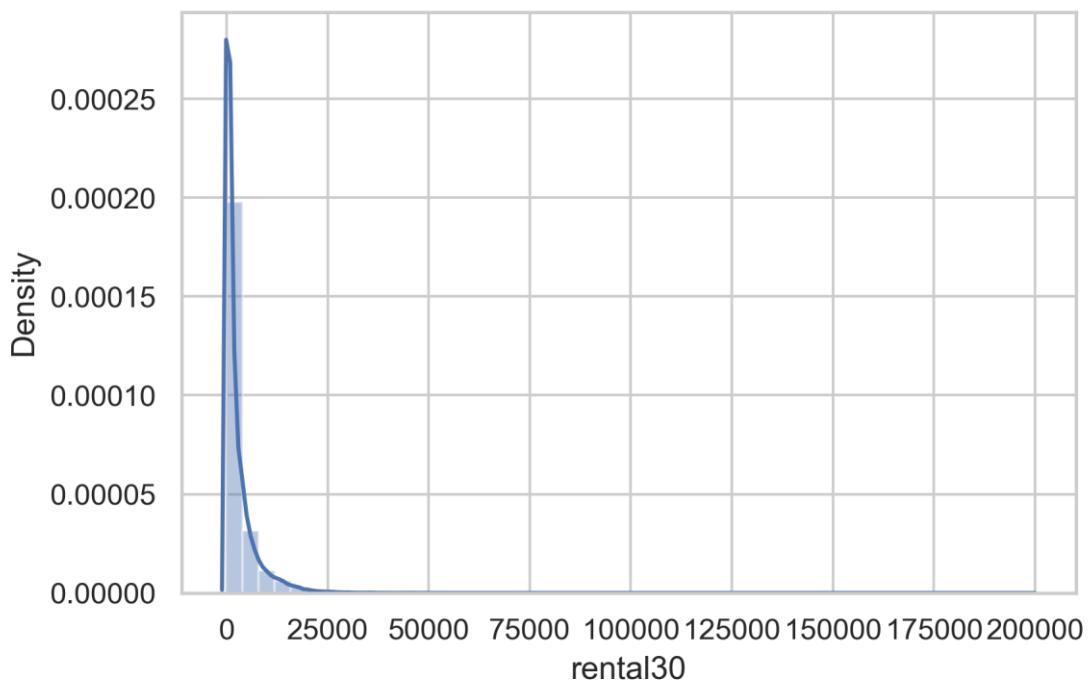


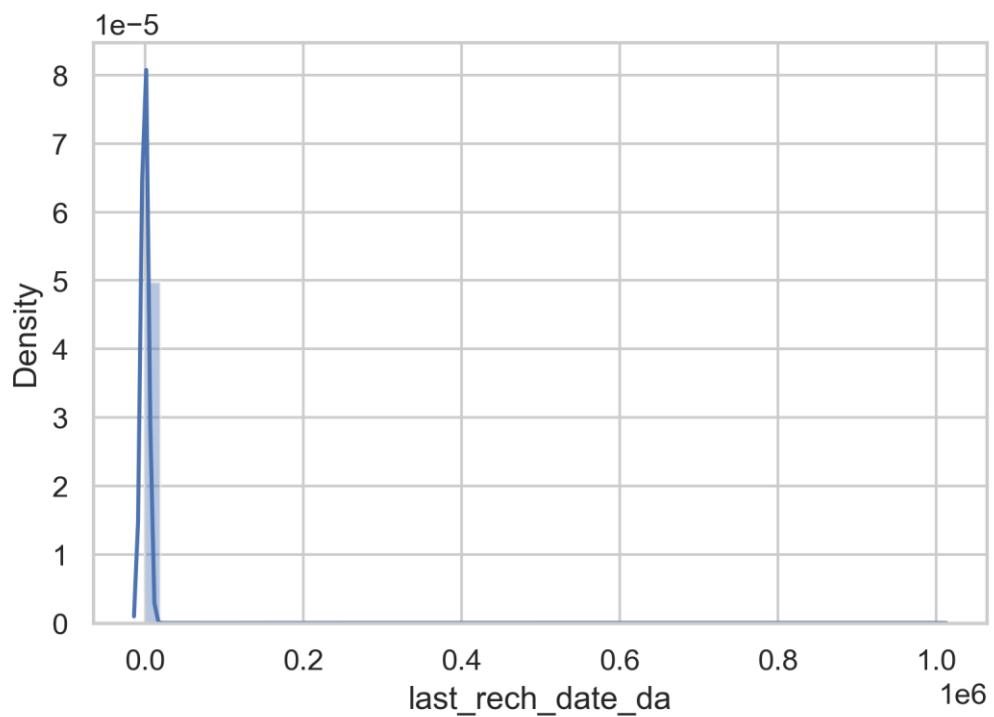
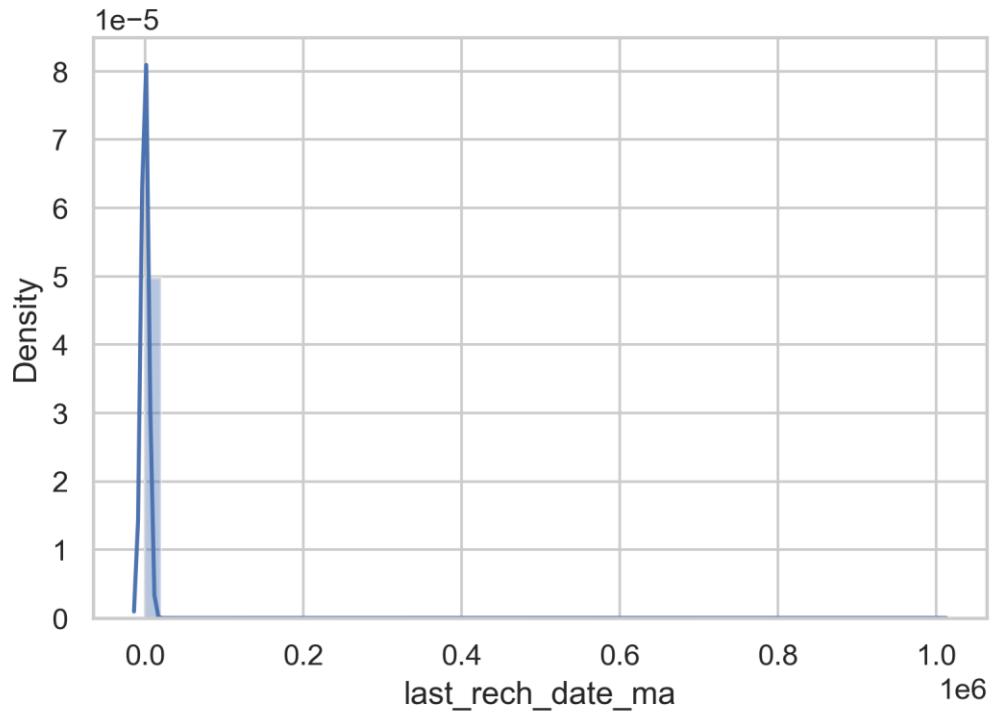


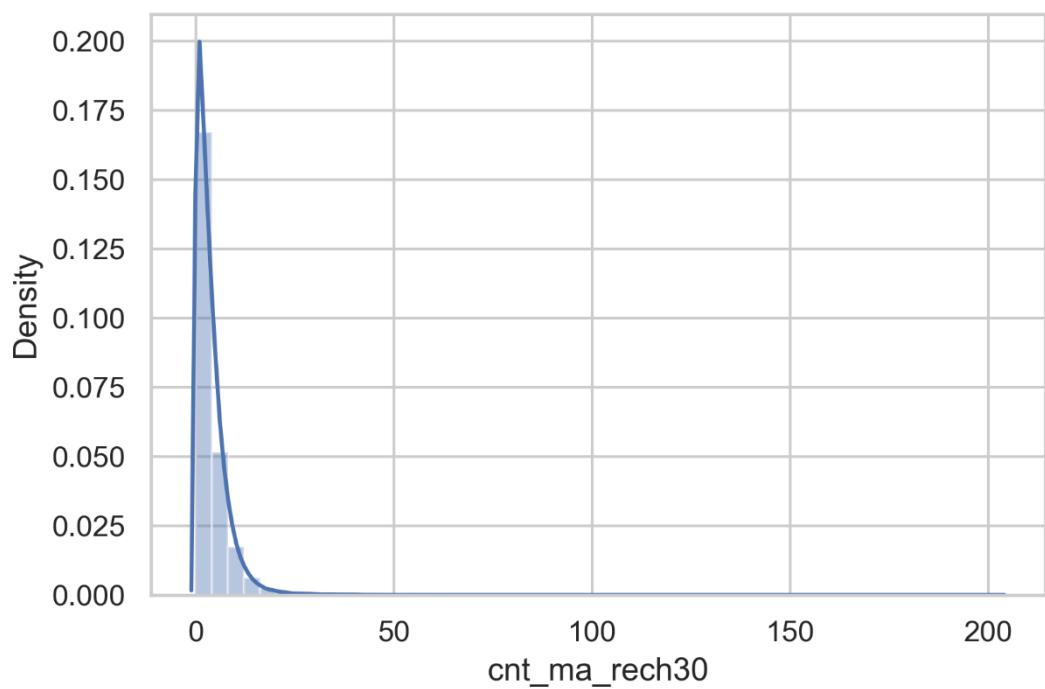
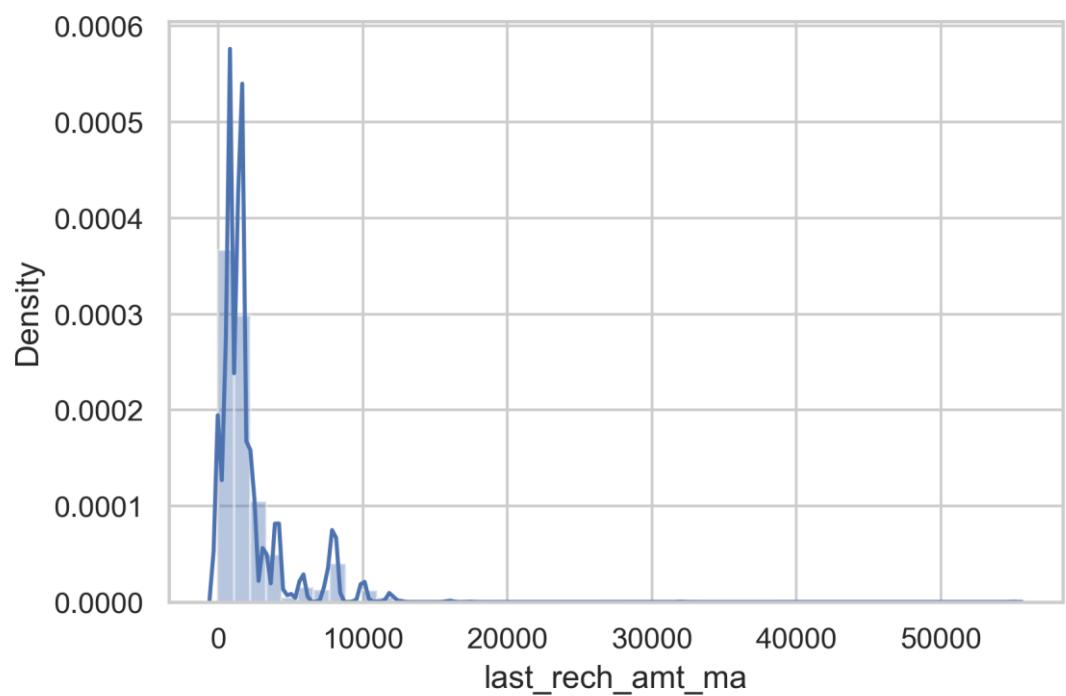


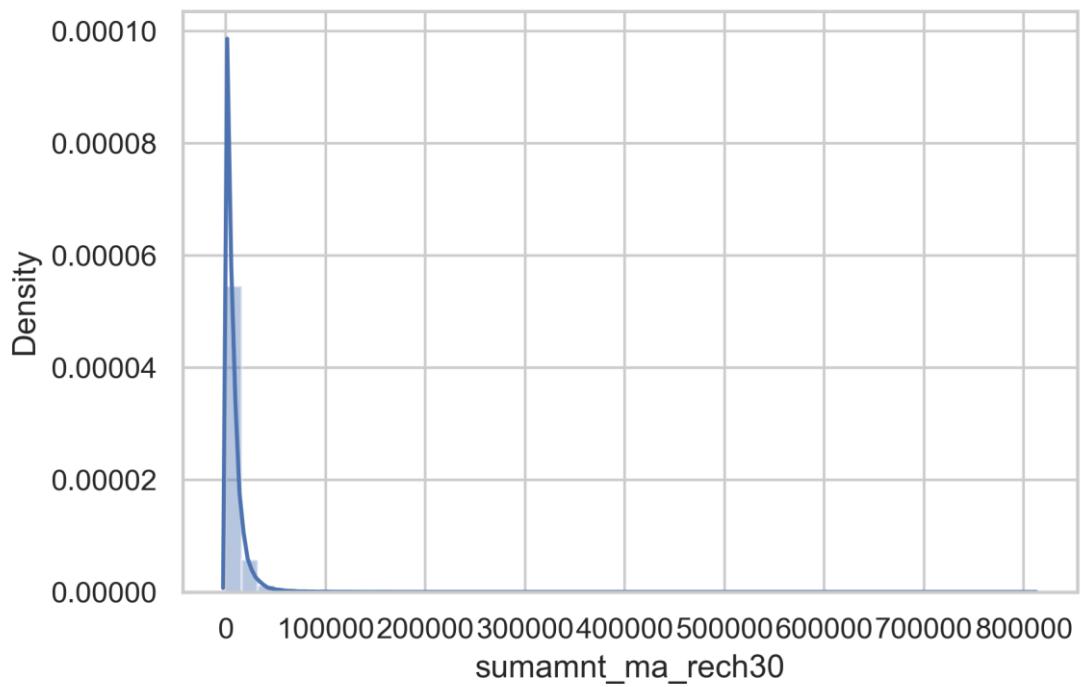
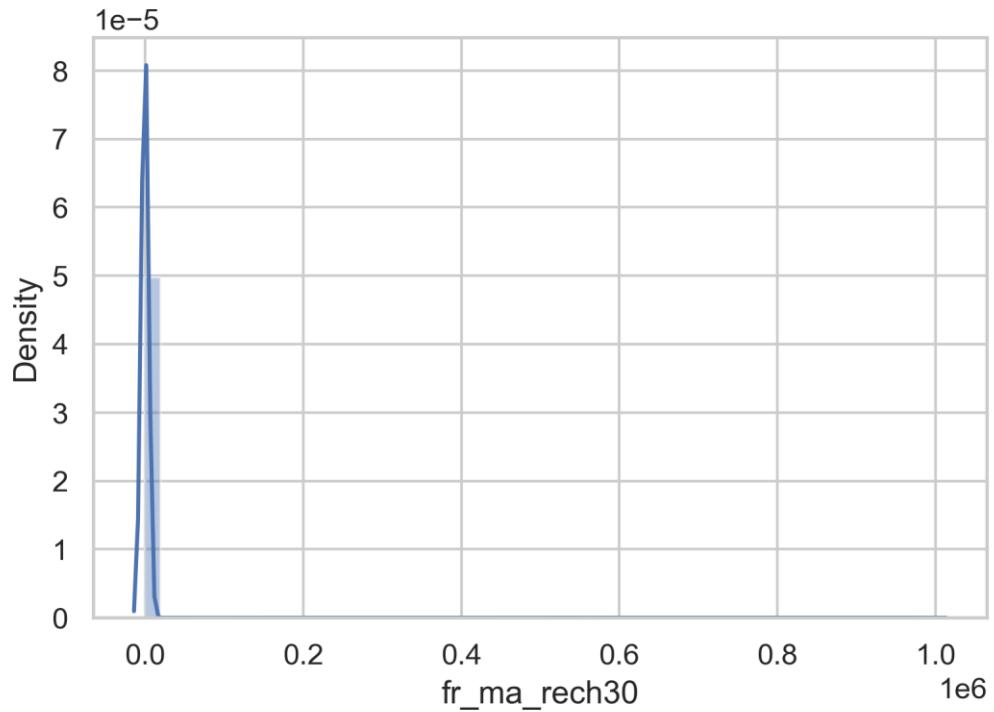


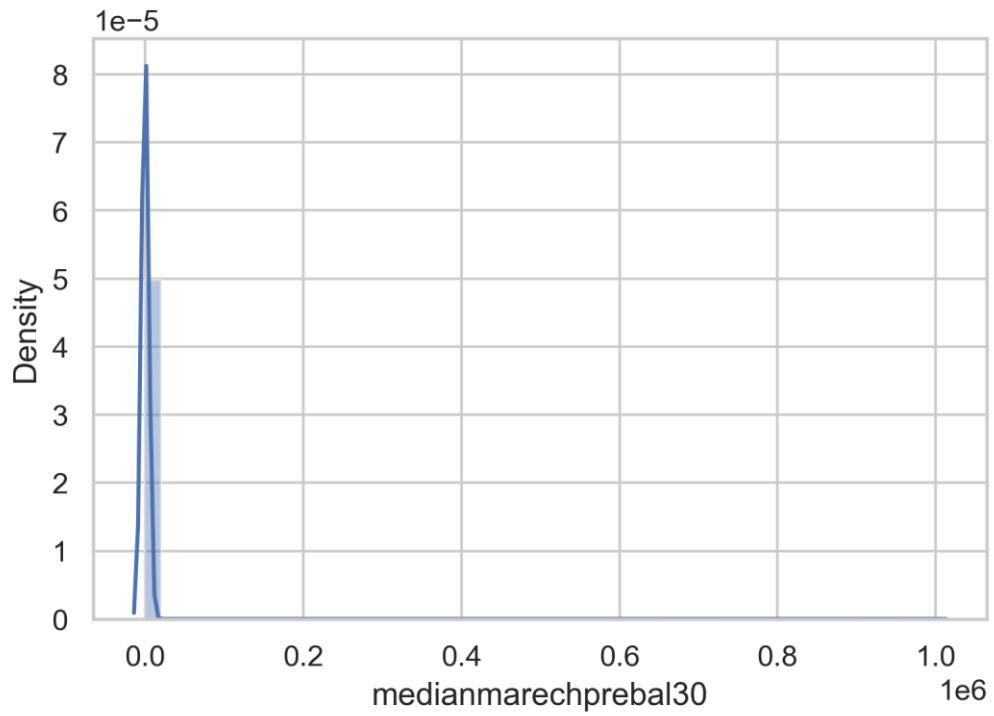
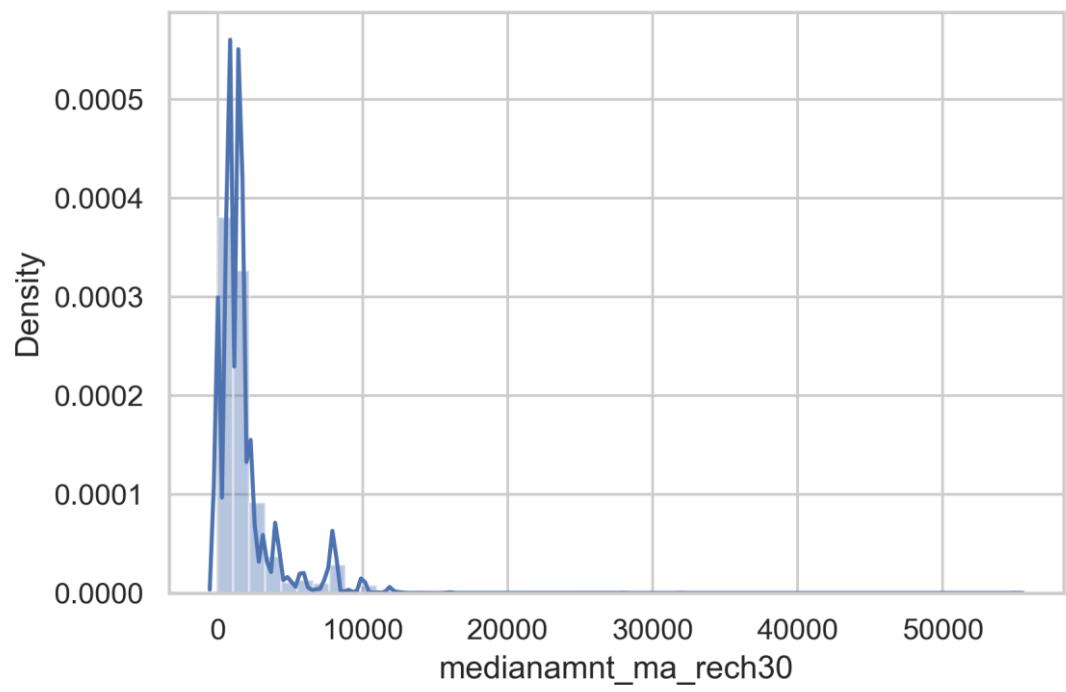


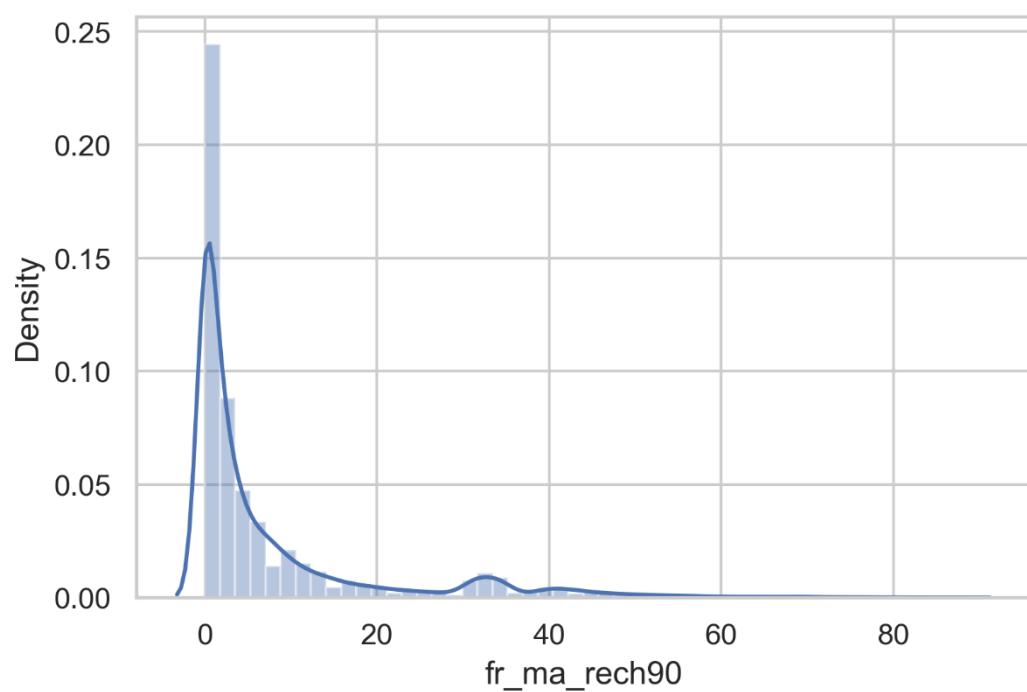
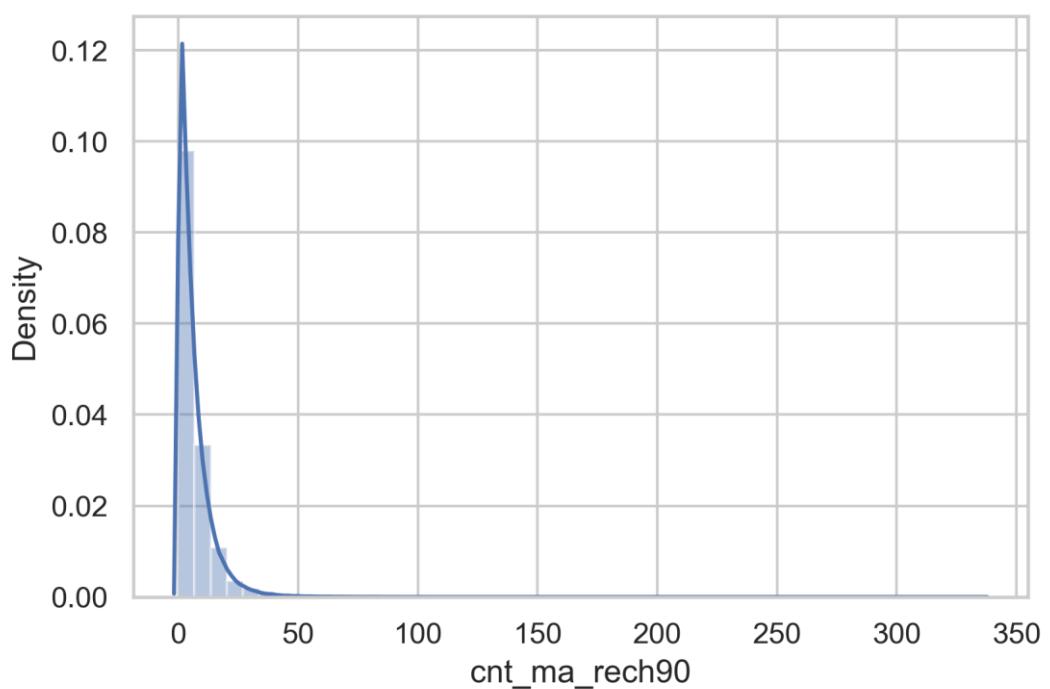


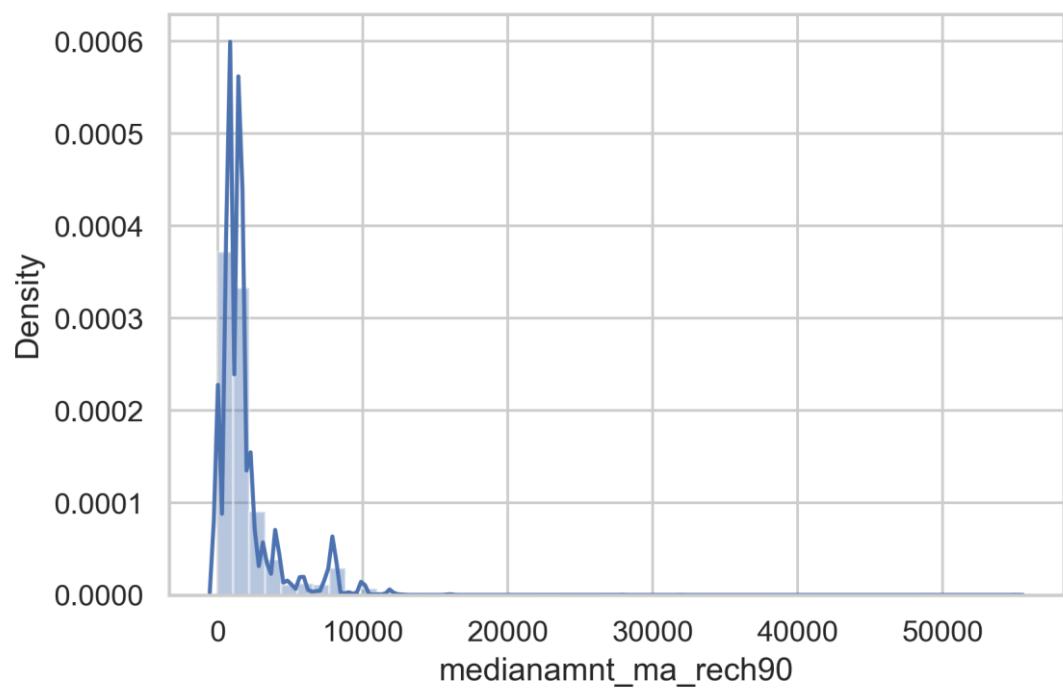
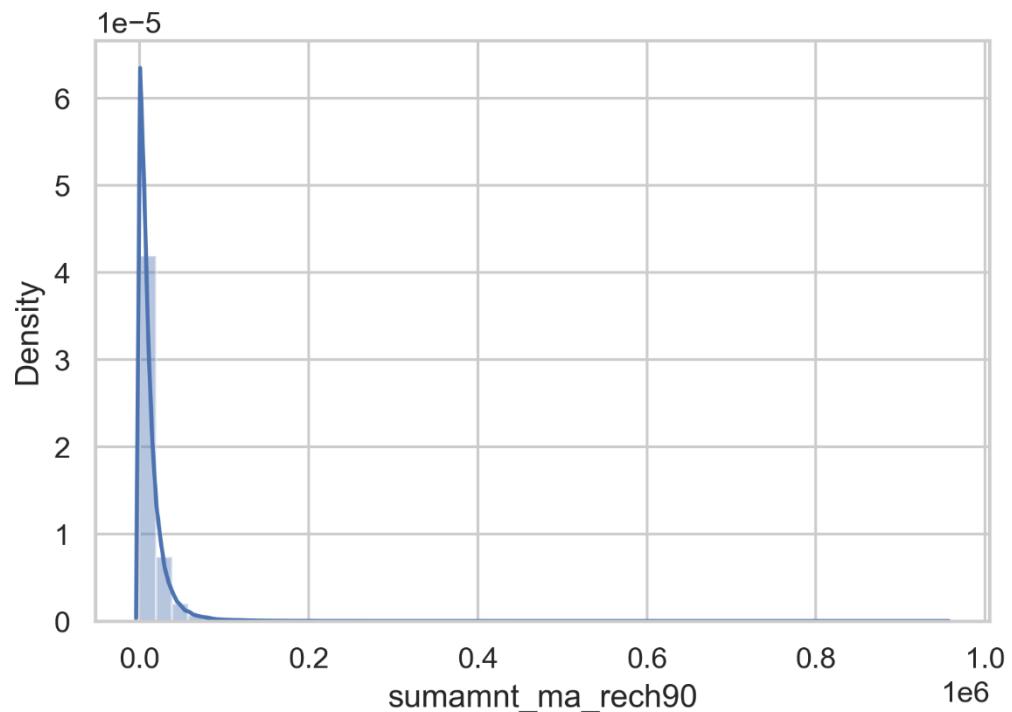


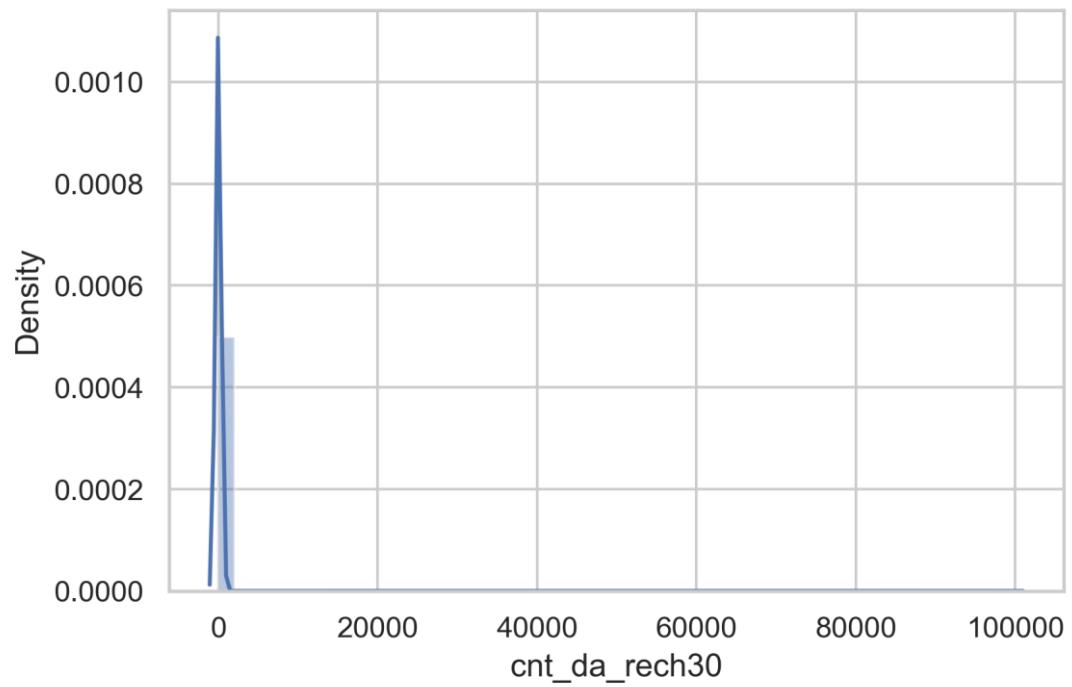
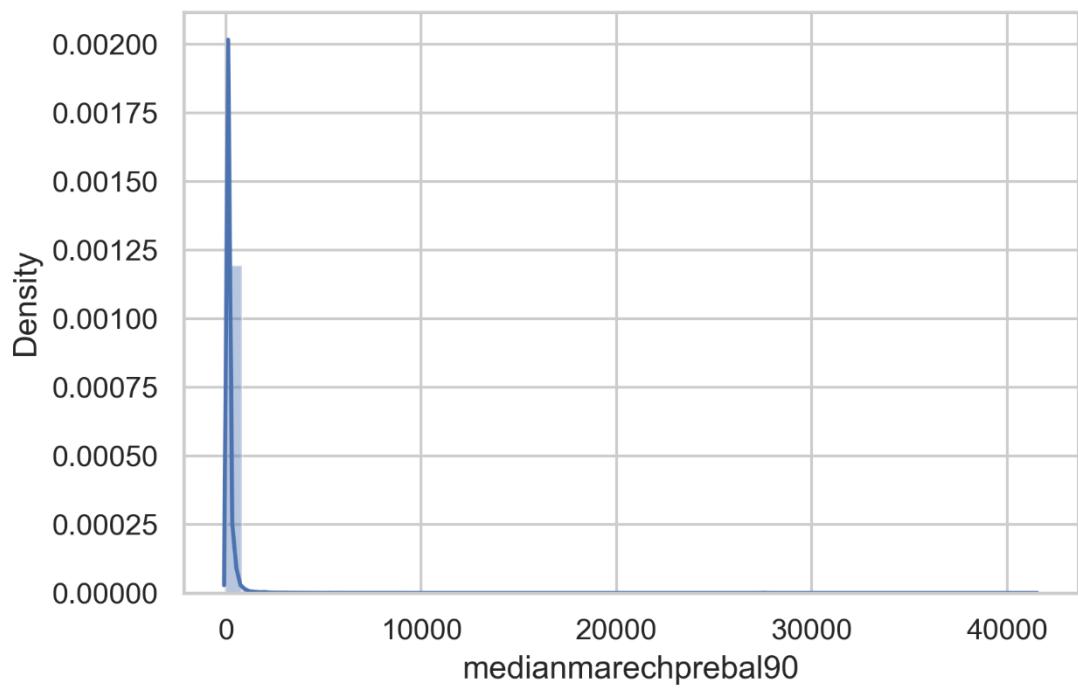


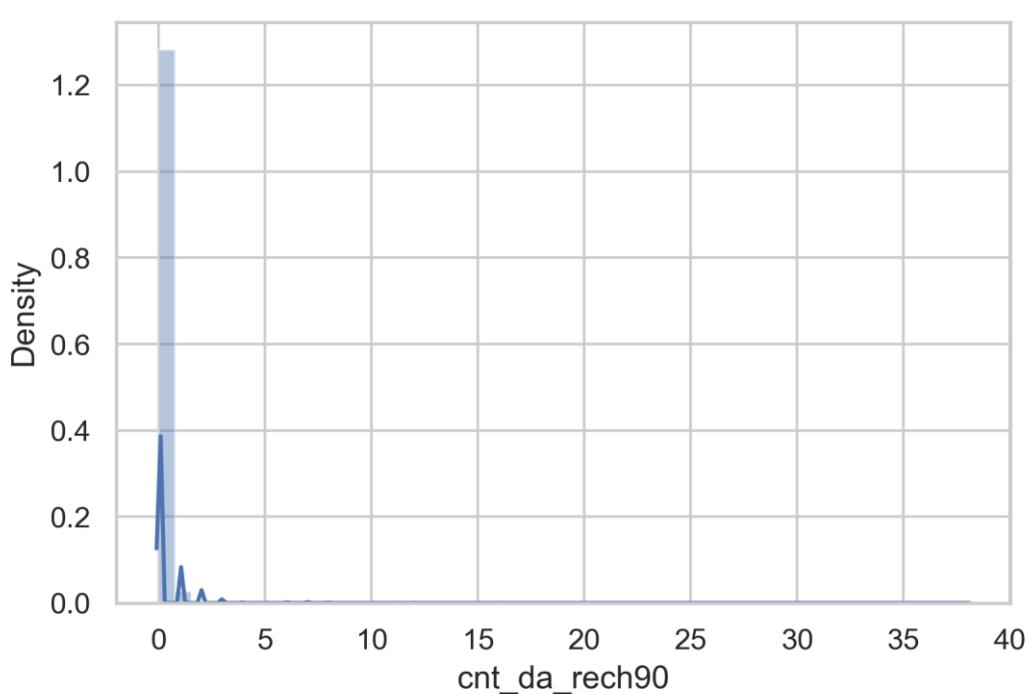
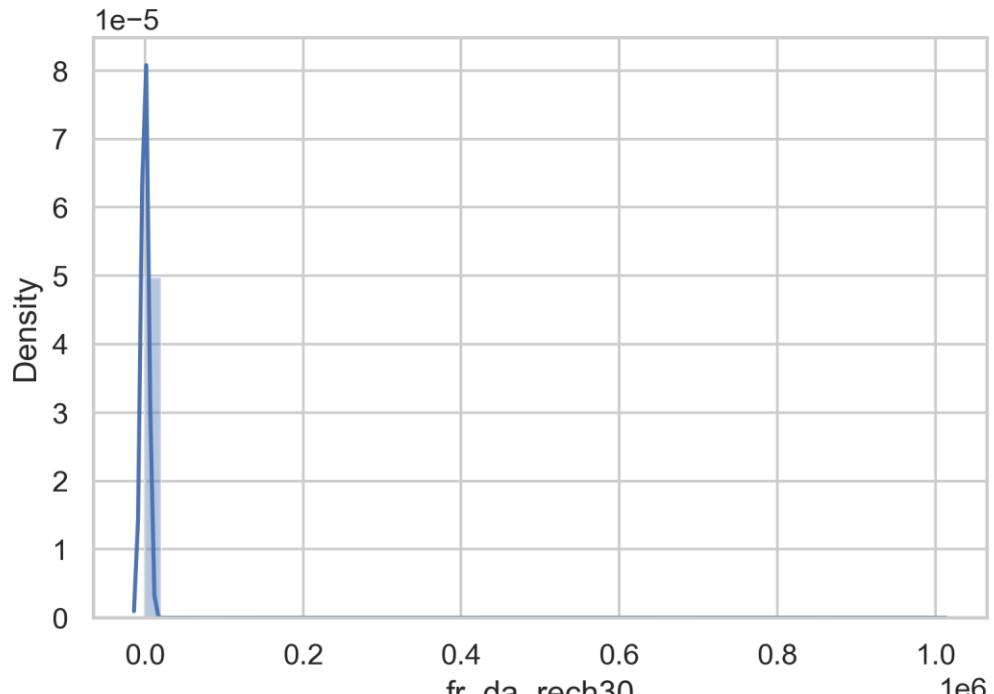


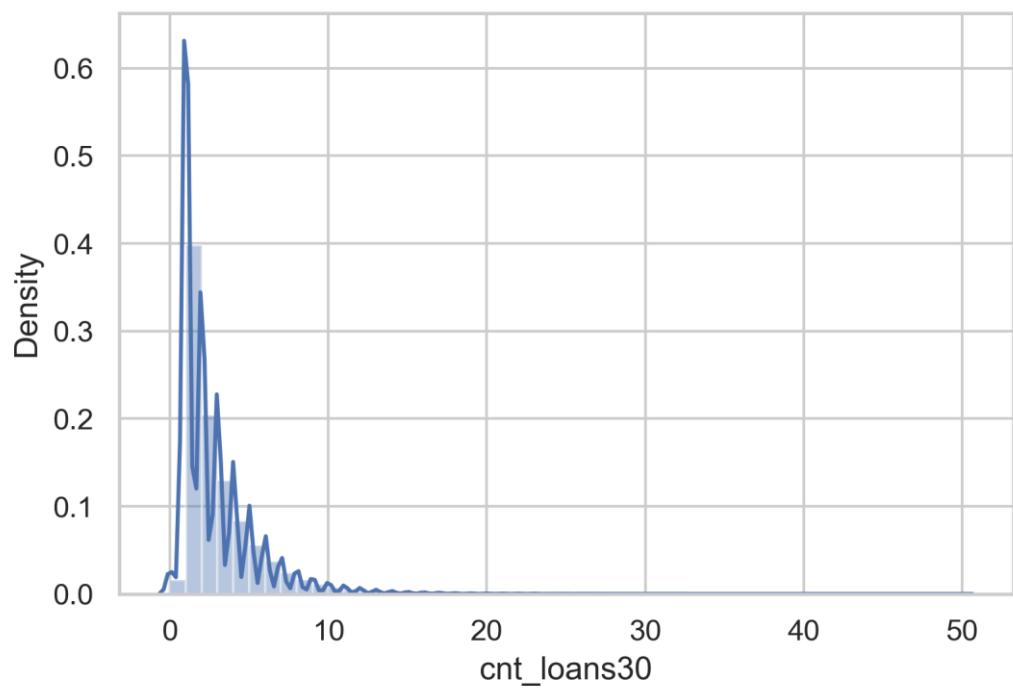
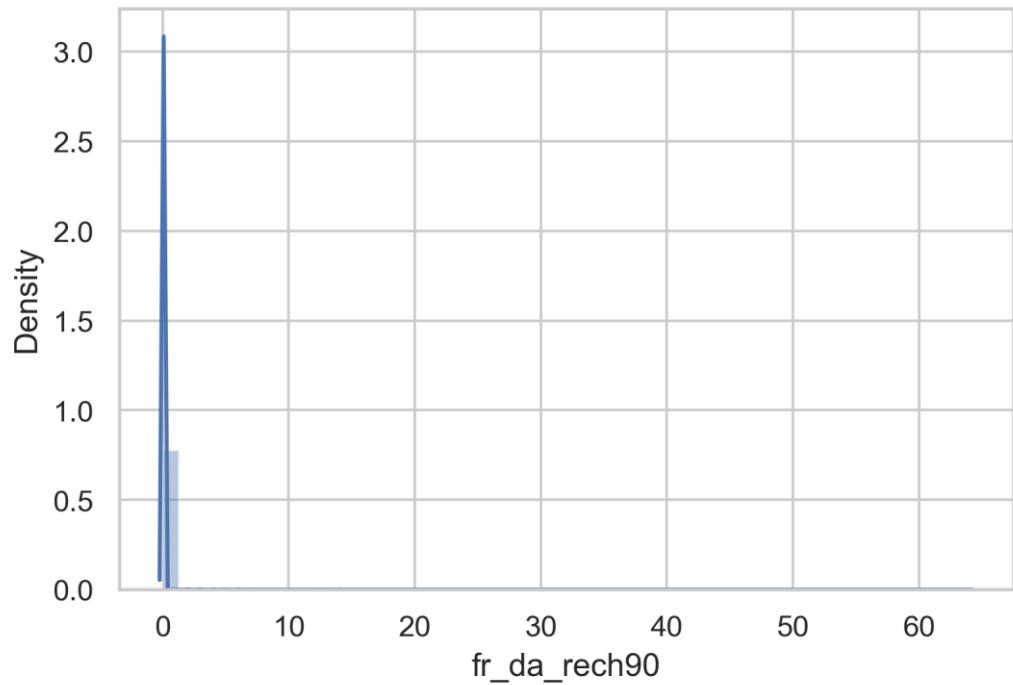


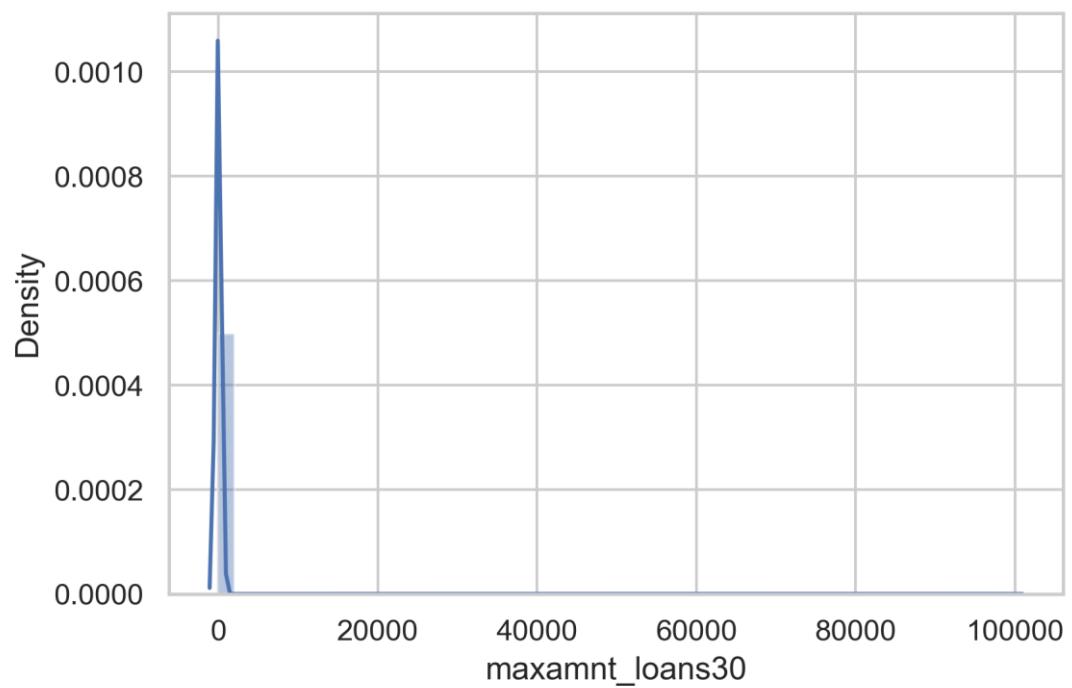
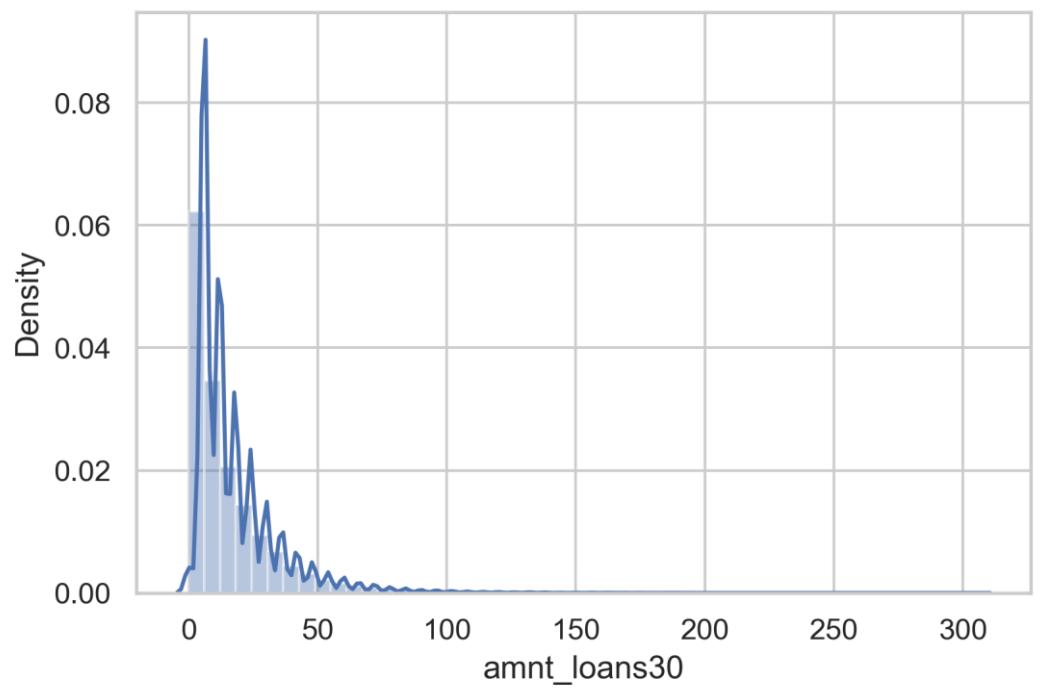


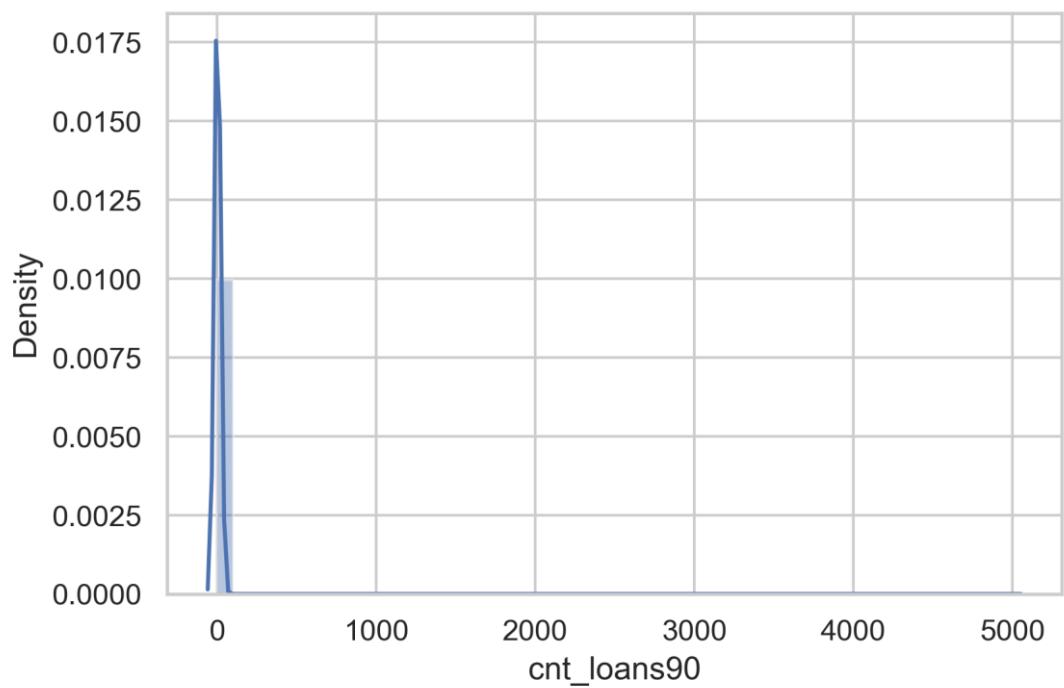
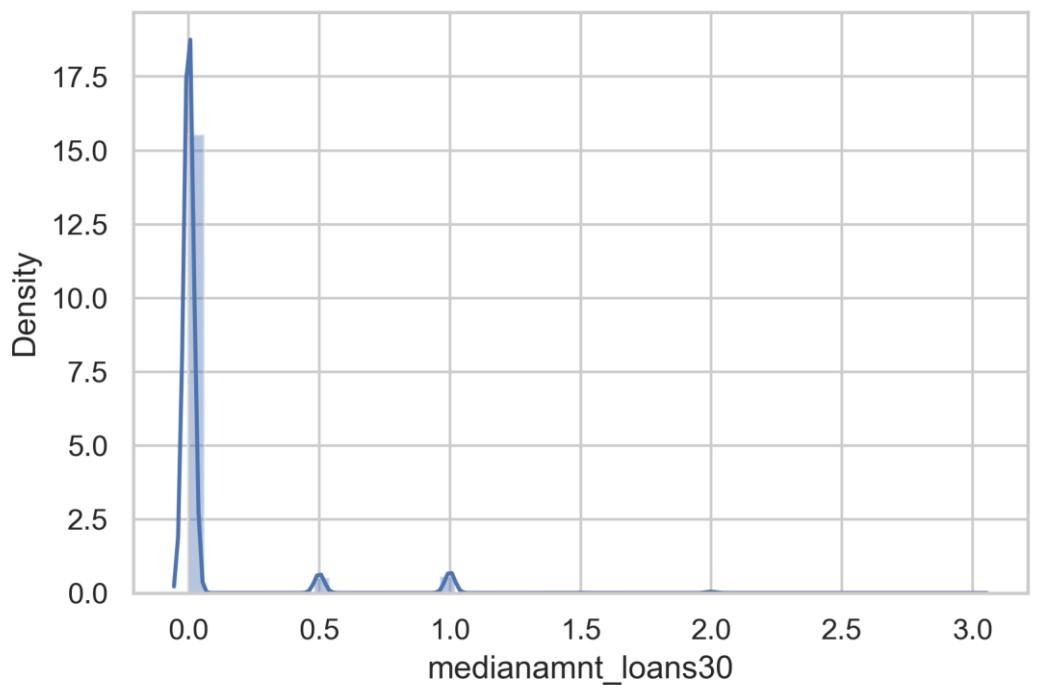


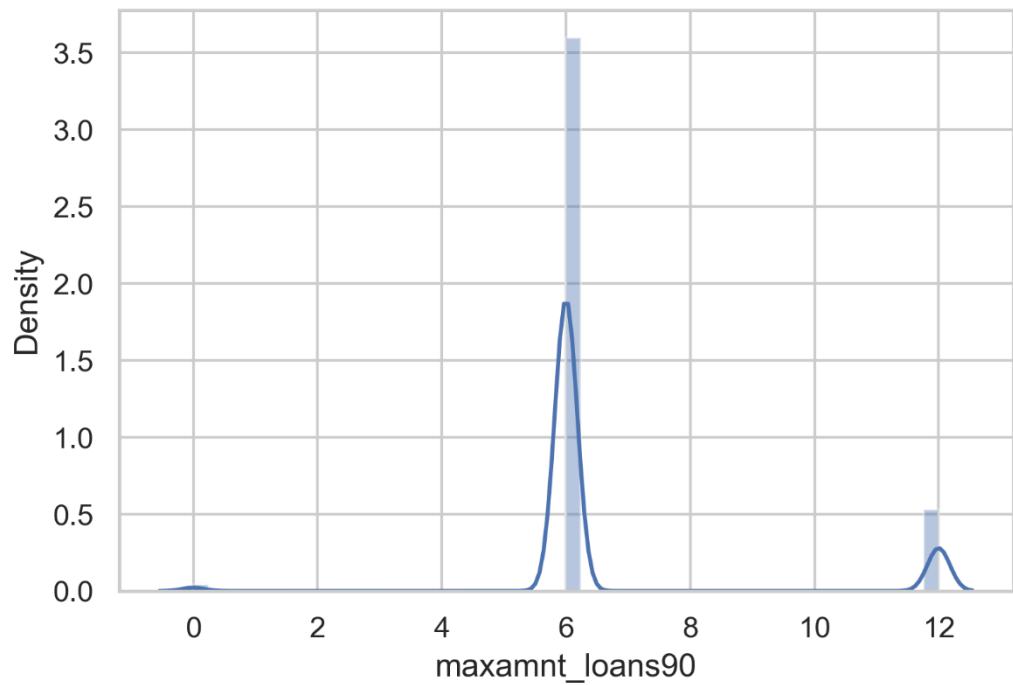
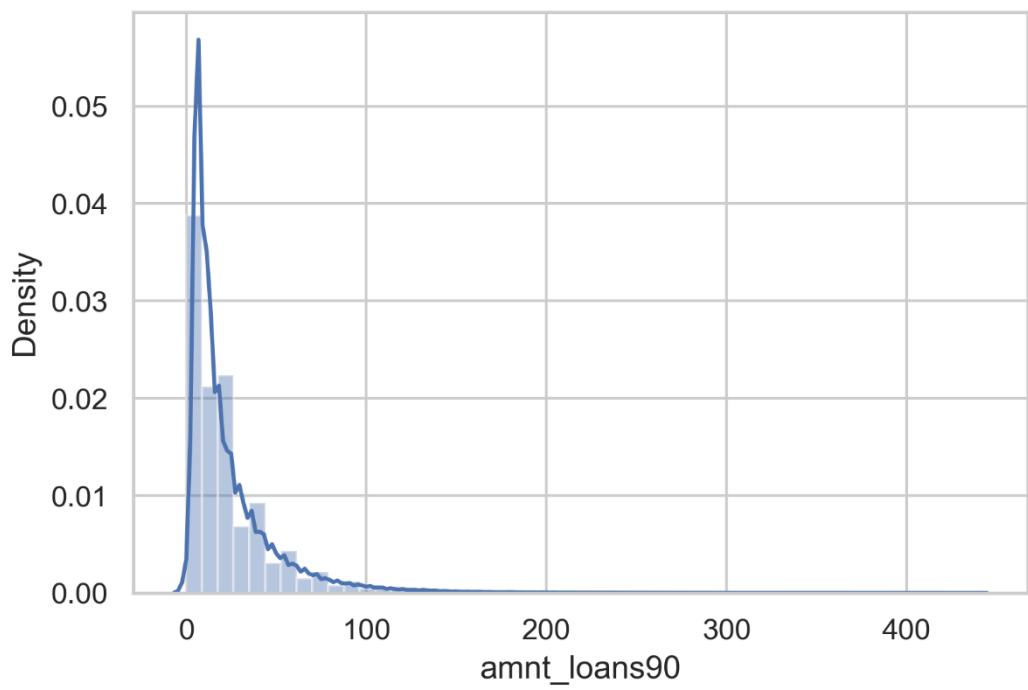


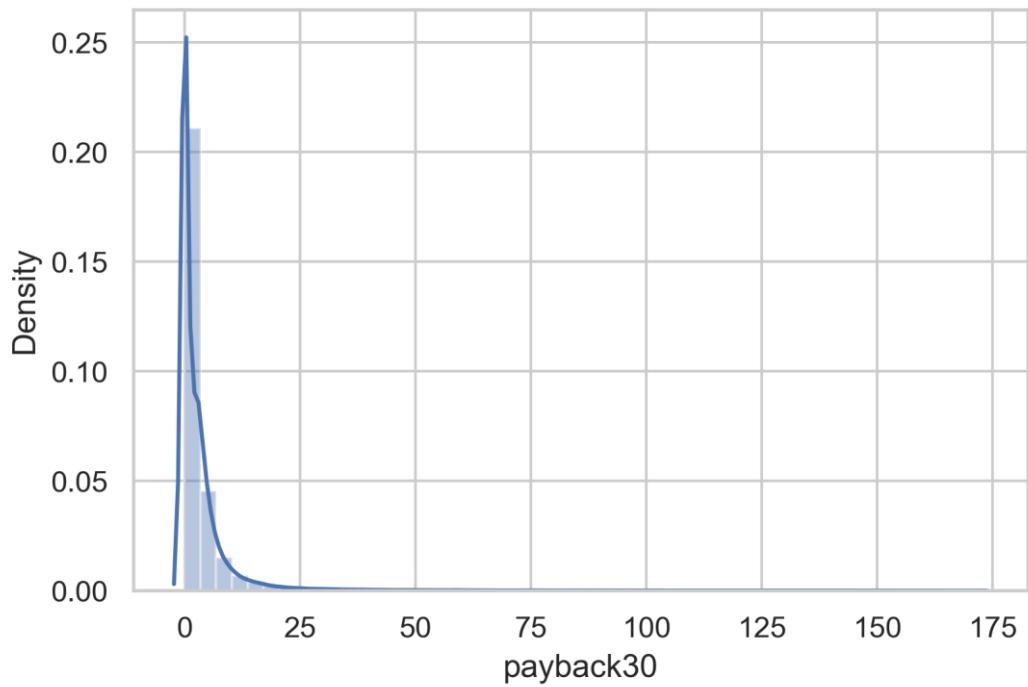
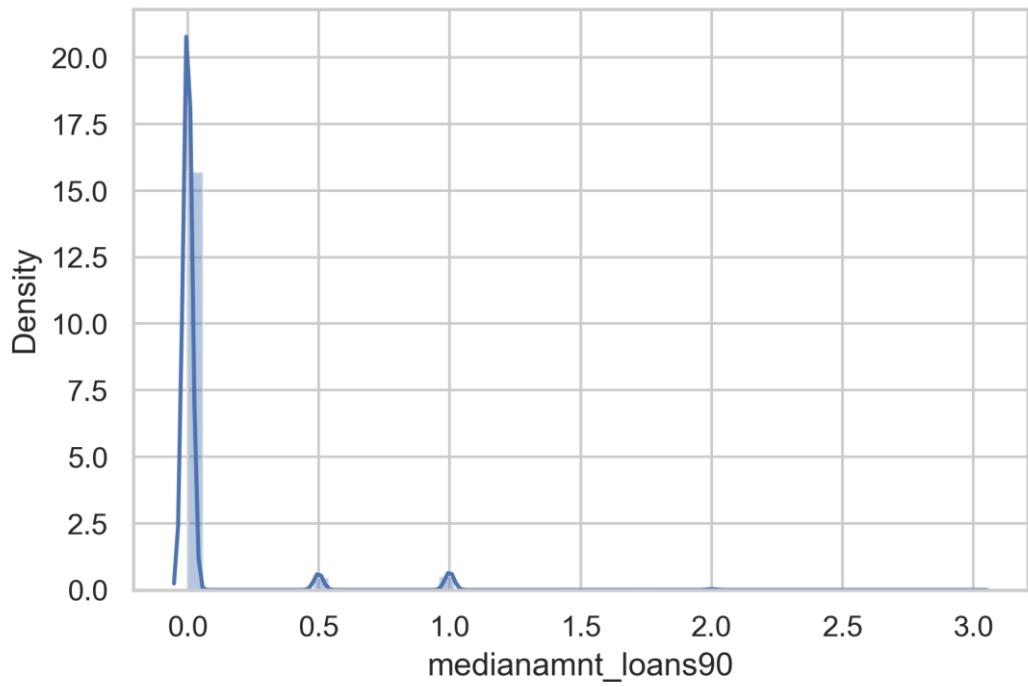


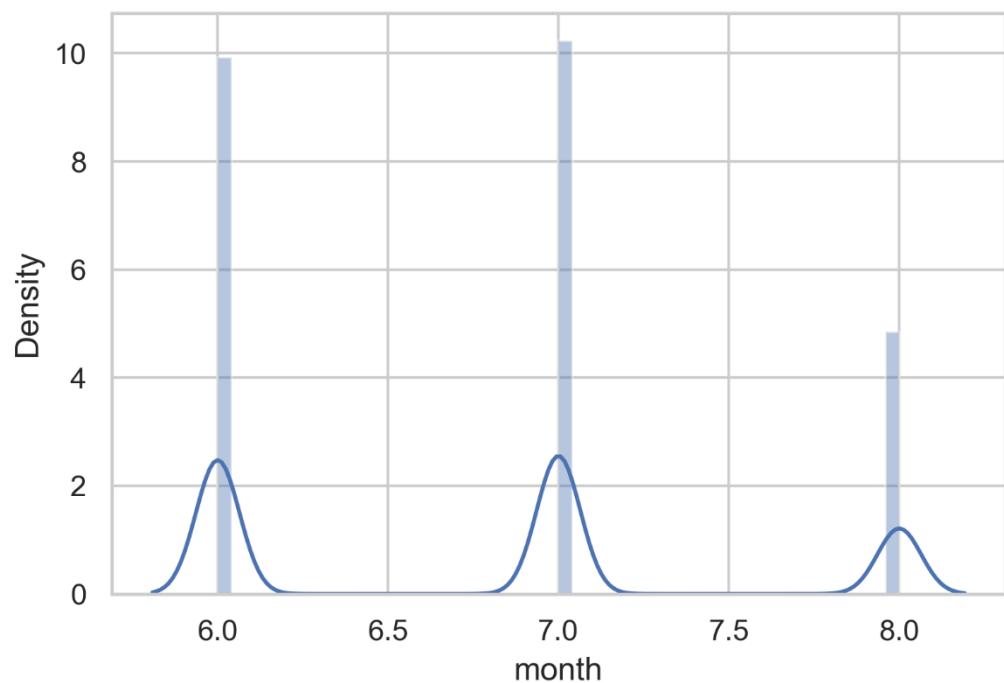
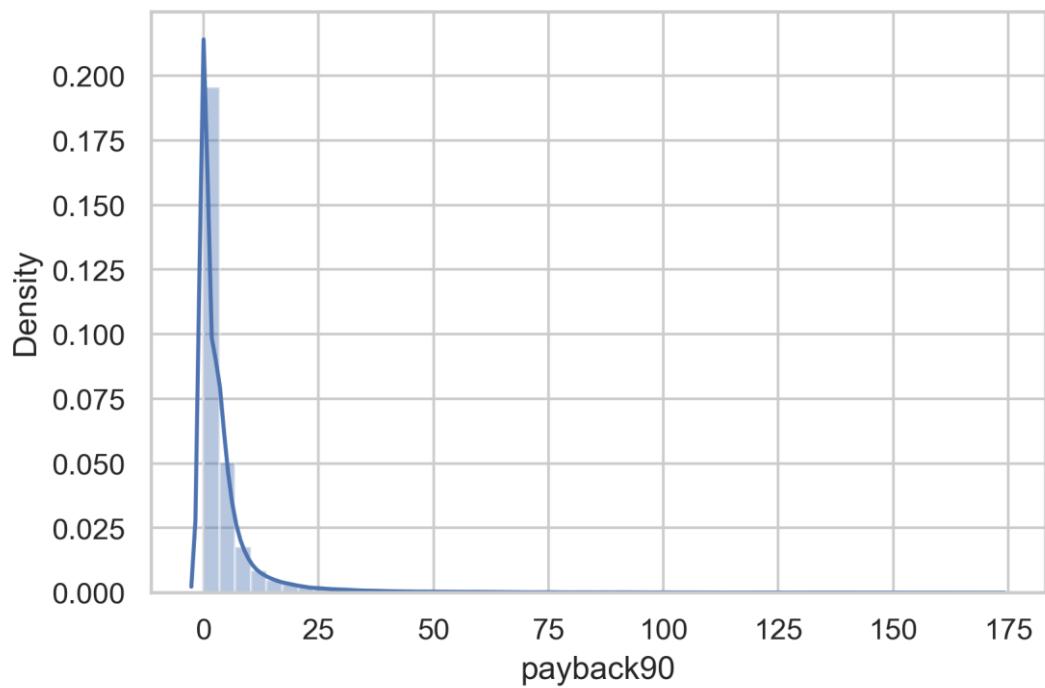


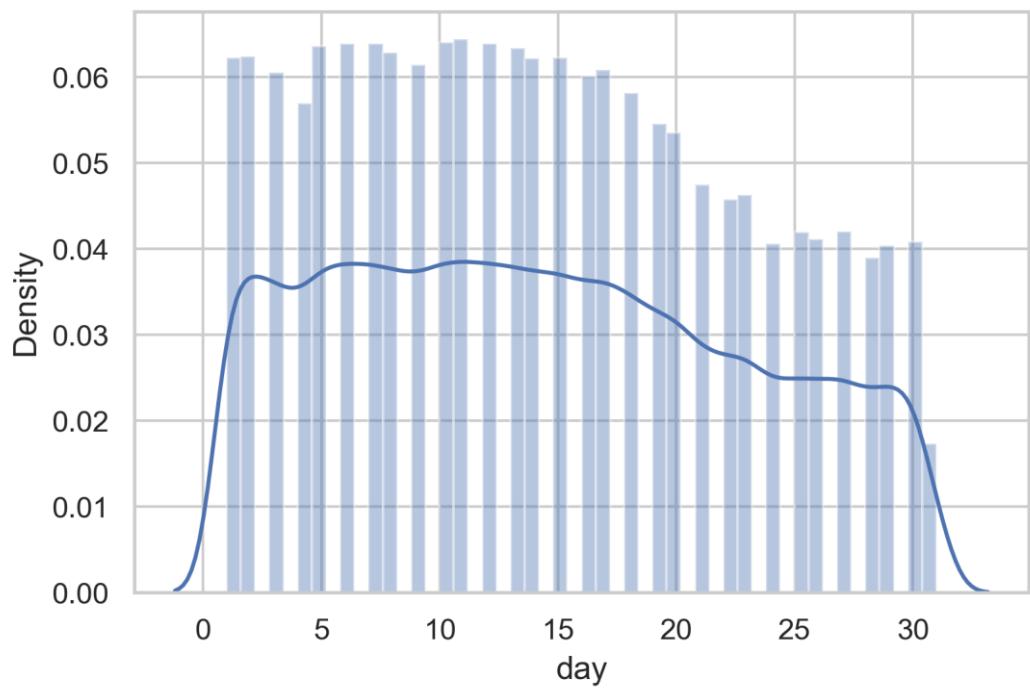


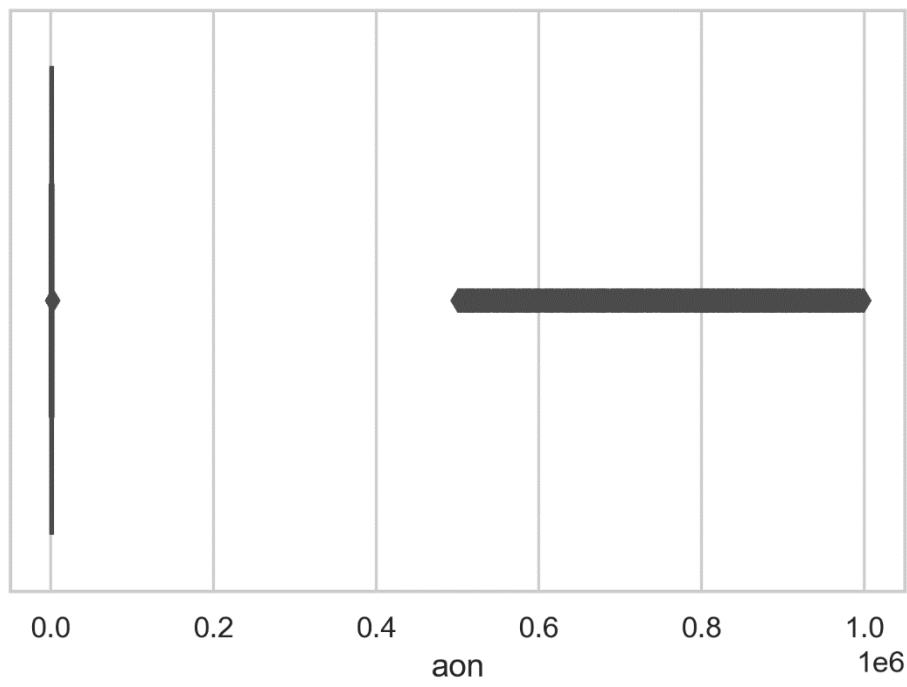
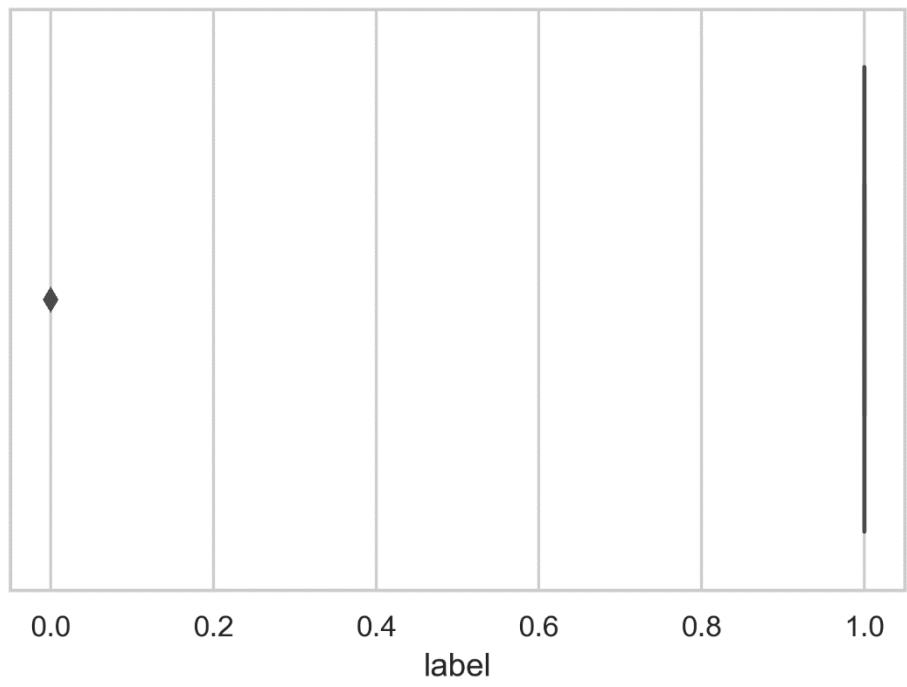


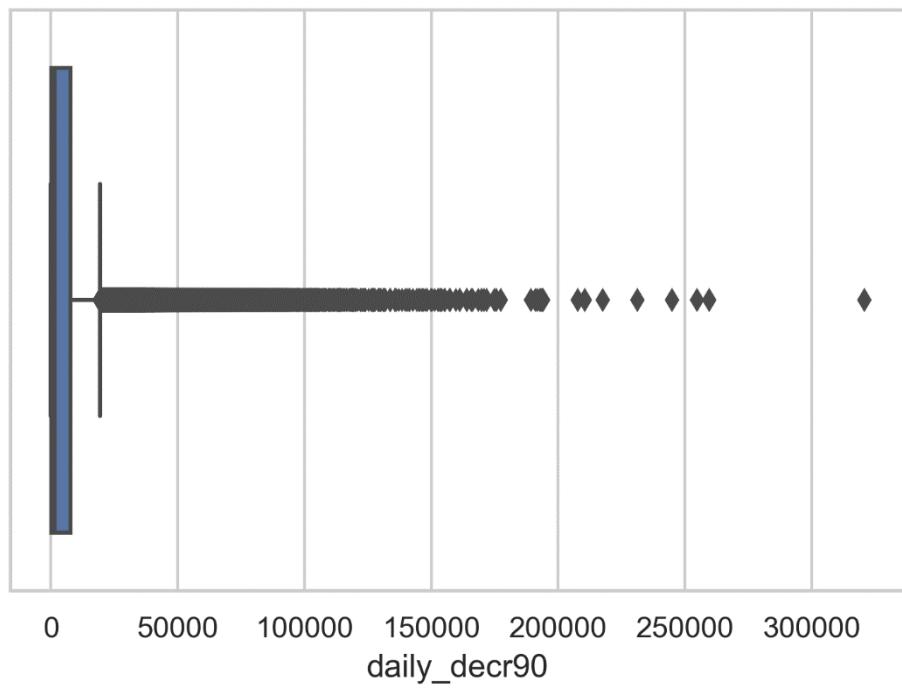
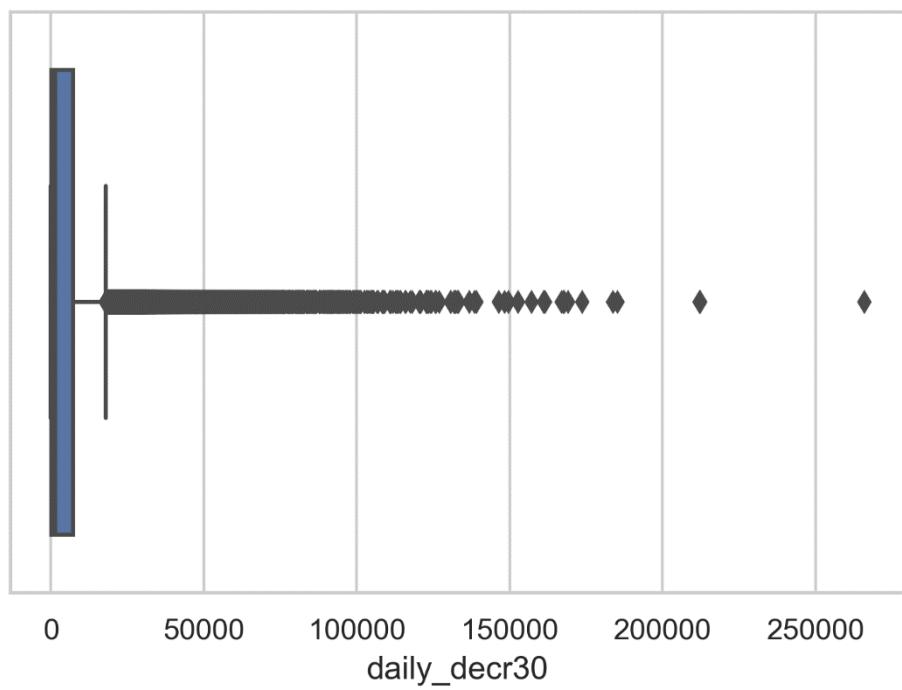


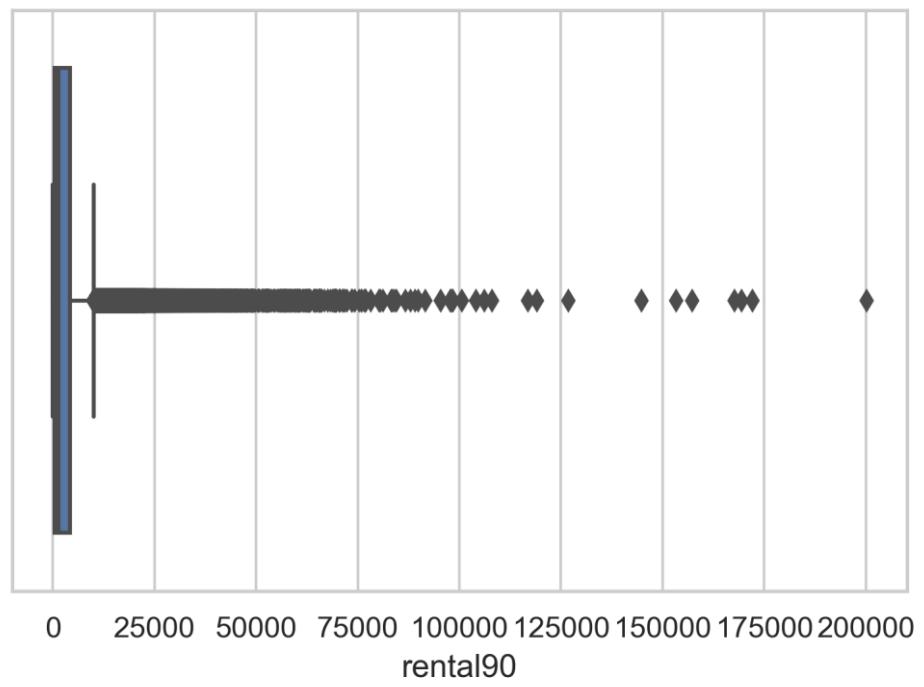
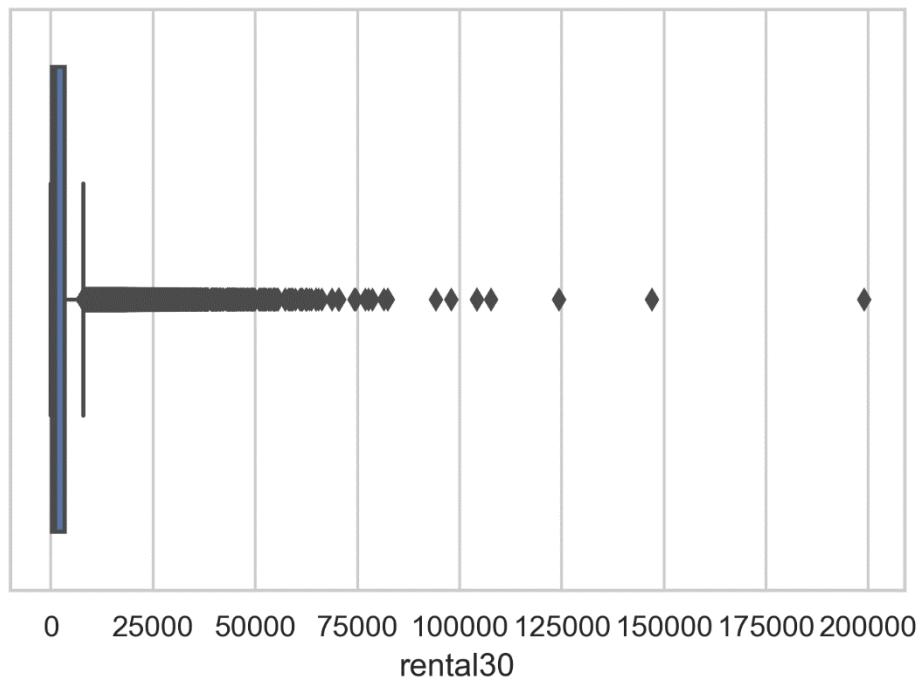


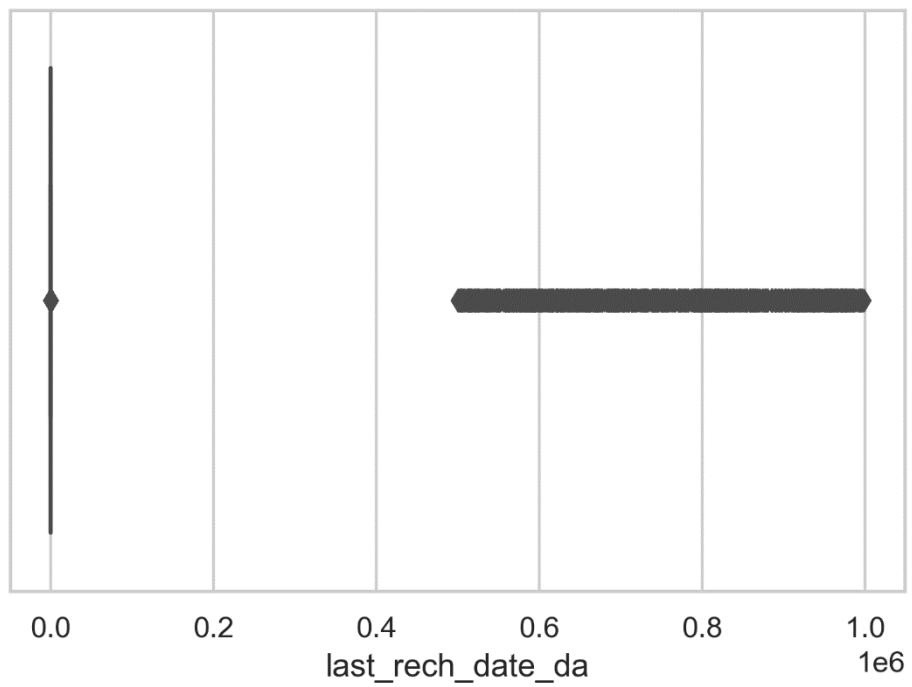
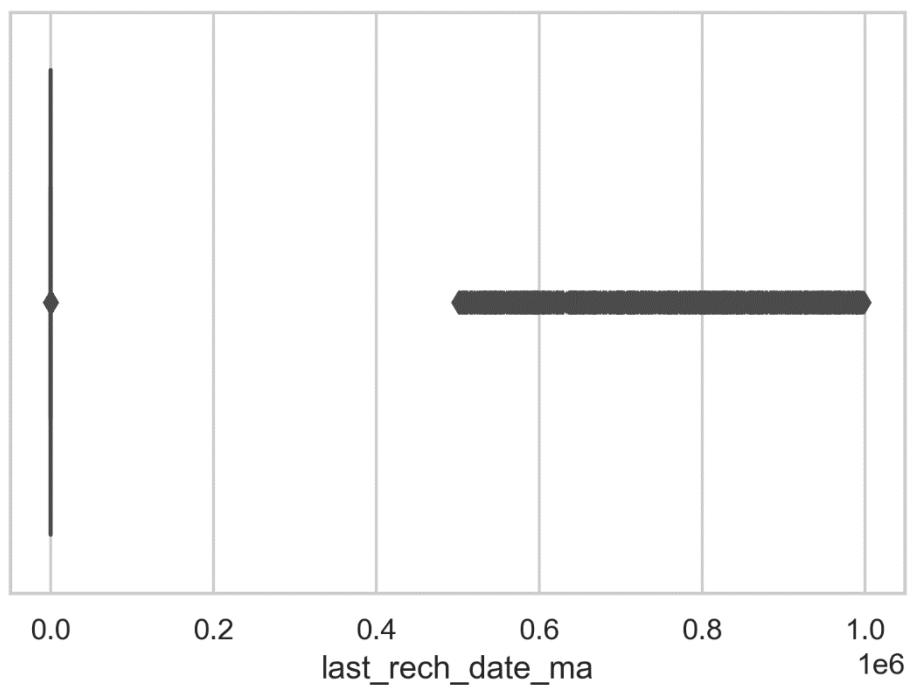


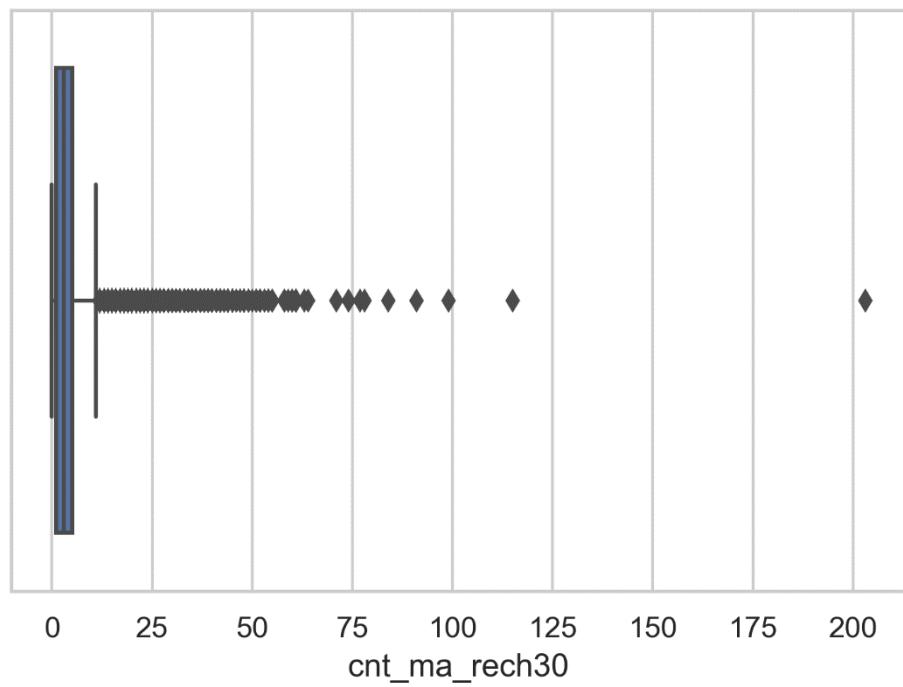
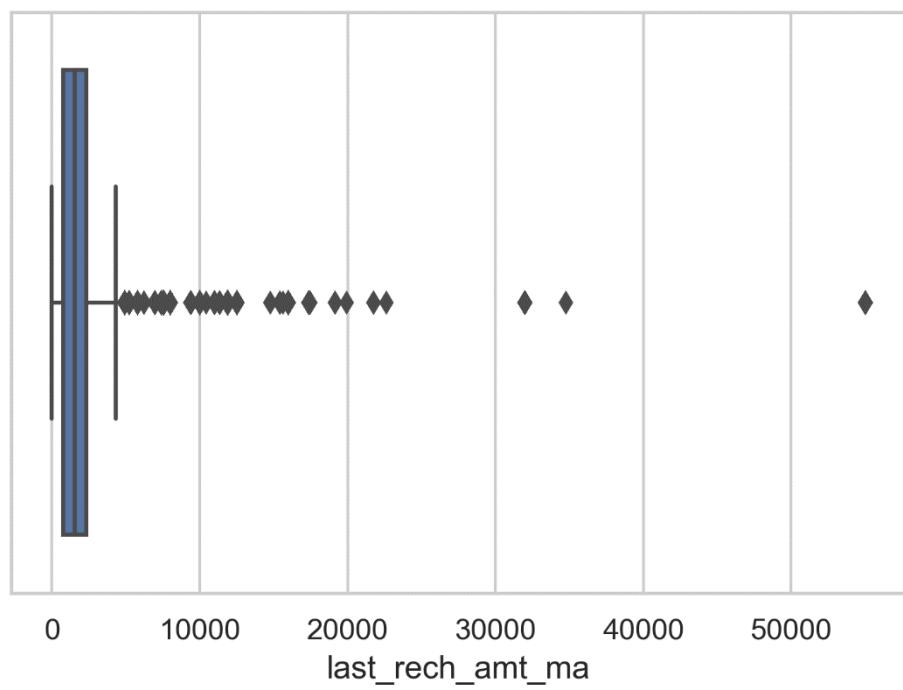


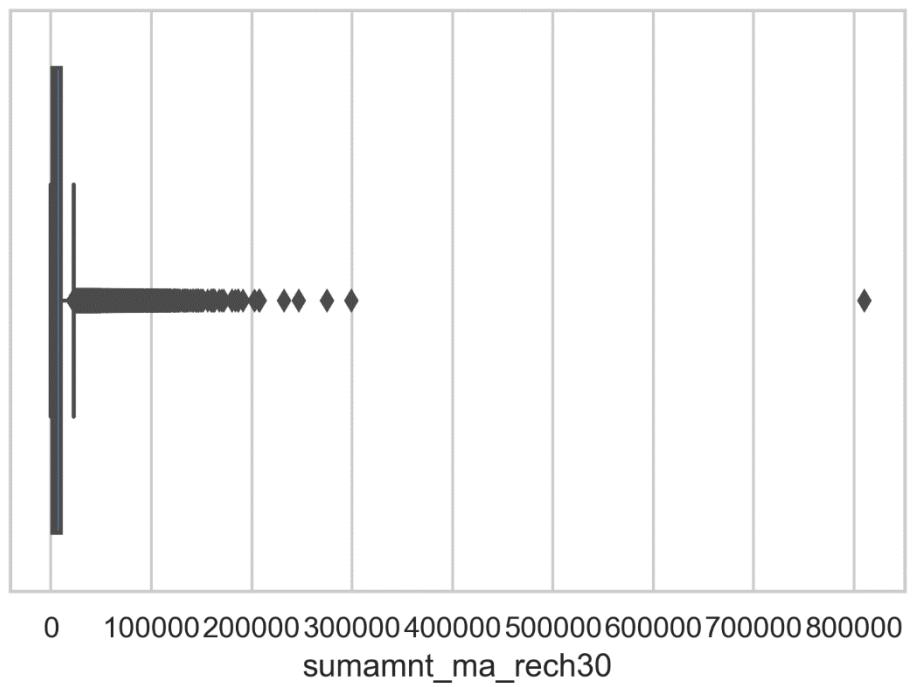
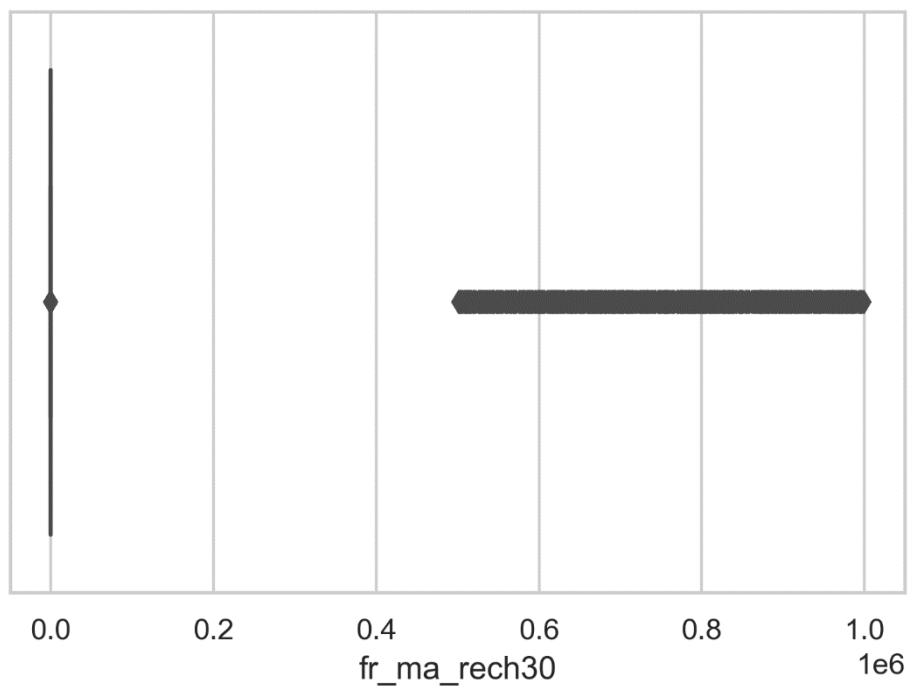


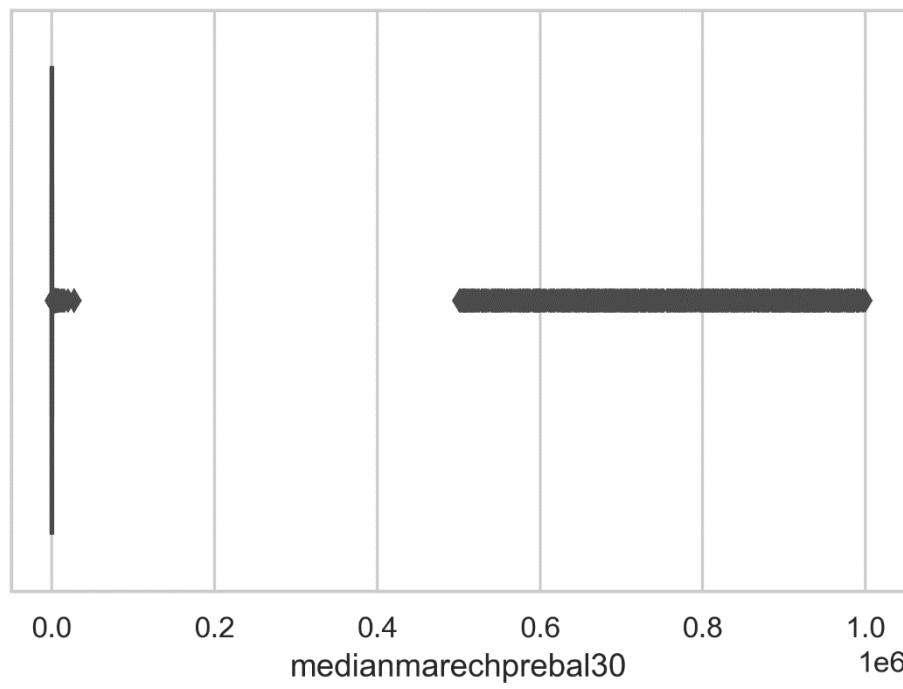
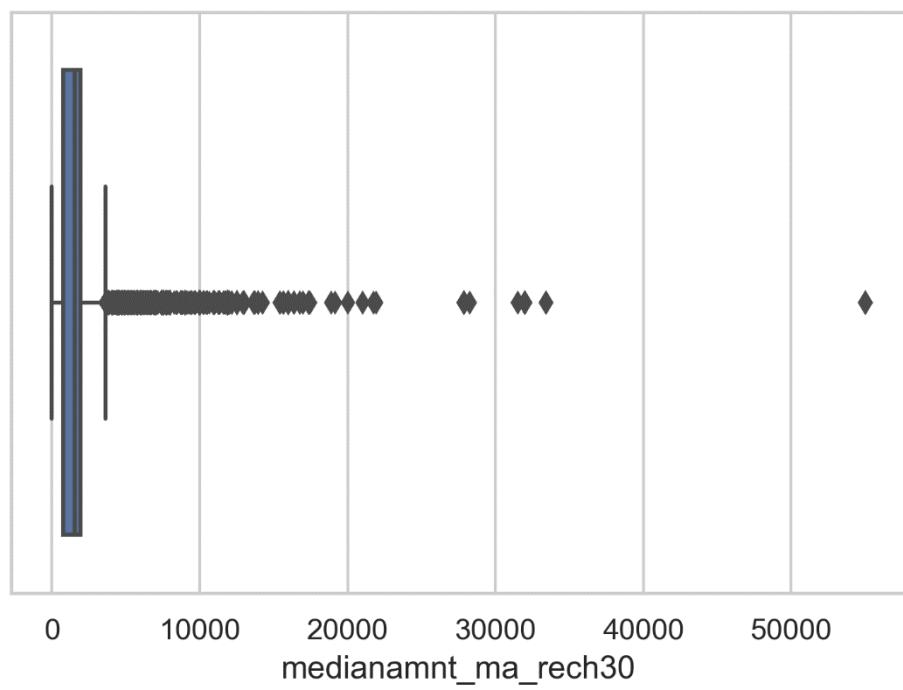


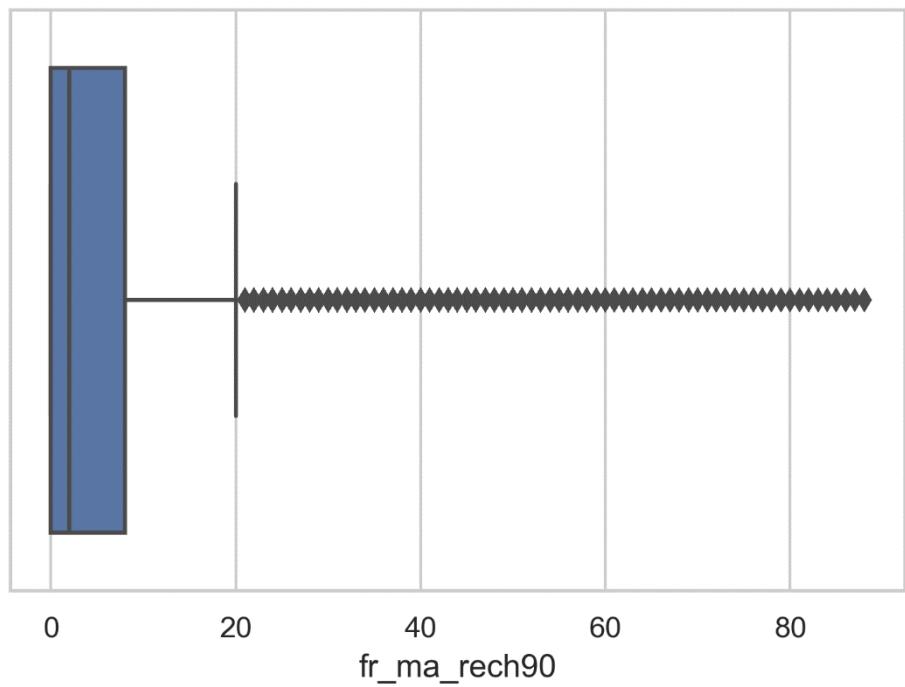
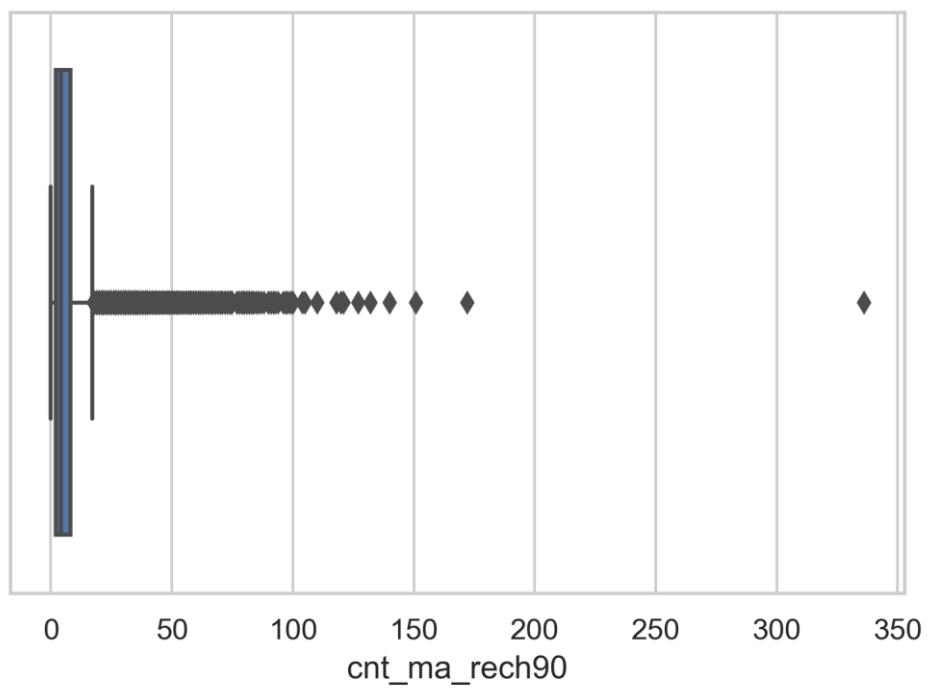


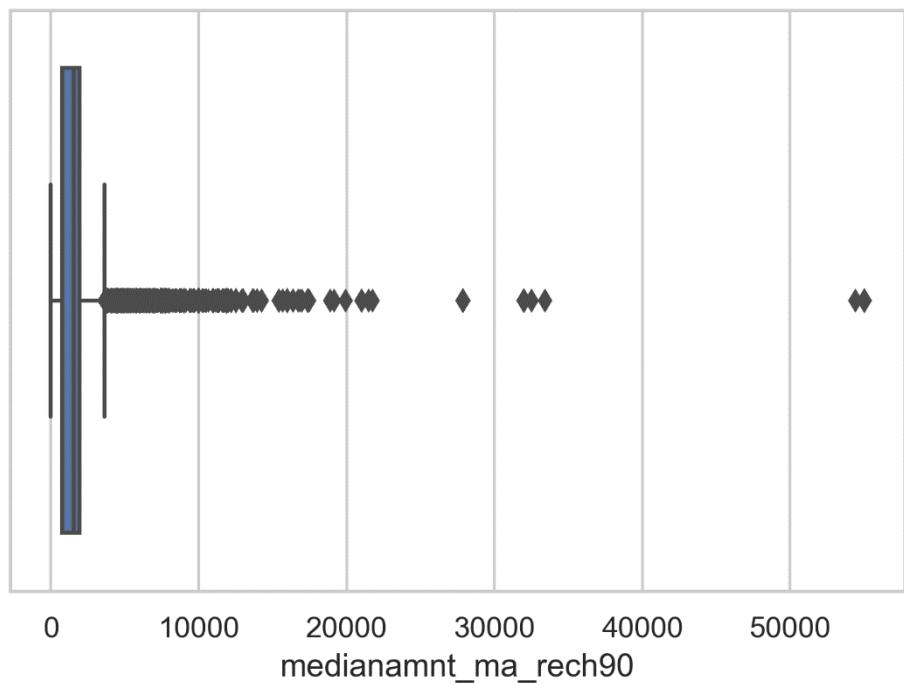
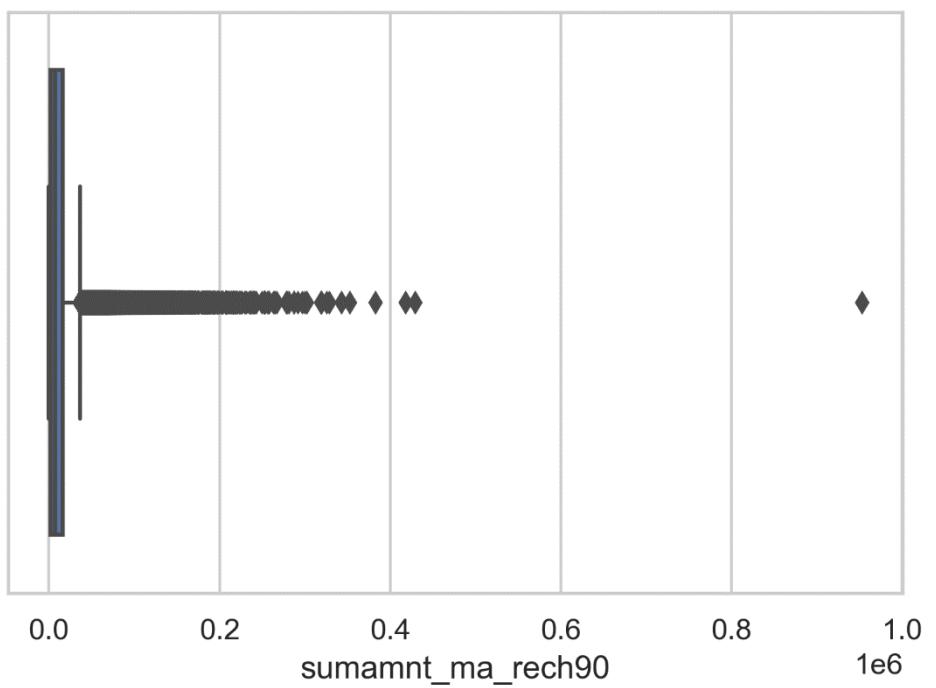


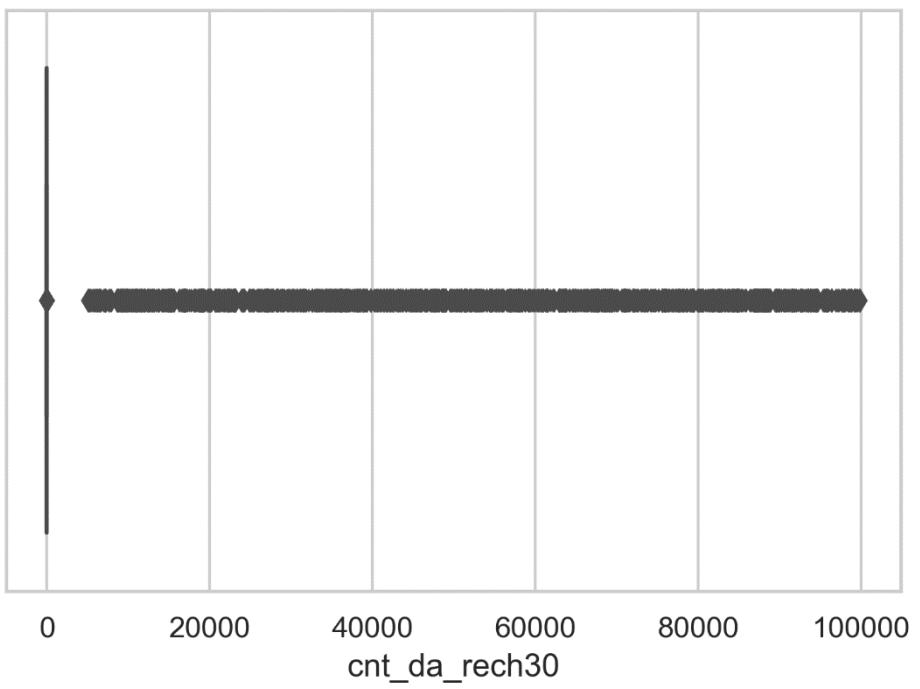
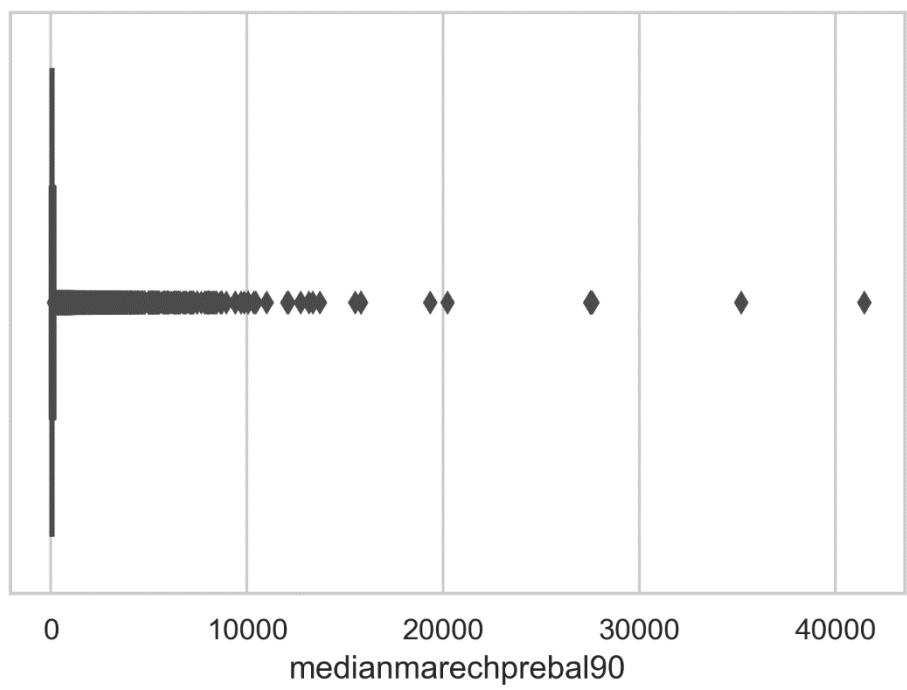


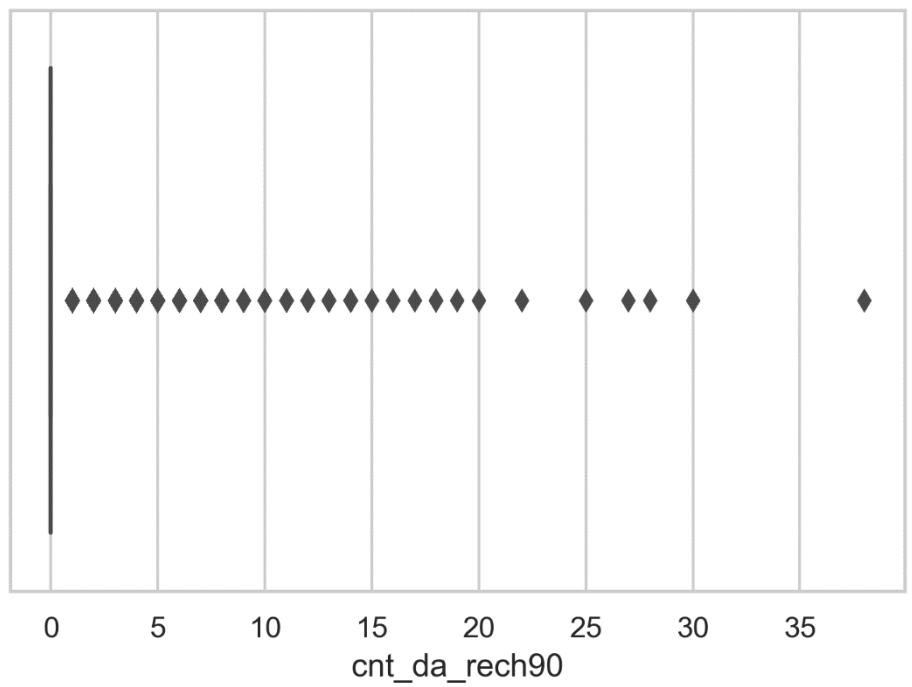
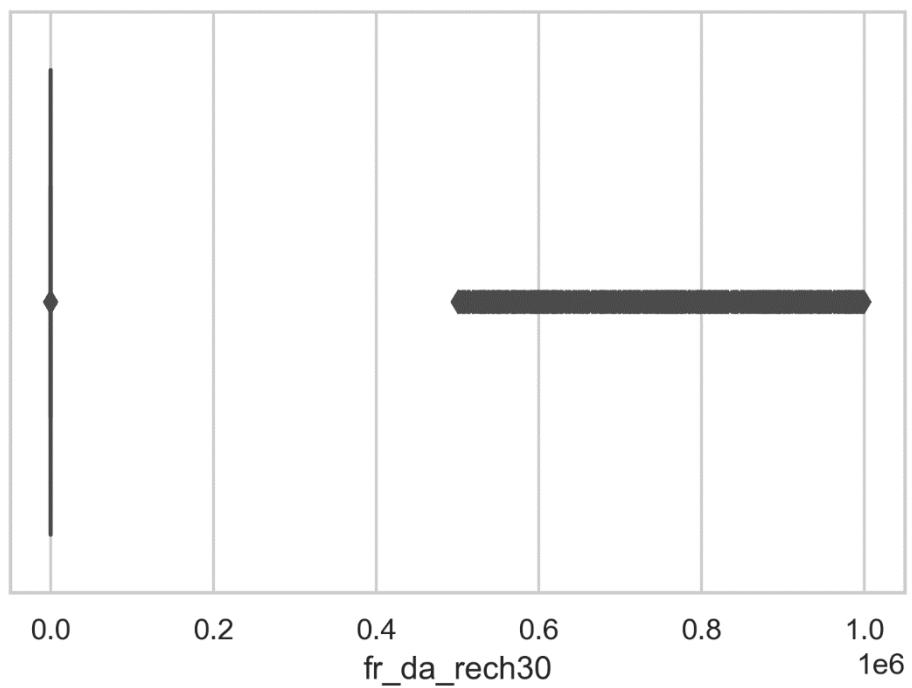


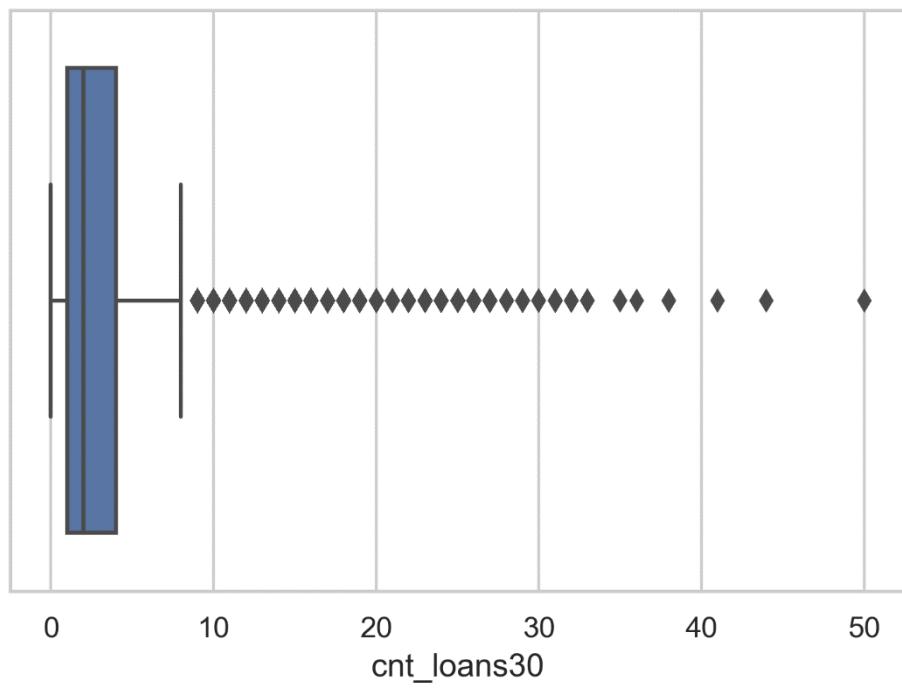
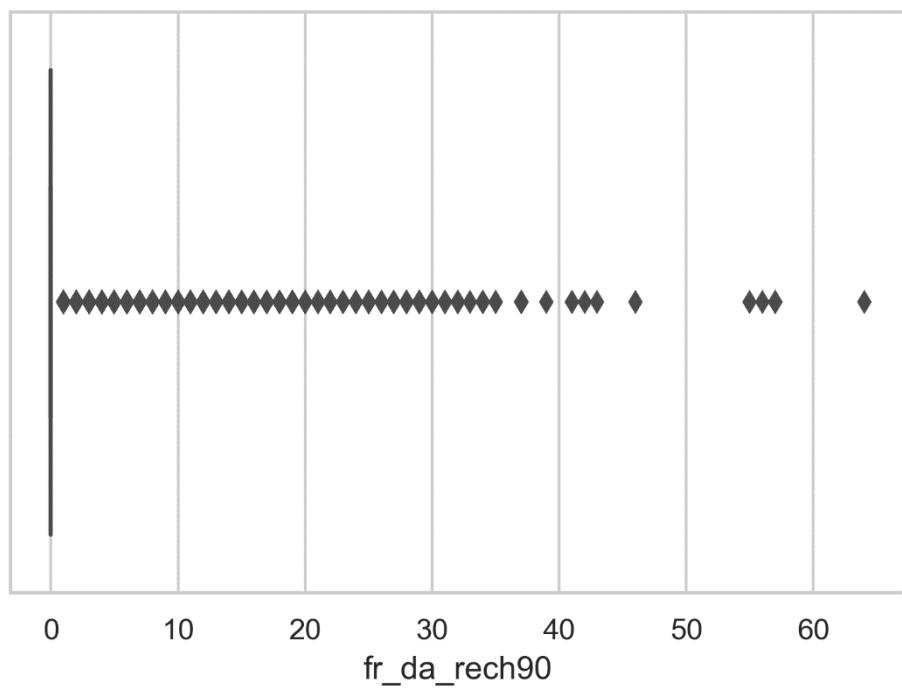


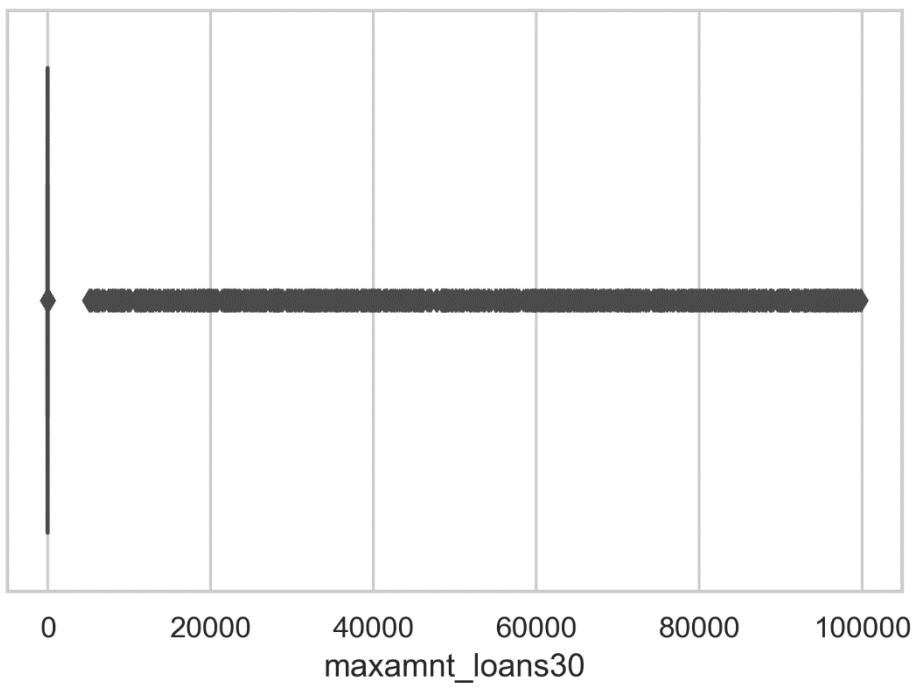
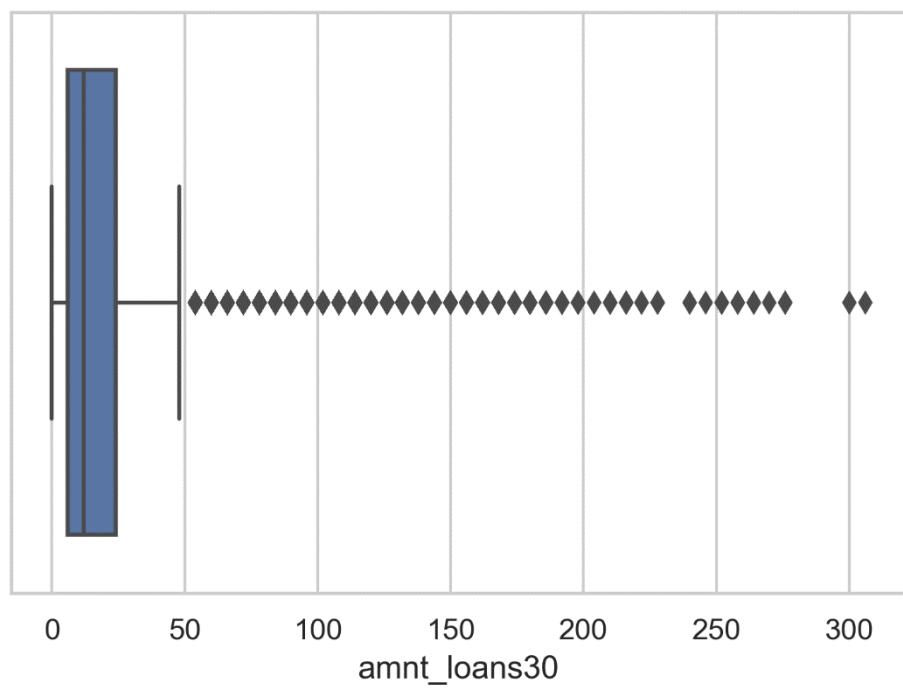


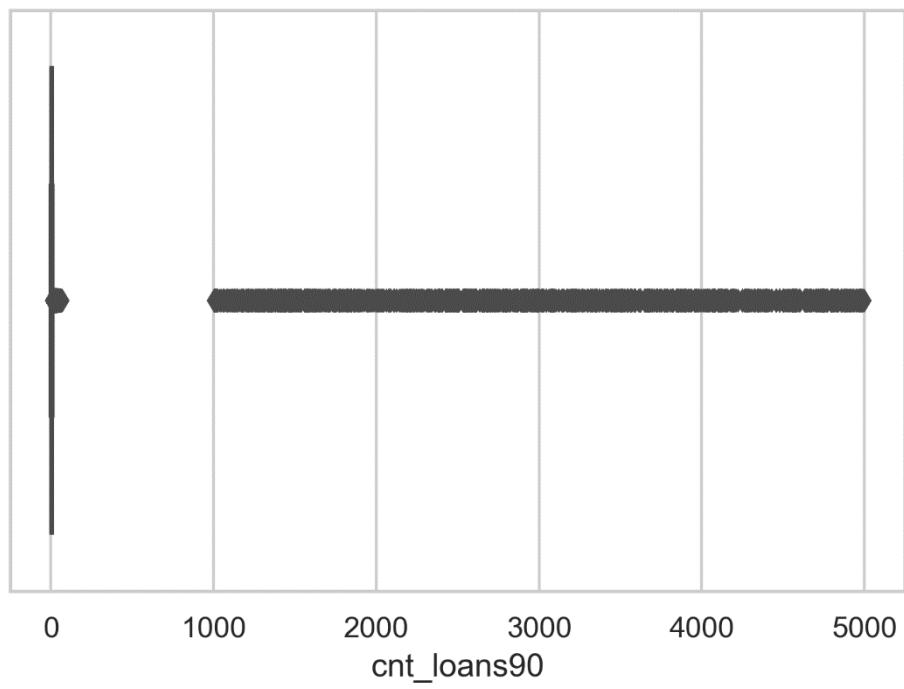
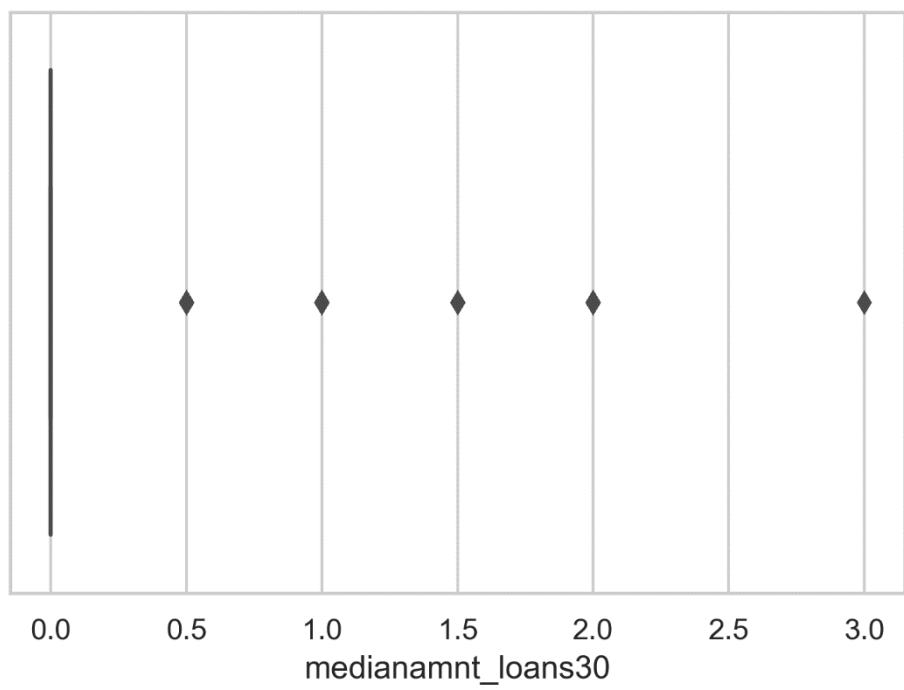


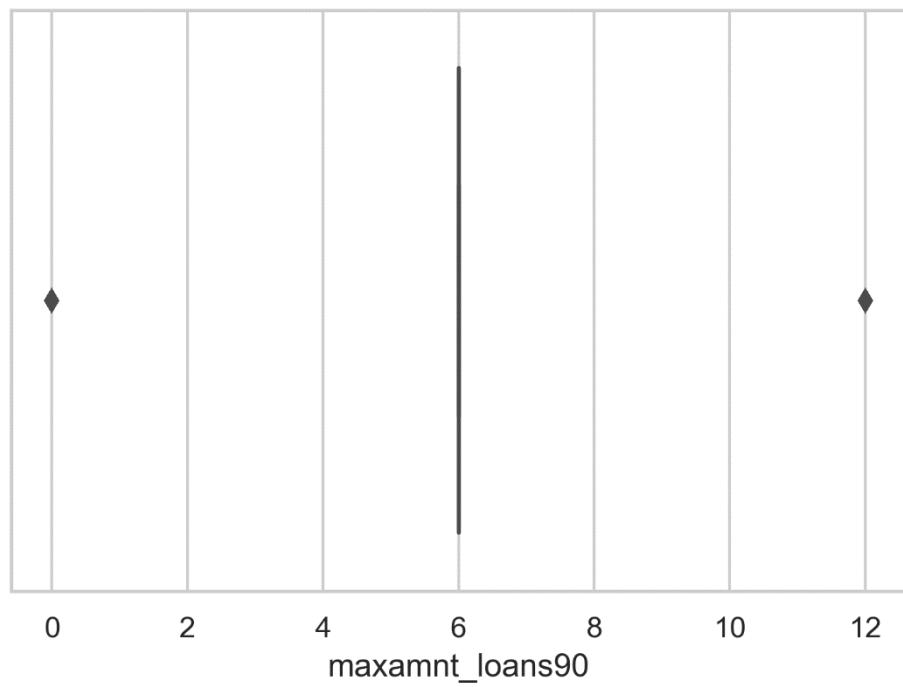
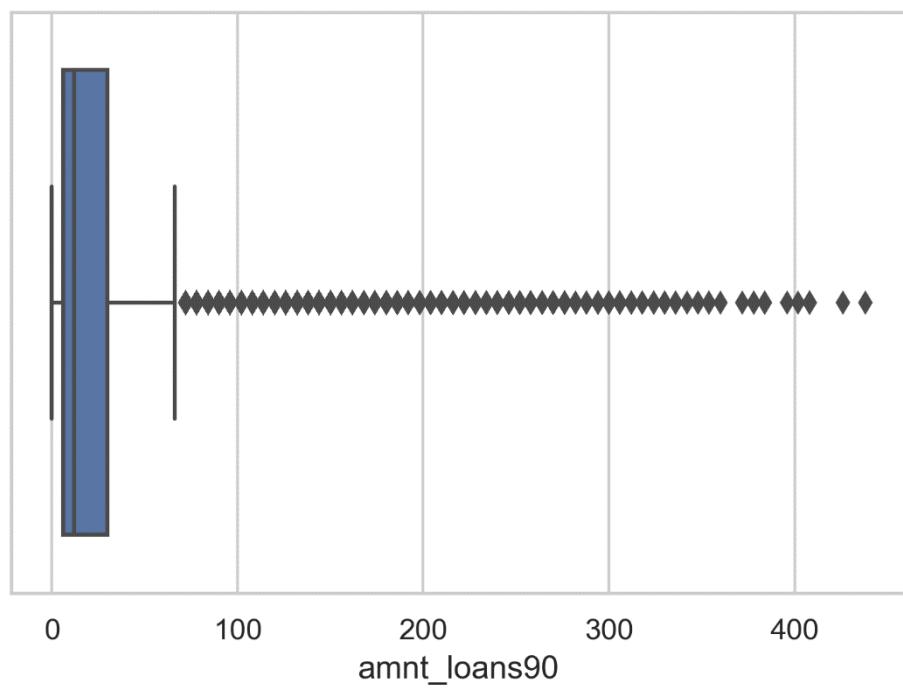


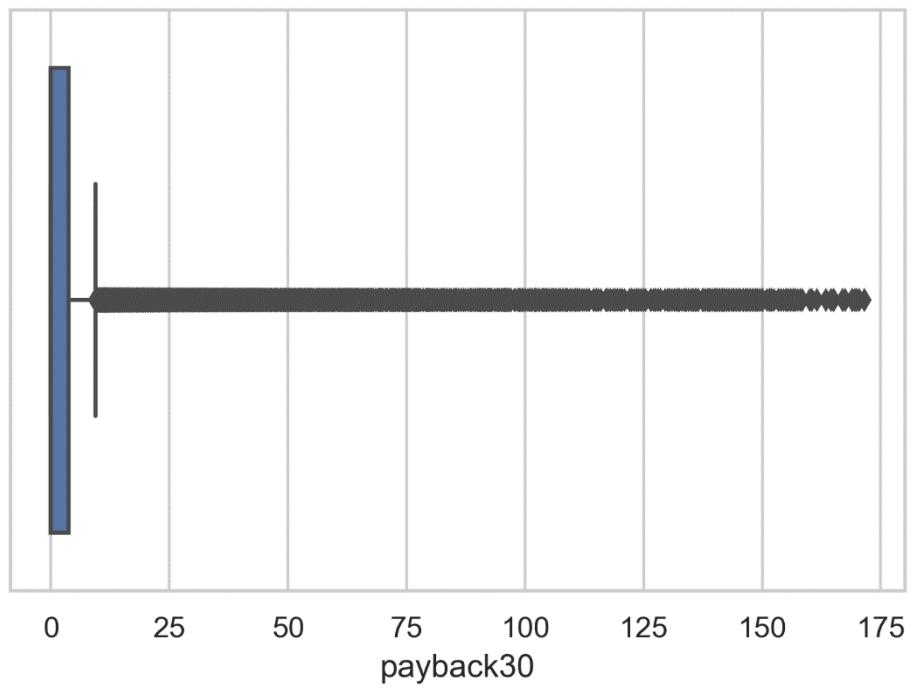
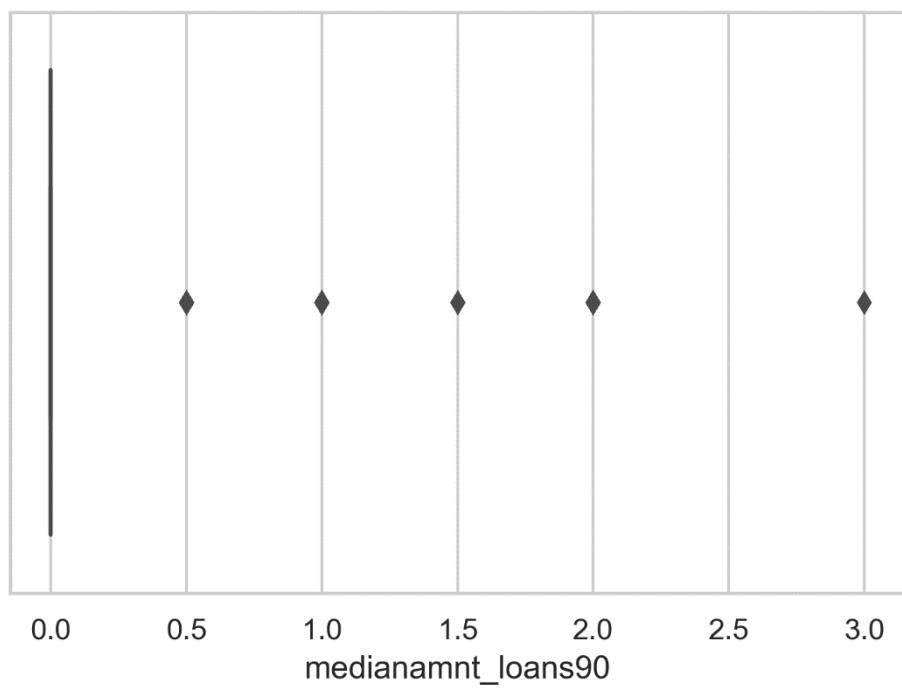


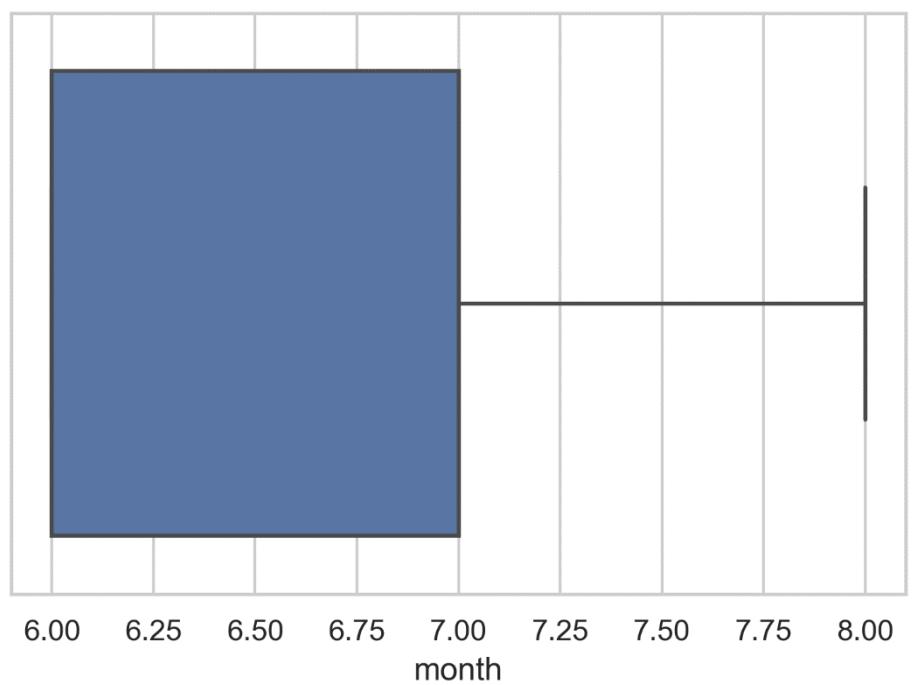
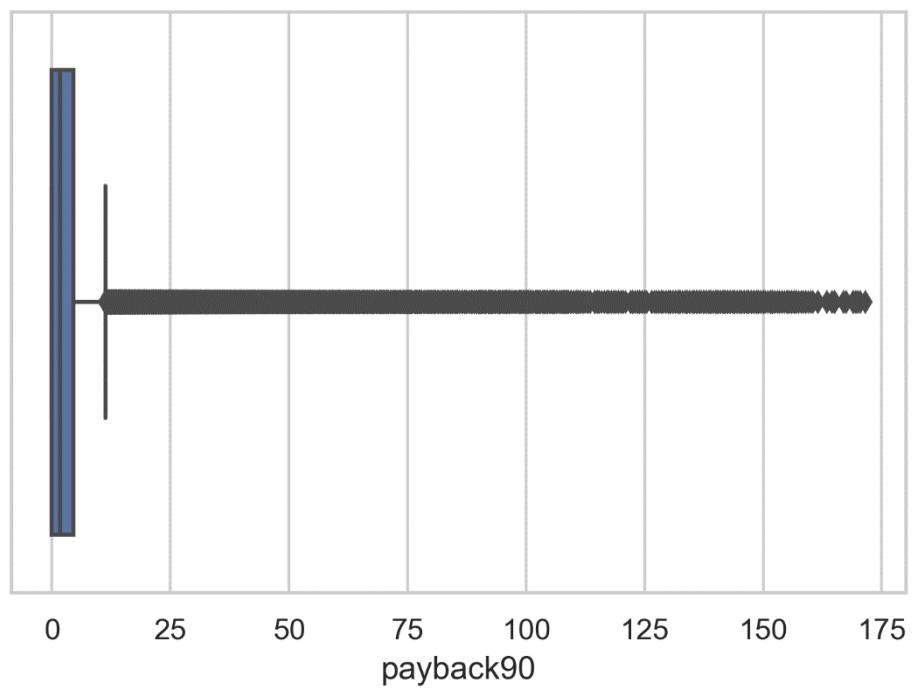


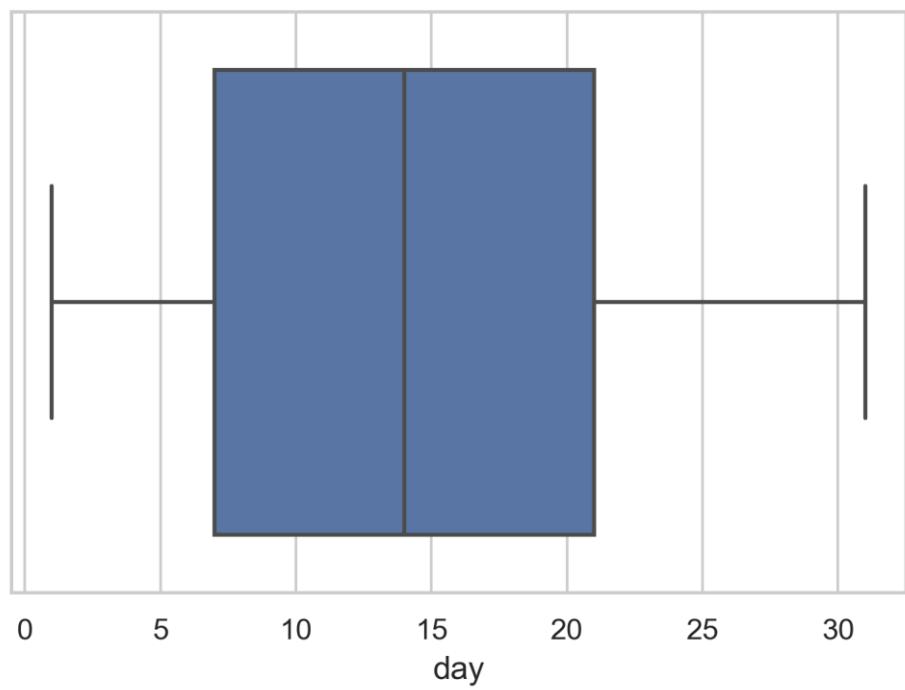


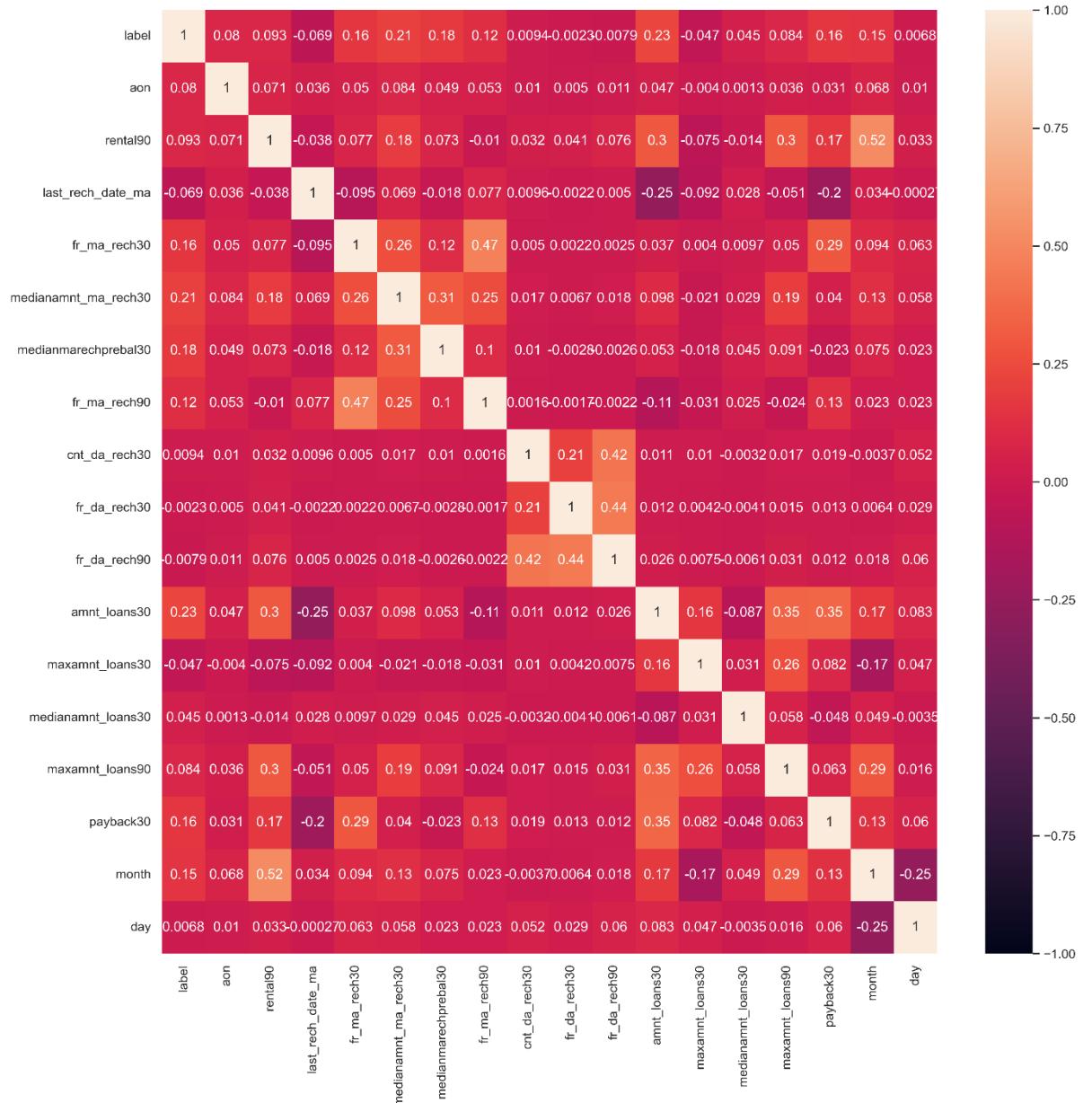


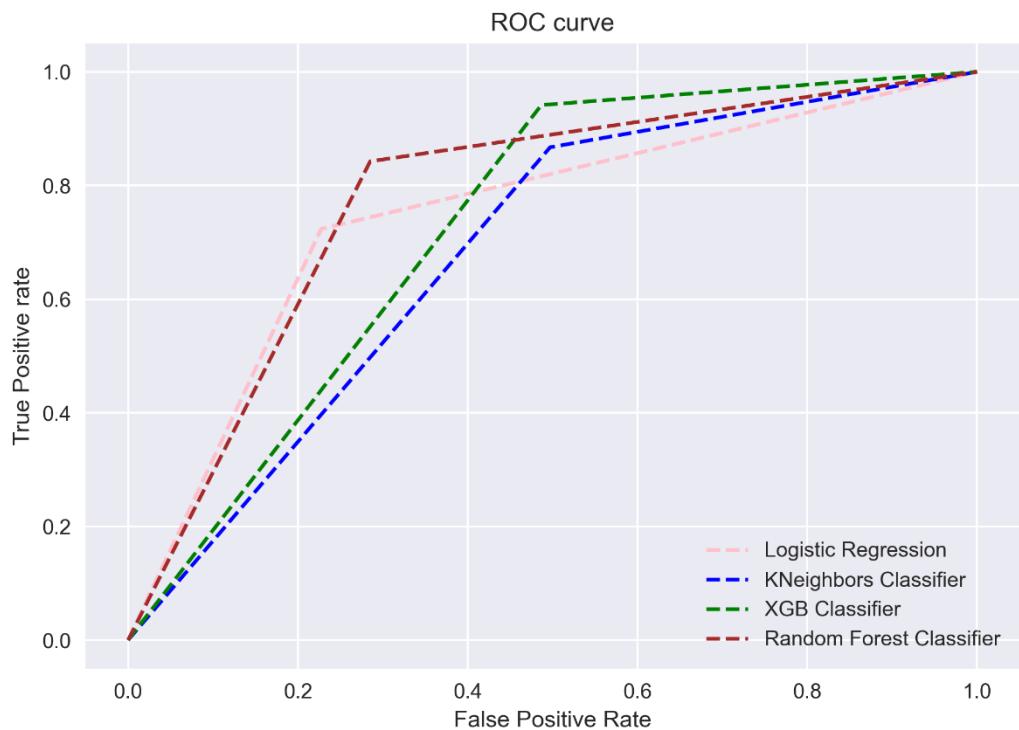
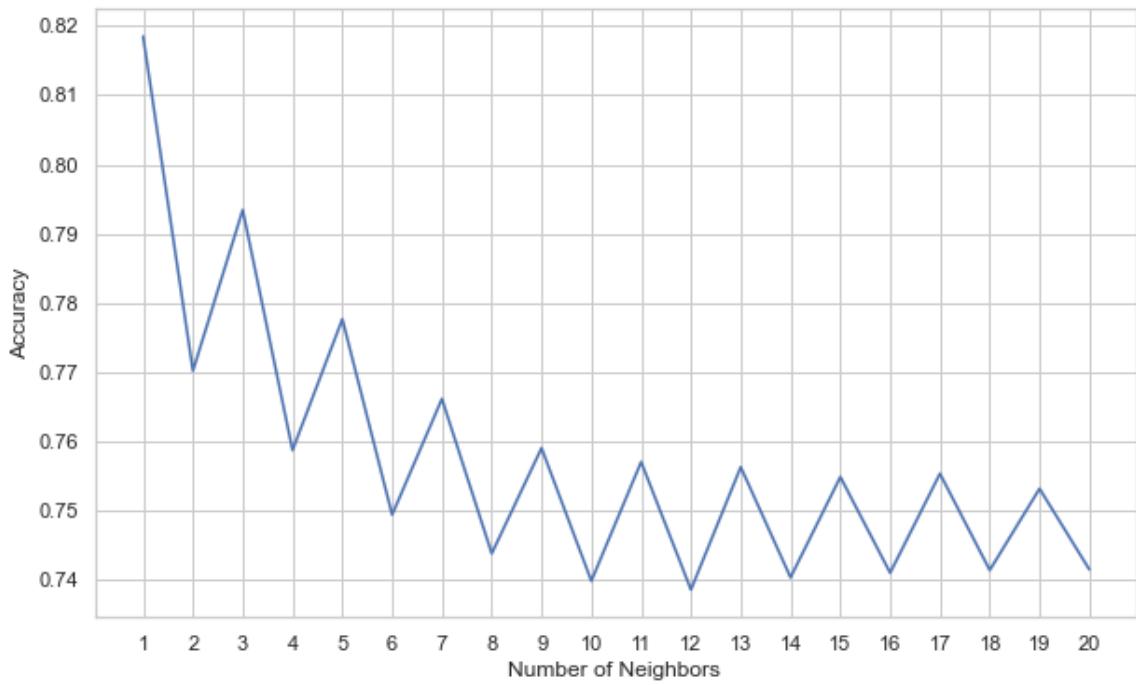












Observations:

1. age of cellular does not seem to have any effect on the label outcome as they are similarly spread for both the labels
2. subscribers with less daily amount spend are the majority ones who are failing to payback

3. subscribers with less average main balance are more unlikely pay the loan within the due date
4. Users who are recharging frequently are the ones who are failing to pay back the loan taken.
5. Users who are recharging less than 100000 Indonesian Rupiah are mostly the user who are not paying back on time.
6. Users who are taking small loans amounts are failing payback the loan amount on time.
7. Users whose median of loan amounts is less than 1.5 are mostly probably the ones who are not paying the loan on time.
8. Users with zero loan amount have 100% success rate.
9. User of august month have 100% success rate.
10. Most of the data right hand skewed.
11. Majority of the columns have outliers in them.

Results:

Out of all the Classifier models Random Forest has performed well, XGB has some over fitting compared all the other model i.e., Random Forest have good precision and recall value compared to all the other models. Also, the predicted values and true value have good fit in case of Random Forest.

Conclusion:

In the whole dataset amount of loan and median of loan amount in last 30 days column has the most influence on the label columns. Payback30 and month columns are the other factors which have good influence on the label prediction.

Outcome of the Study:

Visualizing data helped to negotiate few outliers and biased data. Data Cleaning helps in minimizing the overfitting created during model training and improves the model performance. Random Forest can neglect outliers even when the data is fed with outliers to the machine learning model. XGB Classifier is not worth the use in this type of problem as it required ample amount of time. Simplifying the data was the most challenge part in this project but it overcome if minimal domain knowledge is exploited into this project.

Limitations of this work and Scope for Future Work:

The label predicted are only limited to particular cellular telecom industry. When we try to predict for different cellular network of another country the machine learning model will fail. To improve the model efficiency, we can add more of failure label in order to balance the data and get much improved precision, recall and f1-score for the given model.