

Classification goals

The task is a binary classification problem. The goal is to predict Bad clients (target=1) based on application information provided.

Metrics

We need to choose a metric that will show the rate of right predictions. Common metrics for Binary classification tasks are ROC AUC and Gini.

For the sake of the reducing losses (risk minimization), between two algorithms with close quality metrics the one that will have less False Negative rate will be chosen. This rate shows the number of Bad clients that were misclassified as Good and the greater False Negative rate can lead to future losses when loan is approved.

Target

Apply mapping: Good:0, Bad:1

Result: Column with [0,1] values

Missing values

Only one column 'job type' has missing values (for 234 clients).

To fill missing value new class was added (value=0)

Binary categorical variables

Apply mapping with respect to correlation with Target

Result: Column with [0,1] values

Categorical variables (more than 2 values)

One Hot Encoding was applied to deal with Categorical variables with more than 2 classes. Less frequent classes were dropped.

Ways to study feature influence and importance

1. Comparing means grouped by target
Comparing means of the columns grouped by target shows that some features have quite close mean values for “good” and “bad” clients and that is why they may not have strong predictive power. Among such feature are: “foreign”, “telephon”, “dependents”, “resident”.
2. Categorical-categorical association
Statistics Chi squared, Cramers V, Crosstabs, confusion matrix
3. Visualization
Plots for every column grouped by target:
 - ▶ Stacked Bar Plot (category-target)
 - ▶ Histogram
 - ▶ BoxPlot

Numerical variables

There are 4 numerical variables out of 22 variables: ‘age’, ‘loan duration (m)’, ‘gross payments’ and ‘principal payments’.

'Gross payments' and 'Principal payments'

'Gross payments' and 'principal payments' have strong positive correlation (0.998) and we can drop one of this columns. As 'principal payments' has a slightly higher correlation with target we will keep this column.

There can be identified that among the clients with a large 'principal payments' there is a large percentage of "bad" target.

'Age'

The distribution of the age is skewed to the left, most of the clients are of the age between 22 and 42. There are 6 clients with the age over 70, only one of them is treated as a "bad" client. Younger people tend to have higher ration of "bad" classification.

'loan duration (m)'

The longer the credit the higher the chances to qualify client as "Bad" one.

All 4 are skewed to the right, log-transformation was applied.

Categorical variables

'History (number of loans)' demonstrates negative "correlation" with target.

For clients with 0 and 1 loans the percentage of "Bad" as clients with more loans know that it is worth to be a 'responsible' client to be able to get another loan.

New Features

1. Binarisation. Find values that separate classes.
2. Combining classes to increase the number of clients in less frequent

Model

Logistic Regression from the "sklearn" library.

Validation

To train training and validate models was applied the following procedure:

1. Dataset was splitted into 2 parts (Train/Test ratio=0.3) with stratification to the Target class
2. As minor class "BAD" share is 30% in the dataset the oversampling was applied to the Trainset to balance classes. To balance classes function SMOTE from imblearn.over_sampling library was applied.
3. For the Oversampled Trainset KFold was split into folds (N=6).
4. N models were trained on N-1 fold.
Validation: Predict on a left-off folds and calculate quality metrics. Mean of quality metrics was calculated.
Validation (out of sample): Predict on Testset. Quality metrics calculated. Predictions for class and probability was stored
5. Aggregated prediction:
 - a. Prediction of class – majority rule (mean \geq 0.5)
 - b. Prediction of probability – mean function.

In a nutshell, the quality of the solution is measured on 2 kinds of sets:

1. Left-off folds (balanced data- Valid1)
2. Testset (unbalanced data, 30% of the original dataset – Valid2)

Feature selection

Iterative procedure was applied to eliminate features. On each iteration, the feature that resulted in the largest increase of the quality if not applied was eliminated.

Quality was calculated as minimum from Score_1 and Score_2:

Score_1 = Mean of the Metric score on Left-off folds

Score_2 = Metric score on aggregated prediction

Results

Train/test split=0.7/0.3

Shape_0: 83 (all in)

Shape: 71 (after feature selection)

Folds 6 Seed 101 Selection: 1

Feature selection based on Gini score.

Validation scores on left-off folds (196 clients)

-----min folds-----

GINI 0.5735 Roc AUC 0.7868 FN 15 FP 22

-----mean folds-----

GINI 0.6408 Roc AUC 0.8204 FN 18.7 FP 23.3

-----std folds-----

GINI 0.0432 Roc AUC 0.0216 FN 2.9 FP 1.8

Validation scores on Testset (300 clients)

-----min folds-----

GINI 0.6232 Roc AUC 0.8116 FN 21 FP 45

-----mean folds-----

GINI 0.6316 Roc AUC 0.8158 FN 23.7 FP 48.0

-----std folds-----

GINI 0.0074 Roc AUC 0.0037 FN 1.8 FP 2.2

Aggregated prediction on Test_set (300 clients)

Confusion matrix

	0	1
0	163	47
1	23	67

Class 1 (BAD): 114 Class 0 (Good): 186

Misclassified BAD: 23 (7.667%) !!!

Misclassified GOOD: 47 (15.667%)

Accuracy 76.67%

GINI: 0.6386 Roc AUC 0.8193 FN 23 FP 47

Features

Numerical and “ordinal”

1. checking
2. purpose of loan
3. savings
4. installp
5. marital
6. co-applicant status
7. resident
8. property
9. other

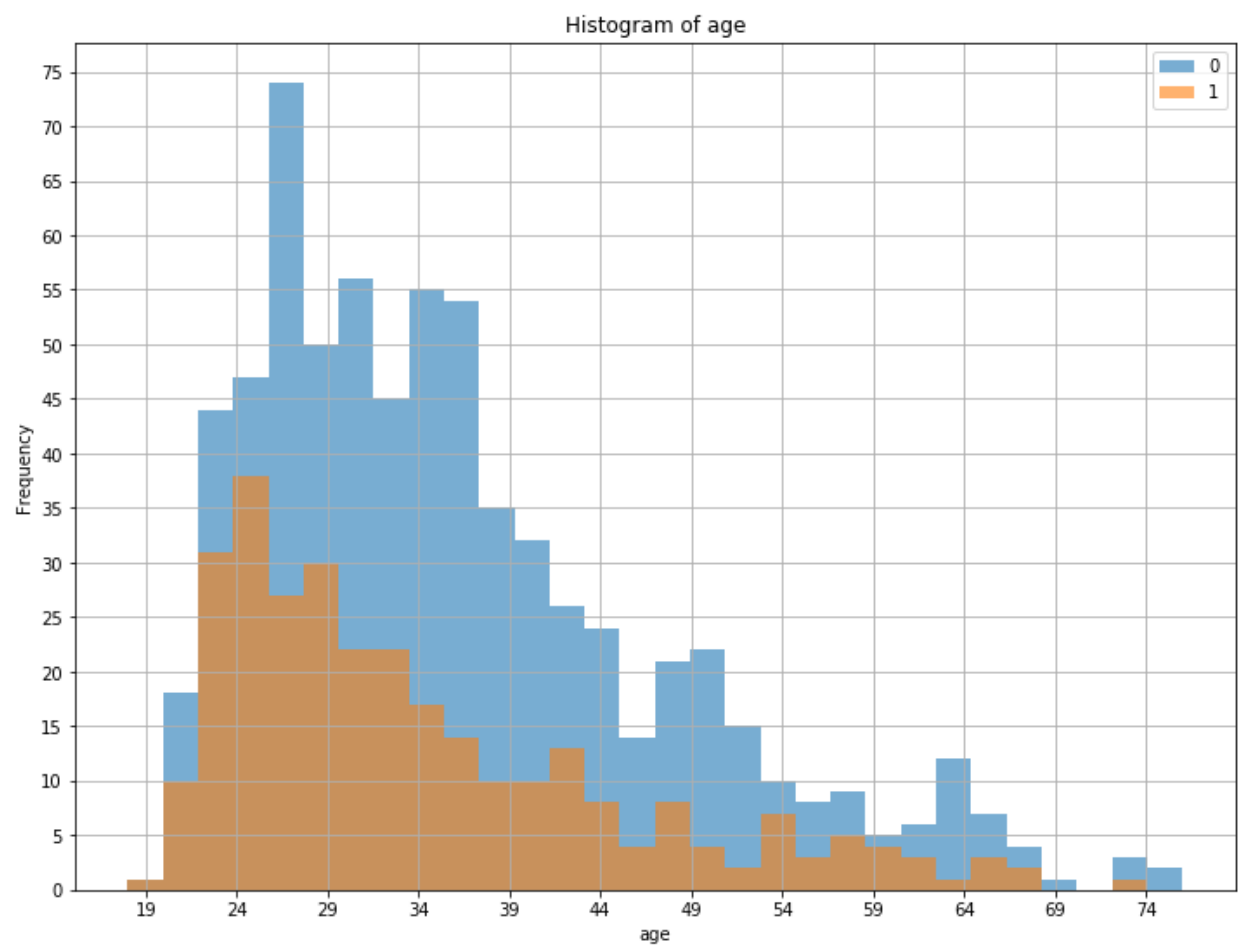
10. dependents
11. telephon
12. foreign
13. job type_m
14. other_bin
15. job_bin
16. housing_bin
17. job type_mod
18. marital_mod
19. age_1
20. history (number of loans)
21. principal payments_1
22. loan duration (m)_1
23. checking_bin
24. marital_bin
25. history (number of loans)_mod
26. savings_bin
27. principal payments_bin
28. loan duration (m)_bin
29. property_bin

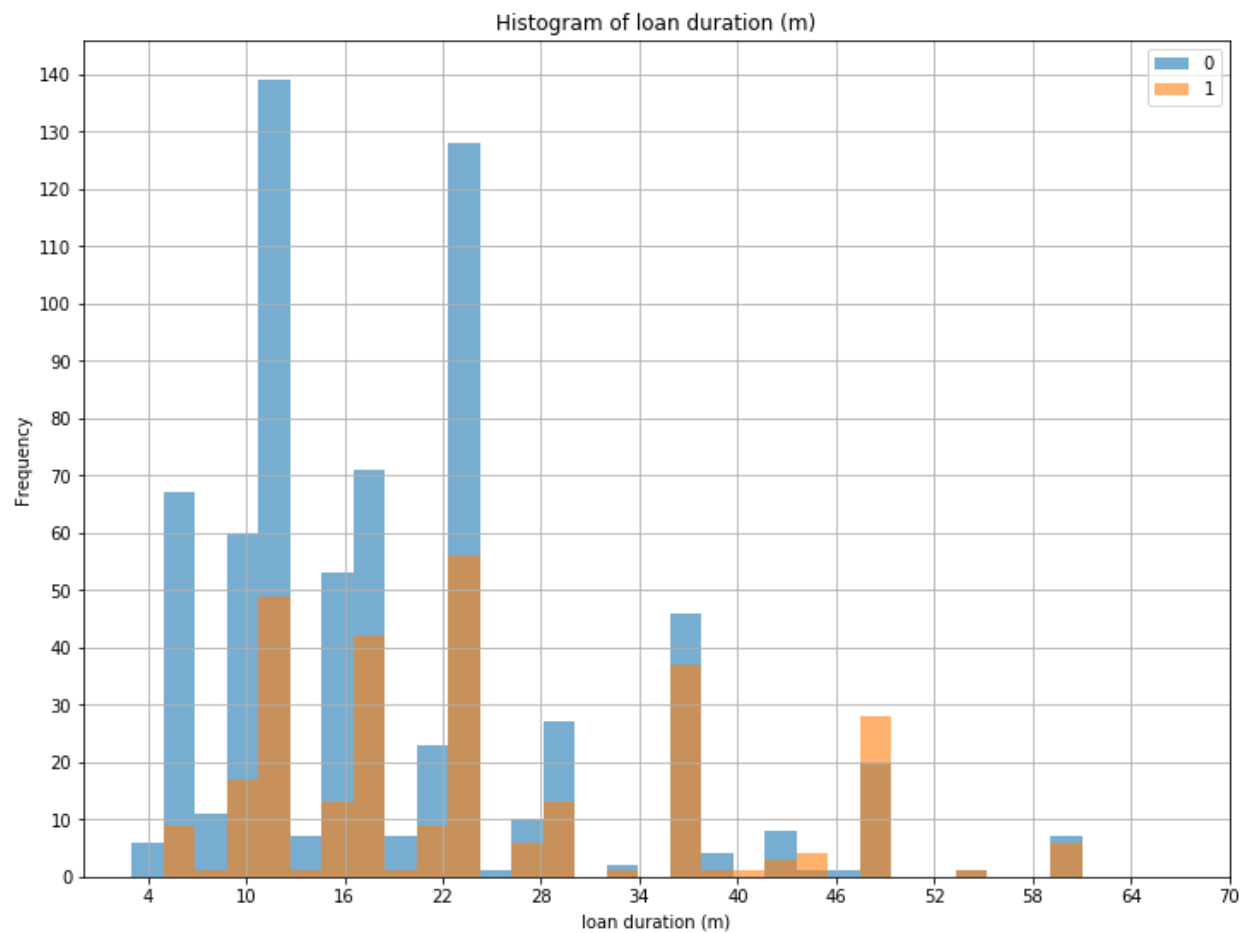
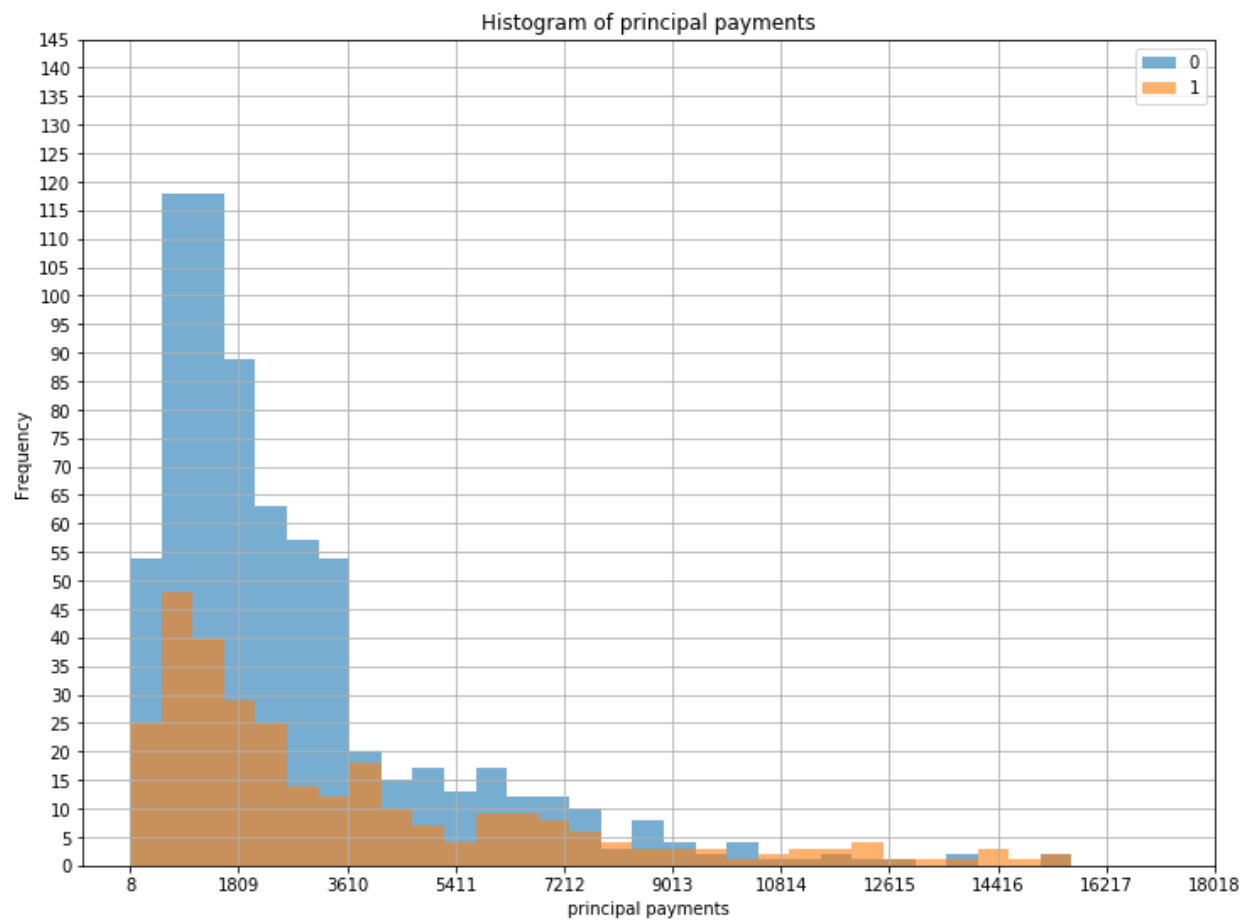
Categorical after One Hot Encoding

1. OHE_checking_4
2. OHE_checking_1
3. OHE_history (number of loans)_2
4. OHE_history (number of loans)_4
5. OHE_history (number of loans)_3
6. OHE_purpose of loan_0
7. OHE_purpose of loan_2
8. OHE_purpose of loan_1
9. OHE_purpose of loan_9
10. OHE_purpose of loan_6
11. OHE_purpose of loan_5
12. OHE_savings_1
13. OHE_savings_5
14. OHE_savings_2
15. OHE_savings_3
16. OHE_employed_5
17. OHE_employed_4
18. OHE_employed_2
19. OHE_installp_4
20. OHE_installp_2
21. OHE_marital_3
22. OHE_marital_2
23. OHE_marital_4
24. OHE_co-applicant status_1
25. OHE_co-applicant status_3
26. OHE_resident_4
27. OHE_resident_2
28. OHE_resident_3
29. OHE_property_3
30. OHE_property_1

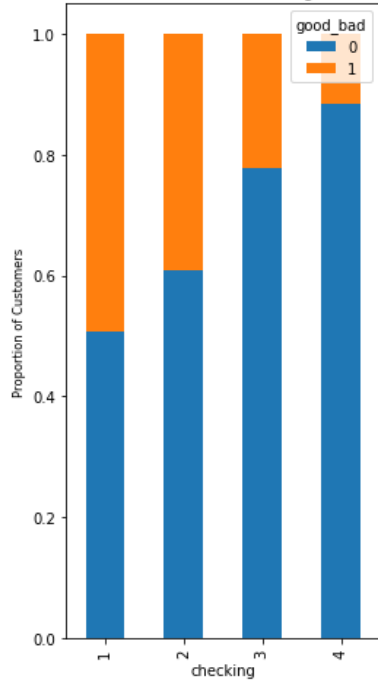
31. OHE_property_2
32. OHE_other_3
33. OHE_other_1
34. OHE_housing_1
35. OHE_exist credit bureau data_1
36. OHE_exist credit bureau data_2
37. OHE_exist credit bureau data_3
38. OHE_job_3
39. OHE_job_2
40. OHE_job_4
41. OHE_job type_3.0
42. OHE_job type_4.0

Visualisation

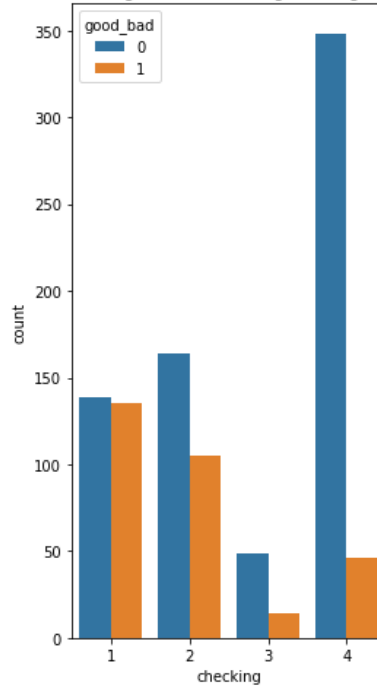




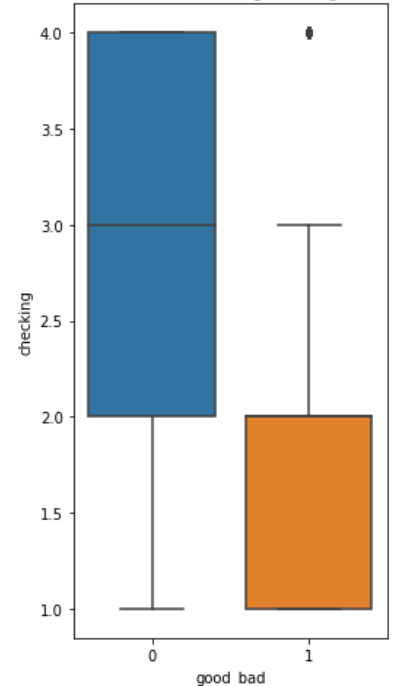
Stacked Bar Chart of "checking" vs target



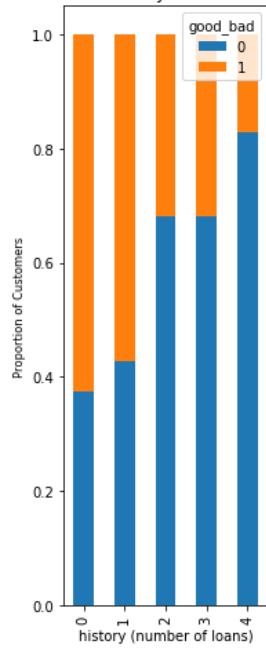
Histogram of "checking" vs target



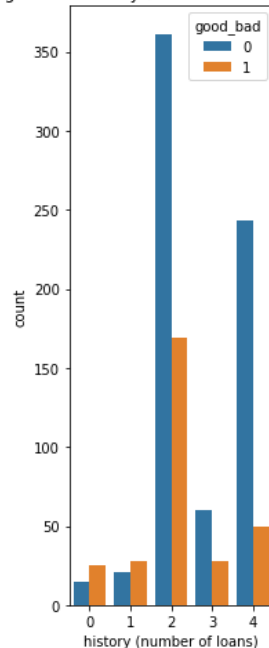
BOX for "checking" vs target



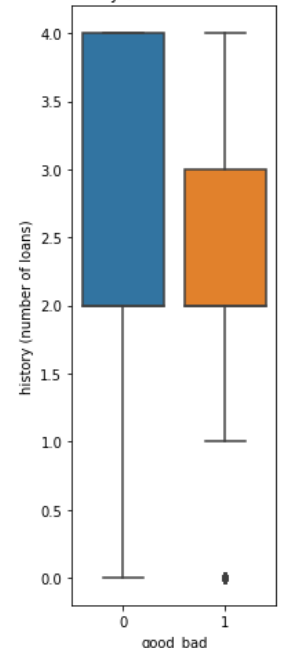
Stacked Bar Chart of "history (number of loans)" vs target



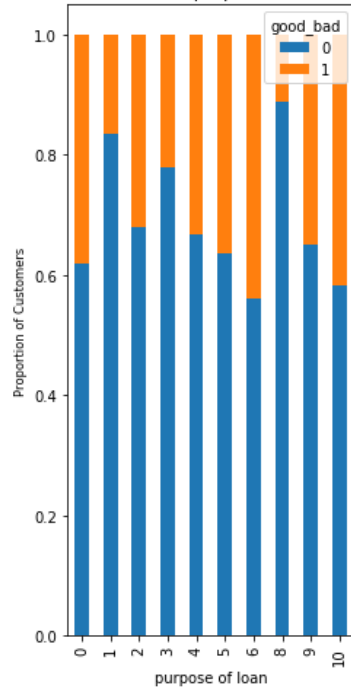
Histogram of "history (number of loans)" vs target



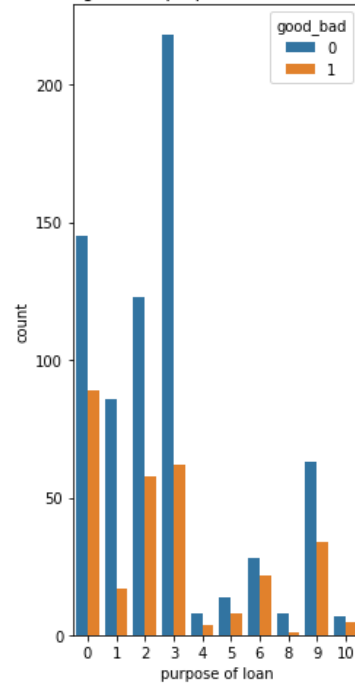
BOX for "history (number of loans)" vs target



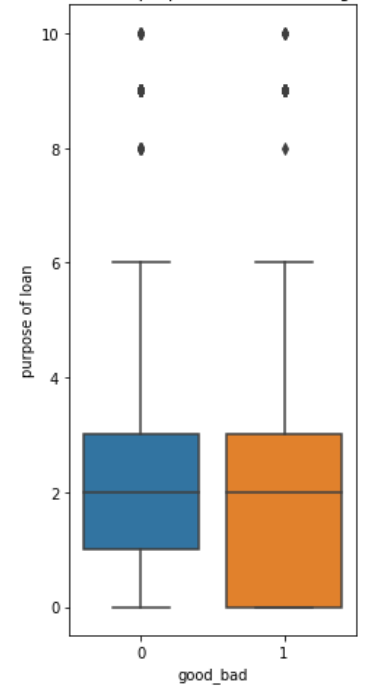
Stacked Bar Chart of "purpose of loan" vs target



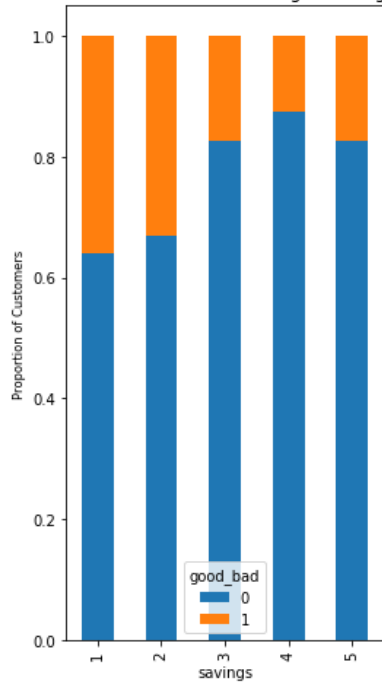
Histogram of "purpose of loan" vs target



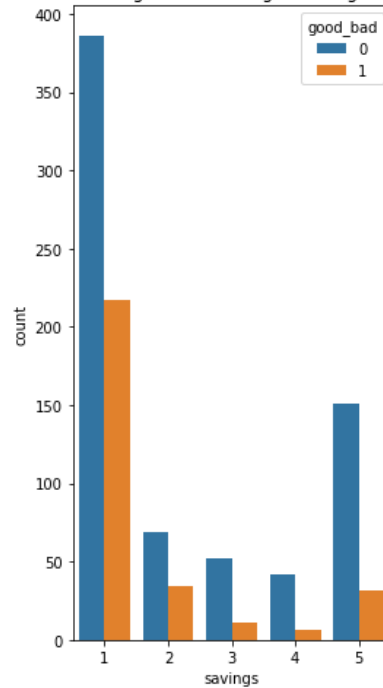
BOX for "purpose of loan" vs target



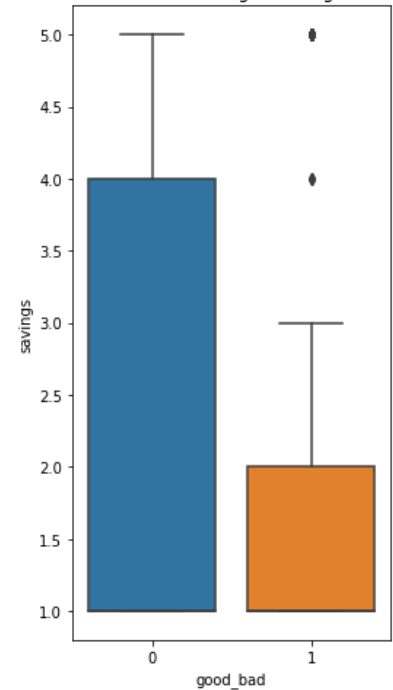
Stacked Bar Chart of "savings" vs target



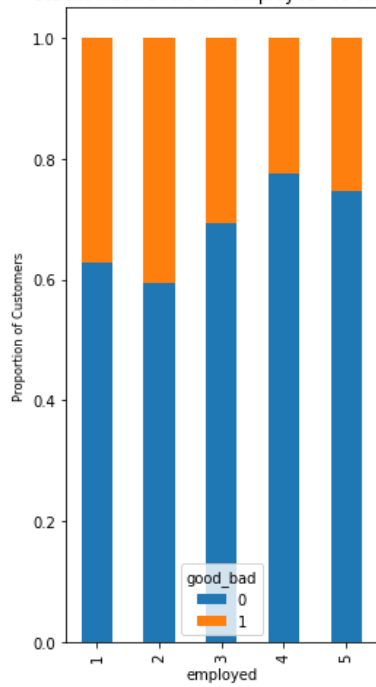
Histogram of "savings" vs target



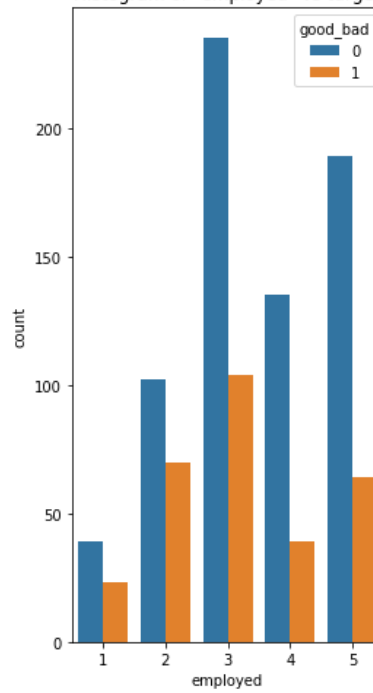
BOX for "savings" vs target



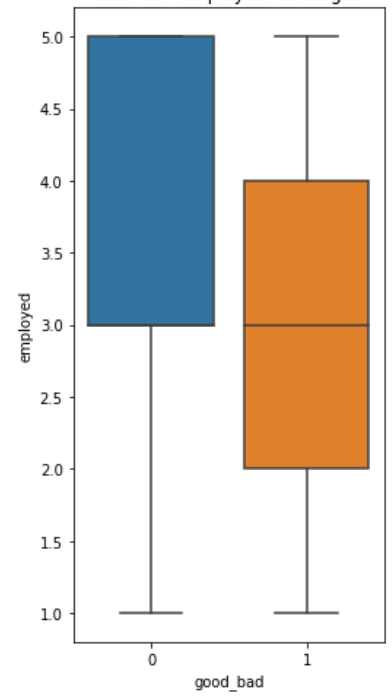
Stacked Bar Chart of "employed" vs target



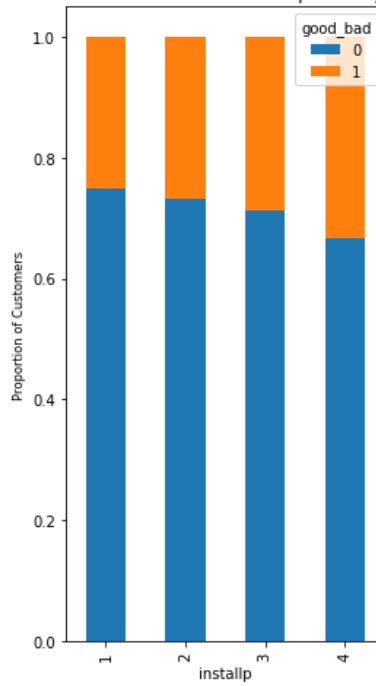
Histogram of "employed" vs target



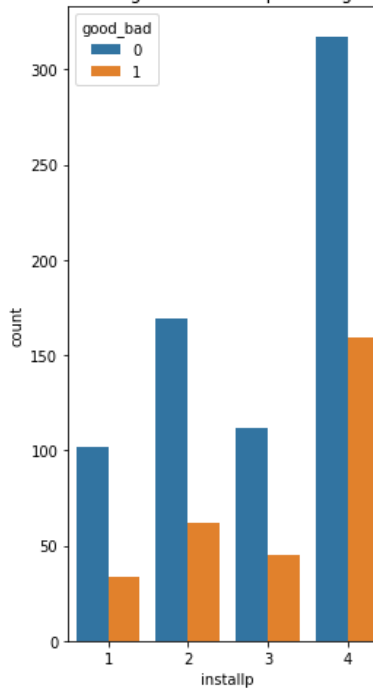
BOX for "employed" vs target



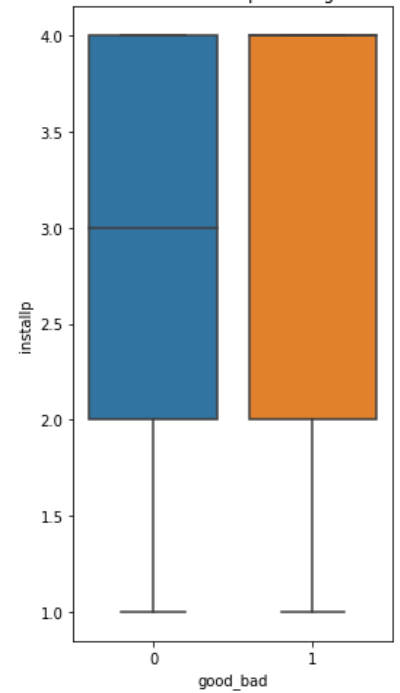
Stacked Bar Chart of "installp" vs target



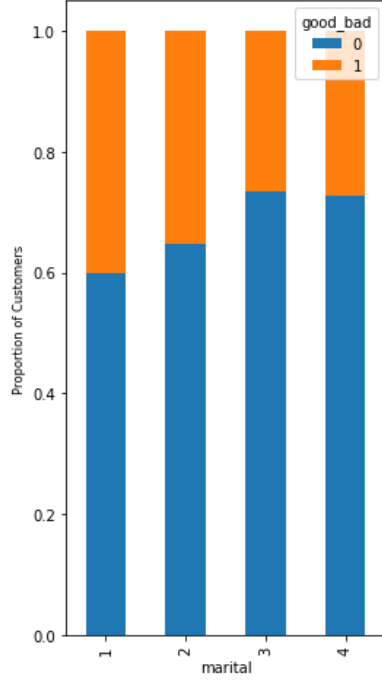
Histogram of "installp" vs target



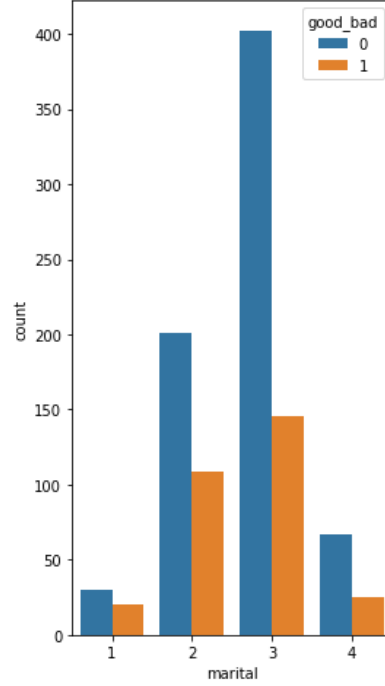
BOX for "installp" vs target



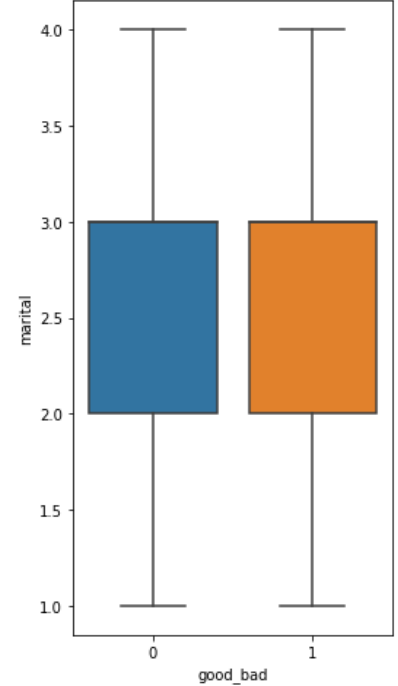
Stacked Bar Chart of "marital" vs target



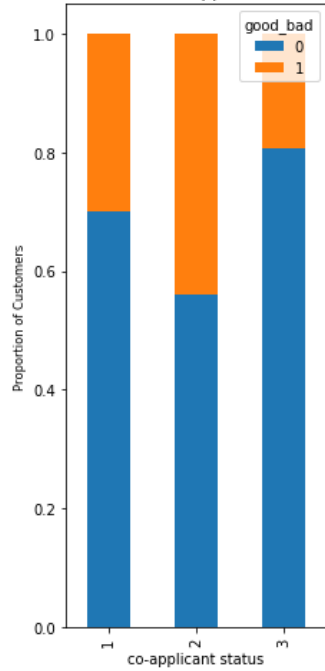
Histogram of "marital" vs target



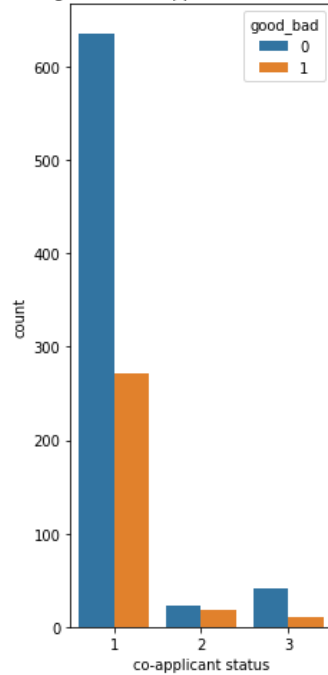
BOX for "marital" vs target



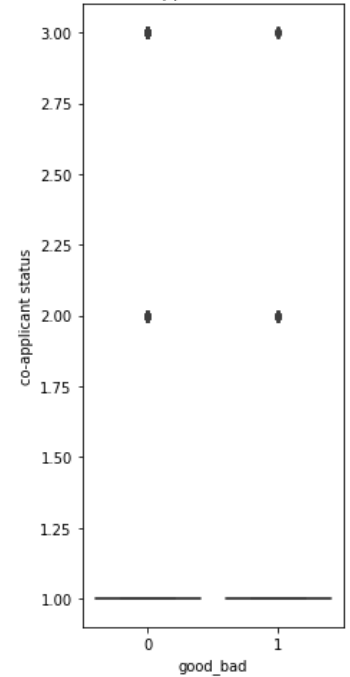
Stacked Bar Chart of "co-applicant status" vs target

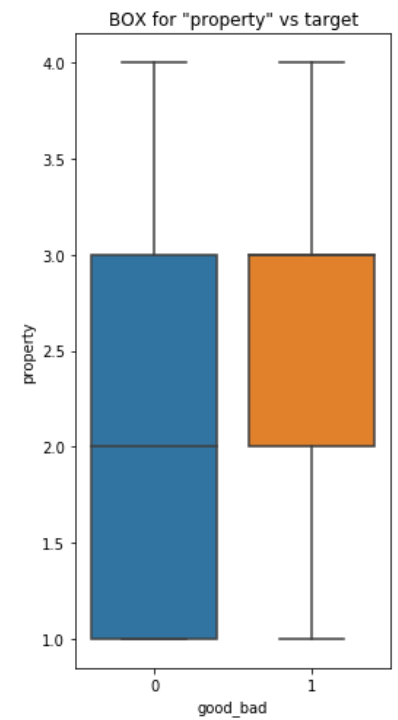
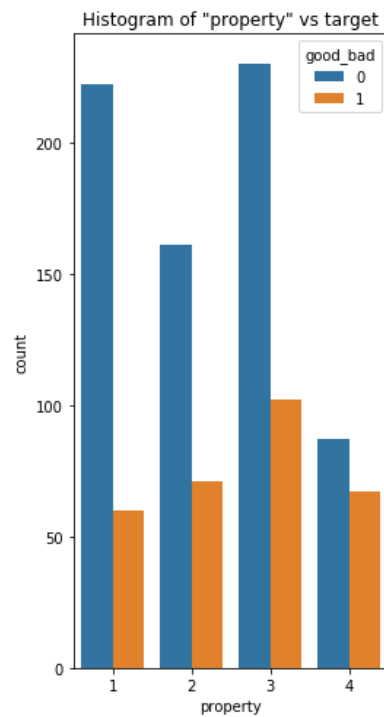
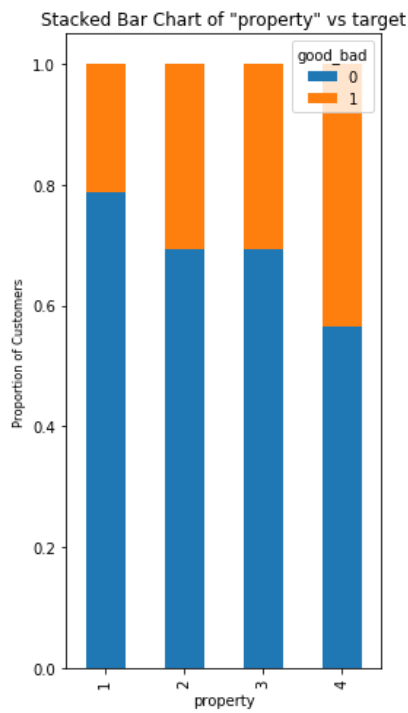
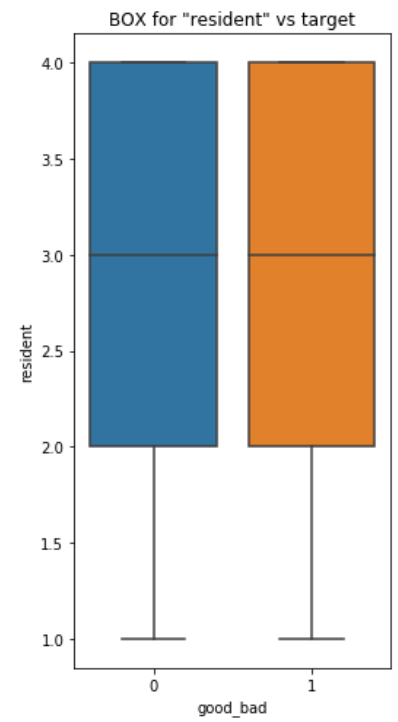
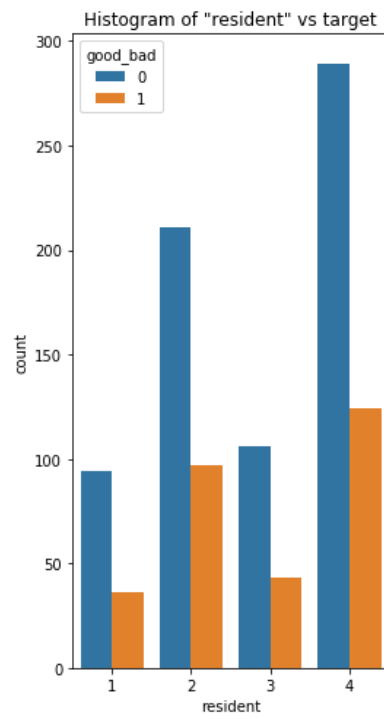
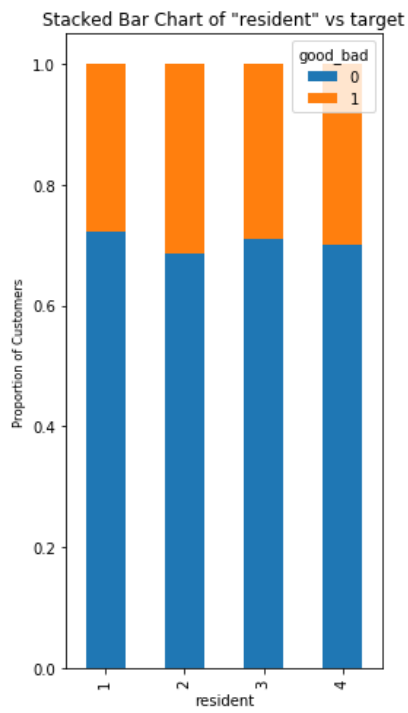


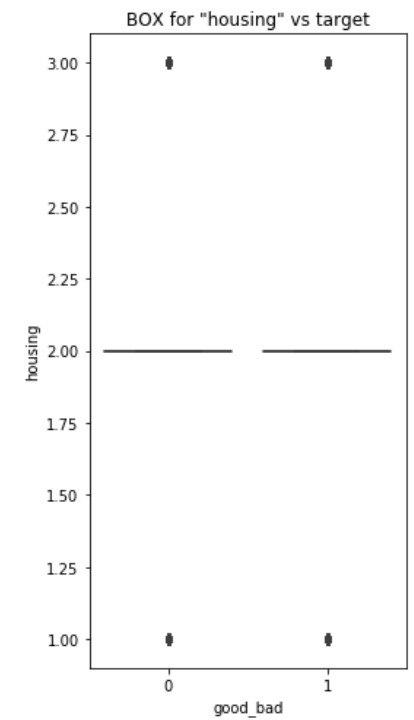
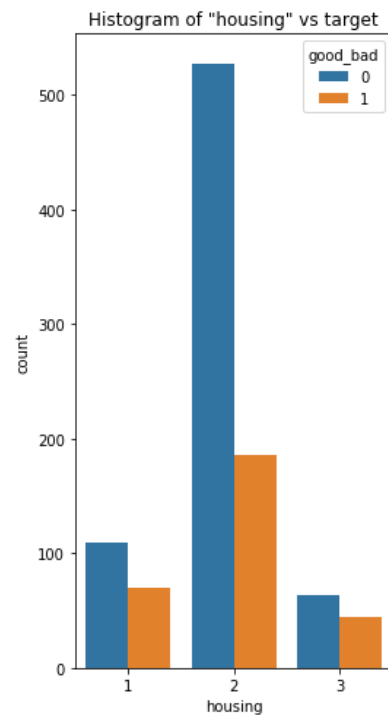
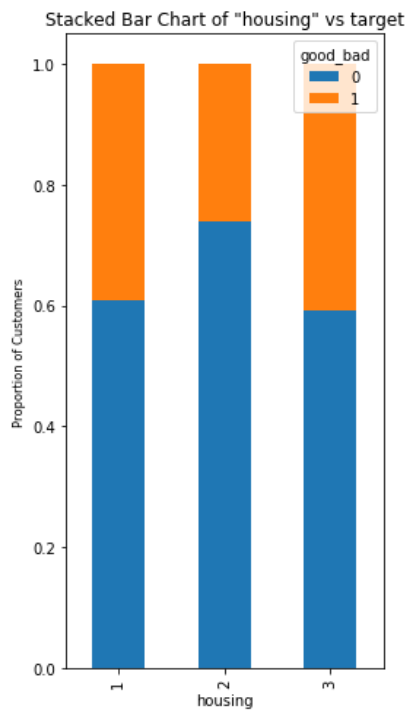
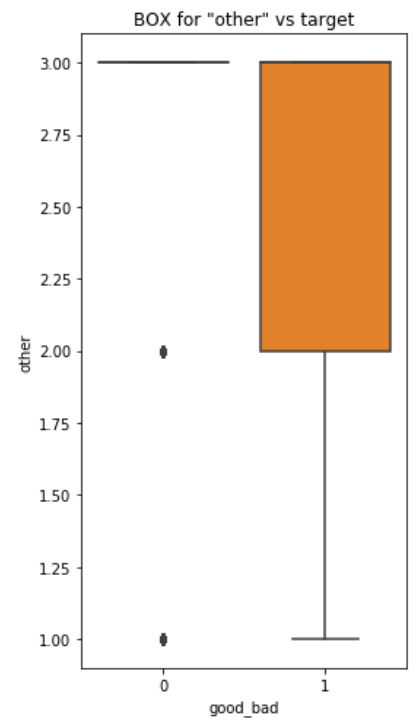
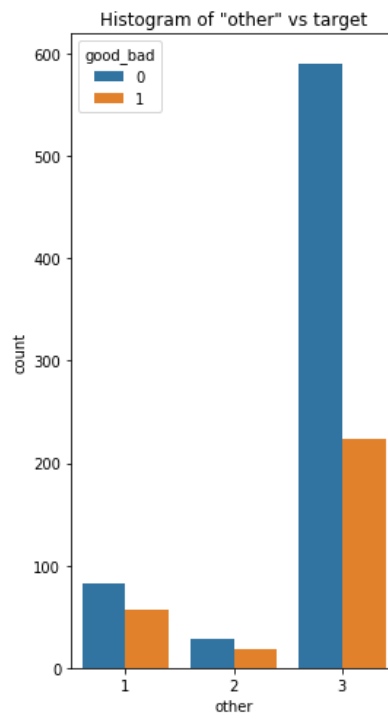
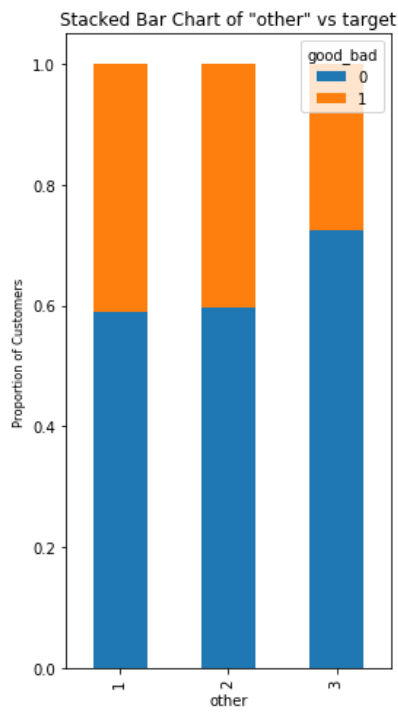
Histogram of "co-applicant status" vs target



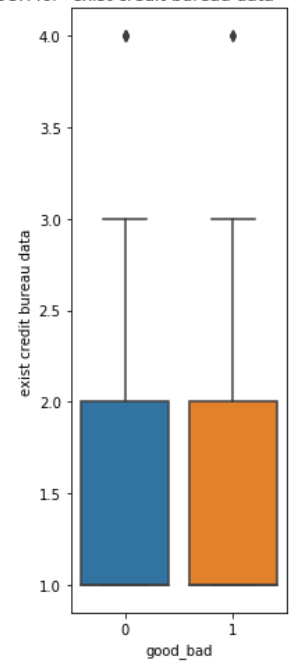
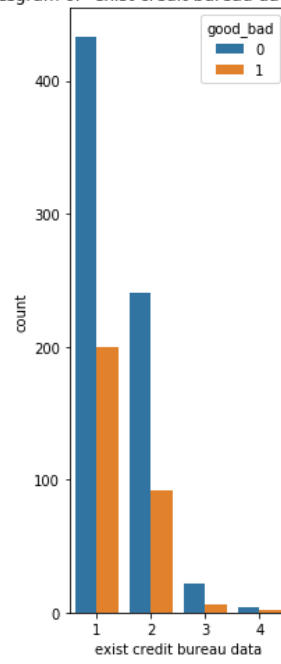
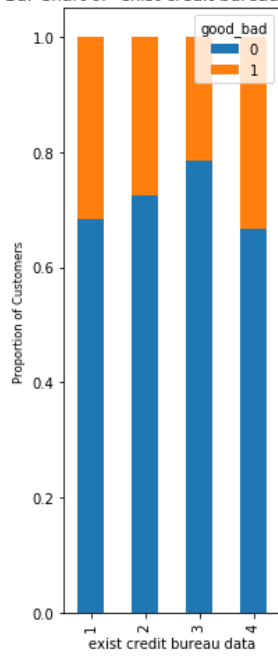
BOX for "co-applicant status" vs target



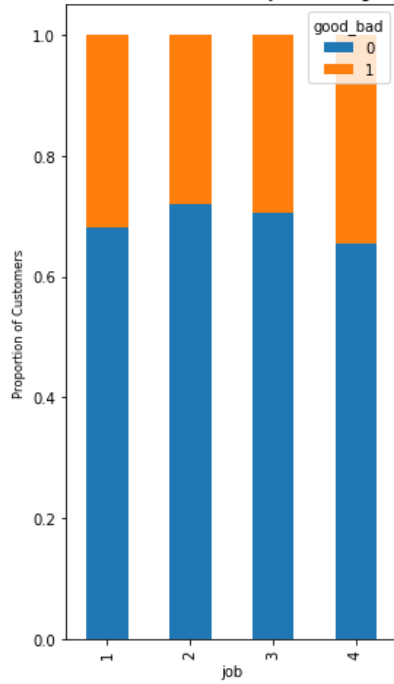




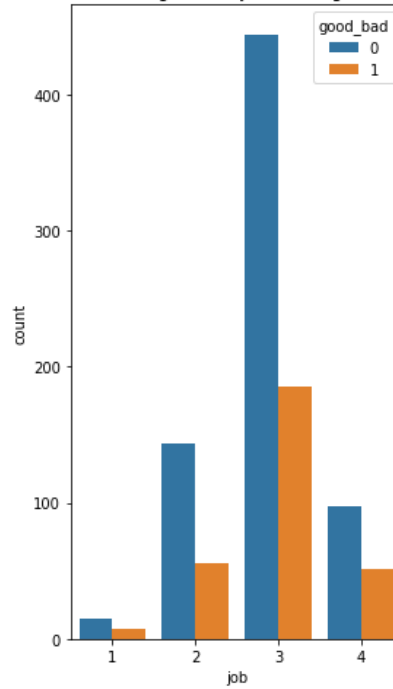
Stacked Bar Chart of "exist credit bureau data" vs target | Histogram of "exist credit bureau data" vs target | BOX for "exist credit bureau data" vs target



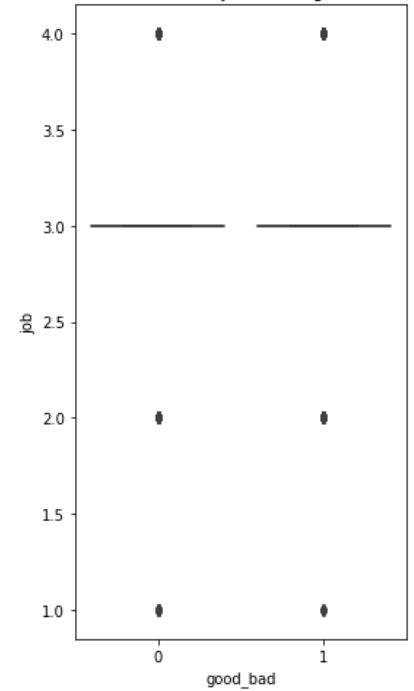
Stacked Bar Chart of "job" vs target



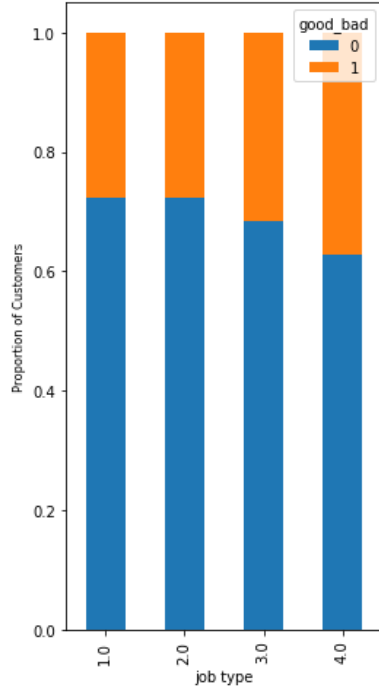
Histogram of "job" vs target



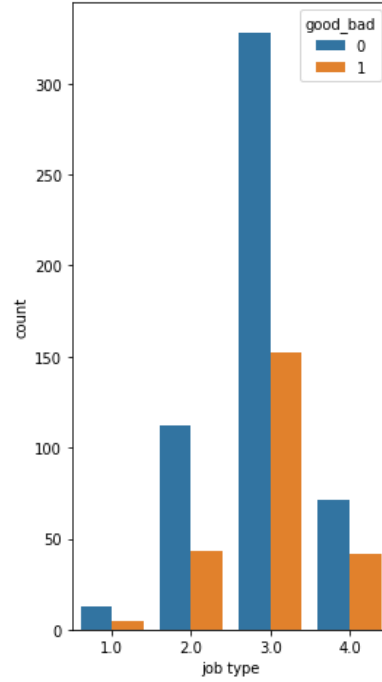
BOX for "job" vs target



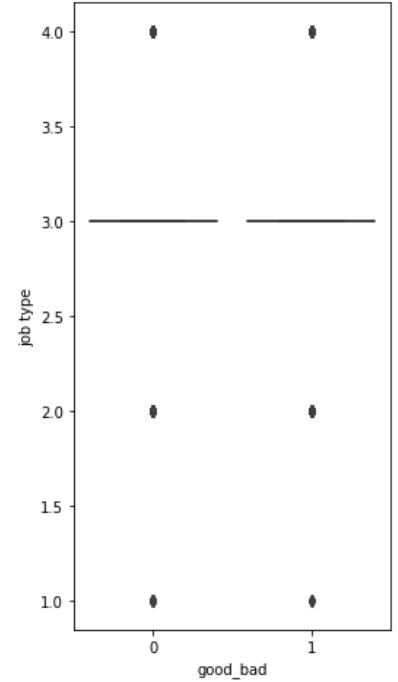
Stacked Bar Chart of "job type" vs target



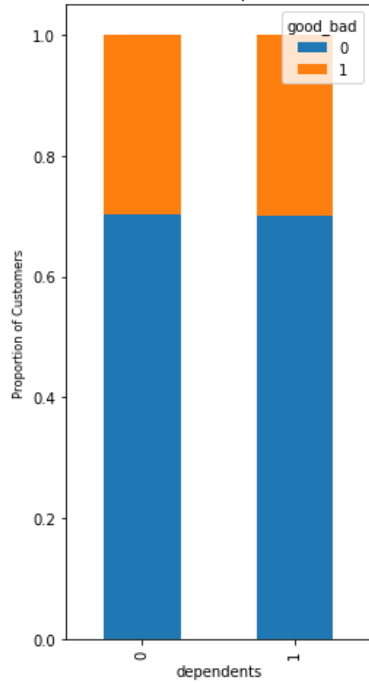
Histogram of "job type" vs target



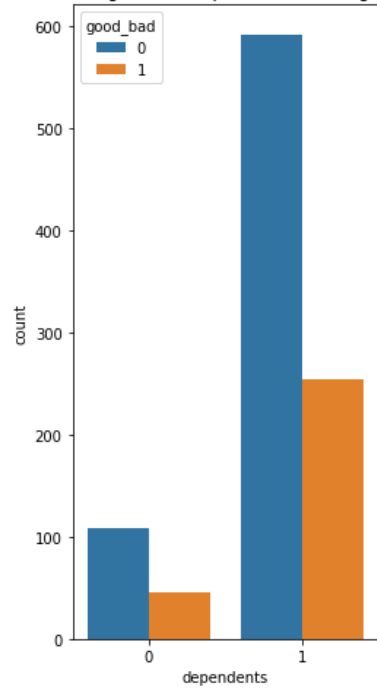
BOX for "job type" vs target



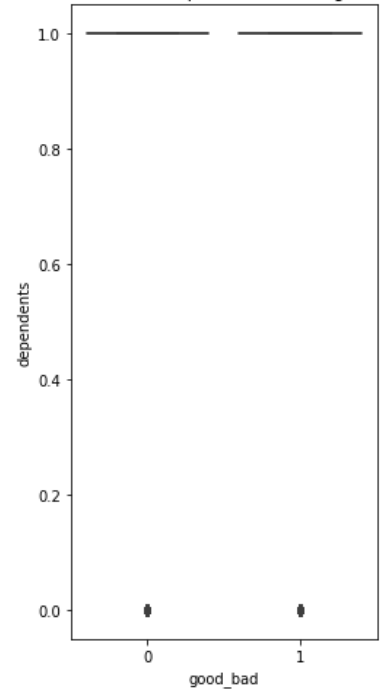
Stacked Bar Chart of "dependents" vs target



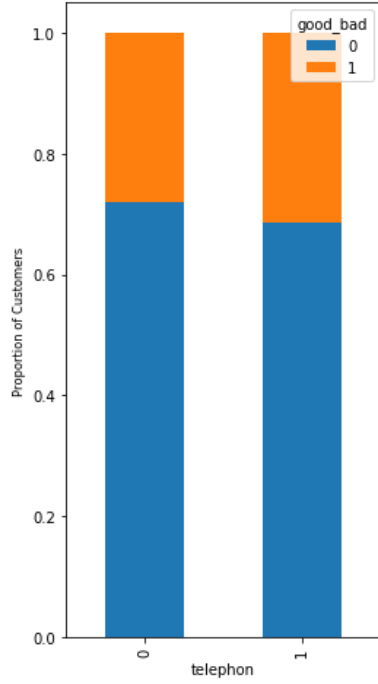
Histogram of "dependents" vs target



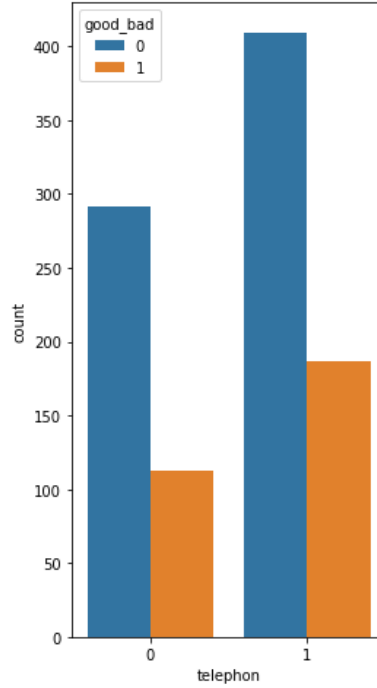
BOX for "dependents" vs target



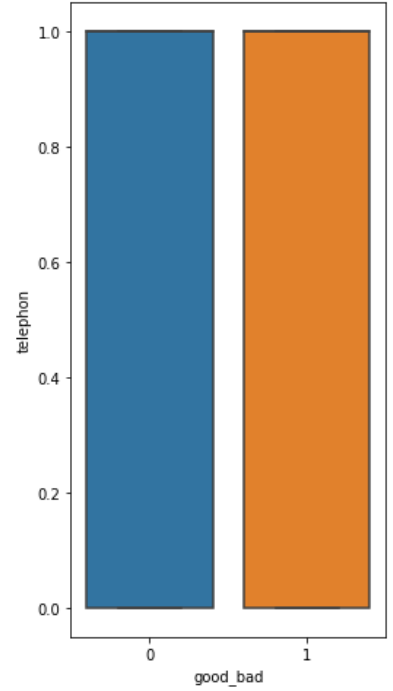
Stacked Bar Chart of "telephon" vs target



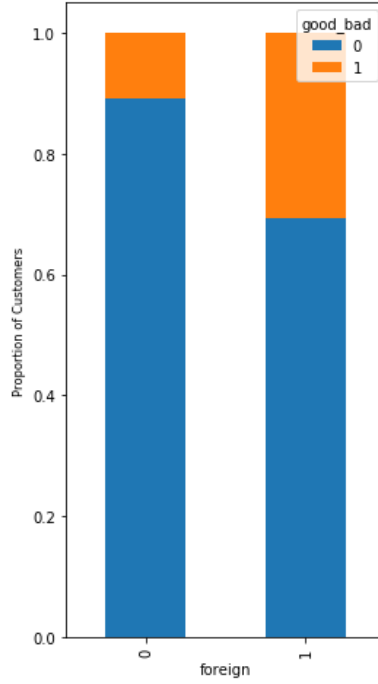
Histogram of "telephon" vs target



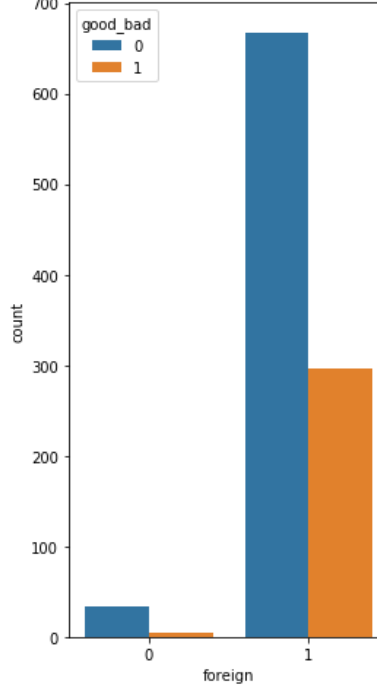
BOX for "telephon" vs target



Stacked Bar Chart of "foreign" vs target



Histogram of "foreign" vs target



BOX for "foreign" vs target

