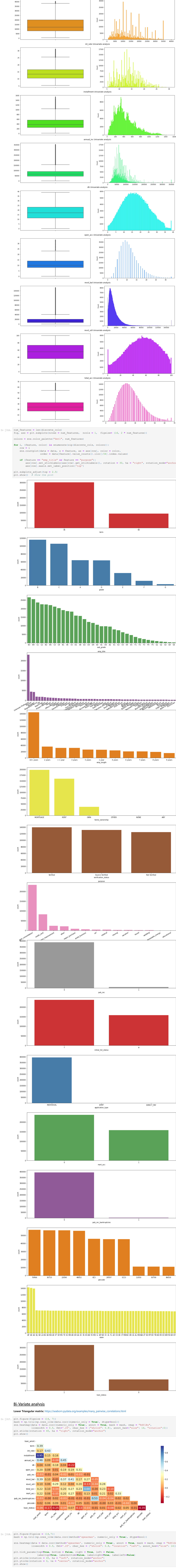


Summary

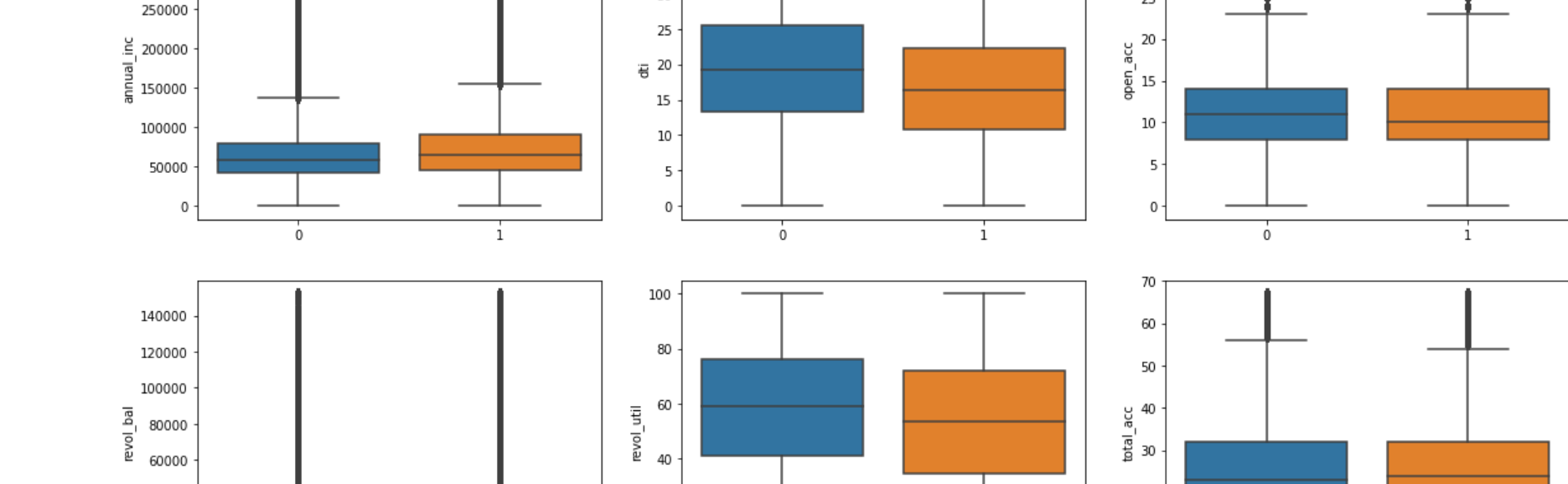
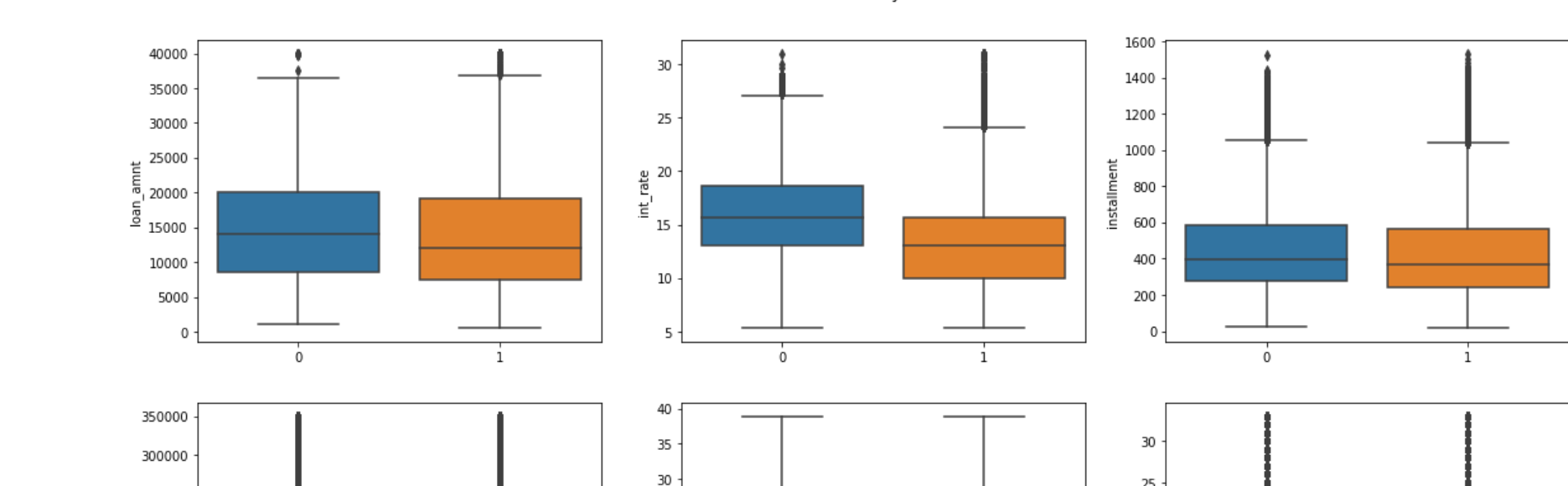
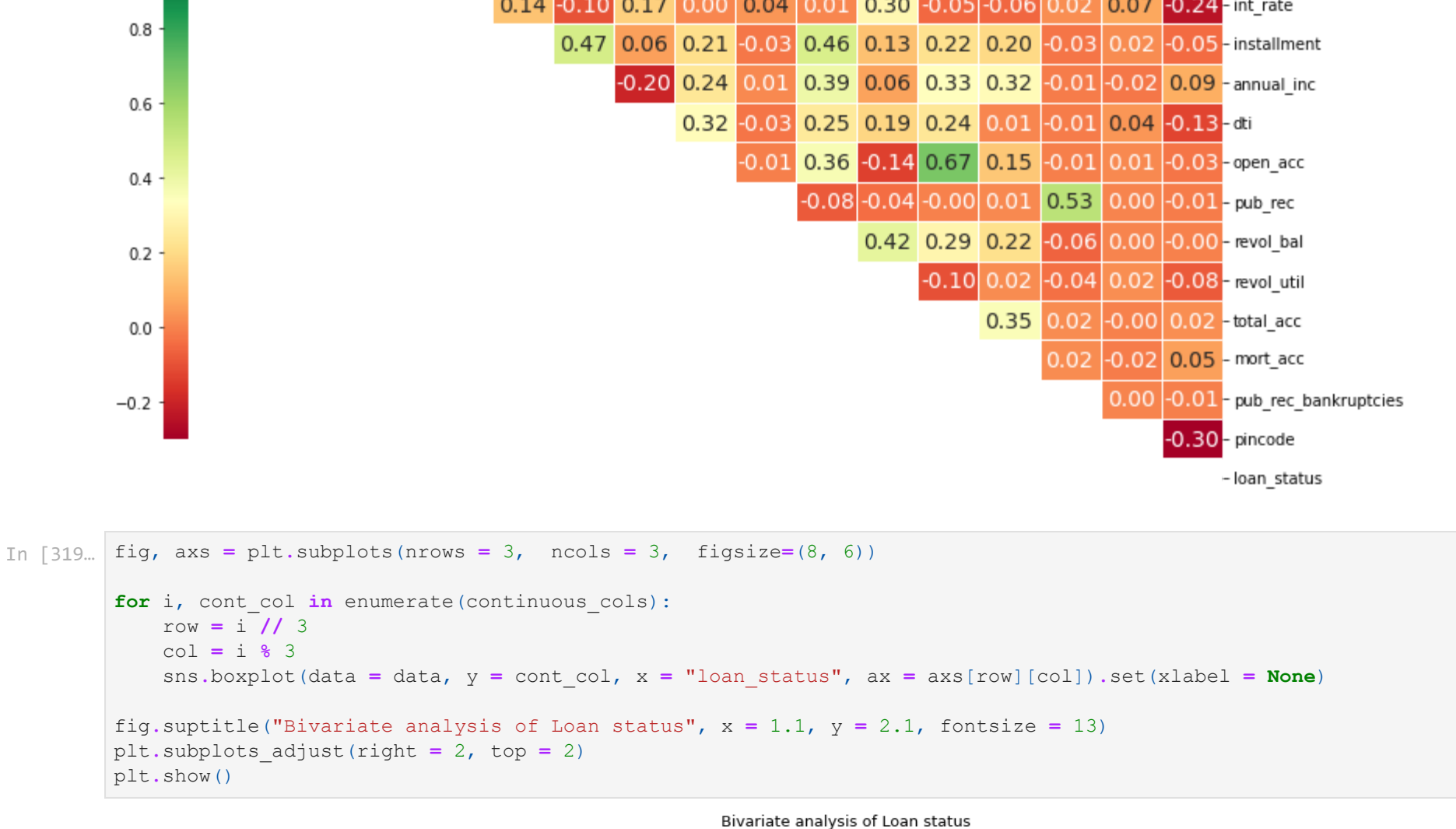
- Observation:** Data has 396030 entries, 25 columns (after initial feature engg). Data had very few null values and they were imputed without losing any data
- Observation:** The column address was extremely granular and has multiple repetitions
- Fix:** Divided the address field into state and pincode to bring about some homogeneity in the data. This will help us check if location is a factor in deciding the loan status
- Observation:** Columns such as Pub_rec, Mort_acc, Pub_rec_bankruptcies are considered negative for a person's creditworthiness. Having more of these can raise concerns for lenders.
Fix: Created flags
- Observation:** Our target value (loan status) has significant class imbalance, this will create issues in prediction.
Fix: We can under technique like undersampling, oversampling and smote to reduce this
- Observation:** Some features such as grade, verification status, initial list status have an approx. even distribution. Whereas, features like purpose, application type and home ownership have skewed data
Recommendation: Loanlap should notice this skew in some data and can gather more insights about their customers to provide services for their target audience. They can also use this data to try to increase their penetration in other markets they aren't dominant in yet
- Observation:** We have very few outliers (as a % of total, but they all represent real valued data. However, some outliers are extremely significant (like for very rich income individuals).
Fix: We will cap the data to reduce the disparity in data. To cap the data we will use 99.7% quantile for upper value and 0.03% quantile for lower value

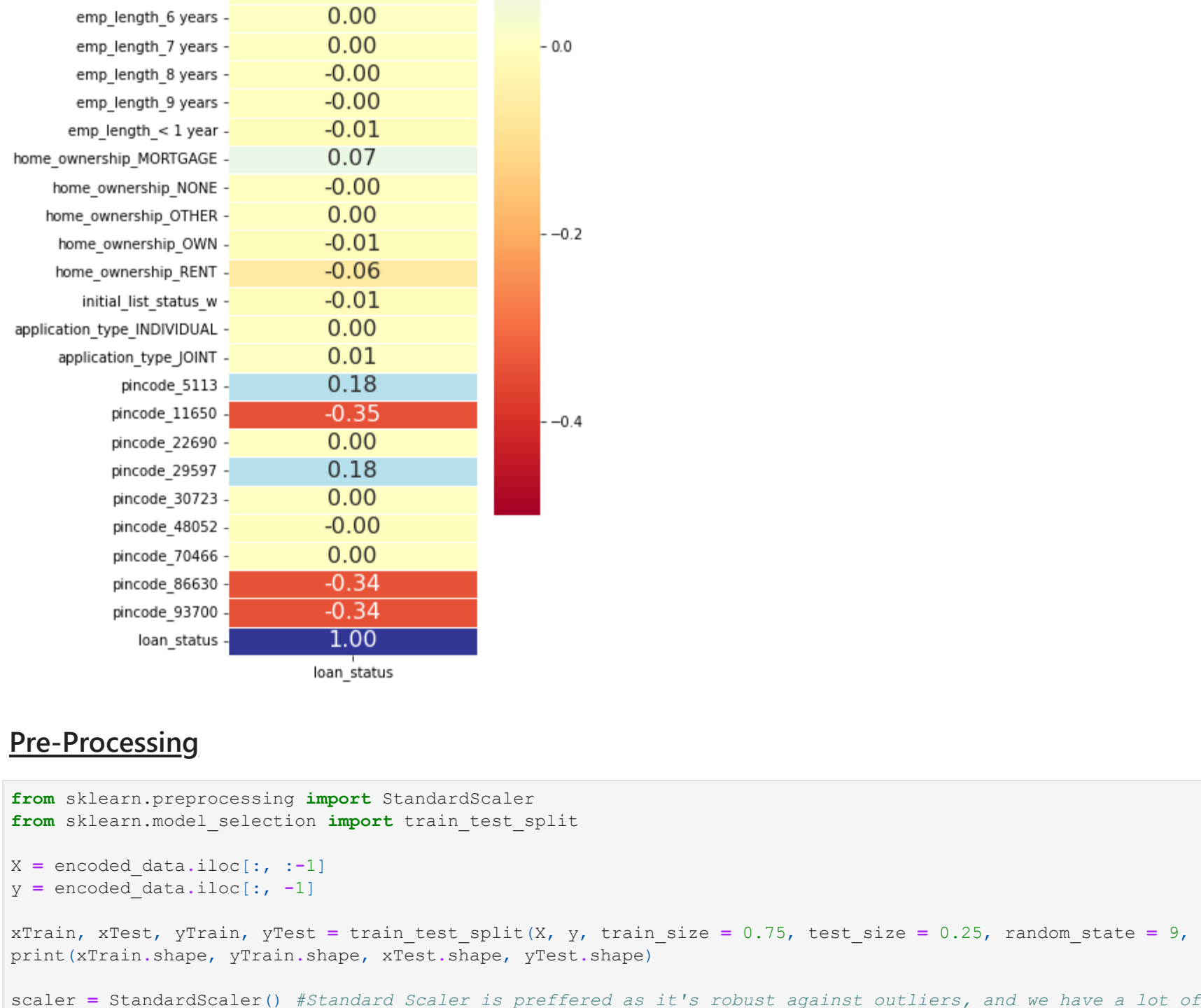
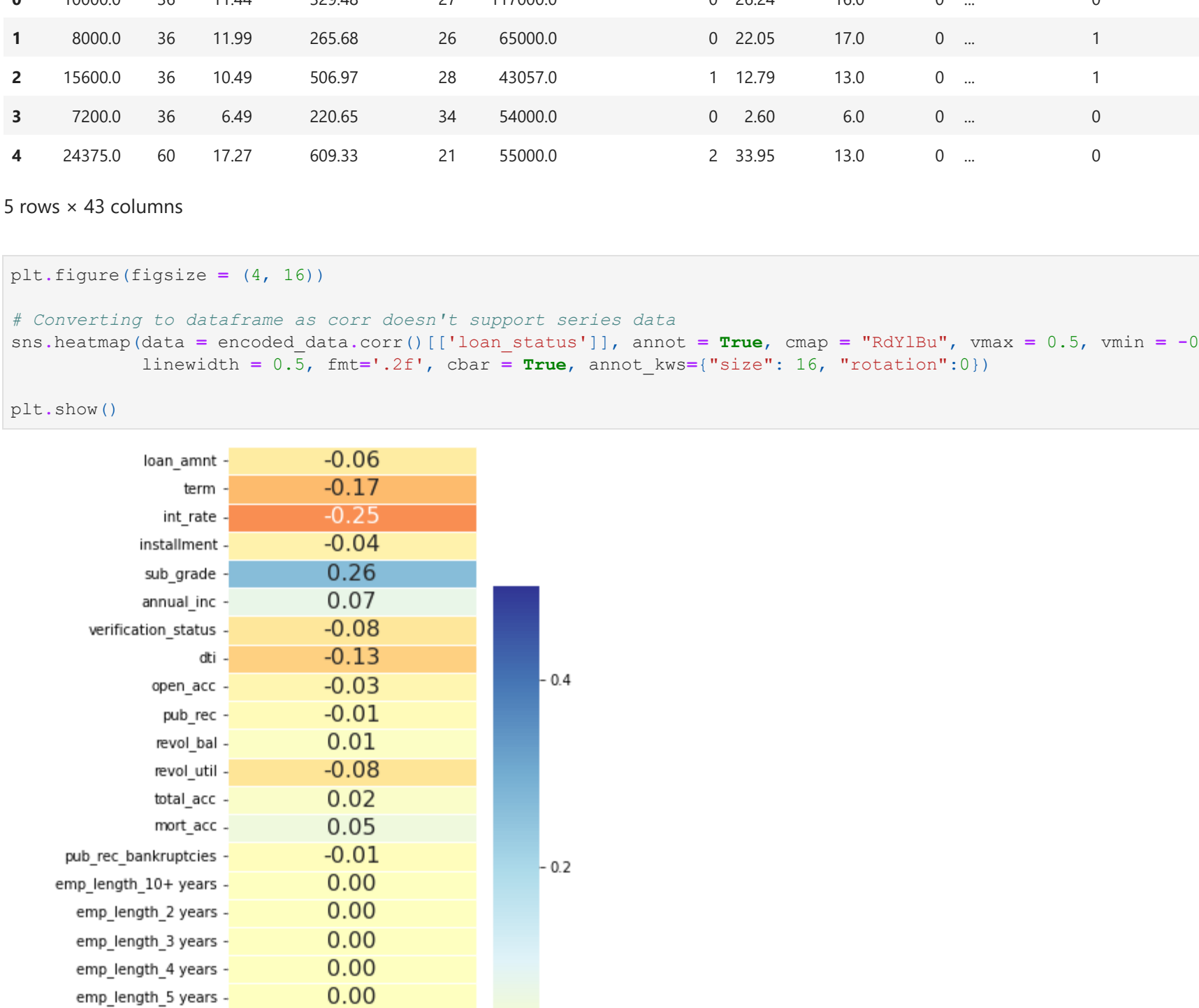
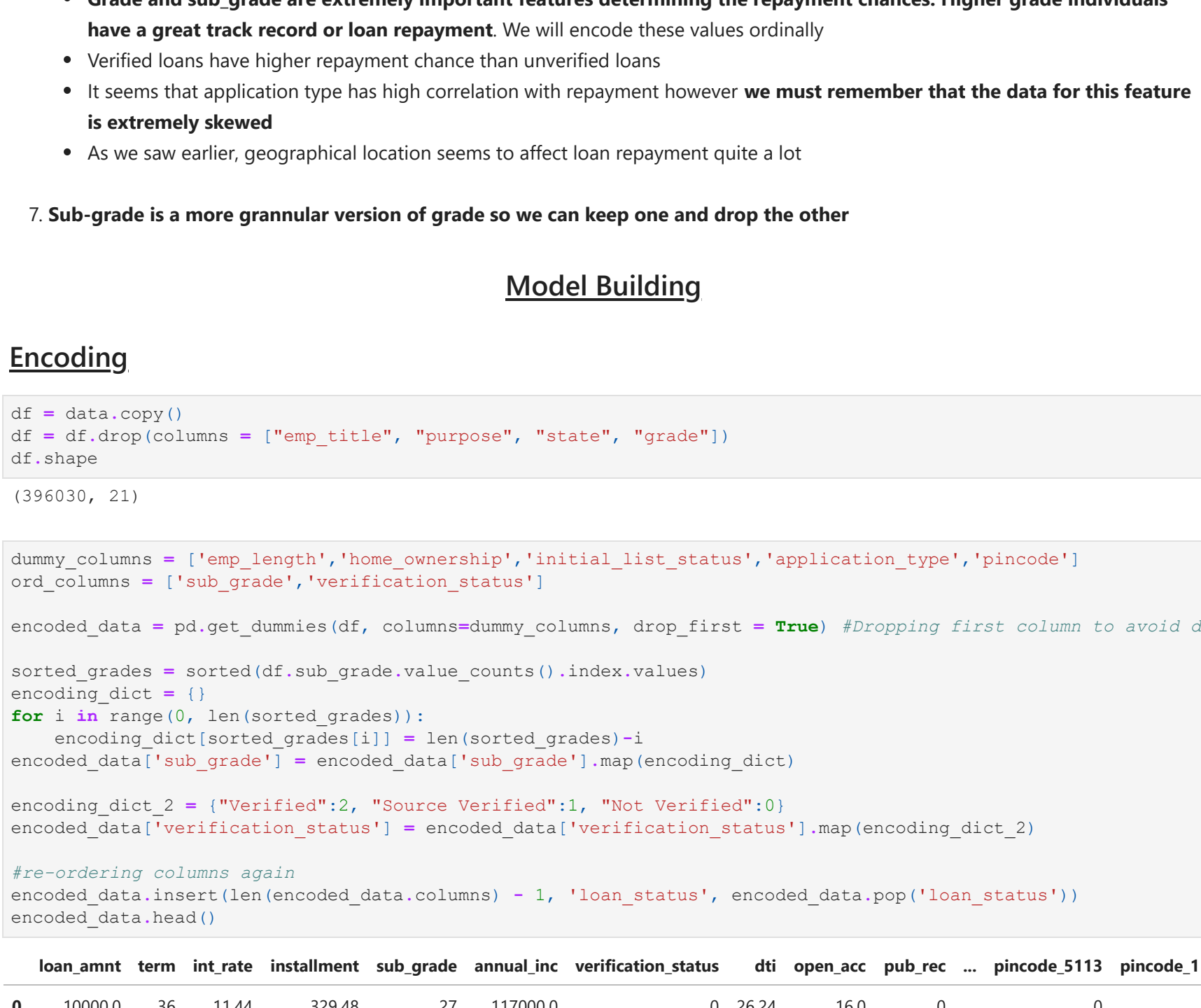
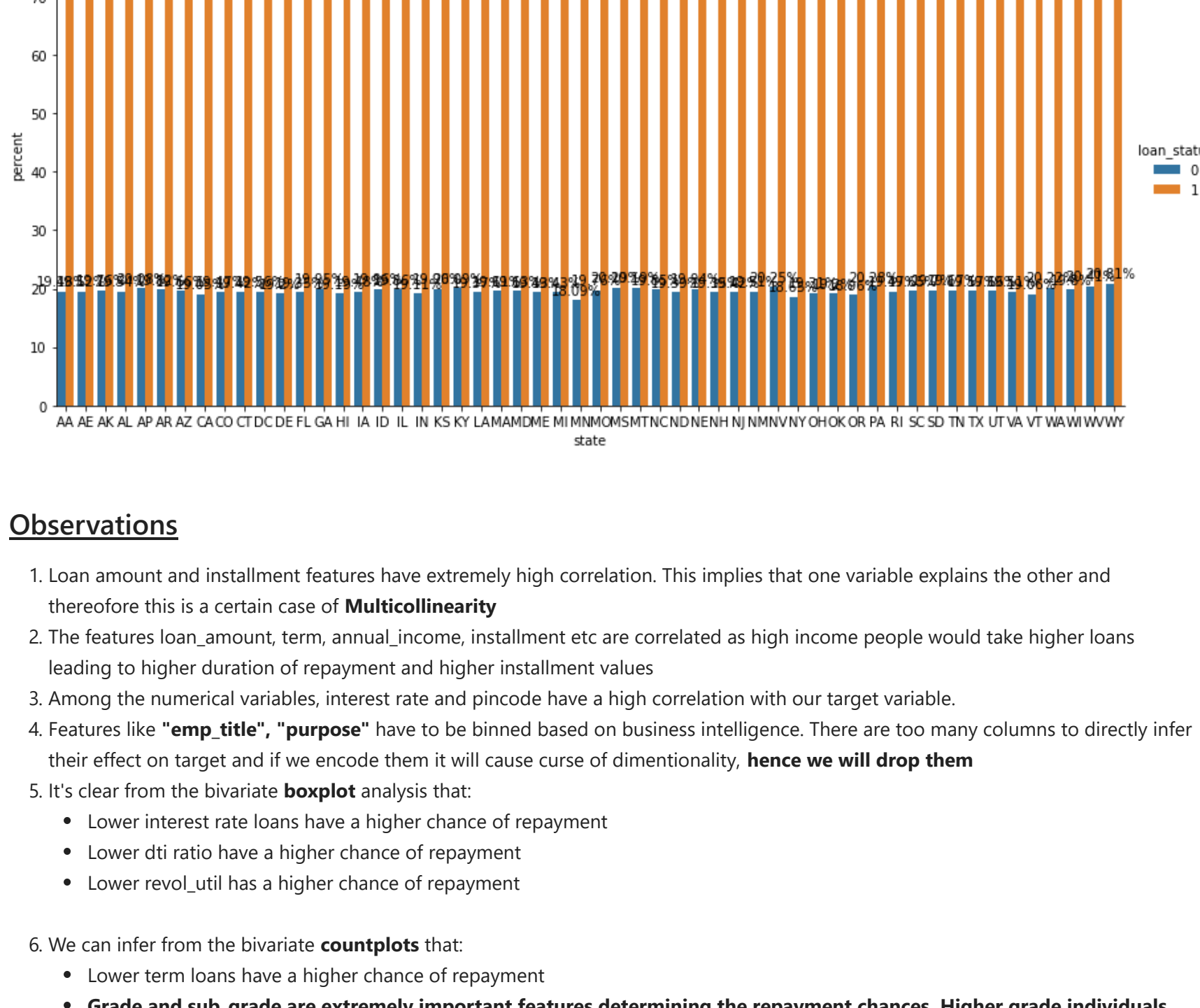
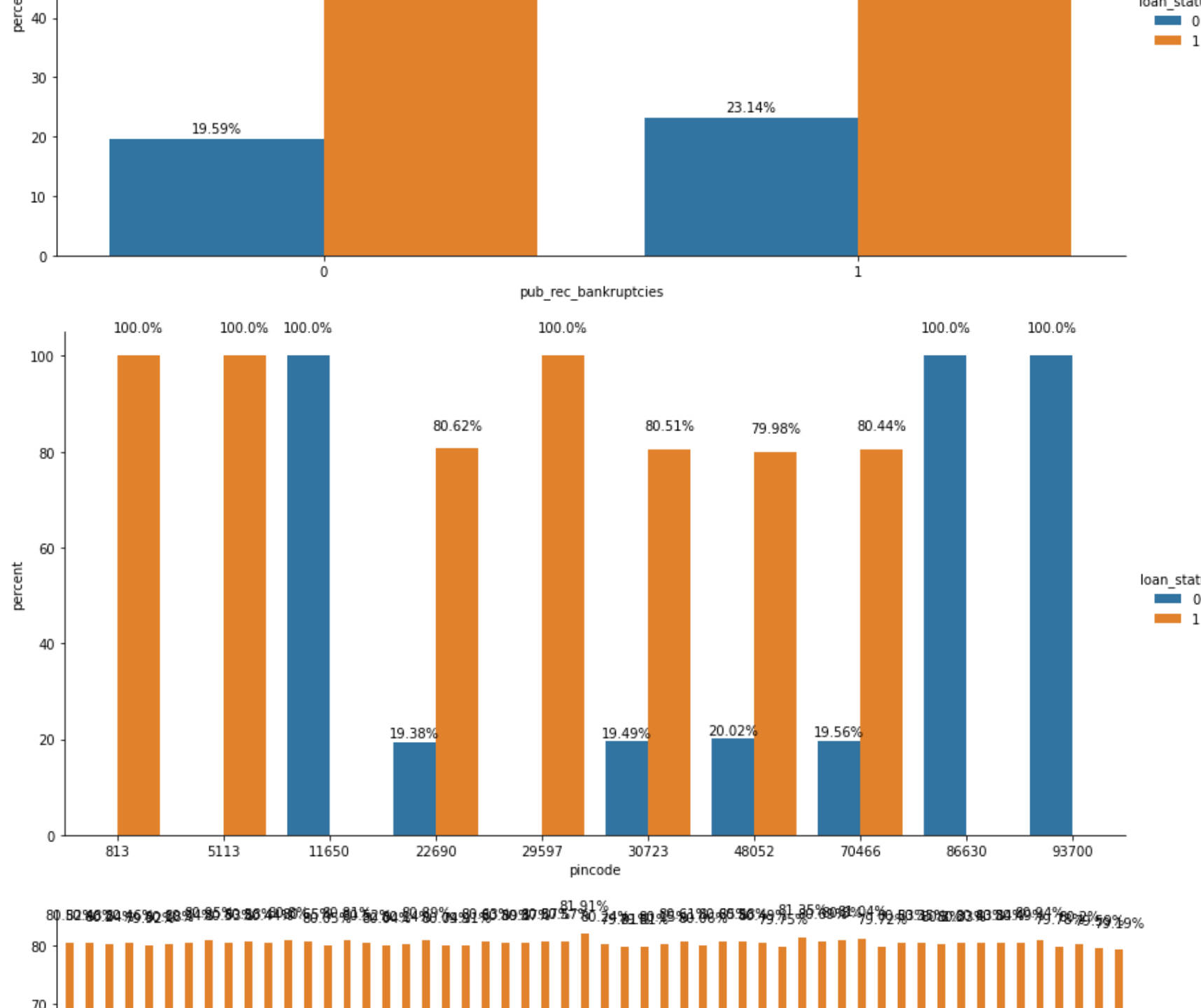
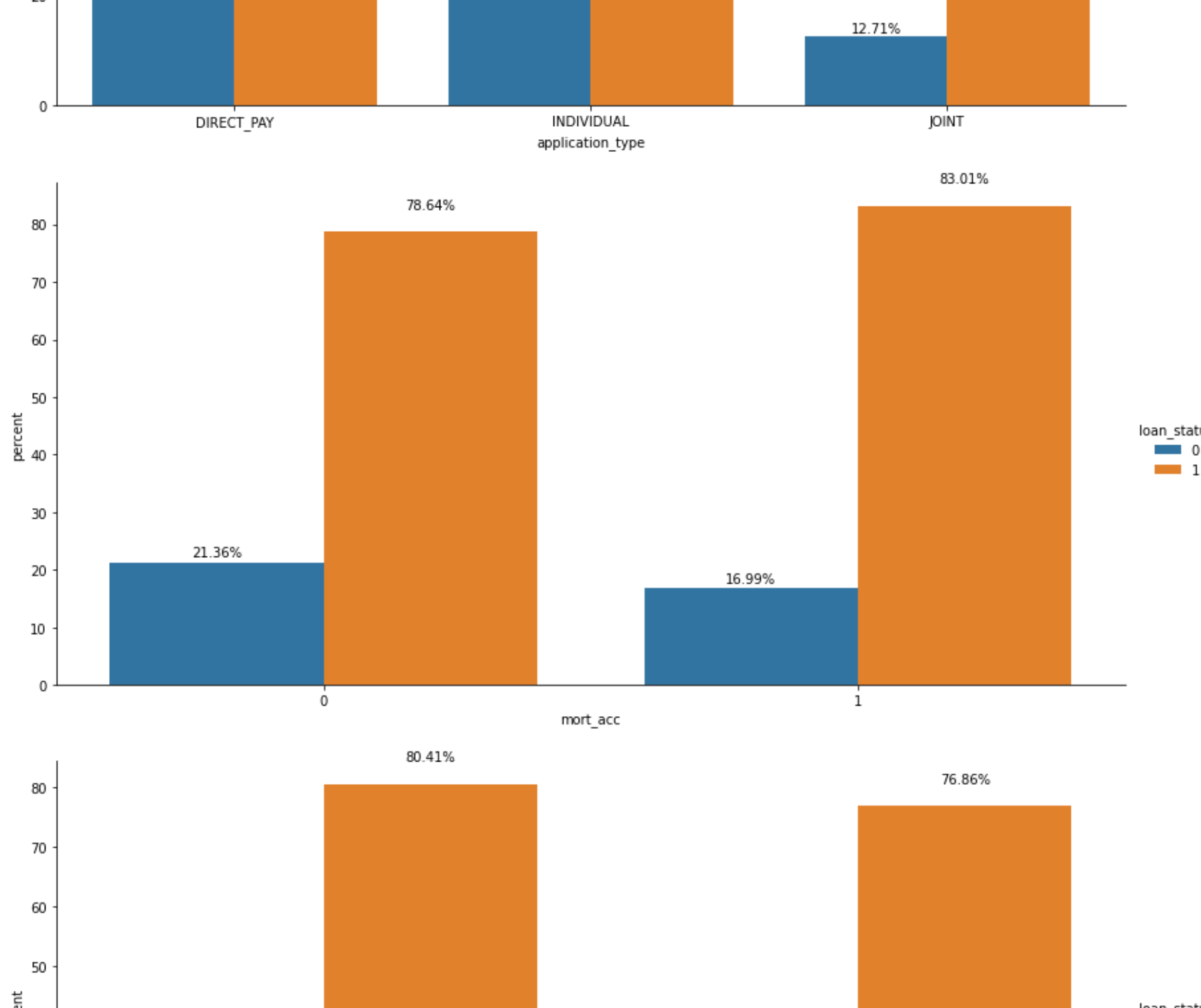
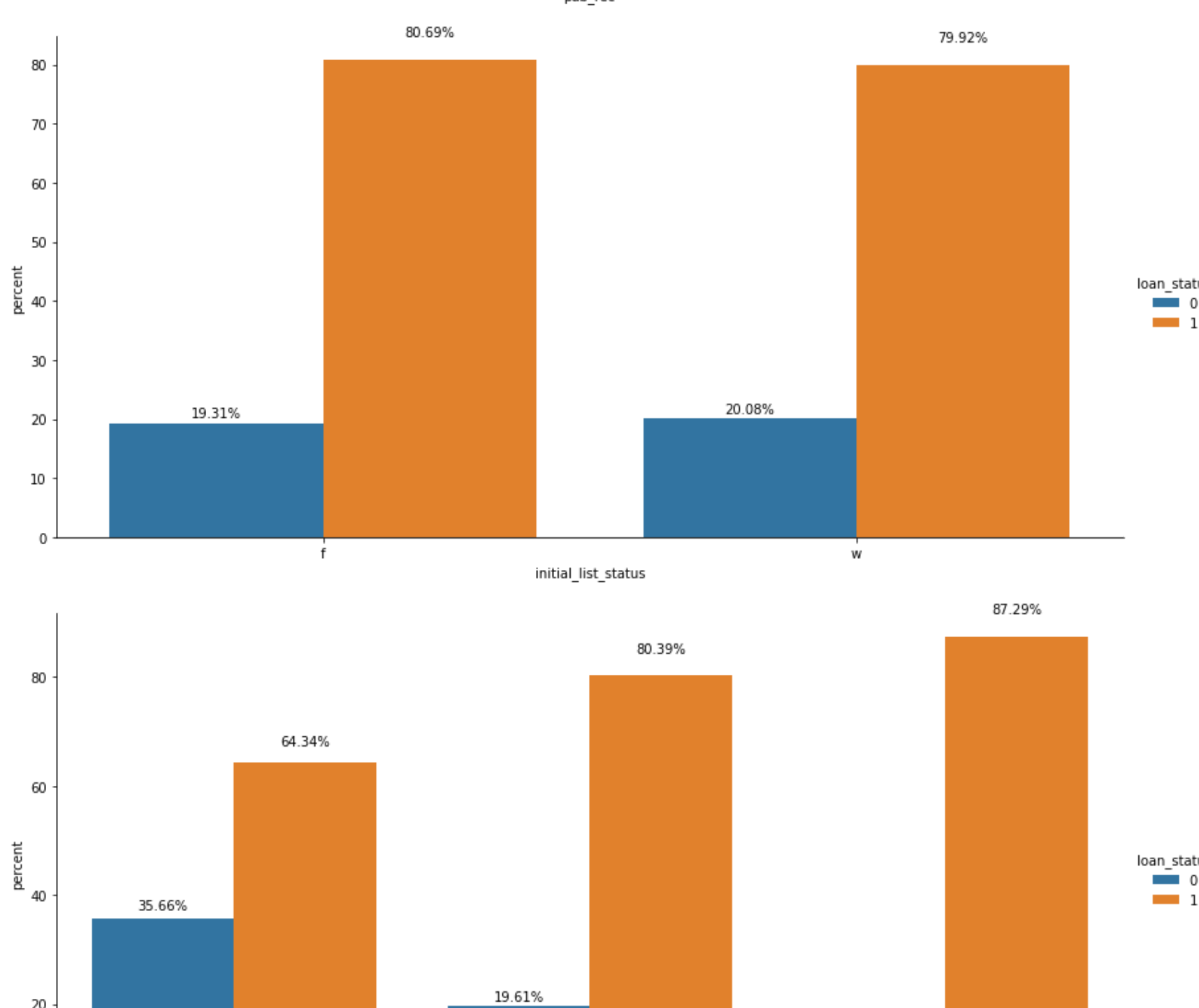
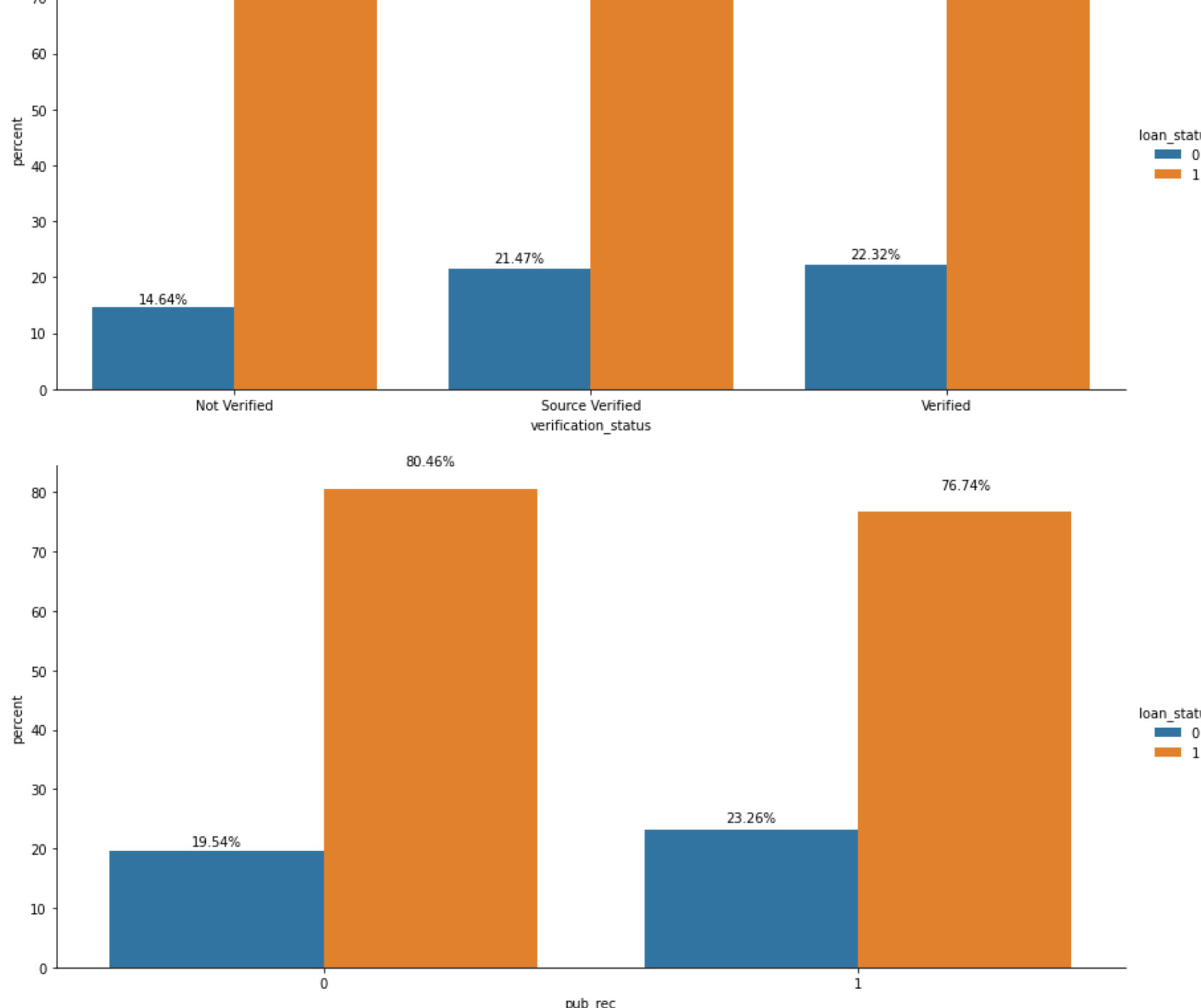
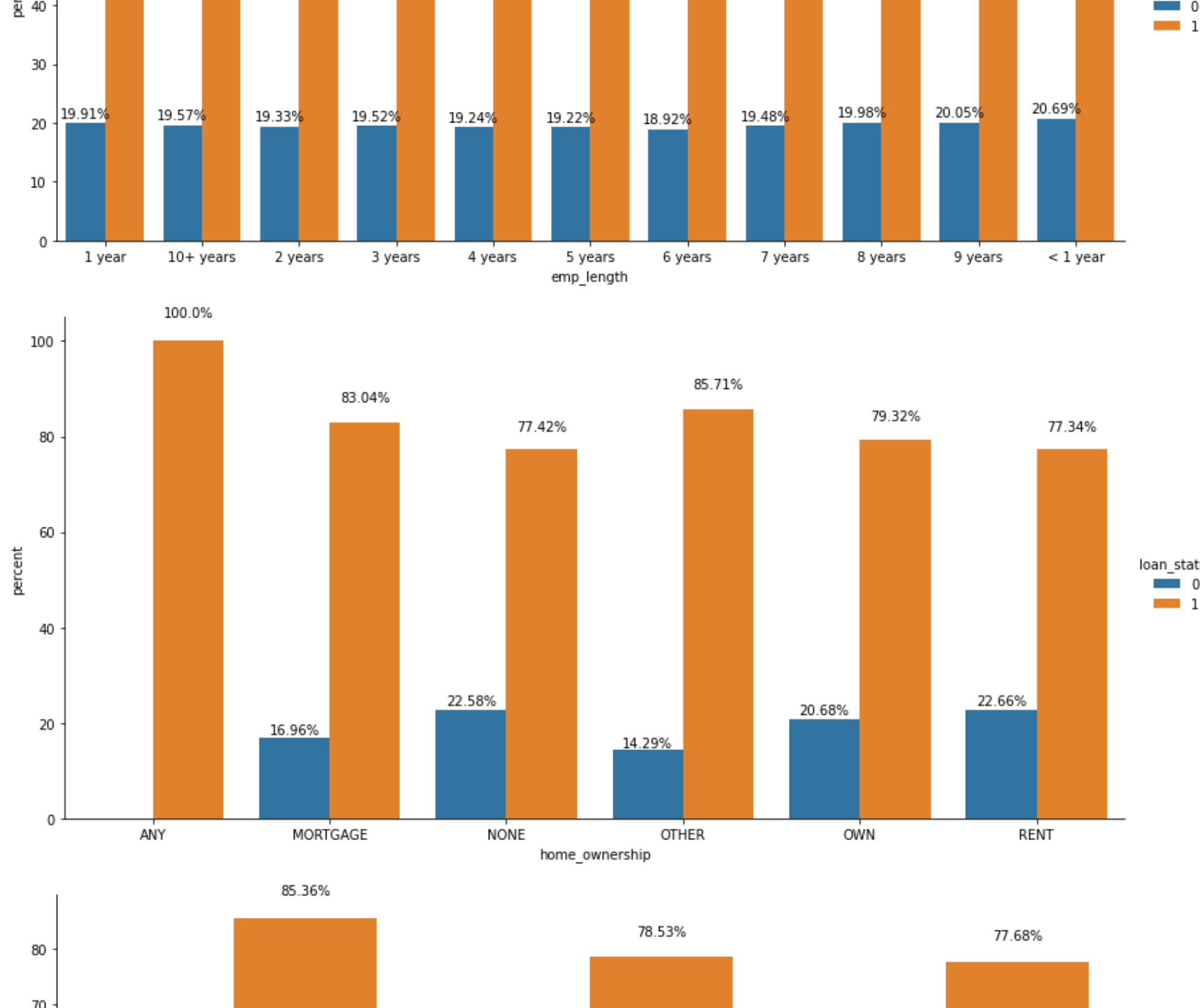
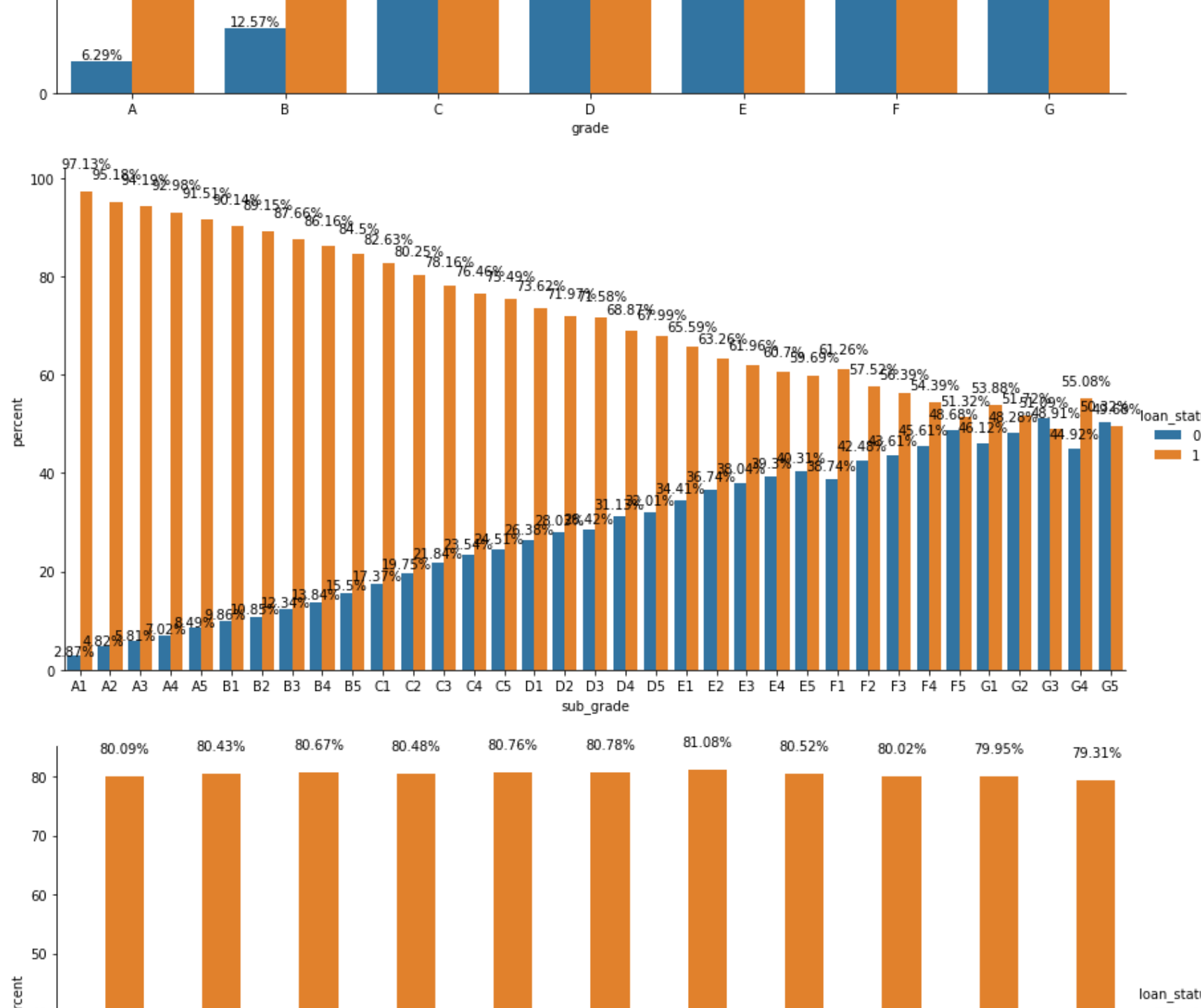
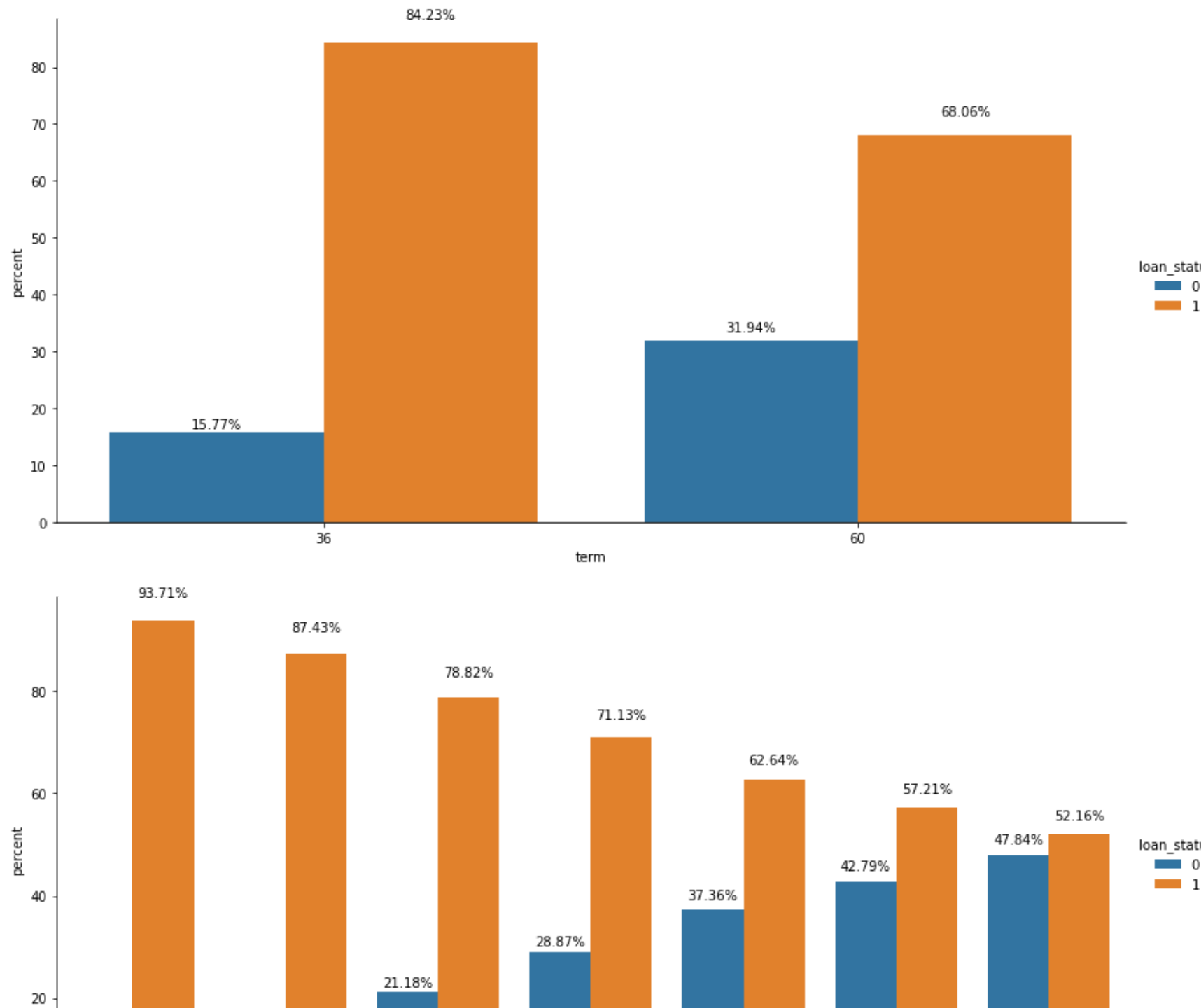
Exploratory Data Analysis



Bi-Variate analysis

Lower Triangular matrix: https://seaborn.pydata.org/examples/many_pairwise_correlations.html





Observations

- Loan amount and installment features have extremely high correlation. This implies that one variable explains the other and therefore this is a certain case of **Multicollinearity**
- The features loan_amount, term, annual_income, installment etc are correlated as high income people would take higher loans leading to higher duration of repayment and higher installment values
- Among the numerical variables, interest rate and pincode have a high correlation with our target variable.
- Features like "term", "purpose" are to be binned based on business intelligence. There are too many columns to directly infer their effect on target and if we encode them it will cause curse of dimensionality, **hence we will drop them**
- It's clear from the bivariate **boxplot** analysis that:
 - Lower interest rate loans have a higher chance of repayment
 - Lower dti ratio have a higher chance of repayment
 - Lower revol_util has a higher chance of repayment
- We can infer from the bivariate **countplots** that:
 - Lower term loans have a higher chance of repayment
 - Grade and sub_grade are extremely important features determining the repayment chances. Higher grade individuals have a great track record or loan repayment.** We will encode these values ordinarily
 - Verified loans have higher repayment chance than unverified loans
 - It seems that application type has high correlation with repayment however **we must remember that the data for this feature is extremely skewed**
 - As we saw earlier, geographical location seems to affect loan repayment quite a lot
- Sub-grade is a more granular version of grade so we can keep one and drop the other

Model Building

Encoding

```
In [32]: df = data.copy()
df = df.drop(columns = ['emp_title', 'purpose', 'state', 'grade'])
df.shape
Out[32]: (396030, 21)
```

```
In [33]: dummy_columns = ['emp_length', 'home_ownership', 'initial_list_status', 'application_type', 'pincode']
ord_columns = ['sub_grade', 'verification_status']
encoded_data = pd.get_dummies(df, columns=dummy_columns, drop_first = True) #Dropping first column to avoid dummy
sorted_grades = sorted(df.sub_grade.value_counts().index.values)
encoding_dict = {}
for i in range(0, len(sorted_grades)):
    encoding_dict[sorted_grades[i]] = len(sorted_grades)-i
encoded_data['sub_grade'] = encoded_data['sub_grade'].map(encoding_dict)
encoding_dict_2 = {"Verified":2, "Source Verified":1, "Not Verified":0}
encoded_data['verification_status'] = encoded_data['verification_status'].map(encoding_dict_2)
```

```
#re-ordering columns again
encoded_data.insert(len(encoded_data.columns) - 1, 'loan_status', encoded_data.pop('loan_status'))
encoded_data.head()
Out[33]:
```

	loan_amnt	term	int_rate	installment	sub_grade	annual_inc	verification_status	dti	open_acc	pub_rec	...	pincode	5113	pincode	1165
0	10000.0	36	11.44	329.48	27	117000.0	0	26.24	16.0	0	...	0	...	0	...
1	8000.0	36	11.99	265.68	26	65000.0	0	22.05	17.0	0	...	0	...	1	...
2	15600.0	36	10.49	506.97	28	43057.0	1	12.79	13.0	0	...	1	...	1	...
3	7200.0	36	6.49	220.65	34	54000.0	0	2.60	6.0	0	...	0	...	0	...
4	24375.0	60	17.27	609.33	21	55000.0	2	31.95	13.0	0	...	0	...	0	...

5 rows × 13 columns

```
In [35]: plt.figure(figsize = (4, 16))
sns.heatmap(data = encoded_data.corr()[['loan_status']], annot = True, cmap = 'magma', vmax = 0.5, vmin = -0.5, linewidth = 0.5, data_cmap = 'magma', bar = True, annot_kws={"size": 16, "rotation":0})
plt.show()
```



Pre-Processing

```
In [34]: from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
X = encoded_data.iloc[:, :-1]
y = encoded_data.iloc[:, -1]
xTrain, xTest, yTrain, yTest = train_test_split(X, y, train_size = 0.75, test_size = 0.25, random_state = 9, sh
print(xTrain.shape, yTrain.shape, xTest.shape, yTest.shape)
```

```
scaler = StandardScaler() #Standard Scaler is preferred as it's robust against outliers, and we have a lot of o
xTrain_Scaled = scaler.fit_transform(xTrain)
xTrain_Scaled_df = pd.DataFrame(data = scaler.fit_transform(xTrain), columns = xTrain.columns, index = xTrain.i
xTest_Scaled = scaler.transform(xTest)
xTest_Scaled_df = pd.DataFrame(data = scaler.transform(xTest), columns = xTest.columns, index = xTest.index)
(297022, 42) (297022, (99008, 42) (99008,)
```

```
In [35]: from statsmodels.stats.outliers_influence import variance_inflation_factor
from statsmodels.tools.tools import add_constant
def dropVIF(df, col_list = None):
    thresh = 10 #V.I.F of > 10 will be dropped
    temp_df = add_constant(df) #statmodel requires constant column separately
    if col_list is None:
        col_list = df.columns.tolist()
        return "one or more element passed in col_list aren't present in given dataframe"
    col_list.insert(0, 'const')
    temp_df = temp_df[col_list]
    flag = 1
```

```
    while flag == 1: # While any column has V.I.F > 10
        flag = 0 #reset flag
        vif_max = 1 #lowest possible VIF score
        col_to_drop = None
        for i in range(len(temp_df.columns)): # Calculating V.I.F for each column
            vif = variance_inflation_factor(temp_df, i)
            vif = np.round(vif, 0)
            if (vif > thresh and vif > vif_max):
                col_to_drop = temp_df.columns[i]
                vif_max = vif
        print(f'Dropped column {col_to_drop} with V.I.F = {vif_max}')
        temp_df.drop(columns = df.columns[df.columns == col_to_drop], inplace = True)
```

```
def getVIF(df):
    temp_df = add_constant(df) #statmodel requires constant column separately
    return pd.DataFrame(data = variance_inflation_factor(temp_df, i) for i in range(0, temp_df.shape[1]),
        index = temp_df.columns.tolist(),
        columns = ['VIF'])
```

```
In [36]: VIF_df = getVIF(xTrain_Scaled_df)
VIF_df
```

	VIF
const	1.00000
loan_amnt	58.698951
term	6.646840
int_rate	21.713933
installment	50.900002
sub_grade	214.54373
annual_inc	1.713175
verification_status	1.174336
dti	1.420700
open_acc	2.137311
pub_rec	1.413369
revol_bal	1.556458
revol_util	1.354573
total_acc	2.191897
mort_acc	1.494894
pub_rec_bankruptcies	1.403470
emp_length_10+ years	4.292723
emp_length_2 years	2.163669
emp_length_3 years	2.040969
emp_length_4 years	1.809903
emp_length_5 years	1.893146
emp_length_6 years	1.712220
emp_length_7 years	1.719341
emp_length_8 years	1.663412
emp_length_9 years	1.538372
emp_length_1 year	2.052605
home_ownership_MORTGAGE	24756.528022
home_ownership_NONE	8.334163
home_ownership_OTHER	29.663377
home_ownership_OWEN	8487.441357
home_ownership_RENT	23846.926783
initial_list_status_w	1.071810
application_type_INDIVIDUAL	2.470954
application_type_JOINT	2.462060
pincode_5113	1.760634
pincode_11650	1.226181
pincode_22690	1.917134
pincode_29597	1.764420
pincode_30723	1.915895
pincode_48052	1.909913
pincode_70466	1.923908
pincode_86630	1.219957
pincode_93700	1.223428

```
In [38]: dropVIF(xTrain_Scaled_df) #encoded_data.columns[0:15].tolist()
Dropped column home_ownership_MORTGAGE with V.I.F = 24757.0
```

```
In [36]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report
lr = LogisticRegression(penalty='l2', max_iter = 1000)
lr.fit(xTrain_Scaled, yTrain)
yPred = lr.predict(xTest_Scaled)
report = classification_report(yTest, yPred)
print(report)
```

	precision	recall	f1-score	support
0	0.92	0.48	0.63	19575
1	0.88	0.99	0.93	79433
accuracy	0.90	0.73	0.89	99008
macro avg	0.90	0.73	0.78	99008
weighted avg	0.89	0.89	0.87	99008

```
In [38]: feature_imp = pd.DataFrame(data = np.round(np.insert(lr.coef_, 0, lr.intercept_), 2),
    https://stackoverflow.com/questions/596442/what-does-the-high-vif-for-the-constant-term-intercept-indicate
https://stats.stackexchange.com/questions/332428/regression-model-constant-causes-multicollinearity-warning-but-not-in-standards
feature_imp['Sorted'] = feature_imp['co-efficients'].abs()
feature_imp.sort_values(by = 'Sorted', axis = 0, ascending=False)
```

```
Out[38]:
```

	Co-efficients	Sorted
Constant	3.87	3.87
pincode_11650	-3.09	3.09
pincode_93700	-3.07	3.07
pincode_86630	-3.04	3.04
pincode_70466	-2.64	2.64
pincode_48052	-2.63	2.63
pincode_22690	-2.62	2.62
pincode_30723	-2.62	2.62
pincode_5113	0.87	0.87
pincode_29597	0.87	0.87
sub_grade	0.75	0.75
int_rate	0.30	0.30
home_ownership_RENT	-0.27	0.27
term	-0.19	0.19
annual_inc	0.18	0.18
home_ownership_MORTGAGE	-0.14	0.14
open_acc	-0.13	0.13
total_acc	0.12	0.12
home_ownership_OWEN	-0.12	0.12
revol_util	-0.09	0.09
revol_bal	0.08	0.08
installment	-0.07	0.07
verification_status	-0.04	0.04
pub_rec	-0.03	0.03
loan_amnt	-0.02	0.02
application_type_JOINT	0.02	0.02
mort_acc	0.02	0.02
application_type_INDIVIDUAL	-0.01	0.01
emp_length_8 years	0.01	0.01
emp_length_1 year	-0.01	0.01
emp_length_9 years	0.01	0.01
emp_length_7 years	0.01	0.01
emp_length_5 years	0.01	0.01
emp_length_4 years	0.01	0.01
emp_length_3 years	0.01	0.01
emp_length_2 years	0.01	0.01
emp_length_10+ years	-0.01	0.01
emp_length_6 years	0.01	0.01
home_ownership_OTHER	0.00	0.00
initial_list_status_w	-0.00	0.00
pub_rec_bankruptcies	-0.00	0.00

```
In [38]: from sklearn.metrics import RocCurveDisplay
RocCurveDisplay.from_estimator(lr, xTest_Scaled, yTest, pos_label = 1)
plt.show()
```



```
In [39]: from sklearn.metrics import RocCurveDisplay
RocCurveDisplay.from_estimator(lr, xTest_Scaled, yTest, pos_label = 0)
plt.show()
```



```
In [39]: from sklearn.metrics import PrecisionRecallDisplay
PrecisionRecallDisplay.from_estimator(lr, xTest_Scaled, yTest, pos_label = 1)
plt.show()
```



```
In [39]: from sklearn.metrics import PrecisionRecallDisplay
PrecisionRecallDisplay.from_estimator(lr, xTest_Scaled, yTest, pos_label = 0)
plt.show()
```


Comments on Results

- Our data had some variables with extremely high multicollinearity, so to reduce it, I have used l2 regularization which will deal with the multicollinearity by constraining the coefficients and by keeping all the variables
- The most important features in our prediction were **pincode** and **sub_grade**. I had noticed this clearly when we performed bi-variate feature analysis
- We can see that while the overall accuracy of our model is good, there exists a clear imbalance in the results achieved for repayed loans and defaulted loans with the former having better results overall (high F1 score). This means our model is **not performing well in detecting the defaulters (class '0')**. This phenomenon might have multiple reasons but as we had already seen that the classes in our data were **heavily imbalanced**, this makes the model focus more on a specific class. To solve this, **techniques like undersampling, oversampling and SMOTE** can be used.
- We can see through our classification metrics that the specificity, the **Recall of target variable is quite bad** for The recall is 0.48, which means that the model is able to identify 48% of the individuals who actually defaulted on their loans. Whereas, since the precision is 0.92, which means that when the model predicts an individual to default on their loan, it is correct 92% of the time. Hence, our **False Negatives for class '0' are quite high, so our model is predicting more people as defaulters than they actually are**
- The AUC of 0.91 suggests that the model has a **high level of discrimination power** in correctly ranking the probabilities of both defaulting individuals (class 0) and individuals who paid their loans back on time (class 1). It indicates that the model is generally effective in distinguishing between the two classes.

Tradeoff Questions

- How can we make sure that our model can detect real defaulters and there are less false positives? This is important as we can lose out on an opportunity to finance more individuals and earn interest on it?
- How **do we explain we can use undersampling and oversampling techniques for the same. We can also play with the threshold values to get the best threshold using which helps us reduce FN**
- Since NPA (non-performing asset) is a real problem in this industry, it's important we play safe and shouldn't disburse loans to anyone.
- This is true, our model is also giving high precision for defaulters which means it's doing well to save LoanTap money by forewarning the lender before default**

Miscellaneous : Ignore

Questionnaire (Answers should present in the text editor along with insights):

- What percentage of customers have fully paid their Loan Amount?
- Comment about the correlation between Loan Amount and Installment features.
- The majority of people have home ownership as ____.
- People with grades 'A' are more likely to fully pay their loan. (T/F)
- Name the top 2 afforded job titles.
- Thinking from a bank's perspective, which metric should our primary focus be on.
- ROC AUC
- Precision
- Recall
- F1 Score
- How does the gap in precision and recall affect the bank?
- Which were the features that heavily affected the outcome?
- Will the results be affected by geographical location? (Yes/No)

- should we add_constant while checking V.I.F, if constant is high then what to do?
- https://stackoverflow.com/questions/596442/what-does-the-high-vif-for-the-constant-term-intercept-indicate
- <https://stats.stackexchange.com/questions/332428/regression-model-constant-causes-multicollinearity-warning-but-not-in-standards>

- PCA + Decision boundary

```
In [ ]: locs = [df.columns.get_loc(c) for c in col_list if c in df]
```

```
In [34]: temp_df = add_constant(encoded_data)
pd.DataFrame(data = [variance_inflation_factor(temp_df, i) for i in range(0, temp_df.shape[1])],
    index = temp_df.columns.tolist(),
    columns = ['VIF'])
```



```
Out[345]:
```

		VIF
	const	134564.130215
	loan_amnt	59.166840
	term	6.688299
	int_rate	21.833780
	installment	51.296394
	sub_grade	21.641137
	annual_inc	1.1715397
	verification_status	1.1745971
	dti	1.424990
	open_acc	2.140934
	pub_rec	1.413320
	revol_bal	1.556564
	revol_util	1.354765
	total_acc	2.193000
	mort_acc	1.494578
	pub_rec_bankruptcies	1.403819
	emp_length_10+ years	4.290924
	emp_length_2 years	2.169255
	emp_length_3 years	2.047375
	emp_length_4 years	1.811337
	emp_length_5 years	1.893157
	emp_length_6 years	1.716832
	emp_length_7 years	1.717539
	emp_length_8 years	1.664160
	emp_length_9 years	1.537277
	emp_length < 1 year	2.046791
	home_ownership_MORTGAGE	33007.202027
	home_ownership_NONE	11.334034
	home_ownership_OTHER	38.327813
	home_ownership_OWN	11385.142601
	home_ownership_RENT	31777.383815
	initial_list_status_w	1.0771600
	application_type_INDIVIDUAL	2.500057
	application_type_JOINT	2.491329
	pincode_5113	1.762635
	pincode_11650	1.516016
	pincode_22690	1.965937
	pincode_29597	1.763699
	pincode_30723	1.966817
	pincode_48052	1.960869
	pincode_70466	1.973256
	pincode_86630	1.505101
	pincode_93700	1.512642
	loan_status	1.850755

```
In [337]: from statsmodels.stats.outliers_influence import variance_inflation_factor
from statsmodels.tools.tools import add_constant

pd.DataFrame(data = [variance_inflation_factor(encoded_data, i) for i in range(0, encoded_data.shape[1])],
              index = encoded_data.columns.tolist(),
              columns = ["VIF{}".format(i)]
)
```

```
Out[337]:
```

		VIF
	loan_amnt	227.841468
	term	118.036616
	int_rate	223.367307
	installment	203.417954
	sub_grade	303.407990
	annual_inc	6.643866
	verification_status	3.067987
	dti	7.942459
	open_acc	12.769271
	pub_rec	1.442525
	revol_bal	2.939612
	revol_util	7.942673
	total_acc	12.399727
	mort_acc	2.487065
	pub_rec_bankruptcies	1.412108
	emp_length_10+ years	6.751358
	emp_length_2 years	2.384889
	emp_length_3 years	2.225101
	emp_length_4 years	1.927844
	emp_length_5 years	2.028771
	emp_length_6 years	1.812169
	emp_length_7 years	1.812804
	emp_length_8 years	1.748755
	emp_length_9 years	1.599071
	emp_length < 1 year	2.226961
	home_ownership_MORTGAGE	1227.366159
	home_ownership_NONE	1.194951
	home_ownership_OTHER	1.683642
	home_ownership_OWN	233.773335
	home_ownership_RENT	987.158600
	initial_list_status_w	1.762007
	application_type_INDIVIDUAL	1377.953528
	application_type_JOINT	2.478282
	pincode_5113	1.990758
	pincode_11650	1.560147
	pincode_22690	2.293134
	pincode_29597	1.992403
	pincode_30723	2.294221
	pincode_48052	2.283125
	pincode_70466	2.304719
	pincode_86630	1.547945
	pincode_93700	1.556365
	loan_status	9.435835

```
In [ ]: # Binning & converting employment length column to integer
# data['emp_length'] = data['emp_length'].apply(lambda yrs: "0 years" if yrs == "0 years" if yrs == "< 1 year" else ("10 years" if
# data['emp_length'] = data['emp_length'].str.split(" ").str[0].astype('float64')
# data['emp_length'] = data['emp_length'].apply(lambda yrs: 0 if yrs < 10 else 1)

In [ ]: # Ordinal encoding verification_status :
# 'Verified' --> verified by LoanTap, 'Source verified' --> income source was verified
```