

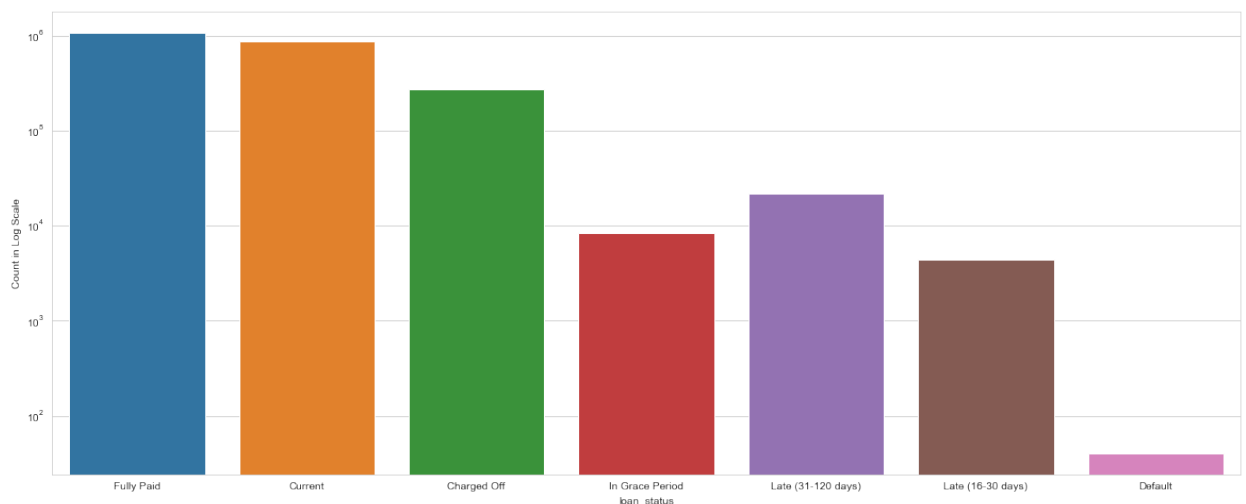
```
In [1]: 1 import pandas as pd
        2 import matplotlib.pyplot as plt
        3 import warnings
        4 warnings.filterwarnings("ignore")
        5 from myEda import EDA #my eda python file
        6
        7 %matplotlib inline
        8
        9 acc_filepath = './archive/accepted_2007_to_2018q4.csv/'
       10 rej_filepath = './archive/rejected_2007_to_2018q4.csv/'
       11 data = pd.read_csv(acc_filepath+"accepted_2007_to_2018q4.csv")
       12 eda = EDA(data)
       13 eda.preprocessing()
```

## Note:

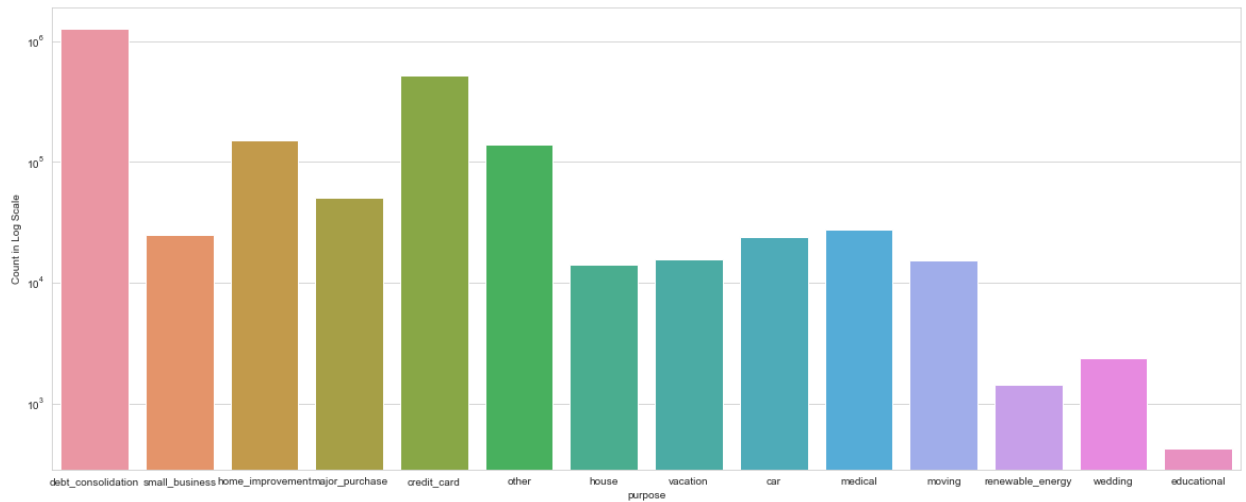
- If the maps are not showing in IPYNB then please check the pdf. The notebook needs to be executed for the maps plots to show.
- All observations are summarized at the end of the notebook.

## Accepted CSV EDA

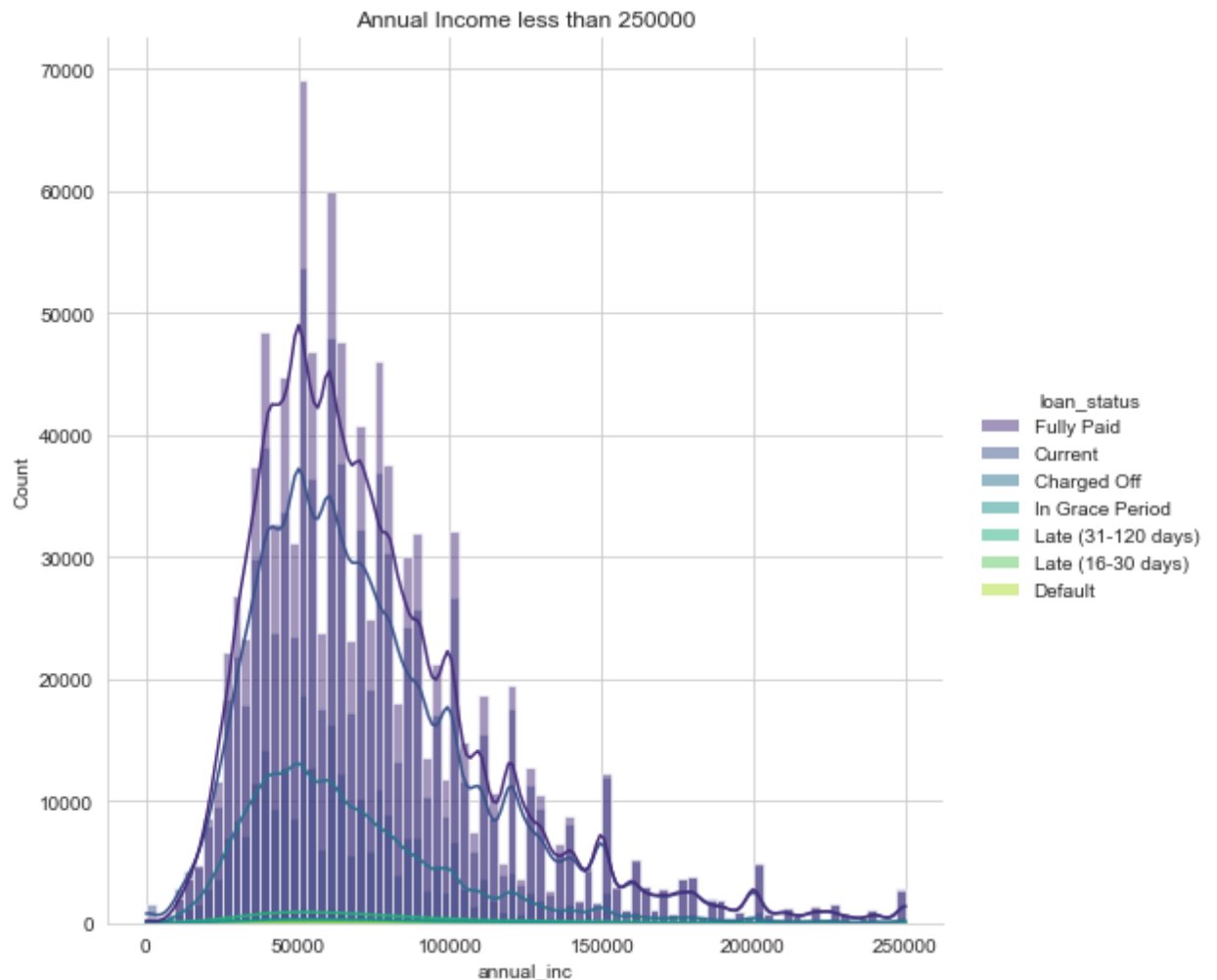
```
In [2]: 1 eda.logcountplot(x="loan_status")
```



```
In [3]: 1 eda.logcountplot("purpose")
```

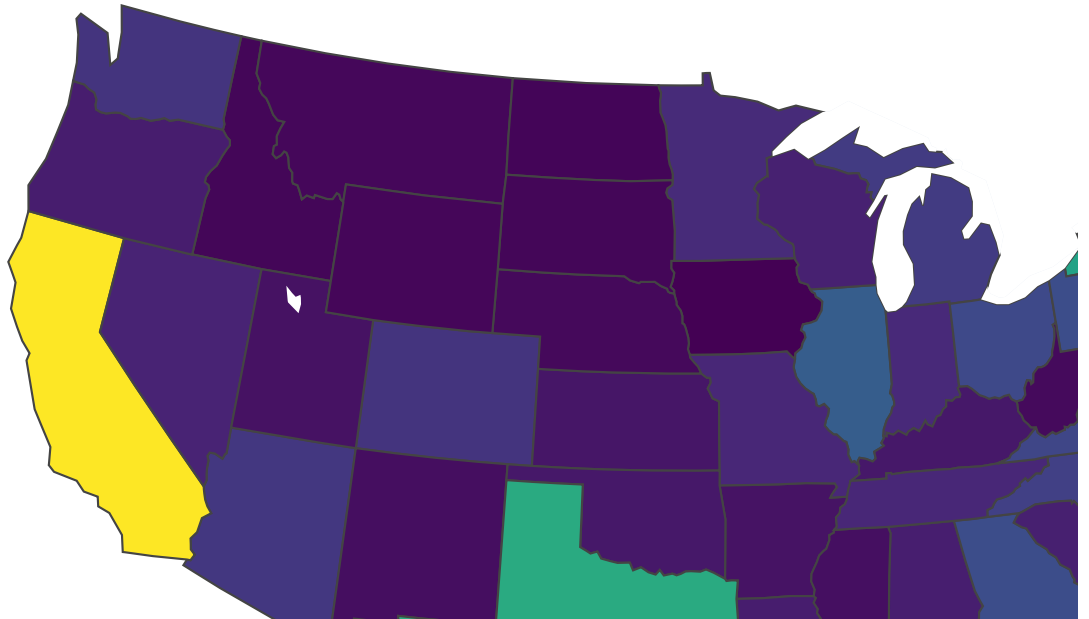


```
In [4]: 1 eda.incomedist()
```



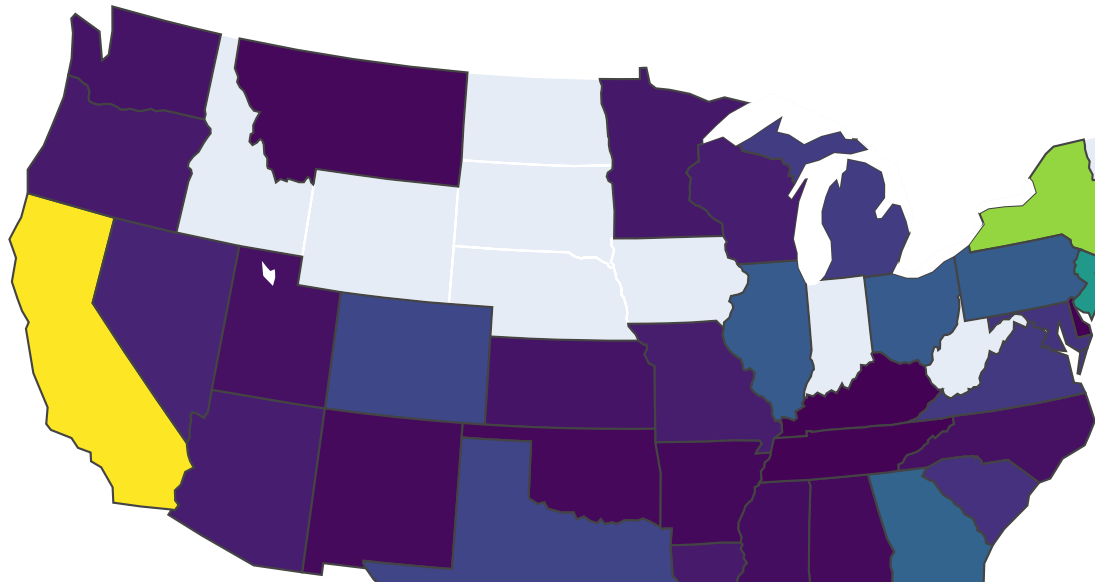
## Total Loan

```
In [5]: 1 eda.mapdist_loan_amt()
```



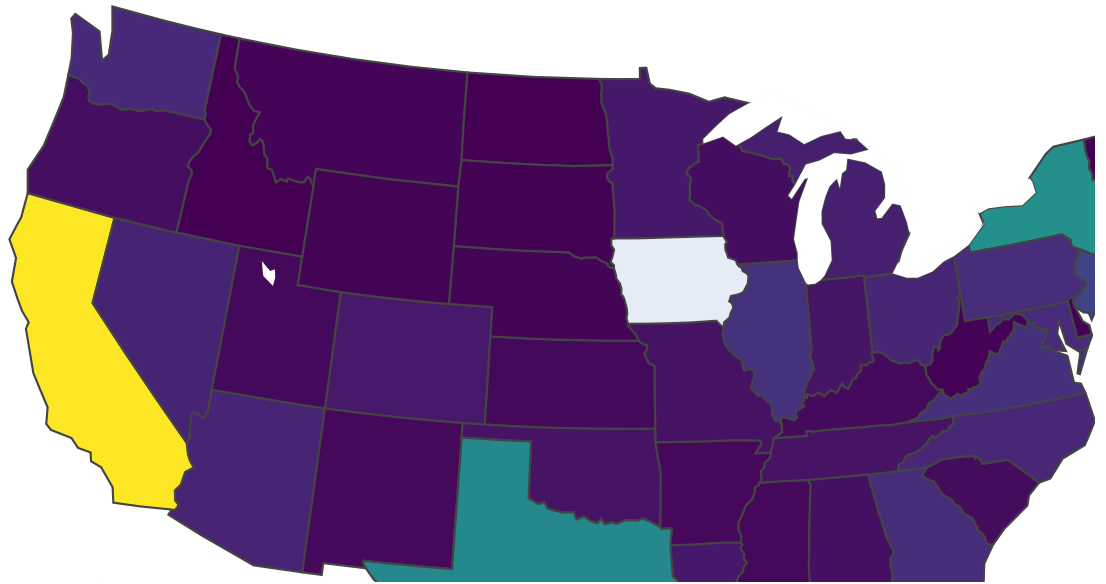
## Educational Loan

```
In [6]: 1 eda.mapdist_pur_loan_amt("educational")
```



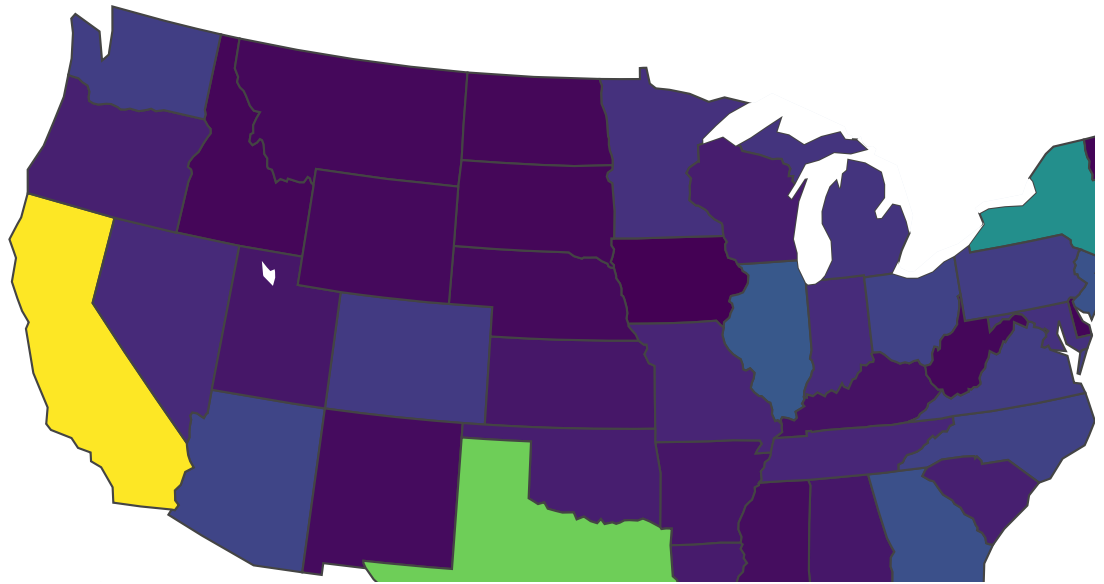
## Vacation Loan

```
In [7]: 1 eda.mapdist_pur_loan_amt("vacation")
```



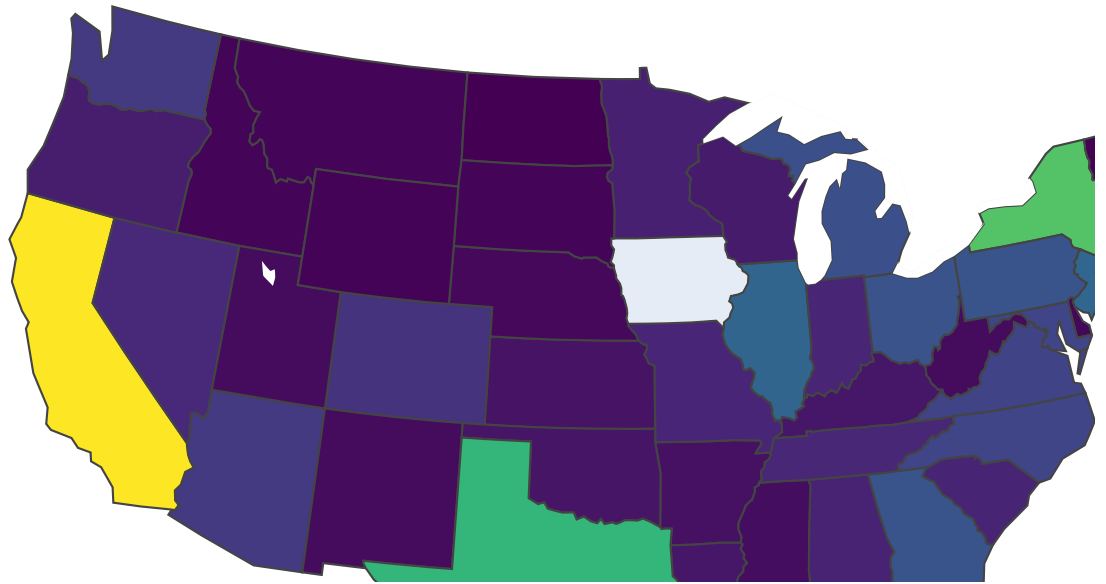
## Medical Loan

```
In [8]: 1 eda.mapdist_pur_loan_amt("medical")
```

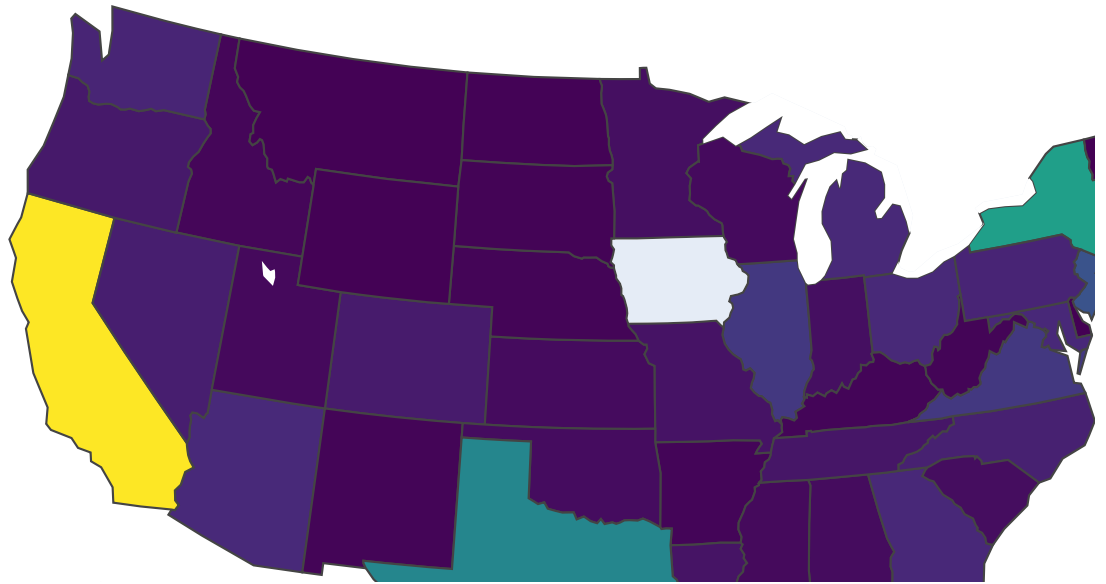


## House Loan

```
In [9]: 1 eda.mapdist_pur_loan_amt("house")
```



```
In [10]: 1 eda.mapdist_pur_loan_amt("renewable_energy")
```

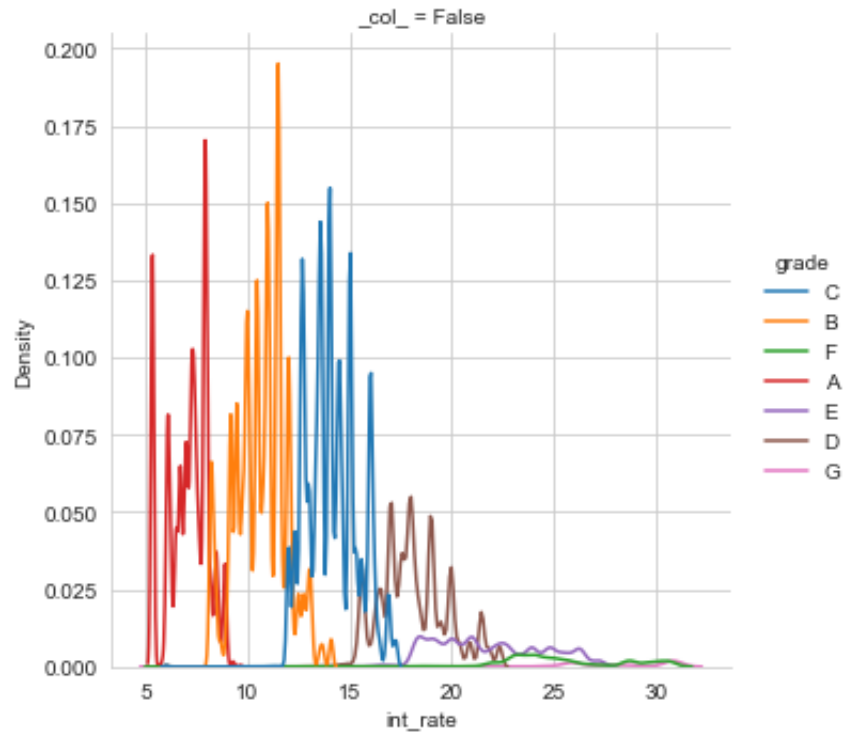


## Loan Grade observations



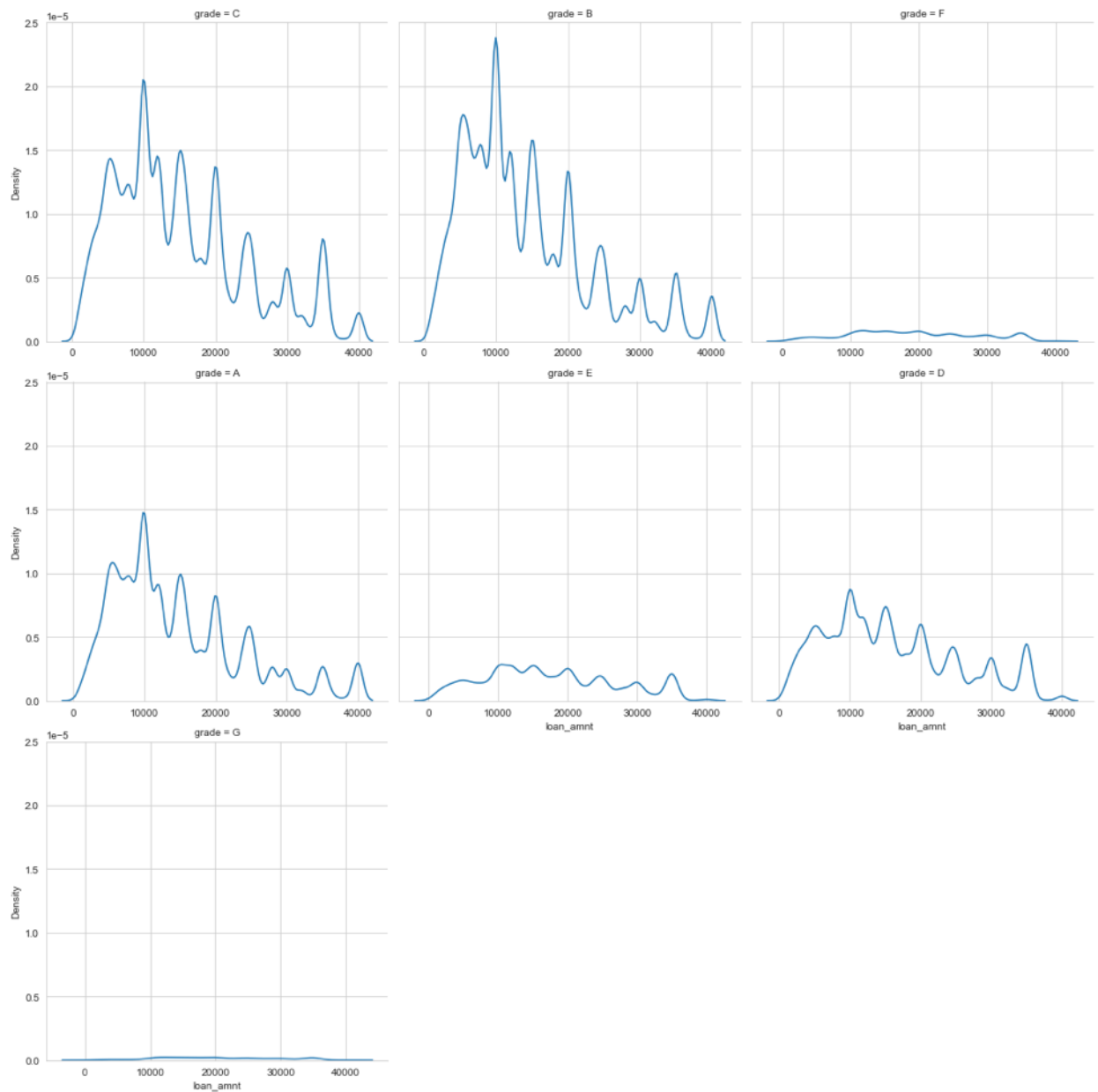
```
In [11]: 1 eda.simpLEDist(x='int_rate', hue ="grade")
```

<Figure size 1440x576 with 0 Axes>



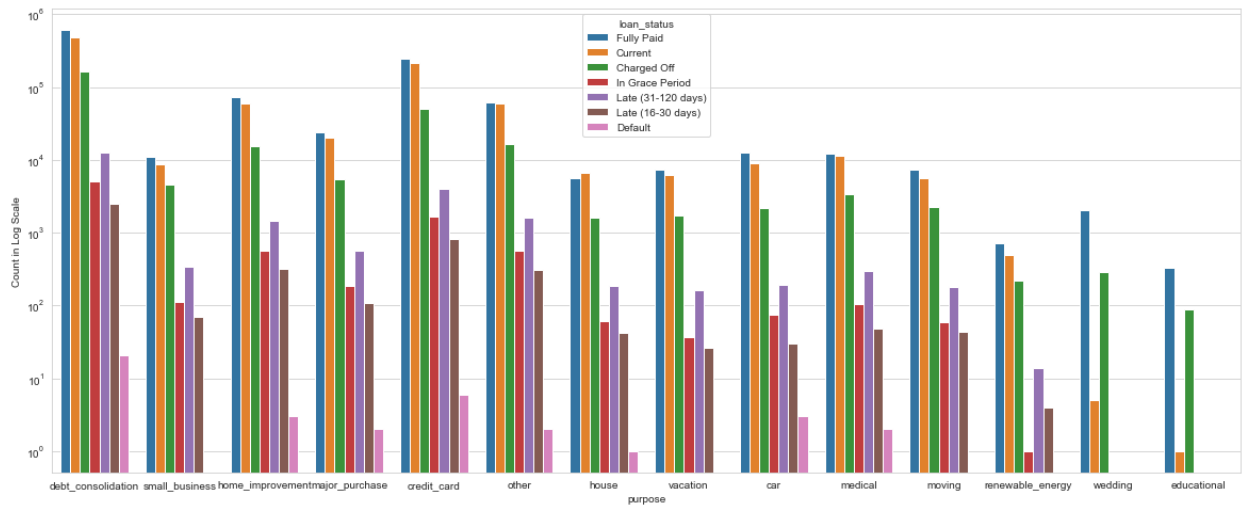
```
In [12]: 1 eda.comparedist(x='loan_amnt', hue =None,col="grade")
```

<Figure size 1440x576 with 0 Axes>

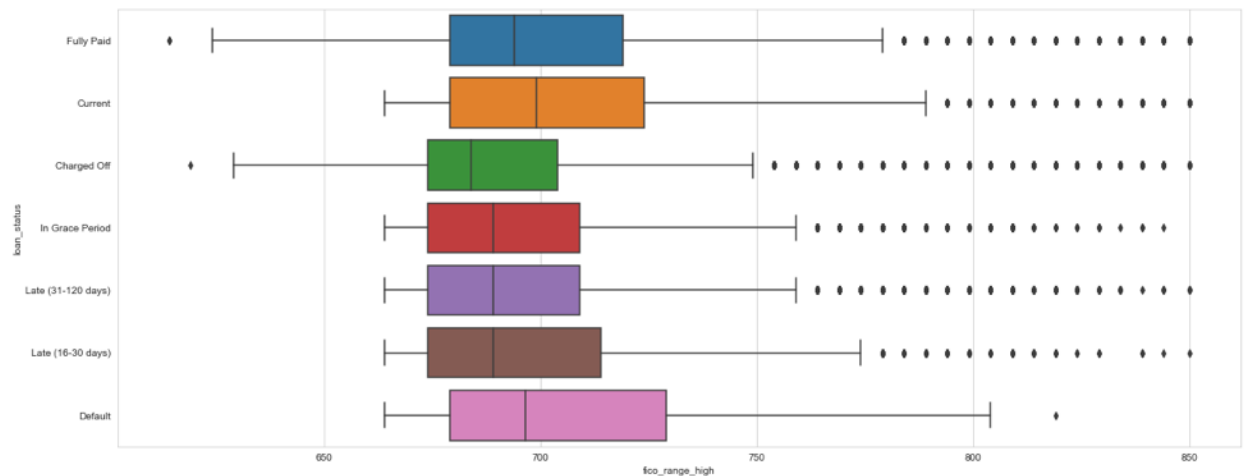


## Loan Status/ FICO Observations

In [13]: `eda.logcountplot(x = "purpose", hue = "loan_status")`

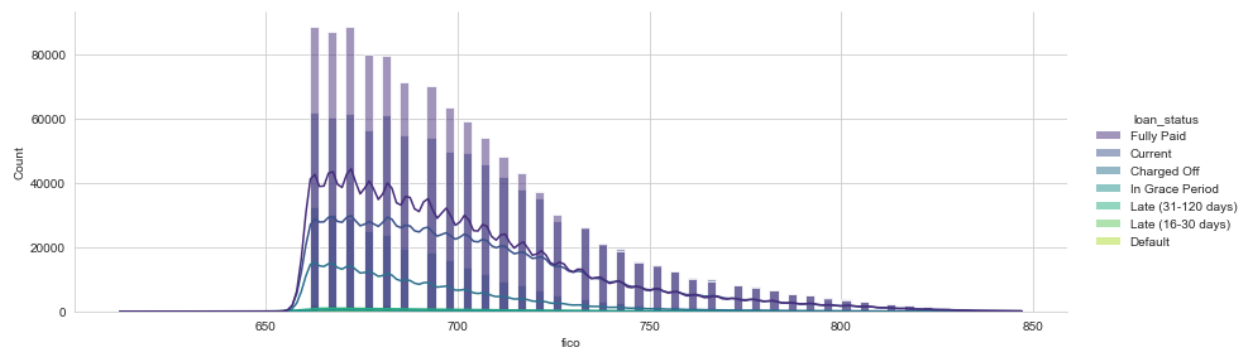


In [14]: `eda.boxplots(y='loan_status', x='fico_range_high')`



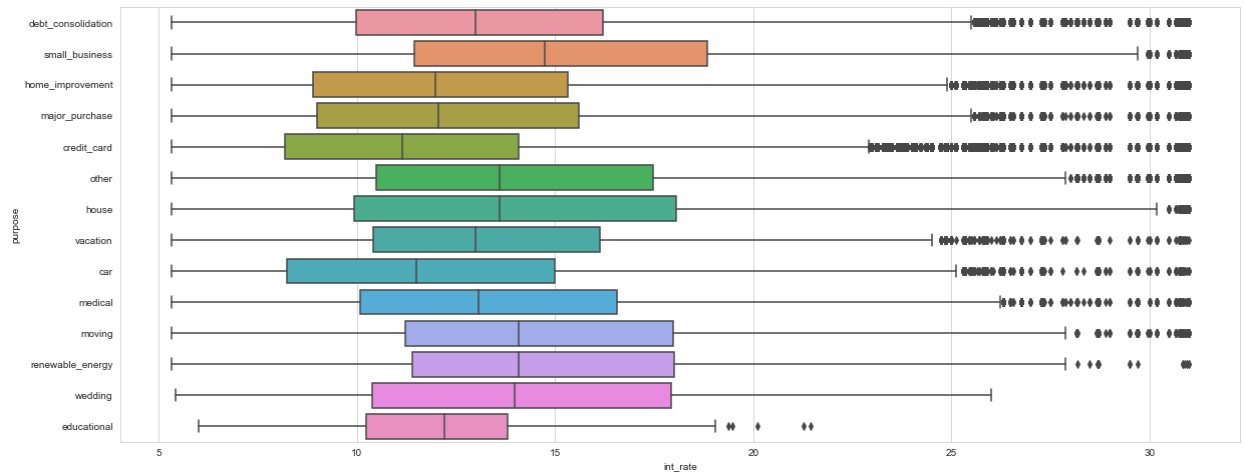
In [15]: `eda.fico_loan()`

<Figure size 6000x2100 with 0 Axes>

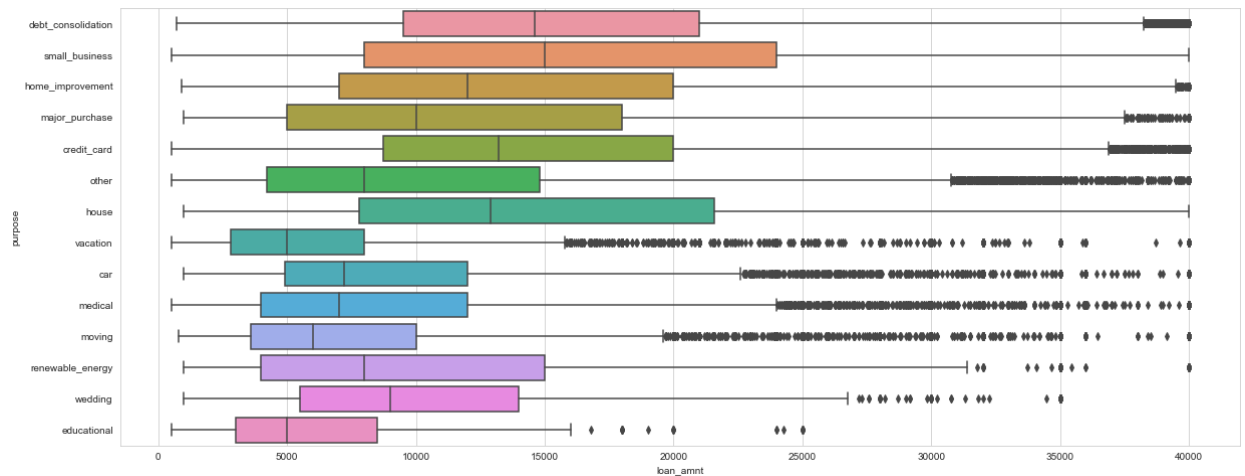


## Loan Purpose Observations

In [16]: `eda.bboxplots(x='int_rate',y='purpose')`

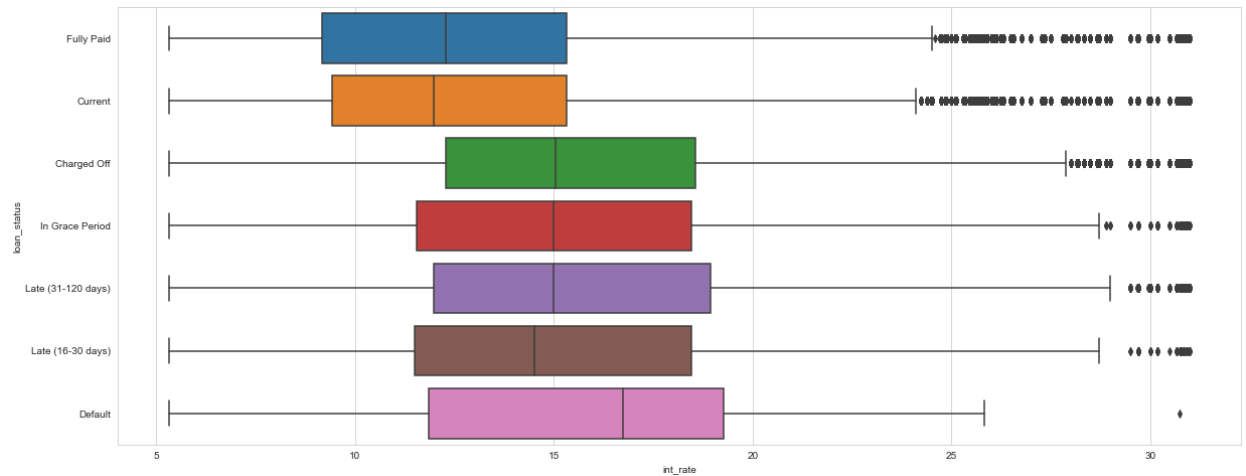


In [17]: `eda.bboxplots(x='loan_amnt',y='purpose')`

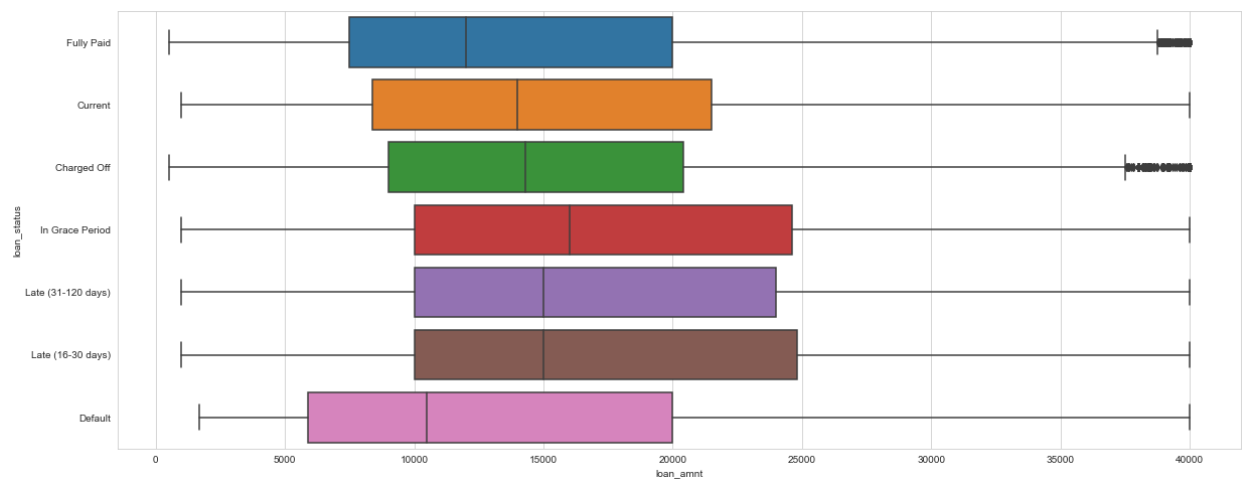


## Loan Status Observations

```
In [18]: 1 eda.boxplots(x='int_rate',y='loan_status')
```

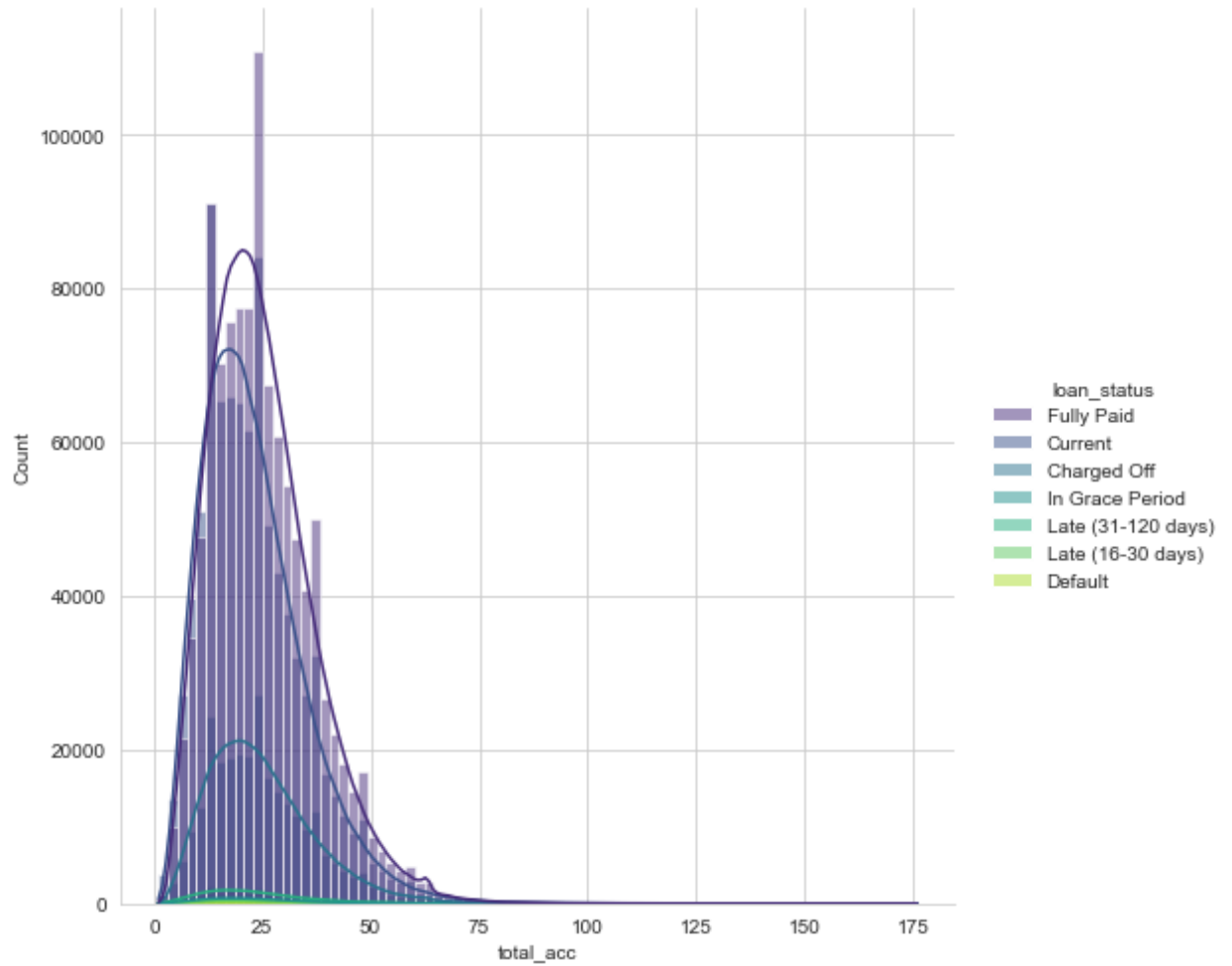


```
In [19]: 1 eda.boxplots(x='loan_amnt',y='loan_status')
```



**Total no. of credit lines**

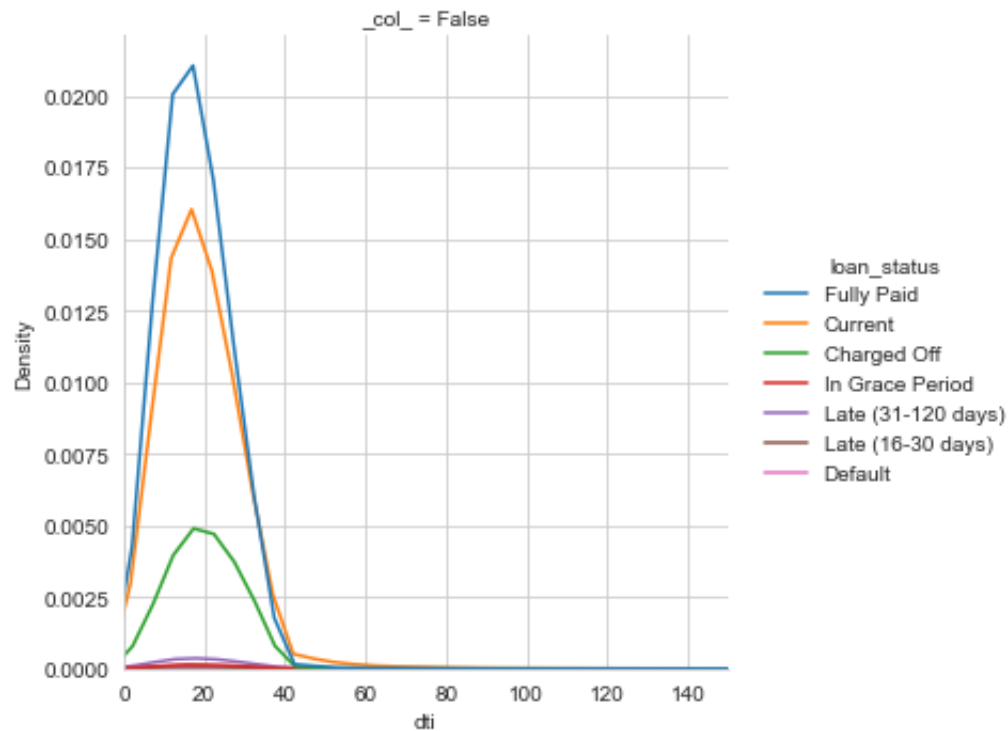
```
In [20]: 1 eda.comparehist(x='total_acc', hue='loan_status')
```



## Debt to Income Ratio

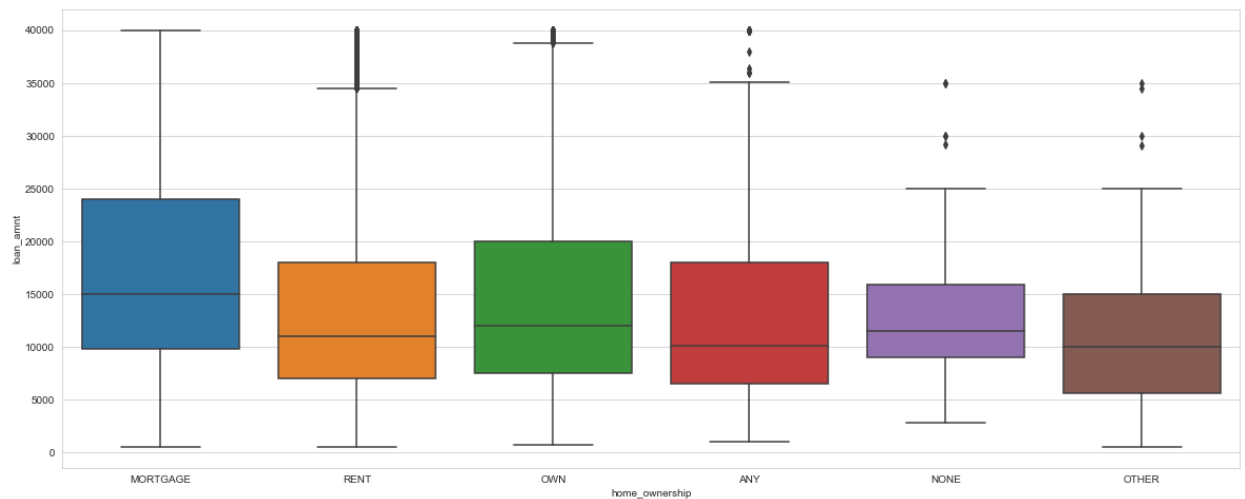
```
In [21]: 1 eda.dti_loanst(x='dti', hue = "loan_status")
```

<Figure size 1440x576 with 0 Axes>



## Home Ownership vs Loan Amount

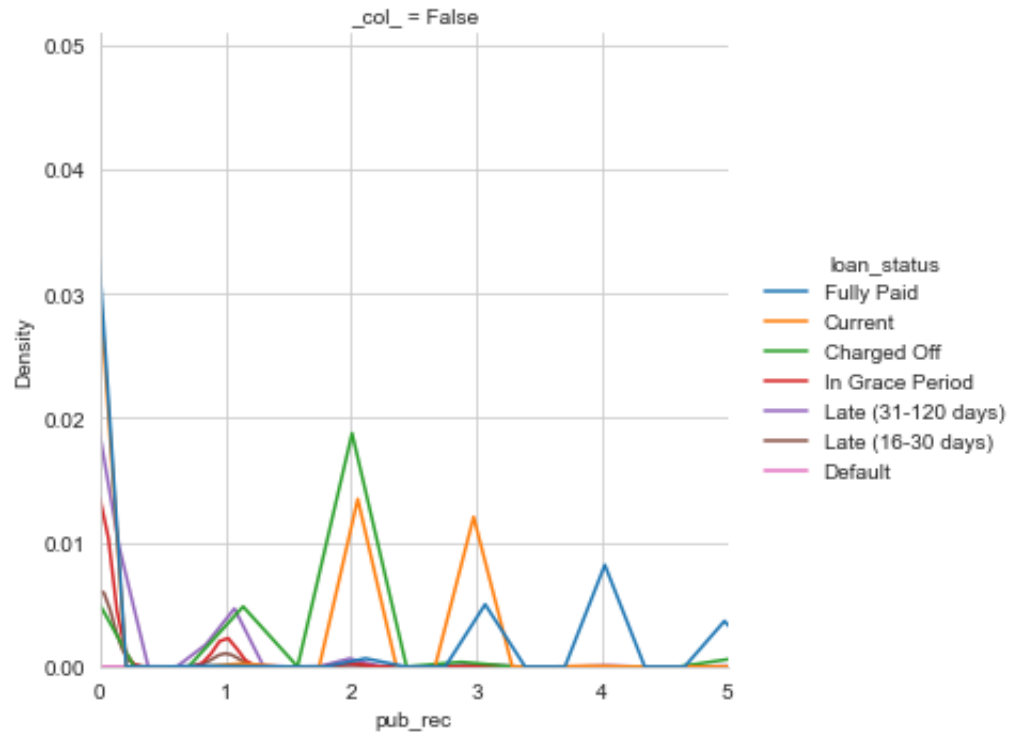
```
In [22]: 1 eda.boxplots(y='loan_amnt', x='home_ownership')
```



## Public Derogatory Records vs Loan Status

```
In [23]: 1 eda.pubrec_simpdist(x='pub_rec',hue='loan_status')
```

<Figure size 1440x576 with 0 Axes>

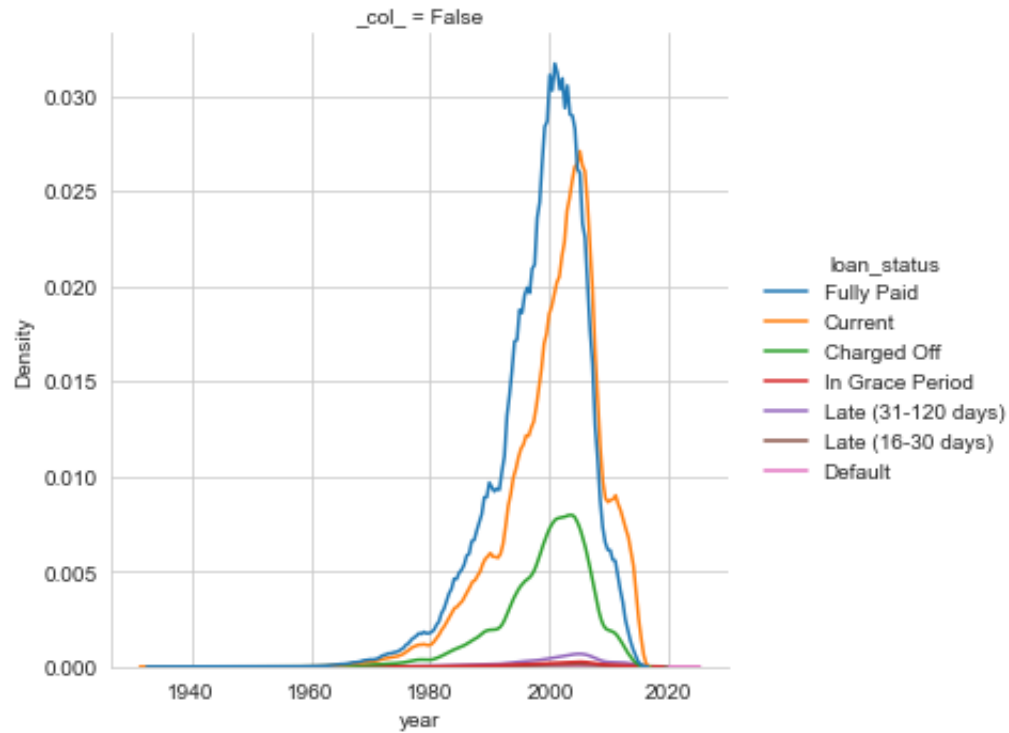


## Earliest Credit Line Opening year vs Loan Status



```
In [24]: 1 eda.datedist(x='earliest_cr_line',hue='loan_status')
```

<Figure size 1440x576 with 0 Axes>



## Observations:

1. Most of the loans were taken in the states of California, Texas, New York and Florida
2. People that paid the loan fully have higher incomes than that of people with loan status: current, charged off, late or default
3. Not just total loan, almost all choropleth maps show that the state of California has the highest loan amount borrowed. It may be the case of biased data sampling too.
4. Also, it is interesting to note that there are a few states without any educational loan taken. Again a possibility of data sampling bias.
5. Renewable energy loans are higher in the outer regions (CA, TX, NY etc) than the inner ones (NE, SD, WY).
6. As the grade increases alphabetically (starting from A), the interest rates go high.
7. These high interest rates lead to very few loans taken as evidenced by the low density of loan amounts borrowed in grades E and above.
8. When the mean of high and low of fico score is considered, the fico score of people who paid fully is higher than the ones that didn't pay yet.
9. Although the interest rates of "Medical" and "House" loans are similar, the loan amount borrowed for houses is far more.
10. Educational Loans have one of the lowest interest rates contrary to popular belief.
11. Although the interest rates for Fully Paid and Current Loans are similar, it appears that lower loan amounts have been paid back but the current, charged off and late ones are large amounts.
12. The total number of credit lines currently in the borrower's credit file (total\_acc) shows that mean total\_acc is higher for people who fully paid off the loan.
13. The debt-to-income (dti) ratio (total monthly debt payments divided by reported monthly income) is highest for people who fully paid the loan.
14. People on mortgage borrow significantly higher amount
15. People who paid off the loan fully have an earlier (mean) credit line opening date.