

# Keras API Project Exercise - Solutions

## The Data

We will be using a subset of the LendingClub DataSet obtained from Kaggle: <https://www.kaggle.com/wordsforthewise/lending-club>

**NOTE: Do not download the full zip from the link! We provide a special version of this file that has some extra feature engineering for you to do. You won't be able to follow along with the original file!**

LendingClub is a US peer-to-peer lending company, headquartered in San Francisco, California.[3] It was the first peer-to-peer lender to register its offerings as securities with the Securities and Exchange Commission (SEC), and to offer loan trading on a secondary market. LendingClub is the world's largest peer-to-peer lending platform.

## Our Goal

Given historical data on loans given with information on whether or not the borrower defaulted (charge-off), can we build a model that can predict whether or not a borrower will pay back their loan? This way in the future when we get a new potential customer we can assess whether or not they are likely to pay back the loan. Keep in mind classification metrics when evaluating the performance of your model!

The 'loan\_status' column contains our label.

## Data Overview

There are many LendingClub data sets on Kaggle. Here is the information on this particular data set:

LoanStatNew	Description
0 loan_amnt	The listed amount of the loan applied for by the borrower. If at some point in time, the credit department reduces the loan amount, then it will be reflected in this value.
1 term	The number of payments on the loan. Values are in months and can be either 36 or 60.
2 int_rate	Interest Rate on the loan
3 installment	The monthly payment owed by the borrower if the loan originates.
4 grade	LC assigned loan grade
5 sub_grade	LC assigned loan subgrade
6 emp_title	The job title supplied by the Borrower when applying for the loan.*
7 emp_length	Employment length in years. Possible values are between 0 and 10 where 0 means less than one year and 10 means ten or more years.
8 home_ownership	The home ownership status provided by the borrower during registration or obtained from the credit report. Our values are RENT, OWN, MORTGAGE, OTHER
9 annual_inc	The self-reported annual income provided by the borrower during registration.
10 verification_status	Indicates if income was verified by LC, not verified, or if the income source was verified
11 issue_d	The month which the loan was funded
12 loan_status	Current status of the loan
13 purpose	A category provided by the borrower for the loan request.
14 title	The loan title provided by the borrower
15 zip_code	The first 3 numbers of the zip code provided by the borrower in the loan application.
16 addr_state	The state provided by the borrower in the loan application
17 dti	A ratio calculated using the borrower's total monthly debt payments on the total debt obligations, excluding mortgage and the requested LC loan, divided by the borrower's self-reported monthly income.
18 earliest_cr_line	The month the borrower's earliest reported credit line was opened
19 open_acc	The number of open credit lines in the borrower's credit file.
20 pub_rec	Number of derogatory public records
21 revol_bal	Total credit revolving balance
22 revol_util	Revolving line utilization rate, or the amount of credit the borrower is using relative to all available revolving credit.
23 total_acc	The total number of credit lines currently in the borrower's credit file
24 initial_list_status	The initial listing status of the loan. Possible values are – W, F
25 application_type	Indicates whether the loan is an individual application or a joint application with two co-borrowers
26 mort_acc	Number of mortgage accounts.
27 pub_rec_bankruptcies	Number of public record bankruptcies

## Starter Code

**Note: We also provide feature information on the data as a .csv file for easy lookup throughout the notebook:**

```
In [1]: import pandas as pd
# KaggleIndex: 396030 entries, 0 to 396029
data_info = pd.read_csv('../DATA/lending_club_info.csv', index_col='LoanStatNew')

In [2]: print(data_info.loc['revol_util']['Description'])
Revolving line utilization rate, or the amount of credit the borrower is using relative to all available revolving credit.

In [3]: def feat_info(col_name):
    print(data_info.loc[col_name]['Description'])

In [4]: feat_info('mort_acc')
Number of mortgage accounts.
```

## Loading the data and other imports

```
In [6]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
# might be needed depending on your version of Jupyter
%matplotlib inline

In [7]: df = pd.read_csv('../DATA/lending_club_loan_two.csv')

In [8]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 396030 entries, 0 to 396029
Data columns (total 27 columns):
loan_amnt      396030 non-null float64
term           396030 non-null object
int_rate       396030 non-null float64
installment    396030 non-null float64
grade          396030 non-null object
sub_grade      396030 non-null object
emp_title      373103 non-null object
emp_length     317729 non-null object
home_ownership 396030 non-null object
annual_inc     396030 non-null float64
verification_status 396030 non-null object
issue_d        396030 non-null object
loan_status    396030 non-null object
purpose        396030 non-null object
title          394275 non-null object
dti            396030 non-null float64
earliest_cr_line 396030 non-null object
open_acc       396030 non-null float64
pub_rec        396030 non-null float64
revol_bal      396030 non-null float64
revol_util     395754 non-null float64
total_acc      396030 non-null float64
initial_list_status 396030 non-null object
application_type 396030 non-null object
mort_acc       358235 non-null float64
pub_rec_bankruptcies 394545 non-null float64
address        396030 non-null object
dtypes: float64(12), object(15)
memory usage: 81.6+ MB
```

## Project Tasks

Complete the tasks below! Keep in mind is usually more than one way to complete the task! Enjoy

## Section 1: Exploratory Data Analysis

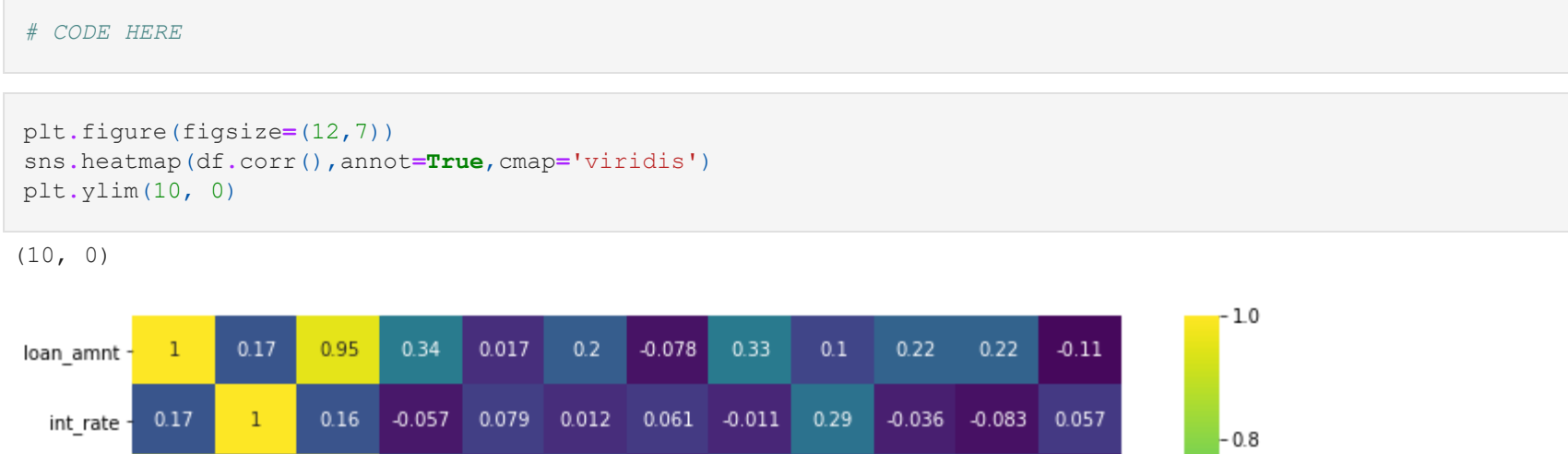
**OVERALL GOAL:** Get an understanding for which variables are important, view summary statistics, and visualize the data

**TASK:** Since we will be attempting to predict loan\_status, create a countplot as shown below.

```
In [9]: # CODE HERE
```

```
sns.countplot(x='loan_status', data=df)
```

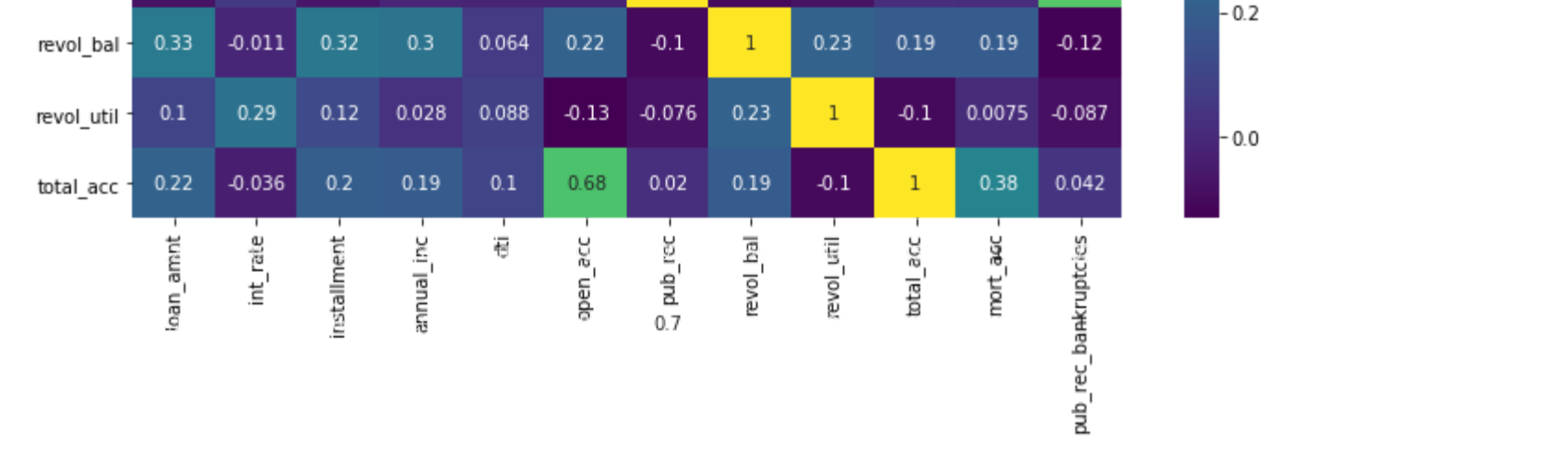
```
Out[10]: <matplotlib.axes._subplots.AxesSubplot at 0x207932022c8>
```



**TASK:** Create a histogram of the loan\_amnt column.

```
In [11]: # CODE HERE
```

```
In [12]: plt.figure(figsize=(12,4))
sns.distplot(df['loan_amnt'], kde=False, bins=40)
plt.xlim(0,45000)
```



**TASK:** Let's explore correlation between the continuous feature variables. Calculate the correlation between all continuous numeric variables using .corr() method.

```
In [13]: # CODE HERE
```

```
In [14]: df.corr()
```

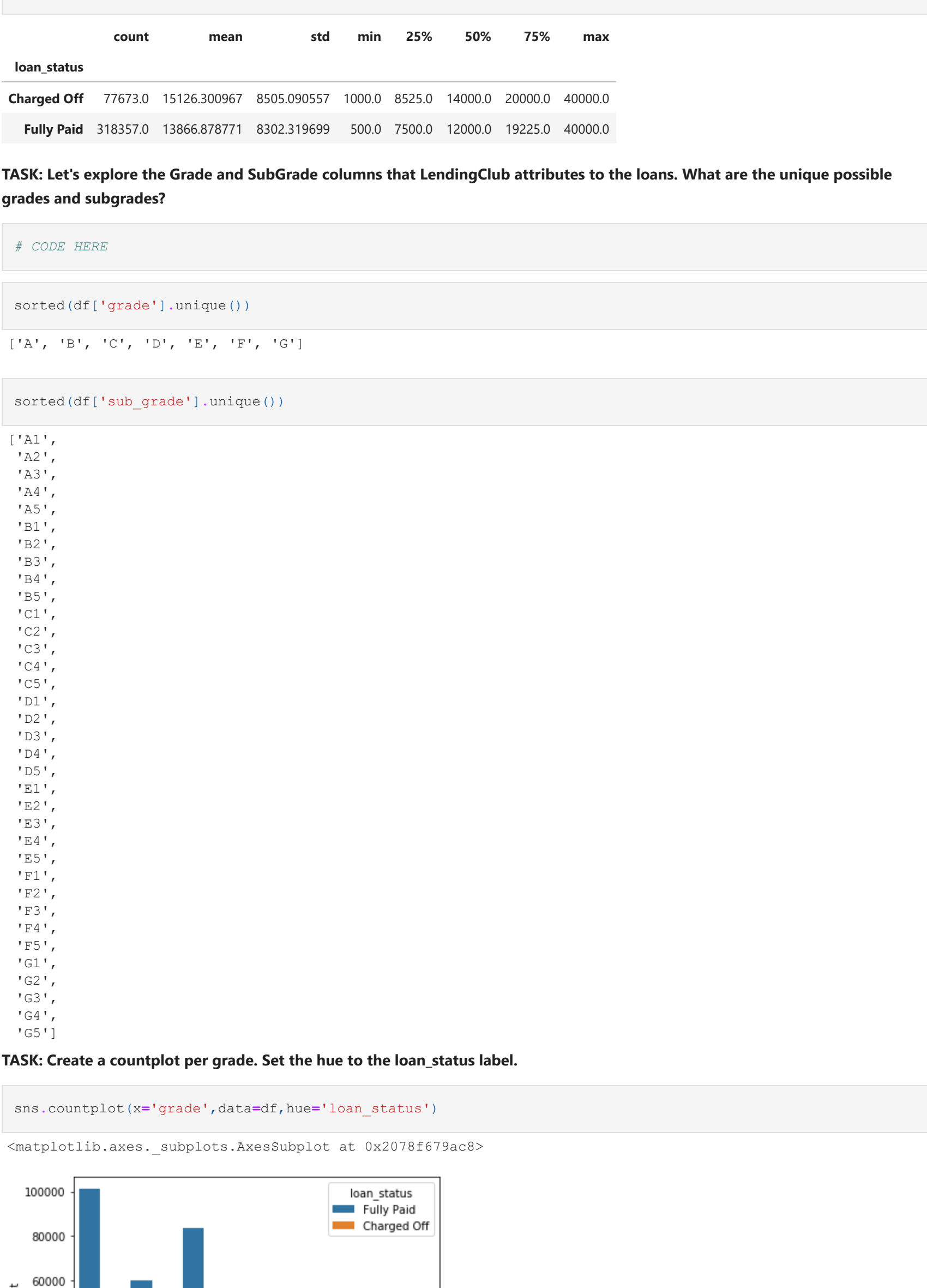
	loan_amnt	int_rate	installment	annual_inc	dti	open_acc	pub_rec	revol_bal	revol_util	total_acc	mort_acc
loan_amnt	1.000000	0.168921	0.953929	0.336887	0.016636	0.198556	-0.077779	0.328320	0.099911	0.223886	0.222315
int_rate	0.168921	1.000000	0.162758	-0.056771	0.079938	0.011649	0.060986	-0.011280	0.293659	-0.036404	-0.082583
installment	0.953929	0.162758	1.000000	0.330381	0.015786	0.188973	-0.067892	0.311645	0.123915	0.202430	0.193694
annual_inc	0.336887	-0.056771	0.330381	1.000000	-0.081685	0.136150	-0.013720	0.299773	0.027871	0.193023	0.236230
dti	0.016636	0.079038	0.015786	-0.081685	1.000000	0.136181	-0.017639	0.063571	0.088375	0.102128	-0.025439
open_acc	0.198556	0.011649	0.188973	0.136150	0.136181	1.000000	-0.018392	0.221192	-0.131420	0.680728	0.109205
pub_rec	-0.077779	0.060986	-0.067892	-0.013720	-0.017639	-0.018392	1.000000	-0.101664	-0.075910	0.019723	0.011552
revol_bal	0.328320	-0.011280	0.311645	0.299773	0.063571	0.221192	-0.101664	1.000000	0.226346	0.191616	0.194925
revol_util	0.099911	0.293659	0.123915	0.027871	0.088375	0.102128	0.680728	0.019723	1.000000	-0.104273	0.075114
total_acc	0.223886	-0.036404	0.202430	0.193023	0.102128	0.680728	0.019723	0.191616	-0.104273	1.000000	0.381072
mort_acc	0.222315	-0.082583	0.193694	0.236230	-0.025439	0.109205	0.011552	0.194925	0.007514	0.381072	1.000000
pub_rec_bankruptcies	-0.106539	-0.057450	-0.098628	-0.050162	-0.014558	-0.027732	0.699408	-0.124532	-0.086751	0.042035	0.002739

**TASK:** Visualize this using a heatmap. Depending on your version of matplotlib, you may need to manually adjust the heatmap.

- Heatmap info
- Help with resizing

```
In [15]: # CODE HERE
```

```
In [16]: plt.figure(figsize=(12,7))
sns.heatmap(df.corr(), annot=True, cmap='viridis')
plt.gcf().set_size_inches(10,10)
```



**TASK:** You should have noticed almost perfect correlation with the "installment" feature. Explore this feature further. Print out their descriptions and perform a scatterplot between them. Does this relationship make sense to you? Do you think there is duplicate information here?

```
In [17]: # CODE HERE
```

```
In [18]: feat_info('installment')
```

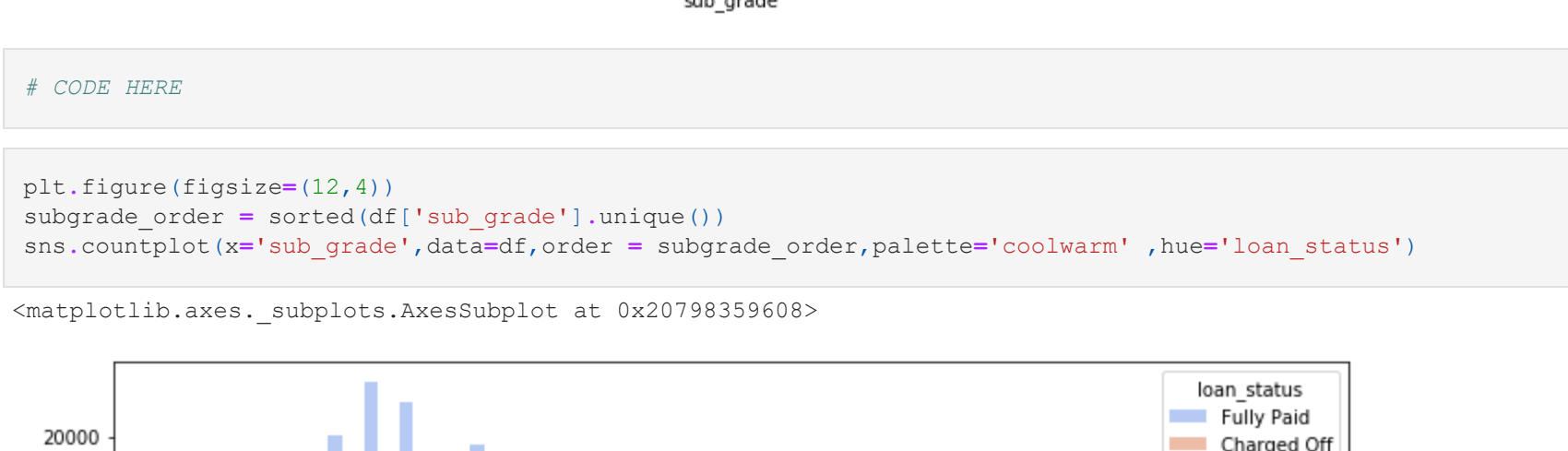
The monthly payment owed by the borrower if the loan originates.

```
In [19]: feat_info('loan_amnt')
```

The listed amount of the loan applied for by the borrower. If at some point in time, the credit department reduces the loan amount, then it will be reflected in this value.

```
In [20]: sns.scatterplot(x='installment', y='loan_amnt', data=df)
```

```
Out[20]: <matplotlib.axes._subplots.AxesSubplot at 0x20798026f48>
```

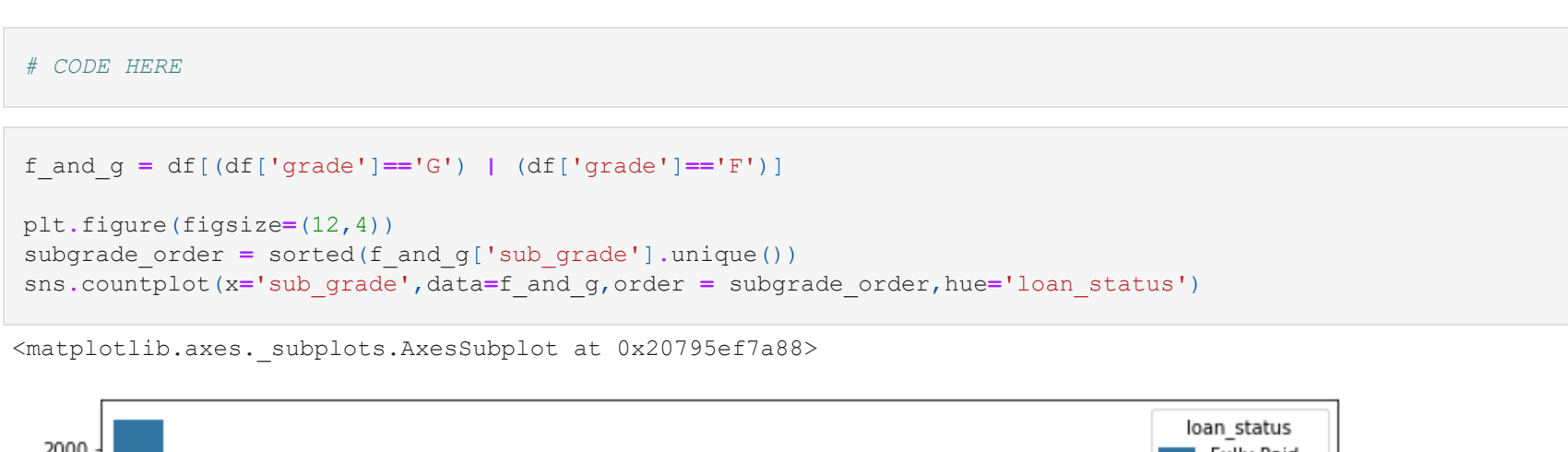


**TASK:** Create a boxplot showing the relationship between the loan\_status and the Loan Amount.

```
In [21]: # CODE HERE
```

```
In [22]: sns.boxplot(x='loan_status', y='loan_amnt', data=df)
```

```
Out[22]: <matplotlib.axes._subplots.AxesSubplot at 0x20798026f48>
```



**TASK:** Calculate the summary statistics for the loan amount, grouped by loan\_status.

```
In [23]: # CODE HERE
```

```
In [24]: df.groupby('loan_status')['loan_amnt'].describe()
```

	count	mean	std	min	25%	50%	75%	max
Charged Off	77673.0	15126.309967	8505.090557	1000.0	8525.0	14000.0	20000.0	40000.0
Fully Paid	318357.0	13866.878771	8302.319699	500.0	7500.0	12000.0	19225.0	40000.0

**TASK:** Let's explore the Grade and SubGrade columns that LendingClub attributes to the loans. What are the unique possible grades and subgrades?

```
In [25]: # CODE HERE
```

```
In [26]: sorted(df['grade'].unique())
```

```
Out[26]: ['A', 'B', 'C', 'D', 'E', 'F', 'G']
```

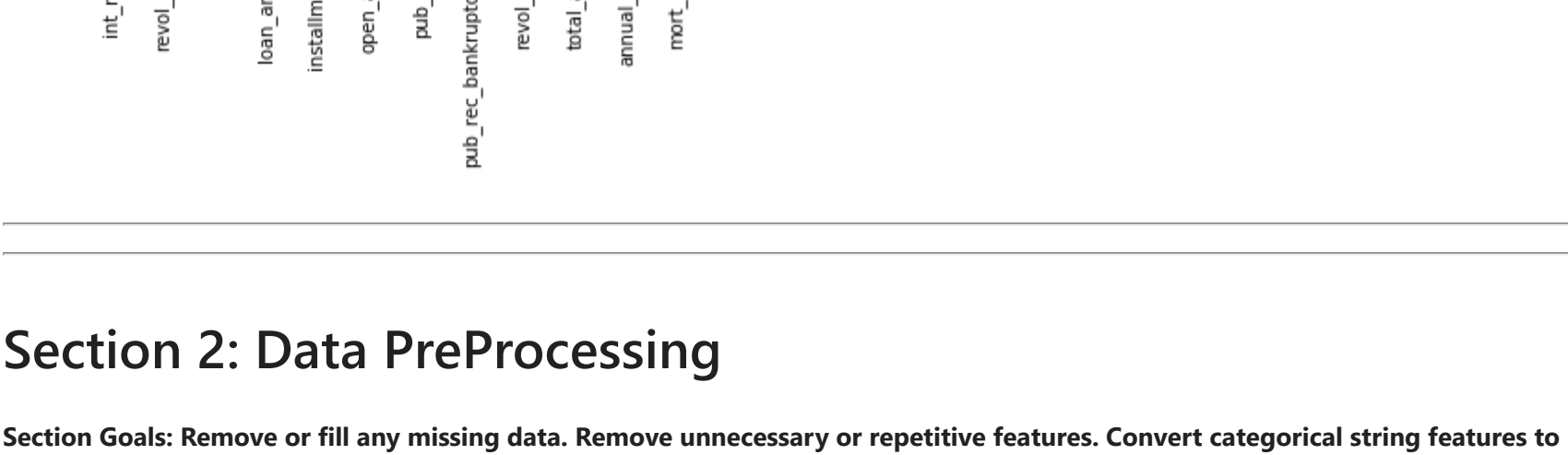
```
In [27]: sorted(df['sub_grade'].unique())
```

```
Out[27]: ['A1',
'A2',
'A3',
'A4',
'A5',
'B1',
'B2',
'B3',
'B4',
'B5',
'C1',
'C2',
'C3',
'C4',
'C5',
'D1',
'D2',
'D3',
'D4',
'D5',
'E1',
'E2',
'E3',
'E4',
'E5',
'E6',
'E7',
'E8',
'E9',
'F1',
'F2',
'F3',
'F4',
'F5',
'F6',
'F7',
'F8',
'F9',
'G1',
'G2',
'G3',
'G4',
'G5']
```

**TASK:** Create a countplot per grade. Set the hue to the loan status label.

```
In [28]: sns.countplot(x='grade', data=df, hue='loan_status')
```

```
Out[28]: <matplotlib.axes._subplots.AxesSubplot at 0x20798359ac8>
```

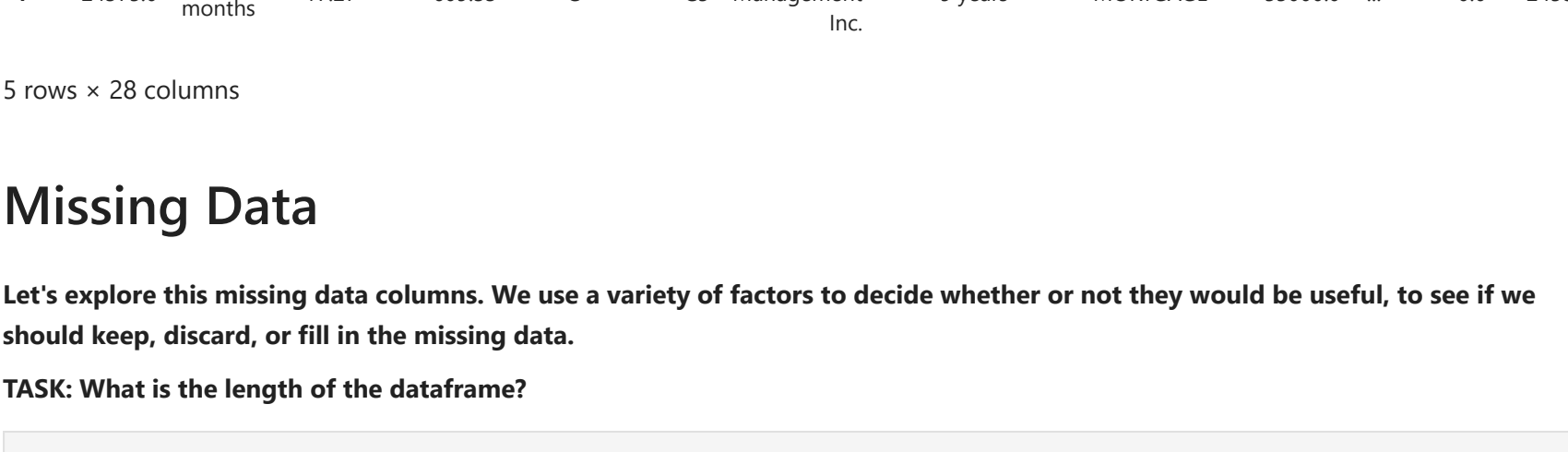


**TASK:** Display a count plot per subgrade. You may need to resize for this plot and reorder the x axis. Feel free to edit the color palette. Explore both all loans made per subgrade as well being separated based on the loan\_status

```
In [29]: #CODE HERE
```

```
In [30]: plt.figure(figsize=(12,4))
subgrade_order = sorted(df['sub_grade'].unique())
sns.countplot(x='sub_grade', data=df, order = subgrade_order, palette='coolwarm' )
```

```
Out[30]: <matplotlib.axes._subplots.AxesSubplot at 0x20798504288>
```



**TASK:** It looks like F and G subgrades don't get paid back that often. Isolate those and recreate the countplot just for those subgrades.

```
In [31]: # CODE HERE
```

```
In [34]: f_and_g = df[(df['grade'] == 'G') | (df['grade'] == 'F')]
plt.figure(figsize=(12,4))
subgrade_order = sorted(f_and_g['sub_grade'].unique())
sns.countplot(x='sub_grade', data=f_and_g, order = subgrade_order, hue='loan_status')
```

```
Out[34]: <matplotlib.axes._subplots.AxesSubplot at 0x20795e7fa8>
```



**TASK:** Create a new column called 'loan\_repaid' which will contain a 1 if the loan status was "Fully Paid" and a 0 if it was "Charged Off".

```
In [35]: # CODE HERE
```

```
In [36]: df['loan_status'].unique()
```

```
Out[36]: array(['Fully Paid', 'Charged Off'], dtype=object)
```

```
In [37]: df['loan_repaid'] = df['loan_status'].map({'Fully Paid':1, 'Charged Off':0})
```

```
In [38]: df[['loan_repaid', 'loan_status']]
```

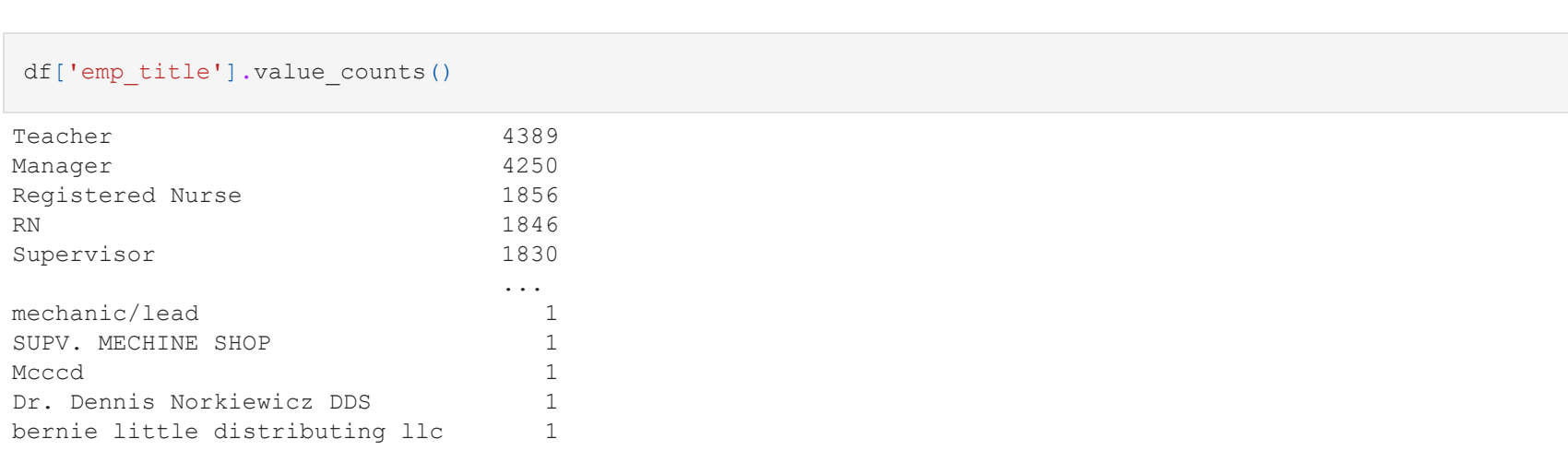
	loan_repaid	loan_status
0	1	Fully Paid
1	1	Fully Paid
2	1	Fully Paid
3	1	Fully Paid
4	0	Charged Off
...	...	...
396025	1	Fully Paid
396026	1	Fully Paid
396027	1	Fully Paid
396028	1	Fully Paid
396029	1	Fully Paid
396030 rows x 2 columns		

**CHALLENGE TASK:** (Note this is hard, but can be done in one line!) Create a bar plot showing the correlation of the numeric features to the new loan\_repaid column. [Helpful Link](#)

```
In [39]: #CODE HERE
```

```
In [40]: df.corr()['loan_repaid'].sort_values().drop('loan_repaid').plot(kind='bar')
```

```
Out[40]: <matplotlib.axes._subplots.AxesSubplot at 0x20798034c8>
```



**TASK:** Let's examine emp\_title and emp\_length to see whether it will be okay to drop them. Print out their feature information using the feat\_info() function from the top of this notebook.

```
In [48]: # CODE HERE
```

```
In [49]: feat_info('emp_title')
print('\n')
feat_info('emp_length')
```

The job title supplied by the Borrower when applying for the loan.\*

Employment length in years. Possible values are between 0 and 10 where 0 means less than one year and 10 means ten or more years.

**TASK:** How many unique employment job titles are there?

```
In [50]: # CODE HERE
```

```
In [51]: df['emp_title'].nunique()
```

```
Out[51]: 173105
```

```
In [52]: df['emp_title'].value_counts()
```

```
Out[52]: Teacher 4389
Manager 4250
Registered Nurse 1856
RN 1846
Supervisor 1830
...
```

**TASK:** Realistically there are too many unique job titles to try to convert this to a dummy variable feature. Let's remove that emp\_title column.

```
In [53]: # CODE HERE
```

```
In [54]: df = df.drop('emp_title', axis=1)
```

**TASK:** Create a count plot of the emp\_length feature column. Challenge: Sort the order of the values.

```
In [55]: # CODE HERE
```

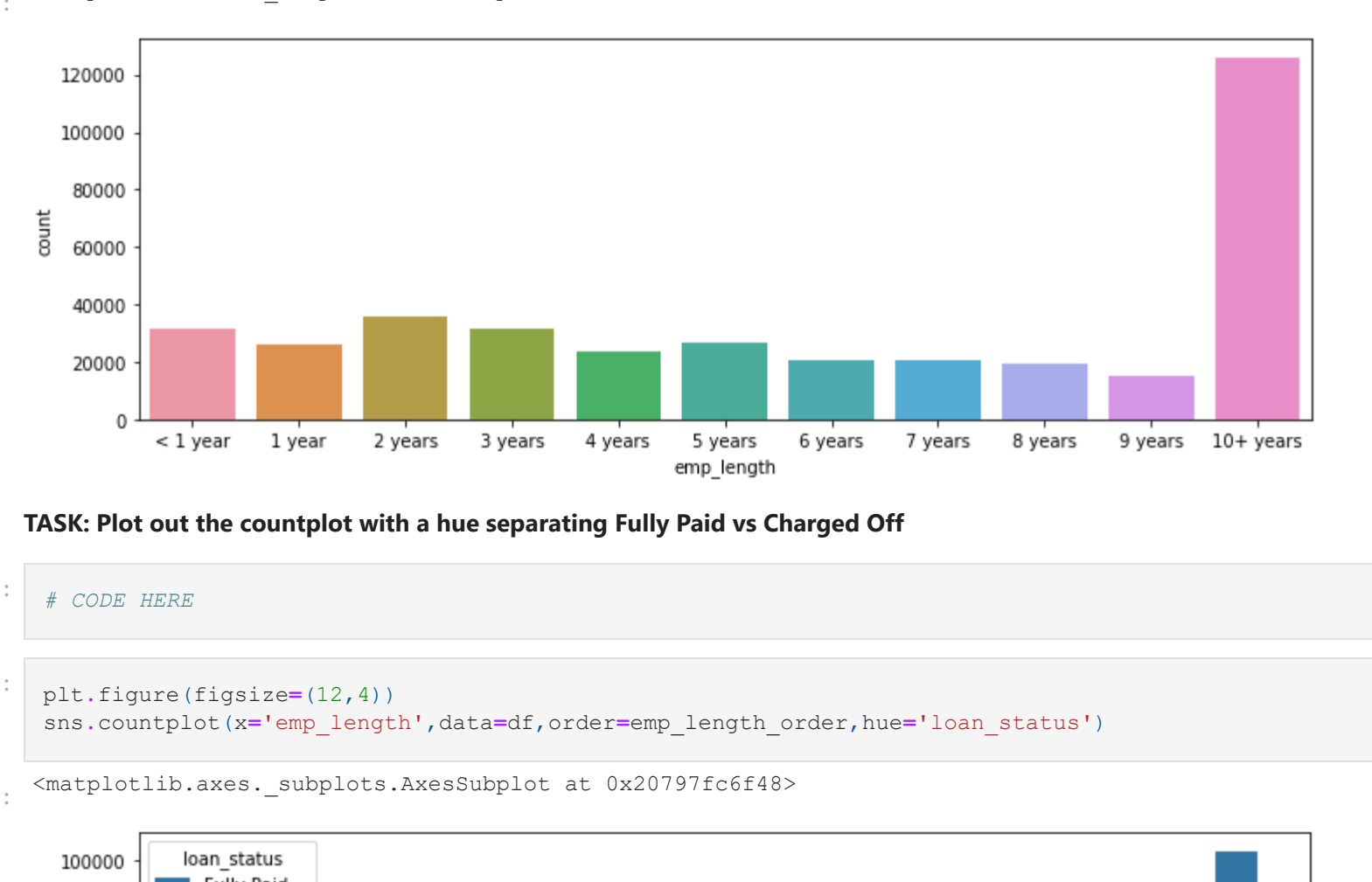
```
In [56]: sorted(df['emp_length'].dropna().unique())
```

```
Out[56]: ['10+ years',
'12 years',
'13 years',
'14 years',
'15 years',
'16 years',
'17 years',
'18 years',
'19 years',
'10+ years']
```

```
In [58]: plt.figure(figsize=(12,4))
sns.countplot(x='emp_length', data=df, order=emp_length_order)
```



<matplotlib.axes.\_subplots.AxesSubplot at 0x20797fc6f48>

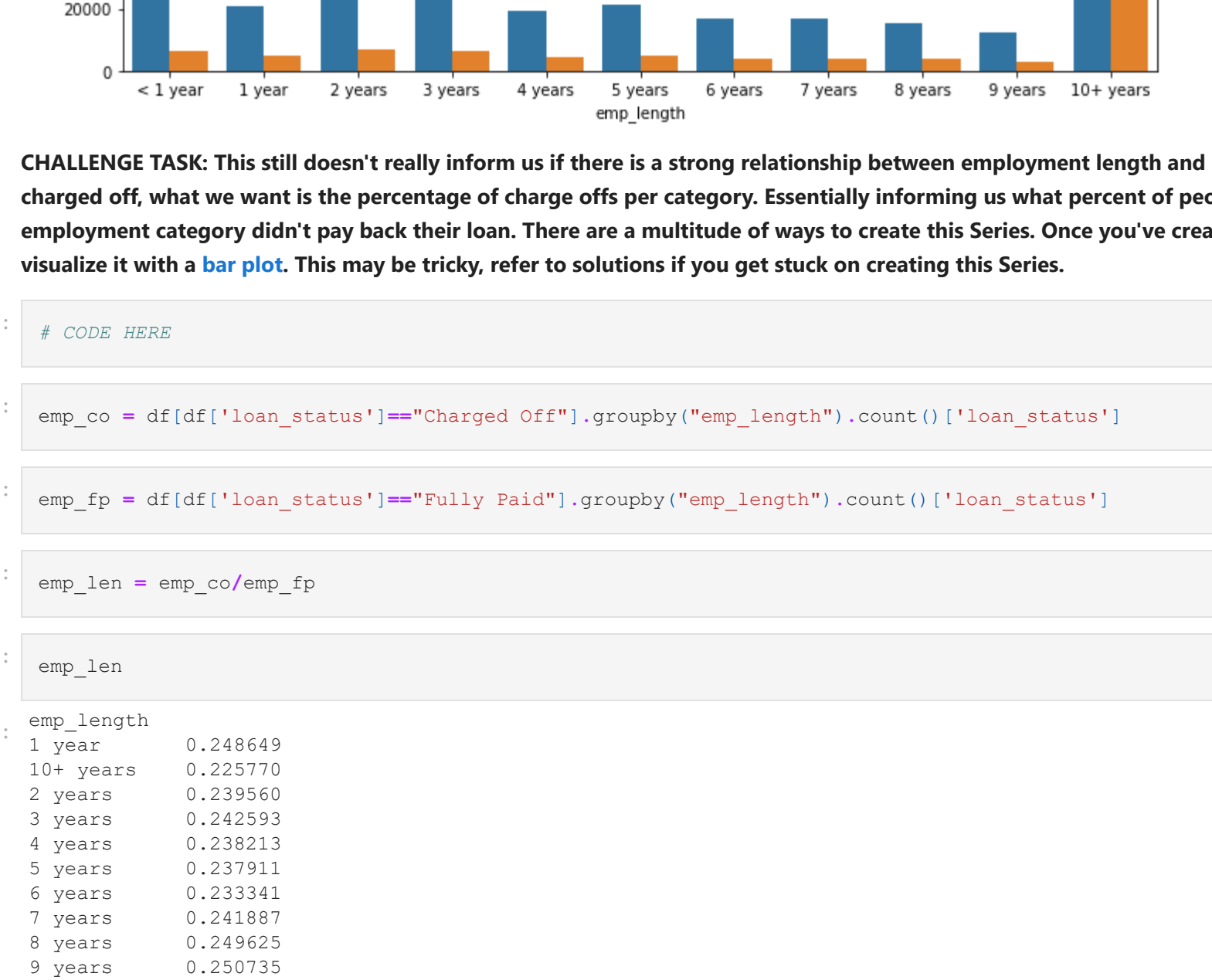


**TASK:** Plot out the countplot with a hue separating Fully Paid vs Charged Off

```
# CODE HERE
```

```
sns.factorplot(x='emp_length', y='count', hue='loan_status', data=df, order=emp_length_order, hue='loan_status')
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x20797fc6f48>



**CHALLENGE TASK:** This still doesn't really inform us if there is a strong relationship between employment length and being charged off, what we want is the percentage of charge offs per category. Essentially informing us what percent of people per employment category didn't pay back their loan. There are a multitude of ways to create this Series. Once you've created it, see if you can visualize it with a **bar plot**. This may be tricky, refer to solutions if you get stuck on creating this Series.

```
# CODE HERE
```

```
emp_co = df[df['loan_status'] == 'Charged Off'].groupby('emp_length').count()['loan_status']
```

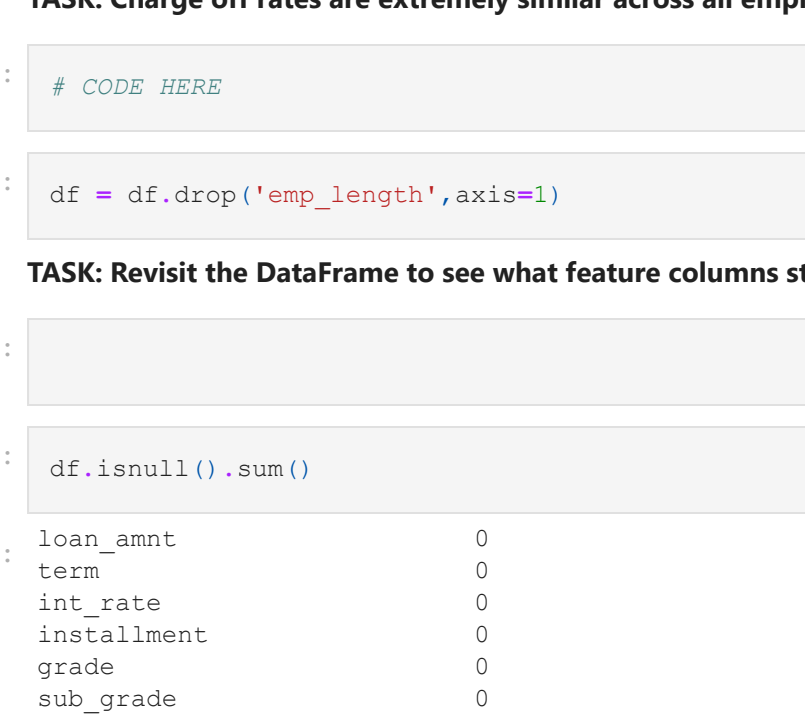
```
emp_fp = df[df['loan_status'] == 'Fully Paid'].groupby('emp_length').count()['loan_status']
```

```
emp_len = emp_co/emp_fp
```

```
emp_len
```

```
emp_len.plot(kind='bar')
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x20797fc6f48>



**TASK:** Charge off rates are extremely similar across all employment lengths. Go ahead and drop the emp\_length column.

```
# CODE HERE
```

```
df = df.drop('emp_length', axis=1)
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x20797fc6f48>

```
# CODE HERE
```

```
df
```

```
df.isnull().sum()
```

```
loan_amnt      0
term           0
int_rate       0
installment    0
grade         0
sub_grade      0
home_ownership 0
annual_inc     0
verification_status 0
issue_d        0
loan_status    0
purpose        0
title          1795
dti            0
earliest_cr_line 0
open_acc       0
pub_rec        0
revol_bal      0
revol_util     276
total_acc      0
mort_acc       0
application_type 0
initial_list_status 0
pub_rec_bankruptcies 535
address        0
loan_repaid    0
dtype: object
```

**TASK:** Review the title column vs the purpose column. Is this repeated information?

```
# CODE HERE
```

```
df['purpose'].head(10)
```

```
0      vacation
1      debt_consolidation
2      credit_card
3      credit_card
4      credit_card
5      debt_consolidation
6      home_improvement
7      credit_card
8      debt_consolidation
9      debt_consolidation
Name: purpose, dtype: object
```

```
df['title'].head(10)
```

```
0      Vacation
1      Debt consolidation
2      Credit card refinancing
3      Credit card refinancing
4      Credit Card Refinancing
5      Debt consolidation
6      Home improvement
7      No More Credit Cards
8      Debt consolidation
9      Debt consolidation
Name: title, dtype: object
```

**TASK:** The title column is simply a string subcategory/description of the purpose column. Go ahead and drop the title column.

```
# CODE HERE
```

```
df = df.drop('title', axis=1)
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x20797fc6f48>

```
# CODE HERE
```

```
df
```

```
df.isnull().sum()
```

```
loan_amnt      0
term           0
int_rate       0
installment    0
grade         0
sub_grade      0
home_ownership 0
annual_inc     0
verification_status 0
issue_d        0
loan_status    0
purpose        0
title          1795
dti            0
earliest_cr_line 0
open_acc       0
pub_rec        0
revol_bal      0
revol_util     276
total_acc      0
mort_acc       0
application_type 0
initial_list_status 0
pub_rec_bankruptcies 535
address        0
loan_repaid    0
dtype: object
```

**TASK:** There are many ways we could deal with this missing data. We could attempt to build a simple model to fill it in, such as a linear model, we could just fill it in based on the mean of the other columns, or you could even bin the columns into categories and then set NaN as its own category. There is no 100% correct approach! Let's review the other columns to see which most highly correlates to mort\_acc

```
# CODE HERE
```

```
feat_info['mort_acc']
```

```
Number of mortgage accounts.
```

**TASK:** Create a value\_counts of the mort\_acc column.

```
# CODE HERE
```

```
df['mort_acc'].value_counts()
```

```
0.0    139777
0.0016  60416
2.0     49948
3.0     38049
4.0     27887
5.0     18194
6.0     11069
7.0      6052
8.0      3121
9.0      1656
10.0     865
11.0     479
12.0     264
13.0     146
14.0     107
15.0      61
16.0      22
17.0      18
18.0      15
19.0      13
20.0      13
21.0      4
25.0      4
27.0      3
29.0      2
32.0      2
36.0      2
38.0      2
39.0      1
40.0      1
44.0      1
Name: mort_acc, dtype: int64
```

**TASK:** There are many ways we could deal with this missing data. We could attempt to build a simple model to fill it in, such as a linear model, we could just fill it in based on the mean of the other columns, or you could even bin the columns into categories and then set NaN as its own category. There is no 100% correct approach! Let's review the other columns to see which most highly correlates to mort\_acc

```
# CODE HERE
```

```
print("Correlation with the mort_acc column")
```

```
dfcorr = df[['mort_acc']].corr()
```

```
Correlation with the mort_acc column
int_rate      -0.025891
dti           -0.025439
revol_util     0.007514
pub_rec       0.011552
pub_rec_bankruptcies 0.027239
loan_repaid    0.073111
mort_acc      0.109205
installment    0.193694
revol_bal     0.194925
loan_amnt     0.222315
annual_inc    0.236320
total_acc     0.381072
mort_acc      1.000000
Name: mort_acc, dtype: float64
```

**TASK:** Looks like the mort\_acc column value correlates with the mort\_acc, this makes sense! Let's try this fillna() approach. We will group the dataframe by the total\_acc and calculate the mean value for the mort\_acc per total\_acc entry. To get the result below:

```
# CODE HERE
```

```
print("Mean of mort_acc column per total_acc")
```

```
df.groupby('total_acc').mean()['mort_acc']
```

```
Mean of mort_acc column per total_acc
```

```
total_acc
```

```
2.0    0.000000
3.0    0.002023
4.0    0.066743
5.0    0.103289
6.0    0.131293
...
```

```
124.0  1.000000
129.0  1.000000
135.0  3.000000
150.0  2.000000
151.0  0.000000
Name: mort_acc, Length: 118, dtype: float64
```

**CHALLENGE TASK:** Let's fill in the missing mort\_acc values based on their total\_acc value. If the mort\_acc is missing, then we will fill in that missing value with the mean value corresponding to its total\_acc value from the Series we created above. This involves using an apply() method with two columns. Check out the link below for more info, or review the solutions video/notebook.

[Helpful Link](#)

```
# CODE HERE
```

```
total_acc_avg = df.groupby('total_acc').mean()['mort_acc']
```

```
total_acc_avg[2.0]
```

```
0.0
```

```
def fill_mort_acc(total_acc, mort_acc):
```

```
    """
```

```
    Accepts the total_acc and mort_acc values for the row.
```

```
    Checks if the mort_acc is NaN - if so, it returns the avg mort_acc value
```

```
    for the corresponding total_acc value for that row.
```

```
    total_acc_avg here should be a Series or dictionary containing the mapping of the
```

```
    groupby averages of mort_acc per total_acc values.
```

```
    """
```

```
    if np.isnan(mort_acc):
```

```
        return total_acc_avg[total_acc]
```

```
    else:
```

```
        return mort_acc
```

```
df['mort_acc'] = df.apply(lambda x: fill_mort_acc(x['total_acc'], x['mort_acc']), axis=1)
```

```
df.isnull().sum()
```

```
loan_amnt      0
term           0
int_rate       0
installment    0
grade         0
sub_grade      0
home_ownership 0
annual_inc     0
verification_status 0
issue_d        0
loan_status    0
purpose        0
title          1795
dti            0
earliest_cr_line 0
open_acc       0
pub_rec        0
revol_bal      0
revol_util     276
total_acc      0
mort_acc       0
application_type 0
initial_list_status 0
pub_rec_bankruptcies 535
address        0
loan_repaid    0
dtype: object
```

**TASK:** The revol\_util and the pub\_rec\_bankruptcies have missing data points, but they account for less than 0.5% of the total data. Go ahead and remove the rows that are missing those values in those columns with dropna().

```
# CODE HERE
```

```
df = df.dropna()
```

```
df.isnull().sum()
```

```
loan_amnt      0
term           0
int_rate       0
installment    0
grade         0
sub_grade      0
home_ownership 0
annual_inc     0
verification_status 0
issue_d        0
loan_status    0
purpose        0
title          1795
dti            0
earliest_cr_line 0
open_acc       0
pub_rec        0
revol_bal      0
revol_util     276
total_acc      0
mort_acc       0
application_type 0
initial_list_status 0
pub_rec_bankruptcies 535
address        0
loan_repaid    0
dtype: object
```

## Categorical Variables and Dummy Variables

We're done working with the missing data! Now we just need to deal with the string values due to the categorical columns.

**TASK:** List all the columns that are currently non-numeric. [Helpful Link](#)

[Another very useful method call](#)

```
# CODE HERE
```

```
df.select_dtypes(['object']).columns
```

```
Index(['term', 'grade', 'sub_grade', 'home_ownership', 'verification_status',
```

```
       'issue_d', 'loan_status', 'purpose', 'earliest_cr_line',
```

```
       'initial_list_status', 'application_type', 'address',
```

```
       dtype='object'])
```

Let's now go through all the string features to see what we should do with them.

### term feature

**TASK:** Convert the term feature into either a 36 or 60 integer numeric data type using .apply() or .map().

```
# CODE HERE
```

```
df['term'].value_counts()
```

```
36 months    301247
60 months    93972
Name: term, dtype: int64
```

```
# Or just use .map()
```

```
df['term'] = df['term'].apply(lambda term: int(term[:3]))
```

```
df['term'].value_counts()
```

```
36 months    301247
60 months    93972
Name: term, dtype: int64
```

### grade feature

**TASK:** We already know grade is part of sub\_grade, so just drop the grade feature.

```
# CODE HERE
```

```
df = df.drop('grade', axis=1)
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x20797fc6f48>

```
# CODE HERE
```

```
df
```

```
df.isnull().sum()
```

```
loan_amnt      0
term           0
int_rate       0
installment    0
grade         0
sub_grade      0
home_ownership 0
annual_inc     0
verification_status 0
issue_d        0
loan_status    0
purpose        0
title          1795
dti            0
earliest_cr_line 0
open_acc       0
pub_rec        0
revol_bal      0
revol_util     276
total_acc      0
mort_acc       0
application_type 0
initial_list_status 0
pub_rec_bankruptcies 535
address        0
loan_repaid    0
dtype: object
```

**TASK:** Convert the subgrade into dummy variables. Then concatenate these new columns to the original dataframe. Remember to drop the original subgrade column and to add drop\_first=True to your get\_dummies call.

```
# CODE HERE
```

```
subgrade_dummies = pd.get_dummies(df['sub_grade'], drop_first=True)
```

```
df = pd.concat([df.drop('sub_grade', axis=1), subgrade_dummies], axis=1)
```

```
df.columns
```

```
Index(['loan_amnt', 'term', 'int_rate', 'installment', 'home_ownership',
```

```
       'annual_inc', 'verification_status', 'issue_d', 'loan_status',
```

```
       'purpose', 'dti', 'earliest_cr_line', 'open_acc', 'pub_rec',
```

```
       'revol_bal', 'revol_util', 'total_acc', 'initial_list_status',
```

```
       'application_type', 'mort_acc', 'pub_rec_bankruptcies', 'address',
```

```
       'loan_repaid', 'A2', 'A3', 'A4', 'A5', 'B1', 'B2', 'B3', 'B4', 'B5',
```

```
       'C1', 'C2', 'C3', 'C4', 'C5', 'D1', 'D2', 'D3', 'D4', 'D5', 'E1', 'E2',
```

```
       'E3', 'E4', 'E5', 'F1', 'F2', 'F3', 'F4', 'F5', 'G1', 'G2', 'G3', 'G4',
```

```
       'G5'],
```

```
       dtype='object'])
```

```
df.select_dtypes(['object']).columns
```

```
Index(['home_ownership', 'verification_status', 'issue_d', 'loan_status',
```

```
       'purpose', 'earliest_cr_line', 'initial_list_status',
```

```
       'application_type', 'address',
```

```
       dtype='object'])
```

### verification\_status, application\_type, initial\_list\_status, purpose

**TASK:** Convert these columns ('verification\_status', 'application\_type', 'initial\_list\_status', 'purpose') into dummy variables and concatenate them with the original dataframe. Remember to set drop\_first=True and to drop the original columns.

```
# CODE HERE
```

```
dummies = pd.get_dummies(df[['verification_status', 'application_type', 'initial_list_status', 'purpose']], drop_first=True)
```

```
df = df.drop(['verification_status', 'application_type', 'initial_list_status', 'purpose'], axis=1)
```

```
df = pd.concat([df, dummies], axis=1)
```

```
df
```

```
df.isnull().sum()
```

```
loan_amnt      0
term           0
int_rate       0
installment    0
grade         0
sub_grade      0
home_ownership 0
annual_inc     0
verification_status 0
issue_d        0
loan_status    0
purpose        0
title          1795
dti            0
earliest_cr_line 0
open_acc       0
pub_rec        0
revol_bal      0
revol_util     276
total_acc      0
mort_acc       0
application_type 0
initial_list_status 0
pub_rec_bankruptcies 535
address        0
loan_repaid    0
dtype: object
```

### home\_ownership

**TASK:** Review the value counts for the home\_ownership column.

```
# CODE HERE
```

```
df['home_ownership'].value_counts()
```

```
MORTGAGE    198022
RENT         159395
OWN          37660
OTHER        110
NONE         29
ANY          0
Name: home_ownership, dtype: int64
```

**TASK:** Convert these to dummy variables, but replace NONE and ANY with OTHER, so that we end up with just 4 categories, MORTGAGE, RENT, OWN, OTHER. Then concatenate them with the original dataframe. Remember to set drop\_first=True and to drop the original columns.

```
# CODE HERE
```

```
df['home_ownership'] = df['home_ownership'].replace({'NONE', 'ANY'}, 'OTHER')
```

```
dummies = pd.get_dummies(df['home_ownership'], drop_first=True)
```

```
df = df.drop('home_ownership', axis=1)
```

```
df = pd.concat([df, dummies], axis=1)
```

```
df
```

```
df.isnull().sum()
```

```
loan_amnt      0
term           0
int_rate       0
installment    0
grade         0
sub_grade      0
home_ownership 0
annual_inc     0
verification_status 0
issue_d        0
loan_status    0
purpose        0
title          1795
dti            0
earliest_cr_line 0
open_acc       0
pub_rec        0
revol_bal      0
revol_util     276
total_acc      0
mort_acc       0
application_type 0
initial_list_status 0
pub_rec_bankruptcies 535
address        0
loan_repaid    0
dtype: object
```

### address

**TASK:** Let's feature engineer a zip code column from the address in the data set. Create a column called 'zip\_code' that extracts the zip code from the address column.

```
# CODE HERE
```

```
df['zip_code'] = df['address'].apply(lambda address: address[-5:])
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x20797fc6f48>

```
# CODE HERE
```

```
df
```

```
df.isnull().sum()
```

```
loan_amnt      0
term           0
int_rate       0
installment    0
grade         0
sub_grade      0
home_ownership 0
annual_inc     0
verification_status 0
issue_d        0
loan_status    0
purpose        0
title          1795
dti            0
earliest_cr_line 0
open_acc       0
pub_rec        0
revol_bal      0
revol_util     276
total_acc      0
mort_acc       0
application_type 0
initial_list_status 0
pub_rec_bankruptcies 535
address        0
loan_repaid    0
dtype: object
```

### OPTIONAL

#### Grabbing a Sample for Training Time

**OPTIONAL:** Use .sample() to grab a sample of the 490k+ entries to save time on training. Highly recommended for lower RAM computers or if you are not using GPU.

```
# CODE HERE
```

```
df = df.sample(frac=0.2, random_state=101)
```

```
print(len(df))
```

```
195219
```

**TASK:** Perform a train/test split with test\_size=0.2 and a random\_state of 101.

```
# CODE HERE
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=101)
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x20797fc6f48>

```
# CODE HERE
```

```
df
```

```
df.isnull().sum()
```

```
loan_amnt      0
term           0
int_rate       0
installment    0
grade         0
sub_grade      0
home_ownership 0
annual_inc     0
verification_status 0
issue_d        0
loan_status    0
purpose        0
title          1795
dti            0
earliest_cr_line 0
open_acc       0
pub_rec        0
revol_bal      0
revol_util     276
total_acc      0
mort_acc       0
application_type 0
initial_list_status 0
pub_rec_bankruptcies 535
address        0
loan_repaid    0
dtype: object
```

### Normalizing the Data

**TASK:** Use a MinMaxScaler to normalize the feature data X\_train and X\_test. Recall we don't want data leakage from the test set so we only fit on the X\_train data.

```
# CODE HERE
```

```
from sklearn.preprocessing import MinMaxScaler
```

```
scaler = MinMaxScaler()
```

```
X_train = scaler.fit_transform(X_train)
```

```
X_test = scaler.transform(X_test)
```

### Creating the Model

**TASK:** Run the cell below to import the necessary Keras functions.

```
# CODE HERE
```

```
import tensorflow as tf
```

```
from tensorflow.keras.models import Sequential
```

```
from tensorflow.keras.layers import Dense, Activation, Dropout
```

```
from tensorflow.keras.constraints import max_norm
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x20797fc6f48>

```
# CODE HERE
```

```
model = Sequential()
```

```
# Choose whatever number of layers/neurons you want.
```

```
# https://stats.stackexchange.com/questions/181/how-to-choose-the-number-of-hidden-layers-and-nodes-in-a-feedf
```

```
# Remember to compile()
```

```
model = Sequential()
```

```
# https://stats.stackexchange.com/questions/181/how-to-choose-the-number-of-hidden-layers-and-nodes-in-a-feedf
```

```
# Input layer
```

```
model.add(Dense(78, activation='relu'))
```

```
model.add(Dropout(0.2))
```

```
# hidden layer
```

```
model.add(Dense(39, activation='relu'))
```

```
model.add(Dropout(0.2))
```

```
# hidden layer
```

```
model.add(Dense(19, activation='relu'))
```

```
model.add(Dropout(0.2))
```

```
# output layer
```

```
model.add(Dense(units=1, activation='sigmoid'))
```

```
# Compile model
```

```
model.compile(loss='binary_crossentropy', optimizer='adam')
```



