



**Copyright by Pierian Data Inc.
Created by Jose Marcial Portilla.**

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 out 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.

	LoanStatNew	Description
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

	LoanStatNew	Description
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

In [2]: data_info = pd.read_csv('../DATA/lending_club_info.csv',index_col='LoanStatNew')

In [3]: 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 [4]: def feat_info(col_name):
    print(data_info.loc[col_name]['Description'])

In [5]: 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):
 #   Column           Non-Null Count  Dtype  
--- 
 0   loan_amnt        396030 non-null  float64 
 1   term              396030 non-null  object  
 2   int_rate          396030 non-null  float64 
 3   installment       396030 non-null  float64 
 4   grade             396030 non-null  object  
 5   sub_grade         396030 non-null  object  
 6   emp_title         373103 non-null  object  
 7   emp_length        377729 non-null  object  
 8   home_ownership    396030 non-null  object  
 9   annual_inc        396030 non-null  float64 
 10  verification_status 396030 non-null  object  
 11  issue_d           396030 non-null  object  
 12  loan_status       396030 non-null  object  
 13  purpose            396030 non-null  object  
 14  title              394274 non-null  object  
 15  dti                396030 non-null  float64 
 16  earliest_cr_line  396030 non-null  object  
 17  open_acc           396030 non-null  float64 
 18  pub_rec            396030 non-null  float64 
 19  revol_bal          396030 non-null  float64 
 20  revol_util         395754 non-null  float64 
 21  total_acc          396030 non-null  float64 
 22  initial_list_status 396030 non-null  object  
 23  application_type   396030 non-null  object  
 24  mort_acc            358235 non-null  float64 
 25  pub_rec_bankruptcies 395495 non-null  float64 
 26  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

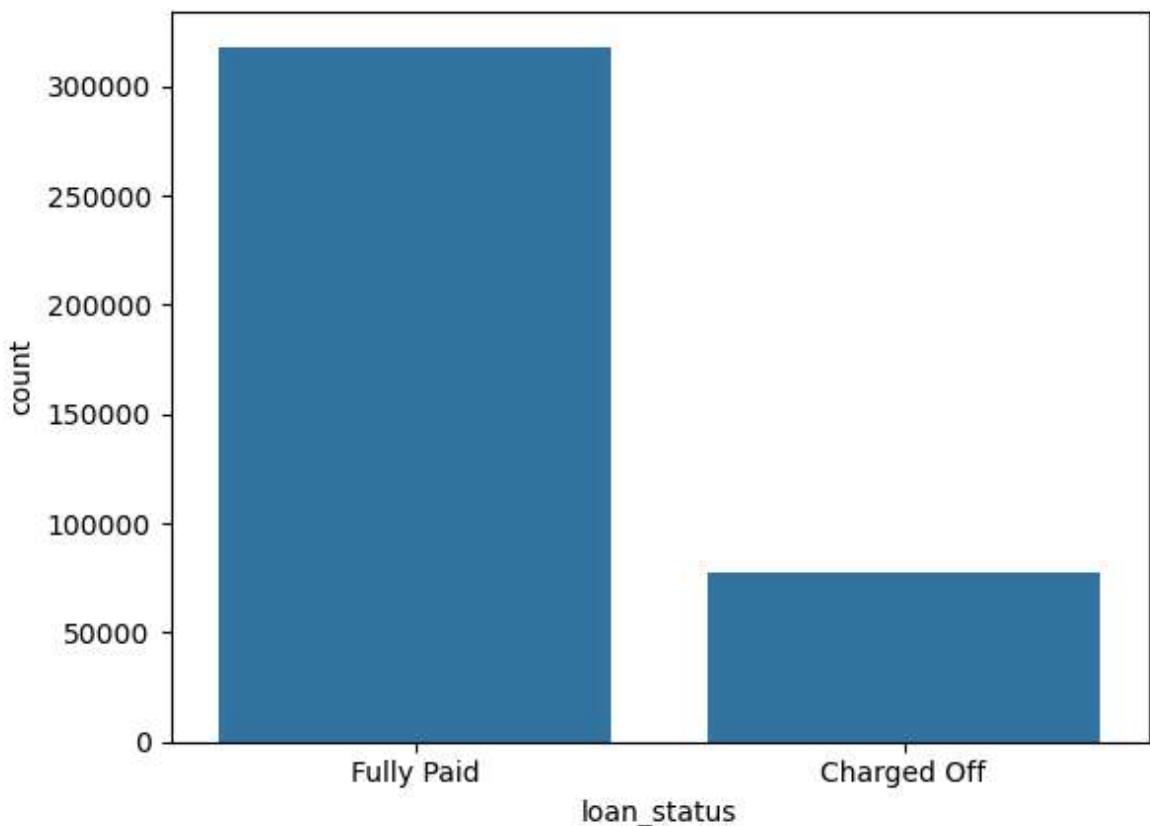
Section 1: Exploratory Data Analysis

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

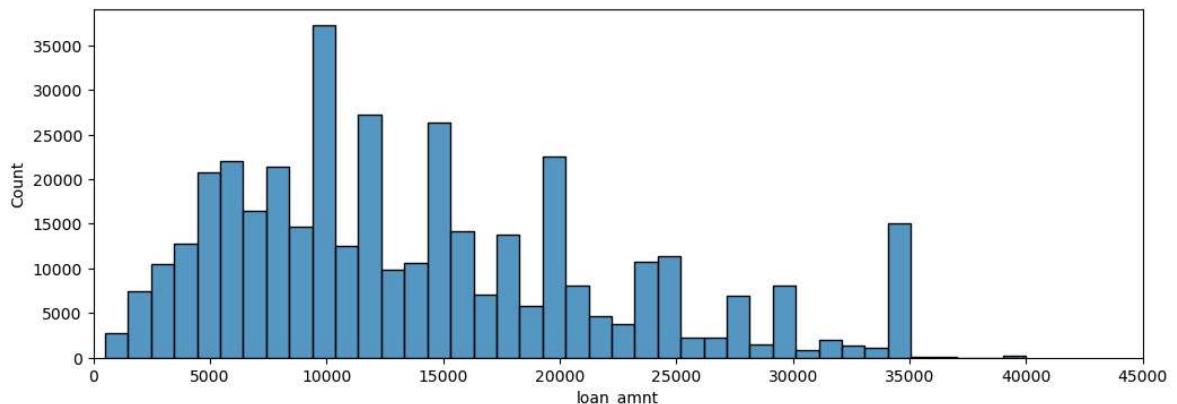
Out[9]: <Axes: xlabel='loan_status', ylabel='count'>



TASK: Create a histogram of the loan_amnt column.

```
In [10]: plt.figure(figsize=(12,4))
sns.histplot(df['loan_amnt'], bins=40, kde=False, edgecolor='black')
plt.xlim(0,45000)
```

Out[10]: (0.0, 45000.0)



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

```
In [11]: df.select_dtypes(include=['number']).corr()
```

Out[11]:

	loan_amnt	int_rate	installment	annual_inc	dti	open_ac
loan_amnt	1.000000	0.168921	0.953929	0.336887	0.016636	0.19855
int_rate	0.168921	1.000000	0.162758	-0.056771	0.079038	0.01164
installment	0.953929	0.162758	1.000000	0.330381	0.015786	0.18897
annual_inc	0.336887	-0.056771	0.330381	1.000000	-0.081685	0.13615
dti	0.016636	0.079038	0.015786	-0.081685	1.000000	0.13618
open_acc	0.198556	0.011649	0.188973	0.136150	0.136181	1.00000
pub_rec	-0.077779	0.060986	-0.067892	-0.013720	-0.017639	-0.01839
revol_bal	0.328320	-0.011280	0.316455	0.299773	0.063571	0.22119
revol_util	0.099911	0.293659	0.123915	0.027871	0.088375	-0.13142
total_acc	0.223886	-0.036404	0.202430	0.193023	0.102128	0.68072
mort_acc	0.222315	-0.082583	0.193694	0.236320	-0.025439	0.10920
pub_rec_bankruptcies	-0.106539	0.057450	-0.098628	-0.050162	-0.014558	-0.02773

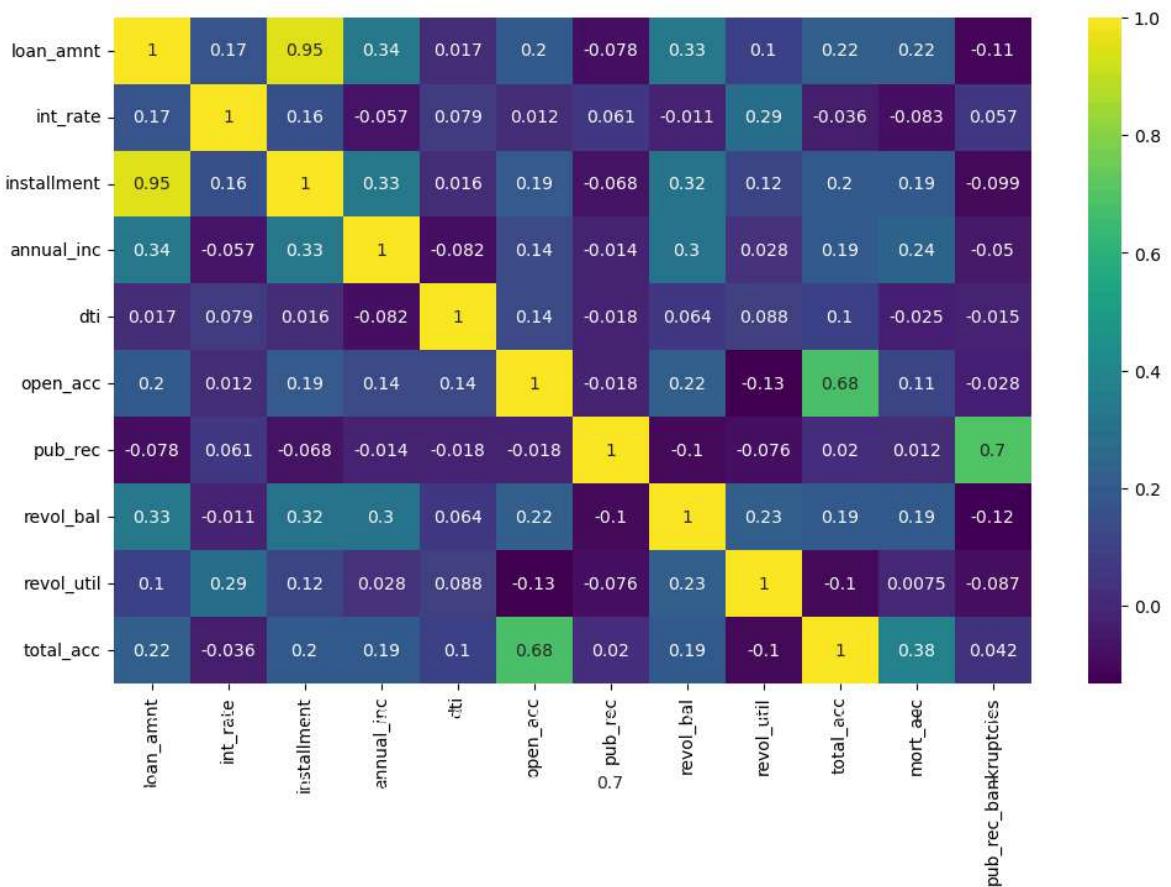
TASK: Visualize this using a heatmap. Depending on your version of matplotlib, you may need to manually adjust the heatmap.

- [Heatmap info](#)
- [Help with resizing](#)

In [12]:

```
plt.figure(figsize=(12,7))
sns.heatmap(df.select_dtypes(include=['number']).corr(), annot=True, cmap='viridis')
plt.ylim(10, 0)
```

Out[12]: (10.0, 0.0)



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 [13]: feat_info('installment')
```

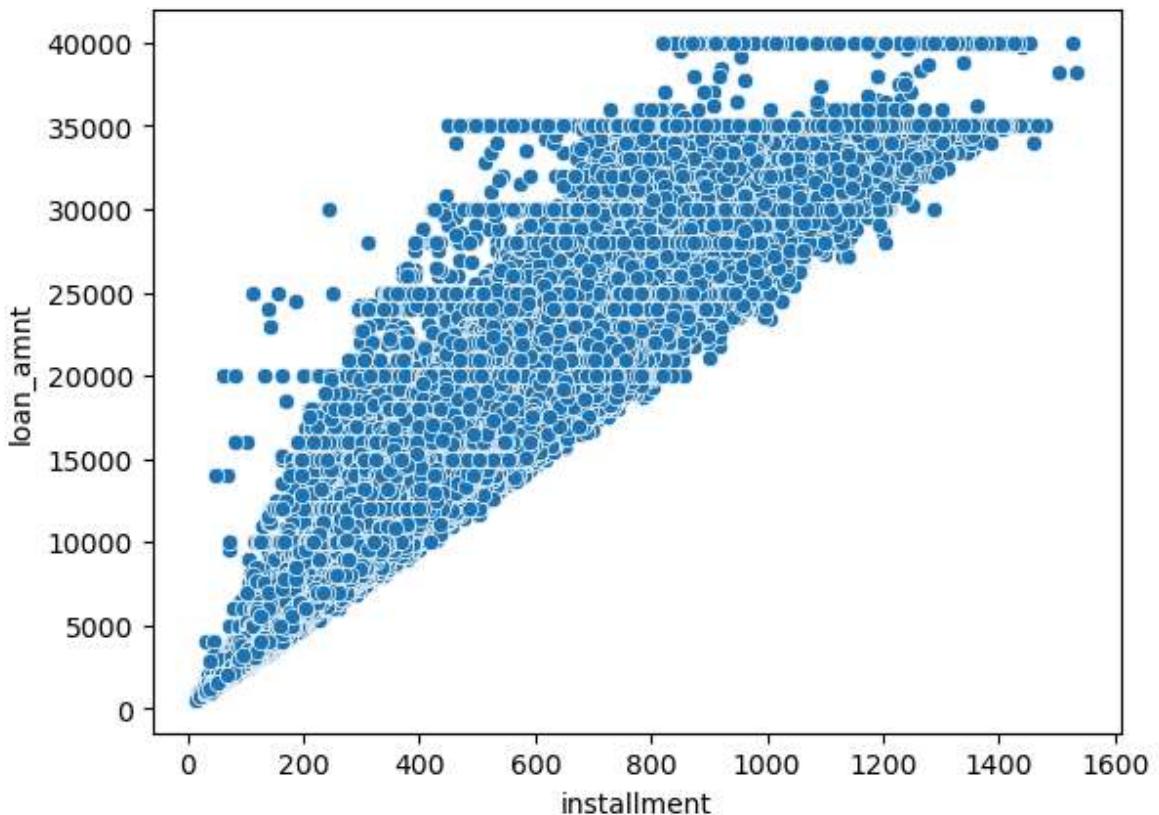
The monthly payment owed by the borrower if the loan originates.

```
In [14]: 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 [15]: sns.scatterplot(x='installment',y='loan_amnt',data=df,)
```

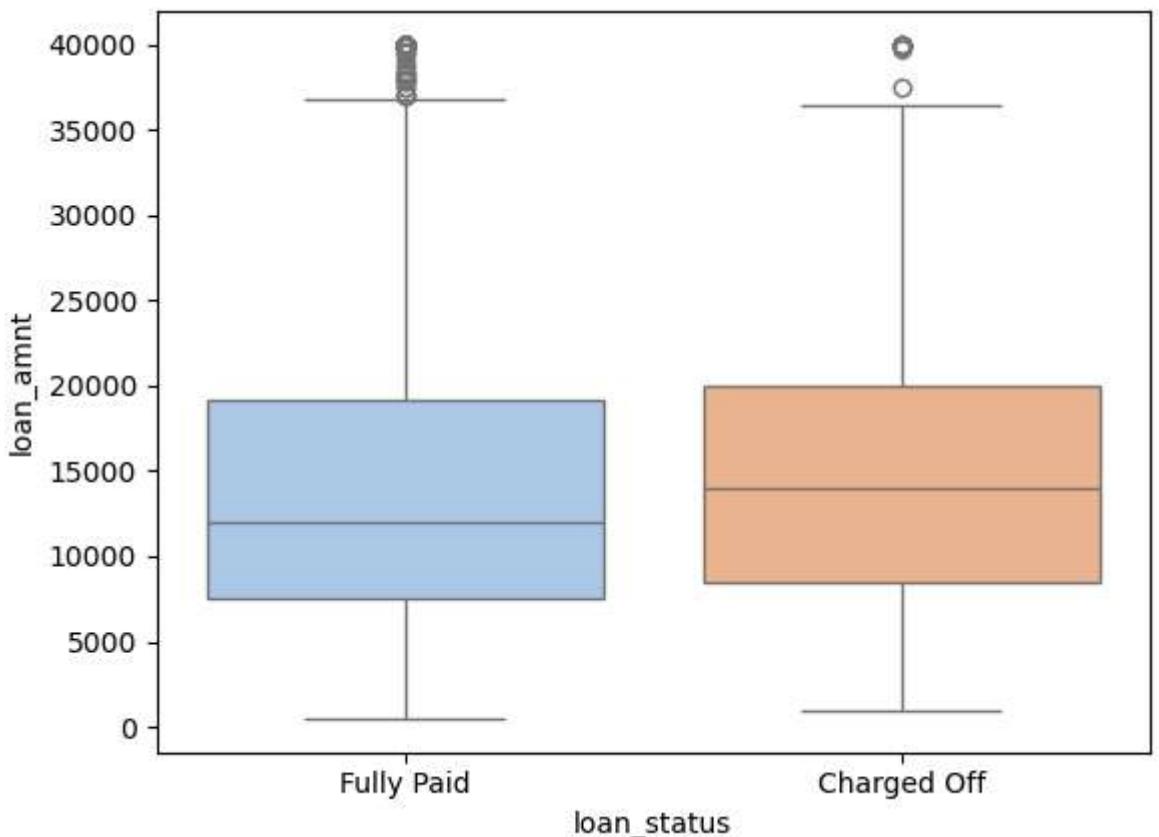
```
Out[15]: <Axes: xlabel='installment', ylabel='loan_amnt'>
```



TASK: Create a boxplot showing the relationship between the loan_status and the Loan Amount.

```
In [16]: sns.boxplot(palette='pastel',x='loan_status',y='loan_amnt',data=df, hue='loan_st
```

```
Out[16]: <Axes: xlabel='loan_status', ylabel='loan_amnt'>
```



TASK: Calculate the summary statistics for the loan amount, grouped by the loan_status.

```
In [17]: df.groupby('loan_status')['loan_amnt'].describe()
```

loan_status	count	mean	std	min	25%	50%	75%	m
Charged Off	77673.0	15126.300967	8505.090557	1000.0	8525.0	14000.0	20000.0	4000
Fully Paid	318357.0	13866.878771	8302.319699	500.0	7500.0	12000.0	19225.0	4000

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

```
In [18]: sorted(df['grade'].unique())
```

```
Out[18]: ['A', 'B', 'C', 'D', 'E', 'F', 'G']
```

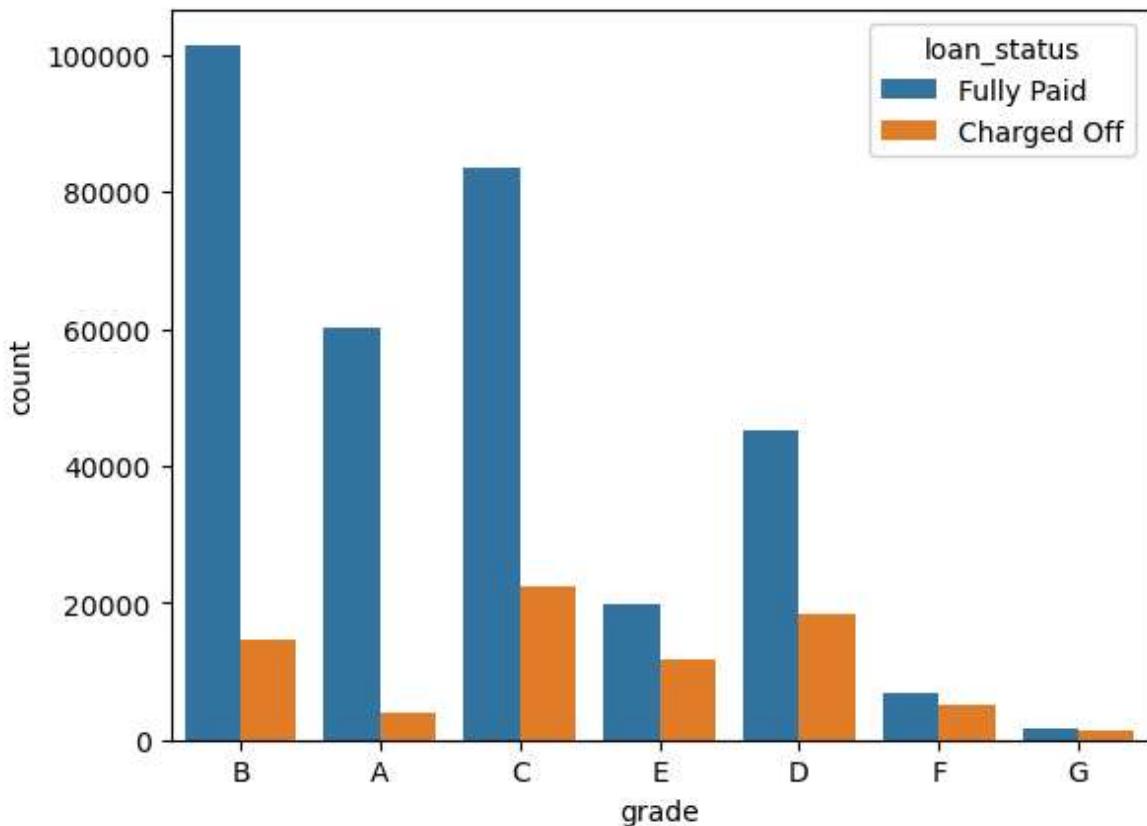
```
In [19]: sorted(df['sub_grade'].unique())
```

```
Out[19]: ['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',
          'F1',
          'F2',
          'F3',
          'F4',
          'F5',
          'G1',
          'G2',
          'G3',
          'G4',
          'G5']
```

TASK: Create a countplot per grade. Set the hue to the loan_status label.

```
In [20]: sns.countplot(x='grade', data=df, hue='loan_status')
```

```
Out[20]: <Axes: xlabel='grade', ylabel='count'>
```



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

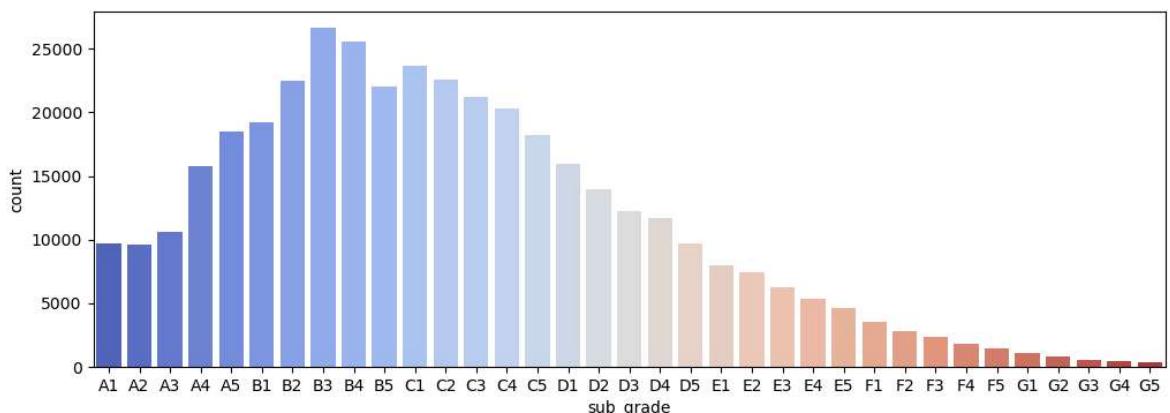
```
In [21]: plt.figure(figsize=(12,4))
subgrade_order = sorted(df['sub_grade'].unique())
sns.countplot(x='sub_grade', data=df, order = subgrade_order, palette='coolwarm')
```

/tmp/ipykernel_7835/3718666630.py:3: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v 0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
sns.countplot(x='sub_grade', data=df, order = subgrade_order, palette='coolwarm')
```

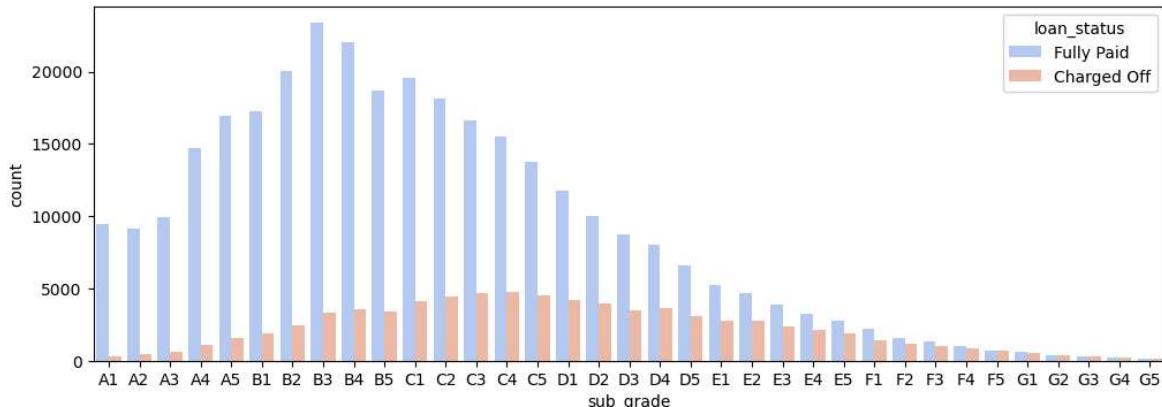
```
Out[21]: <Axes: xlabel='sub_grade', ylabel='count'>
```



```
In [22]: plt.figure(figsize=(12,4))
subgrade_order = sorted(df['sub_grade'].unique())
```

```
sns.countplot(x='sub_grade', data=df, order = subgrade_order, palette='coolwarm' , h
```

Out[22]: <Axes: xlabel='sub_grade', ylabel='count'>

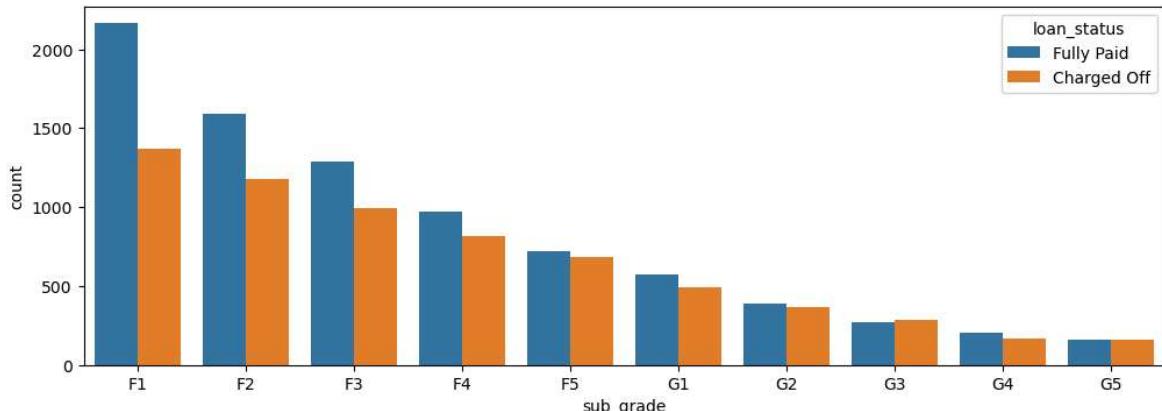


TASK: It looks like F and G subgrades don't get paid back that often. Isolate those and recreate the countplot just for those subgrades.

```
In [23]: 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[23]: <Axes: xlabel='sub_grade', ylabel='count'>



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 [24]: df['loan_status'].unique()
```

Out[24]: array(['Fully Paid', 'Charged Off'], dtype=object)

```
In [25]: df['loan_repaid'] = df['loan_status'].map({'Fully Paid':1,'Charged Off':0})
```

```
In [26]: df[['loan_repaid','loan_status']]
```

Out[26]:

	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

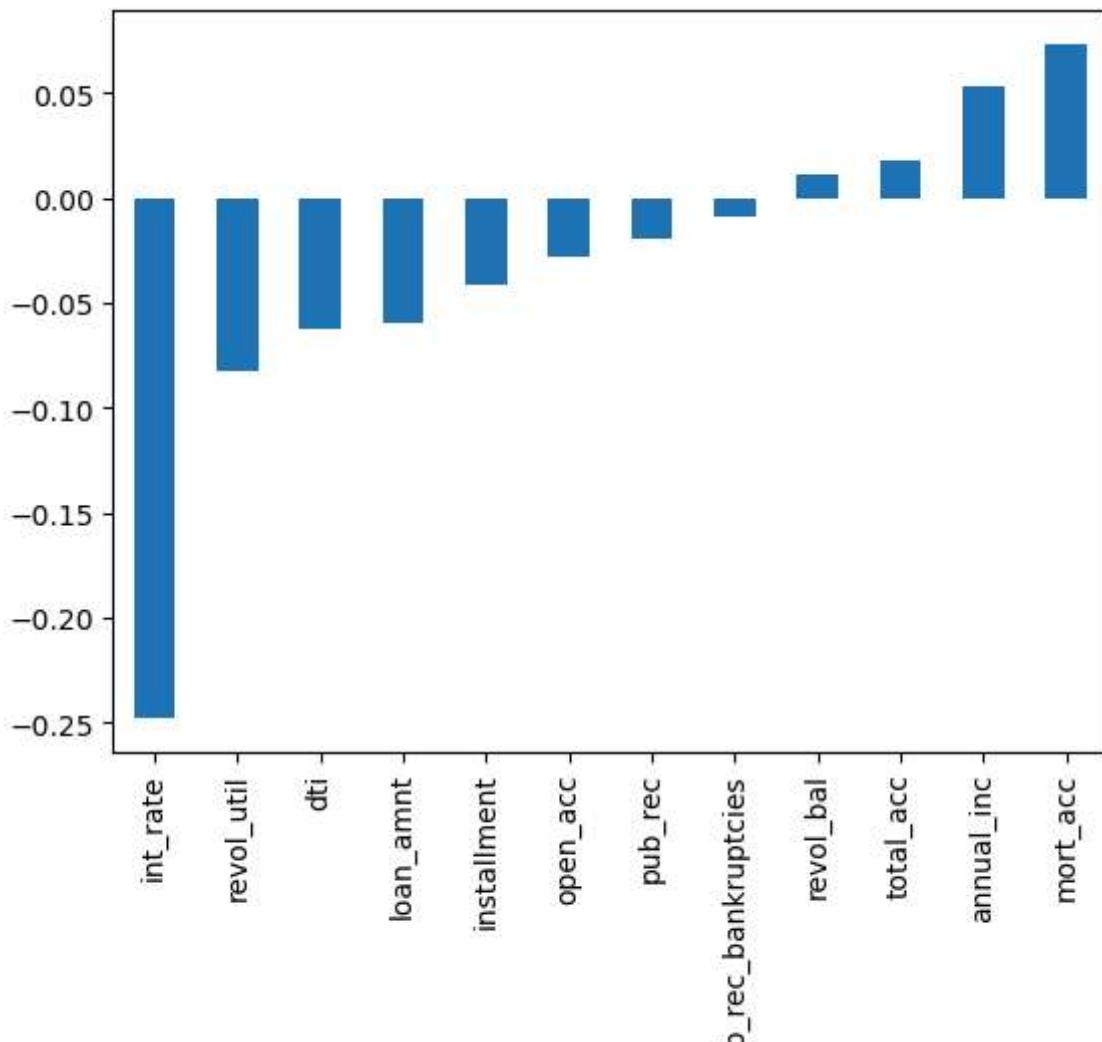
396030 rows × 2 columns

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 [27]: `df.select_dtypes(include=['number']).corr()['loan_repaid'].sort_values().drop('1')`

Out[27]: <Axes: >



Section 2: Data PreProcessing

Section Goals: Remove or fill any missing data. Remove unnecessary or repetitive features. Convert categorical string features to dummy variables.

In [28]: `df.head()`

Out[28]:

	loan_amnt	term	int_rate	installment	grade	sub_grade	emp_title	emp_length
0	10000.0	36 months	11.44	329.48	B	B4	Marketing	10+ years
1	8000.0	36 months	11.99	265.68	B	B5	Credit analyst	4 years
2	15600.0	36 months	10.49	506.97	B	B3	Statistician	< 1 year
3	7200.0	36 months	6.49	220.65	A	A2	Client Advocate	6 years
4	24375.0	60 months	17.27	609.33	C	C5	Destiny Management Inc.	9 years

5 rows × 28 columns



Missing Data

Let's explore this missing data columns. We use a variety of factors to decide whether or not they would be useful, to see if we should keep, discard, or fill in the missing data.

TASK: What is the length of the dataframe?

In [29]: `len(df)`

Out[29]: 396030

TASK: Create a Series that displays the total count of missing values per column.

In [30]: `df.isnull().sum()`

```
Out[30]: loan_amnt      0  
term          0  
int_rate      0  
installment   0  
grade         0  
sub_grade     0  
emp_title    22927  
emp_length   18301  
home_ownership 0  
annual_inc    0  
verification_status 0  
issue_d       0  
loan_status   0  
purpose       0  
title        1756  
dti           0  
earliest_cr_line 0  
open_acc      0  
pub_rec       0  
revol_bal    0  
revol_util   276  
total_acc     0  
initial_list_status 0  
application_type 0  
mort_acc     37795  
pub_rec_bankruptcies 535  
address       0  
loan_repaid   0  
dtype: int64
```

TASK: Convert this Series to be in term of percentage of the total DataFrame

```
In [31]: 100* df.isnull().sum()/len(df)
```

```
Out[31]: loan_amnt      0.000000
          term         0.000000
          int_rate     0.000000
          installment   0.000000
          grade        0.000000
          sub_grade    0.000000
          emp_title    5.789208
          emp_length   4.621115
          home_ownership 0.000000
          annual_inc    0.000000
          verification_status 0.000000
          issue_d       0.000000
          loan_status   0.000000
          purpose       0.000000
          title         0.443401
          dti           0.000000
          earliest_cr_line 0.000000
          open_acc      0.000000
          pub_rec       0.000000
          revol_bal     0.000000
          revol_util    0.069692
          total_acc     0.000000
          initial_list_status 0.000000
          application_type 0.000000
          mort_acc      9.543469
          pub_rec_bankruptcies 0.135091
          address       0.000000
          loan_repaid   0.000000
          dtype: float64
```

Lets examine emp_title and emp_length to see whether it will be okay to drop them. Print out their feature information using feat_info() function

```
In [32]: feat_info('emp_title')
```

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

```
In [33]: feat_info('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.

Lets look at how many unique job titles there are

```
In [34]: df['emp_title'].nunique() # Theres a lot, half of the data set are unique titles
```

```
Out[34]: 173105
```

```
In [35]: df['emp_title'].value_counts()
```

```
Out[35]: emp_title
Teacher           4389
Manager          4250
Registered Nurse 1856
RN              1846
Supervisor        1830
...
lab technologist      1
sutures            1
Ddepartment of defence 1
Sr Sales Analytics Manager 1
UPS Field Service Technician 1
Name: count, Length: 173105, dtype: int64
```

```
In [36]: # way too many titles to convert to dummy variables, so drop it.
```

```
In [37]: df = df.drop('emp_title', axis=1)
```

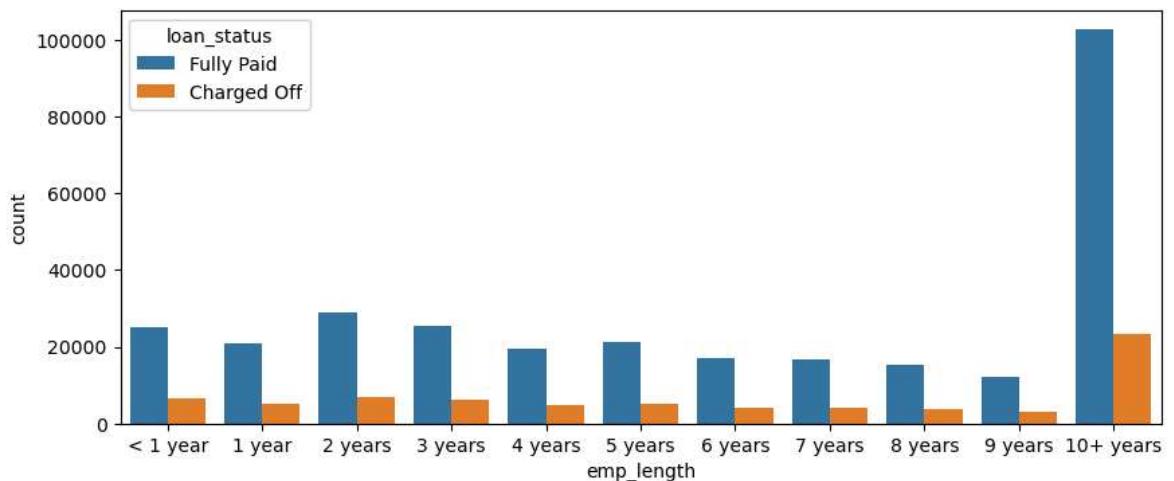
```
In [38]: # Now for emp_length
sorted(df['emp_length'].dropna().unique())
```

```
Out[38]: ['1 year',
 '10+ years',
 '2 years',
 '3 years',
 '4 years',
 '5 years',
 '6 years',
 '7 years',
 '8 years',
 '9 years',
 '< 1 year']
```

```
In [39]: emp_length_order = ['< 1 year',
 '1 year',
 '2 years',
 '3 years',
 '4 years',
 '5 years',
 '6 years',
 '7 years',
 '8 years',
 '9 years',
 '10+ years']
```

```
In [40]: plt.figure(figsize=(10,4))
sns.countplot(x='emp_length', data=df, order=emp_length_order, hue='loan_status')
```

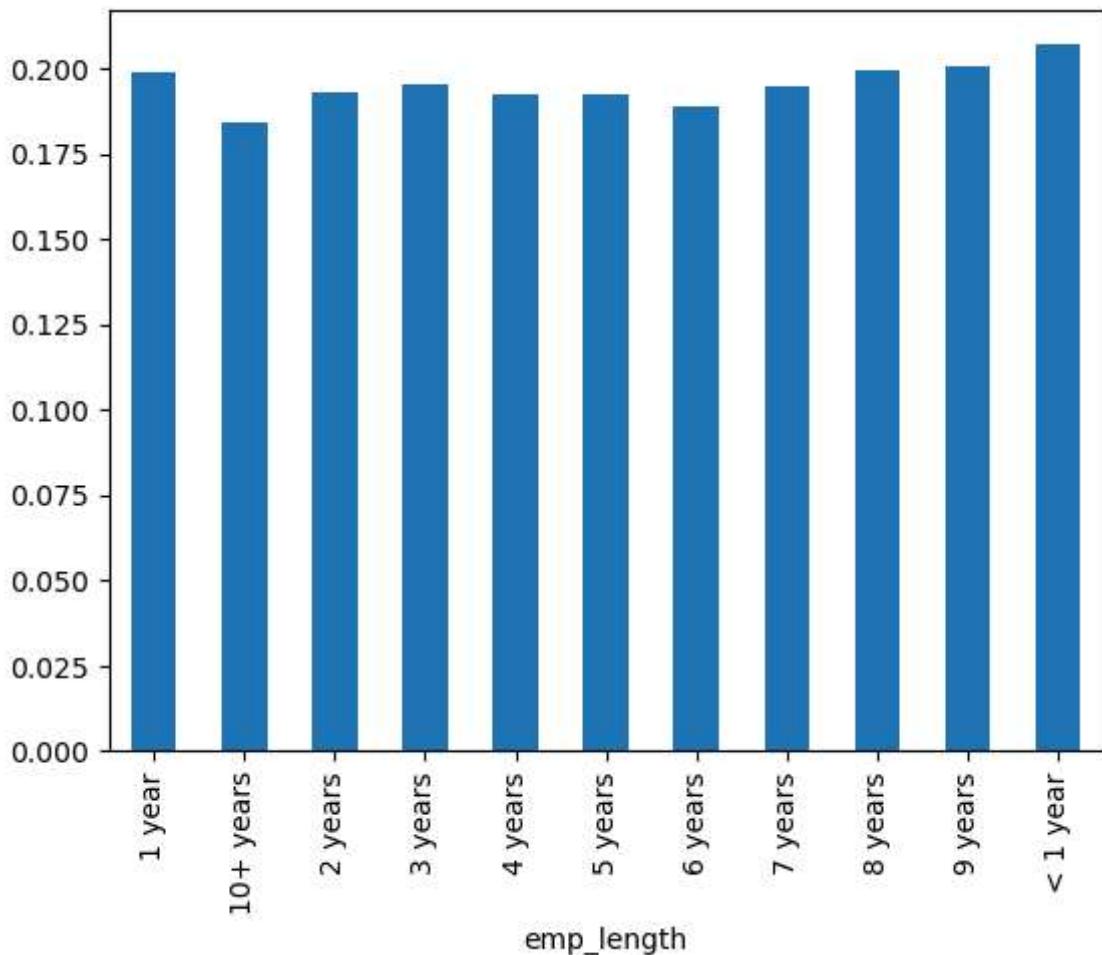
```
Out[40]: <Axes: xlabel='emp_length', ylabel='count'>
```



This current countplot doesn't tell us whether there is a strong relationship between employment length and being charged off. We want to visualise the percentage of charge off's per category. So we will get the percent of people per employer category that didn't pay back their loan.

```
In [41]: emp_co = df[df['loan_status']=='Charged Off'].groupby('emp_length').count()['loan_status']
In [42]: emp_fp = df[df['loan_status']=='Fully Paid'].groupby('emp_length').count()['loan_status']
In [43]: emp_len = emp_co/(emp_fp+emp_co)
emp_len.plot(kind='bar')
```

Out[43]: <Axes: xlabel='emp_length'>



```
In [44]: # They are all about the same...so we should just drop the emp_Length feature  
df = df.drop('emp_length', axis=1)
```

```
In [45]: df.isnull().sum()
```

```
Out[45]: 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                1756
          dti                  0
          earliest_cr_line      0
          open_acc              0
          pub_rec               0
          revol_bal             0
          revol_util            276
          total_acc              0
          initial_list_status    0
          application_type       0
          mort_acc              37795
          pub_rec_bankruptcies   535
          address              0
          loan_repaid            0
          dtype: int64
```

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

```
In [46]: df['purpose'].head(10)
feat_info('purpose')
```

A category provided by the borrower for the loan request.

```
In [47]: df['title'].head(10)
feat_info('title')
```

The loan title provided by the borrower

```
In [48]: #Title is just a string, can be removed
df = df.drop('title',axis=1)
```

NOTE: This is one of the hardest parts of the project! Refer to the solutions video if you need guidance, feel free to fill or drop the missing values of the mort_acc however you see fit! Here we're going with a very specific approach.

TASK: Find out what the mort_acc feature represents

```
In [49]: feat_info('mort_acc')
```

Number of mortgage accounts.

TASK: Create a value_counts of the mort_acc column.

```
In [50]: df['mort_acc'].value_counts()
```

```
Out[50]: mort_acc
0.0      139777
1.0      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     37
17.0     22
18.0     18
19.0     15
20.0     13
24.0     10
22.0     7
21.0     4
25.0     4
27.0     3
26.0     2
32.0     2
31.0     2
23.0     2
34.0     1
28.0     1
30.0     1
Name: count, 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`

```
In [51]: print("Correlation with the mort_acc column")
df.select_dtypes(include=['number']).corr()['mort_acc'].sort_values()
# Here we are checking the correlation between the mort_acc column and the rest
# we can see there is a relatively strong correlation between mort_acc and tota
```

Correlation with the mort_acc column

```
Out[51]: int_rate           -0.082583
          dti                -0.025439
          revol_util          0.007514
          pub_rec              0.011552
          pub_rec_bankruptcies 0.027239
          loan_repaid          0.073111
          open_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 `total_acc` feature 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:

```
In [52]: print("Mean of mort_acc column per total_acc")
df.select_dtypes(include=['number']).groupby('total_acc').mean()['mort_acc']

Mean of mort_acc column per total_acc

Out[52]: total_acc
2.0      0.000000
3.0      0.052023
4.0      0.066743
5.0      0.103289
6.0      0.151293
...
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

```
In [53]: total_acc_avg = df.select_dtypes(include=['number']).groupby('total_acc').mean()
# what this line of code does:

#Groups the DataFrame by each unique value of total_acc, so ALL rows that share
#total_acc value are grouped together.

#Calculates the mean of every numeric column within each total_acc group.
#So for each total_acc value (say 2, 3, 4...), you get the average of mort_acc,
#loan_amnt, and all other numeric features.
```

```
#From that grouped-and-averaged result, you select only the mort_acc column.
```

In [54]: `total_acc_avg[2.0]`

Out[54]: `np.float64(0.0)`

In [55]: `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
#groupby averages of mort_acc per total_acc values.
```

```
if np.isnan(mort_acc):
    return total_acc_avg[total_acc]
else:
    return mort_acc
```

```
# so essentially, this means:
```

```
#Take the total_acc and mort_acc values from this specific row, and feed them in
```

In [56]: `df['mort_acc'] = df.apply(lambda x: fill_mort_acc(x['total_acc'], x['mort_acc']))`

In [57]: `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
dti	0
earliest_cr_line	0
open_acc	0
pub_rec	0
revol_bal	0
revol_util	276
total_acc	0
initial_list_status	0
application_type	0
mort_acc	0
pub_rec_bankruptcies	535
address	0
loan_repaid	0
dtype: int64	

TASK: 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().

```
In [58]: df = df.dropna()
```

```
In [59]: df.isnull().sum()
```

```
Out[59]: 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
dti            0
earliest_cr_line 0
open_acc        0
pub_rec         0
revol_bal       0
revol_util      0
total_acc       0
initial_list_status 0
application_type 0
mort_acc        0
pub_rec_bankruptcies 0
address         0
loan_repaid     0
dtype: int64
```

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

```
In [60]: df.select_dtypes(['object']).columns
```

```
Out[60]: 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().

```
In [61]: df['term'].value_counts()
```

```
Out[61]: term
36 months    301247
60 months     93972
Name: count, dtype: int64
```

```
In [62]: # Or just use .map()
df['term'] = df['term'].apply(lambda term: int(term[:3])) #this takes the first
df['term']
```

```
Out[62]: 0      36
1      36
2      36
3      36
4      60
      ..
396025   60
396026   36
396027   36
396028   60
396029   36
Name: term, Length: 395219, dtype: int64
```

grade feature

TASK: We already know grade is part of sub_grade, so just drop the grade feature.

```
In [63]: df = df.drop('grade', axis=1)
```

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.

```
In [64]: subgrade_dummies = pd.get_dummies(df['sub_grade'], drop_first=True)
```

```
In [65]: df = pd.concat([df.drop('sub_grade', axis=1), subgrade_dummies], axis=1)
# so, Take my original DataFrame, remove the sub_grade column,
# add the new dummy columns instead, and store the result as the updated DataFra
```

```
In [66]: df.columns
```

```
Out[66]: 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')
```

```
In [67]: df.select_dtypes(['object']).columns
```

```
Out[67]: 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.

```
In [68]: ## The bottom code does the same thing as above but for the
#verification_status, application_type,initial_list_status,purpose columns.
#However, our data frame df already has these concatenated, so code wont run
```

```
#CODE
dummies = pd.get_dummies(df[['verification_status', 'application_type','initial_list_status','purpose']])
df = df.drop(['verification_status', 'application_type','initial_list_status','purpose'])
df = pd.concat([df,dummies],axis=1)
```

home_ownership

TASK: Review the value_counts for the home_ownership column.

```
In [69]: df['home_ownership'].value_counts()
```

```
Out[69]: home_ownership
MORTGAGE    198022
RENT        159395
OWN         37660
OTHER        110
NONE         29
ANY           3
Name: count, 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.

```
In [70]: 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)
```

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.

```
In [71]: df['zip_code'] = df['address'].apply(lambda address:address[-5:])
```

TASK: Now make this zip_code column into dummy variables using pandas. Concatenate the result and drop the original zip_code column along with dropping the address column.

```
In [72]: dummies = pd.get_dummies(df['zip_code'],drop_first=True)
df = df.drop(['zip_code','address'],axis=1)
df = pd.concat([df,dummies],axis=1)
```

issue_d

TASK: This would be data leakage, we wouldn't know beforehand whether or not a loan would be issued when using our model, so in theory we wouldn't have an issue_date, drop this feature.

```
In [73]: df = df.drop('issue_d',axis=1)
```

earliest_cr_line

TASK: This appears to be a historical time stamp feature. Extract the year from this feature using a .apply function, then convert it to a numeric feature. Set this new data to a feature column called 'earliest_cr_year'.Then drop the earliest_cr_line feature.

```
In [74]: df['earliest_cr_year'] = df['earliest_cr_line'].apply(lambda date:int(date[-4:]))
df = df.drop('earliest_cr_line',axis=1)
```

```
In [75]: df.select_dtypes(['object']).columns
```

```
Out[75]: Index(['loan_status'], dtype='object')
```

Train Test Split

TASK: Import train_test_split from sklearn.

In [97]: `from sklearn.model_selection import train_test_split`

TASK: drop the load_status column we created earlier, since its a duplicate of the loan_repaid column. We'll use the loan_repaid column since its already in 0s and 1s.

In [98]: `df = df.drop('loan_status', axis=1)`

```

KeyError                                         Traceback (most recent call last)
Cell In[98], line 1
----> 1 df = df.drop(                  ,axis=1)

File ~/anaconda3/envs/Data_Science/lib/python3.12/site-packages/pandas/core/fram
e.py:5588, in DataFrame.drop(self, labels, axis, index, columns, level, inplace,
errors)
    5440 def drop(
    5441     self,
    5442     labels: IndexLabel | None = None,
    (... ) 5449     errors: IgnoreRaise = "raise",
    5450 ) -> DataFrame | None:
    5451     """
    5452     Drop specified labels from rows or columns.
    5453
    (...) 5586             weight  1.0      0.8
    5587     """
-> 5588     return super().drop(
    5589         labels=labels,
    5590         axis=axis,
    5591         index=index,
    5592         columns=columns,
    5593         level=level,
    5594         inplace=inplace,
    5595         errors=errors,
    5596     )

File ~/anaconda3/envs/Data_Science/lib/python3.12/site-packages/pandas/core/gener
ic.py:4807, in NDFrame.drop(self, labels, axis, index, columns, level, inplace, e
rrors)
    4805 for axis, labels in axes.items():
    4806     if labels is not None:
-> 4807         obj = obj._drop_axis(labels, axis, level=level, errors=errors)
    4809 if inplace:
    4810     self._update_inplace(obj)

File ~/anaconda3/envs/Data_Science/lib/python3.12/site-packages/pandas/core/gener
ic.py:4849, in NDFrame._drop_axis(self, labels, axis, level, errors, only_slice)
    4847     new_axis = axis.drop(labels, level=level, errors=errors)
    4848 else:
-> 4849     new_axis = axis.drop(labels, errors=errors)
    4850 indexer = axis.get_indexer(new_axis)
    4852 # Case for non-unique axis
    4853 else:

File ~/anaconda3/envs/Data_Science/lib/python3.12/site-packages/pandas/core/index
es/base.py:7136, in Index.drop(self, labels, errors)
    7134 if mask.any():
    7135     if errors != "ignore":
-> 7136         raise KeyError(f"{labels[mask].tolist()} not found in axis")
    7137 indexer = indexer[~mask]
    7138 return self.delete(indexer)

KeyError: "['loan_status'] not found in axis"

```

TASK: Set X and y variables to the .values of the features and label.

```
In [99]: X = df.drop('loan_repaid',axis=1).values
y = df['loan_repaid'].values
```

TASK: Perform a train/test split with test_size=0.2 and a random_state of 101.

```
In [100... X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20, random
```

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.

```
In [129... from sklearn.preprocessing import MinMaxScaler
```

```
In [130... scaler = MinMaxScaler()
```

```
In [131... X_train = scaler.fit_transform(X_train)
```

```
In [132... X_test = scaler.transform(X_test)
```

Creating the Model

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

```
In [105... 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
```

TASK: Build a sequential model to will be trained on the data. You have unlimited options here, but here is what the solution uses: a model that goes 78 --> 39 --> 19--> 1 output neuron. OPTIONAL: Explore adding Dropout layers 1) 2

```
In [109... #We are going to create our mode:

# Choose whatever number of Layers/neurons you want.

# https://stats.stackexchange.com/questions/181/how-to-choose-the-number-of-hidden-layers-in-a-neural-network

# Remember to compile()
X_train.shape
#This is so we can see how many layers to use.
# Generally, our very first Layer should match the number of features
```

```
Out[109... (316175, 78)
```

```
In [107... model = Sequential()

# input layer
model.add(Dense(78, activation='relu')) # First Layer matches number of feature
model.add(Dropout(0.2)) # To prevent overfitting, 20% of the data is dropped and
```

```
# 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')
```

TASK: Fit the model to the training data for at least 25 epochs. Also add in the validation data for later plotting. Optional: add in a batch_size of 256.

In [110...]

```
model.fit(x=X_train,
           y=y_train,
           epochs=25,
           batch_size=256,
           validation_data=(X_test, y_test),
           )
```

```

Epoch 1/25
1236/1236 2s 772us/step - loss: 0.3639 - val_loss: 0.2648
Epoch 2/25
1236/1236 1s 722us/step - loss: 0.2663 - val_loss: 0.2628
Epoch 3/25
1236/1236 1s 733us/step - loss: 0.2638 - val_loss: 0.2623
Epoch 4/25
1236/1236 1s 760us/step - loss: 0.2620 - val_loss: 0.2620
Epoch 5/25
1236/1236 1s 738us/step - loss: 0.2618 - val_loss: 0.2635
Epoch 6/25
1236/1236 1s 749us/step - loss: 0.2600 - val_loss: 0.2618
Epoch 7/25
1236/1236 1s 722us/step - loss: 0.2592 - val_loss: 0.2620
Epoch 8/25
1236/1236 1s 716us/step - loss: 0.2596 - val_loss: 0.2619
Epoch 9/25
1236/1236 1s 738us/step - loss: 0.2593 - val_loss: 0.2616
Epoch 10/25
1236/1236 1s 724us/step - loss: 0.2582 - val_loss: 0.2620
Epoch 11/25
1236/1236 1s 720us/step - loss: 0.2590 - val_loss: 0.2611
Epoch 12/25
1236/1236 1s 718us/step - loss: 0.2582 - val_loss: 0.2615
Epoch 13/25
1236/1236 1s 734us/step - loss: 0.2571 - val_loss: 0.2611
Epoch 14/25
1236/1236 1s 719us/step - loss: 0.2577 - val_loss: 0.2619
Epoch 15/25
1236/1236 1s 713us/step - loss: 0.2584 - val_loss: 0.2612
Epoch 16/25
1236/1236 1s 717us/step - loss: 0.2568 - val_loss: 0.2615
Epoch 17/25
1236/1236 1s 711us/step - loss: 0.2574 - val_loss: 0.2611
Epoch 18/25
1236/1236 1s 719us/step - loss: 0.2550 - val_loss: 0.2616
Epoch 19/25
1236/1236 1s 717us/step - loss: 0.2572 - val_loss: 0.2613
Epoch 20/25
1236/1236 1s 724us/step - loss: 0.2558 - val_loss: 0.2612
Epoch 21/25
1236/1236 1s 715us/step - loss: 0.2567 - val_loss: 0.2620
Epoch 22/25
1236/1236 1s 729us/step - loss: 0.2569 - val_loss: 0.2616
Epoch 23/25
1236/1236 1s 717us/step - loss: 0.2565 - val_loss: 0.2621
Epoch 24/25
1236/1236 1s 717us/step - loss: 0.2554 - val_loss: 0.2616
Epoch 25/25
1236/1236 1s 717us/step - loss: 0.2545 - val_loss: 0.2616

```

Out[110...]: <keras.src.callbacks.history.History at 0x72ef5c7d6180>

TASK: OPTIONAL: Save your model.

In [111...]: `from tensorflow.keras.models import load_model`

In [112...]: `model.save('full_data_project_model.h5')`

WARNING:absl:You are saving your model as an HDF5 file via `model.save()` or `keras.saving.save_model(model)`. This file format is considered legacy. We recommend using instead the native Keras format, e.g. `model.save('my_model.keras')` or `keras.saving.save_model(model, 'my_model.keras')`.

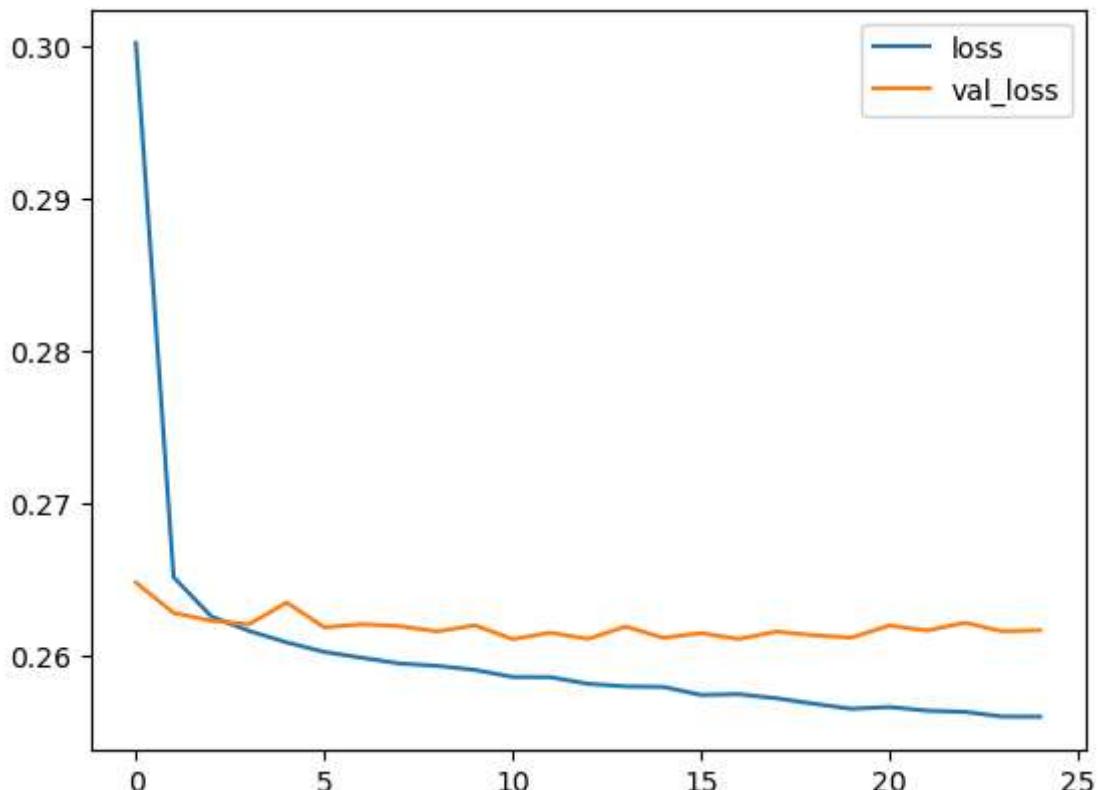
Section 3: Evaluating Model Performance.

TASK: Plot out the validation loss versus the training loss.

```
In [113...]: losses = pd.DataFrame(model.history.history)
```

```
In [114...]: losses[['loss', 'val_loss']].plot()
```

```
Out[114...]: <Axes: >
```



TASK: Create predictions from the X_test set and display a classification report and confusion matrix for the X_test set.

```
In [115...]: from sklearn.metrics import classification_report,confusion_matrix
```

```
In [117...]: predictions = (model.predict(X_test) > 0.5).astype("int32")
```

```
2471/2471 ━━━━━━━━ 1s 218us/step
```

```
In [123...]: print(classification_report(y_test,predictions))
```

	precision	recall	f1-score	support
0	0.97	0.44	0.61	15658
1	0.88	1.00	0.93	63386
accuracy			0.89	79044
macro avg	0.93	0.72	0.77	79044
weighted avg	0.90	0.89	0.87	79044

f1 score on 0 class is a good representation of this model. 80% of the data contains 'loan paid back' values, so the performance of the model in predicting that they paid their loan back is not important as most did pay it back

In [124...]: `confusion_matrix(y_test, predictions)`

Out[124...]: `array([[6918, 8740],
 [186, 63200]])`

TASK: Given the customer below, would you offer this person a loan?

```
import random
random.seed(101)
random_ind = random.randint(0, len(df))

new_customer = df.drop('loan_repaid', axis=1).iloc[random_ind]
new_customer
```

Out[125...]:

loan_amnt	25000.0
term	60
int_rate	18.24
installment	638.11
annual_inc	61665.0
	...
48052	False
70466	False
86630	False
93700	False
earliest_cr_year	1996
Name:	305323, Length: 78, dtype: object

In [134...]: `new_customer = scaler.transform(new_customer.values.reshape(1, 78))`

In [137...]: `prediction = (model.predict(new_customer) > 0.5).astype("int32")`
prediction

1/1 ————— 0s 13ms/step

Out[137...]: `array([[1]], dtype=int32)`

TASK: Now check, did this person actually end up paying back their loan?

In [138...]: `df.iloc[random_ind]['loan_repaid']`

Out[138...]: `np.int64(1)`

GREAT JOB!