

Imports

```
In [21]: from sqlalchemy import create_engine
import pymysql
import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
from matplotlib.colors import mcolorm
from matplotlib import cm
import pandas_profiling
import seaborn as sns
import (color, color_m, font, scale, 2)
# Import the two methods from heatmap library
from heatmap import heatmap, corplot
```

Start with cleaned data from C272A notebook

```
In [22]: df = pd.read_csv('CreditOne_cleaned_formatted.csv', index_col=0)
df.rename(columns={'PAY_0': 'PAY_1', 'formated': 'true'})
df.head()
```

Out [22]:

ID	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_1	PAY_2	PAY_3	PAY_4	...	BILL_AMT1	BILL_AMT2	BILL_AMT3	PAY_5
0	1	20000	female	university	-1	24	2	2	-1	...	0	0	0	0
1	2	120000	female	university	-2	26	-1	0	2	...	3272	3458	3261	0
2	3	80000	female	university	2	34	0	0	0	...	14331	14948	15549	0

3 rows * 25 columns

```
In [23]: df.dtypes.head()
```

Out [23]:

```
ID          int64
LIMIT_BAL  int64
SEX         object
EDUCATION  object
MARRIAGE   int64
dtype: object
```

Pandas profile may be more meaningful now that the data has been cleaned. May give me a good place to start my EDA

```
In [24]: rpt = pandas_profiling.ProfileReport(df)
rpt.to_file(output_file='CreditOne_ProfileReport_CleanedData.html')
```

Use Dummy to split object columns

```
In [25]: df = pd.get_dummies(df)
```

```
In [26]: df.head()
```

Out [26]:

ID	LIMIT_BAL	MARRIAGE	AGE	PAY_1	PAY_2	PAY_3	PAY_4	PAY_5	PAY_6	PAY_7	SEX_female	SEX_male
0	1	20000	1	24	-1	-1	-1	-2	-2	0	0	1
1	2	120000	2	26	-1	-2	0	0	2	0	2000	1
2	3	80000	2	34	0	0	0	0	0	1000	5000	0
3	4	50000	1	37	0	0	0	0	0	1089	1000	1
4	5	50000	1	57	-1	0	-1	0	0	689	679	0

5 rows * 30 columns

```
In [27]: df.dtypes
```

Out [27]:

```
ID          int64
LIMIT_BAL  int64
MARRIAGE   int64
AGE         int64
PAY_1      int64
PAY_2      int64
PAY_3      int64
PAY_4      int64
PAY_5      int64
PAY_6      int64
PAY_7      int64
SEX_female  uint8
SEX_male    uint8
EDUCATION_graduate_school  uint8
EDUCATION_high_school      uint8
EDUCATION_other            uint8
EDUCATION_university        uint8
DEFAULT_default            uint8
DEFAULT_not_default         uint8
dtype: object
```

EDA

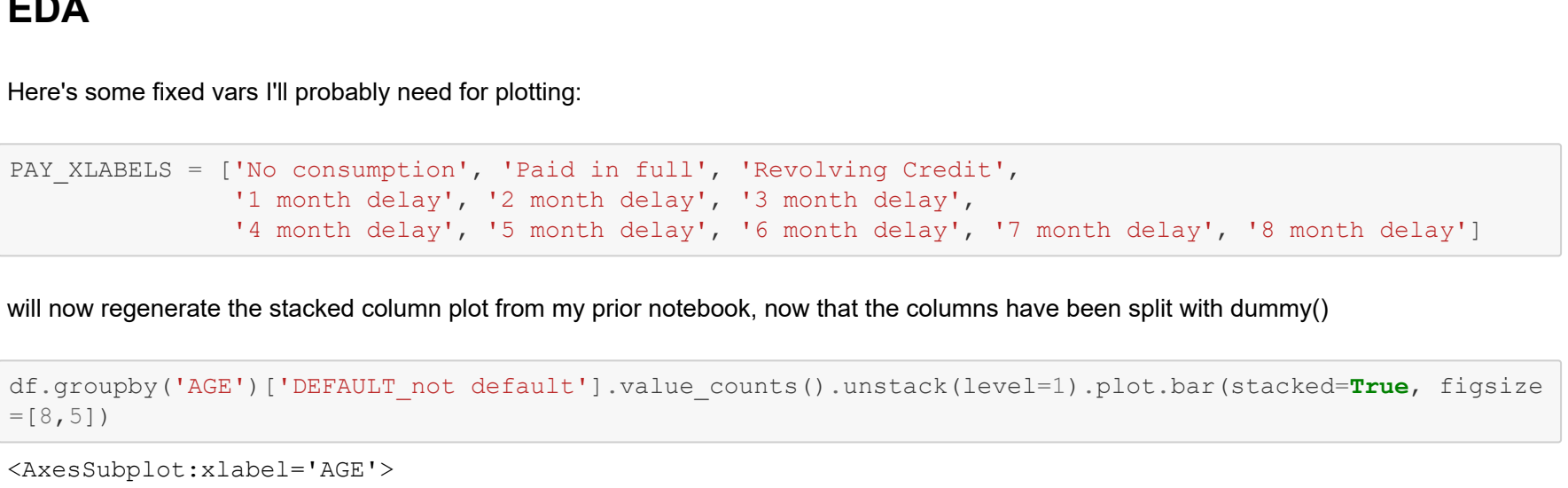
Here's some fixed vars I'll probably need for plotting:

```
In [28]: PAY_VARIABLES = ['No consumption', '95id in full', 'Revolving Credit',
                        '1 month delay', '2 month delay', '3 month delay',
                        '4 month delay', '5 month delay', '6 month delay', '7 month delay', '8 month delay']
```

Will now regenerate the stacked column plot from my prior notebook, now that the columns have been split with dummy()

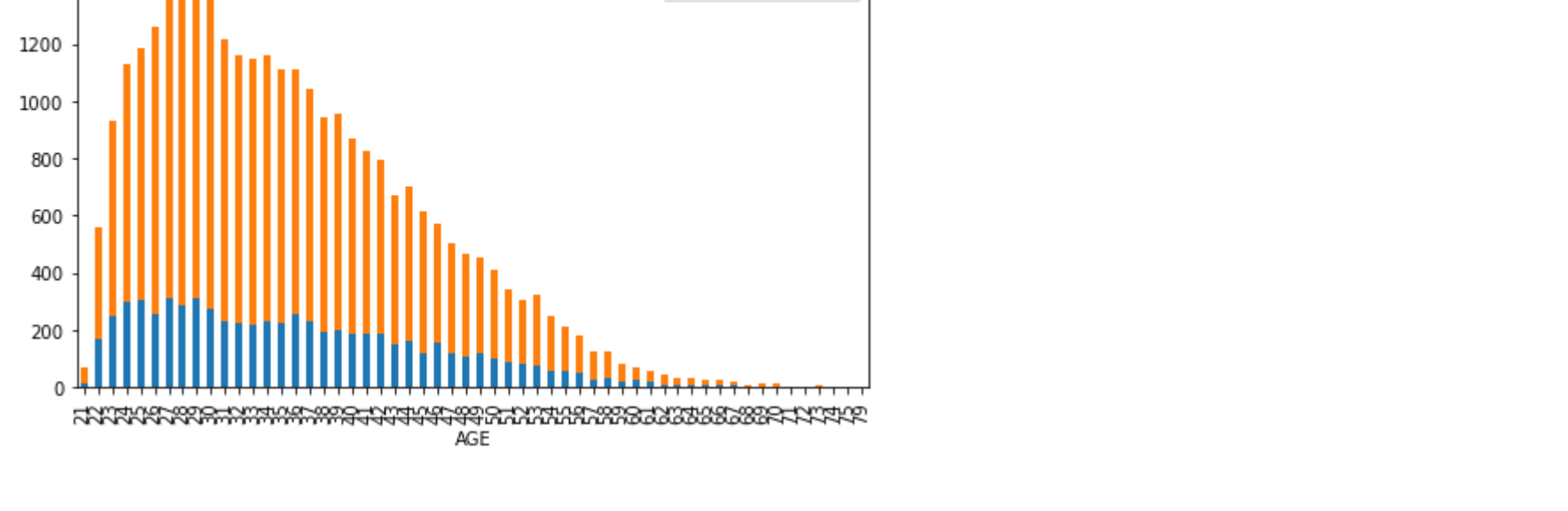
```
In [29]: df.groupby('AGE')['DEFAULT_not default'].value_counts().unstack(level=1).plot.bar(stacked=True, figsize=[8,5])
```

```
Out [29]: %axesSubplot(xlabel='AGE')>
```



Plots from the Project Roadmap, for reference

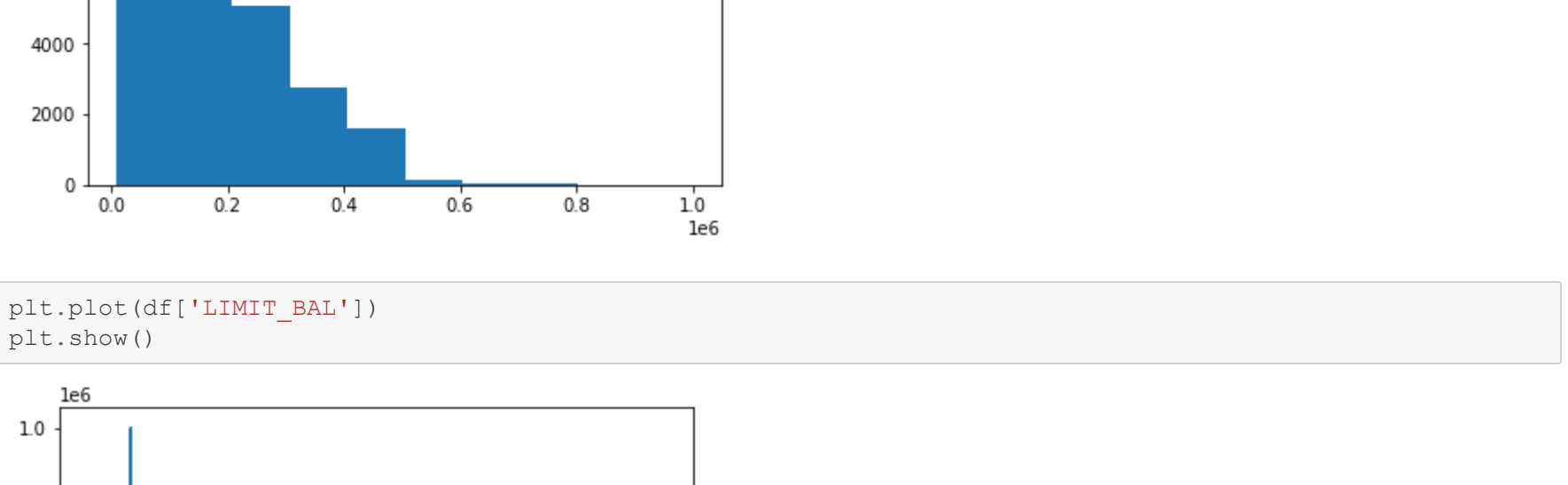
```
In [30]: default_x=6.4
plt.rcParams['figure.figsize'] = (default_x, default_y)
plt.hist(df['LIMIT_BAL'], bins=10)
```



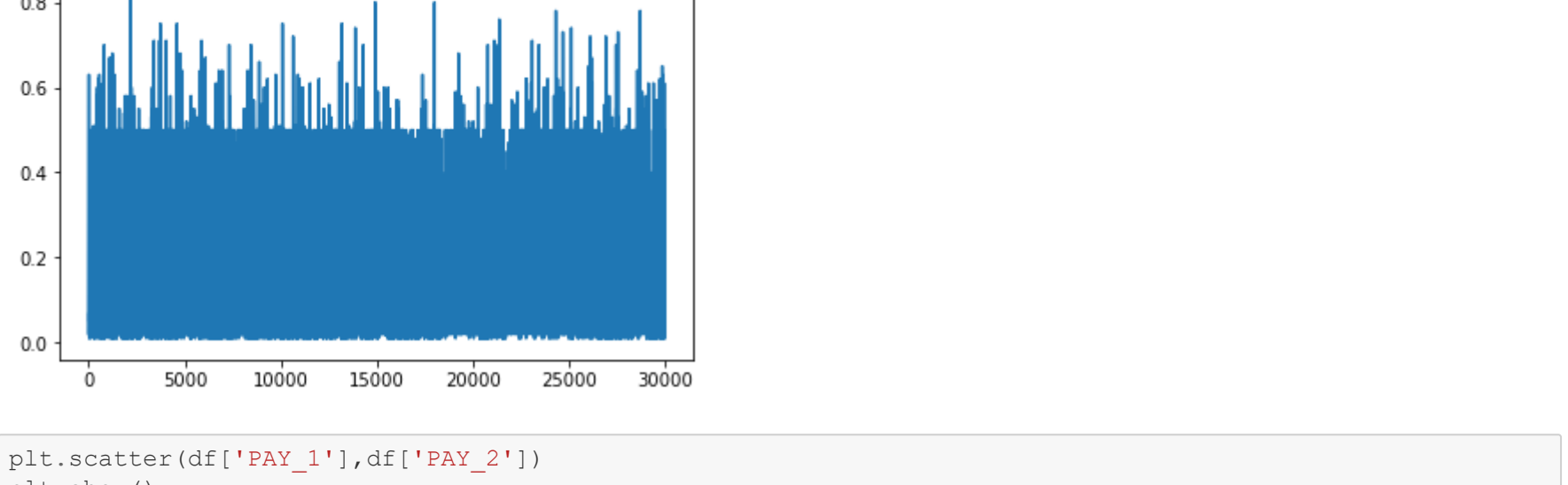
```
In [31]: plt.plot(df['LIMIT_BAL'])
```



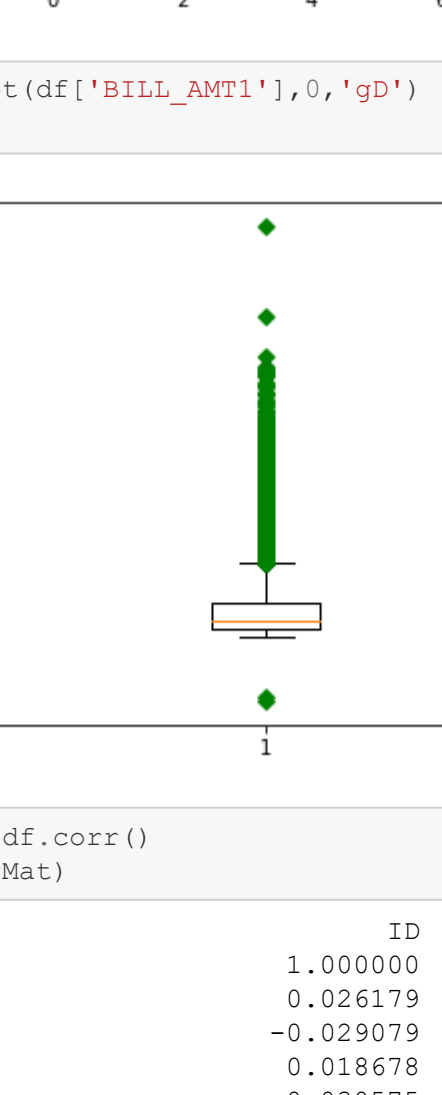
```
In [32]: plt.scatter(df['PAY_1'], df['PAY_2'])
```



```
In [33]: plt.boxplot(df['BILL_AMT1'], 0, 'gp')
```



```
In [34]: corrMat = df.corr()
print(corrMat)
```

<pre>plt.boxplot(df[['BILL_AMT1'],0,'gp']) plt.show()</pre>									
									
<pre>corrMat = df.corr() print(corrMat)</pre>									
ID	LIMIT_BAL	MARRIAGE	AGE	PAY_1 \					
LIMIT_BAL	1.000000	-0.026179	-0.029079	0.018678	-0.030575				
MARRIAGE	0.026179	1.000000	-0.108139	0.144713	-0.271214				
AGE	-0.029079	-0.108139	1.000000	-0.414170	-0.039447				
PAY_1	0.018678	0.144713	-0.414170	1.000000	-0.039447				
PAY_2	-0.030575	-0.271214	0.030197	-0.039447	1.000000				
PAY_3	-0.012115	-0.296382	0.024195	-0.050148	0.072164				
PAY_4	-0.018494	-0.286123	0.032688	-0.053048	0.574245				
PAY_5	-0.002735	-0.267460	0.033122	-0.049722	0.538841				
PAY_6	-0.022139	-0.249411	0.035629	-0.053826	0.509426				
BILL_AMT1	-0.002070	-0.235195	0.034345	-0.048773	0.474553				
BILL_AMT2	0.019389	0.285430	-0.023472	0.056239	-0.187068				
BILL_AMT3	0.007482	0.278314	-0.021602	0.054283	-0.189859				
BILL_AMT4	0.024354	0.283236	-0.024909	0.053710	-0.179785				
BILL_AMT5	0.040351	0.293988	-0.023344	0.051353	-0.179125				
BILL_AMT6	0.016705	0.295562	-0.023393	0.049345	-0.180635				
PAY_AMT1	0.016730	0.290389	-0.021207	0.047613	-0.176980				
PAY_AMT2	0.009742	0.195236	-0.005979	0.026147	-0.079269				
PAY_AMT3	0.008406	0.178408	-0.008093	0.021785	-0.070101				
PAY_AMT4	0.039315	0.210167	-0.003541	0.029247	-0.070561				
PAY_AMT5	0.007793	0.203242	-0.012659	0.012379	-0.064005				
PAY_AMT6	0.000652	0.212702	-0.001205	0.022850	-0.058193				
SEX_female	0.003000	0.215955	-0.006641	0.019478	-0.058673				
SEX_male	0.018497	0.024755	-0.031389	-0.090874	-0.057643				
EDUCATION_graduate_school	-0.018497	-0.024755	-0.031389	0.090874	0.057643				
EDUCATION_high_school	-0.025858	0.258777	0.142129	-0.104923	-0.142720				
EDUCATION_university	0.017149	-0.139686	-0.110845	0.231252	0.058902				
EDUCATION_other	0.037165	-0.013420	-0.009836	0.008982	-0.024937				
EDUCATION_university	0.002569	-0.147530	-0.051797	-0.077626	0.099177				
DEFAULT_default	0.013952	-0.153520	-0.024339	0.013890	0.324794				
DEFAULT_not_default	0.013952	0.153520	0.024339	-0.013890	-0.324794				

ID	PAY_2	PAY_3	PAY_4	PAY_5	PAY_6
LIMIT_BAL	-0.012115	-0.018494	-0.023472	-0.022139	-0.020270
MARRIAGE	-0.296382	-0.286123	-0.267460	-0.249411	-0.235195
AGE	0.024195	0.032688	0.033122	0.035629	0.034345
PAY_1	-0.050419	-0.053048	-0.049722	-0.050826	-0.048773
PAY_2	0.072164	0.574245	0.538841	0.509426	0.474553
PAY_3	1.000000	0.766552	0.662067	0.622780	0.575901
PAY_4	0.766552	1.000000	0.777359	0.681985	0.636186
PAY_5	0.662067	0.777359	1.000000	0.819835	0.716449
PAY_6	0.622780	0.686775	0.819835	1.000000	0.816900
BILL_AMT1	0.575901	0.632684	0.716449	0.104923	0.058673
BILL_AMT2	0.234887	0.208473	0.202812	0.206684	0.207368
BILL_AMT3	0.235257	0.237295	0.225816	0.226913	0.226924
BILL_AMT4	0.224146	0.227494	0.244983	0.243335	0.241181
BILL_AMT5	0.225237	0.227022	0.245917	0.243915	0.246356
BILL_AMT6	0.221348	0.225145	0.242902	0.246783	0.245094
PAY_AMT1	0.008070	0.022327	0.239154	0.262509	0.285091
PAY_AMT2	0.080701	0.001295	-0.003662	-0.006089	-0.001496
PAY_AMT3	-0.058990	-0.067693	-0.001944	-0.003961	-0.005223
PAY_AMT4	-0.053311	-0.069231	-0.069235	0.000602	0.005836
PAY_AMT5	0.046858	-0.046067	-0.043461	-0.059239	-0.057618
PAY_AMT6	-0.037093	-0.035863	-0.033590	-0.031337	-0.046434
SEX_female	-0.036500	-0.035861	-0.026655	-0.023027	-0.025299
SEX_male	-0.070771	-0.066096	-0.060173	-0.050564	-0.044008
EDUCATION_graduate_school	0.070771	0.066096	0.060173	0.050564	0.044008
EDUCATION_high_school	-0.169215	-0.160209	-0.152402	-0.138709	-0.125123
EDUCATION_university	0.064590	0.062461	0.059882	0.049577	0.041370
EDUCATION_university	0.033118	-0.034435	-0.030998	-0.028822	-0.037533
DEFAULT_default	0.122364	0.115644	0.110340	0.103218	0.098013
DEFAULT_not_default	0.263551	0.235253	0.216614	0.204149	0.186666
	-0.263551	-0.235253	-0.216614	-0.204149	-0.186666

ID	PAY_AMT5	PAY_AMT6	SEX_female	SEX_male \	
LIMIT_BAL	
MARRIAGE	
AGE	
PAY_1	
PAY_2	
PAY_3	
PAY_4	
PAY_5	
PAY_6	
BILL_AMT1	
BILL_AMT2	
BILL_AMT3	
BILL_AMT4	
BILL_AMT5	
BILL_AMT6	
PAY_AMT1	
PAY_AMT2	
PAY_AMT3	
PAY_AMT4	
PAY_AMT5	
PAY_AMT6	
SEX_female	
SEX_male	
EDUCATION_graduate_school	
EDUCATION_high_school	
EDUCATION_university	
EDUCATION_university	
DEFAULT_default	
DEFAULT_not_default	

ID	PAY_AMT5	PAY_AMT6	SEX_female	SEX_male \	
LIMIT_BAL	
MARRIAGE	
AGE	
PAY_1	
PAY_2	
PAY_3	
PAY_4	
PAY_5	
PAY_6	
BILL_AMT1	
BILL_AMT2	
BILL_AMT3	
BILL_AMT4	
BILL_AMT5	
BILL_AMT6	
PAY_AMT1	
PAY_AMT2	
PAY_AMT3	
PAY_AMT4	
PAY_AMT5	
PAY_AMT6	
SEX_female	
SEX_male	
EDUCATION_graduate_school	
EDUCATION_high_school	
EDUCATION_university	
EDUCATION_university	
DEFAULT_default	
DEFAULT_not_default	

ID	EDUCATION_graduate_school	EDUCATION_high_school \	
LIMIT_BAL	
MARRIAGE	
AGE	
PAY_1	
PAY_2	
PAY_3	
PAY_4	
PAY_5	
PAY_6	
BILL_AMT1	
BILL_AMT2	
BILL_AMT3	
BILL_AMT4	
BILL_AMT5	
BILL_AMT6	
PAY_AMT1	
PAY_AMT2	
PAY_AMT3	
PAY_AMT4	
PAY_AMT5	
PAY_AMT6	
SEX_female	
SEX_male	
EDUCATION_graduate_school	
EDUCATION_high_school	
EDUCATION_university	
EDUCATION_university	
DEFAULT_default	
DEFAULT_not_default	

ID	EDUCATION_graduate_school	EDUCATION_high_school \	
LIMIT_BAL	
MARRIAGE	
AGE	
PAY_1	
PAY_2	
PAY_3	
PAY_4	
PAY_5	
PAY_6	
BILL_AMT1	
BILL_AMT2	
BILL_AMT3	
BILL_AMT4	
BILL_AMT5	
BILL_AMT6	
PAY_AMT1	
PAY_AMT2	
PAY_AMT3	
PAY_AMT4	
PAY_AMT5	
PAY_AMT6	
SEX_female	
SEX_male	
EDUCATION_graduate_school	
EDUCATION_high_school	
EDUCATION_university	
EDUCATION_university	
DEFAULT_default	
DEFAULT_not_default	

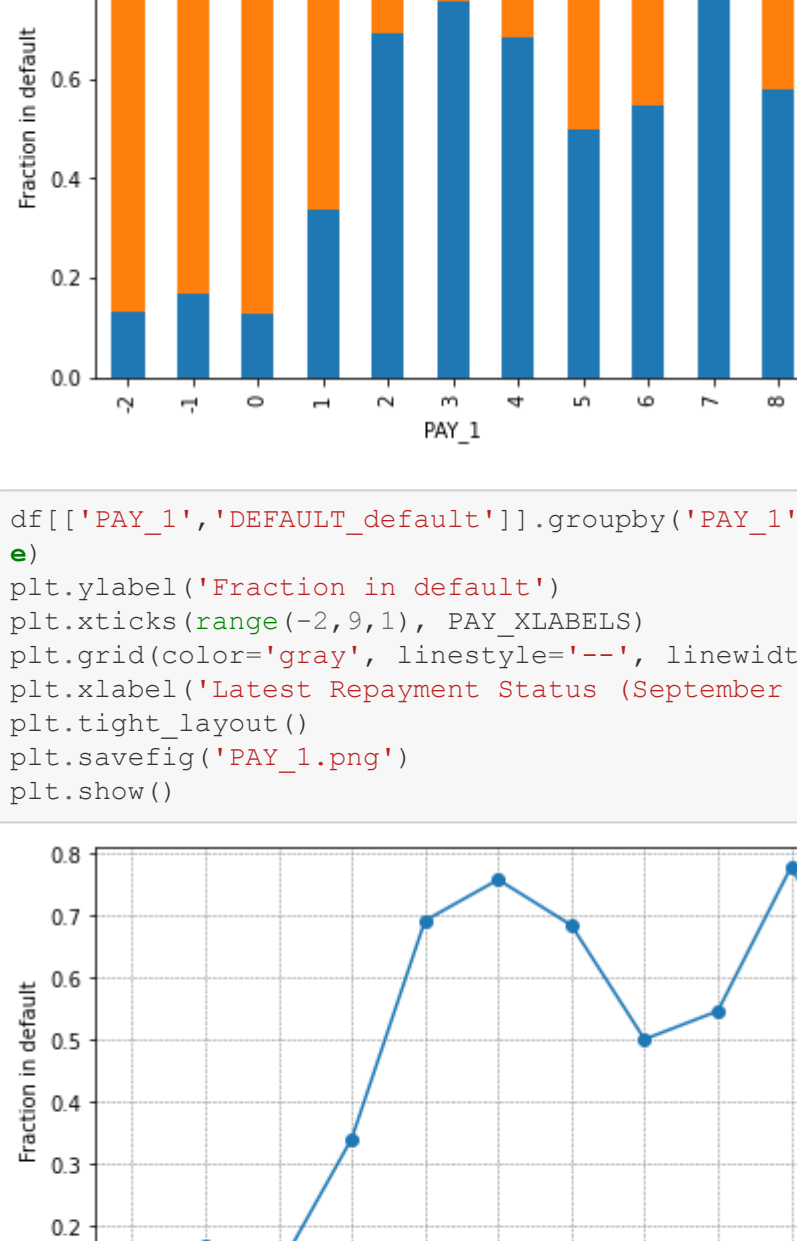
ID	EDUCATION_graduate_school	EDUCATION_high_school \	
LIMIT_BAL	
MARRIAGE	
AGE	
PAY_1	
PAY_2	
PAY_3	
PAY_4	
PAY_5	
PAY_6	
BILL_AMT1	
BILL_AMT2	
BILL_AMT3	
BILL_AMT4	
BILL_AMT5	
BILL_AMT6	
PAY_AMT1	
PAY_AMT2	
PAY_AMT3	
PAY_AMT4	
PAY_AMT5	
PAY_AMT6	
SEX_female	
SEX_male	
EDUCATION_graduate_school	
EDUCATION_high_school	
EDUCATION_university	
EDUCATION_university	
DEFAULT_default	
DEFAULT_not_default	

ID	EDUCATION_graduate_school	EDUCATION_high_school \	
LIMIT_BAL	
MARRIAGE	
AGE	
PAY_1	
PAY_2	
PAY_3	
PAY_4	
PAY_5	
PAY_6	
BILL_AMT1	
BILL_AMT2	
BILL_AMT3	
BILL_AMT4	
BILL_AMT5	
BILL_AMT6	
PAY_AMT1	
PAY_AMT2	
PAY_AMT3	
PAY_AMT4	
PAY_AMT5	
PAY_AMT6	
SEX_female	
SEX_male	
EDUCATION_graduate_school	
EDUCATION_high_school	
EDUCATION_university	
EDUCATION_university	
DEFAULT_default	
DEFAULT_not_default	

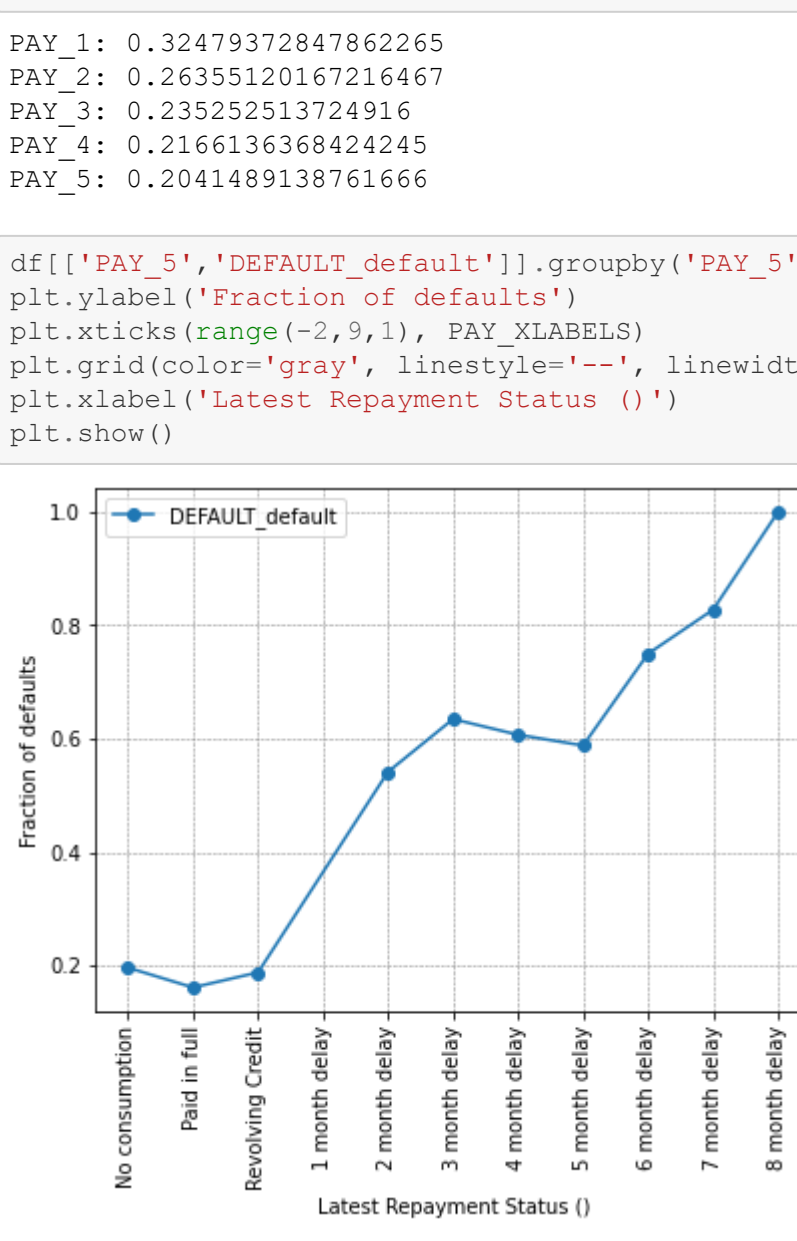
ID	EDUCATION_graduate_school	EDUCATION_high_school \	
LIMIT_BAL	
MARRIAGE	
AGE	
PAY_1	
PAY_2	
PAY_3	
PAY_4	
PAY_5	
PAY_6	
BILL_AMT1	
BILL_AMT2	
BILL_AMT3	
BILL_AMT4	
BILL_AMT5	
BILL_AMT6	
PAY_AMT1	
PAY_AMT2	
PAY_AMT3	
PAY_AMT4	
PAY_AMT5	
PAY_AMT6	
SEX_female	
SEX_male	
EDUCATION_graduate_school	
EDUCATION_high_school	
EDUCATION_university	
EDUCATION_university	
DEFAULT_default	
DEFAULT_not_default	

ID	EDUCATION_graduate_school	EDUCATION_high_school \	
LIMIT_BAL	
MARRIAGE	
AGE	
PAY_1	
PAY_2	
PAY_3	
PAY_4	
PAY_5	
PAY_6	
BILL_AMT1	
BILL_AMT2	
BILL_AMT3	
BILL_AMT4	
BILL_AMT5	
BILL_AMT6	
PAY_AMT1	
PAY_AMT2	
PAY_AMT3	
PAY_AMT4	
PAY_AMT5	
PAY_AMT6	
SEX_female	
SEX_male	
EDUCATION_graduate_school	
EDUCATION_high_school	
EDUCATION_university	
EDUCATION_university			

In [40]:
df[['PAY_1','DEFULT_default']].groupby('PAY_1').mean().plot(kind='bar',stacked=True, rot=90)
plt.ylabel('Fraction in default')
plt.xticks(range(-2,9,1), PAY_XLABELS)
plt.xticks(range(0,11,1), PAY_XLABELS3) #for some reason, adding xticks as bar forces the axis to start at "0"
plt.legend(['PAY_1', 'label1', 'label2', 'label3'])
plt.savefig('PAY_1_bar.png')
plt.show()



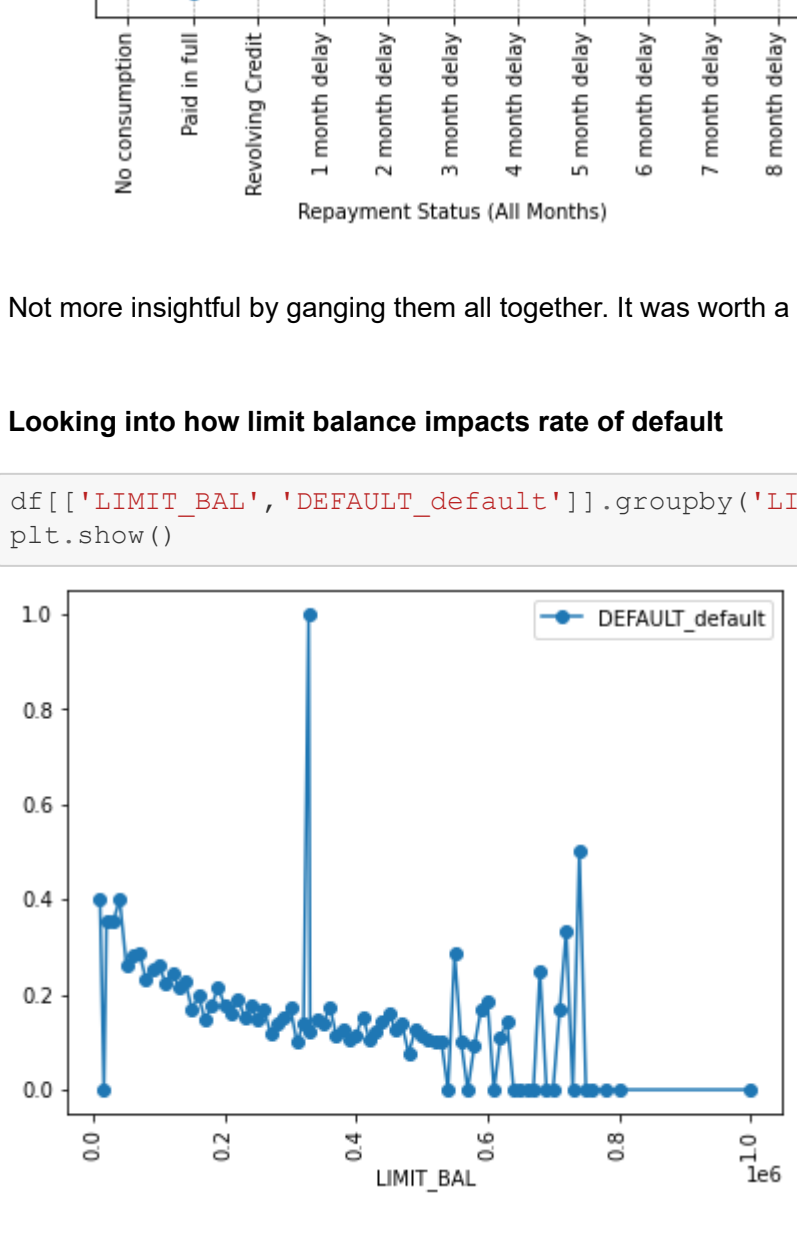
In [41]:
df[['PAY_1','DEFULT_default']].groupby('PAY_1').mean().plot(kind='line', rot=90, style='o-',legend=Non)
plt.ylabel('Fraction in default')
plt.xticks(range(-2,9,1), PAY_XLABELS)
plt.grid(color='gray', linestyle='--', linewidth=0.5)
plt.xlabel('Latest Repayment Status (September 2009)')
plt.tight_layout()
plt.savefig('PAY_1.png')
plt.show()



In [42]:
print("First number is and second number is {}".format(df['DEFAULT_default'].corr(df['PAY_1'])))
First number is and second number is 0.32479372847862265

In [43]:
print("PAY_1: {}".format(df['DEFAULT_default'].corr(df['PAY_1'])))
print("PAY_2: {}".format(df['DEFAULT_default'].corr(df['PAY_2'])))
print("PAY_3: {}".format(df['DEFAULT_default'].corr(df['PAY_3'])))
print("PAY_4: {}".format(df['DEFAULT_default'].corr(df['PAY_4'])))
print("PAY_5: {}".format(df['DEFAULT_default'].corr(df['PAY_5'])))
PAY_1: 0.32479372847862265
PAY_2: 0.26355120167216467
PAY_3: 0.2352513724916
PAY_4: 0.2166136368424245
PAY_5: 0.20241489138761666

In [44]:
df[['PAY_5','DEFULT_default']].groupby('PAY_5').mean().plot(kind='line', rot=90, style='o-')
plt.ylabel('Fraction of default')
plt.xticks(range(-2,9,1), PAY_XLABELS)
plt.grid(color='gray', linestyle='--', linewidth=0.5)
plt.xlabel('Latest Repayment Status (')
plt.show()

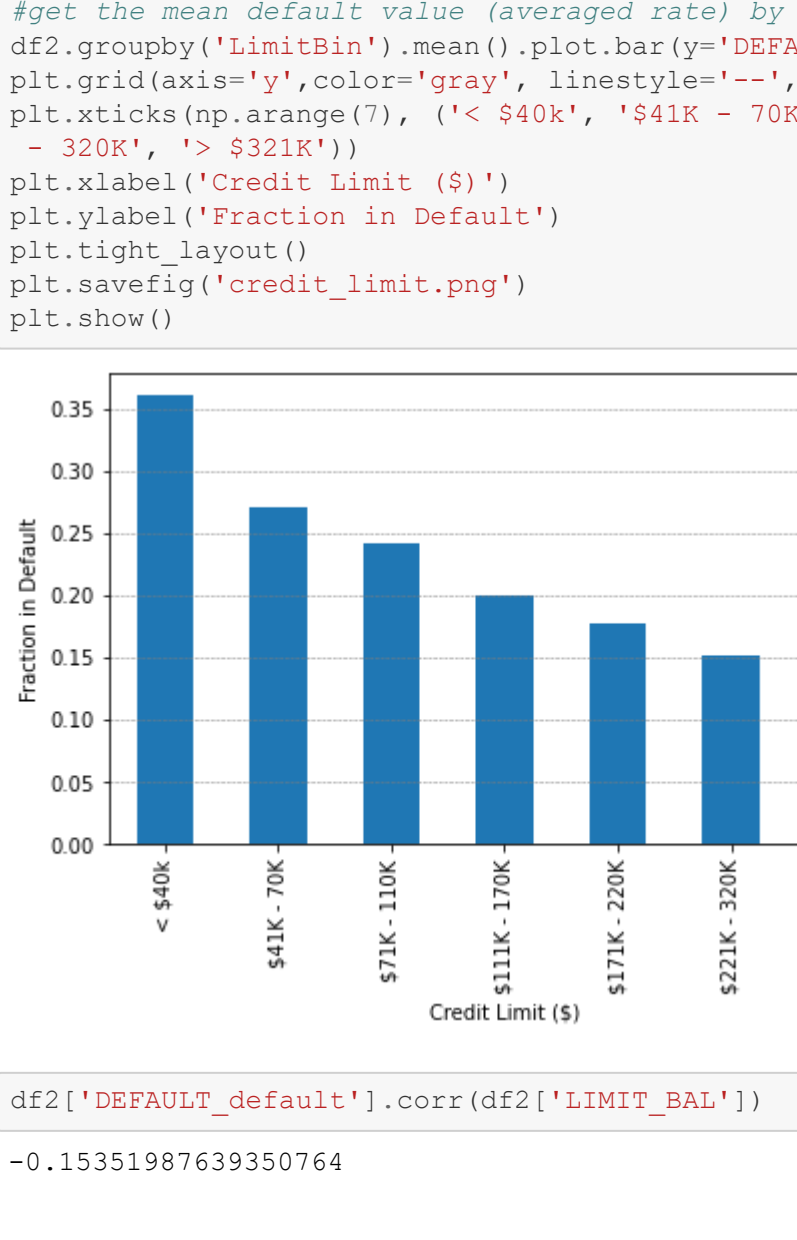


Perhaps try stacking all of the PAY_X columns for a combined result?

In [45]:
df1=df.melt(id_vars=['DEFAULT_default'], value_vars=['PAY_1','PAY_2','PAY_3','PAY_4','PAY_5'], value_name='PAY_ALL')
df1.head()

Out [45]:
DEFAULT_default variable PAY_ALL
0 1 PAY_1 2
1 1 PAY_1 -1
2 0 PAY_1 0
3 0 PAY_1 0
4 0 PAY_1 -1

In [46]:
df1[['PAY_ALL','DEFULT_default']].groupby('PAY_ALL').mean().plot(kind='line', rot=90, style='o-')
plt.ylabel('Fraction in default')
plt.xticks(range(-2,9,1), PAY_XLABELS)
plt.grid(color='gray', linestyle='--', linewidth=0.5)
plt.xlabel('Repayment Status (All Months)')
plt.show()



Not more insightful by ganging them all together. It was worth a shot

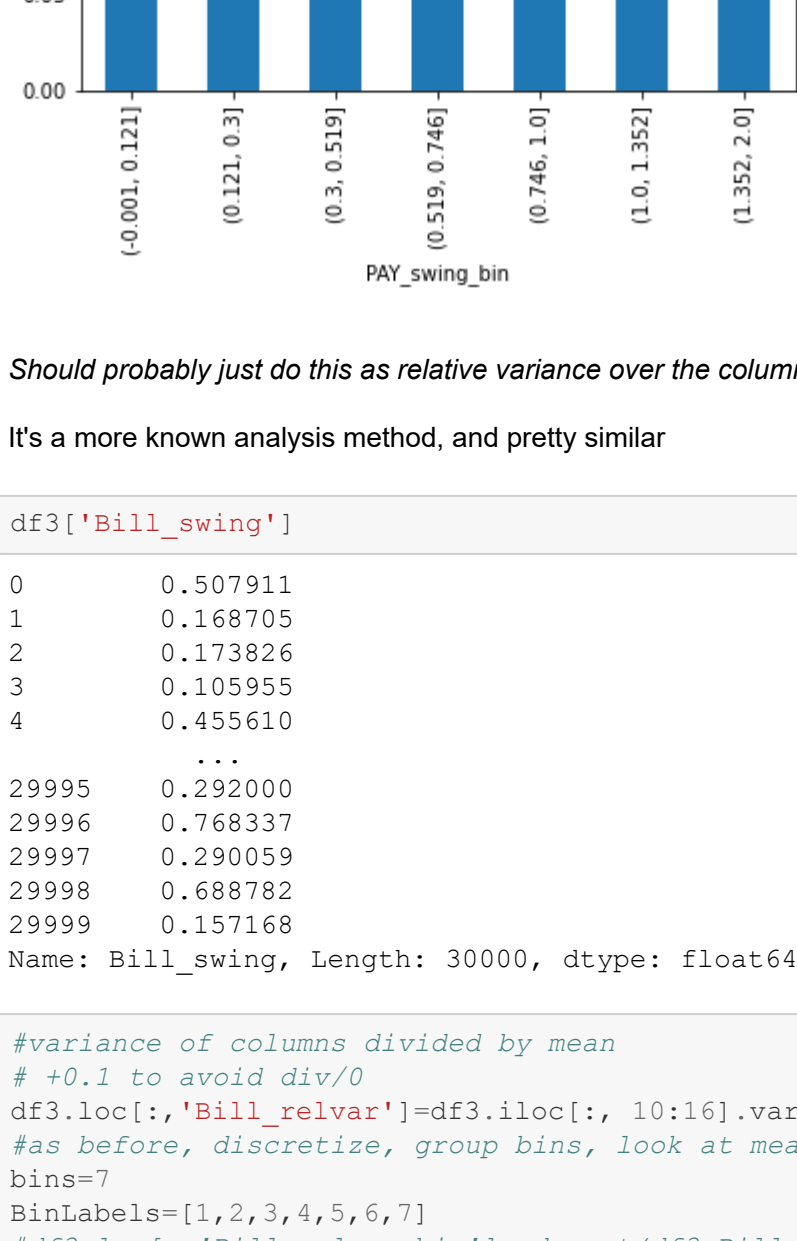
Looking into how limit balance impacts rate of default

In [47]:
df[['LIMIT_BAL','DEFULT_default']].groupby('LIMIT_BAL').mean().plot(kind='line', rot=90, style='o-')
plt.ylabel('Fraction in default')
plt.xticks(range(-2,9,1), PAY_XLABELS)
plt.grid(color='gray', linestyle='--', linewidth=0.5)
plt.xlabel('Repayment Status (All Months)')
plt.show()



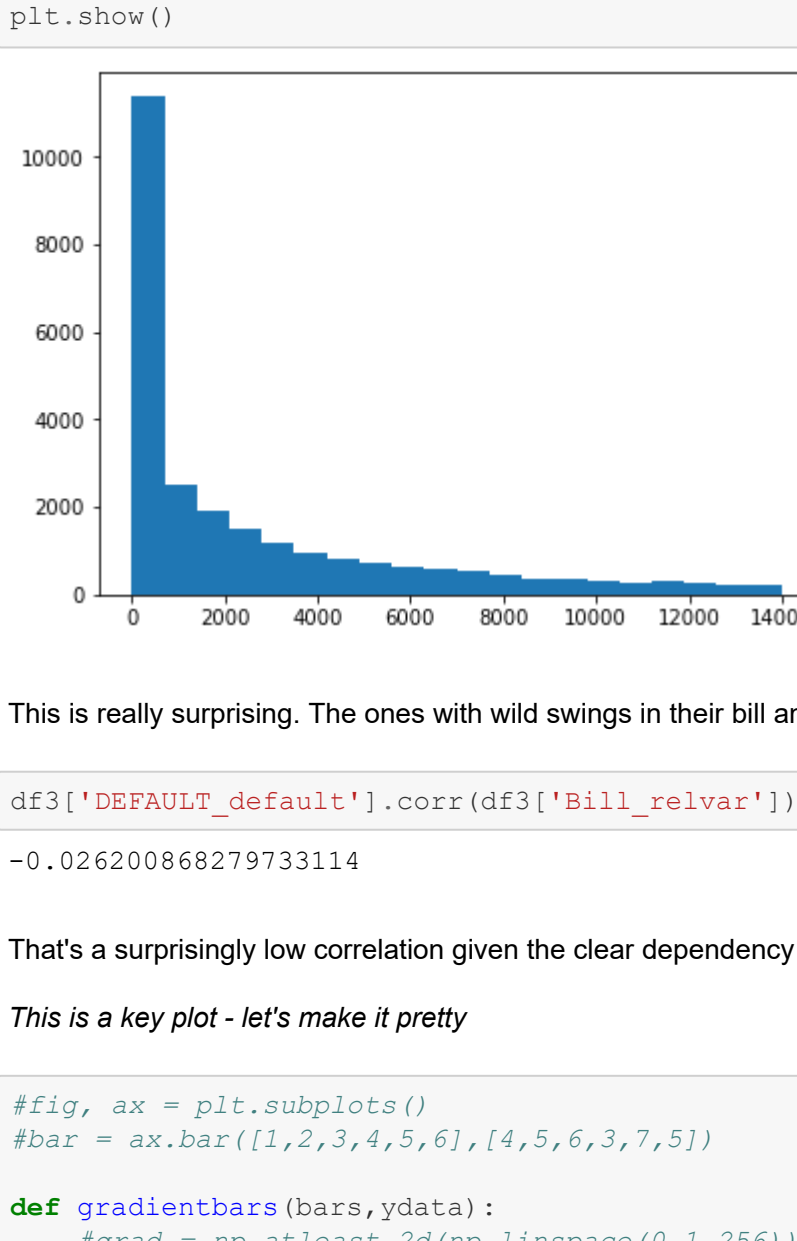
I bet it gets noisy at the end because of fewer customers with LIMIT_BAL that high. Let's check

In [48]:
plt.hist(df['LIMIT_BAL'], bins=30)
plt.show()



I'll try to make the prior one as a bar to blend a bit. Will need to discretize

In [49]:
#make a subset df with the columns we care about
df2 = df.loc[:,['LIMIT_BAL','DEFULT_default']]
#discretize bins
df2.loc[:, 'LimitBin']=pd.qcut(df2.LIMIT_BAL,bins)
#get the mean default value (averaged rate) by bin
df2.groupby('LimitBin').mean().plot.bar(y='DEFAULT_default', legend=None)
plt.xlabel('LimitBin')
plt.grid(color='gray', linestyle='--', linewidth=0.5)
plt.xticks(np.arange(7), ('< \$40k', '\$41K - 70K', '\$71K - 110K', '\$111K - 170K', '\$171K - 220K', '\$221K - 320K', '> \$321K'))
plt.ylabel('Credit Limit (\$)')
plt.tight_layout()
plt.savefig('credit_limit.png')
plt.show()

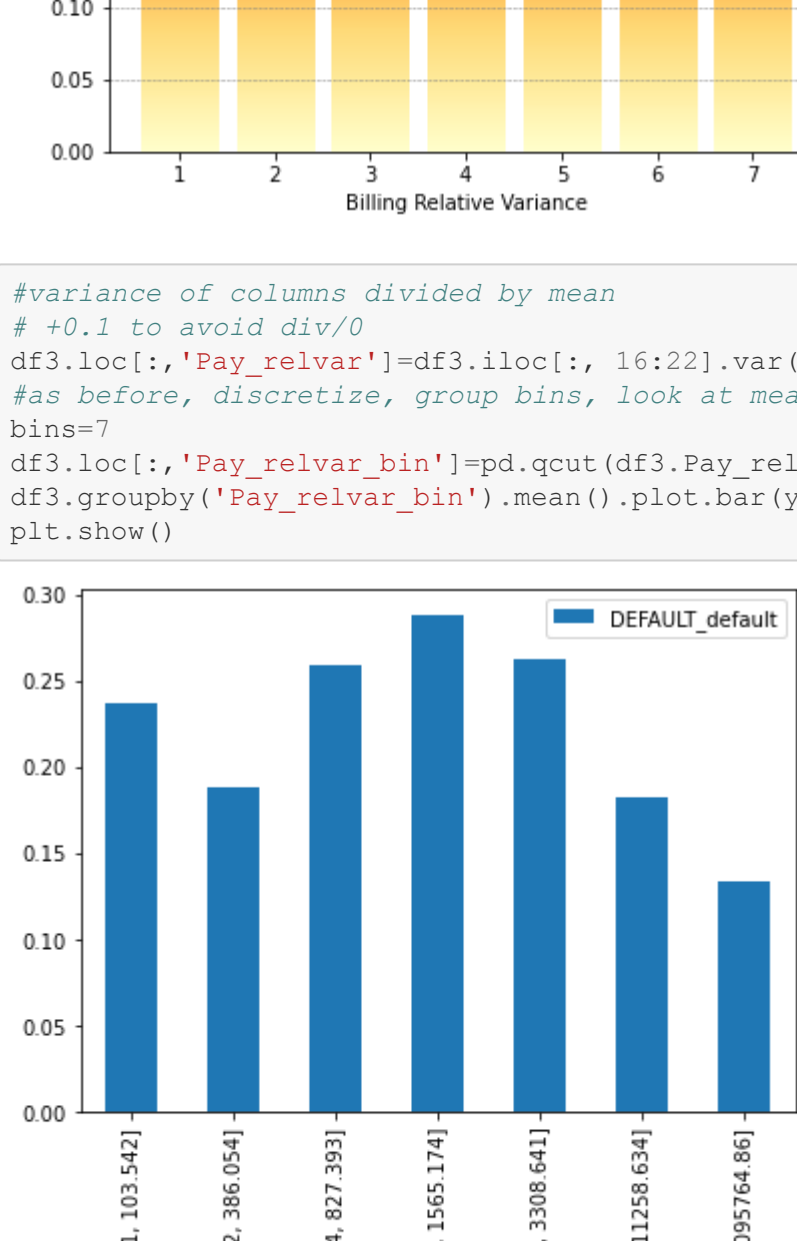


In [50]:
df2['DEFAULT_default'].corr(df2['LIMIT_BAL'])
Out [50]:
-0.15351987639350764

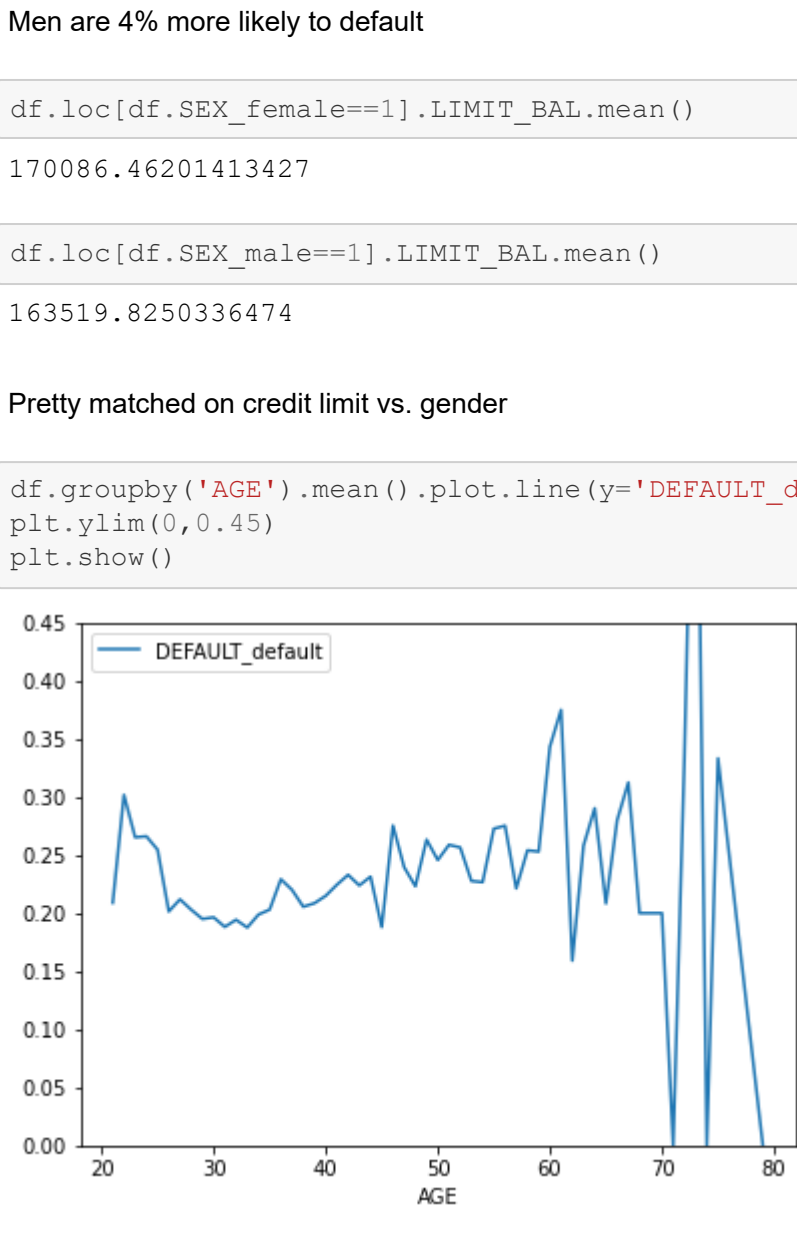
Looking into swings in bills and payments

cumulative add of % change every month

In [51]:
df3 = df.copy()
#Fails really hard for monthly bills of \$0 or really low when doing as a % change.
#Instead, do as a weighted delta. Could do as a for loop with indexes, but would be hard to read later
df3.loc[:, 'Bill_swing'] = (abs(df3.BILL_AMT1 - df3.BILL_AMT2) + \nabs(df3.BILL_AMT2 - df3.BILL_AMT3) + \nabs(df3.BILL_AMT3 - df3.BILL_AMT4) + \nabs(df3.BILL_AMT4 - df3.BILL_AMT5) + \nabs(df3.BILL_AMT5 - df3.BILL_AMT6)) / (abs(df3.iloc[:, 10:16].sum(axis=1)+1)
#had to add a +1 to avoid divide by 0 errors
#as before, discretize, group bins, look at mean defaulting rate
df3.before, df3.after, df3.group = df3.iloc[:, 10:16].var(axis=1) / (df3.iloc[:, 10:16].mean(axis=1)+1)
df3.loc[:, 'Bill_swing_bin']=pd.qcut(df3.Bill_swing,bins)
df3.groupby('Bill_swing_bin').mean().plot.bar(y='DEFAULT_default')
plt.show()



In [52]:
df3.iloc[:, 'PAY_swing'] = (abs(df3.PAY_AMT1 - df3.PAY_AMT2) + \nabs(df3.PAY_AMT2 - df3.PAY_AMT3) + \nabs(df3.PAY_AMT3 - df3.PAY_AMT4) + \nabs(df3.PAY_AMT4 - df3.PAY_AMT5) + \nabs(df3.PAY_AMT5 - df3.PAY_AMT6)) / (abs(df3.iloc[:, 16:22].sum(axis=1)+1))
#had to add a +1 to avoid divide by 0 errors
#as before, discretize, group bins, look at mean defaulting rate
df3.before, df3.after, df3.group = df3.iloc[:, 16:22].var(axis=1) / (df3.iloc[:, 16:22].mean(axis=1)+1)
df3.loc[:, 'PAY_swing_bin']=pd.qcut(df3.PAY_swing,bins)
df3.groupby('PAY_swing_bin').mean().plot.bar(y='DEFAULT_default')
plt.show()

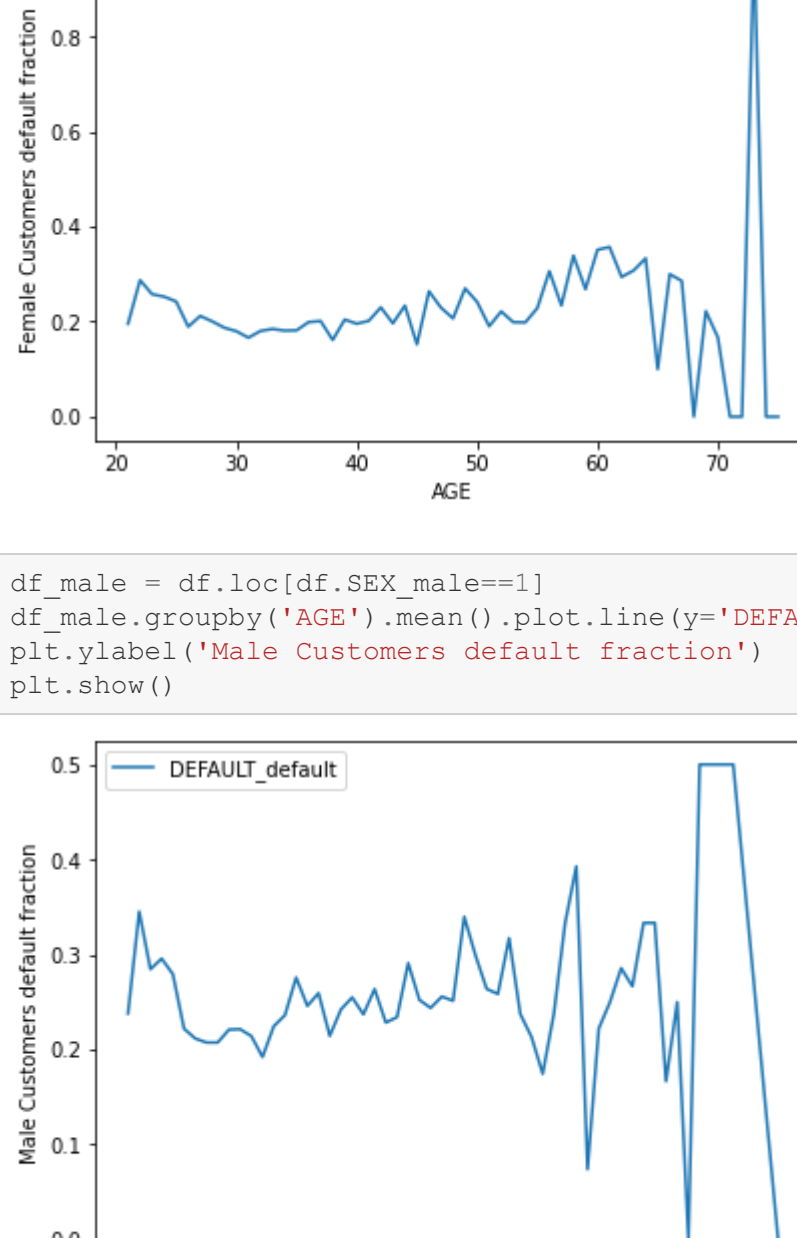


Should probably just do this as relative variance over the columns

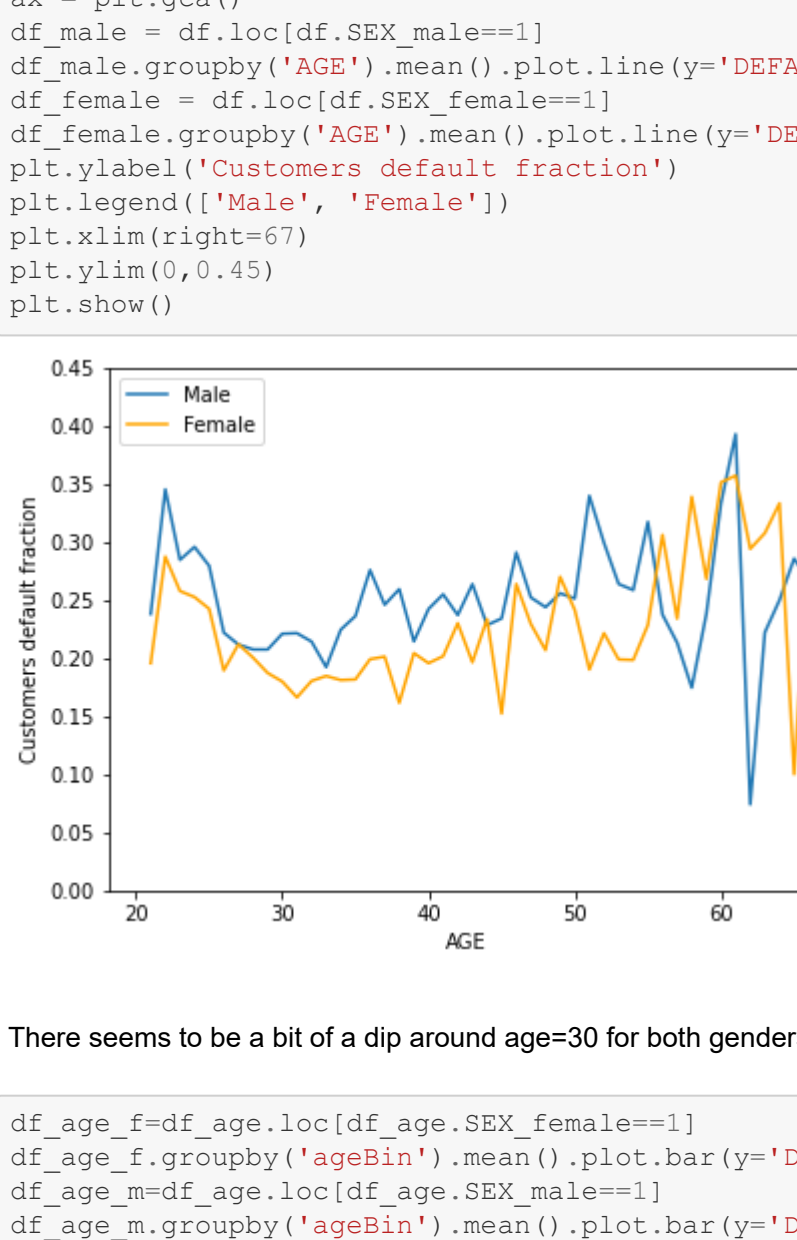
It's a more known analysis method, and pretty similar

In [53]:
df3['Bill_swing']
Out [53]:
0 0.507911
1 0.168705
2 0.173826
3 0.105955
4 0.455610
29995 0.292000
29996 0.768337
29997 0.290938
29998 0.688782
29999 0.157168
Name: Bill_swing, Length: 30000, dtype: float64

In [55]:
#variance of columns divided by mean
+0.1 to avoid div/0
df3.iloc[:, 'Bill_relvar'] = df3.iloc[:, 10:16].var(axis=1) / (df3.iloc[:, 10:16].mean(axis=1)+1)
#as before, discretize, group bins, look at mean defaulting rate
df3.before, df3.after, df3.group = df3.iloc[:, 10:16].var(axis=1) / (df3.iloc[:, 10:16].mean(axis=1)+1)
df3.loc[:, 'Bill_relvar_bin']=pd.qcut(df3.Bill_relvar,bins)
df3.groupby('Bill_relvar_bin').mean().plot.bar(y='DEFAULT_default')
plt.grid(color='gray', linestyle='--', linewidth=0.5)
plt.xticks(np.arange(7), ('< \$40k', '\$41K - 70K', '\$71K - 110K', '\$111K - 170K', '\$171K - 220K', '\$221K - 320K', '> \$321K'))
plt.xlabel('Billing Relative Variance (relative swing in bill amounts month-to-month)')
plt.ylabel('Billing in Default')
plt.show()



In [56]:
plt.hist(df3['Bill_relvar'], bins=20, range=(0, 14000))
plt.show()



This is really surprising. The ones with wild swings in their bill amount tend to default less

In [57]:
df3['DEFAULT_default'].corr(df3['Bill_relvar'])
Out [57]:
-0.026200868279733114

That's a surprisingly low correlation given the clear dependency on the plot.

This is a key plot - let's make it pretty

In [58]:
#fig, ax = plt.subplots()
#bar = ax.bar([1,2,3,4,5,6],[4,5,6,3,7,5])
def gradientbars(bars,ydata):
#grad = np.linspace(0,1,256).T
ax = bars[0].axes
lim = ax.get_xlim()*ax.get_ylim()
for bar in bars:
bar.set_zorder(1)
bar.set_facecolor("none")
x,y = bar.get_xy()
w,h = bar.get_width(), bar.get_height()
grad = np.linspace(2d(np.linspace(0,1,h,max(ydata,256))).T
ax.imshow(grad, extent=(x,x+w,y,y+h), origin='lower', aspect='auto', zorder=0, norm=cm.colors.Norm(vmin=0,vmax=1), cmap=plt.get_cmap('YlOrRd'))
ax.axis(lim)
gradientbars(bar, [4,5,6,3,7,5])

In [59]:
B_relvar_df = df3.groupby('Bill_relvar_bin', as_index=False).mean()
df3.loc[:, 'Pay_relvar'] = df3.iloc[:, 16:22].var(axis=1) / (df3.iloc[:, 16:22].mean(axis=1)+1)
df3.before, df3.after, df3.group = df3.iloc[:, 16:22].var(axis=1) / (df3.iloc[:, 16:22].mean(axis=1)+1)
df3.loc[:, 'Pay_relvar_bin']=pd.qcut(df3.Pay_relvar,bins)
df3.groupby('Pay_relvar_bin').mean().plot.bar(y='DEFAULT_default')
plt.grid(color='gray', linestyle='--', linewidth=0.5)
plt.xticks(np.arange(7), ('< \$40k', '\$41K - 70K', '\$71K - 110K', '\$111K - 170K', '\$171K - 220K', '\$221K - 320K', '> \$321K'))
plt.xlabel('Billing Relative Variance')
plt.ylabel('Fraction in Default')
gradientbars(bars,ydata)
plt.savefig('Bill_relvar.png')
plt.show()



Let's look for anything else... gender perhaps

In [61]:
df.loc[df['SEX_female']==1].DEFAULT_default.mean()
Out [61]:
0.20776280918727916

In [62]:
df.loc[df['SEX_male']==1].DEFAULT_default.mean()
Out [62]:
0.2416722745628841

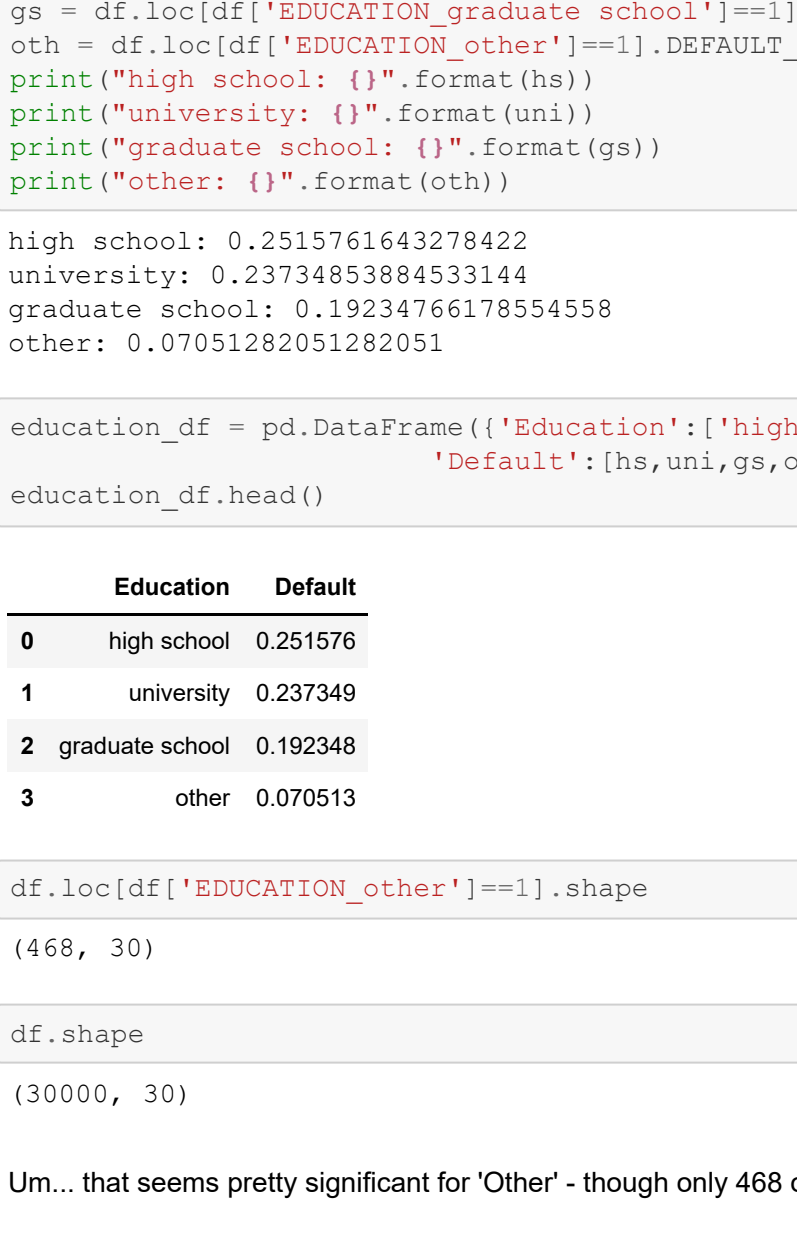
Men are 4% more likely to default

In [63]:
df.loc[df['SEX_female']==1].LIMIT_BAL.mean()
Out [63]:
170086.46201413427

In [64]:
df.loc[df['SEX_male']==1].LIMIT_BAL.mean()
Out [64]:
163519.8250336474

Pretty matched on credit limit vs. gender

In [65]:
df.groupby('AGE').mean().plot.line(y='DEFAULT_default')
plt.ylim(0,0.45)
plt.show()



In [66]:
df_age=df.loc[:,['AGE','SEX_female','SEX_male','DEFAULT_default','DEFAULT_not default']]
df_age_f=df.groupby('AGE').mean().plot.bar(y='DEFAULT_default',subplots=True,color='orange')
df_age_m=df.groupby('AGE').mean().plot.bar(y='DEFAULT_default',subplots=True,color='blue')
df_age_f.groupby('AGE').mean().plot.bar(y='DEFAULT_default',subplots=True,color='orange')
df_age_m.groupby('AGE').mean().plot.bar(y='DEFAULT_default',subplots=True,color='blue')
plt.xticks((0,1,2,3,4,5,6),plt_ageBinLabels)
plt.xlabel('Age Group')
plt.ylabel('Fraction of Default')
plt.show()

In [67]:
df_female = df.loc[df['SEX_female']==1]
df_female.groupby('AGE').mean().plot.line(y='DEFAULT_default')
plt.ylabel('Female Customers default fraction')
plt.show()

In [68]:
df_male = df.loc[df['SEX_male']==1]
df_male.groupby('AGE').mean().plot.line(y='DEFAULT_default')
plt.ylabel('Male Customers default fraction')
plt.show()

Let's overlay these for male and female, and plot in seahorn

In [69]:
ax = plt.gca()
df_male = df.loc[df['SEX_male']==1]
df_male.groupby('AGE').mean().plot.line(y='DEFAULT_default',ax=ax)
df_female = df.loc[df['SEX_female']==1]
df_female.groupby('AGE').mean().plot.line(y='DEFAULT_default',ax=ax,color='orange')
plt.ylabel('Customers default fraction')
plt.legend(['Male', 'Female'])
plt.xlim(right=80)
plt.ylim(0,0.45)
plt.show()

There seems to be a bit of a dip around age=30 for both genders. This may look better as a binned bar chart

In [70]:
df_age_f=df_age.loc[df['AGE','SEX_female']==1]
df_age_f.groupby('AGE').mean().plot.bar(y='DEFAULT_default',subplots=True,color='orange')
df_age_m=df_age.loc[df['AGE','SEX_male']==1]
df_age_m.groupby('AGE').mean().plot.bar(y='DEFAULT_default',subplots=True,color='blue')
plt.xticks((0,1,2,3,4,5,6),plt_ageBinLabels)
plt.xlabel('Age Group')
plt.ylabel('Fraction of Default')
plt.show()

In [71]:
df_template=df_age_f.groupby('ageBin').mean()
df_template.groupby('ageBin').mean().plot.bar(y='DEFAULT_default',subplots=True,color='orange')
df_template.groupby('ageBin').mean().plot.bar(y='DEFAULT_default',subplots=True,color='blue')
df_template.groupby('ageBin').mean().plot.bar(y='DEFAULT_default',subplots=True,color='orange')
df_template.groupby('ageBin').mean().plot.bar(y='DEFAULT_default',subplots=True,color='blue')
plt.xticks((0,1,2,3,4,5,6),plt_ageBinLabels)
plt.xlabel('Age Group')
plt.ylabel('Fraction of Default')
plt.legend(['Male', 'Female'])
plt.tight_layout()
plt.savefig('age_gender_bar.png')
plt.show()

In [72]:
df_male.groupby('AGE').mean().shape
Out [72]:
(55, 29)

In [73]:
df_m = df[['AGE','SEX_female','SEX_male']]

In [74]:
df_m.groupby('AGE')['SEX_female'].value_counts()

Out [74]:
AGE SEX_female
21 1 46
0 21
2 1 421
0 139
23 1 671
73 1 2
74 1 1
75 0 2
421 1
79 0 1
Name: SEX_female, Length: 110, dtype: int64

In [75]:
df_m.groupby('AGE')['SEX_female'].value_counts().unstack(level=1).plot.bar(stacked=True, figsize=(8,5))
plt.legend(['Male', 'Female'])
plt.show()

Let's take a quick look at education to see if there's anything there. Heatmap says no...

In [76]:
h = df.loc[df['EDUCATION_high school']==1].DEFAULT_default.mean()
u = df.loc[df['EDUCATION_university']==1].DEFAULT_default.mean()
g = df.loc[df['EDUCATION_graduate school']==1].DEFAULT_default.mean()
o = df.loc[df['EDUCATION_other']==1].DEFAULT_default.mean()
print("high school: {}".format(h))
print("university: {}".format(u))
print("graduate school: {}".format(g))
print("other: {}".format(o))
high school: 0.2519761643278422
university: 0.23734838484533144
graduate school: 0.19234766178554598
other: 0.07051282051282051

In [77]:
education_df = pd.DataFrame({'Education': ['high school', 'university', 'graduate school', 'other'],
'Default': [h,u,g,o]})
education_df.head()

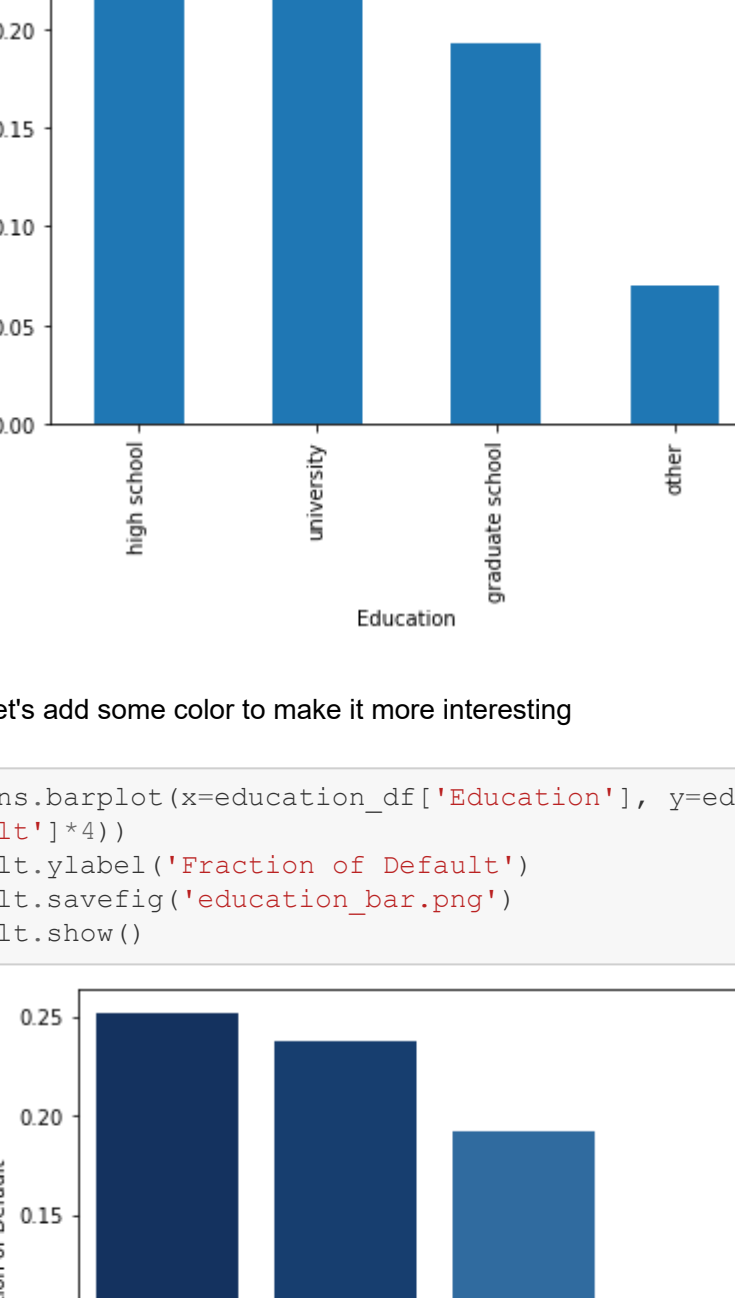
Out [77]:
Education Default
0 high school 0.251976
1 university 0.237348
2 graduate school 0.192348
3 other 0.070513

In [78]:
df.loc[df['EDUCATION_other']==1].shape
Out [78]:
(468, 30)

In [79]:
df.shape
Out [79]:
(30000, 30)

Uhm... that seems pretty significant for 'Other' - though only 468 of them out of 30000! 1.5% of population


```
[80]: education_df.plot.bar(x='Education', y='Default')
      plt.show()
```



Let's add some color to make it more interesting

```
In [81]: sns.barplot(x=education_df['Education'], y=education_df['Default'], palette=cm.Blues(education_df['Defa
      plt.show()
      plt.ylabel('Fraction of Default')
      plt.savefig('education_bar.png')
      plt.show()
```



```
In [ ]:
```