



In [41]: #Turn the Loan Status into a binary variable to see correlation between that and number of days until the loan is paid back

In [42]: #Use this table to work with
model_data = clean_dti_data.copy()
model_data

	Amount_Requested	Risk_Score	debt_to_income_ratio	State	Employment_Length	Purpose	Loan_Status
0	5000.0	669.0	21.80	OK	8	Other	Accepted
1	15000.0	704.0	18.29	FL	2	Debt consolidation	Accepted
2	11200.0	669.0	43.97	NH	0	Medical expenses	Rejected
3	25000.0	669.0	12.89	AL	10	Debt consolidation	Accepted
4	3000.0	764.0	0.58	WA	9	Major purchase	Accepted
...
1081515	3000.0	535.0	1.54	CT	5	Medical expenses	Accepted
1081516	6000.0	580.0	5.55	TX	0	Car financing	Accepted
1081517	3000.0	561.0	0.00	OR	0	Debt consolidation	Accepted
1081518	2500.0	573.0	8.55	NY	0	Other	Accepted
1081519	2500.0	567.0	0.00	AL	2	Car financing	Accepted

1034236 rows x 7 columns

In [43]: columns_to_be_encoded=['Loan_Status', 'State', 'Purpose']
for column in columns_to_be_encoded:
 le=preprocessing.LabelEncoder()
 le.fit(list(model_data[column].values))
 model_data[column]=le.transform(list(model_data[column].values))
with open('{}_{}_pickle'.format(column, 'wb')) as file:
 pickle.dump(le, file, pickle.HIGHEST_PROTOCOL)
model_data

	Amount_Requested	Risk_Score	debt_to_income_ratio	State	Employment_Length	Purpose	Loan_Status
0	5000.0	669.0	21.80	36	8	17	Accepted
1	15000.0	704.0	18.29	9	2	4	Accepted
2	11200.0	669.0	43.97	30	0	14	Rejected
3	25000.0	669.0	12.89	1	10	4	Accepted
4	3000.0	764.0	0.58	47	9	12	Accepted
...
1081515	3000.0	535.0	1.54	6	5	15	Accepted
1081516	6000.0	580.0	5.55	43	0	2	Accepted
1081517	3000.0	561.0	0.00	37	0	4	Accepted
1081518	2500.0	573.0	8.55	34	0	17	Accepted
1081519	2500.0	567.0	0.00	1	2	2	Accepted

1034236 rows x 7 columns

In [44]: #Independent Variable
x=model_data.drop(['Loan_Status'], axis=1)
y=model_data['Loan_Status']

In [45]: new_model_corr = model_data.corr()
plt.figure(figsize=(12, 6))
sns.heatmap(new_model_corr, vmin=-1, vmax=1, annot=True, cmap='BrBG')

Out[45]: <AxesSubplot>

In [46]: X

Out[46]:

	Amount_Requested	Risk_Score	debt_to_income_ratio	State	Employment_Length	Purpose
0	5000.0	669.0	21.80	36	8	17
1	15000.0	704.0	18.29	9	2	4
2	11200.0	669.0	43.97	30	0	14
3	25000.0	669.0	12.89	1	10	4
4	3000.0	764.0	0.58	47	9	12
...
1081515	3000.0	535.0	1.54	6	5	15
1081516	6000.0	580.0	5.55	43	0	2
1081517	3000.0	561.0	0.00	37	0	4
1081518	2500.0	573.0	8.55	34	0	17
1081519	2500.0	567.0	0.00	1	2	2

1034236 rows x 6 columns

In [47]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=.33, random_state=42)

In [48]: #DecisionTreeClassifier
dtree=DecisionTreeClassifier(random_state=42)
dtree.fit(X_train, y_train)

Out[48]: DecisionTreeClassifier(random_state=42)

In [49]: y_pred=dtree.predict(X_test)

In [50]: #What is the perfect prediction
confusion_matrix(y_test, y_pred)

Out[50]: array([[14568, 3656],
[3561, 188423]], dtype=int64)

In [51]: #Precision Score
print("Precision Score", average_precision_score(y_test, y_pred), "--What percentage of positive cases did you catch?
print("Recall", recall_score(y_test, y_pred), "--What percentage of positive cases did you catch?
print(cross_val_score(dtree, X_test, y_test, cv=5, scoring = "accuracy"))
print(f1_score(y_test, y_pred))

Precision Score 0.973204489079665 --What percentage of predictions were correct?
Recall 0.9814515792982749 --What percentage of positive cases did you catch?
[0.97505127 0.97713156 0.97427483 0.97855993 0.97765862]
0.9812088120959321

In []:

In []:

In []: